

```
1 using LinearAlgebra, Random, Plots, Plots.PlotMeasures, PlutoUI, Images, MLDatasets  
, Enzyme , Statistics,LaTeXStrings
```

Implementing a Vision Transformer (ViT) in Julia!

The Transformer architecture, introduced in the paper *Attention Is All You Need* by [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).

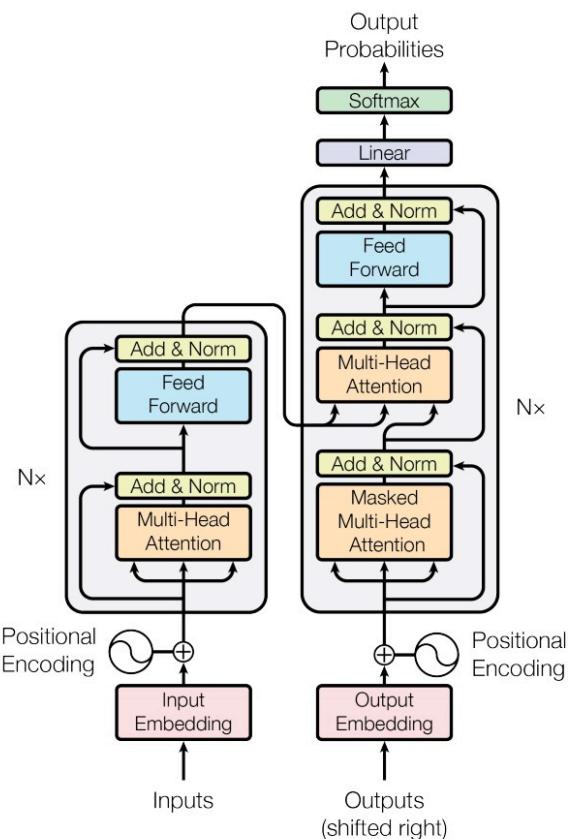
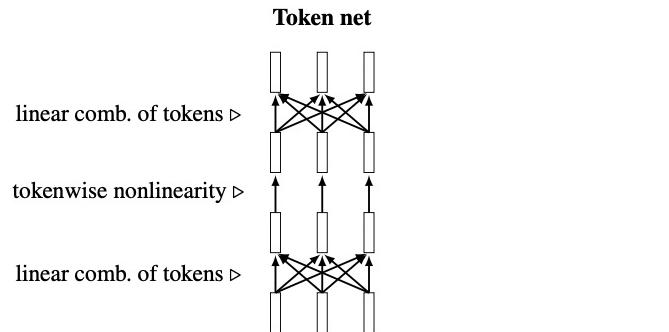


Figure 1: The Transformer - model architecture.

NOTE: We adapt/borrow a lot of material/concepts from [Torralba, A., Isola, P., & Freeman, W. T. \(2021, December 1\). Foundations of Computer Vision. MIT Press; The MIT Press, Massachusetts Institute of Technology.](#)

One thing to keep in mind throughout this notebook is that Transformers operate on **tokens**. Conceptually, the transformer architecture may be thought of as a *token net*, which is abstraction or generalization of the more familiar *neural nets*.

Token nets are just like neural nets, alternating between layers that mix nodes in linear combinations (e.g., fully connected linear layers, convolutional layers, etc.) and layers that apply a pointwise nonlinearity to each node (e.g., relus, per-token MLPs).



Until recently, the best performing models for image classification had been convolutional neural networks (CNNs) introduced in [LeCun et al. \(1998\)](#). Nowadays, transformer architectures have been shown to have similar to better performance. One such model, called Vision Transformer by [Dosovitskiy et al. \(2020\)](#) splits up images into regularly sized patches. The patches are treated as a sequence and attention weights are learned as in a standard transformer model.

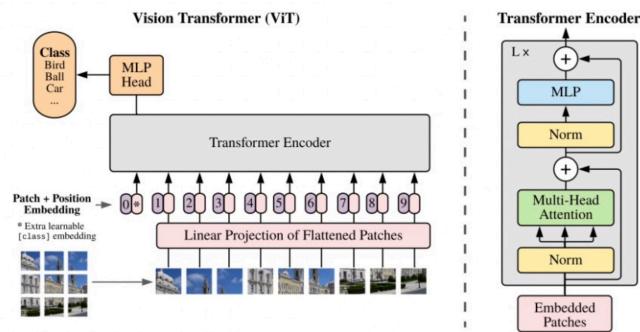


Figure 2: ViT Architecture, Figure from [Dosovitskiy et al. \[2020\]](#)

Transformers basics

In this section, we will learn how a basic transformer attention head works.

First, we must break the image up into tokens that can be processed by the transformer. Our tokens come from patches of the image.

```
visualize_patches (generic function with 1 method)
1 function visualize_patches(image_square, patch_size)
2   # Extract patches from image
3   patches = extract_patches(image_square, patch_size)
4
5   # Calculate grid dimensions (assuming a square grid)
6   n_patches = length(patches)
7   grid_dim = Int(sqrt(n_patches))
8
9   # Initialize a list to hold the individual patch plots
10  plot_list = []
11
12  # Generate the grid of patches
13  for idx in 1:n_patches
14    # Create a heatmap for each patch without axis or colorbar
15    p = heatmap(patches[idx], color=:grays, axis=false, colorbar=false)
16    push!(plot_list, p)
17  end
18
19  # Display the grid of patches
20  plot(plot_list..., layout=(grid_dim, grid_dim), title="Patches of Image",
21       size=(2000, 2000))
22 end
23
```

```
extract_patches (generic function with 1 method)
```

We will use the CIFAR dataset here

```
dataset CIFAR10:
  metadata => Dict{String, Any} with 2 entries
  split => :test
  features => 32×32×3×10000 Array{Float32, 4}
  targets => 10000-element Vector{Int64}
1 begin
2   train_dataset = CIFAR10(dir="cifar/", split=:train)
3   test_dataset = CIFAR10(dir="cifar/", split=:test)
4 end
```



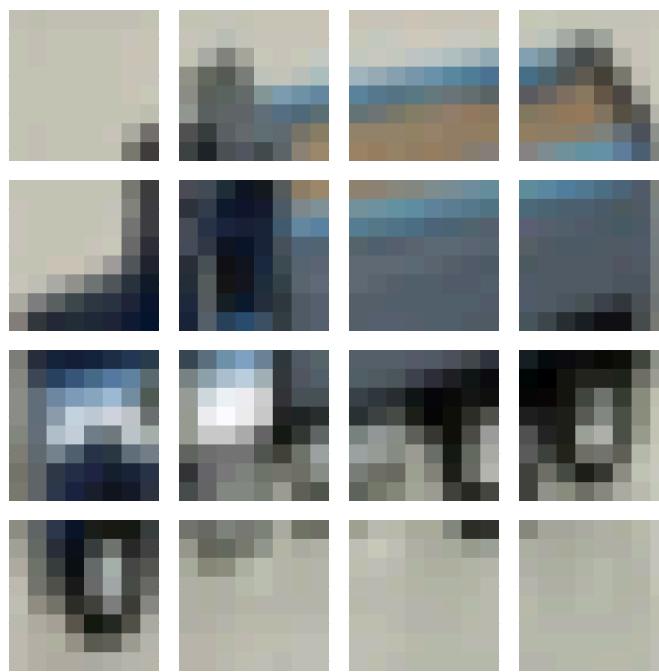
```
1 @bind cifar_img_num Slider(1:1:1000, show_value=true)
```



First, we will break our image up into patches. We can choose different patch sizes.

```
1 # patch size
2 @bind patch_size_slider Slider([2,4,8,16,32], show_value=true)
```

"truck"



These patches are the basis for our tokens (in practice using some sort of embedding discussed later)

`extract_patches_num` (generic function with 1 method)

```

16x8x8x3 Array[Float64, 4]:
[:, :, 1, 1]
 0.768627  0.760784  0.768627  0.768627 ... 0.768627  0.772549  0.780392
 0.780392  0.776471  0.776471  0.776471 ... 0.772549  0.478431  0.239216
 0.501961  0.160784  0.0823259  0.09101961 ... 0.12549  0.14902  0.172549
 0.560874  0.380392  0.188235  0.054902 ... 0.105882  0.0901961  0.223529
 0.780392  0.776471  0.772549  0.772549 ... 0.788235  0.792157  0.796078
 0.294118  0.258824  0.133333  0.0627451 ... 0.345098  0.560784  0.607843
 0.231373  0.368627  0.439216  0.490196 ... 0.184314  0.32549  0.341176
:
 0.333333  0.337255  0.337255  0.337255 ... 0.301961  0.254902  0.180392
 0.627451  0.737255  0.729412  0.67451 ... 0.521569  0.188235  0.2
 0.780392  0.768627  0.772549  0.768627 ... 0.780392  0.780392  0.788235
 0.380392  0.356863  0.329412  0.321569 ... 0.298039  0.313726  0.490196
 0.0941176  0.0313726  0.0431373  0.0627451 ... 0.0666667  0.2666667  0.564701
 0.533333  0.65098  0.654902  0.678431 ... 0.690196  0.72549  0.772549

```

[::, ::, 2, 1] =	0.772549	0.768672	0.772549	0.776471	..	0.776471	0.784314	0.788235
0.784314	0.780392	0.780392	0.780392	0.780392	..	0.780392	0.509804	0.247059
0.541176	0.247059	0.160784	0.219608	0.345098	..	0.376471	0.364706	0.364706
0.631373	0.4	0.164706	0.0580804	0.439216	..	0.160784	0.258824	0.258824
0.788235	0.776471	0.741176	0.776471	0.784314	..	0.788235	0.796078	0.796078
0.188235	0.2	0.0627451	0.0588235	0.309804	..	0.309804	0.501961	0.501961
0.478431	0.560784	0.670588	0.72549	0.196078	..	0.313726	0.329412	0.329412
..
0.321569	0.321569	0.301961	0.247059	0.111674	..	0.0745098	0.043137	0.043137
0.752941	0.784314	0.756863	0.694118	0.658824	..	0.6	0.623529	0.623529
0.764706	0.741176	0.658824	0.537255	0.784314	..	0.792157	0.796078	0.796078
0.298039	0.301961	0.309804	0.309804	0.294118	..	0.313726	0.482353	0.482353
0.0156863	0.0784314	0.0784314	0.141176	0.153333	..	0.435294	0.623529	0.623529

```
1 begin
2     cifar_img_size = 32
3     n_channels = 3
4     dim = 10
5
6     patches = extract_patches_num(train_dataset.features[:, :, :, cifar_img_num]
7         patch_size_slider)
8 end
```

The attention aspect involves queries, keys, and values.

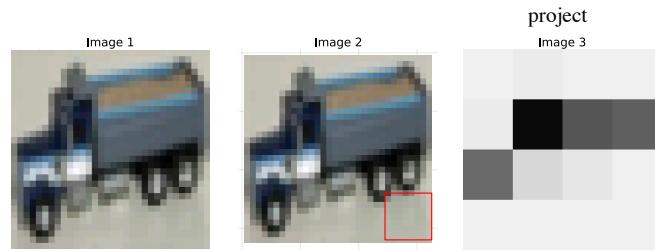
The queries are like the question you are asking, which should direct the attention. The keys tell you how much each token contributes in the query.

The attention matrix is essentially the multiplication of the query and key matrices.

The value matrix maps the attention matrix to the new token values in the next layer.

Example with prefixed keys and queries set to mean color of patch (in reality they are learned during training and we don't know exactly what they are)

16



The attention matrix comes from the product of the key and query matrices, multiplied by the tokens.

Finally, the next layer of tokens is calculated by multiplying the attention matrix and the value matrix.

Building a transformer

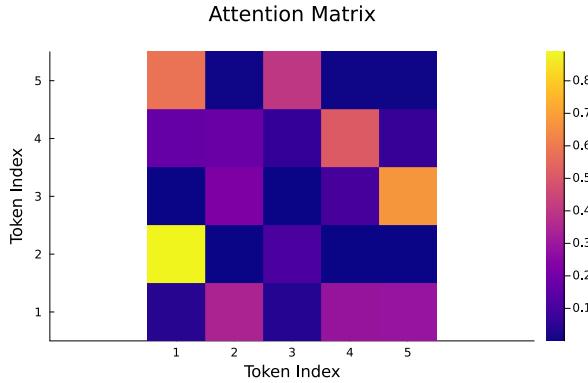
Let's start by defining key components of a **Transformer** model using Julia structs and parametric types, similar to the structure we implemented in *Homework 3*.

We will implement the `AttentionHead`, `MultiHeadedAttention`, and `FeedForwardNetwork` modules as Julia structs. This will set up the parts which get combined together in the `Transformer` model.

```
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::Array{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
39 end
40
```

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



```

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
16         :plasma
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)



To compute the query, key, and value for a set of input tokens, \mathbf{T}_{in} , we apply the same linear transformations to each token in the set, resulting in matrices $\mathbf{Q}_{in}, \mathbf{K}_{in} \in \mathbb{R}^{N \times m}$ and $\mathbf{V}_{in} \in \mathbb{R}^{N \times d}$, where each row is the query/key/value for each token:

$$\mathbf{Q}_{in} = \begin{bmatrix} \mathbf{q}_1^\top \\ \vdots \\ \mathbf{q}_N^\top \end{bmatrix} = \begin{bmatrix} (\mathbf{W}_q \mathbf{t}_1)^\top \\ \vdots \\ (\mathbf{W}_q \mathbf{t}_N)^\top \end{bmatrix} = \mathbf{T}_{in} \mathbf{W}_q^\top \quad \leftarrow \text{query matrix}$$

$$\mathbf{K}_{in} = \begin{bmatrix} \mathbf{k}_1^\top \\ \vdots \\ \mathbf{k}_N^\top \end{bmatrix} = \begin{bmatrix} (\mathbf{W}_k \mathbf{t}_1)^\top \\ \vdots \\ (\mathbf{W}_k \mathbf{t}_N)^\top \end{bmatrix} = \mathbf{T}_{in} \mathbf{W}_k^\top \quad \leftarrow \text{key matrix}$$

$$\mathbf{V}_{in} = \begin{bmatrix} \mathbf{v}_1^\top \\ \vdots \\ \mathbf{v}_N^\top \end{bmatrix} = \begin{bmatrix} (\mathbf{W}_v \mathbf{t}_1)^\top \\ \vdots \\ (\mathbf{W}_v \mathbf{t}_N)^\top \end{bmatrix} = \mathbf{T}_{in} \mathbf{W}_v^\top \quad \leftarrow \text{value matrix}$$

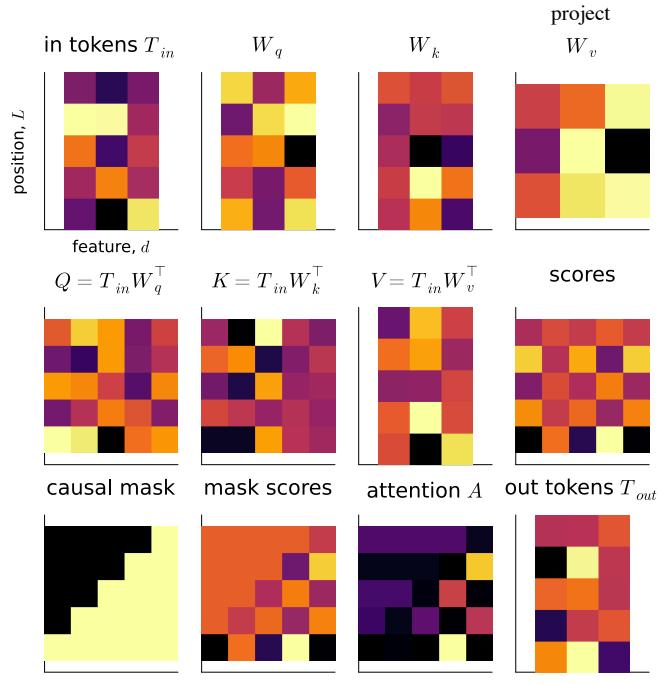
Note that the query and key vectors must have the same dimensionality, m , because we take a dot product between them. Conversely, the value vectors must match the dimensionality of the token code vectors, d , because these are summed up to produce the new token code vectors.

Finally, we have the attention equation:

$$\mathbf{A} = \text{softmax} \left(\frac{\mathbf{Q}_{in} \mathbf{K}_{in}^\top}{\sqrt{m}} \right) \quad \leftarrow \text{attention matrix}$$

$$\mathbf{T}_{out} = \mathbf{A} \mathbf{V}_{in}$$

where the softmax is taken within each row (i.e., over the vector of matches for each separate query vector).



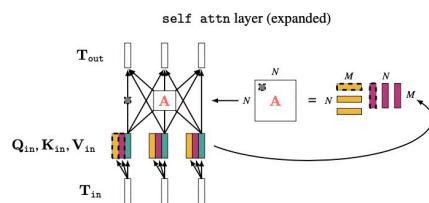
```

1 # Showing causal attention with mask
2 let
3   dim, attn_dim = 3, 8
4   args = (aspect_ratio=1, colorbar=false, grid=false, yticks=false, xticks=false,
5   size=(650,650))
6
7   X = randn(Float64, n_tokens, dim) # example 3-D input of n_tokens
8   W_Q = randn(n_tokens, dim)
9   W_K = randn(n_tokens, dim)
10  W_V = randn(dim, dim)
11
12  p1 = heatmap(X, title=string("in tokens ", L"T_{in}")); args...
13  ylabel=string("position, ", L'L"), xlabel=string("feature, ", L"d"))
14  p2 = heatmap(W_Q, title=L"W_q"; args...)
15  p3 = heatmap(W_K, title=L"W_k"; args...)
16  p4 = heatmap(W_V, title=L"W_v"; args...)
17
18  Q = X * transpose(W_Q)
19  K = X * transpose(W_K)
20  V = X * transpose(W_V)
21
22  p5 = heatmap(Q, title=L"Q = T_{in} W_q^\\intercal"; args...)
23  p6 = heatmap(K, title=L"K = T_{in} W_k^\\intercal"; args...)
24  p7 = heatmap(V, title=L"V = T_{in} W_v^\\intercal"; args...)
25
26  scores = Q * transpose(K) / sqrt(attn_dim) # n_tokens x n_tokens
27  @assert size(scores) == (n_tokens, n_tokens)
28
29  mask = UpperTriangular(ones(n_tokens, n_tokens))
30
31  p9 = heatmap(mask, title="causal mask"; args...)
32
33  masked_scores = scores .* mask .+ (1 .- mask)
34
35  p10 = heatmap(masked_scores, title="mask scores"; args...)
36
37  alpha = softmax(masked_scores, dims=ndims(masked_scores))
38
39  p11 = heatmap(alpha, title=string("attention ", L"A")); args...)
40
41  attn_output = alpha * V
42
43  p12 = heatmap(attn_output, title=string("out tokens ", L"T_{out}"); args...)
44
45  plot!([p1,p2,p3,p4,p5,p6,p7,p8,p9,p10,p11,p12]..., layout=(3,4))
46 end

```

In expanded detail, here are the full mechanics of a self-attention layer, which is the kind of attention layer used in transformers.

Figure 26.10: Self-attention layer expanded.
The nodes with the dashed outline correspond to each other; they represent one query being matched against one key to result in a scalar similarity value, in the gray box, which acts as a weight in the weighted sum computed by A .

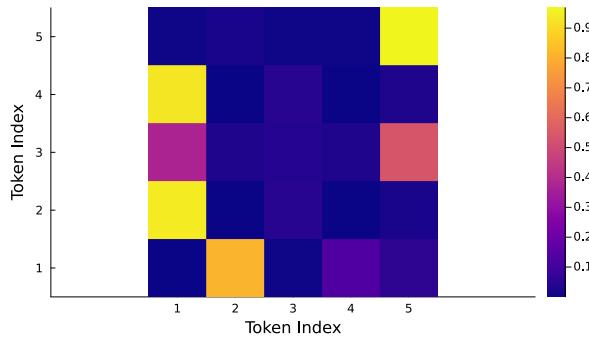


```

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
7     T<:Real
8     # Each head outputs dim-dimensional tokens
9     heads = [AttentionHead{T}(dim, n_hidden) for _ in 1:num_heads]
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::Array{T}, attn_mask::Union{Nothing,
15    Matrix{T}}=nothing) where {T<:Real}
16    outputs, alphas = [], []
17    for head in mha.heads
18      out, alpha = head(X, attn_mask) # Shapes: (N, dim), (N, N)
19      push!(outputs, out)
20      push!(alphas, alpha)
21    end
22    # Concatenate along hidden dimension
23    concatenated = cat(outputs...; dims=2) # Shape: (N, num_heads*dim)
24    attn_output = concatenated * transpose(mha.W_msa) # Shape: (N, dim)
25    attn_alphas = cat(alphas...; dims=3) # Shape: (N, N, num_heads)
26    attn_alphas = permutedims(attn_alphas, (3, 1, 2)) # Shape: (num_heads, N, N)
27  end
28 end

```

Attention Matrix (last head)



```

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 (last head)",
29     c=:plasma, # Choose a colormap, e.g., :viridis or
30     :plasma
31     cLabel="Weight", # Label for the color bar
32     colorbar=true, # Show the color bar
33     grid=false, # Turn off the grid
34     # framestyle=:none, # Removes the axis lines
35     xticks=1:n_tokens, # Ensures ticks are at each integer index
36     yticks=1:n_tokens,
37   )
38 end

```

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

Figure 26.12: Example of self-attention maps where each token is an image patch and the query and key vectors are both set to the mean color of the patch, normalized to be a unit vector.



```

1 ### 3. Feed-Forward Network (FFN)
2 struct Linear{T<:Real}
3   W::Matrix{T} # Weight matrix
4   b::Vector{T} # Bias vector
5
6   # Constructor: Initialize weights and biases
7   function Linear{!}(in_features::Int, out_features::Int) where T<:Real
8     W = randn(T, out_features, in_features) / sqrt(in_features) # Xavier
9     initialization
10    b = zeros(T, out_features)
11    return new{T}(W, b)
12  end
13
14  # Apply the linear transformation
15  function (linear::Linear{T})(X::Array{T}) where T<:Real
16    @assert size(X, 2) == size(linear.W, 2) "Input dimension must match the
17    weight matrix"
18    return X * transpose(linear.W) .+ linear.b
19  end
20 end

```

```

1 # Test 'Linear' implementation
2 let
3   in_features = 64 # Input dimension
4   out_features = 10 # Output dimension
5
6   # Initialize the Linear module
7   linear_layer = Linear{Float64}(in_features, out_features)
8
9   # Example input matrix (T, in_features), where T is the number of tokens
10  n_tokens = 20 # Number of tokens
11  X = randn(Float64, n_tokens, in_features)
12
13  # Apply the linear transformation
14  output = linear_layer(X)
15
16  # Verify dimensions
17  println("Input shape: ", size(X)) # Should be (n_tokens, in_features)
18  println("Output shape: ", size(output)) # Should be (n_tokens, out_features)
19 end

```

Input shape: (20, 64)
Output shape: (20, 10)



```

1 ### 4. Feed-Forward Network (FFN)
2 struct FeedForwardNetwork{T<:Real}
3   layer1::Linear{T} # First linear transformation
4   layer2::Linear{T} # Second linear transformation
5
6   function FeedForwardNetwork{T}(dim::Int, n_hidden::Int) where T<:Real
7     # Initialize the two linear layers
8     layer1 = Linear{!}(dim, n_hidden) # Shape: (n_hidden, dim)
9     layer2 = Linear{T}(n_hidden, dim) # Shape: (dim, n_hidden)
10    return new{T}(layer1, layer2)
11  end
12
13  function (ffn::FeedForwardNetwork{T})(X::Array{T}) where {T<:Real}
14    # Apply the first linear transformation
15    X = ffn.layer1(X) # Shape: (N, n_hidden)
16    X = max.(0, X) # ReLU activation
17    # Apply the second linear transformation
18    return ffn.layer2(X) # Shape: (N, dim)
19  end
20 end

```

```

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 input with n_tokens tokens of dim-
6   dimensional vectors
6   ffn_output = ffn(X)
7
8   # Verify output shape
9   println("Input shape: ", size(X)) # Should be (n_tokens, dim)
10  println("Output shape: ", size(ffn_output)) # Should be (n_tokens, dim)
11 end
12

```

Input shape: (5, 3)
Output shape: (5, 3)



Recap so far

1. **AttentionHead Implementation:**
 - o Projects the input token matrix X (shape: (N, dim)) to query, key, and value matrices.
 - o Computes scaled dot-product attention, applies an optional mask, and then applies softmax to get attention weights alpha with shape (N, N) .
 - o Returns the attention output attn_output (shape: (N, dim)) and attention weights alpha.
2. **MultiHeadedAttention Implementation:**
 - o Creates multiple AttentionHead instances and collects their outputs.
 - o 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:**
 - o A two-layer feed-forward network with an intermediate hidden layer of size n_{hidden} .
 - o 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 `W2` and `b2`.
- Returns the output with shape `(N, dim)`, maintaining the same dimension as the input tokens.

```

1  ### 5. 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::Array{T}, attn_mask::Union{Nothing,
16    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
28

```

```

1  ### 6. Transformer
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::Array{T}; attn_mask::Union{Nothing,
15    Matrix{T}}=nothing, return_attn::Bool=false) 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      if return_attn
20        push!(collected_alphas, alphas) # Collect attention weights
21      end
22    end
23    # Return the final output and collected attention weights from all layers
24    (if required)
25      if return_attn
26        return X, collected_alphas
27      else
28        return X, nothing
29      end
30    end
31  end
32 end
33

```

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 # 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'
6   dimensions
7   output, alphas = residual_block(X)
8   println("AttentionResidual output shape: ", size(output))
9   println("Attention weights shape (from one layer): ", size(alphas))
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'
6   dimensions
7   output, collected_alphas = transformer(X; return_attn=true)
8   println("Transformer output shape: ", size(output))
9   println("Collected attention weights shape ($num_layers layers): ",
10   [size(collected_alphas[i]) for i in 1:length(collected_alphas)])
end

```

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

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.

Using callable structs, parametric types, and matrix operations, we set up the basic components and combined them to create a Transformer module.

Let's view our Julia callable struct implementations of the `AttentionHead` and `Transformer` modules side-by-side with a bare-bones implementation in PyTorch.

Show side-by-side comparison:

AttentionHead implemented in Python (left) and Julia (right)

```
class AttentionHead(nn.Module):
    def __init__(self, dim_int, n_hidden: int):
        super().__init__()
        self.W_q = nn.Linear(dim_int, n_hidden) # Q weight matrix
        self.W_k = nn.Linear(dim_int, n_hidden) # K weight matrix
        self.W_v = nn.Linear(dim_int, n_hidden) # V weight matrix
        self.n_hidden = n_hidden
        n_hidden

    def forward(self, x: torch.Tensor, attn_mask: Optional[torch.Tensor] = None, attn_head: torch.Tensor, attn_kv: torch.Tensor, attn_v: torch.Tensor):
        # x: Input tensor of shape (B x T x C)
        # attn_kv: optional mask tensor of shape (B x T x T)
        # attn_head: optional mask tensor of shape (B x T x T)

        # ===== Answer: START =====
        Q = self.W_q(x) # shape: (B x T x n_hidden)
        K = self.W_k(x) # shape: (B x T x n_hidden)
        V = self.W_v(x) # shape: (B x T x n_hidden)

        # Compute attention scores (softmax)
        K_T = K.permute(0, 2, 1, 3) # b -> b x t
        scores = F.softmax(K @ V / self.n_hidden ** 0.5) # (B x T x T)

        # Apply attention mask if provided
        if attn_mask is not None:
            scores = scores.masked_fill(attn_mask == 0, float("-inf"))

        # Compute attention weights (softmax)
        alpha = F.softmax(scores, dim=-1) # (B x T x T)

        # Compute the attention output as a weighted sum of values
        attn_output = torch.matmul(alpha, V) # (B x T x n_hidden)
        # ===== Answer: END =====
        return attn_output, alpha
```

```
## 1. Attention Head
struct AttentionHead{T<:Real}
    W_q::Matrix{T} # Shape: (n_hidden, dim)
    W_k::Matrix{T} # Shape: (n_hidden, dim)
    W_v::Matrix{T} # Shape: (n_hidden, dim)
    n_hidden::Int # dimensionality of key and query vectors
end

function AttentionHead{T}(dim:Int, n_hidden:Int) where {T<:Real}
    Matrix{T}(randn(T, n_hidden, dim), randn(T, n_hidden, dim), randn(T, n_hidden, dim))
end

function (head::AttentionHead{T})(x::Matrix{T}, attn_mask::Union{Nothing,
    Matrix{T}}) where {T<:Real}
    # x is expected to be an input token matrix with shape (B, dim)
    # Project input tokens to query, key, and value representations
    # x: (B x T x C) -> (B x T x n_hidden)
    # K = x * transpose(head.W_k) # Shape: (B, n_hidden)
    # V = x * transpose(head.W_v) # Shape: (B, n_hidden)

    # Compute scaled dot-product attention
    scores = 0.0 * x * transpose(x) / sqrt(head.n_hidden) # Shape: (B, B)
    # Apply attention mask if provided
    if attn_mask != nothing
        scores = scores .* attn_mask .+ (I - attn_mask) * -Inf
    end

    # Apply softmax along the last dimension
    alpha = softmax(scores, dims=last(scores)) # Shape: (B, B)

    # Compute attention output as weighted sum of values
    attn_output = alpha * V # Shape: (B, dim)

    # attn_output is the (B, dim) output token matrix
    # alpha is the (B, B) attention matrix
    return attn_output, alpha
end
```

We don't handle batched inputs in the Julia implementation so assume `batch_size` (B) is 1.

Transformer implemented in Python (left) and Julia (right)

```
class Transformer(nn.Module):
    def __init__(self, dim_int, n_hidden, dim_attn, n_heads, num_heads: int,
                 mlp_dim: int, num_layers: int):
        super().__init__()

        # ===== Answer: START =====
        self.dim_int = dim_int
        self.n_hidden = n_hidden
        self.dim_attn = dim_attn
        self.n_heads = n_heads
        self.mlp_dim = mlp_dim
        self.num_heads = num_heads
        self.num_layers = num_layers

        self.layers = nn.ModuleList([AttentionResidualBlock(dim_int, dim_attn, n_heads, mlp_dim, num_heads) for _ in range(num_layers)])
        # ===== Answer: END =====
        self.dim_int = dim_int
        self.n_hidden = n_hidden
        self.dim_attn = dim_attn
        self.n_heads = n_heads
        self.mlp_dim = mlp_dim
        self.num_heads = num_heads
        self.num_layers = num_layers

    def forward(self, x: torch.Tensor, attn_mask: torch.Tensor, attn_kv: torch.Tensor, attn_v: torch.Tensor):
        # x: Input tensor of shape (B x T x C)
        # attn_kv: optional mask tensor of shape (B x T x T)
        # attn_head: optional mask tensor of shape (B x T x T)

        # ===== Answer: START =====
        x = self.layers[0](x, attn_kv, attn_v, attn_head, attn_mask)
        for layer in self.layers[1:-1]:
            x = layer(x, attn_kv, attn_v, attn_head, attn_mask)
        x = self.layers[-1](x, attn_kv, attn_v, attn_head, attn_mask)
        # ===== Answer: END =====
        return x
```

```
## 2. Transformer
struct Transformer{T<:Real}
    layers::Vector{AttentionResidual{T}} # Sequence of AttentionResidual blocks
end

# Constructor: initializes a sequence of attention residual blocks
# AttentionResidual{T}(dim:Int, attn_dim:Int, mlp_dim:Int, num_heads:Int,
# num_layers:Int) where {T<:Real}
#     layers = [AttentionResidual{T}(dim, attn_dim, mlp_dim, num_heads) for _ in 1:num_layers]
#     return new{T}(layers)
end

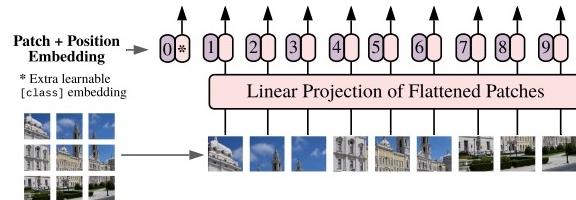
# Apply the Transformer model to input X
function (transformer::Transformer{T})(x::Matrix{T}, attn_mask::Union{Nothing,
    Matrix{T}}) where {T<:Real}
    collected_alphas = [] # To store attention weights from each layer
    for layer in transformer.layers
        x, alphas = layer(x, attn_mask) # Apply each residual block
        push!(collected_alphas, alphas) # Collect attention weights
    end
    # Return the final output and collected attention weights from all layers
    return x, collected_alphas
end

# ===== Answer: START =====
# Collect attention weights from each layer
# x: (B x T x C) -> (B x T x n_heads)
# x = x * transpose(head.W_k) # Shape: (B x T x n_heads)
# collected_attn = []
# for layer in layers:
#     x, alpha = layer(x, attn_kv, attn_v, attn_head, attn_mask)
#     collected_attn.append(alpha)
# if return_alpha is True, return the alphas
# else, return the output
# ===== Answer: END =====
# return output, collected_attn if return_alpha else None
end
```

We don't handle batched inputs in the Julia implementation so assume `batch_size` (B) is 1.

But recall that what we want is to make is a **Vision Transformer**. This requires some additional layers for tokenizing an image. These are

- (1) patch embedding; and
- (2) positional encoding



We will explore these in detail next.

Patch Embedding

It turns out the patch embedding 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 the first step of *tokenizing* our image data.

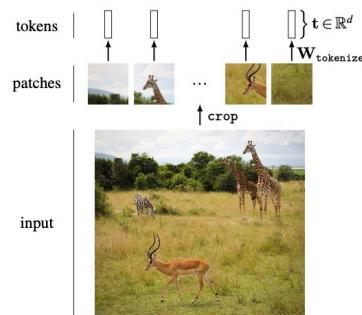
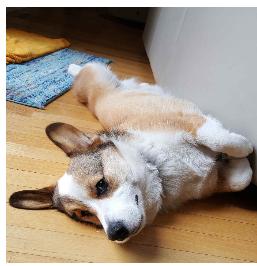


Figure 26.2: Tokenization: converting an image to a set of vectors. $\mathbf{W}_{\text{tokenize}}$ is a learnable linear projection from the dimensionality of the vectorized crops to d dimensions. This is just one of many possible ways to tokenize an image.



```

1 # Load and preprocess an example image
2 begin
3   # Download image of Philip
4   url = "https://user-images.githubusercontent.com/6933510/107239146-dcc3fd00-6a28-
5   11eb-8c7b-41aaef6618935.png" # high-res
6   philip_filename = download(url)
7   philip = load(philip_filename)
8
9   # Resize the image to a square (e.g., 256x256)
10  img_size = 256
11  image_square = imresize(philip, (img_size, img_size))
12
13 end

```

[16, 32, 64, 128, 256]

```

1 begin
2   # Calculate divisors of the image size
3   function divisors_of_half_image_size(img_size)
4     divisors = []
5     for i in div(img_size, 16):div(img_size, 1)
6       if img_size % i == 0
7         push!(divisors, i)
8     end
9   end
10  return divisors
11
12
13 # Get the divisors of the image size
14 patch_size_options = divisors_of_half_image_size(img_size)
15
16

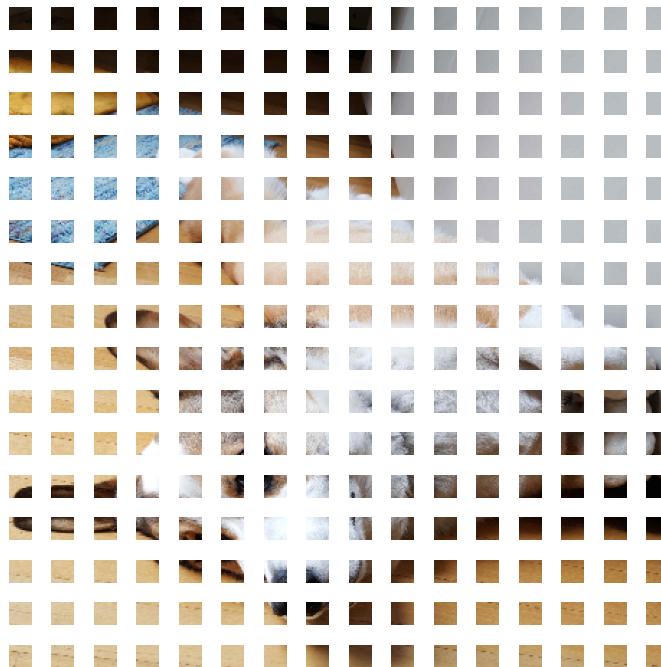
```

16

```

1 # Create the slider with the patch size options
2 @bind patch_size Slider(patch_size_options, show_value=true,
3 default=patch_size_options[1])

```



```

1 # Visualize patches for the chosen 'patch_size'
2 visualize_patches(image_square, patch_size)

```

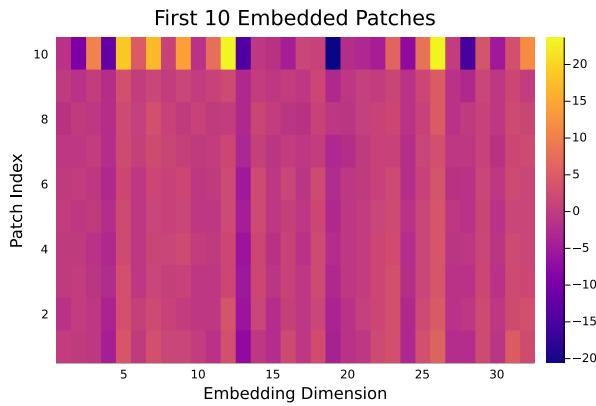
A patch embedding layer is one that takes each of the image patches, like those displayed above, and then projects that into a vector. One approach is to simply flatten each patch and use a linear projection (using a matrix multiplication) to convert this into a vector. Since we are working on RGB images (3 channels), we define a linear projection for each channel independently and then combine them.

```

1  ### 7. PatchEmbedLinear
2  struct PatchEmbedLinear{T<:Real}
3      img_size::Int
4      patch_size::Int
5      nin::Int # Number of input channels (e.g., RGB → nin = 3)
6      nout::Int # Desired output dimensionality for each patch
7      num_patches::Int
8      W::Vector{Matrix{T}} # Linear projection weights, one for each channel
9
10     function PatchEmbedLinear{T}(img_size::Int, patch_size::Int, nin::Int,
11         nout::Int) where T<:Real
12         @assert img_size % patch_size == 0 "img_size must be divisible by patch_size"
13         num_patches = (img_size ÷ patch_size)^2
14         # Create a distinct weight matrix for each channel
15         W = [randn(T, nout, patch_size^2) for _ in 1:nin]
16         return new{T}(img_size, patch_size, nin, nout, num_patches, W)
17     end
18
19     function (embed::PatchEmbedLinear{T})(image::Matrix{<:RGB}) where T<:Real
20         # Ensure image size matches expected dimensions
21         img_size = size(image, 1)
22         @assert img_size == embed.img_size "Image size does not match module
23         configuration"
24
25         # Split the RGB image into three separate channel matrices
26         channels = channelview(image) # Shape: (3, 256, 256)
27         @assert size(channels, 1) == embed.nin "Number of image channels does not
28         match nin"
29
30         # Extract patches and project for each channel
31         projected_channels = []
32         for c in 1:embed.nin
33             # Extract the channel matrix (2D slice)
34             channel_matrix = Matrix{T}(channels[c, :, :]) # Shape: (256, 256)
35             # Extract patches
36             patches = extract_patches(channel_matrix, embed.patch_size)
37             # Flatten patches and project
38             patch_matrix = hcat([vec(patch) for patch in patches]...) # Shape:
39             (num_patches, patch_size^2)
40             projected = patch_matrix * transpose(embed.W[c]) # Shape: (num_patches,
41             nout)
42             push!(projected_channels, projected)
43         end
44
45     end
46 end

```

 32
1 @bind nout Slider([16:16:96...], show_value=true, default=32)



```

1 # Test 'PatchEmbedLinear' implementation for RGB images
2 let
3     img_size = size(image_square, 1) # Image size (assumes square image: img_size x
4         img_size)
5     nin = size(channelview(image_square), 1) # Number of input channels (e.g., RGB)
6     # nout = 64 # Desired output dimensionality for each patch embedding
7
8     # Create a 'PatchEmbedLinear' instance
9     patch_embed = PatchEmbedLinear{Float64}(img_size, patch_size, nin, nout)
10
11    # Use the color image of Philip the dog (loaded as 'image_square')
12    embedded_patches = patch_embed(image_square)
13
14    # Verify output dimensions
15    println("Input image shape: ", size(image_square))
16    println("Number of patches: ", patch_embed.num_patches)
17    println("Embedded patches shape: ", size(embedded_patches))
18
19    # Visualize the embedded patches (first few)
20    npatch = min(10, size(embedded_patches, 1))
21    heatmap(embedded_patches[1:patch, :], :
22        title="First 10 Embedded Patches",
23        xlabel="Embedding Dimension",
24        ylabel="Patch index",
25        c=:plasma,
26        clabel="Value")
27 end
27

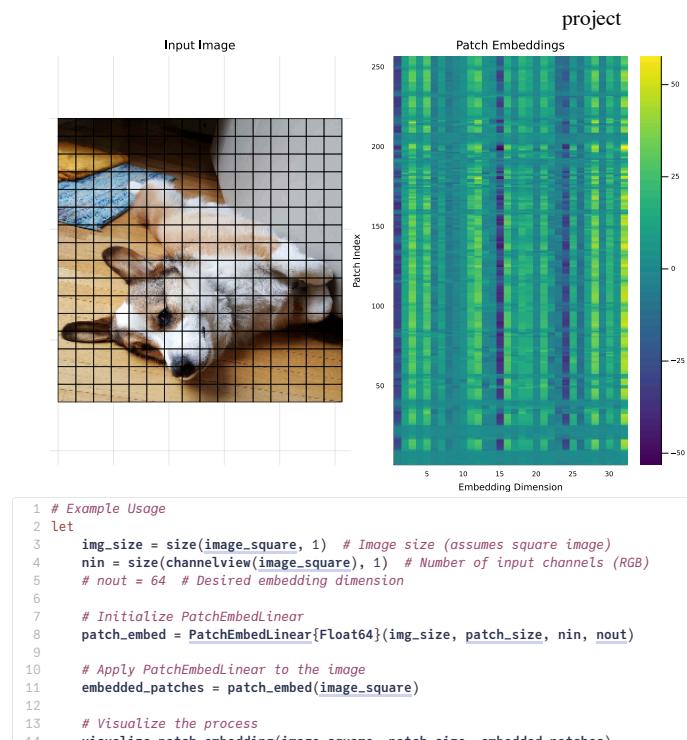
```

Input image shape: (256, 256)
Number of patches: 256
Embedded patches shape: (256, 32)

```

visualize_patch_embedding (generic function with 1 method)
1 function visualize_patch_embedding(image::Matrix{<:RGB}, patch_size::Int,
2     embedded_patches::Matrix{Float64})
3     # Extract patches for visualization
4     patches = extract_patches(image, patch_size)
5
6     # Calculate grid dimensions for patches
7     n_patches = length(patches)
8     grid_dim = div(size(image, 1), patch_size) # Grid dimensions (e.g., 16x16 for
9     256x256 with patch_size=16)
10
11    # Start with the original image
12    p1 = plot(image, title="Input Image", size=(800, 800), color=:grays, axis=false)
13
14    # Overlay transparent patches
15    for i in 0:grid_dim-1
16        for j in 0:grid_dim-1
17            # Draw rectangles for patches
18            rectangle = Shape([j*patch_size, (j+1)*patch_size, (j+1)*patch_size,
19                j*patch_size],
20                [i*patch_size, i*patch_size, (i+1)*patch_size,
21                (i+1)*patch_size])
22            plot!(rectangle, lw=1.5, linealpha=0.7, fillalpha=0.0, color=:red,
23            legend=false)
24        end
25    end
26
27    # Visualize embeddings as a heatmap
28    p2 = heatmap(embedded_patches, title="Patch Embeddings", xlabel="Embedding
29        Dimension",
30        ylabel="Patch Index", color=:viridis, size=(800, 800))
31
32    # Combine the two plots
33    plot(p1, p2, layout=(1, 2), size=(1200, 800))
34 end
34

```



Positional Encoding

One reason why CNNs worked so well for image recognition is because they have an inductive bias for local structure. In a Transformer, every token can attend to every other token in the sequence. Because self-attention operation is permutation invariant, it is important to use proper positional encoding to provide order information to the model. The positional encoding $\mathbf{P} \in \mathbb{R}^{L \times d}$ has the same dimension as the input embedding, so it can be added on the input directly.

The vanilla Transformer considered two types of encodings:

- (1) *Sinusoidal positional encoding*: Each dimension of the positional encoding corresponds to a sinusoid of different wavelengths in different dimensions.

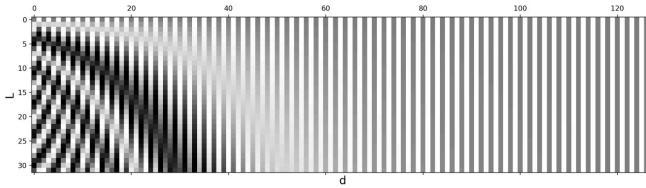
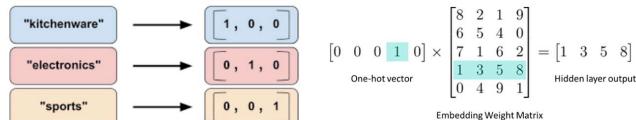


Fig. 3. Sinusoidal positional encoding with $L = 32$ and $d = 128$. The value is between -1 (black) and 1 (white) and the value 0 is in gray.

- (2) *Learned positional encoding*: As its name suggests, assigns each element in a sequence with a learned column vector which encodes its absolute position.

We will implement the latter (2) by implementing an Embedding module since it is straightforward and because embedding layers are extremely useful and ubiquitous in deep learning code.



1) Convert indices of 3 training examples to one-hot encoding

$$\begin{bmatrix} 2 \\ 3 \\ 1 \end{bmatrix} \rightarrow \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \end{bmatrix}$$

2) Multiply one-hot encoded inputs with weight matrix

$$\begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \end{bmatrix} \times \begin{bmatrix} 0.3374 & -0.1778 & -0.3035 & -0.5880 & 1.5810 \\ 1.3010 & 1.2753 & -0.2010 & -0.1606 & -0.4015 \\ 0.6957 & -1.8061 & -1.1589 & 0.3255 & -0.6315 \\ -2.8400 & -0.7849 & -1.4096 & -0.4076 & 0.7953 \end{bmatrix}$$

$$= \begin{bmatrix} 0.6957 & -1.8061 & -1.1589 & 0.3255 & -0.6315 \\ -2.8400 & -0.7849 & -1.4096 & -0.4076 & 0.7953 \\ 1.3010 & 1.2753 & -0.2010 & -0.1606 & -0.4015 \end{bmatrix}$$

```

1  ### 8. Embedding
2  struct Embedding{T<:Real}
3      emb::Matrix{T} # Shape: (num_embeddings, embedding_dim)
4
5      function Embedding{T}(num_embeddings::Int, embedding_dim::Int) where T<:Real
6          emb = randn(T, num_embeddings, embedding_dim) # Randomly initialize the
7          embedding matrix
8          return new{T}(emb)
9      end
10
11     function (embedding::Embedding{T})(indices::Vector{Int}) where T<:Real
12         # Ensure indices are valid
13         @assert all(1 .≤ indices .≤ size(embedding.emb, 1)) "Indices out of range"
14         # Perform lookup for each index
15         return embedding.emb[indices, :]
16     end
17 end
18
19

```

```

1 # Test Embedding module
2 let
3     num_embeddings = 100 # Number of patches
4     embedding_dim = 64 # Embedding dimension
5
6     # Create the embedding module
7     embedding = Embedding{Float64}(num_embeddings, embedding_dim)
8
9     # Sample indices to lookup
10    indices = [1, 10, 50, 100]
11
12    # Lookup embeddings for the indices
13    embeddings = embedding(indices)
14
15    # Verify output
16    println("Indices: ", indices)
17    println("Embedding shape: ", size(embeddings)) # Should be (length(indices),
18    embedding_dim)
19 end

```

Indices: [1, 10, 50, 100]
Embedding shape: (4, 64)



Almost there!

We just need to define a few more layers to put together the ViT.

```

1 struct Sequential{T<:Real}
2   # This is an array of modules where each module is applied sequentially to the
3   # input.
4   seq::AbstractVector
5
6   function Sequential{T}(seq::AbstractVector) where T<:Real
7     return new{T}(seq)
8   end
9
10  function (sequential::Sequential{T})(x::Array{T}) where T<:Real
11    for mod in sequential.seq
12      x = mod(x)
13    end
14    return x
15  end

```

```

1 # Test 'Sequential' implementation
2 let
3   dim, mlp_dim, nout = 3, 8, 5
4   layer1 = Linear{Float64}(dim, mlp_dim)
5   layer2 = Linear{Float64}(mlp_dim, nout)
6
7   # Initialize 'Sequential' with two layers
8   sequential = Sequential{Float64}([layer1, layer2])
9
10  # Create a sample input
11  X = randn(Float64, n_tokens, dim) # Input with n_tokens tokens of dim-
12  # dimensional vectors
13
14  # Apply 'Sequential'
15  output = sequential(X)
16
17  # Verify output shape
18  println("Input shape: ", size(X))           # Should be (10, dim)
19  println("Output shape: ", size(output))       # Should be (10, nout)
end

```

Input shape: (5, 3)
Output shape: (5, 5)



```

1 struct LayerNorm{T<:Real}
2   dim::Int
3
4   function LayerNorm{T}(dim::Int) where T<:Real
5     return new{T}(dim)
6   end
7
8   function (layernorm::LayerNorm{T})(X::Array{T}) where T<:Real
9     mean_val = mean(X; dims=layernorm.dim)
10    std_val = std(X; dims=layernorm.dim)
11    return (X .- mean_val) ./ (std_val .+ eps(T))
12  end
13 end
14

```

```

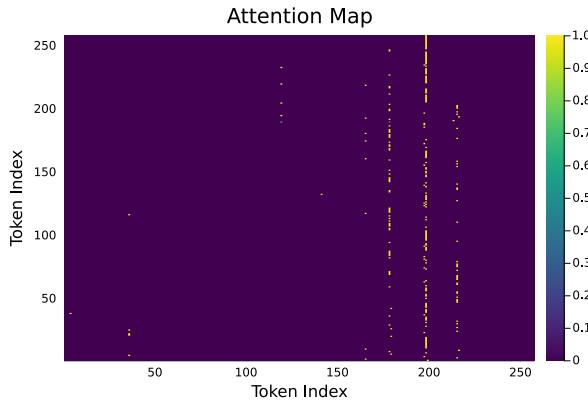
1 struct Parameter{T<:Real}
2   param::Vector{T}
3
4   function Parameter{T}(dim::Int) where T<:Real
5     param = randn(T, dim)
6     return new{T}(param)
7   end
8 end

```

```

1 struct VisionTransformer{T<:Real}
2   patch_embed::PatchEmbedLinear{T} # Patch embedding layer
3   pos_E::Embedding{T} # Positional encoding
4   cls_token::Parameter{T} # Learned class embedding token
5   transformer::Transformer{T} # Transformer encoder
6   head::Sequential{T} # Classification head
7
8   function VisionTransformer{T}(
9     n_channels::Int, nout::Int, img_size::Int, patch_size::Int, dim::Int,
10    attn_dim::Int, mlp_dim::Int, num_heads::Int, num_layers::Int
11  ) where T<:Real
12    # Initialize each component
13    patch_embed = PatchEmbedLinear{T}(img_size, patch_size, n_channels, dim)
14    pos_E = Embedding{T}((img_size ÷ patch_size)^2, dim)
15    cls_token = Parameter{T}(dim)
16    transformer = Transformer{T}(dim, attn_dim, mlp_dim, num_heads, num_layers)
17    head = Sequential{T}([LayerNorm{T}(dim), Linear{T}(dim, nout)]) # LayerNorm
18  along embedding dim
19
20  return new{T}(patch_embed, pos_E, cls_token, transformer, head)
21 end
22
23 function (vt::VisionTransformer{T})(img::Matrix{<:RGB}); return_attn=false)
24 where T<:Real
25   # Generate patch embeddings
26   embs = vt.patch_embed(img) # Shape: (num_patches, dim)
27
28   # Add positional encoding
29   N, D = size(embs) # Number of patches (N) and embedding dimension (D)
30   pos_ids = collect(1:N) # Generate positional indices
31   embs .= vt.pos_E.(pos_ids) # Add positional encodings to embeddings
32
33   # Add the class token
34   cls_token = vt.cls_token.param # Shape: (dim,)
35   x = vcat(cls_token', embs) # Shape: (N+1, dim)
36
37   # Apply the transformer
38   x, alphas = vt.transformer(x; attn_mask=nothing, return_attn=return_attn)
39
40   # Pass through the classification head
41   cls_token_out = reshape(x[1, :], 1, :) # Reshape into a matrix for
42   'Sequential'
43   out = vt.head(cls_token_out)[1, :] # Final output as a vector
44
45   return out, alphas
46 end
47

```



```

1 # Test Vision Transformer implementation for single-image inputs
2 let
3     # Define hyperparameters
4     img_size = 256          # Image size (assumes square image)
5     patch_size = 16         # Patch size
6     n_channels = 3          # Number of input channels (RGB)
7     dim = 64                # Embedding dimension
8     attn_dim = 128           # Attention hidden dimension
9     mlp_dim = 256            # Feed-forward network hidden dimension
10    num_heads = 4             # Number of attention heads
11    num_layers = 6            # Number of transformer layers
12    # nout = 10               # Number of output classes (e.g., for classification)
13
14    # Initialize the Vision Transformer
15    vt = VisionTransformer{Float64}(
16        n_channels, nout, img_size, patch_size, dim, attn_dim, mlp_dim, num_heads,
17        num_layers
18    )
19
20    # Use the color image of Philip the dog (loaded as 'image_square')
21    out, alphas = vt(image_square, return_attn=true)
22
23    # Verify output dimensions
24    println("Output shape: ", size(out)) # Should be (nout,)
25    println("Attention weights shape: ", size(alphas)) # Should be (num_layers,
26        num_heads, N+1, N+1)
27
28    # Visualize one attention map (optional)
29    heatmap(
30        alphas[1][1, :, :], # Visualize the attention weights for the first layer
31        and first head
32        title="Attention Map",
33        xlabel="Token Index",
34        ylabel="Token Index",
35        c=viridis,
36        clabel="Attention Weight (first layer)"
37    )
38 end
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
559

```

Output shape: (32,)
Attention weights shape: (6,)



Loading an image dataset (CIFAR10)

[6, 9, 9, 4, 1, 1, 2, 7, 8, 3, 4, 7, 7, 2, 9, 9, 9, 3, 2, 6, more ,2, 1, 3, 5, 7, 3, 5, 1, 3, !

```

1 begin
2     # Load CIFAR-10 training data (already loaded above)
3     # train_dataset = CIFAR10(dir="cifar/", split=:train)
4     # test_dataset = CIFAR10(dir="cifar/", split=:test)
5
6     subset_size = 1000
7     train_subset = train_dataset.features[:, :, :, 1:subset_size]
8     train_labels = train_dataset.targets[1:subset_size]
9 end

```



```

1 let
2     # Extract a single CIFAR-10 image and label
3     img_data = train_dataset.features[:, :, :, 1] # Shape: 32x32x3 (HWC format)
4     img_data_permuted = permutedims(img_data, (3, 1, 2)) # CHW format for the
5     VisionTransformer
6     rgb_image = colorview(RGB, img_data_permuted)
7     target_label = train_dataset.targets[1] # Class label (1-based index)
8
9     # Display the CIFAR-10 image
10    rgb_image
11 end

```

Cross-Entropy Loss Function

We will define the cross-entropy loss function. For multi-class classification with C classes, the cross-entropy loss for a single sample is given by:

$$\text{Loss} = - \sum_{c=1}^C y_c \log(\hat{y}_c)$$

Here:

- y_c is 1 if the sample belongs to class c , otherwise 0.
- \hat{y}_c is the predicted probability for class c .

```
cross_entropy_loss (generic function with 1 method)
1 function cross_entropy_loss(predictions::Matrix{Float64}, targets::Vector{Int})
2     # Convert targets to one-hot encoding
3     num_samples, num_classes = size(predictions)
4     one_hot_targets = zeros(Float64, num_samples, num_classes)
5     for i in 1:num_samples
6         one_hot_targets[i, targets[i]] = 1.0
7     end
8     # Compute log of softmax predictions
9     log_probs = log.(softmax(predictions, dims=2))
10    # Compute the loss
11    loss = -sum(one_hot_targets .* log_probs) / num_samples
12    return loss
13 end
```

```
compute_gradient (generic function with 1 method)
1 # Compute gradient for the VisionTransformer
2 function compute_gradient(vit_model::VisionTransformer{T}, img::Array{T},
3     target::Int) where T<:Real
4     function loss_fn()
5         # Forward pass through the model
6         output, _ = vit_model(img)
7         # Compute the loss (batch size = 1)
8         y_pred = reshape(output, 1, :) # Ensure output is a matrix
9         return cross_entropy_loss(y_pred, [target])
10    end
11    # Compute gradients of the loss function with respect to model parameters
12    grad = Enzyme.gradient(loss_fn, vit_model)
13    return grad
14 end
```

```
1 # Test the VisionTransformer with gradient computation
2 # let
3 #     # Initialize the Vision Transformer
4 #     img_size = 32 # CIFAR-10 image size
5 #     patch_size = 4 # Patch size
6 #     n_channels = 3 # Number of input channels (RGB)
7 #     dim = 64 # Embedding dimension
8 #     attn_dim = 128 # Attention hidden dimension
9 #     mlp_dim = 256 # Feed-forward network hidden dimension
10 #    num_heads = 4 # Number of attention heads
11 #    num_layers = 6 # Number of transformer layers
12 #    nout = 10 # Number of output classes (CIFAR-10 has 10 classes)
13 #
14 #    vit_model = VisionTransformer{Float32}(
15 #        n_channels, nout, img_size, patch_size, dim, attn_dim, mlp_dim, num_heads,
16 #        num_layers
16 #    )
17 #
18 #    # Compute gradients
19 #    img_data = train_dataset.features[:, :, :, 1] # Shape: 32x32x3 (HWC format)
20 #    img_data_permuted = permutedims(img_data, (3, 1, 2))
21 #    rgb_image = colorview(RGB, img_data_permuted)
22 #    target_label = train_dataset.targets[1] # Class label (1-based index)
23 #    grad = compute_gradient(vit_model, img_data_permuted, target_label)
24 #
25 #    # Print gradient details
26 #    println("Gradient: ", grad)
27 # end
28
```

1 Enter cell code...