# FSAN/ELEG815: Statistical Learning

Gonzalo R. Arce

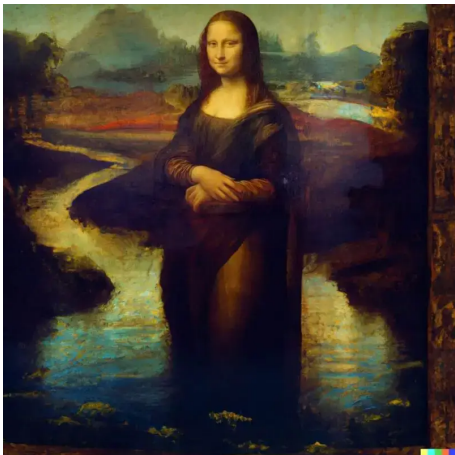**Department of Electrical and Computer Engineering**
**University of Delaware**
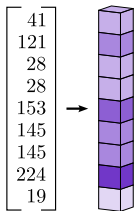
Transformers

# Transformers - Revolutionary Architecture

► ChatGPT is based on the GPT (Generative Pretrained Transformer) architecture.

► Introduced in the paper "Attention is All You Need" by Vaswani et al. in 2017.

► Excel in NLP and Imaging tasks thanks to their capacity to incorporate extensive context. Outperforms in image classification, segmentation, and machine translation.

► The name "transformer" reflects the ability to seamlessly *transform* one sequence of data into another, thanks to its sophisticated self-attention mechanisms.

# Transformers - Image Generation
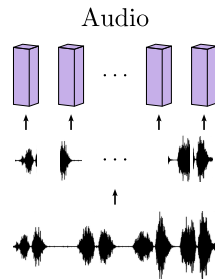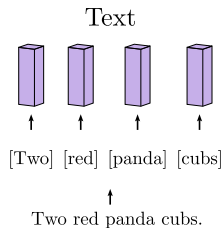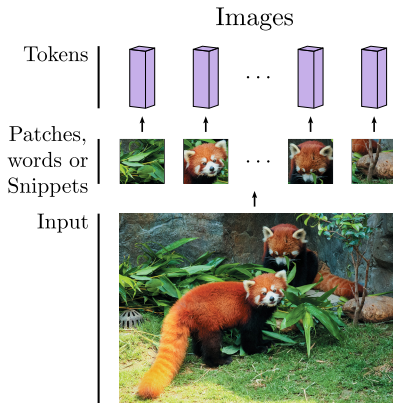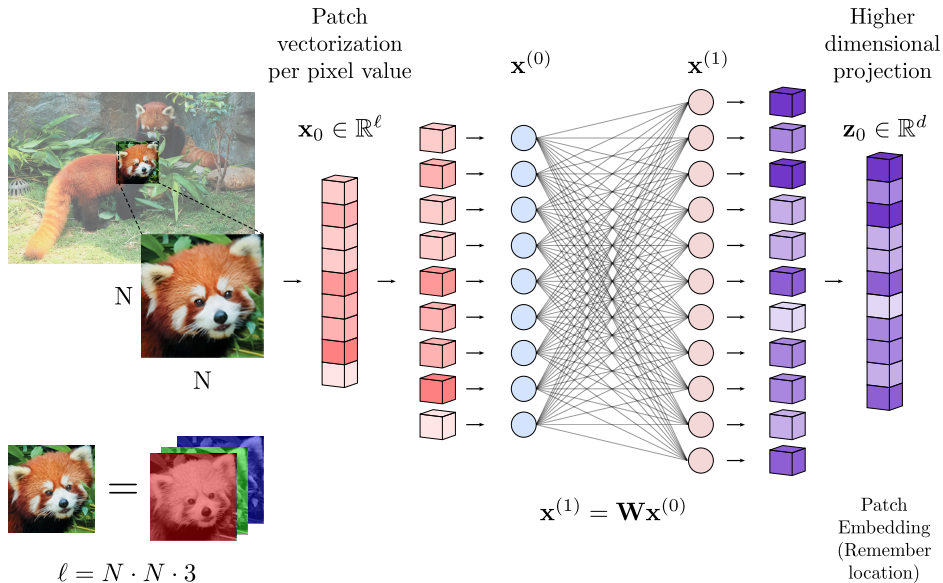
# Tokens and Input Tokenizing

# Image Tokenization and Linear Transformation



Patch vectorization per pixel value

$\mathbf{x}_0 \in \mathbb{R}^\ell$

$\mathbf{x}^{(0)}$

$\mathbf{x}^{(1)}$

Higher dimensional projection

$\mathbf{z}_0 \in \mathbb{R}^d$

$\mathbf{x}^{(1)} = \mathbf{W}\mathbf{x}^{(0)}$

Patch Embedding (Remember location)

$\ell = N \cdot N \cdot 3$

# Patch Embedding Has Meaning



$Z$ Space

# Coloring Problem



Missing color



$= 0.29$  $+\ 0.58$  $+\ 0.11$

# Coloring Problem - Query Vector



Patch projection

$\mathbf{z}_0 \in \mathbb{R}^d$

(learned)

$\mathbf{W}_q \in \mathbb{R}^{d' \times d}$    $\mathbf{z}_0$    $\mathbf{q}_0 \in \mathbb{R}^{d'}$

"Looking for panda's head color"

In general, $d' \geq d$

# Coloring Problem



"Looking for panda's head color"

query

# Coloring Problem - Key Matrix



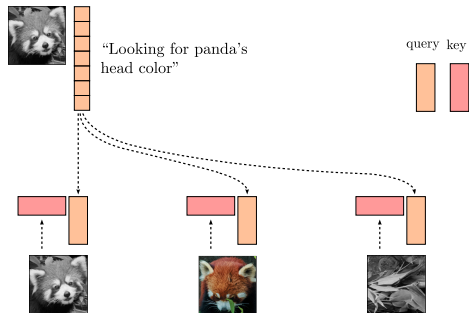$$\mathbf{k}_i = \mathbf{W}_k \mathbf{z}_i$$

$$\mathbf{K} = \mathbf{Z}\mathbf{W}_k$$

# Coloring Problem

# Coloring Problem - Self-Attention



$$\mathbf{a}_0 = \mathrm{softmax} \left( \frac{\mathbf{K}^T \mathbf{q}_0}{\sqrt{d'}} \right)$$

$\sqrt{d'}$ Scales the dot product for numerical stability on the softmax function and balancing signal magnitudes with respect to the dimensionality

$$\mathrm{softmax}(x_i) = \frac{e^{x_i}}{\sum_{j=0}^{t} e^{x_j}}$$

# Coloring Problem - Self-Attention Heat Map



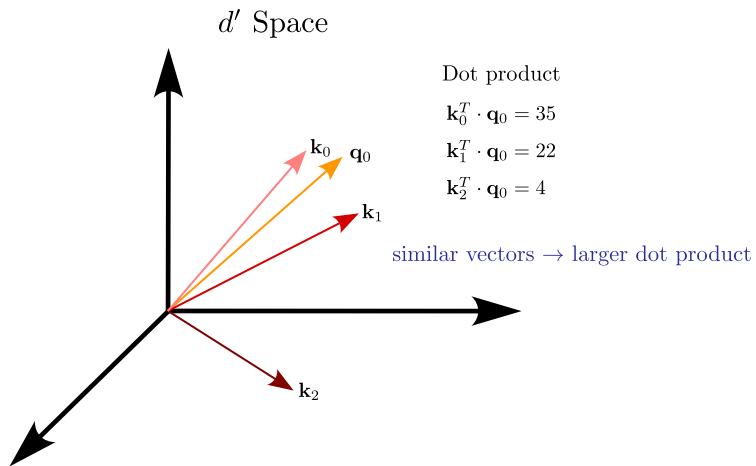Attention to each patch

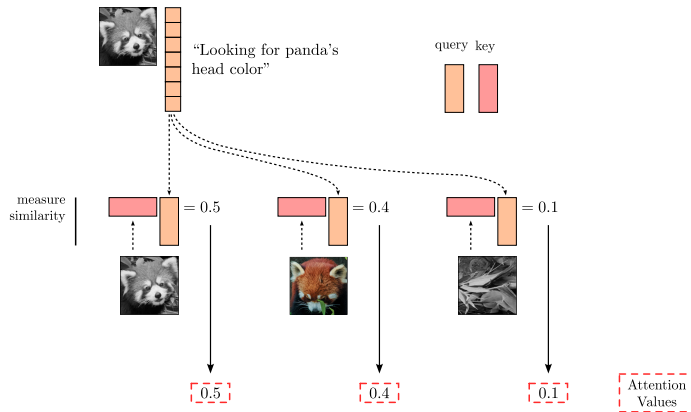To generate vegetation's color

To generate panda's color
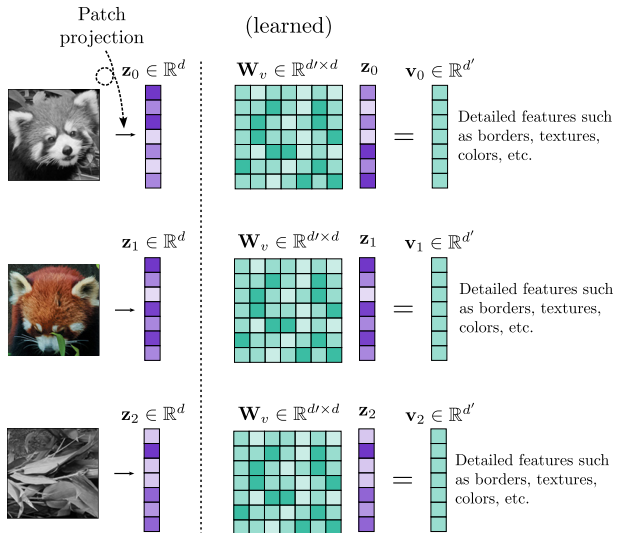
High attention

Low attention
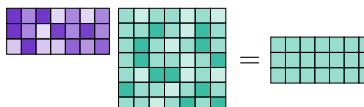
# Why Dot Product for Similarity?



$d'$ Space

Dot product

$\mathbf{k}_0^T \cdot \mathbf{q}_0 = 35$

$\mathbf{k}_1^T \cdot \mathbf{q}_0 = 22$

$\mathbf{k}_2^T \cdot \mathbf{q}_0 = 4$

similar vectors → larger dot product

# Coloring Problem

# Coloring Problem - Value Matrix



Patch projection

$\mathbf{z}_0 \in \mathbb{R}^d$

(learned)

$\mathbf{W}_v \in \mathbb{R}^{d' \times d}$ $\quad \mathbf{z}_0 \quad$ $\mathbf{v}_0 \in \mathbb{R}^{d'}$

Detailed features such as borders, textures, colors, etc.

$\mathbf{z}_1 \in \mathbb{R}^d$

$\mathbf{W}_v \in \mathbb{R}^{d' \times d}$ $\quad \mathbf{z}_1 \quad$ $\mathbf{v}_1 \in \mathbb{R}^{d'}$

Detailed features such as borders, textures, colors, etc.

$\mathbf{z}_2 \in \mathbb{R}^d$

$\mathbf{W}_v \in \mathbb{R}^{d' \times d}$ $\quad \mathbf{z}_2 \quad$ $\mathbf{v}_2 \in \mathbb{R}^{d'}$

Detailed features such as borders, textures, colors, etc.
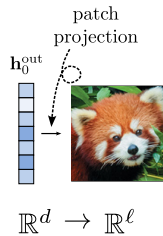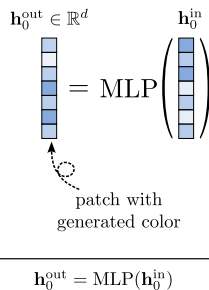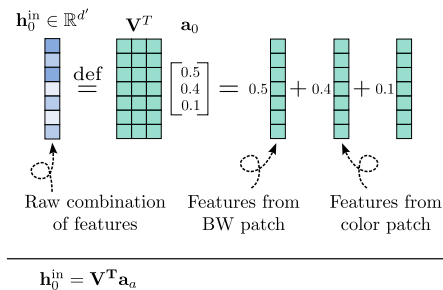
$$\mathbf{k}_i = \mathbf{W}_k \mathbf{z}_i$$

$\mathbf{Z} \in \mathbb{R}^{t \times d}$ $\quad \mathbf{W}_v \in \mathbb{R}^{d \times d'}$ $\quad \mathbf{V} \in \mathbb{R}^{t \times d'}$
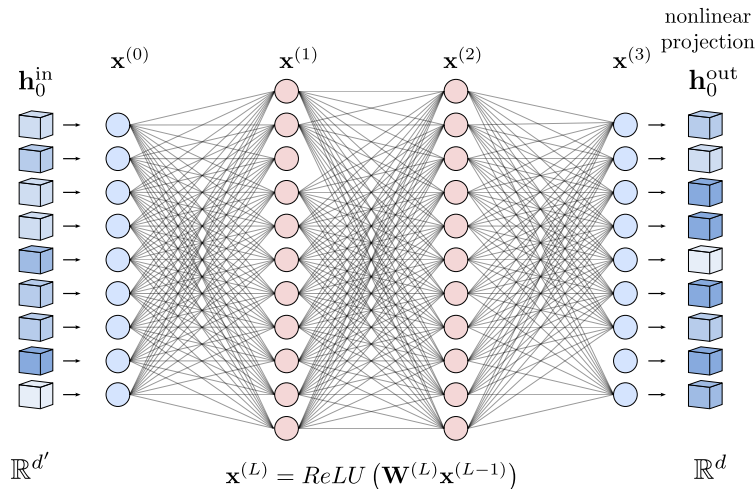
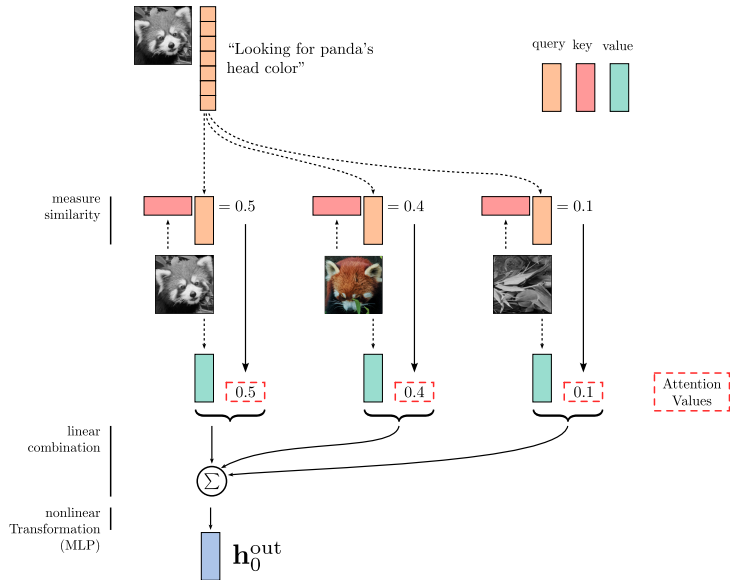$$\mathbf{V} = \mathbf{Z} \mathbf{W}_v$$

# Coloring Problem

# Coloring Problem - Hidden Representation



$$\mathbf{h}_0^{\text{in}} \in \mathbb{R}^{d'} \qquad \mathbf{V}^T \qquad \mathbf{a}_0$$

$$\mathbf{h}_0^{\text{in}} \overset{\text{def}}{=} \mathbf{V}^T \begin{bmatrix} 0.5 \\ 0.4 \\ 0.1 \end{bmatrix} = 0.5 \square + 0.4 \square + 0.1 \square$$

Raw combination of features    Features from BW patch    Features from color patch

$$\mathbf{h}_0^{\text{in}} = \mathbf{V^T a}_a$$

$$\mathbf{h}_0^{\text{out}} \in \mathbb{R}^d \qquad \mathbf{h}_0^{\text{in}}$$

$$\mathbf{h}_0^{\text{out}} = \text{MLP}\left( \mathbf{h}_0^{\text{in}} \right)$$

patch with generated color

$$\mathbf{h}_0^{\text{out}} = \text{MLP}(\mathbf{h}_0^{\text{in}})$$

patch projection

$$\mathbf{h}_0^{\text{out}}$$

$$\mathbb{R}^d \;\rightarrow\; \mathbb{R}^\ell$$

# Nonlinear Transformation of Tokens - MLP



$$\mathbf{x}^{(L)} = ReLU\left(\mathbf{W}^{(L)}\mathbf{x}^{(L-1)}\right)$$

# Coloring Problem

# Coloring Problem - $\mathbf{Q}, \mathbf{K}$ and $\mathbf{V}$ Matrices



$$\mathbf{Q} = \mathbf{Z}\mathbf{W}_q \qquad \mathbf{K} = \mathbf{Z}\mathbf{W}_k \qquad \mathbf{V} = \mathbf{Z}\mathbf{W}_v$$

$$\mathbf{A} = \mathrm{softmax}\left(\frac{\mathbf{Q}^T \mathbf{K}}{\sqrt{d'}}\right)$$

# Coloring Problem - Attention Map

Attention map



□ High attention

■ Low attention

$$\text{MLP} \left( \underset{\mathbf{A}}{\boxed{\phantom{xx}}} \; \underset{\mathbf{V}}{\boxed{\phantom{xxxx}}} \right) = \underset{\mathbf{H}^{\text{out}} \in \mathbb{R}^{t \times d}}{\boxed{\phantom{xxxx}}}$$

patch projection

Patch Projection $\overset{\text{def}}{=} \text{MLP}(\mathbf{H}^{out})$

$$\mathbf{H}^{\text{out}} = \text{MLP}(\mathbf{AV})$$

# Position in the Matrix Has no Influence

# Positional Embedding - Encode Position

▶ Solution: Add a positional vector to each patch projection.

# Positional embedding - Sinusoidal Waves

$$\text{PE}(i, 2j) = \sin\left(\frac{i}{10000^{2j/d}}\right)$$

$$\text{PE}(i, 2j+1) = \cos\left(\frac{i}{10000^{2j/d}}\right)$$

where $i$ is the position of the patch in the sequence, $j$ for $j = 0, \ldots, d$ is the dimension within the positional vector and $d$ is the size of the positional vector.

For each dimension of $\mathbf{t}_i$, there is a sinusoidal wave with different frequency

$$\mathbf{t}_{i,0} = \sin\left(\frac{i}{10000^{2(0)/d}}\right)$$

$$\mathbf{t}_{i,1} = \cos\left(\frac{i}{10000^{2(0)/d}}\right)$$

$$\mathbf{t}_{i,2} = \sin\left(\frac{i}{10000^{2(1)/d}}\right)$$

$\mathbf{t}_i$

$j = 0$

$j = 2$

$j = 4$

$j = 0 \rightarrow$
$j = 2 \rightarrow$
$j = 4 \rightarrow$

$\cdots$

$\mathbf{t}_0$          $\mathbf{t}_5$

# Positional embedding - Sinusoidal Waves

▶ This function provides a large number of vectors with a constant distance between them, i.e. $\|\mathbf{t}_i - \mathbf{t}_{i+1}\|_2$ is constant for all $i$.



Positional Encoding Vector distance

# Attention Layer - Normalization



$$\mathbf{h}_0^{\text{in}} = f_{\text{norm}}(\mathbf{V}^T \mathbf{a}_0 + \hat{\mathbf{z}}_0)$$

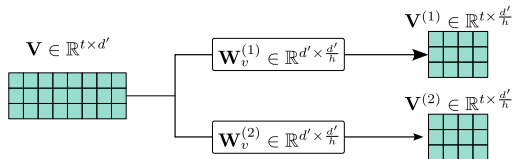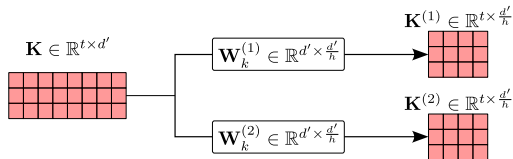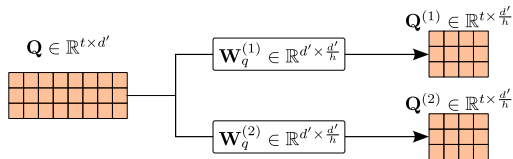$$\mathbf{h}_0^{\text{out}} = f_{\text{norm}}(\text{MLP}(\mathbf{h}_0^{\text{in}}) + \mathbf{h}_0^{\text{in}})$$

# Single Attention Layer - Single-Head and Multi Head

# Single Attention Layer - $\mathbf{Q}, \mathbf{K}$ and $\mathbf{V}$ split



for $h = 2$ (two heads)

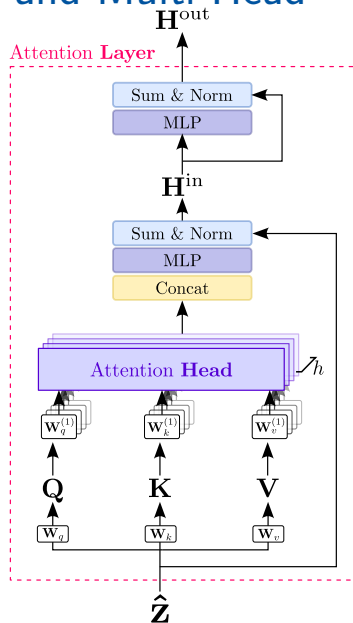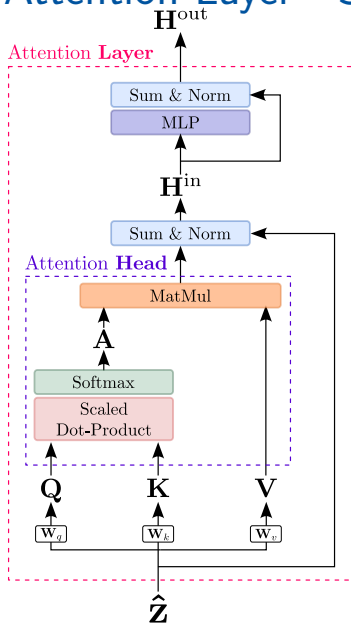# Single Attention Layer - Multi-Head Attention



for $h = 2$ (two heads)

# Single Attention Layer - Multi-Head Concat



for $h = 2$ (two heads)

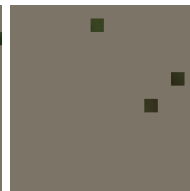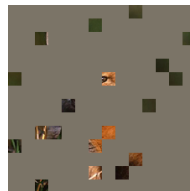# Single Attention Layer - Single Head and Multi Head

# Transformer

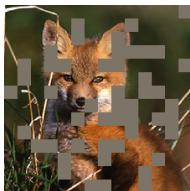# Vision Transformer - Reconstructing Masked Image
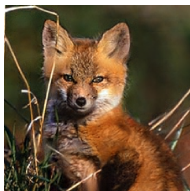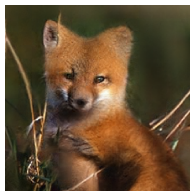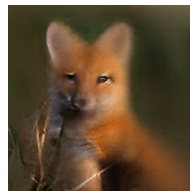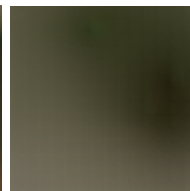


Ground Truth

Masked Image

Reconstruction

25%　　　75%　　　90%　　　98%

Mask Ratio

# Vision Transformer - Reconstructing Masked Image