

Computer Vision - 2025

Lecture #03. Visual Transformers

Lectures by Alexei Kornaev ^{1,2,3}

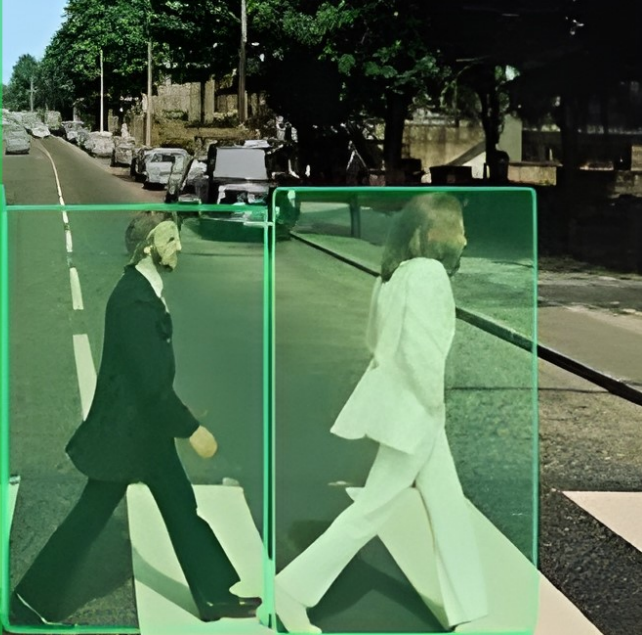
Practical sessions by Kirill Yakovlev ²

¹AI Institute, Innopolis University (IU), Innopolis

²Robotics & CV Master's Program, IU, Innopolis

³RC for AI, National RC for Oncology, Moscow

February 2, 2025



Agenda

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

① Outcomes

② Transformer Architecture Overview

③ Vision Transformers (ViTs)

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

Section 1. Outcomes

Outcomes

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

This week's lecture on Vision Transformers (ViTs) aims to provide an understanding of transformer-based architectures in computer vision. By the end of this week, students will be able to:

- 1 Understand the core concepts of Transformer architecture: embeddings, positional encoding, self-attention, and multi-head attention.
- 2 Explain Vision Transformers (ViTs), including patch embeddings, transformer blocks, and classification heads.
- 3 Compare CNNs and ViTs for vision tasks, highlighting their strengths and weaknesses.
- 4 Analyze modern ViT architectures such as DeiT and hybrid CNN-ViT models.

Key Takeaway: Vision Transformers (ViTs) provide an alternative to CNNs by leveraging self-attention, enabling global context modeling in images .

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

Section 2. Transformer Architecture Overview

Transformer Architecture

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

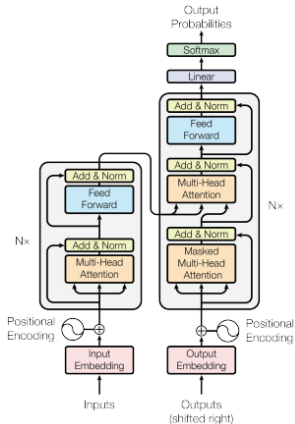


Figure: High-level Transformer architecture, consisting of an encoder and decoder stack [Vaswani et al., 2023].

Self-Attention Mechanism

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

Self-attention allows the model to weigh different parts of the input sequence when making predictions.

Scaled Dot-Product Attention

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V$$

where:

- Q, K, V are the query, key, and value matrices.
- d_k is the dimensionality of the key vectors.

Physiology of cats

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

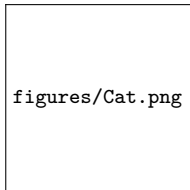


Figure: Responses of the cat's visual cortex cell.

Implementation in math.

The horizontal derivative kernel approximates $\frac{\partial I}{\partial x}$ using **finite differences**. For images, this translates to computing intensity changes along the x-axis. The Sobel kernel for horizontal derivative is:

$$\begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix}$$

This kernel computes:

$$\frac{\partial I}{\partial x} \approx (I(x+1, y) - I(x-1, y))$$

and incorporates smoothing to reduce noise. It detects horizontal edges by highlighting intensity changes from left to right.

Positional Embedding in Transformers

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

Why Positional Embedding?

Transformers, unlike CNNs or RNNs, do not inherently understand the order of input tokens. Positional embeddings are added to input embeddings to inject information about the position of tokens in the sequence.

- **Input Embedding:** Represents the content of each token.
- **Positional Embedding:** Represents the position of each token.
- **Combined:** Input embedding + Positional embedding.

$$\mathbf{E} = \mathbf{E}_{\text{input}} + \mathbf{E}_{\text{position}} \quad (1)$$

Key Takeaway: Positional embeddings enable Transformers to process sequences with order-awareness.

Sinusoidal Positional Encoding

CV-2025

A.Korbaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

Sinusoidal Encoding

In the original Transformer paper [Vaswani et al., 2023], sinusoidal functions are used to generate positional encodings. These encodings are deterministic and do not require learning.

$$PE_{(pos, 2i)} = \sin \left(\frac{pos}{10000^{\frac{2i}{d}}} \right) \quad (2)$$

$$PE_{(pos, 2i+1)} = \cos \left(\frac{pos}{10000^{\frac{2i}{d}}} \right) \quad (3)$$

- pos : Position of the token in the sequence.
- i : Dimension index of the embedding.
- d : Dimensionality of the embedding.

Key Takeaway: Sinusoidal encoding allows the model to generalize to sequences longer than those seen during training.

Positional Embedding in Vision Transformers (ViTs)

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

Patch Embeddings and Positional Encoding

In Vision Transformers (ViTs), the input image is divided into patches, and each patch is treated as a token. Positional embeddings are added to these patch embeddings to preserve spatial information.

$$\mathbf{E}_{\text{patch}} = \text{Linear}(\text{Flatten}(\text{Patches})) \quad (4)$$

$$\mathbf{E} = \mathbf{E}_{\text{patch}} + \mathbf{E}_{\text{position}} \quad (5)$$

- **Patches:** Image divided into fixed-size patches (e.g., 16x16).
- **Positional Embedding:** Added to patch embeddings to encode spatial location.

Key Takeaway: Positional embeddings in ViTs help the model understand the spatial arrangement of patches in the image.

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

Section 3. Vision Transformers (ViTs)

Image to Patch Splitting

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

Key Idea

Instead of processing the image as a whole, ViTs divide it into **non-overlapping patches** of fixed size.

- Given an image $x \in \mathbb{R}^{H \times W \times C}$ (height H , width W , channels C), we divide it into N patches of size $P \times P$.
- The number of patches is:

$$N = \frac{HW}{P^2}$$

Linear Projection of Flattened Patches

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

Patch Embeddings

Each patch is **flattened** into a vector and projected into a higher-dimensional space using a linear layer.

- Each patch $x_p \in \mathbb{R}^{P \times P \times C}$ is **flattened** into a vector $x_p \in \mathbb{R}^{P^2 C}$.
- A trainable weight matrix $E \in \mathbb{R}^{D \times P^2 C}$ projects it into a **D-dimensional** embedding:

$$z_i = Ex_i, \quad i = 1, \dots, N$$

Positional Embeddings

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

Why Positional Encoding?

Transformers process input as a **set of tokens** without spatial order. Positional embeddings restore spatial information.

- Each patch embedding receives a **learned** positional encoding E_p .
- The final input to the Transformer is:

$$z_i = E x_i + E_p(i), \quad i = 1, \dots, N$$

- A **[CLS]** token is prepended for classification tasks.

Learned Positional Embeddings in ViTs

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

Instead of fixed sinusoidal functions, Vision Transformers often use learned positional embeddings:

$$\mathbf{z}_0 = \text{class token}, \quad \mathbf{z}_i = E(x_i) + E_p(i), \quad i = 1, \dots, N$$

where:

- $E(x_i)$ is the patch embedding,
- $E_p(i)$ is the learned positional embedding for patch i ,
- \mathbf{z}_0 is an extra **classification token**.

Transformer Encoder

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

Self-Attention for Global Feature Learning

ViTs apply the standard Transformer encoder from NLP to image patches.

- Each layer consists of **Multi-Head Self-Attention (MSA)** and **MLP blocks**.
- Layer normalization (LN) is applied before each block.

Multi-Head Self-Attention

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V$$

where Q, K, V are query, key, and value matrices.

MLP Head for Classification

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

Final Prediction Layer

After Transformer encoding, the **[CLS]** token representation is passed through an MLP for classification.

- The final representation of the **[CLS]** token z_0 is used for classification:

$$y = \text{MLP}(z_0)$$

- The MLP consists of two fully connected layers with a non-linearity (ReLU/GELU).



Comparison: CNNs vs. ViTs

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

Key Differences

- CNNs learn spatial hierarchies via local receptive fields and weight sharing.
- ViTs model long-range dependencies using self-attention but lack inductive biases like locality and translation invariance.

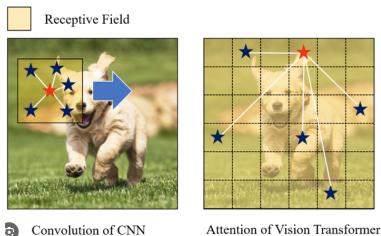


Figure: Comparison of CNNs and ViTs in feature extraction and learning patterns [Baek et al., 2022]

Strengths and Weaknesses of ViTs

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

Strengths

- **Long-range dependencies** captured effectively.
- **Scalability** to large datasets.
- **Flexibility** to adapt across domains (text, vision, multimodal).

Weaknesses

- **Data-hungry** - requires large-scale training.
- **Computationally expensive** due to self-attention complexity ($O(N^2)$).
- **Lack of inductive bias** â struggles with small datasets.

Conclusion

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

Key Takeaways

- Vision Transformers (ViTs) treat images as sequences of patches and use **self-attention** for feature extraction.
- Unlike CNNs, ViTs do not rely on convolution but instead learn **global relationships** from data.
- While powerful, ViTs require **large datasets and high computational resources** to generalize well.
- Hybrid architectures (CNN + ViT) help balance efficiency and performance.

Looking Ahead: Modern architectures like **DeiT, Swin Transformer, and hybrid models** address ViT limitations.

Bibliography

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

- A. Dosovitskiy et al., "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale," NeurIPS 2020. [arXiv:2010.11929](https://arxiv.org/abs/2010.11929)
- A. Vaswani et al., "Attention Is All You Need," NeurIPS 2017. [arXiv:1706.03762](https://arxiv.org/abs/1706.03762)
- H. Touvron et al., "Training Data-efficient Image Transformers Distillation through Attention," ICML 2021. [arXiv:2012.12877](https://arxiv.org/abs/2012.12877)
- Z. Liu et al., "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows," ICCV 2021. [arXiv:2103.14030](https://arxiv.org/abs/2103.14030)
- S. Prince, *Understanding Deep Learning*, MIT Press, 2023.

Hands-on coding: **Kernels** [**CV-2025**]

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

Hands-on the book by Howard and Gugger [2020]

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

ViT Models Zoo

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

The following resources should be met:

- ① ViTs, Community Computer Vision Course by Hugging Face
- ② Vision Transformer (ViT) by Hugging Face

Bibliography

CV-2025

A.Kornaev,
K.Yakovlev

Outcomes

Transformer
Architecture
Overview

Vision
Transformers
(ViTs)

Sihun Baek, Jihong Park, Praneeth Vepakomma, Ramesh Raskar, Mehdi Bennis, and Seong-Lyun Kim. Visual transformer meets cutmix for improved accuracy, communication efficiency, and data privacy in split learning, 2022. URL <https://arxiv.org/abs/2207.00234>.

CV-2025. The course repository. <https://github.com/Mechanics-Mechatronics-and-Robotics/CV-2025>, 2025. Accessed: February 2, 2025.

J. Howard and S. Gugger. **Deep Learning for Coders with Fastai and Pytorch: AI Applications Without a PhD**. O'Reilly Media, Incorporated, 2020. ISBN 9781492045526. URL <https://books.google.no/books?id=xd6LxgEACAAJ>.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023. URL <https://arxiv.org/abs/1706.03762>.