

```
1 using LinearAlgebra, Random, Plots, PlutoUI, Images, Interact#, MLDatasets
```

Our project is to implement a simple Vision Transformer in Julia!

The Transformer architecture, introduced in the paper *Attention Is All You Need* (Vaswani et al., 2017), is the most ubiquitous neural network architecture in modern machine learning. Its parallelism and scalability to large problems has seen it adopted in domains beyond those it was traditionally considered for (sequential data) and it quickly replaced convolutional neural networks for image-based tasks.

```
1 md"""
2 ### Our project is to implement a simple Vision Transformer in Julia!
3
4 The Transformer architecture, introduced in the paper _Attention Is All You Need_
  (Vaswani et al., 2017), is the most ubiquitous neural network architecture in modern
  machine learning. Its parallelism and scalability to large problems has seen it
  adopted in domains beyond those it was traditionally considered for (sequential
  data) and it quickly replaced convolutional neural networks for image-based tasks.
5 """
```

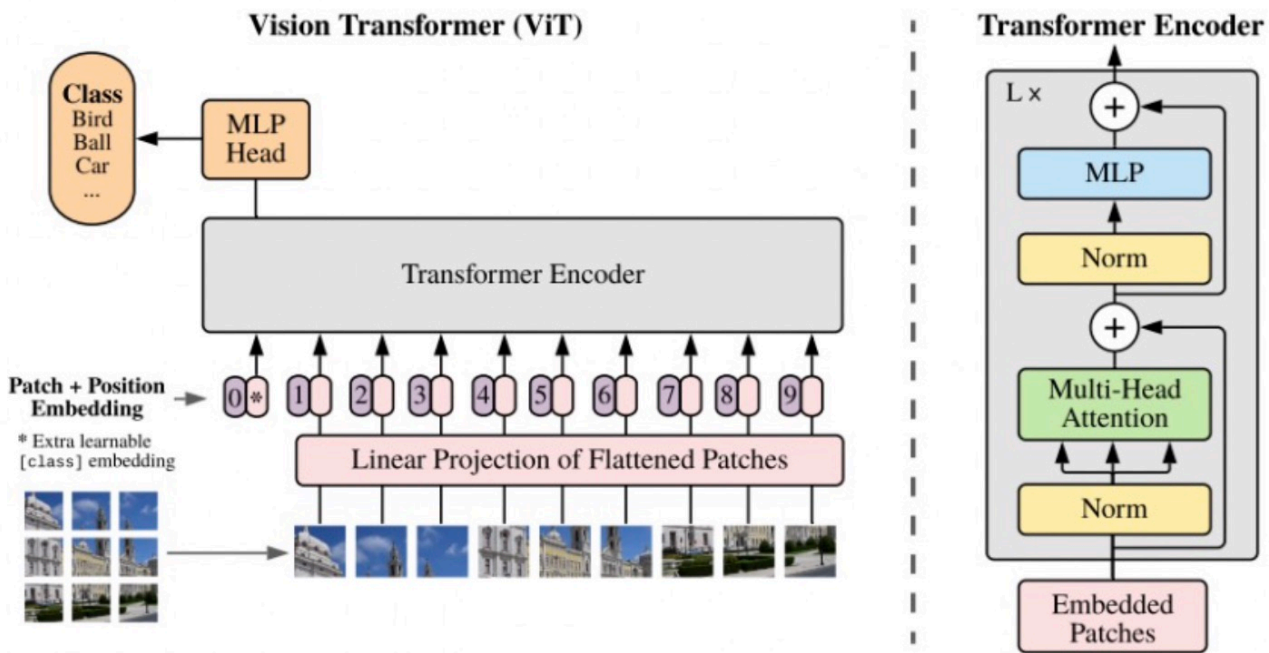


Figure 2: ViT Architecture, Figure from Dosovitskiy et al. [2020]

```
1 md"""
2 ![ViT Model]
  (https://github.com/qsimeon/julia_class_project/blob/e698587c2c2b7455404e6126c06f4ec0
  4c463032/vit_arch.jpg?raw=true)
3 """
```

Let's start by defining key components of a Vision Transformer (ViT) model using Julia structs and parametric types, similar to the structure we implemented in Homework 3. We will implement the `AttentionHead`, `MultiHeadedAttention`, and `FeedForwardNetwork` layers as Julia structs. This will set up the parts which get combined together in the `Transformer` model.

```
1 md"""
2 Let's start by defining key components of a Vision Transformer (ViT) model using
  Julia structs and parametric types, similar to the structure we implemented in
  Homework 3. We will implement the `AttentionHead`, `MultiHeadedAttention`, and
  `FeedForwardNetwork` layers as Julia structs. This will set up the parts which get
  combined together in the `Transformer` model.
3 """
```

`softmax` (generic function with 1 method)

```
1 # Stable softmax implementation
2 function softmax(x; dims=1)
3     exp_x = exp.(x .- maximum(x, dims=dims)) # stability trick
4     return exp_x ./ sum(exp_x, dims=dims)
5 end
```

```

1  ### 1. Attention Head
2  struct AttentionHead{T<:Real}
3      W_K::Matrix{T} # Shape: (n_hidden, dim)
4      W_Q::Matrix{T} # Shape: (n_hidden, dim)
5      W_V::Matrix{T} # Shape: (dim, dim)
6      n_hidden::Int # dimensionality of key and query vectors
7
8      function AttentionHead{T}(dim::Int, n_hidden::Int) where T<:Real
9          return new{T}(randn(T, n_hidden, dim), randn(T, n_hidden, dim), randn(T,
10 dim, dim), n_hidden)
11 end
12
13 function (head::AttentionHead{T})(X::Matrix{T}, attn_mask::Union{Nothing,
14 Matrix{T}}=nothing) where {T<:Real}
15     # X is expected to be an input token matrix with shape (N, dim)
16     # Project input tokens to query, key, and value representations
17     Q = X * transpose(head.W_Q) # Shape: (N, n_hidden)
18     K = X * transpose(head.W_K) # Shape: (N, n_hidden)
19     V = X * transpose(head.W_V) # Shape: (N, dim)
20
21     # Compute scaled dot-product attention
22     scores = Q * transpose(K) / sqrt(head.n_hidden) # Shape: (N, N)
23
24     # Apply attention mask if provided
25     if attn_mask != nothing
26         scores = scores .* attn_mask .+ (1 .- attn_mask) * -Inf
27     end
28
29     # Apply softmax along the last dimension
30     alpha = softmax(scores, dims=ndims(scores)) # Shape: (N, N)
31
32     # Compute attention output as weighted sum of values
33     attn_output = alpha * V # Shape: (N, dim)
34
35     # attn_output is the (N, dim) output token matrix
36     # alpha is the (N, N) attention matrix
37     return attn_output, alpha
38 end
end

```

 5

```
1 @bind n_tokens Slider(5:20, show_value=true)
```

Attention Matrix



```
1 # Test 'AttentionHead' implementation
2 let
3   dim, attn_dim = 3, 8
4   head = AttentionHead{Float64}(dim, attn_dim)
5   X = randn(Float64, n_tokens, dim) # example 3-D input of n_tokens
6   attn_output, alpha = head(X)
7   println("attention output shape: ", size(attn_output))
8   println("attention weight shape: ", size(alpha))
9   heatmap(
10     alpha,
11     aspect_ratio=:equal,
12     xlabel="Token Index",
13     ylabel="Token Index",
14     title="Attention Matrix",
15     c=:plasma, # Choose a colormap, e.g., :viridis or :plasma
16
17     clabel="Weight", # Label for the color bar
18     colorbar=true, # Show the color bar
19     grid=false, # Turn off the grid
20     # framestyle=:none, # Removes the axis lines
21     xticks=1:n_tokens, # Ensures ticks are at each integer index
22     yticks=1:n_tokens,
23   )
24
25 end
```

```
attention output shape: (5, 3)
attention weight shape: (5, 5)
```

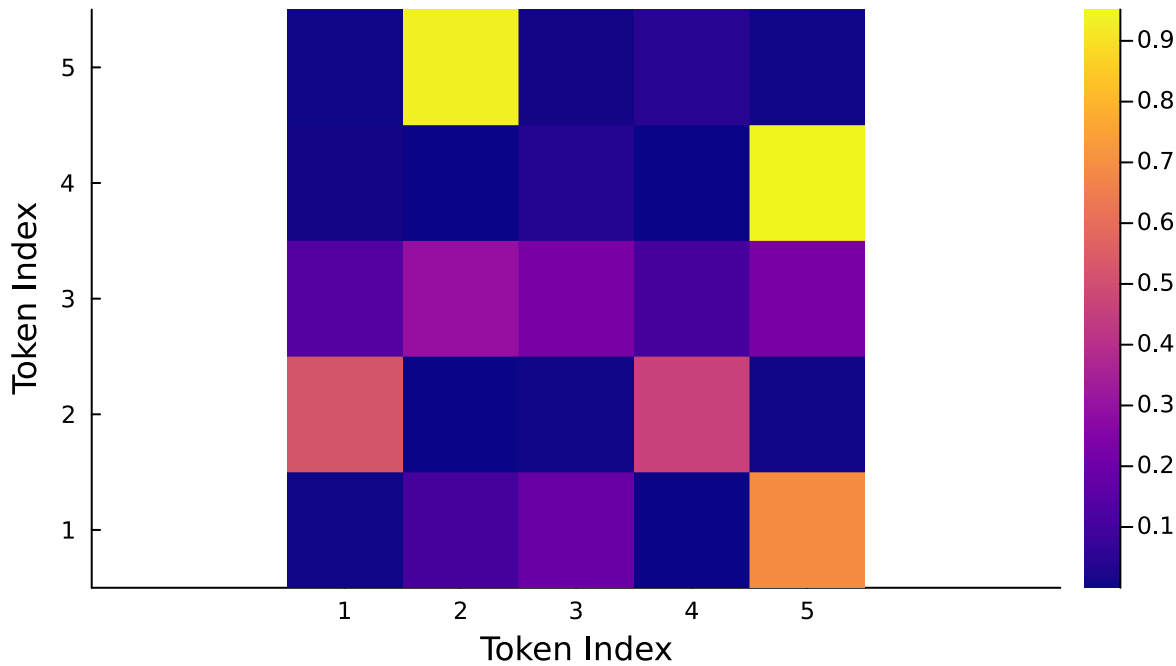


```

1  ### 2. Multi-Headed Attention
2  struct MultiHeadedAttention{T<:Real}
3      heads::Vector{AttentionHead{T}}
4      W_msa::Matrix{T} # Shape: (dim, num_heads*dim)
5
6      function MultiHeadedAttention{T}(dim::Int, n_hidden::Int, num_heads::Int) where
T<:Real
7          # Each head outputs dim-dimensional tokens
8          heads = [AttentionHead{T}(dim, n_hidden) for _ in 1:num_heads]
9          # Our MHA outputs tokens with the same dimension as the input tokens
10         W_msa = randn(T, dim, num_heads * dim)
11         return new{T}(heads, W_msa)
12     end
13
14     function (mha::MultiHeadedAttention{T})(X::Matrix{T}, attn_mask::Union{Nothing,
Matrix{T}}=nothing) where {T<:Real}
15         outputs, alphas = [], []
16         for head in mha.heads
17             out, alpha = head(X, attn_mask) # Shapes: (N, dim), (N, N)
18             push!(outputs, out)
19             push!(alphas, alpha)
20         end
21         # Concatenate along hidden dimension
22         concatenated = cat(outputs...; dims=2) # Shape: (N, num_heads*dim)
23         attn_output = concatenated * transpose(mha.W_msa) # Shape: (N, dim)
24         attn_alphas = cat(alphas...; dims=3) # Shape: (N, N, num_heads)
25         attn_alphas = permutedims(attn_alphas, (3, 1, 2)) # Shape: (num_heads, N, N)
26         return attn_output, attn_alphas
27     end
28 end

```

Attention Matrix



```
1 # Test 'MultiHeadedAttention' implementation
2 let
3   dim, attn_dim, num_heads = 3, 8, 5
4   heads = [AttentionHead{Float64}(dim, attn_dim) for _ in 1:num_heads]
5   W_msa = randn(Float64, dim, num_heads * dim)
6   X = randn(Float64, n_tokens, dim) # example 3-D input of n_tokens
7
8   outputs, alphas = [], []
9   for head in heads
10     attn_out, alpha = head(X)
11     push!(outputs, attn_out)
12     push!(alphas, alpha)
13   end
14
15   concatenated = cat(outputs...; dims=2) # Shape: (N, num_heads*dim)
16   attn_alphas = cat(alphas...; dims=3) # Shape: (N, N, num_heads)
17   attn_alphas = permutedims(attn_alphas, (3, 1, 2)) # Shape: (num_heads, N, N)
18   attn_output = concatenated * transpose(W_msa) # Shape: (N, dim)
19
20   println("attention output shape: ", size(attn_output))
21   println("attention weights shape: ", size(attn_alphas))
22   # plot the attention mask from the last head
23   heatmap(
24     attn_alphas[end,:,:],
25     aspect_ratio=:equal,
26     xlabel="Token Index",
27     ylabel="Token Index",
28     title="Attention Matrix",
29     c=:plasma, # Choose a colormap, e.g., :viridis or :plasma
30     clabel="Weight", # Label for the color bar
31     colorbar=true, # Show the color bar
32     grid=false, # Turn off the grid
33     # framestyle=:none, # Removes the axis lines
34     xticks=1:n_tokens, # Ensures ticks are at each integer index
35     yticks=1:n_tokens,
```

```
36 )
37 end
```

```
attention output shape: (5, 3)
attention weights shape: (5, 5, 5)
```



```
1 ### 3. Feed-Forward Network (FFN)
2 struct FeedForwardNetwork{T<:Real}
3     W1::Matrix{T} # Shape: (n_hidden, dim)
4     W2::Matrix{T} # Shape: (dim, n_hidden)
5     b1::Vector{T} # Shape: (n_hidden,)
6     b2::Vector{T} # Shape: (dim,)
7
8     function FeedForwardNetwork{T}(dim::Int, n_hidden::Int) where T<:Real
9         # Our FFN outputs tokens with the same dimension as the input tokens
10        return new{T}(randn(T, n_hidden, dim), randn(T, dim, n_hidden), randn(T,
n_hidden), randn(T, dim))
11    end
12
13    function (ffn::FeedForwardNetwork{T})(X::Matrix{T}) where {T<:Real}
14        # X is expected to be an input token matrix with shape (N, dim)
15        X = X * transpose(ffn.W1) .+ ffn.b1' # Shape: (N, n_hidden)
16        X = max.(0, X) # ReLU activation
17        return X * transpose(ffn.W2) .+ ffn.b2' # Shape: (N, dim)
18    end
19 end
20
```

```
1 # Test 'FeedForwardNetwork' implementation
2 let
3     dim, mlp_dim = 3, 8
4     ffn = FeedForwardNetwork{Float64}(dim, mlp_dim)
5     X = randn(Float64, n_tokens, dim) # example 3-D input of n_tokens
6     ffn_output = ffn(X)
7     println("feedforward output shape: ", size(ffn_output))
8 end
```

```
feedforward output shape: (5, 3)
```



Recap so far

1. AttentionHead Implementation:

- Projects the input token matrix X (shape: (N, dim)) to query, key, and value matrices.
- Computes scaled dot-product attention, applies an optional mask, and then applies softmax to get attention weights α with shape (N, N) .
- Returns the attention output attn_output (shape: (N, dim)) and attention weights α .

2. MultiHeadedAttention Implementation:

- Creates multiple AttentionHead instances and collects their outputs.
- Concatenates these outputs along the hidden dimension, applies a linear transformation (W_{msa}), and stacks the attention weights from each head into a 3D tensor with shape $(\text{num_heads}, N, N)$.

3. FeedForwardNetwork (FFN) Implementation:

- A two-layer feed-forward network with an intermediate hidden layer of size n_{hidden} .
- Projects the input token matrix X (shape: (N, dim)) to an intermediate hidden representation (shape: (N, n_{hidden})) using W_1 and b_1 , followed by a ReLU activation.
- Transforms the hidden representation back to the original input dimension dim using W_2 and b_2 .
- Returns the output with shape (N, dim) , maintaining the same dimension as the input tokens.

```
1 md"""
2 ### Recap so far
3
4 1. AttentionHead Implementation:
5     - Projects the input token matrix 'X' (shape:  $(N, \text{dim})$ ) to query,
6       key, and value matrices.
7     - Computes scaled dot-product attention, applies an optional mask, and then
8       applies softmax to get attention weights 'alpha' with shape  $(N, N)$ .
9     - Returns the attention output 'attn_output' (shape:  $(N, \text{dim})$ ) and
10       attention weights 'alpha'.
11
12
13 2. MultiHeadedAttention Implementation:
14     - Creates multiple 'AttentionHead' instances and collects their outputs.
15     - Concatenates these outputs along the hidden dimension, applies a linear
16       transformation ('W_msa'), and stacks the attention weights from each head into a 3D
17       tensor with shape '(num_heads, N, N)'.
18
19
20 3. FeedForwardNetwork (FFN) Implementation:
21     - A two-layer feed-forward network with an intermediate hidden layer of size
22       'n_hidden'.
23     - Projects the input token matrix 'X' (shape:  $(N, \text{dim})$ ) to an
24       intermediate hidden representation (shape:  $(N, n_{\text{hidden}})$ ) using 'W1' and
25       'b1', followed by a ReLU activation.
26     - Transforms the hidden representation back to the original input dimension
27       'dim' using 'W2' and 'b2'.
28     - Returns the output with shape  $(N, \text{dim})$ , maintaining the same
29       dimension as the input tokens.
30 """
31
```



```

1  ### 4. Attention Residual
2  struct AttentionResidual{T<:Real}
3      attn::MultiHeadedAttention{T}  # Multi-headed attention mechanism
4      ffn::FeedForwardNetwork{T}     # Feed-forward network
5
6      # Constructor: initializes attention and feed-forward sub-layers
7      function AttentionResidual{T}(dim::Int, attn_dim::Int, mlp_dim::Int,
8          num_heads::Int) where T<:Real
9          attn_layer = MultiHeadedAttention{T}(dim, attn_dim, num_heads)
10         ffn_layer = FeedForwardNetwork{T}(dim, mlp_dim)
11         return new{T}(attn_layer, ffn_layer)
12     end
13
14     # Apply the AttentionResidual block to input x
15     function (residual::AttentionResidual{T})(X::Matrix{T},
16         attn_mask::Union{Nothing, Matrix{T}}=nothing) where {T<:Real}
17         # Apply the multi-headed attention layer
18         attn_out, alphas = residual.attn(X, attn_mask) # attn_out: (N, dim),
19         # alphas: (num_heads, N, N)
20         # First residual connection with attention output
21         X = X .+ attn_out
22         # Apply the feed-forward network and add the second residual connection
23         X = X .+ residual.ffn(X)
24         # Return the final output and attention weights
25         return X, alphas
26     end
27 end

```

```

1  ### 5. Attention Residual
2  struct Transformer{T<:Real}
3      layers::Vector{AttentionResidual{T}} # Sequence of AttentionResidual blocks
4
5      # Constructor: initializes a sequence of attention residual blocks
6      function Transformer{T}(dim::Int, attn_dim::Int, mlp_dim::Int, num_heads::Int,
7          num_layers::Int) where T<:Real
8          layers = [AttentionResidual{T}(dim, attn_dim, mlp_dim, num_heads) for _ in
9              1:num_layers]
10         return new{T}(layers)
11     end
12
13     # Apply the Transformer model to input X
14     function (transformer::Transformer{T})(X::Matrix{T}, attn_mask::Union{Nothing,
15         Matrix{T}}=nothing) where {T<:Real}
16         collected_alphas = [] # To store attention weights from each layer
17         for layer in transformer.layers
18             X, alphas = layer(X, attn_mask) # Apply each residual block
19             push!(collected_alphas, alphas) # Collect attention weights
20         end
21         # Return the final output and collected attention weights from all layers
22         return X, collected_alphas
23     end
24 end

```

Testing the AttentionResidual and Transformer

Let's test the AttentionResidual and Transformer structs to confirm that they work as expected with the previously implemented components.

```
1 md"""
2 ### Testing the AttentionResidual and Transformer
3
4 Let's test the 'AttentionResidual' and 'Transformer' structs to confirm that they
  work as expected with the previously implemented components.
5 """
```

```
1 # Test 'AttentionResidual' implementation
2 let
3   dim, attn_dim, mlp_dim, num_heads = 8, 16, 32, 3
4   residual_block = AttentionResidual{Float64}(dim, attn_dim, mlp_dim, num_heads)
5   X = randn(Float64, n_tokens, dim) # example input with n_tokens, each of 'dim'
    dimensions
6   output, alphas = residual_block(X)
7   println("AttentionResidual output shape: ", size(output))
8   println("Attention weights shape (from one layer): ", size(alphas))
9 end
```

```
AttentionResidual output shape: (5, 8)
Attention weights shape (from one layer): (3, 5, 5)
```

```
1 # Test 'Transformer' implementation
2 let
3   dim, attn_dim, mlp_dim, num_heads, num_layers = 8, 16, 32, 3, 6
4   transformer = Transformer{Float64}(dim, attn_dim, mlp_dim, num_heads, num_layers)
5   X = randn(Float64, n_tokens, dim) # example input with n_tokens, each of 'dim'
    dimensions
6   output, collected_alphas = transformer(X)
7   println("Transformer output shape: ", size(output))
8   println("Collected attention weights shape: ", size(collected_alphas[1]), " for
    ", num_layers, " layers")
9 end
```

```
Transformer output shape: (5, 8)
Collected attention weights shape: (3, 5, 5) for 6 layers
```

Our modules so far build up the Transformer

- **AttentionHead**: Implements a single attention head, creating query, key, and value projections, computing the attention scores, and applying a softmax.
- **MultiHeadedAttention**: Combines multiple `AttentionHead`s, concatenates their outputs, and applies a final linear transformation.
- **FeedForwardNetwork**: A simple feed-forward network with two linear layers and a ReLU activation in between.
- **AttentionResidual**: Combines multi-head attention and feed-forward network layers with residual connections.
- **Transformer**: Stacks multiple `AttentionResidual` layers to form the complete Transformer encoder.

We've set up a basic structure of a Transformer using callable structs, parametric types, and matrix operations.

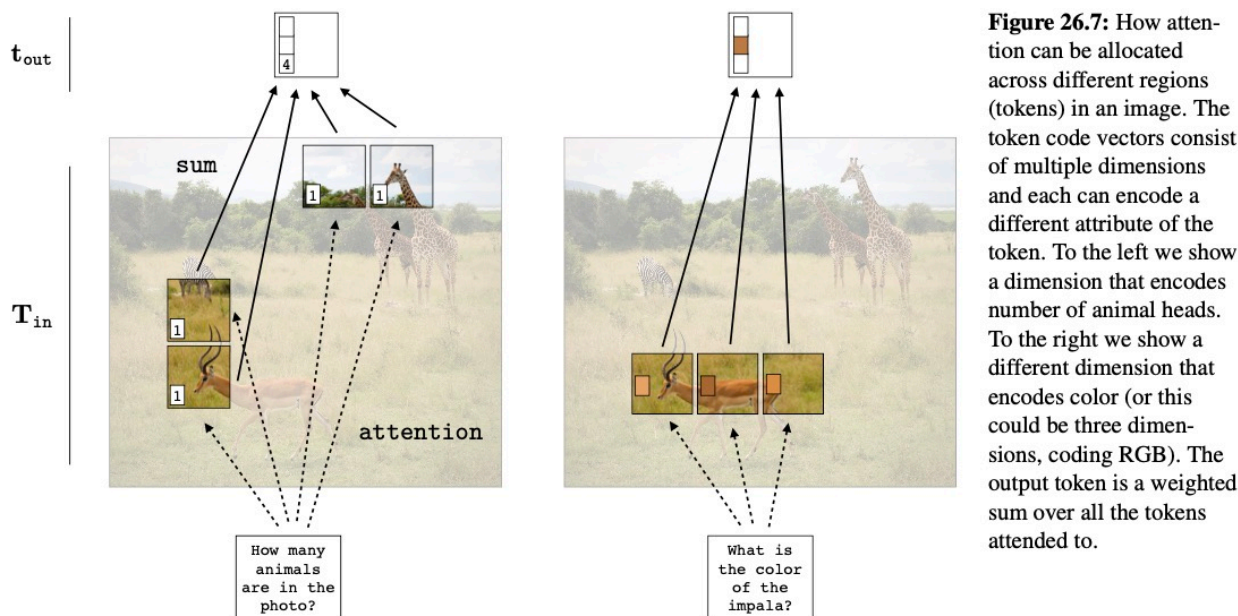
Let's view our Julia callable struct implementation side-by-side with a bare-bones implementation in PyTorch.

TODO: Side-by-side comparison.

```
1 md"""
2 ##### Our modules so far build up the Transformer
3
4 - **AttentionHead**: Implements a single attention head, creating query, key, and
  value projections, computing the attention scores, and applying a softmax.
5 - **MultiHeadedAttention**: Combines multiple `AttentionHead`s, concatenates their
  outputs, and applies a final linear transformation.
6 - **FeedForwardNetwork**: A simple feed-forward network with two linear layers and a
  ReLU activation in between.
7 - **AttentionResidual**: Combines multi-head attention and feed-forward network
  layers with residual connections.
8 - **Transformer**: Stacks multiple `AttentionResidual` layers to form the complete
  Transformer encoder.
9
10 ---
11
12 We've set up a basic structure of a Transformer using callable structs, parametric
  types, and matrix operations.
13
14 Let's view our Julia callable struct implementation side-by-side with a bare-bones
  implementation in PyTorch.
15
16 **TODO:** Side-by-side comparison.
17 """
```

We want to make a Vision Transformer. This requires some additional layers for image processing: patch embedding and positional encoding.

TODO: Implement PatchEmbed in Julia and make a visual example of applying it to some image like this:



```
1 md"""
2 We want to make a Vision Transformer. This requires some additional layers for image
3 processing: patch embedding and positional encoding.
4
5 **TODO:** Implement 'PatchEmbed' in Julia and make a visual example of applying it
6 to some image like this:
7
8 ![Patch Image]
9 (https://github.com/qsimeon/julia\_class\_project/blob/main/patch\_embed.jpg?raw=true)
10 """
```

It turns out the patch embedding is can be implemented by applying a strided convolution. However, we will take the more direct and visualizable approach of chopping up an image into patches and linearly projecting the vector that is the flattened patch to the desired dimensionality.

Remember Transformers operate on tokens i.e. transformations of tokens. What we are doing here is essentially *tokenizing* our image data.

```
1 md"""
2 It turns out the patch embedding is can be implemented by applying a strided
3 convolution. However, we will take the more direct and visualizable approach of
4 chopping up an image into patches and linearly projecting the vector that is the
5 flattened patch to the desired dimensionality.
6
7 Remember Transformers operate on tokens i.e. transformations of tokens. What we are
8 doing here is essentially *tokenizing* our image data.
9 """
```

```

1  ### 5. PatchEmbed struct
2  begin
3      struct PatchEmbed{T<:Real}
4          img_size::Int
5          patch_size::Int
6          nin::Int
7          nout::Int
8          num_patches::Int
9          W::Matrix{T} # Linear projection weights
10
11      function PatchEmbed{T}(img_size::Int, patch_size::Int, nin::Int, nout::Int)
12  where T<:Real
13          @assert img_size % patch_size == 0 "img_size must be divisible by
14  patch_size"
15
16          num_patches = (img_size ÷ patch_size)^2
17          W = randn(T, nout, patch_size^2 * nin) # Linear projection matrix for
18  each patch
19          return new{T}(img_size, patch_size, nin, nout, num_patches, W)
20      end
21  end
22  end

```



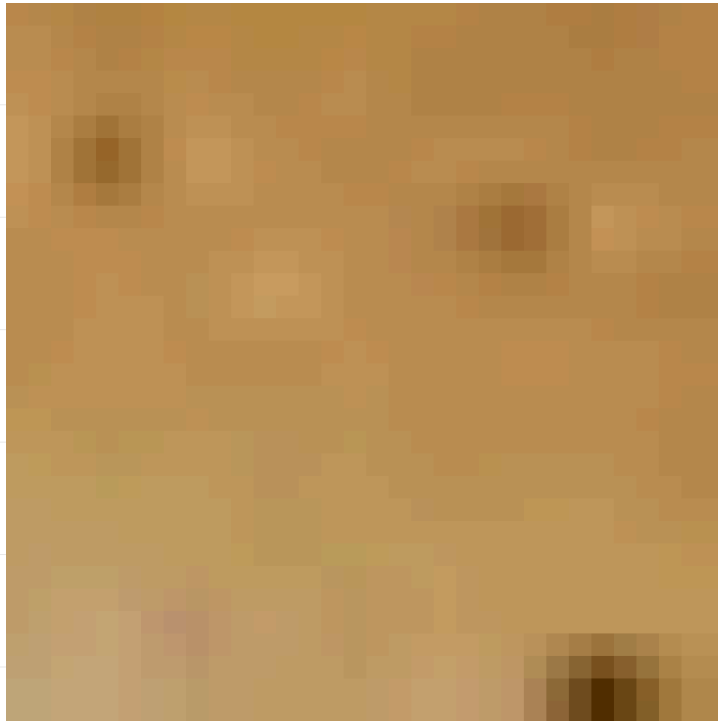
```

1  # Load and preprocess the image
2  begin
3      image_url =
4      "https://github.com/qsimeon/julia_class_project/blob/e698587c2c2b7455404e6126c06f4ec0
5      4c463032/reduced_phil.png?raw=true"
6      image_file = download(image_url)
7      image = load(image_file)
8
9      # Resize the image to a square (e.g., 256x256)
10     img_size = 256
11     image_square = imresize(image, (img_size, img_size))
12 end

```

extract_patches (generic function with 1 method)

```
1 # Define function to extract patches
2 function extract_patches(image, patch_size)
3     patches = []
4     for i in 1:patch_size:size(image, 1)
5         for j in 1:patch_size:size(image, 2)
6             push!(patches, view(image, i:i+patch_size-1, j:j+patch_size-1))
7         end
8     end
9     return patches
10 end
```



```
1 begin
2     # Parameters
3     patch_size = 32 # Size of each patch (32x32)
4
5     # Extract patches
6     patches = extract_patches(image_square, patch_size)
7
8     # Visualize patches in a grid
9     n_patches = length(patches)
10    grid_dim = Int(sqrt(n_patches)) # Assumes square grid for simplicity
11
12    # Create a grid plot using 'plot(grid=true)'
13    plot(layout=(grid_dim, grid_dim), title="Patches of Image", margin=5mm)
14
15    for idx in 1:n_patches
16        plot!(heatmap(patches[idx], color=:grays, axis=false, colorbar=false),
17            layout=(grid_dim, grid_dim), subplot=idx)
18    end
19
20    # Display the grid of patches
21    plot!()
22 end
```

