# EECE 571F: Advanced Topics in Deep Learning

## Lecture 2: Invariance, Equivariance, and Deep Learning Models for Sets/Sequences

Renjie Liao

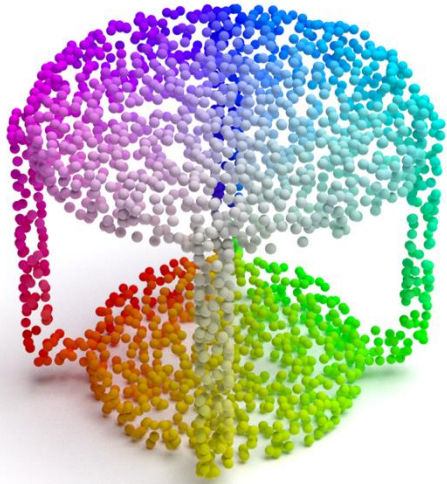University of British Columbia

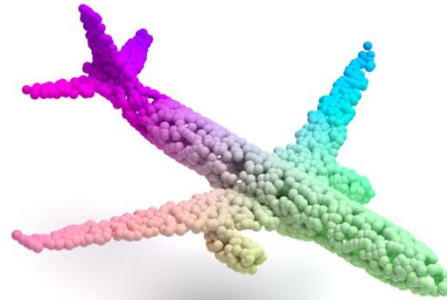Winter, Term 2, 2025

# Outline

- Invariance & Equivariance Principle
  - Translation equivariance in convolutions
  - Permutation equivariance and invariance

- Models for Sets
  - DeepSets: representation theorem of permutation-invariant set functions & architecture
  - DeepSets: permutation-equivariant linear mapping & architecture

- Models for Sequences
  - Transformers
  - Positional encoding vs. Rotary Positional Embeddings (RoPE)
  - Attention & Flash Attention
  - Pre-norm vs. post-norm
  - Vision Transformers (ViT) & Swin Transformers
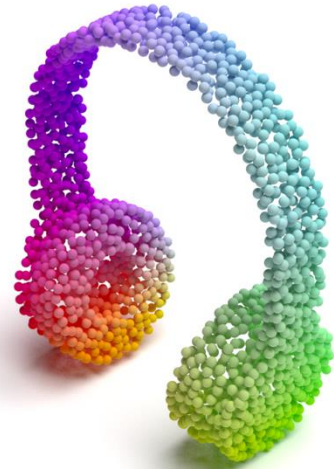
# Motivating Applications for Sets

- Population Statistics
- Point Cloud Classification



Table

Airplane

Earphone

Image Credit: https://github.com/AnTao97/PointCloudDatasets

# Invariance & Equivariance

- Invariance:

    A mathematical object (or a class of mathematical objects) remains unchanged after operations or transformations of a certain type are applied to the objects

$$f(X) = f(g(X))$$

# Invariance & Equivariance

- Invariance:

    A mathematical object (or a class of mathematical objects) remains unchanged after operations or transformations of a certain type are applied to the objects

$$f(X) = f(g(X))$$

# Invariance & Equivariance

- Invariance:

    A mathematical object (or a class of mathematical objects) remains unchanged after operations or transformations of a certain type are applied to the objects

$$f(X) = f(g(X))$$

- Equivariance:

    Applying a transformation and then computing the function produces the same result as computing the function and then applying the transformation

$$g(f(X)) = f(g(X))$$

# Outline

- Invariance & Equivariance Principle
  - **Translation equivariance in convolutions**
  - Permutation equivariance and invariance

- Models for Sets
  - DeepSets: representation theorem of permutation-invariant set functions & architecture
  - DeepSets: permutation-equivariant linear mapping & architecture

- Models for Sequences
  - Transformers
  - Positional encoding vs. Rotary Positional Embeddings (RoPE)
  - Attention & Flash Attention
  - Pre-norm vs. post-norm
  - Vision Transformers (ViT) & Swin Transformers

# Revisit Convolution

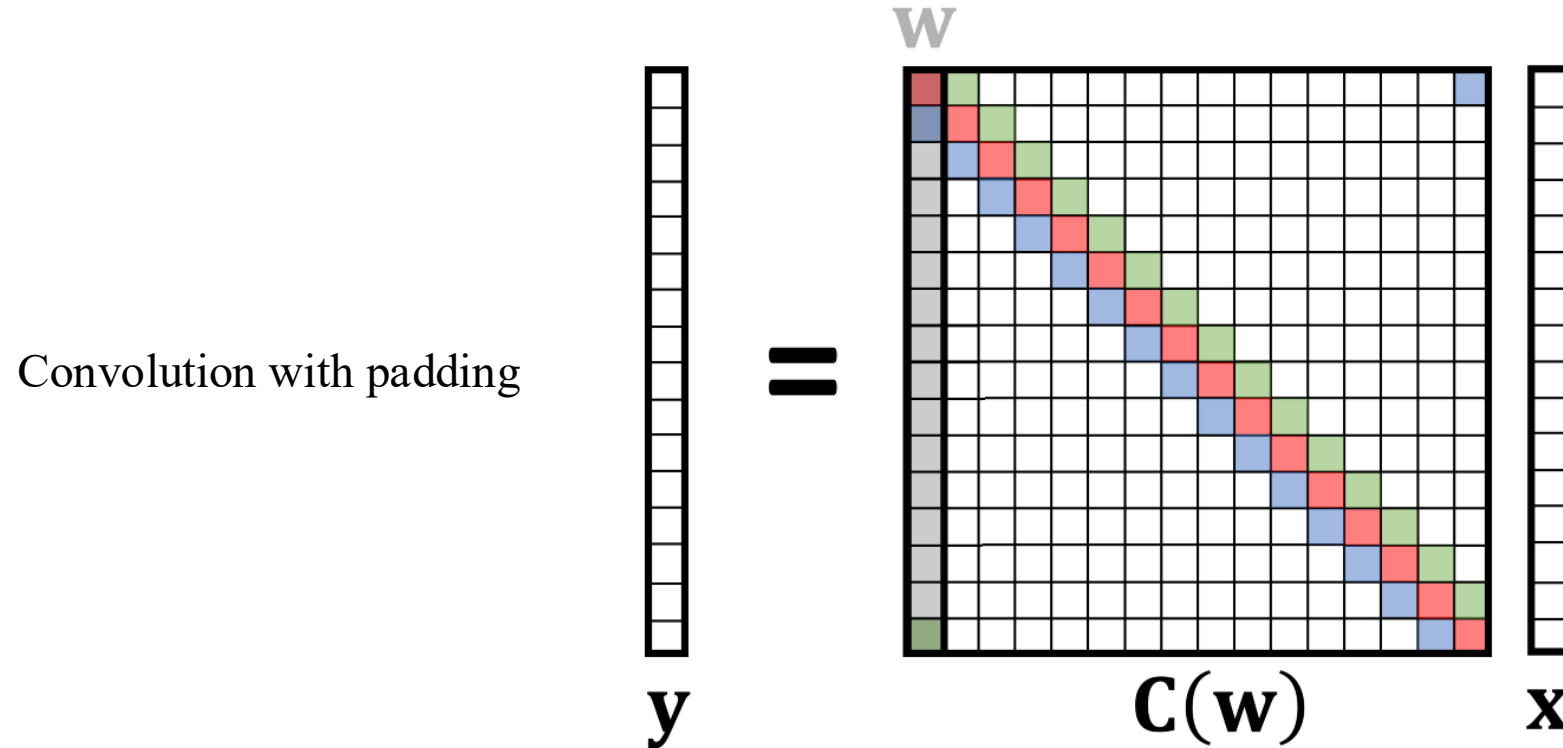Matrix multiplication views of (discrete) convolution:

- Filter => Toeplitz matrix
- Data => Toeplitz matrix

# Revisit Convolution

Matrix multiplication views of (discrete) convolution:

- Filter => Toeplitz matrix
- Data => Toeplitz matrix

Consider a special Toeplitz matrix: circulant matrix (must be square!)

Convolution with padding

Image Credit: https://towardsdatascience.com/deriving-convolution-from-first-principles-4ff124888028

# Translation/Shift Operator

# Translation/Shift Operator

Shift operator is also a circulant matrix!

# Translation/Shift Equivariance

Matrix multiplication is non-commutative. But not for circulant matrices!



$$C(w) \qquad S^T \text{ shift operator} \qquad = \qquad S^T \text{ shift operator} \qquad C(w)$$

Image Credit: https://towardsdatascience.com/deriving-convolution-from-first-principles-4ff124888028

# Translation/Shift Equivariance

Matrix multiplication is non-commutative. But not for circulant matrices!



$$C(w) \quad S^T \quad = \quad S^T \quad C(w)$$

**C(w)**     **Sᵀ** shift operator     **Sᵀ** shift operator     **C(w)**

Convolution is translation equivariant, i.e., Conv(Shift(X)) = Shift(Conv(X))!

Image Credit: https://towardsdatascience.com/deriving-convolution-from-first-principles-4ff124888028

# Translation/Shift Invariance

Global pooling gives you shift-invariance!

Image Credit: https://towardsdatascience.com/covolutional-neural-network-cb0883dd6529

# Translation/Shift Equivariance Invariance

Yann LeCun's LeNet Demo:

# Outline

- Invariance & Equivariance Principle
  - Translation equivariance in convolutions
  - **Permutation equivariance and invariance**

- Models for Sets
  - DeepSets: representation theorem of permutation-invariant set functions & architecture
  - DeepSets: permutation-equivariant linear mapping & architecture

- Models for Sequences
  - Transformers
  - Positional encoding vs. Rotary Positional Embeddings (RoPE)
  - Attention & Flash Attention
  - Pre-norm vs. post-norm
  - Vision Transformers (ViT) & Swin Transformers

# Permutation Invariance



Table

Point Clouds $\qquad$ $X \in \mathbb{R}^{n \times 3}$

Probability of Classes $\qquad$ $Y \in \mathbb{R}^{1 \times K}$

Permutation / Shuffle $\qquad$ $P \in \mathbb{R}^{n \times n}$

# Permutation Invariance



Table

Point Clouds $\qquad$ $X \in \mathbb{R}^{n \times 3}$

Probability of Classes $\qquad$ $Y \in \mathbb{R}^{1 \times K}$

Permutation / Shuffle $\qquad$ $P \in \mathbb{R}^{n \times n}$

$$
\begin{bmatrix} 2 \\ 5 \\ 3 \\ 1 \\ 4 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{bmatrix}
$$

# Geometric Interpretation of Permutation Matrix

Birkhoff Polytope

$$B_n = \{P \in \mathbb{R}^{n \times n} | \forall i \forall j \ P_{ij} \geq 0, \forall i \ \sum_j P_{ij} = 1, \forall j \ \sum_i P_{ij} = 1\}$$

<span style="color:red">Doubly Stochastic Matrix</span>

# Geometric Interpretation of Permutation Matrix

Birkhoff Polytope

$$B_n = \{P \in \mathbb{R}^{n \times n} | \forall i \forall j \ P_{ij} \geq 0, \forall i \ \sum_j P_{ij} = 1, \forall j \ \sum_i P_{ij} = 1\}$$

<span style="color:red">Doubly Stochastic Matrix</span>

Birkhoff–von Neumann Theorem:

1. Birkhoff Polytope is the convex hull of permutation matrices

2. Permutation matrices = Vertices of Birkhoff Polytope $S_n$

# Geometric Interpretation of Permutation Matrix

Birkhoff Polytope

$$B_n = \{P \in \mathbb{R}^{n \times n} | \forall i \forall j \ P_{ij} \geq 0, \forall i \ \sum_j P_{ij} = 1, \forall j \ \sum_i P_{ij} = 1\}$$

<span style="color:red">Doubly Stochastic Matrix</span>

Birkhoff–von Neumann Theorem:

1. Birkhoff Polytope is the convex hull of permutation matrices

2. Permutation matrices = Vertices of Birkhoff Polytope $S_n$

# Permutation Invariance



Table

Point Clouds $\qquad X \in \mathbb{R}^{n \times 3}$

Probability of Classes $\qquad Y \in \mathbb{R}^{1 \times K}$

Permutation / Shuffle $\qquad P \in \mathbb{R}^{n \times n}$

$$Y = f(PX) \qquad \forall P \in S_n$$

Image Credit: https://github.com/AnTao97/PointCloudDatasets

# Permutation Equivariance



Table

| | |
|---|---|
| Point Clouds | $X \in \mathbb{R}^{n \times 3}$ |
| Probability of Classes | $Y \in \mathbb{R}^{1 \times K}$ |
| Permutation / Shuffle | $P \in \mathbb{R}^{n \times n}$ |
| Point Representations | $H \in \mathbb{R}^{n \times d}$ |

# Permutation Equivariance



Table

Point Clouds $\qquad$ $X \in \mathbb{R}^{n \times 3}$

Probability of Classes $\qquad$ $Y \in \mathbb{R}^{1 \times K}$

Permutation / Shuffle $\qquad$ $P \in \mathbb{R}^{n \times n}$

Point Representations $\qquad$ $H \in \mathbb{R}^{n \times d}$

$$H = f(X)$$

# Permutation Equivariance



Table

Point Clouds $\qquad$ $X \in \mathbb{R}^{n \times 3}$

Probability of Classes $\qquad$ $Y \in \mathbb{R}^{1 \times K}$

Permutation / Shuffle $\qquad$ $P \in \mathbb{R}^{n \times n}$

Point Representations $\qquad$ $H \in \mathbb{R}^{n \times d}$

$$H = f(X)$$

$$PH = Pf(X) = f(PX)$$

# Permutation Equivariance



Table

Point Clouds $\qquad X \in \mathbb{R}^{n \times 3}$

Probability of Classes $\qquad Y \in \mathbb{R}^{1 \times K}$

Permutation / Shuffle $\qquad P \in \mathbb{R}^{n \times n}$

Point Representations $\qquad H \in \mathbb{R}^{n \times d}$

$$H = f(X)$$

$$PH = Pf(X) = f(PX)$$

Image Credit: https://github.com/AnTao97/PointCloudDatasets

# More on Invariance & Equivariance

- What about other transformations, e.g., scaling, 2D/3D rotations, Gauge transformation?

Image Credit: http://yann.lecun.com/exdb/lenet/scale.html

# More on Invariance & Equivariance

- What about other transformations, e.g., scaling, 2D/3D rotations, Gauge transformation?



- Generalize to Group Invariance & Equivariance

    Recommend Taco Cohen's PhD Thesis: https://pure.uva.nl/ws/files/60770359/Thesis.pdf

Image Credit: http://yann.lecun.com/exdb/lenet/scale.html

# Outline

- Invariance & Equivariance Principle
  - Translation equivariance in convolutions
  - Permutation equivariance and invariance

- Models for Sets
  - **DeepSets: representation theorem of permutation-invariant set functions & architecture**
  - DeepSets: permutation-equivariant linear mapping & architecture

- Models for Sequences
  - Transformers
  - Positional encoding vs. Rotary Positional Embeddings (RoPE)
  - Attention & Flash Attention
  - Pre-norm vs. post-norm
  - Vision Transformers (ViT) & Swin Transformers

# Deep Learning for Sets

- Point-level Tasks

  Input: a vector per point

  Output: a label/vector per point

  Predictions of individual points are independent, e.g., image classification

# Deep Learning for Sets

- Point-level Tasks

  Input: a vector per point

  Output: a label/vector per point

  Predictions of individual points are independent, e.g., image classification

- Set-level Tasks

  Input: a set of vectors, each corresponds to a point

  Output: a label/vector per set

  Prediction of a set depends on all points, e.g., point cloud classification

# Deep Learning for Sets

Key Challenges:

- Varying-sized input sets

- Permutation equivariant and invariant models

- Expressive models

# Deep Learning for Sets

- Deep Sets [1]

  **Theorem 2** *A function $f(X)$ operating on a set $X$ having elements from a countable universe, is a valid set function,* i.e., ***invariant*** *to the permutation of instances in $X$, iff it can be decomposed in the form $\rho\left(\sum_{x \in X} \phi(x)\right)$, for suitable transformations $\phi$ and $\rho$.*

[1] Zaheer, Manzil, et al. "Deep sets." Advances in neural information processing systems 30 (2017).

# Deep Learning for Sets

- Deep Sets [1]

**Theorem 2** *A function $f(X)$ operating on a set $X$ having elements from a countable universe, is a valid set function,* i.e., ***invariant*** *to the permutation of instances in $X$, iff it can be decomposed in the form $\rho\left(\sum_{x \in X} \phi(x)\right)$, for suitable transformations $\phi$ and $\rho$.*

**Sketch of Proof**

Sufficiency: summation is permutation invariant!

[1] Zaheer, Manzil, et al. "Deep sets." Advances in neural information processing systems 30 (2017).

# Deep Learning for Sets

- Deep Sets [1]

  **Theorem 2** *A function $f(X)$ operating on a set $X$ having elements from a countable universe, is a valid set function,* i.e., ***invariant*** *to the permutation of instances in $X$, iff it can be decomposed in the form $\rho\left(\sum_{x \in X} \phi(x)\right)$, for suitable transformations $\phi$ and $\rho$.*

  **Sketch of Proof**

  Sufficiency: summation is permutation invariant!

  Necessity: any valid set function is essentially a function (absorbed in $\rho$) that takes an injective representation of a set as input. Then we just need to construct an injective set representation in the form of $\sum_{x \in X} \phi(x)$ since any other injective set representations can be obtained via some suitable transformation (absorbed in $\rho$) from $\sum_{x \in X} \phi(x)$.

[1] Zaheer, Manzil, et al. "Deep sets." Advances in neural information processing systems 30 (2017).

# Deep Learning for Sets

- Deep Sets [1]

**Theorem 2** *A function $f(X)$ operating on a set $X$ having elements from a countable universe, is a valid set function, i.e.,* ***invariant*** *to the permutation of instances in $X$, iff it can be decomposed in the form $\rho\left(\sum_{x \in X} \phi(x)\right)$, for suitable transformations $\phi$ and $\rho$.*

**Sketch of Proof**

Sufficiency: summation is permutation invariant!

Necessity: any valid set function is essentially a function (absorbed in $\rho$) that takes an injective representation of a set as input. Then we just need to construct an injective set representation in the form of $\sum_{x \in X} \phi(x)$ since any other injective set representations can be obtained via some suitable transformation (absorbed in $\rho$) from $\sum_{x \in X} \phi(x)$.

1. Construct a mapping $c : \mathfrak{X} \to \mathbb{N}$

[1] Zaheer, Manzil, et al. "Deep sets." Advances in neural information processing systems 30 (2017).

# Deep Learning for Sets

- Deep Sets [1]

  **Theorem 2** *A function $f(X)$ operating on a set $X$ having elements from a countable universe, is a valid set function,* i.e., ***invariant*** *to the permutation of instances in $X$, iff it can be decomposed in the form $\rho\left(\sum_{x \in X} \phi(x)\right)$, for suitable transformations $\phi$ and $\rho$.*

  **Sketch of Proof**

  Sufficiency: summation is permutation invariant!

  Necessity: any valid set function is essentially a function (absorbed in $\rho$) that takes an injective representation of a set as input. Then we just need to construct an injective set representation in the form of $\sum_{x \in X} \phi(x)$ since any other injective set representations can be obtained via some suitable transformation (absorbed in $\rho$) from $\sum_{x \in X} \phi(x)$.

  1. Construct a mapping $c : \boxed{\mathfrak{X}} \rightarrow \mathbb{N}$

  <span style="color:red">Countable Universe</span>

[1] Zaheer, Manzil, et al. "Deep sets." Advances in neural information processing systems 30 (2017).

# Deep Learning for Sets

- Deep Sets [1]

  **Theorem 2** *A function $f(X)$ operating on a set $X$ having elements from a countable universe, is a valid set function, i.e., **invariant** to the permutation of instances in $X$, iff it can be decomposed in the form $\rho\left(\sum_{x \in X} \phi(x)\right)$, for suitable transformations $\phi$ and $\rho$.*

  **Sketch of Proof**

  Sufficiency: summation is permutation invariant!

  Necessity: any valid set function is essentially a function (absorbed in $\rho$) that takes an injective representation of a set as input. Then we just need to construct an injective set representation in the form of $\sum_{x \in X} \phi(x)$ since any other injective set representations can be obtained via some suitable transformation (absorbed in $\rho$) from $\sum_{x \in X} \phi(x)$.

  1. Construct a mapping $\qquad c : \mathfrak{X} \to \mathbb{N}$

  2. Let $\qquad\qquad\qquad \phi(x) = 4^{-c(x)}$

[1] Zaheer, Manzil, et al. "Deep sets." Advances in neural information processing systems 30 (2017).

# Deep Learning for Sets

- Deep Sets [1]

  **Theorem 2** *A function $f(X)$ operating on a set $X$ having elements from a countable universe, is a valid set function, i.e., **invariant** to the permutation of instances in $X$, iff it can be decomposed in the form $\rho \left( \sum_{x \in X} \phi(x) \right)$, for suitable transformations $\phi$ and $\rho$.*

  **Sketch of Proof**

  Sufficiency: summation is permutation invariant!

  Necessity: any valid set function is essentially a function (absorbed in $\rho$) that takes an injective representation of a set as input. Then we just need to construct an injective set representation in the form of $\sum_{x \in X} \phi(x)$ since any other injective set representations can be obtained via some suitable transformation (absorbed in $\rho$) from $\sum_{x \in X} \phi(x)$.

  1. Construct a mapping $\qquad c : \mathfrak{X} \to \mathbb{N}$

  2. Let $\qquad \phi(x) = 4^{-c(x)}$

  3. Injection $\qquad X \in 2^{\mathfrak{X}} \to \sum_{x \in X} \phi(x)$

[1] Zaheer, Manzil, et al. "Deep sets." Advances in neural information processing systems 30 (2017).

# Deep Learning for Sets

- Deep Sets [1]

  **Theorem 2** *A function $f(X)$ operating on a set $X$ having elements from a countable universe, is a valid set function, i.e., **invariant** to the permutation of instances in $X$, iff it can be decomposed in the form $\rho\left(\sum_{x \in X} \phi(x)\right)$, for suitable transformations $\phi$ and $\rho$.*

  **Sketch of Proof**

  Sufficiency: summation is permutation invariant!

  Necessity: any valid set function is essentially a function (absorbed in $\rho$) that takes an injective representation of a set as input. Then we just need to construct an injective set representation in the form of $\sum_{x \in X} \phi(x)$ since any other injective set representations can be obtained via some suitable transformation (absorbed in $\rho$) from $\sum_{x \in X} \phi(x)$.

  1. Construct a mapping $\qquad\qquad c : \mathfrak{X} \to \mathbb{N}$

  2. Let $\qquad\qquad\qquad\qquad \phi(x) = 4^{-c(x)}$

  3. Injection $\qquad\qquad\qquad X \in \boxed{2^{\mathfrak{X}}} \to \sum_{x \in X} \phi(x)$ $\qquad$ <span style="color:red">Power Set</span>

[1] Zaheer, Manzil, et al. "Deep sets." Advances in neural information processing systems 30 (2017).

# Deep Learning for Sets

- Deep Sets [1]

  **Theorem 2** *A function $f(X)$ operating on a set $X$ having elements from a countable universe, is a valid set function, i.e., **invariant** to the permutation of instances in $X$, iff it can be decomposed in the form $\rho\left(\sum_{x \in X} \phi(x)\right)$, for suitable transformations $\phi$ and $\rho$.*

  **Sketch of Proof**

  Sufficiency: summation is permutation invariant!

  Necessity: any valid set function is essentially a function (absorbed in $\rho$) that takes an injective representation of a set as input. Then we just need to construct an injective set representation in the form of $\sum_{x \in X} \phi(x)$ since any other injective set representations can be obtained via some suitable transformation (absorbed in $\rho$) from $\sum_{x \in X} \phi(x)$.

  1. Construct a mapping $\qquad c : \mathfrak{X} \to \mathbb{N}$

  2. Let $\qquad\qquad\qquad \phi(x) = 4^{-c(x)}$

  3. Injection $\qquad\qquad X \in 2^{\mathfrak{X}} \to \sum_{x \in X} \phi(x)$

  *Why base 4?*

[1] Zaheer, Manzil, et al. "Deep sets." Advances in neural information processing systems 30 (2017).

# Deep Learning for Sets

Necessity:

1. Construct a mapping $\qquad c : \mathfrak{X} \to \mathbb{N}$

2. Let $\qquad \phi(x) = 4^{-c(x)}$

3. Injection $\qquad X \in 2^{\mathfrak{X}} \to \sum_{x \in X} \phi(x)$

For better illustrate the problem, let us switch to base 2, i.e., $\phi(x) = 2^{-c(x)}$.

# Deep Learning for Sets

Necessity:

Why base 4?

1. Construct a mapping $\quad c : \mathfrak{X} \to \mathbb{N}$

2. Let $\qquad\qquad\qquad \phi(x) = 4^{-c(x)}$

3. Injection $\qquad\qquad X \in 2^{\mathfrak{X}} \to \sum_{x \in X} \phi(x)$

For better illustrate the problem, let us switch to base 2, i.e., $\phi(x) = 2^{-c(x)}$.

Therefore, the above "injection" essentially maps the binary string to a real number via the binary expansion.

# Deep Learning for Sets

Necessity:

Why base 4?

1. Construct a mapping $\qquad c : \mathfrak{X} \to \mathbb{N}$

2. Let $\qquad\qquad\qquad \phi(x) = 4^{-c(x)}$

3. Injection $\qquad\qquad X \in 2^{\mathfrak{X}} \to \sum_{x \in X} \phi(x)$

For better illustrate the problem, let us switch to base 2, i.e., $\phi(x) = 2^{-c(x)}$.

Therefore, the above "injection" essentially maps the binary string to a real number via the binary expansion.

For example, suppose $\mathfrak{X} = \{1, 2, \dots\}$ and the size is $|\mathfrak{X}|$

Then the size-$|\mathfrak{X}|$ binary string of set $X_1 = \{1\}$ is $b_1 = 10 \dots$ and its binary expansion is $\sum_{x \in X_1} \phi(x) = \sum_{i=1} \frac{b_1[i]}{2^i} = \frac{1}{2} = 0.5$

# Deep Learning for Sets

Necessity:

1. Construct a mapping $\quad c : \mathfrak{X} \to \mathbb{N}$

2. Let $\quad \phi(x) = 4^{-c(x)}$

3. Injection $\quad X \in 2^{\mathfrak{X}} \to \sum_{x \in X} \phi(x)$

**Why base 4?**

For better illustrate the problem, let us switch to base 2, i.e., $\phi(x) = 2^{-c(x)}$.

Therefore, the above "injection" essentially maps the binary string to a real number via the binary expansion.

For example, suppose $\mathfrak{X} = \{1, 2, \dots\}$ and the size is $|\mathfrak{X}|$

Then the size-$|\mathfrak{X}|$ binary string of set $X_1 = \{1\}$ is $b_1 = 10 \dots$ and its binary expansion is $\sum_{x \in X_1} \phi(x) = \sum_{i=1} \frac{b_1[i]}{2^i} = \frac{1}{2} = 0.5$

Then the binary string of set $X_2 = \{2, 3, \dots\}$ is $b_2 = 011 \dots$ and its binary expansion is $\sum_{x \in X_2} \phi(x) = \sum_{i=1}^{\infty} \frac{b_2[i]}{2^i} = \sum_{i=2}^{\infty} \frac{1}{2^i} = 0.5$

# Deep Learning for Sets

Necessity:

**Why base 4?**

1. Construct a mapping $c : \mathfrak{X} \to \mathbb{N}$

2. Let $\phi(x) = 4^{-c(x)}$

3. Injection $X \in 2^{\mathfrak{X}} \to \sum_{x \in X} \phi(x)$

For better illustrate the problem, let us switch to base 2, i.e., $\phi(x) = 2^{-c(x)}$.

Therefore, the above "injection" essentially maps the binary string to a real number via the binary expansion.

For example, suppose $\mathfrak{X} = \{1, 2, \dots\}$ and the size is $|\mathfrak{X}|$

Then the size-$|\mathfrak{X}|$ binary string of set $X_1 = \{1\}$ is $b_1 = 10 \dots$ and its binary expansion is $\sum_{x \in X_1} \phi(x) = \sum_{i=1} \frac{b_1[i]}{2^i} = \frac{1}{2} = 0.5$

Then the binary string of set $X_2 = \{2, 3, \dots\}$ is $b_2 = 011 \dots$ and its binary expansion is $\sum_{x \in X_2} \phi(x) = \sum_{i=1}^{\infty} \frac{b_2[i]}{2^i} = \sum_{i=2}^{\infty} \frac{1}{2^i} = 0.5$

**Dyadic rationals do not have unique binary expansions!**

# Deep Learning for Sets

Suppose we use base B, where B > 1. the value of a tail of a geometric series starting from index n+1 is:

$$\sum_{i=n+1}^{\infty} B^{-i} = \frac{B^{-(n+1)}}{1 - B^{-1}} = \frac{B^{-(n+1)}}{\frac{B-1}{B}} = \frac{1}{B^n(B-1)}$$

# Deep Learning for Sets

Suppose we use base B, where B > 1. the value of a tail of a geometric series starting from index n+1 is:

$$\sum_{i=n+1}^{\infty} B^{-i} = \frac{B^{-(n+1)}}{1 - B^{-1}} = \frac{B^{-(n+1)}}{\frac{B-1}{B}} = \frac{1}{B^n(B-1)}$$

We want to ensure that even if a set X contains every single element from index n+1 onwards, its sum still cannot "reach" the value of the n-th element alone. This requires:

$$\phi(x_n) > \sum_{i=n+1}^{\infty} \phi(x_i)$$

$$B^{-n} > \frac{1}{B^n(B-1)}$$

# Deep Learning for Sets

Suppose we use base B, where B > 1. the value of a tail of a geometric series starting from index n+1 is:

$$\sum_{i=n+1}^{\infty} B^{-i} = \frac{B^{-(n+1)}}{1 - B^{-1}} = \frac{B^{-(n+1)}}{\frac{B-1}{B}} = \frac{1}{B^n(B-1)}$$

We want to ensure that even if a set X contains every single element from index n+1 onwards, its sum still cannot "reach" the value of the n-th element alone. This requires:

$$\phi(x_n) > \sum_{i=n+1}^{\infty} \phi(x_i)$$

$$B^{-n} > \frac{1}{B^n(B-1)}$$

If we simplify this inequality, we get: $\quad 1 > \dfrac{1}{B-1} \implies B-1 > 1 \implies B > 2$

Therefore, any base greater than 2 works!

# Deep Learning for Sets

- Deep Sets [1]

Invariant Architecture

Image Credit: [1] Zaheer, Manzil, et al. "Deep sets." Advances in neural information processing systems 30 (2017).

# Outline

- Invariance & Equivariance Principle
  - Translation equivariance in convolutions
  - Permutation equivariance and invariance

- Models for Sets
  - DeepSets: representation theorem of permutation-invariant set functions & architecture
  - **DeepSets: permutation-equivariant linear mapping & architecture**

- Models for Sequences
  - Transformers
  - Positional encoding vs. Rotary Positional Embeddings (RoPE)
  - Attention & Flash Attention
  - Pre-norm vs. post-norm
  - Vision Transformers (ViT) & Swin Transformers

# Deep Learning for Sets

- Deep Sets [1]       $\mathbf{f}_\Theta(\mathbf{x}) \doteq \boldsymbol{\sigma}(\Theta\mathbf{x}) \quad \Theta \in \mathbb{R}^{M \times M}$

**Lemma 3** *The function $\mathbf{f}_\Theta : \mathbb{R}^M \to \mathbb{R}^M$ defined above is permutation* **equivariant** *iff all the off-diagonal elements of $\Theta$ are tied together and all the diagonal elements are equal as well. That is,*

$$\Theta = \lambda\mathbf{I} + \gamma\left(\mathbf{1}\mathbf{1}^\mathsf{T}\right) \qquad \lambda, \gamma \in \mathbb{R} \quad \mathbf{1} = [1, \dots, 1]^\mathsf{T} \in \mathbb{R}^M \qquad \mathbf{I} \in \mathbb{R}^{M \times M} \text{ is the identity matrix}$$

[1] Zaheer, Manzil, et al. "Deep sets." Advances in neural information processing systems 30 (2017).

# Deep Learning for Sets

- Deep Sets [1]
$$\mathbf{f}_\Theta(\mathbf{x}) \doteq \boldsymbol{\sigma}(\Theta\mathbf{x}) \quad \Theta \in \mathbb{R}^{M \times M}$$

**Lemma 3** *The function* $\mathbf{f}_\Theta : \mathbb{R}^M \to \mathbb{R}^M$ *defined above is permutation **equivariant** iff all the off-diagonal elements of* $\Theta$ *are tied together and all the diagonal elements are equal as well. That is,*

$$\Theta = \lambda\mathbf{I} + \gamma\left(\mathbf{1}\mathbf{1}^\mathsf{T}\right) \qquad \lambda, \gamma \in \mathbb{R} \quad \mathbf{1} = [1, \dots, 1]^\mathsf{T} \in \mathbb{R}^M \qquad \mathbf{I} \in \mathbb{R}^{M \times M} \text{ is the identity matrix}$$

**Sketch of Proof**

Permutation Equivariance $\sigma(\Theta\pi\mathbf{x}) = \pi\sigma(\Theta\mathbf{x})$ (w. element-wise nonlinearity) reduces to $\qquad \pi\Theta\mathbf{x} = \Theta\pi\mathbf{x}$

[1] Zaheer, Manzil, et al. "Deep sets." Advances in neural information processing systems 30 (2017).

# Deep Learning for Sets

- Deep Sets [1]  $\mathbf{f}_\Theta(\mathbf{x}) \doteq \boldsymbol{\sigma}(\Theta\mathbf{x}) \quad \Theta \in \mathbb{R}^{M \times M}$

**Lemma 3** *The function* $\mathbf{f}_\Theta : \mathbb{R}^M \to \mathbb{R}^M$ *defined above is permutation* ***equivariant*** *iff all the off-diagonal elements of* $\Theta$ *are tied together and all the diagonal elements are equal as well. That is,*

$\Theta = \lambda\mathbf{I} + \gamma\left(\mathbf{1}\mathbf{1}^\mathsf{T}\right) \qquad \lambda, \gamma \in \mathbb{R} \quad \mathbf{1} = [1, \dots, 1]^\mathsf{T} \in \mathbb{R}^M \qquad \mathbf{I} \in \mathbb{R}^{M \times M}$ *is the identity matrix*

**Sketch of Proof**

Permutation Equivariance $\sigma(\Theta\pi\mathbf{x}) = \pi\sigma(\Theta\mathbf{x})$ (w. element-wise nonlinearity) reduces to $\qquad \pi\Theta\mathbf{x} = \Theta\pi\mathbf{x}$

Sufficiency: $\Theta$ is commutable with permutation matrix

[1] Zaheer, Manzil, et al. "Deep sets." Advances in neural information processing systems 30 (2017).

# Deep Learning for Sets

- Deep Sets [1]
$$\mathbf{f}_\Theta(\mathbf{x}) \doteq \boldsymbol{\sigma}(\Theta\mathbf{x}) \quad \Theta \in \mathbb{R}^{M \times M}$$

**Lemma 3** *The function* $\mathbf{f}_\Theta : \mathbb{R}^M \to \mathbb{R}^M$ *defined above is permutation* ***equivariant*** *iff all the off-diagonal elements of* $\Theta$ *are tied together and all the diagonal elements are equal as well. That is,*

$$\Theta = \lambda\mathbf{I} + \gamma\left(\mathbf{1}\mathbf{1}^\mathsf{T}\right) \qquad \lambda, \gamma \in \mathbb{R} \quad \mathbf{1} = [1, \dots, 1]^\mathsf{T} \in \mathbb{R}^M \qquad \mathbf{I} \in \mathbb{R}^{M \times M} \text{ is the identity matrix}$$

**Sketch of Proof**

Permutation Equivariance $\sigma(\Theta\pi\mathbf{x}) = \pi\sigma(\Theta\mathbf{x})$ (w. element-wise nonlinearity) reduces to $\qquad \pi\Theta\mathbf{x} = \Theta\pi\mathbf{x}$

Sufficiency: $\Theta$ is commutable with permutation matrix

Necessity: consider a special permutation (i.e., transposition / swap) $\qquad \pi_{i,j}^\mathsf{T} = \pi_{i,j}^{-1} = \pi_{j,i}$

    1. All diagonal elements are identical

$$\pi_{k,l}\Theta = \Theta\pi_{k,l} \implies \pi_{k,l}\Theta\pi_{l,k} = \Theta \implies (\pi_{k,l}\Theta\pi_{l,k})_{l,l} = \Theta_{l,l} \implies \Theta_{k,k} = \Theta_{l,l}$$

[1] Zaheer, Manzil, et al. "Deep sets." Advances in neural information processing systems 30 (2017).

# Deep Learning for Sets

- Deep Sets [1]
$$\mathbf{f}_\Theta(\mathbf{x}) \doteq \boldsymbol{\sigma}(\Theta\mathbf{x}) \quad \Theta \in \mathbb{R}^{M \times M}$$

**Lemma 3** *The function* $\mathbf{f}_\Theta : \mathbb{R}^M \to \mathbb{R}^M$ *defined above is permutation* ***equivariant*** *iff all the off-diagonal elements of* $\Theta$ *are tied together and all the diagonal elements are equal as well. That is,*

$$\Theta = \lambda \mathbf{I} + \gamma \left(\mathbf{1}\mathbf{1}^\mathsf{T}\right) \qquad \lambda, \gamma \in \mathbb{R} \quad \mathbf{1} = [1, \ldots, 1]^\mathsf{T} \in \mathbb{R}^M \qquad \mathbf{I} \in \mathbb{R}^{M \times M} \text{ is the identity matrix}$$

**Sketch of Proof**

Permutation Equivariance $\sigma(\Theta\pi\mathbf{x}) = \pi\sigma(\Theta\mathbf{x})$ (w. element-wise nonlinearity) reduces to $\qquad \pi\Theta\mathbf{x} = \Theta\pi\mathbf{x}$

Sufficiency: $\Theta$ is commutable with permutation matrix

Necessity: consider a special permutation (i.e., transposition / swap) $\qquad \pi_{i,j}^\mathsf{T} = \pi_{i,j}^{-1} = \pi_{j,i}$

1. All diagonal elements are identical

$$\pi_{k,l}\Theta = \Theta\pi_{k,l} \;\Rightarrow\; \pi_{k,l}\Theta\pi_{l,k} = \Theta \;\Rightarrow\; (\pi_{k,l}\Theta\pi_{l,k})_{l,l} = \Theta_{l,l} \;\Rightarrow\; \Theta_{k,k} = \Theta_{l,l}$$

2. All off-diagonal elements are identical

$$\pi_{j',j}\pi_{i,i'}\Theta = \Theta\pi_{j',j}\pi_{i,i'} \;\Rightarrow\; \pi_{j',j}\pi_{i,i'}\Theta(\pi_{j',j}\pi_{i,i'})^{-1} = \Theta \qquad\qquad\Rightarrow$$

$$\pi_{j',j}\pi_{i,i'}\Theta\pi_{i',i}\pi_{j,j'} = \Theta \;\Rightarrow\; (\pi_{j',j}\pi_{i,i'}\Theta\pi_{i',i}\pi_{j,j'})_{i,j} = \Theta_{i,j} \;\Rightarrow\; \Theta_{i',j'} = \Theta_{i,j}$$

[1] Zaheer, Manzil, et al. "Deep sets." Advances in neural information processing systems 30 (2017).

# Deep Learning for Sets

- Deep Sets [1]

  Equivariant Architecture

$$f(\mathbf{x}) = \sigma\left(\mathbf{x}\Lambda - \mathbf{1}\mathbf{1}^\mathsf{T}\mathbf{x}\Gamma\right)$$



Optional conditioning based on meta-information

Image Credit: [1] Zaheer, Manzil, et al. "Deep sets." Advances in neural information processing systems 30 (2017).

# Deep Learning for Sets

- Deep Sets [1]

Recipe for making the model deep:

Stack multiple equivariant layers (+ invariant layer at the end), e.g., PointNet [2]

Image Credit: [1] Zaheer, Manzil, et al. "Deep sets." Advances in neural information processing systems 30 (2017).

# Outline

- Invariance & Equivariance Principle
  - Translation equivariance in convolutions
  - Permutation equivariance and invariance

- Models for Sets
  - DeepSets: representation theorem of permutation-invariant set functions & architecture
  - DeepSets: permutation-equivariant linear mapping & architecture

- Models for Sequences
  - **Transformers**
  - Positional encoding vs. Rotary Positional Embeddings (RoPE)
  - Attention & Flash Attention
  - Pre-norm vs. post-norm
  - Vision Transformers (ViT) & Swin Transformers

# Deep Learning for Sequences

- Language Models

$$P(\boldsymbol{x}^{(t+1)} \mid \boldsymbol{x}^{(t)}, \ldots, \boldsymbol{x}^{(1)})$$

the students opened their _____

books

laptops

exams

minds

Image Credit: http://web.stanford.edu/class/cs224n/

# Deep Learning for Sequences

- Language Models

$$P(\boldsymbol{x}^{(t+1)} | \boldsymbol{x}^{(t)}, \ldots, \boldsymbol{x}^{(1)})$$

the students opened their _____

books

laptops

exams

minds

- Machine Translation

INPUT

Je    suis    étudiant

THE TRANSFORMER

OUTPUT

I    am    a    student

Image Credit: http://web.stanford.edu/class/cs224n/ https://jalammar.github.io/illustrated-transformer/

# Deep Learning for Sequences

Key Challenges:

- Varying-sized input sequences

# Deep Learning for Sequences

Key Challenges:

- Varying-sized input sequences
- Orders "may" be crucial for cognition

Aoccdrnig to a rscheearch at Cmabrigde Uinervtisy, it deosn't mttaer in waht oredr the ltteers in a wrod are, the olny iprmoetnt tihng is taht the frist and lsat ltteer be at the rghit pclae. The rset can be a toatl mses and you can sitll raed it wouthit porbelm. Tihs is bcuseae the huamn mnid deos not raed ervey lteter by istlef, but the wrod as a wlohe.

Image Credit: https://www.mrc-cbu.cam.ac.uk/people/matt.davis/cmabridge

# Deep Learning for Sequences

Key Challenges:

- Varying-sized input sequences
- Orders "may" be crucial for cognition

  **Aoccdrnig to a rscheearch at Cmabrigde Uinervtisy, it deosn't mttaer in waht oredr the ltteers in a wrod are, the olny iprmoetnt tihng is taht the frist and lsat ltteer be at the rghit pclae. The rset can be a toatl mses and you can sitll raed it wouthit porbelm. Tihs is bcuseae the huamn mnid deos not raed ervey lteter by istlef, but the wrod as a wlohe.**

- Complex statistical dependencies (e.g. long-range ones)

As aliens entered our planet

# Deep Learning for Sequences

Key Challenges:

- Varying-sized input sequences
- Orders "may" be crucial for cognition

> **Aoccdrnig to a rscheearch at Cmabrigde Uinervtisy, it deosn't mttaer in waht oredr the ltteers in a wrod are, the olny iprmoetnt tihng is taht the frist and lsat ltteer be at the rghit pclae. The rset can be a toatl mses and you can sitll raed it wouthit porbelm. Tihs is bcuseae the huamn mnid deos not raed ervey lteter by istlef, but the wrod as a wlohe.**

- Complex statistical dependencies (e.g. long-range ones)



As aliens entered our planet

LSTM [1]
GRU [2]
Seq2Seq [3]
Transformer [4]

Image Credit: https://www.mrc-cbu.cam.ac.uk/people/matt.davis/cmabridge/ https://towardsdatascience.com/illustrated-guide-to-transformers-step-by-step-explanation-f74876522bc0
[1] Hochreiter, S. "Long Short-term Memory." Neural Computation MIT-Press (1997). [2] Cho, Kyunghyun. "Learning phrase representations using RNN encoder-decoder for statistical machine translation." arXiv preprint arXiv:1406.1078 (2014). [3] Sutskever, I. "Sequence to Sequence Learning with Neural Networks." arXiv preprint arXiv:1409.3215 (2014). [4] Vaswani, A. "Attention is all you need." Advances in Neural Information Processing Systems (2017).

# Transformers

Image Credit: https://jalammar.github.io/illustrated-transformer/

# Transformers



OUTPUT: I am a student

INPUT: Je suis étudiant

# Transformers

Image Credit: Vaswani, A. "Attention is all you need." Advances in Neural Information Processing Systems (2017). https://jalammar.github.io/illustrated-transformer/

# Transformers

Image Credit: Vaswani, A. "Attention is all you need." Advances in Neural Information Processing Systems (2017). https://jalammar.github.io/illustrated-transformer/

# Transformers

# Transformers

# Transformers

# Transformers

Image Credit: Vaswani, A. "Attention is all you need." Advances in Neural Information Processing Systems (2017). https://jalammar.github.io/illustrated-transformer/

# Input Encoding



Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

$N\times$

Add & Norm

Masked
Multi-Head
Attention

Add & Norm

Feed
Forward

$N\times$

Add & Norm

Multi-Head
Attention

Positional
Encoding

Positional
Encoding

Input
Embedding

Output
Embedding

Inputs

Outputs
(shifted right)

Image Credit: Vaswani, A. "Attention is all you need." Advances in Neural Information Processing Systems (2017). https://jalammar.github.io/illustrated-transformer/

# Input Embedding

# Outline

- Invariance & Equivariance Principle
  - Translation equivariance in convolutions
  - Permutation equivariance and invariance

- Models for Sets
  - DeepSets: representation theorem of permutation-invariant set functions & architecture
  - DeepSets: permutation-equivariant linear mapping & architecture

- Models for Sequences
  - Transformers
  - **Positional encoding vs. Rotary Positional Embeddings (RoPE)**
  - Attention & Flash Attention
  - Pre-norm vs. post-norm
  - Vision Transformers (ViT) & Swin Transformers

# Positional Encoding



$$PE_{(pos, 2i)} = sin(pos/10000^{2i/d_{model}})$$

$$PE_{(pos, 2i+1)} = cos(pos/10000^{2i/d_{model}})$$

# Positional Encoding



$$PE_{(pos,2i)} = sin(pos/10000^{2i/d_{model}})$$

$$PE_{(pos,2i+1)} = cos(pos/10000^{2i/d_{model}})$$

Image Credit: https://lilianweng.github.io/lil-log/2020/04/07/the-transformer-family.html

# Absolute vs. Relative Positional Encoding

Encode relative position information could help better model the dependency among tokens.

How to encode relative positions?
- We can inject the relative position into the bias of attention.
- We can use Rotary Position Embedding (RoPE) [1], which is more effective empirically.

To understand RoPE, let us recap how to rotate a 2D vector:

$$\begin{bmatrix} x_1' \\ x_2' \end{bmatrix} = \underbrace{\begin{bmatrix} \cos m\theta & -\sin m\theta \\ \sin m\theta & \cos m\theta \end{bmatrix}}_{\mathbf{R}_{\theta,m}} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

<span style="color:red">Rotation matrix is orthogonal and preserves the norm!</span>

[1] Su, Jianlin, et al. "RoFormer: Enhanced Transformer with Rotary Position Embedding." arXiv preprint arXiv:2104.09864 (2021).

# Rotary Positional Embedding

RoPE first divide d-dimension vector space in d/2 subspaces and then rotate them based on the position:

$$
\begin{bmatrix} x'_1 \\ \vdots \\ x'_d \end{bmatrix} = \underbrace{\begin{bmatrix} \cos m\theta_1 & -\sin m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ \sin m\theta_1 & \cos m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cos m\theta_2 & -\sin m\theta_2 & \cdots & 0 & 0 \\ 0 & 0 & \sin m\theta_2 & \cos m\theta_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \cos m\theta_{d/2} & -\sin m\theta_{d/2} \\ 0 & 0 & 0 & 0 & \cdots & \sin m\theta_{d/2} & \cos m\theta_{d/2} \end{bmatrix}}_{\mathbf{R}^d_{\Theta, m}} \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}
$$

Here $\quad \Theta = \{\theta_i = 10000^{-2(i-1)/d}, i \in [1, 2, ..., d/2]\}$

In practice, we can apply 2D rotations to pairs $(x_1, x_{1+d/2}), (x_2, x_{2+d/2}), \ldots, (x_{d/2}, x_d)$

# Rotary Positional Embedding

RoPE first divide d-dimension vector space in d/2 subspaces and then rotate them based on the position:

Image Credit: Su, Jianlin, et al. "RoFormer: Enhanced Transformer with Rotary Position Embedding." arXiv preprint arXiv:2104.09864 (2021).

# Rotary Positional Embedding

What do we gain in RoPE?

- Inner product depends on the relative position

Let us look at the case of 2D:

$$\left\langle \begin{bmatrix} x_1' \\ x_2' \end{bmatrix}, \begin{bmatrix} y_1' \\ y_2' \end{bmatrix} \right\rangle = \left\langle \underbrace{\begin{bmatrix} \cos m\theta & -\sin m\theta \\ \sin m\theta & \cos m\theta \end{bmatrix}}_{\mathbf{R}_{\theta,m}} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \underbrace{\begin{bmatrix} \cos n\theta & -\sin n\theta \\ \sin n\theta & \cos n\theta \end{bmatrix}}_{\mathbf{R}_{\theta,n}} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \right\rangle$$

$$= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^\top \begin{bmatrix} \cos m\theta & -\sin m\theta \\ \sin m\theta & \cos m\theta \end{bmatrix}^\top \begin{bmatrix} \cos n\theta & -\sin n\theta \\ \sin n\theta & \cos n\theta \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^\top \begin{bmatrix} \cos m\theta \cos n\theta + \sin m\theta \sin n\theta & -\cos m\theta \sin n\theta + \sin m\theta \cos n\theta \\ -\sin m\theta \cos n\theta + \cos m\theta \sin n\theta & \sin m\theta \sin n\theta + \cos m\theta \cos n\theta \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^\top \begin{bmatrix} \cos(m-n)\theta & \sin(m-n)\theta \\ \sin(n-m)\theta & \cos(m-n)\theta \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^\top \underbrace{\begin{bmatrix} \cos(m-n)\theta & -\sin(m-n)\theta \\ \sin(m-n)\theta & \cos(m-n)\theta \end{bmatrix}}_{\mathbf{R}_{\theta,m-n}}^\top \underbrace{\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}}_{\mathbf{R}_{\theta,0}} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \left\langle \mathbf{R}_{\theta,m-n} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \mathbf{R}_{\theta,0} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \right\rangle$$

# Rotary Positional Embedding

What do we gain in RoPE?

- Inner product depends on the relative position

Let us look at the case of 2D:

<span style="color:red">This holds for d-dimension as we construct a block-diagonal matrix with 2D rotation matrices!</span>

$$\left\langle \begin{bmatrix} x_1' \\ x_2' \end{bmatrix}, \begin{bmatrix} y_1' \\ y_2' \end{bmatrix} \right\rangle = \left\langle \underbrace{\begin{bmatrix} \cos m\theta & -\sin m\theta \\ \sin m\theta & \cos m\theta \end{bmatrix}}_{\mathbf{R}_{\theta,m}} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \underbrace{\begin{bmatrix} \cos n\theta & -\sin n\theta \\ \sin n\theta & \cos n\theta \end{bmatrix}}_{\mathbf{R}_{\theta,n}} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \right\rangle$$

$$= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^\top \begin{bmatrix} \cos m\theta & -\sin m\theta \\ \sin m\theta & \cos m\theta \end{bmatrix}^\top \begin{bmatrix} \cos n\theta & -\sin n\theta \\ \sin n\theta & \cos n\theta \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^\top \begin{bmatrix} \cos m\theta \cos n\theta + \sin m\theta \sin n\theta & -\cos m\theta \sin n\theta + \sin m\theta \cos n\theta \\ -\sin m\theta \cos n\theta + \cos m\theta \sin n\theta & \sin m\theta \sin n\theta + \cos m\theta \cos n\theta \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^\top \begin{bmatrix} \cos(m-n)\theta & \sin(m-n)\theta \\ \sin(n-m)\theta & \cos(m-n)\theta \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^\top \underbrace{\begin{bmatrix} \cos(m-n)\theta & -\sin(m-n)\theta \\ \sin(m-n)\theta & \cos(m-n)\theta \end{bmatrix}}_{\mathbf{R}_{\theta,m-n}}^\top \underbrace{\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}}_{\mathbf{R}_{\theta,0}} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \left\langle \mathbf{R}_{\theta,m-n} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \mathbf{R}_{\theta,0} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \right\rangle$$

# Rotary Positional Embedding

What do we gain in RoPE?

- Long-term decay of inner product w.r.t. relative positions

Image Credit: Su, Jianlin, et al. "RoFormer: Enhanced Transformer with Rotary Position Embedding." arXiv preprint arXiv:2104.09864 (2021).

# Outline

- Invariance & Equivariance Principle
  - Translation equivariance in convolutions
  - Permutation equivariance and invariance

- Models for Sets
  - DeepSets: representation theorem of permutation-invariant set functions & architecture
  - DeepSets: permutation-equivariant linear mapping & architecture

- Models for Sequences
  - Transformers
  - Positional encoding vs. Rotary Positional Embeddings (RoPE)
  - **Attention & Flash Attention**
  - Pre-norm vs. post-norm
  - Vision Transformers (ViT) & Swin Transformers

# Encoder

# Multi-Head Attention

Image Credit: Vaswani, A. "Attention is all you need." Advances in Neural Information Processing Systems (2017). https://theaisummer.com/transformer/

# Multi-Head Attention

# Multi-Head Attention

# Multi-Head Attention

# Multi-Head Attention

# Multi-Head Attention



Why square root?

# Multi-Head Attention

# Multi-Head Attention

# Layer Norm [1] & Residual Connection



$$\mu_i = \frac{1}{K} \sum_{k=1}^{K} x_{i,k}$$

$$\sigma_i^2 = \frac{1}{K} \sum_{k=1}^{K} (x_{i,k} - \mu_i)^2$$

$$\hat{x}_{i,k} = \frac{x_{i,k} - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}$$

$$y_i = \gamma \hat{x}_i + \beta \equiv \text{LN}_{\gamma,\beta}(x_i)$$

[1] Ba, Jimmy Lei. "Layer normalization." arXiv preprint arXiv:1607.06450 (2016). Image Credit: Vaswani, A. "Attention is all you need." Advances in Neural Information Processing Systems (2017).
https://towardsdatascience.com/illustrated-guide-to-transformers-step-by-step-explanation-f74876522bc0

# Decoder



For certain applications like language models, decoder should be autoregressive!

# Masked Multi-Head Attention



Prevent attending from future!

# Masked Multi-Head Attention

# Masked Multi-Head Attention

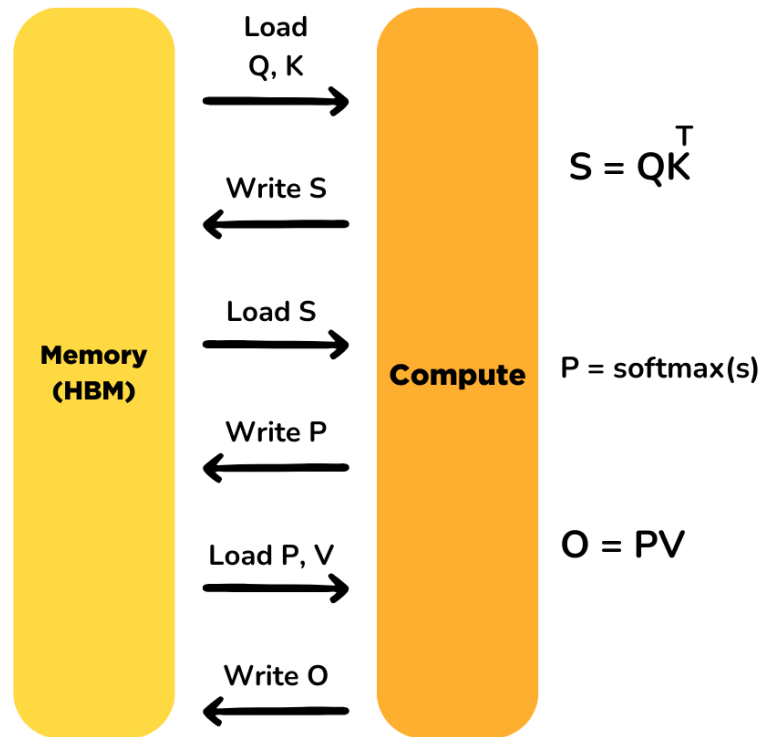# Limitations



- O(L^2) time/memory cost for self-attention

  Methods like Reformer [1] speed up attention to O(L log L) using locality-sensitive hashing techniques

- How can we incorporate prior knowledge into attention rather than having a fully connected attention?

  - Encourage sparse attention

  - Inject known graph structures

  - ……

[1] Kitaev, Nikita, Łukasz Kaiser, and Anselm Levskaya. "Reformer: The efficient transformer." arXiv preprint arXiv:2001.04451 (2020). Image Credit: Vaswani, A. "Attention is all you need." Advances in Neural Information Processing Systems (2017).
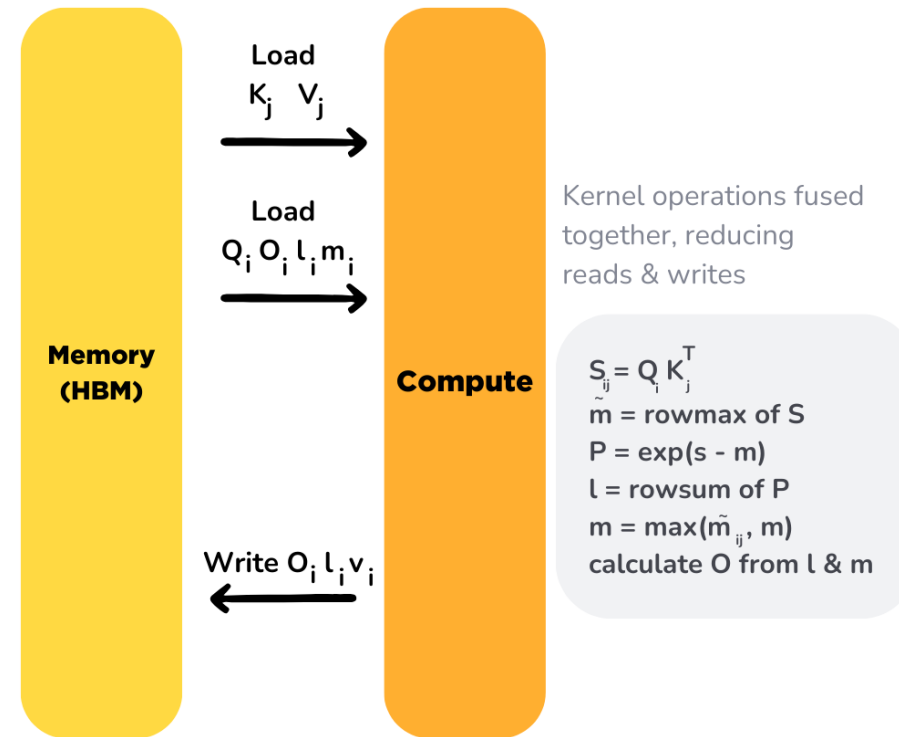
# Flash Attention [1]

Flash attention accelerates attention by using on-chip static random-access memory (SRAM, small memory but fast) to reduce the IO with high bandwidth memory (HBM, large memory but slow).



Standard Attention Implementation

Flash Attention

$$S = QK^T$$

$$P = \text{softmax}(s)$$

$$O = PV$$

Kernel operations fused together, reducing reads & writes

$S_{ij} = Q_i K_j^T$
$\tilde{m} = \text{rowmax of } S$
$P = \exp(s - m)$
$l = \text{rowsum of } P$
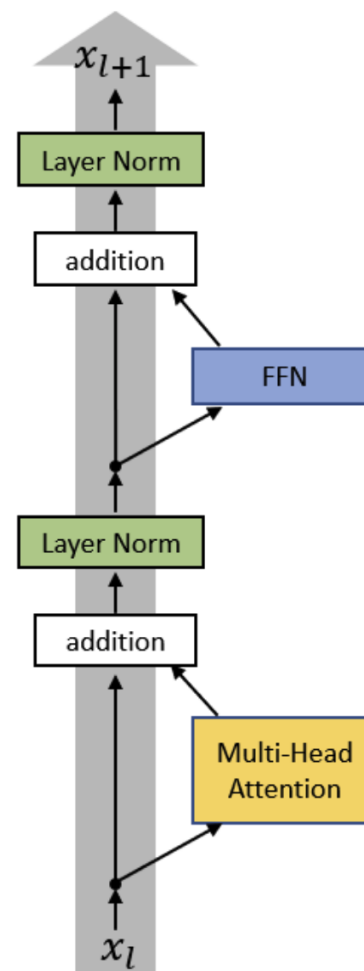$m = \max(\tilde{m}_{ij}, m)$
calculate O from l & m

Initialize O, l and m matrices with zeroes. m and l are used to calculate cumulative softmax. Divide Q, K, V into blocks (due to SRAM's memory limits) and iterate over them, for i is row & j is column.

[1] Dao, Tri, et al. "Flashattention: Fast and memory-efficient exact attention with io-awareness." Advances in Neural Information Processing Systems 35 (2022): 16344-16359. Image Credit: https://huggingface.co/docs/text-generation-inference/en/conceptual/flash_attention
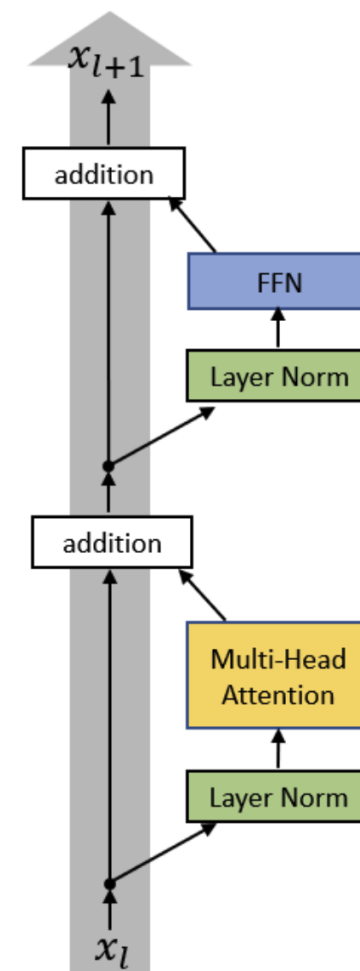
# Outline

- Invariance & Equivariance Principle
  - Translation equivariance in convolutions
  - Permutation equivariance and invariance

- Models for Sets
  - DeepSets: representation theorem of permutation-invariant set functions & architecture
  - DeepSets: permutation-equivariant linear mapping & architecture

- Models for Sequences
  - Transformers
  - Positional encoding vs. Rotary Positional Embeddings (RoPE)
  - Attention & Flash Attention
  - **Pre-norm vs. post-norm**
  - Vision Transformers (ViT) & Swin Transformers

# Pre-Norm vs. Post-Norm

Where to place the Layer Normalization?
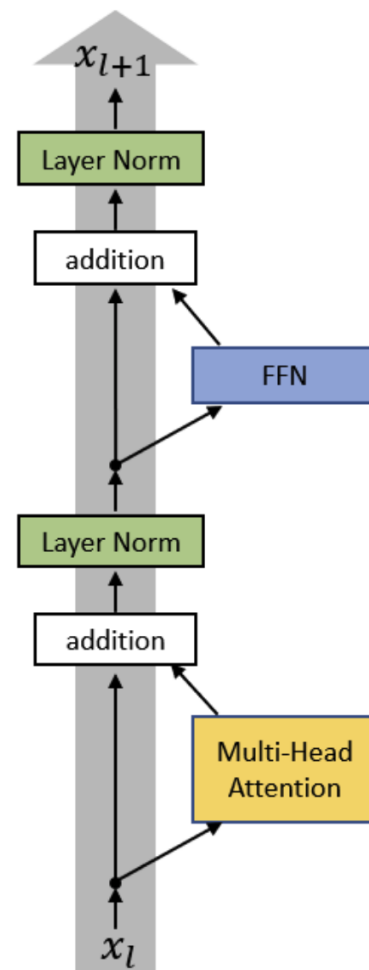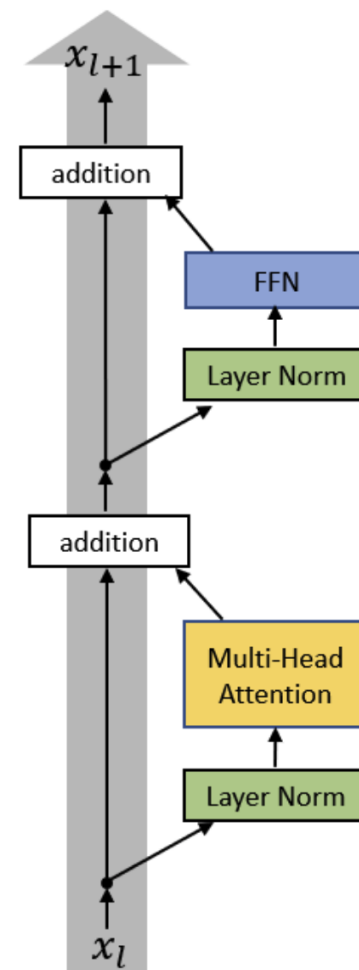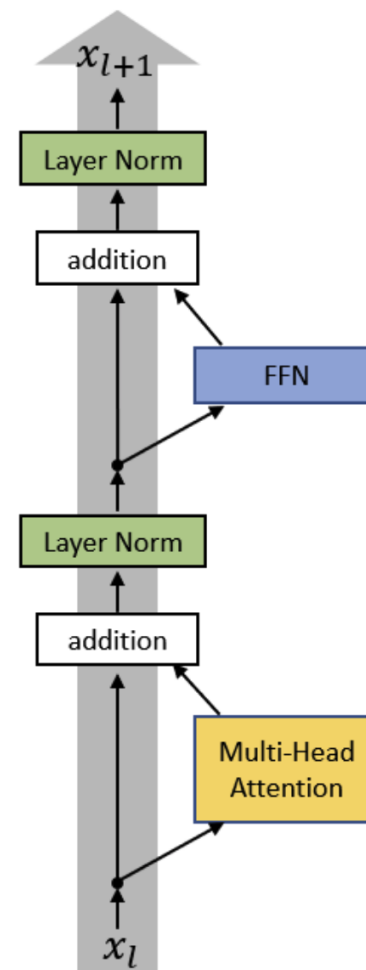


Post-Norm          Pre-Norm

# Pre-Norm vs. Post-Norm

Where to place the Layer Normalization?

- Gradient norm in the Post-Norm Transformer is large for parameters near the output and will be likely to decay as the layer gets closer to input



Post-Norm            Pre-Norm

Image Credit: Xiong, Ruibin, et al. "On layer normalization in the transformer architecture." International Conference on Machine Learning. PMLR, 2020.
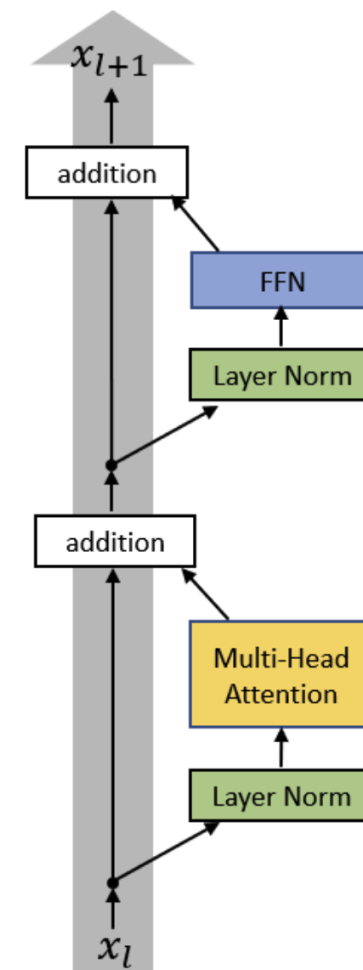
# Pre-Norm vs. Post-Norm

Where to place the Layer Normalization?

- Gradient norm in the Post-Norm Transformer is large for parameters near the output and will be likely to decay as the layer gets closer to input

- Training the Pre-Norm Transformer does not rely on the learning rate warm-up stage and can be trained much faster than the Post-Norm



Post-Norm

Pre-Norm

# Outline

- Invariance & Equivariance Principle
  - Translation equivariance in convolutions
  - Permutation equivariance and invariance

- Models for Sets
  - DeepSets: representation theorem of permutation-invariant set functions & architecture
  - DeepSets: permutation-equivariant linear mapping & architecture

- Models for Sequences
  - Transformers
  - Positional encoding vs. Rotary Positional Embeddings (RoPE)
  - Attention & Flash Attention
  - Pre-norm vs. post-norm
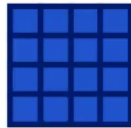  - **Vision Transformers (ViT) & Swin Transformers**

# Extensions: Vision Transformers [1]

[1] Dosovitskiy, Alexey, et al. "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale." International Conference on Learning Representations. 2020. Image Credit: https://github.com/lucidrains/vit-pytorch

# Extensions: Swin Transformers [1]

**Standard MSA**

Attention for each patch is computed against all patches, resulting in quadratic complexity

[1] Liu, Ze, et al. "Swin transformer: Hierarchical vision transformer using shifted windows." Proceedings of the IEEE/CVF international conference on computer vision. 2021. Image Credit: https://towardsdatascience.com/a-comprehensive-guide-to-swin-transformer-64965f89d14c
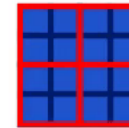
# Extensions: Swin Transformers



**Standard MSA**
Attention for each patch is computed against all patches, resulting in quadratic complexity
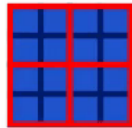
**Window-based MSA**
Attention for each patch is only computed within its own window (drawn in red). Window size is 2x2 in this example.

Image Credit: https://towardsdatascience.com/a-comprehensive-guide-to-swin-transformer-64965f89d14c
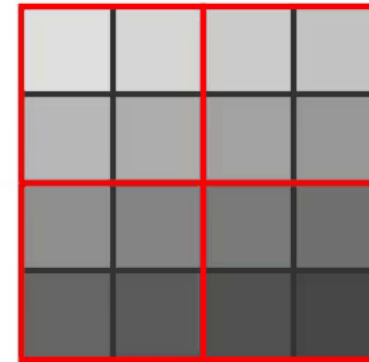
# Extensions: Swin Transformers

**Window-based MSA**

Attention for each patch is only computed within its own window (drawn in red).
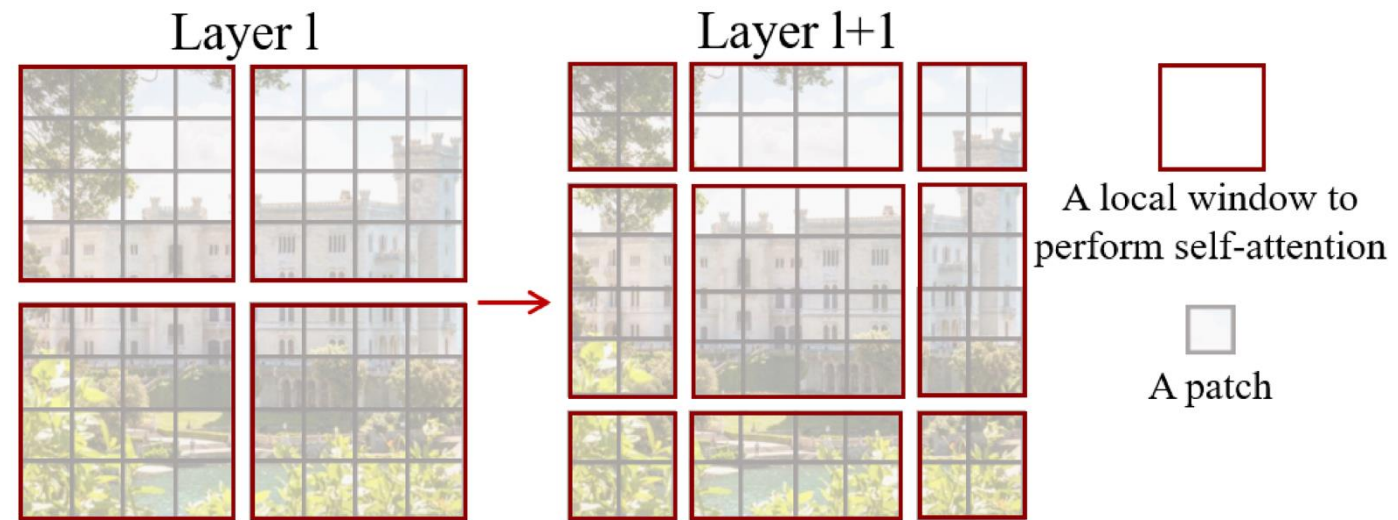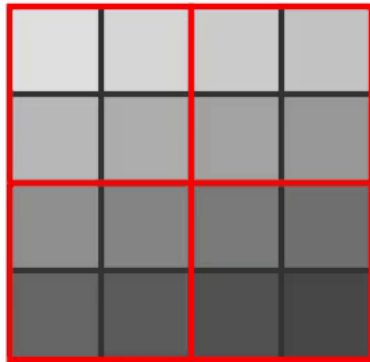Window size is 2x2 in this example.

**Shifted Window MSA**

Step 1: Shift window by a factor of M/2, where M = window size
Step 2: For efficient batch computation, move patches into empty slots to create a complete window.
This is known as 'cyclic shift' in the paper.

Image Credit: https://towardsdatascience.com/a-comprehensive-guide-to-swin-transformer-64965f89d14c

# Extensions: Swin Transformers



**Shifted Window MSA**

Step 1: Shift window by a factor of M/2, where M = window size
Step 2: For efficient batch computation, move patches into empty slots to create a complete window.
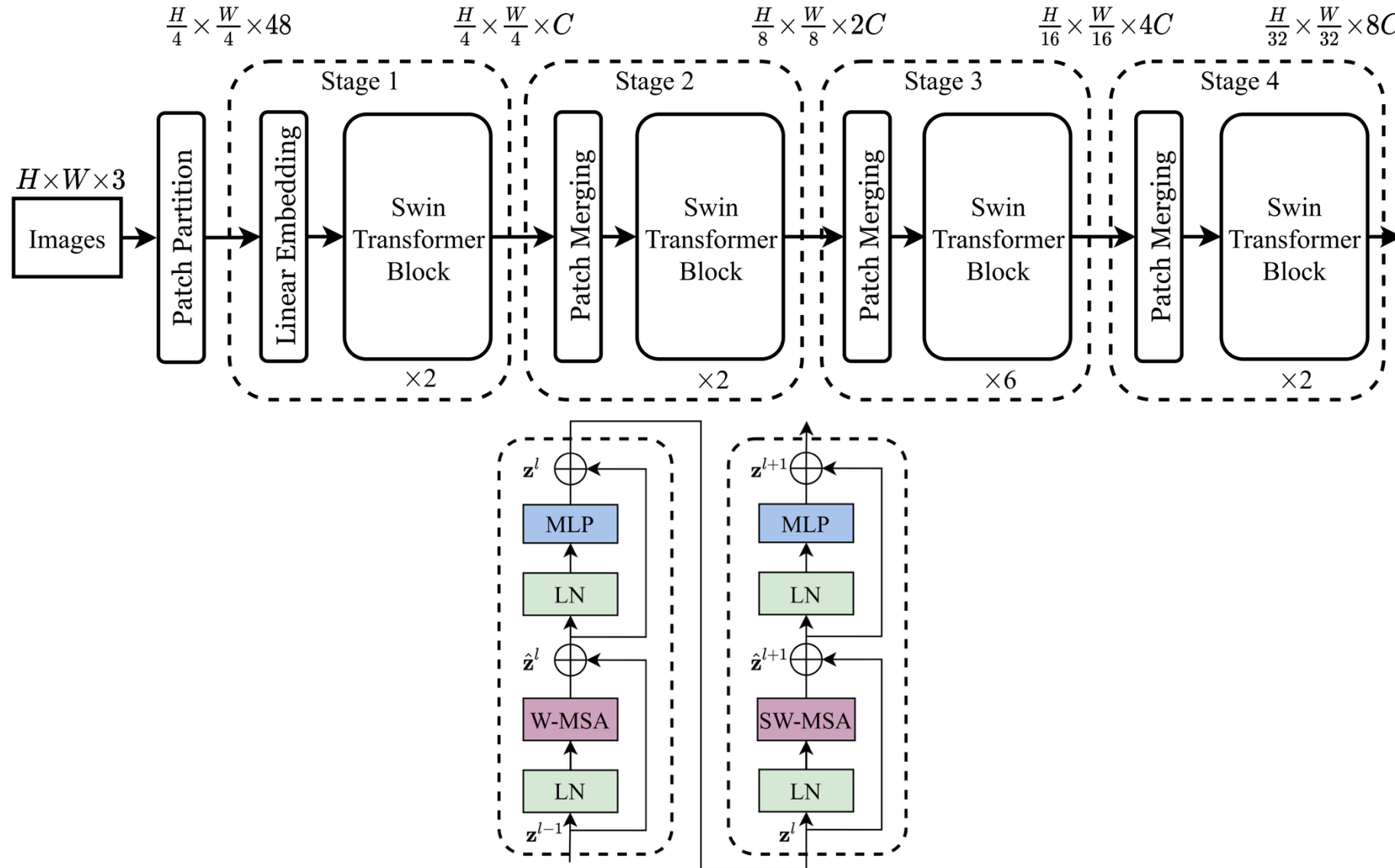This is known as 'cyclic shift' in the paper.

Layer l

Layer l+1

A local window to perform self-attention

A patch

# Extensions: Swin Transformers

# Questions?