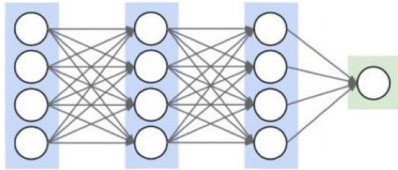# Chapter 7    Tensors for deep learning theory

### Analyzing deep learning architectures via tensorization
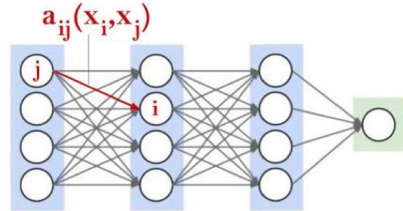
## 7.1  Introduction

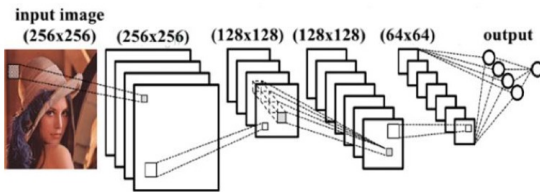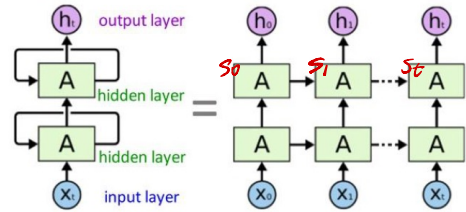There are several prominent deep learning architecture classes.



**FIGURE 7.1**

Prominent deep learning architecture classes. At each layer, fully connected networks connect all outputs to all inputs via learned weights, while in self-attention networks these weights are input-dependent. See a thorough presentation and analysis of self-attention in Section 7.3.1. Convolutional and recurrent networks are discussed in Section 7.4.

a) Fully-Connected Networks: MLP $\quad x \rightarrow (w_0 x + b_0 \rightarrow f) \rightarrow (w_1 x + b_1 \rightarrow f) \rightarrow \cdots \rightarrow y$

b) Convolutional Networks:  convolution kernel

c) Recurrent Networks: $\cdots \rightarrow x_t, x_{t-1}, S_{t-2} \rightarrow x_t, S_{t-1} \rightarrow S_t \rightarrow h_t$, sequence

d) Self-Attention Networks: $\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V$

Three high-level topics:

(1) Expressivity    (2) Optimization    (3) Generalization

## 7.2 Bounding a function's expressivity via tensorization

A scalar function $y$ over $N$ input vectors $\{x^j\}_{j=1}^N$ in dimension $d_x$   $y: (\mathbb{R}^{d_x})^N \to \mathbb{R}$

▷ A measure of capacity for modeling input dependencies

A partition $(A, B)$   $A$ and $B$ are disjoint subsets of $[N] := \{1, \dots, N\}$

Separation rank of $y: (\mathbb{R}^{d_x})^N \to \mathbb{R}$ w.r.t. $(A, B)$:

$$\text{Sep}_{(A,B)}(y) := \min \Big\{ R \in \mathbb{N} \cup \{0\} :$$
$$\exists g_1^A, \dots, g_R^A, g_1^B, \dots, g_R^B : (\mathbb{R}^{d_x})^{N/2} \to \mathbb{R},$$
$$y(x^1, \dots, x^N) =$$
$$\sum_{v=1}^R g_v^A (\{x^i : i \in A\}) \, g_v^B (\{x^i : i \in B\}) \Big\}.$$

If $\text{Sep}_{(A,B)}(y) = 1$, then $y$ is seperable. ( If $y$ is pdf, independent )

$\text{Sep}_{(A,B)}(y) \uparrow$, dependency between $\{x^j\}_{j \in A}$ and $\{x^j\}_{j \in B}$ $\uparrow$

An upper bound of it limits its ability to model input dependencies.

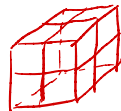A lower bound of it guarantees that this ability is above a threshold.

▽ Upper bound

$$y = \sum_{i=1}^R y_i , \text{ the } \text{sep}_{(A,B)}(y) \leq R$$

▽ Lower bound ( Bounding correlations with tensor matricization ranks )

• A tensor $\mathcal{A}$ of order $N$ and dimension $M_i$ in each mode $i \in [N]$

$\mathcal{A}_{j_1, \dots, j_N}$, where $j_i \in [M_i]$    $N = 3, M_1 = M_2 = M_3 = 2$

- **Grid tensor** ( tensor-based multivariate function discretization )

  A set of points on an exponentially large grid $\implies$ a tensor

  Fixing a set of template vectors $x^{(1)}, \dots, x^{(M)} \in \mathbb{R}^{d_x}$,

  the point on the grid are the set $\left\{ \left( x^{(j_1)}, \dots, x^{(j_N)} \right) \right\}_{j_1, \dots, j_N = 1}^{M}$ $\quad \underbrace{y(x^1, \dots, x^N)}_{}$

  given function $\;\; \mathcal{A}(y) \in (\mathbb{R}^M)^N$
  grid tensor
  induced by $y$

  $$ \mathcal{A}(y)_{j_1 \dots j_N} \equiv y\left( x^1 = x^{(j_1)}, \dots, x^N = x^{(j_N)} \right) $$

  A tensor $\mathcal{A}$ of order $N$ and dimension $M$ in each mode $i \in [N]$

- **Matricization of $\mathcal{A}$** w.r.t. the balanced partition $(A, B)$   $[\![ \mathcal{A} ]\!]_{A,B} \in \mathbb{R}^{M^{N/2} \times M^{N/2}}$

  The arrangement of the tensor elements as a matrix:

  rows $\rightarrow A$, columns $\rightarrow B$

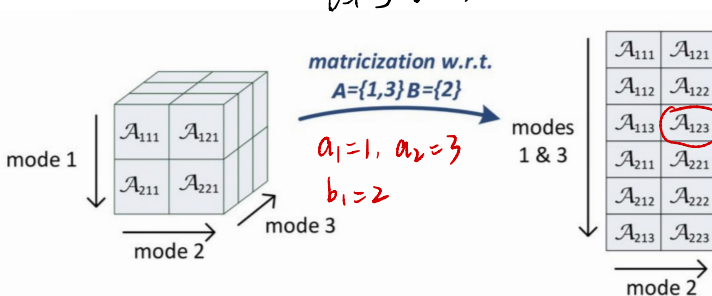  $A = \{ a_1, \dots, a_{N/2} \}, \quad a_1 < \dots < a_{N/2}$

  $B = \{ b_1, \dots, b_{N/2} \}, \quad b_1 < \dots < b_{N/2}$

  The entry $\mathcal{A}_{j_1 \dots j_N}$ is stored in the following matrix entry:

  $$ \left( [\![ \mathcal{A} ]\!]_{A,B} \right)_{row, col} = \mathcal{A}_{j_1 \dots j_N}, $$

  $$ row = 1 + \sum_{t=1}^{N/2} (j_{a_t} - 1) \, M^{N/2 - t} $$

  $$ col = 1 + \sum_{t=1}^{N/2} (j_{b_t} - 1) \, M^{N/2 - t} $$



$j_1 = 1, j_2 = 2, j_3 = 3$

$row = 1 + \sum_{t=1}^{2} (j_{a_t} - 1) \cdot 2^{2-t}$

$= 1 + (3 - 1) \cdot 2^0 = 3$

$col = 1 + \sum_{t=1}^{1} (j_{b_t} - 1) \cdot 2^{1-t}$

$= 1 + (2 - 1) \cdot 2^0 = 2$

matricization w.r.t.
A={1,3} B={2}

$a_1 = 1, \; a_2 = 3$
$b_1 = 2$

modes 1 & 3

**FIGURE 7.3**

An illustrative example of a matricization of an order-3 tensor.

The following claim establishes a fundamental relation between a function's separation rank and the rank of the matrix obtained by the corresponding grid tensor matricization, which lower bounds it.

**Claim 7.1.** Let $y$ be a scalar function over $N$ inputs $\{x^j \in \mathbb{R}^{d_x}\}_{j=1}^N$ and let $(A, B)$ be any partition of $N$. For any integer $M$ and any set of template vectors $x^{(1)}, \ldots, x^{(M)} \in \mathbb{R}^{d_x}$, we have

$$sep_{(A,B)}(y) \geq rank\left( \llbracket A(y) \rrbracket_{A,B} \right),$$

where $A(y)$ is the grid tensor of $y$ with respect to the above template vectors.

**Proof.**   ① $sep_{(A,B)}(y) = \infty$  ✓

② $sep_{(A,B)}(y) = R \in \mathbb{N}$ and $\{g_\nu^A, g_\nu^B\}_{\nu=1}^R$ the functions of decomposition

$$y(x^1, \ldots, x^N) = \sum_{\nu=1}^R g_\nu^A(x^i : i \in A) \cdot g_\nu^B(x^i : i \in B) \quad \text{holds}$$

By def. of the grid tensor,

$$A(y)_{j_1 \cdots j_N} \equiv y(x^1 = x^{(j_1)}, \ldots, x^N = x^{(j_N)})$$
$$= \sum_{\nu=1}^R g_\nu^A(x^{(j_\nu)}, i \in A) \cdot g_\nu^B(x^{(j_i)} : i \in B)$$
$$\equiv \sum_{\nu=1}^R V_{j_i : i \in [A]}^\nu \, U_{j_i : i \in [B]}^\nu \quad \text{holds}$$

where $V^\nu$ and $U^\nu$ are the tensors holding the values of $g_\nu^A, g_\nu^B$, at the input points defined by the template vectors.

$\llbracket V^\nu \rrbracket_{A,B} \rightarrow$ column vectors (only have A part) $\rightarrow \mathbf{r}_\nu$

$\llbracket U^\nu \rrbracket_{A,B} \rightarrow$ row vectors (only have B part) $\rightarrow \mathbf{u}_\nu^T$

$$\llbracket A_{(y)} \rrbracket_{A,B} = \sum_{\nu=1}^{R} v_\nu u_\nu^T, \qquad \text{holds}$$

$$\text{Rank}(v_\nu), \quad \text{Rank}(u_\nu^T) \leq 1$$

$$\text{Rank}(v_\nu u_\nu^T) \leq \min(\text{Rank}(v_\nu), \text{Rank}(u_\nu^T)) \leq 1$$

$$\text{Rank}\left(\sum_{\nu=1}^{R} v_\nu u_\nu^T\right) \leq \sum_{\nu=1}^{R} \text{Rank}(v_\nu u_\nu^T) = R = \text{sep}_{(A,B)}(y). \qquad \#$$

The entries of analyzed functions' grid tensors vary polynomially with learned weight $\Theta$. Finding a single network weight configuration with the above property guarantees the bound for all configurations but a set of measure zero.

**Claim 7.2.** Let $M, N, K \in \mathbb{N}$, $1 \leq R \leq \min\{M, N\}$, and a polynomial mapping $A : \mathbb{R}^K \rightarrow \mathbb{R}^{M \times N}$ (i.e. for every $i \in [M]$, $j \in [N]$, $A_{ij} : \mathbb{R}^K \rightarrow \mathbb{R}$ is a polynomial function). If there exists a point $\Theta \in \mathbb{R}^K$ s.t. $\text{rank}(A(\Theta)) \geq R$, then the set $\{\Theta \in \mathbb{R}^K \mid \text{rank}(A(\Theta)) < R\}$ has zero measure.

**Proof.** $\text{rank}(A(\Theta)) \geq R \iff \exists$ a nonzero $R \times R$ 余子式 minor of $A(\Theta)$

is polynomial in the entries of $A(\Theta) / \Theta$

$c = \binom{M}{R} \cdot \binom{N}{R}$ be the number of minors in $A$, minors: $\{f_i(\Theta)\}_{i=1}^{c}$

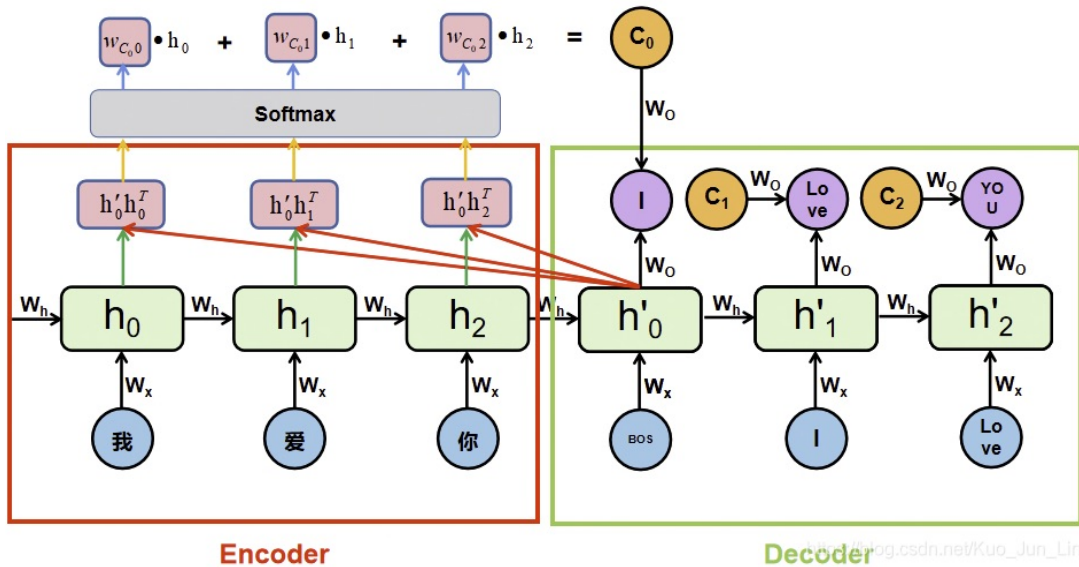Define the polynomial function $f(\Theta) = \sum_{i=1}^{c} f_i(\Theta)^2$

$f(\Theta) = 0 \iff$ For all $i \in [c]$, $f_i(\Theta) = 0 \iff \text{rank}(A(\Theta)) < R$

By relying on Claim 7.2 above, we establish lower bounds on $\text{rank}(\llbracket A(y) \rrbracket_{A,B})$ by choosing a simple-to-analyze assignment of the network's learned weights $\Theta$ and showing that for this value of the weights, there exist template vectors for which the rank of the grid tensor's matricization reaches a certain value. By relying on Claim 7.1 this in turn implies a lower bound on the separation rank of the analyzed architecture.
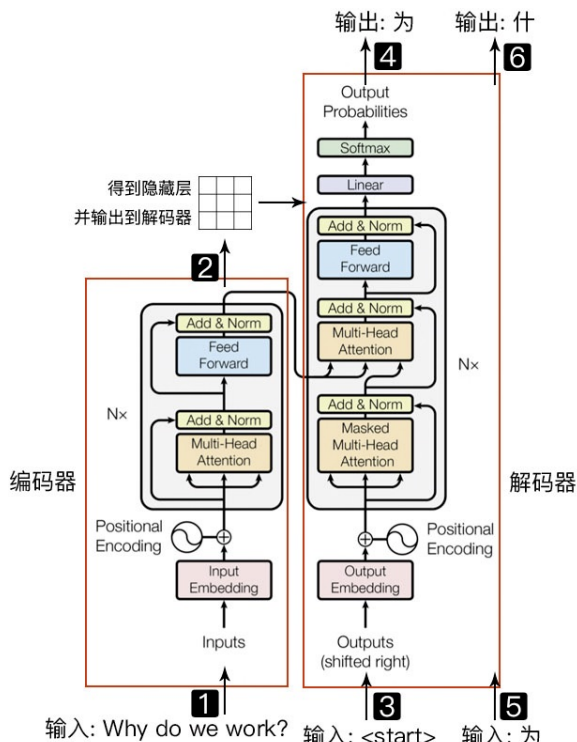
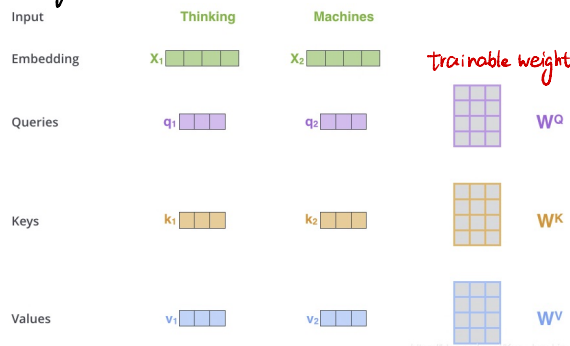# 7.3 A case study: self-attention networks
## 7.3.1 The self-attention mechanism
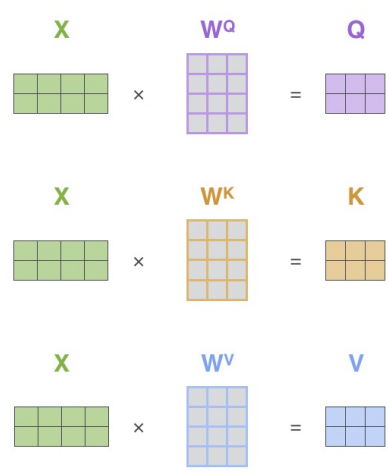▷ Attention in RNN



▷ Transformer

# Single-head Attention

Input    Thinking    Machines

| Input | Thinking | Machines |
|---|---|---|
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ |
| Divide by $\sqrt{d\_k}$ | 14 | 12 |
| Softmax | 0.88 | 0.12 |
| Softmax X Value | $v_1$ | $v_2$ |
| Sum | $z_1$ | $z_2$ |

Embedding   $X_1$   $X_2$

Queries   $q_1$   $q_2$

trainable weight

$W^Q$

Keys   $k_1$   $k_2$

$W^K$

Values   $v_1$   $v_2$

$W^V$

matrix

$X$   $W^Q$   $Q$

$\times$   $=$

$X$   $W^K$   $K$

$\times$   $=$

$X$   $W^V$   $V$

$\times$   $=$
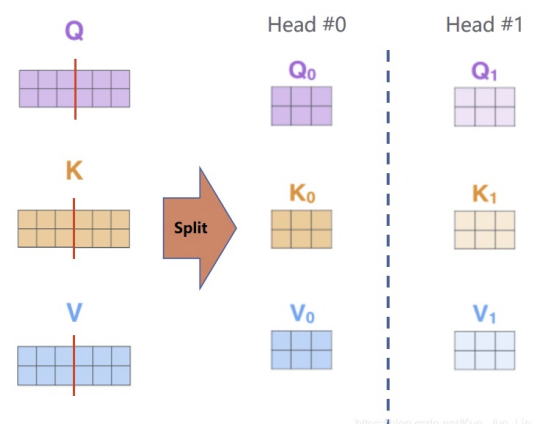
matrix

$$\text{softmax}\left( \frac{Q \times K^T}{\sqrt{d_k}} \right) V$$

$$= Z$$

# Multi-head Attention

$X$   $W^Q$   $Q$

$\times$   $=$

$X$   $W^K$   $K$

$\times$   $=$

$X$   $W^V$   $V$

$\times$   $=$

$Q$    Head #0    Head #1

$K$   Split   $Q_0$   $Q_1$

$V$   $K_0$   $K_1$

$V_0$   $V_1$

[ # of words × C ]    [ C × embedding length ]    [ # of words × embedding length ]

▷ The operation of a self-attention layer

depth $L$ : the number of concatenated layers

$N$ inputs : $\{ x^j \in \mathbb{R}^{d_x} \}_{j=1}^N \xrightarrow[\text{layer}]{\text{first self-attention}} N$ outputs : $\{ y^{l,j} \in \mathbb{R}^{d_x} \}_{j=1}^N$ ← width

$\bar{i}$-th output of layer $l \in [L]$ : $y^{l,i} \in \mathbb{R}^{d_x}$

- Single-headed self-attention

$$\mathbf{y}^{l+1,i}\left(\mathbf{y}^{l,1}, ..., \mathbf{y}^{l,N}\right) = \sum_{j=1}^N a_j^i \left( W^{\mathrm{V},l} \mathbf{y}^{l,j} \right), \qquad (7.5)$$

$$a_j^i = \left\langle W^{\mathrm{Q},l} \mathbf{y}^{l,i}, W^{\mathrm{K},l} \mathbf{y}^{l,j} \right\rangle, \quad \in \mathbb{R}^{d_x \times d_x}$$

- Multi-headed self-attention

$\in \mathbb{R}^{d_x \times d_a}$

$$\mathbf{y}^{l+1,i}\left(\mathbf{y}^{l,1}, ..., \mathbf{y}^{l,N}\right) = \sum_{h=1}^H W^{\mathrm{O},l,h} \sum_{j=1}^N a_{hj}^i \left( W^{\mathrm{V},l,h} \mathbf{y}^{l,j} \right), \qquad (7.6)$$

$$a_{hj}^i = \left\langle W^{\mathrm{Q},l,h} \mathbf{y}^{l,i}, W^{\mathrm{K},l,h} \mathbf{y}^{l,j} \right\rangle, \quad \in \mathbb{R}^{d_a \times d_x}, \; d_a = d_x / H$$

By recursively applying Eq. (7.6) $L$ times we attain a depth-$L$ self-attention network. In Section 7.3.4 we prove that the recursive relation in Eq. (7.6) implies that the function realized by a network with representation dimension $d_x$ and $H$ attention heads per layer at output location $i \in [N]$ can be written as (Corollary 7.1)

$$\mathbf{y}^{i,L,d_x,H,\Theta}(\mathbf{x}^1, \ldots, \mathbf{x}^N) := \sum_{j_1,\ldots,j_C=1}^{N} \mathbf{g}^L(\mathbf{x}^i, \mathbf{x}^{j_1}, \ldots, \mathbf{x}^{j_C}), \qquad (7.7)$$

where $\Theta$ stands for all $4LH$ learned weight matrices: $\forall(l,h) \in [L] \otimes [H]$: $W^{K,l,h}$, $W^{Q,l,h}$, $W^{V,l,h} \in \mathbb{R}^{d_a \times d_x}$, and $W^{O,l,h} \in \mathbb{R}^{d_x \times d_a}$, and the function $\mathbf{g}^L$ is a placeholder which integrates $C + 1 := \frac{3^L - 1}{2} + 1$ different network input vectors.

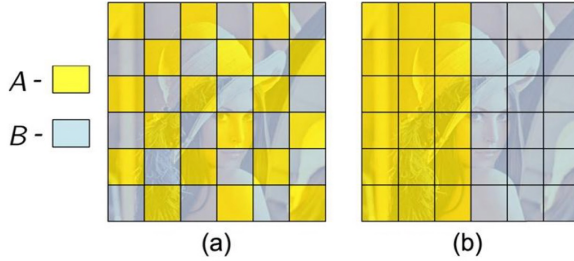▷ Partition invariance of the self-attention separation rank



(a)          (b)

**FIGURE 7.4**

The separation rank of the function realized by self-attention networks is agnostic to the input partition. In contrast, the functions realized by convolutional and recurrent networks have exponentially higher separation ranks w.r.t. the interleaved partition (a) than w.r.t. the left-right partition (b).

**Proposition 7.1.** *For $p \in [d_x]$, let $y_p^{i,L,d_x,H,\Theta}$ be the scalar function computing the $p$-th entry of an output vector at position $i \in [N]$ of the depth-$L$ self-attention network with embedding dimension $d_x$ and $H$ attention heads per layer, defined in Eqs. (7.6) and (7.7). Then, its separation rank w.r.t. balanced partitions, which obey $A \cup B = [N]$, $|A|, |B| = N/2$, is invariant to the identity of the partition, i.e., $\forall A \cup B = [N]$, $\tilde{A} \cup \tilde{B} = [N]$, s.t. $|A|, |B|, |\tilde{A}|, |\tilde{B}| = N/2$:*

$$\mathrm{sep}_{(A,B)}\left(y_p^{i,L,d_x,H,\Theta}\right) = \mathrm{sep}_{(\tilde{A},\tilde{B})}\left(y_p^{i,L,d_x,H,\Theta}\right). \qquad (7.8)$$

Accordingly, for the discussion on self-attention, we will omit the specification of the partition, denoting $\mathrm{sep}(y_p^{i,L,d_x,H,\Theta})$ as the separation rank of $y_p^{i,L,d_x,H,\Theta}$ w.r.t. any balanced partition of the inputs.

**Proof.** $A = (a_1, \ldots, a_{N/2})$, $B = (b_1, \ldots, b_{N/2})$, $\breve{A} = (\breve{a}_1, \ldots, \breve{a}_{N/2})$, $\breve{B} = (\breve{b}_1, \ldots, \breve{b}_{N/2})$

By $\pi \in S_N$ the unique permutation that satisfies

$$\forall m \in \left[\tfrac{N}{2}\right], \quad \pi(a_m) = \breve{a}_m \ \wedge \ \pi(b_m) = \breve{b}_m$$

Assume that $a_1 = \breve{a}_1 = i$, where $i$ is the output location.

Assume that $sep(y^{i,L,d_x,H,\Theta}; A, B) = R$, there exist $g_1^A, \ldots, g_R^A ; g_1^B, \ldots, g_R^B$, s.t.

$$\forall \mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(N)} \in \mathbb{R}^{d_x}:$$

$$y_p^{i,L,d_x,H,\Theta}\big(\mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(N)}\big) = \sum_{v=1}^{R} g_v^A\big(\mathbf{x}^{(a_1)}, \ldots, \mathbf{x}^{(a_{\frac{N}{2}})}\big) g_v^B\big(\mathbf{x}^{(b_1)}, \ldots, \mathbf{x}^{(b_{\frac{N}{2}})}\big).$$

Since both $a_1$ and $\pi(a_1)$ are equal to $i$, the summations over $j_1, \ldots, j_C$ in Eq. (7.7) imply that for any $\mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(N)} \in \mathbb{R}^{d_x}$ we have

$$y_p^{i,L,d_x,H,\Theta}\big(\mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(N)}\big) = y_p^{i,L,d_x,H,\Theta}\big(\mathbf{x}^{(\pi(1))}, \ldots, \mathbf{x}^{(\pi(N))}\big), \quad x^{(i)} = x^{(\pi(i))}$$

and therefore

$$= \sum_{v=1}^{R} g_v\big(\mathbf{x}^{(\pi(a_1))}, \ldots, \mathbf{x}^{(\pi(a_{\frac{N}{2}}))}\big) g_v'\big(\mathbf{x}^{(\pi(b_1))}, \ldots, \mathbf{x}^{(\pi(b_{\frac{N}{2}}))}\big)$$

$$= \sum_{v=1}^{R} g_v\big(\mathbf{x}^{(\breve{a}_1)}, \ldots, \mathbf{x}^{(\breve{a}_{\frac{N}{2}})}\big) g_v'\big(\mathbf{x}^{(\breve{b}_1)}, \ldots, \mathbf{x}^{(\breve{b}_{\frac{N}{2}})}\big).$$

So we proved that

$$sep\big(y_p^{i,L,d_x,H,\Theta}; \tilde{A}, \tilde{B}\big) \le sep\big(y_p^{i,L,d_x,H,\Theta}; A, B\big).$$

Finally, by switching the roles of $\tilde{A}, \tilde{B}$ and $A, B$, we can get the inverse inequality so we conclude that

$$sep\big(y_p^{i,L,d_x,H,\Theta}; \tilde{A}, \tilde{B}\big) = sep\big(y_p^{i,L,d_x,H,\Theta}; A, B\big). \qquad \square$$

It does not integrate the input in a predefined pattern like convolutional networks, but rather computes the input integration in an input-dependent manner $(a_j^i)$

## 7.3.2 Self-attention architecture expressivity questions

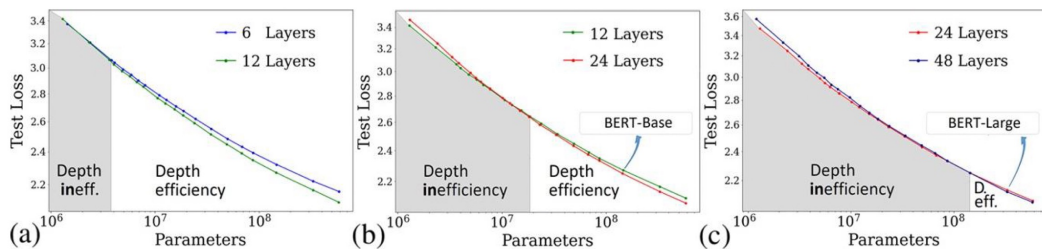- The depth-to-width interplay in self attention



**FIGURE 7.5**

From [47], an experimental validation of the two depth efficiency/**in**efficiency regimes predicted by the tensorial framework of this chapter. [47] further showed that the transition between regimes occurs in exponentially larger network sizes as the networks gets deeper, in agreement with the theory in Section 7.3.3.1.

- The input embedding rank bottleneck in self-attention

When the dimensionality of the inputs is lower than the network width, the network width's contribution to expressivity is capped.

$\forall j$ the input $x^j \in \mathbb{R}^{d_x}$ can be written as

$$\mathbf{x}^j = M\mathbf{w}^j \tag{7.9}$$
$$\text{for } M \in \mathbb{R}^{d_x \times r} \; ; \; \mathbf{w}^j \in R^r$$
$$\text{s.t. } r < d_x,$$

where $r$ is referred to as the input embedding rank. For example, a naive representation of an image as a collection of $N$ RGB vectors in $\mathbb{R}^3$ would imply that $r = 3$, and the identified bottleneck would come into play for $d_x > 3$.

- Mid-architecture rank bottlenecks in self-attention

  eg. T5-11B, with $H \cdot d_a \gg d_x$

  Low representation dimension caps the ability an excessive parameter increase in the self-attention operation. (T5-11B is ~50% redundant.)

## 7.3.3 Results on the operation of self-attention

- The effect of depth $L$

  $L < \log_3(d_x)$, deepen > widen (theorem 7.1)

  $L > \log_3(d_x)$, deepen $\widehat{\backsim}$ widen (theorem 7.2)

**Theorem 7.1.** *For $p \in [d_x]$, let $y_p^{i,L,d_x,H,\Theta}$ be the scalar function computing the p-th entry of an output vector at position $i \in [N]$ of the depth-L self-attention network with embedding dimension $d_x$ and $H$ attention heads per layer, defined in Eqs. (7.6) and (7.7). Let $\text{sep}(y_p^{i,L,d_x,H,\Theta})$ be its separation rank (Section 7.2). If $L, d_x$ obey $L < \log_3(d_x)$, then the following holds almost everywhere in the network's learned parameter space, i.e., for all values of the weight matrices (represented by $\Theta$) but a set of Lebesgue measure zero,*

$$3^{L-2}\big(\log_3(d_x - H) + a\big) \leq \log_3\big(\text{sep}\big(y_p^{i,L,d_x,H,\Theta}\big)\big) \leq \frac{3^L - 1}{2} \log_3(d_x + H) \quad (7.10)$$

*with $a = -L + [2 - \log_3 2]$ (note that $\log_3(d_x - H) + a > 0$ in this regime of $L < \log_3(d_x)$).*

depth $\backsim$ sep : double exponentially , width $\backsim$ sep : polynomially

$H \backsim$ width : linear

**Theorem 7.2.** *For $y_p^{i,L,d_x,H,\Theta}$ as defined in Theorem 7.1, if $L > \log_3(d_x)$, then the following holds almost everywhere in the network's learned parameter space, i.e., for all values of the weight matrices (represented by $\Theta$) but a set of Lebesgue measure zero:*

$$\frac{1}{2}d_x \cdot L + b_1 + b_2 \leq \log_3\big(\text{sep}\big(y_p^{i,L,d_x,H,\Theta}\big)\big) \leq 2d_x \cdot L + c_1 + c_2, \quad (7.11)$$

*with corrections on the order of $L$: $b_1 = -L(\frac{H}{2} + 1)$, $c_1 = L$, and on the order of $d_x \log_3(d_x)$: $b_2 = -d_x(1 + \frac{1}{2}\log_3(\frac{d_x - H}{2}))$, $c_2 = -2d_x \cdot \log_3 d_x/2\sqrt{2e} + \log_3 d_x$.*

depth $\backsim$ sep & width $\backsim$ sep : exponetially

- **The effect of bottlenecks**

The following theorem states that the network's capacity to model dependencies is harmed by a low-rank input embedding.

**Theorem 7.3.** *Let $y_p^{i,L,d_x,H,r}$ be the scalar function computing the pth entry of an output vector at position $i \in [N]$ of the $H$-headed depth-$L$ width-$d_x$ self-attention network defined in Eq. (7.6), where the embedding rank $r$ is defined by Eq. (7.9). Let $sep(y_p^{i,L,d_x,H,r})$ denote its separation rank (Section 7.2.1). Then the following holds (upper bound):*

$$\log\left(sep\left(y_p^{i,L,d_x,H,r}\right)\right) = \tilde{O}\left(L \cdot \min\{r, d_x\}\right). \tag{7.12}$$

*Further assume that $L > \log_3 d_x$, $H < r$. Then for all values of the network weights but a set of Lebesgue measure zero, the following holds (lower bound):*
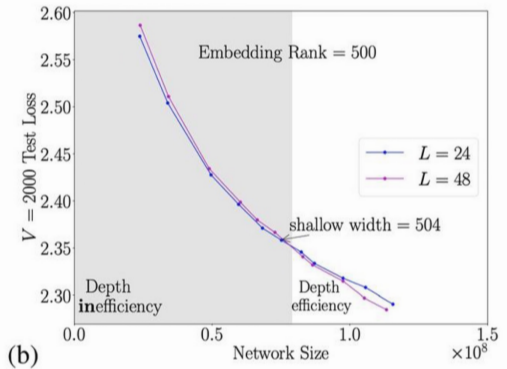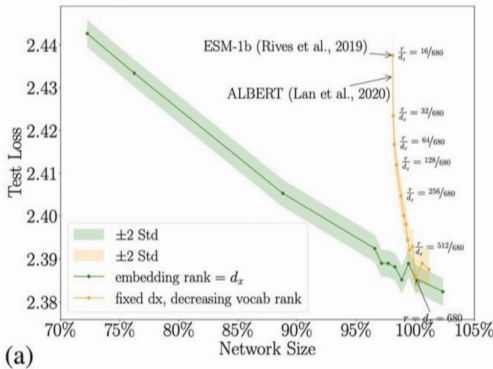
$$\log\left(sep\left(y_p^{i,L,d_x,H,r}\right)\right) = \tilde{\Omega}\left(L \cdot \left(\min\{r, d_x\} - H\right)\right). \tag{7.13}$$

① The input embedding rank bottleneck.

A network with a low-rank embedding $r < d_x$ cannot express the operation of a full-rank $r = d_x$ network.

② Effect on the depth-to-width interplay

$d_x$ (shallow) $> r$, the deeper network is better

(a)

(b)

③ A mid-architecture bottleneck — width caps the internal attention dimension.

Upper bound does not depend on $H$.