

## Appendix

### A. Tensors and Higher Order Singular Value Decomposition

As a higher order analogue of Singular Value Decomposition (SVD) for matrices, Higher Order Singular Value Decomposition (HOSVD) for tensors is the main tool to develop our theorems. In this section, we give a quick review of tensors, and introduce the essential part of HOSVD which helps better understand the theorems. The contents here are based on (De Lathauwer et al., 2000) and (Kolda & Bader, 2009).

**Basic Notations** Below is a table for different notations.

Type	Notation	Examples
Tensor	Boldface Euler script letter	$\mathcal{A}$
Matrix	Boldface capital letter	$\mathbf{A}$
Vectors	Boldface lowercase letters	$\mathbf{a}$
Scalars	Lowercase letters	$a$

The *order* of a tensor is the number of its dimensions, which is also called modes. We use subscripts on corresponding notations to denote a specific slice of a tensor. For example,  $\mathcal{A} \in \mathbb{R}^{4 \times 2 \times 4 \times 5}$  is an order 4 tensor, and  $a_{ijkl}$  denotes the  $(i, j, k, l)$ -element of  $\mathcal{A}$ . The vector coordinated along the first dimension is denoted as  $\mathbf{a}_{:jkl}$ , and the matrix coordinated along the first and the second dimension is denoted as  $\mathbf{A}_{::kl}$ .

The norm of a tensor  $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_K}$  is defined by

$$\|\mathcal{A}\| = \sqrt{\sum_{i_1=1}^{I_1} \dots \sum_{i_K=1}^{I_K} a_{i_1 \dots i_K}^2}, \quad (30)$$

Let  $\mathcal{A} \in \mathbb{R}^{I_1 \times \dots \times I_K}$  and  $\mathcal{B} \in \mathbb{R}^{J_1 \times \dots \times J_{K'}}$  be two tensors, then the outer product of  $\mathcal{A}$ ,  $\mathcal{B}$  is a tensor in  $\mathbb{R}^{I_1 \times \dots \times I_K \times J_1 \times \dots \times J_{K'}}$ , which is defined as

$$(\mathcal{A} \otimes \mathcal{B})_{i_1 \dots i_K j_1 \dots j_{K'}} = a_{i_1 \dots i_K} b_{j_1 \dots j_{K'}}. \quad (31)$$

**Tensor reshaping** The tensor reshaping refers that one can reshape a tensor into a matrix, or a matrix into a tensor based on some specific rules. The process of reshaping a tensor into a matrix is called tensor flattening, while reshaping a matrix into a tensor is called tensorisation. Tensor reshaping is the core idea of HOSVD.

**Definition A.1.** Let  $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_K}$  be an order  $K$  tensor. The mode- $k$  flattening of  $\mathcal{A}$  is denoted as  $\mathbf{A}_{(k)} \in \mathbb{R}^{I_k \times (I_1 \dots I_{k-1} I_{k+1} \dots I_K)}$ , where the  $(i_1, i_2, \dots, i_K)$ -element of  $\mathcal{A}$  is mapped to the  $(i_k, j)$ -element of the matrix  $\mathbf{A}_{(k)}$  with

$$j = 1 + \sum_{\substack{s=1 \\ s \neq k}}^K \left( (i_s - 1) \prod_{\substack{s'=1 \\ s' \neq k}}^{s-1} I_{s'} \right). \quad (32)$$

That is, the columns of  $\mathbf{A}_{(k)}$  are actually the vectors  $\mathbf{a}_{i_1 \dots i_{k-1} : i_{k+1} \dots i_K}$ .

**Definition A.2.** Let  $\mathbf{A} \in \mathbb{R}^{I_k \times (I_1 I_2 \dots I_{K-1})}$  be a matrix. Then  $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_K}$  is the tensorisation of  $\mathbf{A}$  if its mode- $K$  flattening equals to  $\mathbf{A}$ .

We illustrate the above two definitions by an example.

**Example A.1.** Consider an order 3 tensor  $\mathcal{A} \in \mathbb{R}^{4 \times 3 \times 2}$  such that

$$\mathbf{A}_{::1} = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \\ 10 & 11 & 12 \end{pmatrix}, \quad \mathbf{A}_{::2} = \begin{pmatrix} 13 & 14 & 15 \\ 16 & 17 & 18 \\ 19 & 20 & 21 \\ 22 & 23 & 24 \end{pmatrix}. \quad (33)$$

Then

$$\begin{aligned}\mathbf{A}_{(1)} &= \begin{pmatrix} 1 & 2 & 3 & 13 & 14 & 15 \\ 4 & 5 & 6 & 16 & 17 & 18 \\ 7 & 8 & 9 & 19 & 20 & 21 \\ 10 & 11 & 12 & 22 & 23 & 24 \end{pmatrix}, \\ \mathbf{A}_{(2)} &= \begin{pmatrix} 1 & 4 & 7 & 10 & 13 & 16 & 19 & 22 \\ 2 & 5 & 8 & 11 & 14 & 17 & 20 & 23 \\ 3 & 6 & 9 & 12 & 15 & 18 & 21 & 24 \end{pmatrix}, \\ \mathbf{A}_{(3)} &= \begin{pmatrix} 1 & 4 & 7 & 10 & \cdots & 3 & 6 & 9 & 12 \\ 13 & 16 & 19 & 22 & \cdots & 15 & 18 & 21 & 24 \end{pmatrix}.\end{aligned}\tag{34}$$

Conversely,  $\mathcal{A}$  is the tensorisation of  $\mathbf{A}_{(K)}$  in  $\mathbb{R}^{4 \times 3 \times 2}$ .

Recall the tensorisation  $T_{l^K}(\mathbf{a})$  we defined in the paper, where  $\mathbf{a} \in \mathbb{R}^{l^K}$  is a vector. In this case, we first rearrange the vector  $\mathbf{a} \in \mathbb{R}^{l^K}$  into a matrix  $\mathbf{A} \in \mathbb{R}^{l \times l^{K-1}}$  according to row major ordering, then  $T_{l^K}(\mathbf{a})$  is the tensor in  $\mathbb{R}^{l \times l \times \cdots \times l}$  defined as the tensorisation of  $\mathbf{A}$ .

### Singular values and the rank of tensors

**Definition A.3.** The singular values of a tensor  $\mathcal{A} \in \mathbb{R}^{I_1 \times \cdots \times I_K}$  are defined as

$$\left\{ \sigma_{i_k}^{(k)} \right\}, \quad k = 1, \dots, K,\tag{35}$$

where  $\left\{ \sigma_{i_k}^{(k)} \right\}$  are the singular values of  $\mathbf{A}_{(k)}$  arising from the SVD for matrices. That is, the singular values of  $\mathcal{A}$  is a collection of all the singular values of its mode- $k$  flattening.

**Remark A.1.** The norm of  $\mathcal{A}$  satisfies

$$\|\mathcal{A}\| = \|\mathbf{A}_{(k)}\|_F = \sqrt{\sum_{i_k=1}^{r_k} \left( \sigma_{i_k}^{(k)} \right)^2}, \quad k = 1, 2, \dots, K,\tag{36}$$

where  $r_k$  is the matrix rank of  $\mathbf{A}_{(k)}$ .

There are various ways to define the rank of a tensor. Here we use the following definition.

**Definition A.4.** The rank of a tensor  $\mathcal{A}$  is defined as

$$\text{rank } \mathcal{A} := \sum_{k=1}^K r_k.\tag{37}$$

Here  $r_k$  is the matrix rank of  $\mathbf{A}_{(k)}$ , which is also called the  $k$ -rank of  $\mathcal{A}$ . Recall that the matrix rank equals to the number of its non-zero singular values, it follows from Definition A.3 that the tensor rank also equals to the number of its non-zero singular values.

**Remark A.2.** Recall that the matrix SVD can be written in the form of outer product:  $\mathbf{A} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$  with  $r = \text{rank } \mathbf{A}$ .

This can be generalised to HOSVD. For any  $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_K}$ , we have

$$\mathcal{A} = \sum_{i_1=1}^{r_1} \sum_{i_2=1}^{r_2} \cdots \sum_{i_K=1}^{r_K} s_{i_1 i_2 \dots i_K} \mathbf{u}_{i_1}^{(1)} \otimes \mathbf{u}_{i_2}^{(2)} \otimes \cdots \otimes \mathbf{u}_{i_K}^{(K)},\tag{38}$$

where  $r_k$  is the  $k$ -rank of  $\mathcal{A}$ ,  $s_{i_1 i_2 \dots i_K} \in \mathbb{R}$  and  $\mathbf{u}_{i_k}^{(k)} \in \mathbb{R}^{I_k}$ .

We end up with the following approximation property of HOSVD, which can be viewed as an analogue of Eckart-Young-Mirsky theorem for matrix SVD.

**Proposition A.5.** Let  $\mathcal{A}, \hat{\mathcal{A}} \in \mathbb{R}^{I_1 \times \dots \times I_K}$  with the  $k$ -ranks denoted as  $r_1, \dots, r_K$  and  $r'_1, \dots, r'_K$  respectively. Let  $\sigma_1^{(k)} \geq \sigma_2^{(k)} \geq \dots \geq \sigma_{r_k}^{(k)} \geq 0$  be the singular values of  $\mathcal{A}_{(k)}$ , we have

$$\inf_{\hat{\mathcal{A}}} \|\mathcal{A} - \hat{\mathcal{A}}\|^2 \leq \sum_{i_1=r'_1+1}^{r_1} (\sigma_{i_1}^{(1)})^2 + \sum_{i_2=r'_2+1}^{r_2} (\sigma_{i_2}^{(2)})^2 + \dots + \sum_{i_K=r'_K+1}^{r_K} (\sigma_{i_K}^{(K)})^2, \quad (39)$$

where the infimum is taken over  $\hat{\mathcal{A}} \in \mathbb{R}^{I_1 \times \dots \times I_K}$  such that  $r'_k \leq r_k, k = 1, \dots, K$ .

**Remark A.3.** Note that for matrix SVD, Eckart-Young-Mirsky theorem shows that the approximation error equals to the tail sum of singular values, but for HOSVD only the upper bound holds.

## B. Proofs

Now we get down to prove all the results shown in the text.

Recall that the input space  $\mathcal{X}$  is a Hilbert space, one can apply the standard representation theorem.

**Theorem B.1.** (Riesz Representation Theorem) For any continuous linear functional  $H$  defined on  $\mathcal{X}$ , there exists a unique  $\rho \in \mathcal{X}$  such that

$$H(\mathbf{x}) = \sum_{s=-\infty}^{\infty} \rho(s)^\top x(s), \quad (40)$$

and

$$\|H\| = \|\rho\|_{\mathcal{X}}. \quad (41)$$

*Proof.* See Bramwell & Kreyszig (1979), Theorem 3.8-1. □

Based on this, we prove Lemma 1.

*Proof of Lemma 1.* By Riesz Representation Theorem, for any  $t \in \mathbb{Z}$  and  $H_t \in \mathbf{H}$ , there exists a unique  $\rho_t \in \mathcal{X}$  such that

$$H_t(\mathbf{x}) = \sum_{s=-\infty}^{\infty} \rho_t(s)^\top x(s). \quad (42)$$

With the fact that  $\mathbf{H}$  is causal, we have

$$H_t(\mathbf{x}) = \sum_{s=-\infty}^t \rho_t(s)^\top x(s). \quad (43)$$

By the time homogeneity  $H_t(\mathbf{x}) = H_{t+\tau}(\mathbf{x}^{(\tau)})$  with  $\tau = -t$ , we get

$$H_t(\mathbf{x}) = \sum_{s=-\infty}^t \rho_t(s)^\top x(s) = \sum_{s=-\infty}^0 \rho_0(s)^\top x(s+t). \quad (44)$$

The conclusion follows by taking  $\rho^{(\mathbf{H})}(s) = \rho_0(-s)$ . □

Recall the example in section 4.2, where we showed that for any matrix  $\mathbf{A}$  with the rank no more than 2, there exists  $\rho^{(\hat{\mathbf{H}})} \in \mathcal{H}_{\text{CNN}}^{(2,2,\{M_k\})}$  such that  $T_{2^2}(\rho^{(\hat{\mathbf{H}})}) = \mathbf{A}$ . Now we extend this to the general  $l$  and  $K$ .

**Proposition B.2.** For any  $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K \in \mathbb{R}^l$ , we have

$$T_{lK} \left( \mathbf{w}_K \underset{l^{K-1}}{*} \mathbf{w}_{K-1} \underset{l^{K-2}}{*} \dots \underset{l}{*} \mathbf{w}_1 \right) = \mathbf{w}_K \otimes \mathbf{w}_{K-1} \otimes \dots \otimes \mathbf{w}_1. \quad (45)$$

*Proof.* We prove this by induction. When  $K = 2$ , it is the usual outer product for vectors. Suppose the conclusion holds for  $K$ , we prove for  $K + 1$ . Let  $\mathbf{f}_K := \mathbf{w}_K \ast_{l^{K-1}} \mathbf{w}_{K-1} \ast_{l^{K-2}} \cdots \ast_l \mathbf{w}_1 = (c_1, c_2, \dots, c_{l^K})$  and  $\mathbf{w}_{K+1} = (w_1, w_2, \dots, w_l)$ , then

$$\begin{aligned} (\mathbf{w}_{K+1} \ast_{l^K} \mathbf{f}_K)(t) &= \sum_{s \in \mathbb{N}} w_{K+1}(s) f_K(t - l^K s) \\ &= w(0) f_K(t) + w(1) f_K(t - l^K) + \cdots + w(l-1) f_K(t - (l-1)l^K) \\ &= w_1 f_K(t) + w_2 f_K(t - l^K) + \cdots + w_l f_K(t - (l-1)l^K). \end{aligned} \quad (46)$$

The right hand side is non-zero only when  $t = 0, 1, 2, \dots, l^{K+1} - 1$ , which gives

$$(\mathbf{w}_{K+1} \ast_{l^K} \mathbf{f}_K) = (w_1 c_1, w_1 c_2, \dots, w_1 c_{l^K}, w_2 c_1, w_2 c_2, \dots, w_2 c_{l^K}, \dots, w_l c_1, w_l c_2, \dots, w_l c_{l^K}) \in \mathbb{R}^{l^{K+1}}. \quad (47)$$

By the tensorisation for vectors discussed above,  $\mathbf{w}_{K+1} \ast_{l^K} \mathbf{f}_K$  can be rearranged into a matrix according to row major ordering:

$$\begin{pmatrix} w_1 c_1 & w_1 c_2 & \cdots & w_1 c_{l^K} \\ w_2 c_1 & w_2 c_2 & \cdots & w_2 c_{l^K} \\ \vdots & \vdots & \ddots & \vdots \\ w_l c_1 & w_l c_2 & \cdots & w_l c_{l^K} \end{pmatrix} \in \mathbb{R}^{l \times l^K}. \quad (48)$$

This is in fact the mode-K flattening of  $\mathbf{w}_{K+1} \otimes T_{l^K}(\mathbf{f}_K)$ . By the induction hypothesis, we conclude that

$$T_{l^{K+1}}(\mathbf{w}_{K+1} \ast_{l^K} \mathbf{f}_K) = \mathbf{w}_{K+1} \otimes T_{l^K}(\mathbf{f}_K) = \mathbf{w}_{K+1} \otimes \mathbf{w}_K \otimes \mathbf{w}_{K-1} \otimes \cdots \otimes \mathbf{w}_1. \quad (49)$$

□

**Proposition B.3.** Let  $\mathcal{A} \in \mathbb{R}^{l \times l \times \cdots \times l}$  be an order  $K$  tensor with the  $k$ -rank  $r_k$ ,  $k = 1, \dots, K$ . There exists  $\rho^{(\hat{H})} \in \mathcal{H}_{\text{CNN}}^{(l)}$  such that  $T_{l^K}(\rho^{(\hat{H})}) = \mathcal{A}$ .

*Proof.* Recall the linear CNN model

$$\begin{aligned} \mathbf{h}_{0,i} &= \mathbf{x}_i, \\ \mathbf{h}_{k+1,i} &= \sum_{j=1}^{M_k} \mathbf{w}_{kj} \ast_{l^k} \mathbf{h}_{k,j}, \\ \hat{\mathbf{y}} &= \mathbf{h}_{K+1}. \end{aligned} \quad (50)$$

By linearity, there exists  $\rho^{(\hat{H})} \in \mathcal{H}_{\text{CNN}}^{(l, K, \{M_k\})}$  such that

$$\rho^{(\hat{H})} = \sum_{i_K=1}^{M_K} \sum_{i_{K-1}=1}^{M_{K-1}} \cdots \sum_{i_1=1}^{M_1} \mathbf{w}_{K, i_K} \ast_{l^K} \mathbf{w}_{K-1, i_{K-1}} \ast_{l^{K-1}} \cdots \ast_l \mathbf{w}_{1, i_1}, \quad (51)$$

where  $\mathbf{w}_{k, i_k} \in \{\mathbf{w}_{kij}\}$  is a filter at layer  $k$ . According to Proposition B.2, we have

$$T_{l^K}(\rho^{(\hat{H})}) = \sum_{i_K=1}^{M_K} \sum_{i_{K-1}=1}^{M_{K-1}} \cdots \sum_{i_1=1}^{M_1} \mathbf{w}_{K, i_K} \otimes \mathbf{w}_{K-1, i_{K-1}} \otimes \cdots \otimes \mathbf{w}_{1, i_1}. \quad (52)$$

The conclusion now follows from Remark A.2 for sufficient large  $M_k$ . □

Now we can prove the first main theorem in the main text.

*Proof of Theorem 2.* By Lemma 1 and (12) in the main text, we have

$$\begin{aligned}
 \|\mathbf{H} - \hat{\mathbf{H}}\|^2 &= \sup_{t \in \mathbb{Z}} \|H_t - \hat{H}_t\|^2 = \sup_{t \in \mathbb{Z}} \sup_{\|\mathbf{x}\|_X \leq 1} \left| \sum_{s \in \mathbb{N}} (\boldsymbol{\rho}^{(\mathbf{H})}(s) - \boldsymbol{\rho}^{(\hat{\mathbf{H}})}(s))^\top \mathbf{x}(t-s) \right|^2 \\
 &\leq \sum_{s \in \mathbb{N}} |\boldsymbol{\rho}^{(\mathbf{H})}(s) - \boldsymbol{\rho}^{(\hat{\mathbf{H}})}(s)|^2 \\
 &= \sum_{s=0}^{l^K-1} |\boldsymbol{\rho}^{(\mathbf{H})}(s) - \boldsymbol{\rho}^{(\hat{\mathbf{H}})}(s)|^2 + \sum_{s=l^K}^{\infty} |\boldsymbol{\rho}^{(\mathbf{H})}(s)|^2 \\
 &= \left\| T_{l^K}(\boldsymbol{\rho}^{(\hat{\mathbf{H}})}) - T_{l^K}(\boldsymbol{\rho}^{(\mathbf{H})}) \right\|^2 + \sum_{s=l^K}^{\infty} |\boldsymbol{\rho}^{(\mathbf{H})}(s)|^2. \tag{53}
 \end{aligned}$$

Since  $\boldsymbol{\rho}^{(\mathbf{H})} \in \ell^2$ , we can choose  $K$  appropriately large such that the second term is less than  $\epsilon$ . According to Proposition B.3, there exists  $\mathcal{H}_{\text{CNN}}^{(l)}$  such that the first term is zero. The proof is completed.  $\square$

Lemma 3 follows from Proposition A.5 and the definition of rank for tensors.

This was denoted by  $\|\cdot\|_{l,g}$  in the main text. Here we amend the notation as we do not discuss the norm aspects.

**Proposition B.4.** *Let  $\boldsymbol{\rho}$  be a finitely supported sequence such that  $r(\boldsymbol{\rho}) \leq l^K - 1$ . Denote the singular values of  $T_{l^K}(\boldsymbol{\rho})$  by  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{l^K}$ . For any sequence  $\hat{\boldsymbol{\rho}} \in \mathbb{R}^{l^{K+1}}$  with  $\hat{\boldsymbol{\rho}}_{[0, l^K-1]} = \boldsymbol{\rho}$  and  $\hat{\boldsymbol{\rho}}_{[l^K, l^{K+1}]} = 0$ , the singular values of  $T_{l^{K+1}}(\hat{\boldsymbol{\rho}})$  are*

$$\|T_{l^K}(\boldsymbol{\rho})\| = \sigma_{l^{K+1}} \geq \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{l^K} \geq 0 = \sigma_{l^{K+1}} = \sigma_{l^{K+1}} = \dots = \sigma_{l^{K+1}-1}. \tag{54}$$

*Proof.* The singular values of  $T_{l^{K+1}}(\hat{\boldsymbol{\rho}})$  arising from mode-1 to mode- $K$  flattening are not changed, since there are only additions of zero columns. Now we consider the mode- $(K+1)$  flattening. We add zeros to an additional dimension of the tensor such that  $\hat{a}_{i_1 i_2 \dots i_{K+1}} = a_{i_1 i_2 \dots i_K}$  and  $\hat{a}_{i_1 i_2 \dots i_{K+1} j} = 0$  for  $2 \leq j \leq l$ . Then the columns of  $T_{l^{K+1}}(\hat{\boldsymbol{\rho}})_{(K+1)}$  are the vectors along the new dimension  $(\hat{a}_{i_1 i_2 \dots i_{K+1}}, 0, 0, \dots, 0) = (a_{i_1 i_2 \dots i_K}, 0, 0, \dots, 0)$ , which gives a unique non-zero singular value  $\sigma_{l^{K+1}} = \|T_{l^K}(\boldsymbol{\rho})\|$ , and  $(l-1)$  zero singular values.  $\square$

**Proposition B.5.** *Suppose the function  $g$  is monotonously decreasing and strictly positive. Then for any  $\mathbf{H} \in \mathcal{C}$  such that  $\boldsymbol{\rho}^{(\mathbf{H})}$  is finitely supported, we have  $\mathbf{H} \in \mathcal{C}^{(l,g)}$ .*

*Proof.* Let  $K' = \inf\{K \in \mathbb{N}_+ : l^K \geq r(\boldsymbol{\rho}^{(\mathbf{H})})\}$ . Then for any  $s \geq l^{K'}$ ,  $\sum_{i=s+K'+k}^{l^{K'+k}} |\sigma_i^{(K'+k)}|^2 = 0$  for all  $k \in \mathbb{N}_+$  according to Proposition B.4, which completes the proof.  $\square$

**Proposition B.6.** *Suppose  $\mathbf{H} \in \mathcal{C}$  with  $\boldsymbol{\rho}^{(\mathbf{H})}$  finitely supported. Then there exists a finitely supported decreasing  $g$  such that  $\mathbf{H} \in \mathcal{C}^{(l,g)}$ .*

*Proof.* The proof is a straightforward application of Proposition B.4.  $\square$

Next we prove the main theorem on error bound.

*Proof of theorem 4.* (i) Lower bound. Since

$$\|\mathbf{H} - \hat{\mathbf{H}}\| = \sup_{t \in \mathbb{Z}} \sup_{\|\mathbf{x}\|_X \leq 1} \left| \sum_{s \in \mathbb{N}} (\boldsymbol{\rho}^{(\mathbf{H})}(s) - \boldsymbol{\rho}^{(\hat{\mathbf{H}})}(s))^\top \mathbf{x}(t-s) \right|,$$

by taking a specific  $\mathbf{x}$  with  $x_i(0) = \frac{1}{\sqrt{d}} \operatorname{sgn}(\rho_i^{(\mathbf{H})}(t) - \rho_i^{(\hat{\mathbf{H}})}(t))$  and  $x_i(s) = 0$  otherwise, we have

$$\begin{aligned} &\geq \frac{1}{\sqrt{d}} \sup_{t \in \mathbb{N}} \left| \rho^{(\mathbf{H})}(t) - \rho^{(\hat{\mathbf{H}})}(t) \right| \\ &\geq \frac{1}{\sqrt{d}} \sup_{t \in [l^K, \infty]} \left| \rho^{(\mathbf{H})}(t) \right| \\ &\geq \frac{1}{\sqrt{d}} \sup_{t \in [l^K, \infty]} \|\rho^{(\mathbf{H})}(t)\|_2, \end{aligned} \quad (55)$$

where the inequality holds by taking a specific  $\mathbf{x}$  with  $x_i(0) = 1$  and  $x_i(s) = 0$  otherwise.

(ii) Upper bound. Following the proof of Theorem 2, or (53) gives

$$\|\mathbf{H} - \hat{\mathbf{H}}\| \leq \sum_{i=1}^d \left\| T_{l^K}(\rho_i^{(\hat{\mathbf{H}})}) - T_{l^K}(\rho_i^{(\mathbf{H})}) \right\| + \|\rho_{[l^K, \infty]}^{(\mathbf{H})}\|_2. \quad (56)$$

The remaining task is to bound the first term. Based on Lemma 3, we only need to calculate the maximum possible rank of  $T_{l^K}(\rho_i^{(\hat{\mathbf{H}})})$ . Let  $r_k$  denote the  $k$ -rank of  $T_{l^K}(\rho_i^{(\hat{\mathbf{H}})})$ . From Remark A.2, by absorbing the scalar  $s_{i_1 i_2 \dots i_K}$  into any of the vector  $\mathbf{u}_{i_k}^{(k)}$ , we have the following relationship:

$$d \prod_{k=1}^K r_k + lK \geq d \prod_{k=1}^K r_k + \sum_{k=1}^K r_k \geq \sum_{k=2}^K M_k M_{k-1}, \quad (57)$$

thus we have  $\prod_{k=1}^K r_k \geq \frac{1}{d} (\sum_{k=2}^K M_k M_{k-1} - lK) = M$ , which implies  $\operatorname{rank} T_{l^K}(\rho_i^{(\hat{\mathbf{H}})}) = \sum_{k=1}^K r_k \geq KM^{\frac{1}{K}}$ . Combined with Definition 3 gives the conclusion.  $\square$

Next we look in details the two examples about comparison between RNNs and CNNs in details.

**Example where RNNs out-perform CNNs.** We take a scalar input with  $d = 1$ . Consider a target  $\mathbf{H} \in \mathcal{C}$  with the representation  $\rho^{(\mathbf{H})}(t) = \gamma^t$ , where  $0 < \gamma < 1$ . It is easy for RNNs to approximate this target, since the representation has a power form. In fact, we have  $\mathbf{H} \in \mathcal{H}_{\text{RNN}}^{(1)}$ , i.e. a RNN with one hidden unit is sufficient to achieve an exact representation for any  $\gamma \in (0, 1)$ .

For any CNN model  $\hat{\mathbf{H}} \in \mathcal{H}_{\text{CNN}}^{(l, K, \{M_k\})}$ , based on the lower bound of Theorem 4, we have that

$$\|\mathbf{H} - \hat{\mathbf{H}}\|^2 \geq \frac{1}{\sqrt{d}} \sup_{t \in [l^K, \infty]} \|\rho^{(\mathbf{H})}(t)\|_2. \quad (58)$$

Thus, in order to achieve an approximation error with  $\|\mathbf{H} - \hat{\mathbf{H}}\| < \epsilon$ , we have

$$\sup_{t \in [l^K, \infty]} \|\rho^{(\mathbf{H})}(t)\|_2 = \gamma^{l^K} < \epsilon. \quad (59)$$

This implies  $l^K \geq \frac{\log(\epsilon)}{\log(\gamma)}$ . That is, the number of layers necessary to achieve an approximation error smaller than  $\epsilon$  diverges to infinity as  $\gamma$  approaches 1.

**Example where CNNs out-perform RNNs.** We still take a scalar input with  $d = 1$ . Consider a target  $\mathbf{H} \in \mathcal{C}$  with the representation

$$\rho^{(\mathbf{H})}(t) = \begin{cases} 1, & t = 2^K \\ 0, & \text{otherwise} \end{cases}, \quad K \in \mathbb{N}_+. \quad (60)$$

We have  $\mathbf{H} \in \mathcal{H}_{\text{CNN}}^{(2, K, \{1\})}$ . That is, a  $K$ -layer CNN with one channel per layer is sufficient to achieve an exact representation.

Recall that RNN approximates the target  $\rho^{(H)}$  with a power sum  $\rho^{(\hat{H})}(s) = c^\top W^{s-1}U$ . Suppose here  $W \in \mathbb{R}^{m \times m}$  is a diagonalisable matrix with negative eigenvalues. It has some special structures which are summarised in the following theorem.

**Theorem B.7.** (*Borwein & Erdélyi, 1996*) Let  $E_m := \left\{ u : u(t) = c_0 + \sum_{i=0}^m c_i \gamma_i^t, c_i \in \mathbb{R}, \gamma_i > 0 \right\}$ , then

$$\sup_{u \in E_m} \frac{|u'(y)|}{\sup_{s \in [a, b]} u(s)} \leq \frac{2m-1}{\min\{y-a, b-y\}}, \quad y \in (a, b). \quad (61)$$

We rewrite this theorem into a discrete form.

**Corollary B.8.** Let  $u \in E_m$ . Then

$$|u(t+1) - u(t)| \leq \frac{2m}{t} \sup_{s \geq 0} u(s). \quad (62)$$

*Proof.* By the mean value theorem, there exists  $y \in [t, t+1]$  such that  $|u(t+1) - u(t)| = |u'(y)|$ . The corollary then follows from Theorem B.7.  $\square$

For a fixed  $m$ , the changes between  $u(t+1)$  and  $u(t)$  approaches zero as  $t$  goes to infinity. This implies that if there is a sudden change in  $u$  far from the origin, the number of terms  $m$  must be large.

In order to achieve an approximation error with  $\|H - \hat{H}\| < \epsilon$ , by taking a specific  $x$  as the unit sample function where  $x(0) = 1$  and  $x(s) = 0$  otherwise, we have

$$\epsilon > \sup_t \sup_{\|x\| \leq 1} |H_t(x) - \hat{H}_t(x)| \geq \sup_t |H_t(x) - \hat{H}_t(x)| \quad (63)$$

$$= \sup_t |\rho^{(H)}(t) - c^\top W^{t-1}U| \quad (64)$$

$$= \sup_t |\rho^{(H)}(t) - u(t)|, \quad u \in E_{m^2}. \quad (65)$$

Since

$$|u(2^K + 1)| < \epsilon, \quad (66)$$

$$|u(2^K) - 1| < \epsilon, \quad (67)$$

we have

$$|u(2^K + 1) - u(2^K)| = |u(2^K + 1) - 1 - u(2^K) + 1| \quad (68)$$

$$> 1 - |u(2^K + 1)| - |u(2^K) - 1| \quad (69)$$

$$> 1 - 2\epsilon \quad (70)$$

Combining with Corollary B.8 gives

$$m^2 > 2^{K-1} \frac{1 - 2\epsilon}{1 + \epsilon}. \quad (71)$$

As  $K$  increases, the number of parameters needed for RNNs to achieve an error less than  $\epsilon$  increases exponentially, while this increment is linear for CNNs.

### C. Special structures of dilated convolutions

In this section, we discuss an interesting structure of dilated convolutions.

**Proposition C.1.** Let  $w_1, \dots, w_K$  be  $K$  filters with the same filter size  $l$ ,  $\mathbf{s} = (s_1, s_2, \dots, s_K)$  with  $0 \leq s_k \leq l - 1$ ,  $k = 1, \dots, K$ . Suppose all entries of  $w_k$  are zero except  $w_k(s_k) = 1$ , then

$$(w_K \underset{l^{K-1}}{*} w_{K-1} \underset{l^{K-2}}{*} \dots \underset{l}{*} w_1)(t) = \begin{cases} 1, & t = \hat{t} \\ 0, & \text{otherwise} \end{cases}, \quad (72)$$

where  $\hat{t} = (s_K s_{K-1} \dots s_1)_l := \sum_{i=0}^{K-1} s_{i+1} l^i$ . That is,  $\hat{t}$  can be written as a base  $l$  expansion with digits  $s_K$  to  $s_1$ .

*Proof.* We prove this by induction. When  $K = 1$ , the conclusion is obvious. Suppose the conclusion holds for  $K$ , then by (47),

$$(w_{K+1} \underset{l^K}{*} \mathbf{f}_K)(t) = (w_1 c_1, w_1 c_2, \dots, w_1 c_{l^K}, w_2 c_1, w_2 c_2, \dots, w_2 c_{l^K}, \dots, w_l c_1, w_l c_2, \dots, w_l c_{l^K}). \quad (73)$$

Suppose  $w_{K+1}(s_{K+1}) = c_m = 1$ , then the position of 1 in the above vector is  $s_{K+1} l^K + c_m$ , which means the results also holds for  $K + 1$ .  $\square$

This result allows us to construct a filter with value 1 at any specific position  $\hat{t}$ , by choosing filters  $\{w_k\}$  according to the base  $l$  expansion of  $\hat{t}$ . We illustrate this by an example.

**Example C.1.** Let  $l = 4, K = 3$  and  $w_1 = (0, 0, 0, 1), w_2 = (1, 0, 0, 0), w_3 = (0, 1, 0, 0)$ . The positions of value 1 are recorded in  $\mathbf{s} = (3, 0, 1)$ . Then

$$(w_3 \underset{4^2}{*} w_2 \underset{4^1}{*} w_1)(t) = \begin{cases} 1, & t = 19 = (103)_4 \\ 0, & \text{otherwise} \end{cases} \quad (74)$$

Based on above, one can define a notion of sparsity, which gives another sufficient condition for the exact representation.

**Definition C.2.** Suppose  $\rho^{(H)}$  is a finitely supported sequence. The sparsity of  $\rho^{(H)}$  is defined as the number of its non-zero elements, which is denoted by  $\|\rho^{(H)}\|_0$ .

**Corollary C.3.** Suppose  $\rho^{(H)}$  is a finitely supported sequence with  $r(\rho^{(H)}) \leq l^K - 1$ . Then  $K \|\rho^{(H)}\|_0$  filters are sufficient to achieve an exact representation.

*Proof.* This follows from Theorem C.1, where for each non-zero element, we can use  $K$  filters to generate it.  $\square$

In the main text, we use the condition that  $K(M + 1)^{\frac{1}{K}} \geq \text{rank } T_{l^K}(\rho^{(H)})$  to ensure an exact representation. Instead of calculating the rank of  $T_{l^K}(\rho^{(H)})$ , Corollary C.3 gives us another way to decide the number of filters sufficient to have an exact representation. This gives us the insight that if a target is sparse in the sense that  $\|\rho^{(H)}\|_0$  is small, it can also be efficiently approximated by CNNs.

**Remark C.1.** Notice that the rank of a tensor not only depends on its sparsity, but also depends on specific positions of the non-zero elements. To illustrate, consider the following example

$$\rho_1 = (1 \ 0 \ 1 \ 0) \text{ and } \rho_2 = (1 \ 0 \ 0 \ 1). \quad (75)$$

Both of them have a sparsity 2, but  $\text{rank } T_{2^2}(\rho_1) = 2$  while  $\text{rank } T_{2^2}(\rho_2) = 4$ .