

APPENDIX

A SEQUENCE OF KRONECKER PRODUCTS

The Kronecker product between a sequence of factor tensors is given by

$$\left(\mathcal{A}^{(1)} \otimes \dots \otimes \mathcal{A}^{(S)}\right)_{i_1 \dots i_N} \triangleq \mathcal{A}_{j_1^{(1)} \dots j_N^{(1)}}^{(1)} \dots \mathcal{A}_{j_1^{(S)} \dots j_N^{(S)}}^{(S)}, \quad (15)$$

where

$$j_n^{(k)} = \begin{cases} i_n - \sum_{t=1}^{k-2} j_n^{(t)} \prod_{l=t+1}^S a_n^{(l)} \bmod a_n^{(S)} & k = S, \\ \left\lfloor \frac{i_n - \sum_{t=1}^{k-1} j_n^{(t)} \prod_{l=t+1}^S a_n^{(l)}}{\prod_{l=k+1}^S a_n^{(l)}} \right\rfloor & \text{otherwise,} \end{cases} \quad (16)$$

and $\mathcal{A}^{(k)} \in \mathbb{R}^{a_1^{(k)} \times \dots \times a_N^{(k)}}$.

B ALTERNATIVE EXPANSION DIRECTIONS OF SEKRON

The proposed SeKron structure represents a given tensor $\mathcal{W} \in \mathbb{R}^{w_1 \times \dots \times w_n}$ using a sequence of Kronecker products as follows:

$$\mathcal{W} = \sum_{r_1=1}^{R_1} \mathcal{A}_{r_1}^{(1)} \otimes \sum_{r_2=1}^{R_2} \mathcal{A}_{r_1 r_2}^{(2)} \otimes \dots \otimes \sum_{r_{S-1}=1}^{R_{S-1}} \mathcal{A}_{r_1 \dots r_{S-1}}^{(S-1)} \otimes \mathcal{A}_{r_1 \dots r_{S-1}}^{(S)}. \quad (4 \text{ revisited})$$

While this decomposition structure is obtained by recursively finding the Kronecker decomposition of the right-most tensor, many alternative sequential Kronecker structures can be obtained as illustrated in Figure 4. However, such alternative structures do not fall within our SeKron framework as they cannot make use of our convolution algorithm (Algorithm 2)

C THEOREM PROOFS

Theorem 1 (Tensor Decomposition using a Sequence of Kronecker Products). *Any tensor $\mathcal{W} \in \mathbb{R}^{w_1 \times \dots \times w_N}$ can be represented by a sequence of Kronecker products between $S \in \mathbb{N}$ factors:*

$$\mathcal{W} = \sum_{r_1=1}^{R_1} \mathcal{A}_{r_1}^{(1)} \otimes \sum_{r_2=1}^{R_2} \mathcal{A}_{r_1 r_2}^{(2)} \otimes \dots \otimes \sum_{r_{S-1}=1}^{R_{S-1}} \mathcal{A}_{r_1 \dots r_{S-1}}^{(S-1)} \otimes \mathcal{A}_{r_1 \dots r_{S-1}}^{(S)}, \quad (4)$$

where $R_i \in \mathbb{N}$ and $\mathcal{A}^{(k)} \in \mathbb{R}^{R_1 \times \dots \times R_k \times a_1^{(k)} \times \dots \times a_N^{(k)}}$.

Proof. First, we define intermediate tensors

$$\mathcal{B}_{r_1 \dots r_k}^{(k)} \triangleq \sum_{r_{k+1}=1}^{R_{k+1}} \mathcal{A}_{r_1 \dots r_{k+1}}^{(k+1)} \otimes \sum_{r_{k+2}=1}^{R_{k+2}} \mathcal{A}_{r_1 \dots r_{k+2}}^{(k+2)} \otimes \dots \otimes \sum_{r_{S-1}=1}^{R_{S-1}} \mathcal{A}_{r_1 \dots r_{S-1}}^{(S-1)} \otimes \mathcal{A}_{r_1 \dots r_{S-1}}^{(S)} \quad (5 \text{ revisited})$$

Then the reconstruction error can be written as

$$\left\| \mathcal{W}_{r_1 \dots r_{k-1}}^{(k)} - \sum_{r_k=1}^{\hat{R}_k} \mathcal{A}_{r_1 \dots r_k}^{(k)} \otimes \mathcal{B}_{r_1 \dots r_k}^{(k)} \right\|_{\text{F}}^2 \quad (17)$$

where $\mathcal{W}^{(1)}$ is the initial tensor being decomposed. As described in Section 3.2, using reshaping operations

$$\mathcal{W}_{r_1 \dots r_{k-1}}^{(k)} = \text{MAT}(\text{UNFOLD}(\mathcal{W}_{r_1 \dots r_{k-1}}^{(k)}, \mathbf{d}_{\mathcal{B}_{r_1 \dots r_k}^{(k)}})), \quad (8 \text{ revisited})$$

$$\mathbf{a}_{r_1 \dots r_k}^{(k)} = \text{VEC}(\text{UNFOLD}(\mathcal{A}_{r_1 \dots r_k}^{(k)}, \mathbf{d}_{\mathcal{A}_{r_1 \dots r_k}^{(k)}})), \quad \mathbf{b}_{r_1 \dots r_k}^{(k)} = \text{VEC}(\mathcal{B}_{r_1 \dots r_k}^{(k)}), \quad (9 \text{ revisited})$$

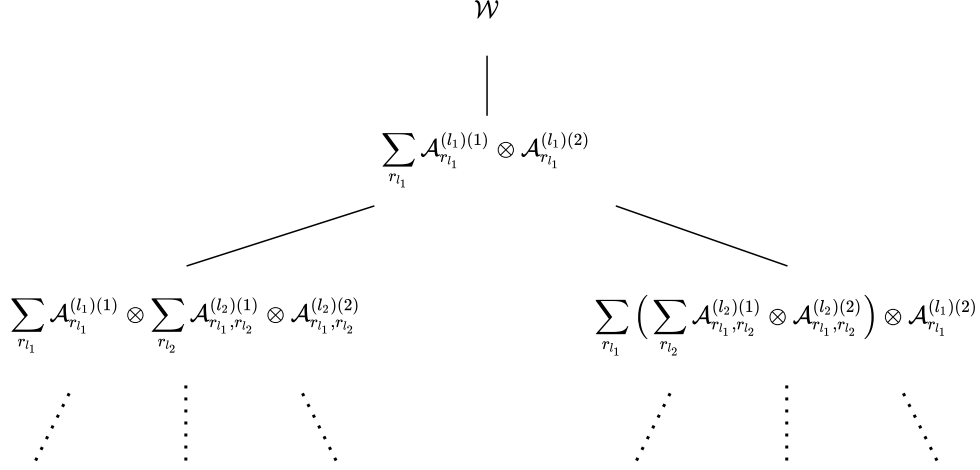


Figure 4: Illustration of alternative expansion directions using sequences of Kronecker products. SeKron structures are those which are leftmost on each level of the tree. Each node is obtained through the decomposition of a single tensor present in its parent node.

that preserve the sum of squares allows us to equivalently write the reconstruction error as

$$\left\| \mathbf{W}_{r_1 \dots r_{k-1}}^{(k)} - \sum_{r_k=1}^{\hat{R}_k} \mathbf{a}_{r_1 \dots r_k}^{(k)} \mathbf{b}_{r_1 \dots r_k}^{(k)\top} \right\|_{\text{F}}^2. \quad (18)$$

Now consider the singular value decomposition of matrix $\mathbf{W}_{r_1 \dots r_{k-1}}^{(k)}$ and let $\mathbf{u}_{r_1 \dots r_k}^{(k)}, \mathbf{v}_{r_1 \dots r_k}^{(k)}$ denote its left and right singular vectors, respectively (with the right singular vector scaled according to its corresponding singular value). Set $\mathbf{a}_{r_1 \dots r_k}^{(k)} = \mathbf{u}_{r_k}^{(k)}$ and define and define error terms

$$\delta_{r_1 \dots r_k}^{(k)} = \mathbf{v}_{r_1 \dots r_k}^{(k)} - \mathbf{b}_{r_1 \dots r_k}^{(k)}, \quad \epsilon_{r_1 \dots r_k}^{(k)} = \|\delta_{r_1 \dots r_k}^{(k)}\|. \quad (19)$$

Expanding out equation 18 reveals its recursive form

$$\left\| \mathbf{W}_{r_1 \dots r_{k-1}}^{(k)} - \sum_{r_k=1}^{\hat{R}_k} \mathbf{a}_{r_1 \dots r_k}^{(k)} \mathbf{b}_{r_1 \dots r_k}^{(k)\top} \right\|_{\text{F}}^2 = \left\| \mathbf{W}_{r_1 \dots r_{k-1}}^{(k)} - \sum_{r_k=1}^{\hat{R}_k} \mathbf{a}_{r_1 \dots r_k}^{(k)} (\mathbf{v}_{r_k}^{(k)} - \delta_{r_1 \dots r_k}^{(k)})^\top \right\|_{\text{F}}^2 \quad (20)$$

$$= \left\| \mathbf{W}_{r_1 \dots r_{k-1}}^{(k)} - \sum_{r_k=1}^{\hat{R}_k} \mathbf{a}_{r_1 \dots r_k}^{(k)} \mathbf{v}_{r_1 \dots r_k}^{(k)\top} + \sum_{r_k=1}^{\hat{R}_k} \mathbf{a}_{r_1 \dots r_k}^{(k)} \delta_{r_1 \dots r_k}^{(k)\top} \right\|_{\text{F}}^2 \quad (21)$$

$$\leq \left\| \mathbf{W}_{r_1 \dots r_{k-1}}^{(k)} - \sum_{r_k=1}^{\hat{R}_k} \mathbf{a}_{r_1 \dots r_k}^{(k)} \mathbf{v}_{r_1 \dots r_k}^{(k)\top} \right\|_{\text{F}}^2 + \sum_{r_k=1}^{\hat{R}_k} \left\| \mathbf{a}_{r_1 \dots r_k}^{(k)} \delta_{r_1 \dots r_k}^{(k)\top} \right\|_{\text{F}}^2 \quad (22)$$

$$\leq \left\| \mathbf{W}_{r_1 \dots r_{k-1}}^{(k)} - \sum_{r_k=1}^{\hat{R}_k} \mathbf{a}_{r_1 \dots r_k}^{(k)} \mathbf{v}_{r_1 \dots r_k}^{(k)\top} \right\|_{\text{F}}^2 + \sum_{r_k=1}^{\hat{R}_k} d^{(k)} \epsilon_{r_1 \dots r_k}^{(k)} \quad (23)$$

$$= \left(\sum_{r_k=\hat{R}_k+1}^{R_k} \sigma_{r_k}^2(\mathbf{W}_{r_1 \dots r_{k-1}}^{(k)}) \right) + \left(\sum_{r_k=1}^{\hat{R}_k} d^{(k)} \epsilon_{r_1 \dots r_k}^{(k)} \right) \quad (24)$$

$$= \sum_{r_k=\hat{R}_k+1}^{R_k} \sigma_{r_k}^2(\mathbf{W}_{r_1 \dots r_{k-1}}^{(k)}) + \sum_{r_k=1}^{\hat{R}_k} d^{(k)} \left\| \mathbf{v}_{r_1 \dots r_k}^{(k)} - \mathbf{b}_{r_1 \dots r_k}^{(k)} \right\|_{\text{F}}^2 \quad (25)$$

where $d^{(k)} \in \mathbb{N}$ is the number of dimensions of vector $\mathbf{a}_{r_1 \dots r_k}^{(k)}$ and R_k is the rank of matrix $\mathbf{W}_{r_1 \dots r_{k-1}}^{(k)}$, $\sigma_{r_k}(\mathbf{W}_{r_1 \dots r_{k-1}}^{(k)})$ denotes the r_k^{th} singular value of tensor $\mathcal{W}_{r_1 \dots r_{k-1}}^{(k)}$. By reshaping

vectors $\mathbf{v}_{r_1 \dots r_k}^{(k)}$, $\mathbf{b}_{r_1 \dots r_k}^{(k)}$ to matrices according to

$$\mathbf{V}_{r_1 \dots r_k}^{(k)} = \text{MAT} \left(\text{UNFOLD} \left(\text{VEC}^{-1} \left(\mathbf{v}_{r_1 \dots r_k}^{(k)}, \prod_{s=k+1}^S \mathbf{d}^{(s)} \right), \prod_{s=k+2}^S \mathbf{d}^{(s)} \right) \right), \quad (26)$$

$$\mathbf{B}_{r_1 \dots r_k}^{(k)} = \text{MAT} \left(\text{UNFOLD} \left(\text{VEC}^{-1} \left(\mathbf{b}_{r_1 \dots r_k}^{(k)}, \prod_{s=k+1}^S \mathbf{d}^{(s)} \right), \prod_{s=k+2}^S \mathbf{d}^{(s)} \right) \right), \quad (27)$$

where $\mathbf{d}^{(s)} = (a_1^{(s)}, \dots, a_N^{(s)})$ describes the dimensions of the s^{th} factor, we can re-write equation 25 as

$$\sum_{r_k=\hat{R}_k+1}^{R_k} \sigma_{r_k}^2(\mathbf{W}_{r_1 \dots r_{k-1}}^{(k)}) + \sum_{r_k=1}^{\hat{R}_k} d^{(k)} \left\| \mathbf{v}_{r_1 \dots r_k}^{(k)} - \mathbf{b}_{r_1 \dots r_k}^{(k)} \right\|_{\text{F}}^2 \quad (28)$$

$$= \sum_{r_k=\hat{R}_k+1}^{R_k} \sigma_{r_k}^2(\mathbf{W}_{r_1 \dots r_{k-1}}^{(k)}) + \sum_{r_k=1}^{\hat{R}_k} d^{(k)} \left\| \mathbf{v}_{r_1 \dots r_k}^{(k)} - \mathbf{B}_{r_1 \dots r_k}^{(k)} \right\|_{\text{F}}^2 \quad (29)$$

$$= \sum_{r_k=\hat{R}_k+1}^{R_k} \sigma_{r_k}^2(\mathbf{W}_{r_1 \dots r_{k-1}}^{(k)}) + \sum_{r_k=1}^{\hat{R}_k} d^{(k)} \left\| \mathbf{v}_{r_1 \dots r_k}^{(k)} - \sum_{r_{k+1}=1}^{\hat{R}_{k+1}} \mathbf{a}_{r_1 \dots r_{k+1}}^{(k+1)} \mathbf{b}_{r_1 \dots r_{k+1}}^{(k+1)\top} \right\|_{\text{F}}^2. \quad (30)$$

The last line reveals the recursive nature of the formula (compare with equation 20). Unrolling the recursive formula for $k = 1, \dots, S-1$, by setting $\mathbf{W}_{r_1 \dots r_k}^{(k+1)} \leftarrow \mathbf{V}_{r_1 \dots r_k}^{(k)}$, leads to the following formula for the reconstruction error:

$$\begin{aligned} \varepsilon_{\text{SeKron}}(\mathbf{W}, \mathbf{r}, \mathbf{D}) = & \sum_{r_1=\hat{R}_1+1}^{R_1} \sigma_{r_1}^2(\mathbf{W}^{(1)}) + d^{(1)} \sum_{r_1=1}^{\hat{R}_1} \sum_{r_2=\hat{R}_2+1}^{R_2} \sigma_{r_2}^2(\mathbf{W}_{r_1}^{(2)}) + \dots \\ & + d^{(1)} d^{(2)} \dots d^{(S-2)} \sum_{r_1, r_2, \dots, r_{S-2}=1}^{\hat{R}_1, \dots, \hat{R}_{S-2}} \sum_{r_{S-1}=\hat{R}_{S-1}+1}^{R_{S-1}} \sigma_{r_{S-1}}^2(\mathbf{W}_{r_1 \dots r_{S-2}}^{(S-1)}) \end{aligned} \quad (31)$$

where $\mathbf{r} = (\hat{R}_1, \dots, \hat{R}_{S-1})$ contains the rank values, $\mathbf{D}_s = \mathbf{d}^{(s)}$ contains the Kronecker factor shapes and is referred to as the **Dr-SeKron** approximation error (note that the dependency of intermediate matrices $\mathbf{W}_{r_1 \dots r_{k-1}}^{(k)}$ on Kronecker factor shapes \mathbf{D} is implied). Selecting $\hat{R}_i = R_i \forall i$ in equation 31 results in zero reconstruction error. \square

Theorem 2. *The factorization structure imposed by CP, Tucker, TT and TR when decomposing a given tensor $\mathcal{W} \in \mathbb{R}^{w_1 \times \dots \times w_N}$ can be achieved using SeKron.*

Proof. The SeKron decomposition of tensor \mathcal{W} is given by

$$\mathcal{W}_{i_1 \dots i_N} = \sum_{r_1, \dots, r_S=1}^{R_1, \dots, R_S} \mathcal{A}_{r_1 j_1^{(1)} \dots j_N^{(1)}}^{(1)} \dots \mathcal{A}_{r_1 \dots r_{S-1} j_1^{(S)} \dots j_N^{(S)}}^{(S)} \quad (32)$$

where $\mathcal{A}^{(k)} \in \mathbb{R}^{R_1 \times \dots \times R_k \times a_1^{(k)} \times \dots \times a_N^{(k)}}$ and

$$j_n^{(k)} = \begin{cases} i_n - \sum_{t=1}^{k-2} j_n^{(t)} \prod_{l=t+1}^S a_n^{(l)} \bmod a_n^{(S)} & k = S, \\ \left\lfloor \frac{i_n - \sum_{t=1}^{k-1} j_n^{(t)} \prod_{l=t+1}^S a_n^{(l)}}{\prod_{l=k+1}^S a_n^{(l)}} \right\rfloor & \text{otherwise,} \end{cases} \quad (16 \text{ revisited})$$

The CP decomposition of tensor \mathcal{W} in scalar form is

$$\mathcal{W}_{i_1 \dots i_N} = \sum_{r=1}^{R^{(\text{CP})}} \mathcal{A}_{r i_1}^{(\text{CP}_1)} \dots \mathcal{A}_{r i_N}^{(\text{CP}_N)} \quad (33)$$

where $\mathcal{A}^{(\text{CP}_k)} \in \mathbb{R}^{R^{(\text{CP})} \times w_k}$. Configuring the SeKron decomposition in equation 32 such that $S = N$; $R_1 = R^{(\text{CP})}$; $R_2, \dots, R_N = 1$ and $a_n^{(n)} = w_n$ for $n = 1, \dots, N$ leads to the equivalent form

$$\mathcal{W}_{i_1 \dots i_N} = \sum_{r_1=1}^{R^{(\text{CP})}} \mathcal{A}_{r_1 i_1 1 \dots 1}^{(1)} \dots \mathcal{A}_{r_1 1 \dots 1 i_N}^{(N)}. \quad (34)$$

The Tucker decomposition of tensor \mathcal{W} is given by

$$\mathcal{W}_{i_1 \dots i_N} = \sum_{r_1=1, \dots, r_N}^{R_1^{(\text{T})}, \dots, R_N^{(\text{T})}} \mathcal{G}_{r_1 \dots r_N} \mathcal{A}_{i_1 r_1}^{(\text{T}_1)} \dots \mathcal{A}_{i_N r_N}^{(\text{T}_N)} \quad (35)$$

where $\mathcal{G} \in \mathbb{R}^{R_1^{(\text{T})} \times \dots \times R_N^{(\text{T})}}$ and $\mathcal{A}^{(\text{T}_k)} \in \mathbb{R}^{w_k \times R_k^{(\text{T})}}$. The SeKron decomposition of tensor \mathcal{W} , with $S = N + 1$, $R_n = R_n^{(\text{T})}$ and $a_n^{(n)} = w_n$ for $n = 1, \dots, N$ yields

$$\mathcal{W}_{i_1 \dots i_N} = \sum_{r_1, \dots, r_N=1}^{R_1^{(\text{T})}, \dots, R_N^{(\text{T})}} \mathcal{A}_{r_1 i_1 1 \dots 1}^{(1)} \dots \mathcal{A}_{r_1 \dots r_N 1 \dots 1 i_N}^{(N)} \mathcal{A}_{r_1 \dots r_N 1 \dots 1}^{(N+1)}, \quad (36)$$

which is equivalent to equation 35 in the special case where there are nullity constraints on some elements in the Kronecker factors, such that for $k = 2, \dots, N$

$$\mathcal{A}_{r_1 \dots r_k 1 \dots 1 i_k 1 \dots 1}^{(k)} = 0 \quad \text{when} \quad r_j \in \{x \in \mathbb{N} \mid x \leq R_j^{(\text{T})}, x \neq R_j^{(\text{T}^*)}\} \quad j = 1, \dots, k-1 \quad (37)$$

for any choice of $R_j^{(\text{T}^*)} \in \{x \in \mathbb{N} \mid x \leq R_j^{(\text{T})}\}$. The Tensor Ring (TR) decomposition of \mathcal{W} is given by

$$\mathcal{W}_{i_1 \dots i_N} = \sum_{r_1=1, \dots, r_N}^{R_1^{(\text{TR})}, \dots, R_N^{(\text{TR})}} \mathcal{A}_{i_1 r_1 r_2}^{(\text{TR}_1)} \dots \mathcal{A}_{i_N r_N r_{N+1}}^{(\text{TR}_N)} \quad (38)$$

where $\mathcal{A}^{(\text{TR}_k)} \in \mathbb{R}^{w_k \times R_k^{(\text{TR})} \times R_{k+1}^{(\text{TR})}}$, and $R_1^{(\text{TR})} = R_{N+1}^{(\text{TR})}$. As the Tensor Train decomposition can be viewed as a special case of the Tensor Ring decomposition (with $R_1^{(\text{TR})} = R_{N+1}^{(\text{TR})} = 1$), it suffices to show that SeKron generalizes Tensor Ring. The SeKron decomposition of tensor \mathcal{W} , with $S = N + 1$; $R_k = R_k^{(\text{TR})}$ for $k = 1, \dots, N - 1$ and $a_n^{(n+1)} = w_n$ for $n = 1, \dots, N$ leads to

$$\mathcal{W}_{i_1 \dots i_N} = \sum_{r_1, \dots, r_N=1}^{R_1^{(\text{TR})}, \dots, R_N^{(\text{TR})}} \mathcal{A}_{r_1 1 \dots 1}^{(1)} \mathcal{A}_{r_1 r_2 i_1 1 \dots 1}^{(2)} \dots \mathcal{A}_{r_1 \dots r_{N+1} 1 \dots 1 i_N}^{(N+1)}, \quad (39)$$

which is equivalent to equation 38 in the special case where some elements in the Kronecker factors are constrained, such that all elements in tensor $\mathcal{A}^{(1)}$ are constrained to one and

$$\mathcal{A}_{r_1 \dots r_k 1 \dots 1 i_k 1 \dots 1}^{(k)} = 0 \quad \forall r_j \in \{x \in \mathbb{N} \mid x \leq R_j^{(\text{TR})}, x \neq R_j^{(\text{TR}^*)}\} \quad (40)$$

for

$$j = \begin{cases} 1, \dots, k-2 & k = 2, \dots, N \\ 2, \dots, k-1 & k = N+1 \end{cases} \quad (41)$$

for any choice of $R_j^{(\text{TR}^*)} \in \{x \in \mathbb{N} \mid x \leq R_j^{(\text{TR})}\}$. \square

Theorem 3 (Linear Mappings with Sequences of Kronecker Products). *Any linear mapping using a given tensor \mathcal{W} can be written directly in terms of its Kronecker factors $\mathcal{A}^{(k)} \in \mathbb{R}^{R_1 \times \dots \times R_N \times a_1^{(k)} \times \dots \times a_N^{(k)}}$. That is:*

$$\mathcal{W}_{i_1 \dots i_N} \mathcal{X}_{i_1+z_1, \dots, i_N+z_N} = \sum_{r_1, \dots, r_N}^{R_1, \dots, R_N} \mathcal{A}_{r_1 j_1^{(1)} \dots j_N^{(1)}}^{(1)} \dots \mathcal{A}_{r_1 \dots r_{S-1} j_1^{(S)} \dots j_N^{(S)}}^{(S)} \mathcal{X}_{f(j_1)+z_1, \dots, f(j_N)+z_N}$$

where $j_n^{(k)} \in \mathbb{N}$ is a function of input indices (see Appendix A) and $f(\mathbf{j}_n) = \sum_{k=1}^S j_n^{(k)} \prod_{l=k+1}^S a_n^{(l)}$

Proof. First we bring out the summations in the SeKron representaion of \mathcal{W}

$$\mathcal{W} = \sum_{r_1}^{R_1} \mathcal{A}_{r_1}^{(1)} \otimes \sum_{r_2}^{R_2} \mathcal{A}_{r_1 r_2}^{(2)} \otimes \cdots \otimes \sum_{r_{S-1}}^{R_{S-1}} \mathcal{A}_{r_1 \cdots r_{S-1}}^{(S-1)} \otimes \mathcal{A}_{r_1 \cdots r_{S-1}}^{(S)}, \quad (4 \text{ revisited})$$

such that

$$\mathcal{W} = \sum_{r_1, \dots, r_{S-1}}^{R_1, \dots, R_{S-1}} \mathcal{A}_{r_1}^{(1)} \otimes \cdots \otimes \mathcal{A}_{r_1 r_2 \cdots r_{S-1}}^{(S)}. \quad (42)$$

Then, using the scalar form definition of sequences of kronecker products in equation 16

$$j_n^{(k)} = \begin{cases} i_n - \sum_{t=1}^{k-2} j_n^{(t)} \prod_{l=t+1}^S a_n^{(l)} \bmod a_n^{(S)} & k = S, \\ \left\lfloor \frac{i_n - \sum_{t=1}^{k-1} j_n^{(t)} \prod_{l=t+1}^S a_n^{(l)}}{\prod_{l=k+1}^S a_n^{(l)}} \right\rfloor & \text{otherwise,} \end{cases} \quad (16 \text{ revisited})$$

allows us to re-write equation 42 in scalar form as

$$\mathcal{W}_{i_1 \cdots i_N} = \sum_{r_1 \cdots r_{S-1}}^{R_1} \mathcal{A}_{r_1 j_1^{(1)} \cdots j_N^{(1)}}^{(1)} \cdots \mathcal{A}_{r_1 \cdots r_{S-1} j_1^{(S)} \cdots j_N^{(S)}}^{(S)} \quad (43)$$

As the $j_n^{(k)}$ terms decompose i_n into an integer weighted sum, we can recover i_n using

$$i_n = f(\mathbf{j}_n) \triangleq \sum_{k=1}^S j_n^{(k)} \prod_{l=k+1}^S a_n^{(l)}, \quad (44)$$

where $\mathbf{j}_n = (j_n^{(1)}, \dots, j_n^{(S)})$. Thus, we can write

$$\mathcal{X}_{i_1+z_1, \dots, i_N+z_N} = \mathcal{X}_{f(\mathbf{j}_1)+z_1, \dots, f(\mathbf{j}_N)+z_N}. \quad (45)$$

Finally, combining equations equation 43 and equation 45 leads to

$$\mathcal{W}_{i_1 \cdots i_N} \mathcal{X}_{i_1+z_1, \dots, i_N+z_N} = \sum_{r_1, \dots, r_N}^{R_1, \dots, R_k} \mathcal{A}_{r_1 j_1^{(1)} \cdots j_N^{(1)}}^{(1)} \cdots \mathcal{A}_{r_1 \cdots r_{S-1} j_1^{(S)} \cdots j_N^{(S)}}^{(S)} \mathcal{X}_{f(\mathbf{j}_1)+z_1, \dots, f(\mathbf{j}_N)+z_N}$$

□

Theorem 4. (Universal approximation via shallow SeKron networks) Any shallow SeKron factorized neural network $\hat{f}^{(s)}$ with an L -Lipschitz activation function a , is dense in the class of continuous functions $C(X)$ for any compact subset X of \mathbb{R}^d

Proof. Let \hat{f} denote a shallow neural network, and $f \in C(X)$. Then,

$$\|f - \hat{f}^{(s)}\|_2^2 \triangleq \int_X \left(f(x) - \hat{f}^{(s)}(x) \right)^2 d\mu \quad (46)$$

$$= \int_X \left(f(x) - \hat{f}(x) \right)^2 d\mu \quad (47)$$

$$+ \int_X \left(\hat{f}(x) - \hat{f}^{(s)}(x) \right)^2 d\mu \quad (48)$$

$$+ 2 \int_X \left(f(x) - \hat{f}(x) \right) \left(\hat{f}(x) - \hat{f}^{(s)}(x) \right) d\mu \quad (49)$$

According to Hornik (1991), equation 47 is dense in $C(X)$; therefore, it suffices to show that equation 48 is bounded as well.

$$\int_X \left(\hat{f}(x) - \hat{f}^{(s)}(x) \right)^2 d\mu = \int_X \left(\mathbf{w}^\top \mathbf{a}(\mathbf{W}\mathbf{x}) - \mathbf{w}^\top \mathbf{a}(\mathbf{W}^{(s)}\mathbf{x}) \right)^2 d\mu \quad (50)$$

$$\leq L \|\mathbf{w}\|_2^2 \|X\|_2^2 \varepsilon_{\text{SeKron}}(\mathbf{W}, \mathbf{r}, \mathbf{D}) \quad (51)$$

where ε denotes the **Dr**-SeKron approximation error as in equation 31, with matrix **D** and vector **r** describing the shapes of the Kronecker factors the ranks used in the SeKron decomposition of **W**, respectively. \square

D IMPLEMENTATION DETAILS

In all of our experiments we use 4 NVIDIA Tesla V100 SXM2 32 GB GPUs during training and evaluate run time on a single core of Intel(R) Xeon(R) Gold 6148 CPU @ 2.40GHz.

D.1 IMAGENET EXPERIMENTS

We train all models using stochastic gradient descent for 100 epochs using a batch size of 256. The learning rate is initially set to 0.1 and reduced by a factor of $10\times$ at epochs number 30, 60 and 90. We also use a 0.0001 weight decay.

D.2 CIFAR-10 EXPERIMENTS

We train all models using using stochastic gradient descent for 200 epochs using a batch size of 128. The learning rate is initially set to 0.1 and is reduced by a factor of $5\times$ at epochs number 60, 120 and 160. We use nestrov momentum set to 0.9 and weight decay set to 0.0005.

D.3 DIV2K

We train all models using using the ADAM optimizer for 300 epochs using a batch size of 16. The optimizer’s learning rate is set to 0.0001 and β_1, β_2 are set to 0.9, 0.999 respectively.