



# Vision Transformers for Plant Traits Prediction

## Lisa Fung & Annabelle Aurelia Jayadinata

### AFFILIATIONS

Stanford University, Department of Computer Science & Electrical Engineering, CS231N

### 01. Problem

- Climate change is expected to cause approximately 250,000 additional deaths per year in 2030-2050, from under nutrition, malaria, diarrhea and heat stress alone.
- Rapidly transforming biosphere causes shifts in ecosystem species distributions and plant adaptations, measured by **changes in plant traits**
- Predicting the global-scale impact of climate change on ecosystems is **difficult to quantify** due to **insufficient data on plant traits**

**Objective:** Employ deep learning-based regression models, including Convolutional Neural Networks (CNNs) and Vision Transformers to predict plant traits from photographs and ancillary geodata



trait.ID	trait.name
X4	Stem specific density (SSD) or wood density (stem dry mass per stem fresh volume)
X11	Leaf area per leaf dry mass (specific leaf area, SLA or LMA)
X18	Plant height
X26	Seed dry mass
X50	Leaf nitrogen (N) content per leaf area
X312	Leaf area (in the case of compound leaves: leaf, undefined if petiole in- or excluded)

### 02. Background

- Datasets:** Data is a pre-labeled link between the TRY database (plant trait information) and the iNaturalist database (55,500 citizen science plant photographs) provided through Kaggle.
  - Additional location-based geo-datasets for each plant photograph, including temperature and precipitation data, and satellite data.
  - Use a 4:1 split for training and validation.
- Performance Metrics:** The evaluation metric is the mean R2 (R-Squared) over all 6 traits.

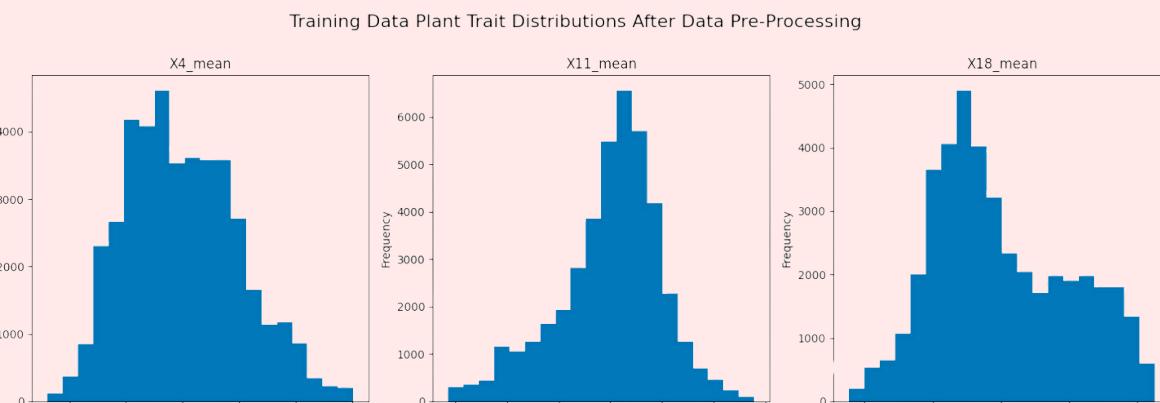
$$R^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

#### Loss Function:

$$\text{SmoothL1}(x) = \begin{cases} 0.5x^2 & \text{if } |x| < 1 \\ |x| - 0.5 & \text{otherwise.} \end{cases}$$

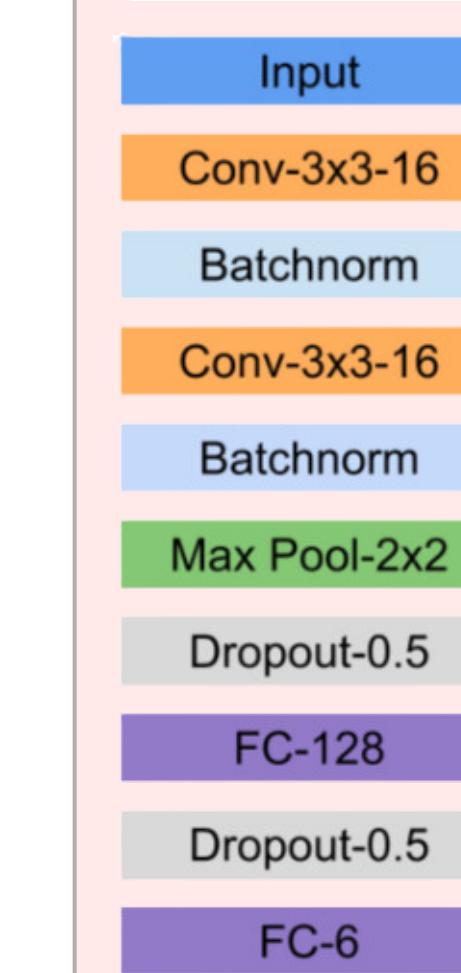
#### Data Preprocessing:

- Filter outliers using a quantile range (0.005, 0.995) for all traits
- Apply a logarithmic base 10 transformation to right-skewed trait data, which includes the five plant traits except trait X4
- Normalize data to achieve a zero mean and unit standard deviation across all traits



### 03. Methods

#### BASELINE MODEL

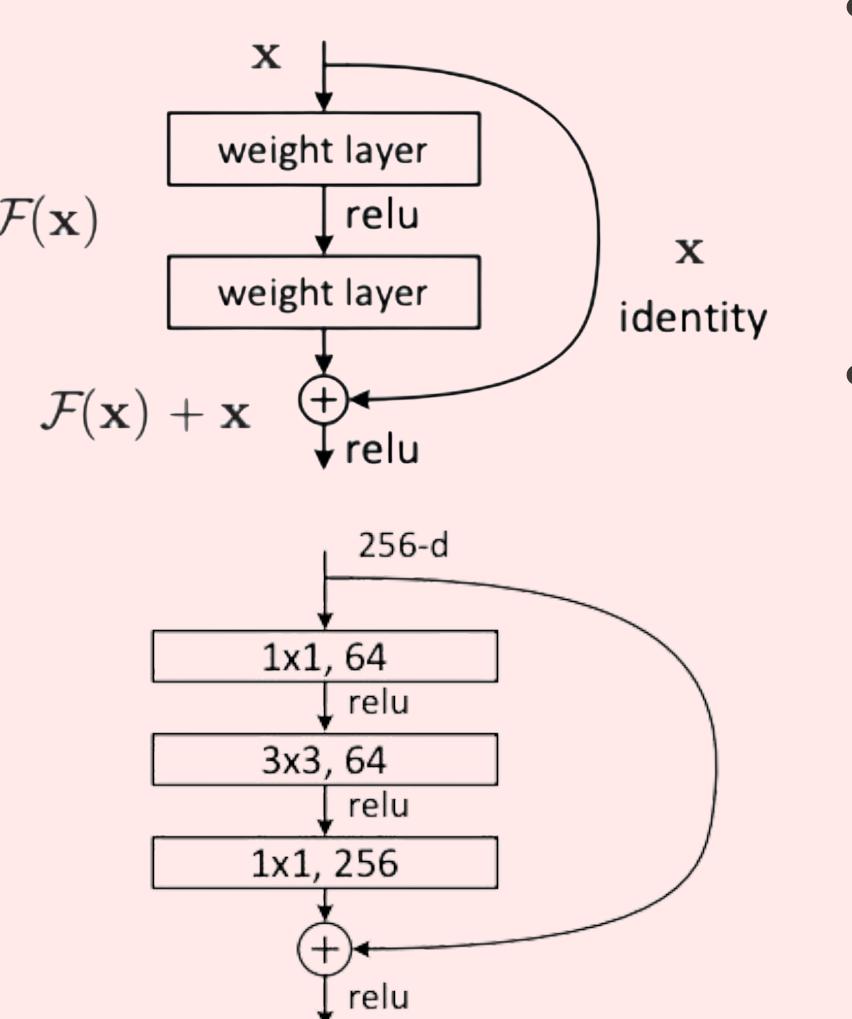


CNN model with two 3x3 convolutional layers, one 2x2 max pooling layer, and two fully connected layers.

- ReLU activations, batch normalization, and dropout. Training uses Adam optimizer with a step learning rate scheduler initialized to 1e-4.

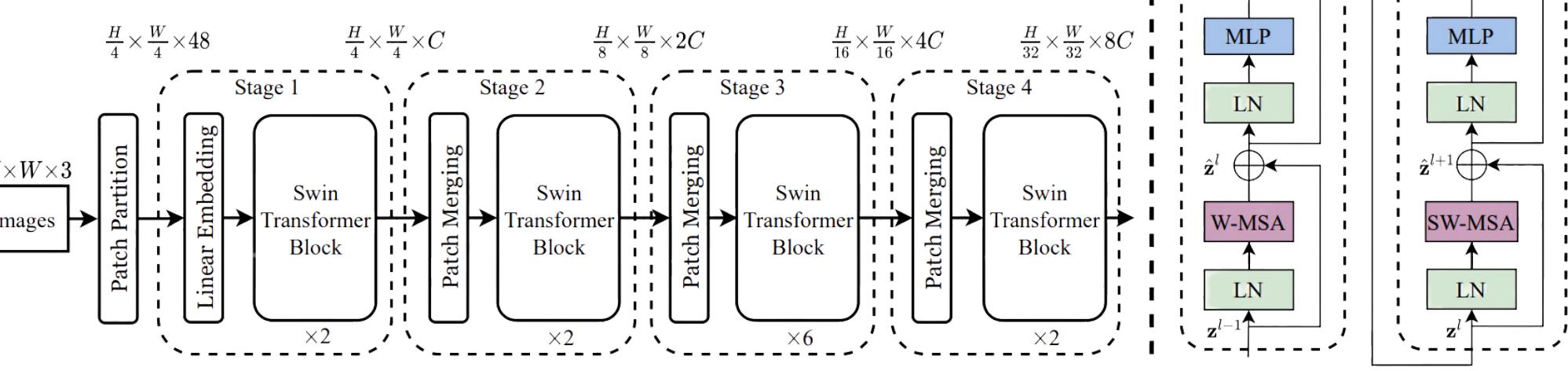
**VGG-16** A simple CNN architecture consisting of 16 layers (13 convolutional and 3 fully connected layers)

#### RESNET-50



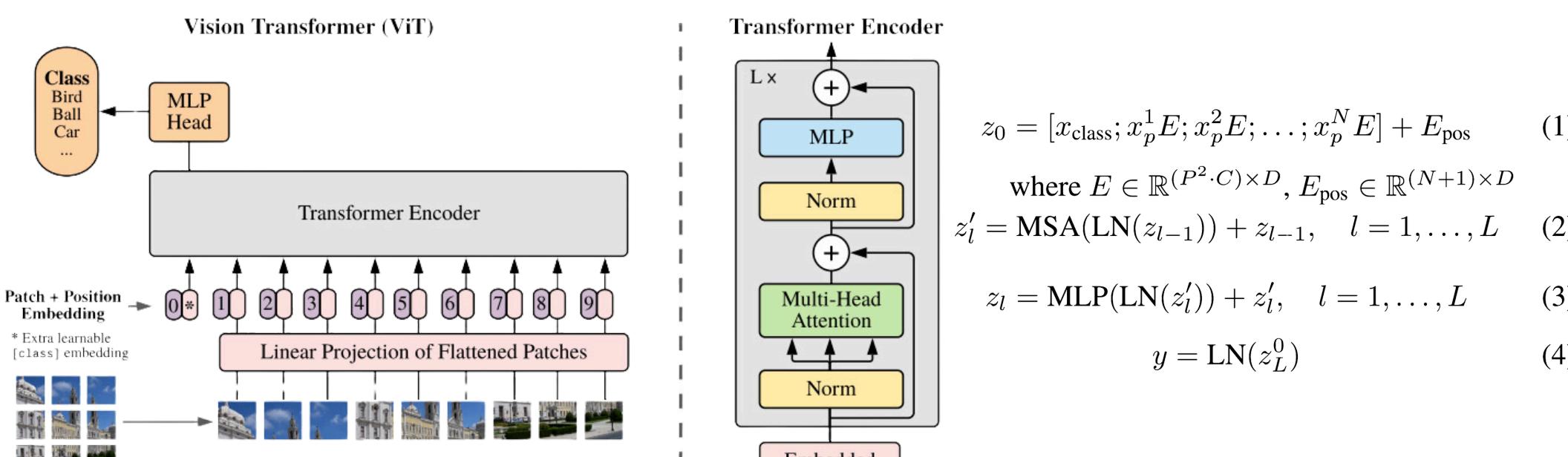
- ResNet's **residual mappings** allow the stacked nonlinear layers to fit another mapping of  $F(x) := H(x) - x$ . The original mapping is recast into  $F(x) + x$ .
- In ResNet-50, for each residual function  $F$ , we use a stack of 3 layers which are 1x1, 3x3, and 1x1 convolutions
  - 1x1 layers reduce and then increase (restore) dimensions
  - 3x3 layer is a bottleneck with smaller input and output dimensions

#### SWIN TRANSFORMER



- Swin Transformers rely on **hierarchical feature extraction**
  - Features extracted at different scales, then combined to form a holistic representation of the image
  - At each level, features are split into non-overlapping patches and processed by a series of transformer blocks

#### VISION TRANSFORMER (ViT)



- The Vision Transformer (ViT) encoder consists of alternating layers of multiheaded self-attention and MLP blocks. LayerNorm (LN) is applied before every block, and residual connections placed after every block. The MLP contains two layers with a GELU non-linearity.

### 04. Experiments & Results

#### SWIN TRANSFORMER

We fine-tune three different pre-trained Swin Transformer models that take in plant images as input and output 6 plant trait values:

1. SwinV2 Tiny Window Size  
16x16, Image Size 256 x 256:  
28.3M parameters

2. SwinV2 Small Window Size 16x16, Image Size 256x256: 49.7M parameters

3. Swin Large Window Size 12x12, Image Size 384x384: 196.7M parameters

Swin Transformer Model	Train $R^2$	Validation $R^2$	Test $R^2$	
$R_{\log}^2$	$R_{\text{orig}}^2$	$R_{\log}^2$	$R_{\text{orig}}^2$	
SwinV2 Tiny baseline	0.5726	X	X	<b>0.26691</b>
SwinV2 Tiny second (12 epochs)	0.7448	X	X	0.21284
SwinV2 Tiny second (8 epochs)	0.6317	X	X	0.24936
SwinV2 Tiny log-adjusted	X	0.0795	X	0.12825
SwinV2 Small	0.6593	X	0.364	<b>0.27308</b>
Swin Large	0.8812	X	X	0.053
				<b>0.27745</b>

Table 4. Summary of Swin Transformer Models Performance

(a) Validation Sample R2 Scores, Average R2 = 0.26694

X4: 0.40, X11: 0.11, X18: 0.12, X26: 0.12, X50: 0.12, X312: 0.22

(b) Validation Sample R2 Scores, Average R2 = 0.21284

X4: 0.9021, X11: 0.0367, X18: 0.21284, X26: 0.21284, X50: 0.21284, X312: 0.21284

Part Test Sample: 0.21284

Part Validation: 0.21284

Part Test Sample: 0.21284