

Computer Vision - 2025

Week #13. Multi-Modal Data Processing - II

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

Practical sessions by Kirill Yakovlev ²

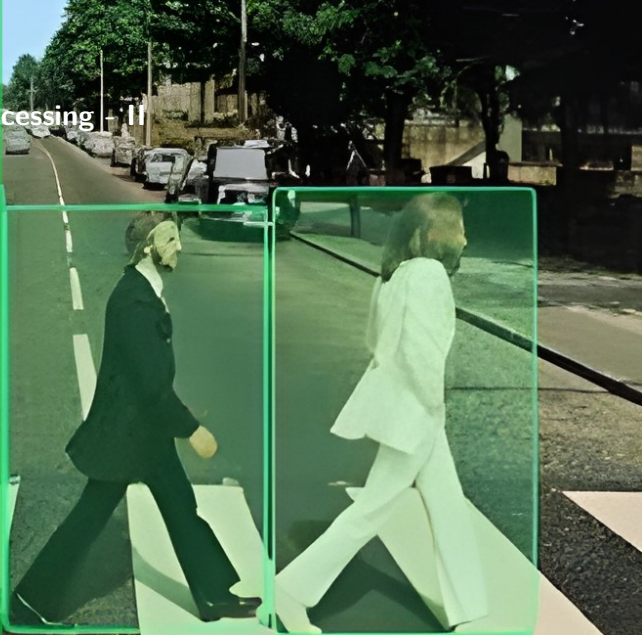
¹AI Institute, Innopolis University (IU), Innopolis

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

³Dept. of M^2R , Orel State University, Orel

⁴RC for AI, National RC for Oncology, Moscow

April 14, 2025



Agenda

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

- 1 Recap: CLIP architecture and loss
- 2 DNNs Training Paradigms
- 3 VLM's Objectives (losses)
- 4 Transfer Learning of Large Models
- 5 Prospects
- 6 VLM + Control System = VLAM
- 7 Other Types of Multimodalities
- 8 Meta-Learning

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Multi-

Section 1. Recap: CLIP architecture and loss

CLIP Loss: Core Mechanism

CV-2025

A.Korbaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs

Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

Definition

Align image-text pairs in a shared space using symmetric contrastive loss:

$$\mathcal{L}_{\text{CLIP}} = -\frac{1}{2N} \sum_{i=1}^N \left[\log \frac{e^{s_{ii}/\tau}}{\sum_j e^{s_{ij}/\tau}} + \log \frac{e^{s_{ii}/\tau}}{\sum_j e^{s_{ji}/\tau}} \right]$$

Handwritten red notes: "I, I, I" above the first log term and "I, I, I" above the second log term.

where $s_{ij} = \text{cos_sim}(I_i, T_j)$ for image and text embeddings, τ is the temperature parameter (learned or fixed) to scale logits.

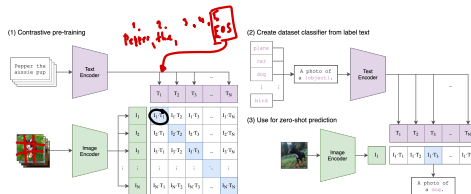


Figure: Contrastive Language-Image Pre-training (CLIP) [Radford et al., 2021]

Pre-test: CLIP Loss Example with a Small Batch (N=2)

CV-2025

A.Korbaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs

Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

Setup

Batch size $N = 2$, temperature $\tau = 0.07$; **image embeddings**: $I_1 = [0.8, 0.6]$, $I_2 = [0.6, 0.8]$; **text embeddings**: $T_1 = [0.7, 0.7]$, $T_2 = [0.7, -0.7]$; normalized embeddings: $\|I\| = \|T\| = 1$; **similarity Matrix**:

$$S = \begin{bmatrix} I_1 \cdot T_1 & I_1 \cdot T_2 \\ I_2 \cdot T_1 & I_2 \cdot T_2 \end{bmatrix} = \begin{bmatrix} 0.98 & 0.14 \\ 0.98 & -0.14 \end{bmatrix}.$$

Step 1: Image \rightarrow Text Loss (for I_1)

$$\text{Softmax}_{\tau}(S_{I_1}) = \frac{e^{X.XX/X.XX}}{e^{X.XX/X.XX} + e^{X.XX/X.XX}} \approx 1.0,$$

$$\mathcal{L}_{CE}(I_1) = -\log(1.0) \approx 0.$$

Step 2: Text \rightarrow Image Loss (for T_2)

$$\text{Softmax}_{\tau}(S_{T_2}) = \frac{e^{X.XX/X.XX}}{e^{X.XX/X.XX} + e^{X.XX/X.XX}} \approx 0.018,$$

$$\mathcal{L}_{CE}(T_2) = -\log(0.018) \approx 4.0.$$

Total Loss

$$\mathcal{L}_{CLIP} = \frac{1}{2 \times 2} (0 + 4.0 + \dots) \quad (\text{Sum over all pairs}).$$

Pre-test: CLIP Loss Example with a Small Batch (N=2)

CV-2025

A.Korbaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs

Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

Setup

Batch size $N = 2$, temperature $\tau = 0.07$; **image embeddings**: $I_1 = [0.8, 0.6]$, $I_2 = [0.6, 0.8]$; **text embeddings**: $T_1 = [0.7, 0.7]$, $T_2 = [0.7, -0.7]$; normalized embeddings: $\|I\| = \|T\| = 1$; **similarity Matrix**:

$$S = \begin{bmatrix} I_1 \cdot T_1 & I_1 \cdot T_2 \\ I_2 \cdot T_1 & I_2 \cdot T_2 \end{bmatrix} = \begin{bmatrix} 0.98 & 0.14 \\ 0.98 & -0.14 \end{bmatrix}.$$

Step 1: Image \rightarrow Text Loss (for I_1)

$$\text{Softmax}_{\tau}(S_{I_1}) = \frac{e^{0.98/0.07}}{e^{0.98/0.07} + e^{0.14/0.07}} = \frac{e^{14}}{e^{14} + e^2} \approx 1.0,$$
$$\mathcal{L}_{CE}(I_1) = -\log(1.0) \approx 0.$$

Step 2: Text \rightarrow Image Loss (for T_2)

$$\text{Softmax}_{\tau}(S_{T_2}) = \frac{e^{-0.14/0.07}}{e^{0.14/0.07} + e^{-0.14/0.07}} = \frac{e^{-2}}{e^2 + e^{-2}} \approx 0.018,$$
$$\mathcal{L}_{CE}(T_2) = -\log(0.018) \approx 4.0.$$

Total Loss

$$\mathcal{L}_{\text{CLIP}} = \frac{1}{2 \times 2} (0 + 4.0 + \dots) \quad (\text{Sum over all pairs}).$$

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

Section 2. DNNs Training Paradigms

CLIP is Followed by...

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

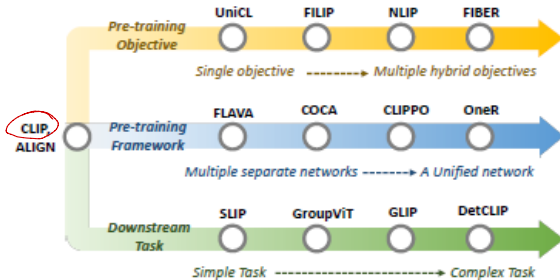


Figure: Illustration of development of VLMs for visual recognition [Zhang et al., 2024]

Three DNNs Training Paradigms

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

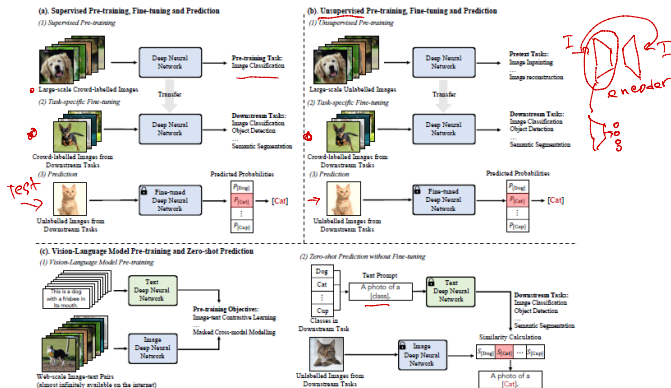


Figure: Three DNN training paradigms in visual recognition. Compared with the paradigms in (a) and (b) that requires fine-tuning for each specific task with task-specific labelled data, the new learning paradigm with VLMs in (c) enables effective usage of web data and zero-shot predictions without task-specific fine-tuning [Zhang et al., 2024]

Pre-Training → Transfer → Distillation

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

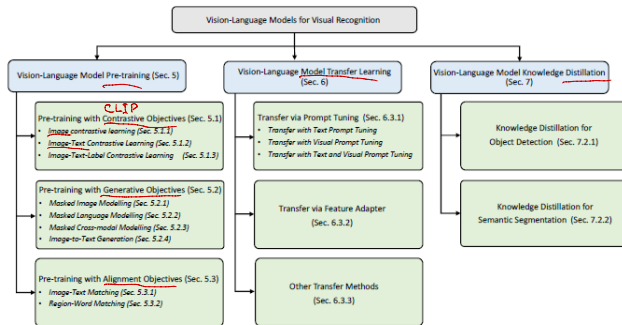


Figure: Typology of vision-language models for visual recognition [Zhang et al., 2024]

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Multi-

Section 3. VLM's Objectives (losses)

VLM Pretraining Objectives

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs

Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

Three Categories

- **Contrastive**: Align visual and textual embeddings via pairwise similarity
- **Generative**: Predict masked content or generate text given vision
- **Alignment**: Match global (image-text) or local (region-word) pairs

Setup

Let $\mathcal{D} = \{x_n^I, x_n^T\}_{n=1}^N$, $z^I = \underline{f_\theta(x^I)}$, $z^T = \underline{f_\phi(x^T)}$.

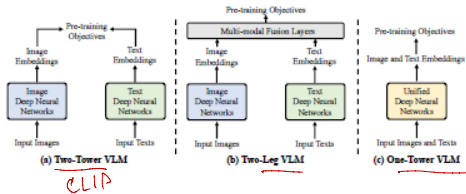


Figure: Typical VLM frameworks [Zhang et al., 2024]

Contrastive Objectives (1)

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Multi-

Image Contrastive Learning (InfoNCE)

$$\mathcal{L}_{\text{InfoNCE}}^I = -\frac{1}{B} \sum_{i=1}^B \log \frac{\exp(z_i^I \cdot z_i^I / \tau)}{\sum_{j \neq i} \exp(z_i^I \cdot z_j^I / \tau)}$$

Image-Text Contrastive Learning (CLIP)

$$\mathcal{L}_{\text{InfoNCE}}^{I \rightarrow T} = -\frac{1}{B} \sum_i \log \frac{\exp(z_i^I \cdot z_i^T / \tau)}{\sum_j \exp(z_i^I \cdot z_j^T / \tau)}$$

$$\mathcal{L}_{\text{InfoNCE}}^{T \rightarrow I} = -\frac{1}{B} \sum_i \log \frac{\exp(z_i^T \cdot z_i^I / \tau)}{\sum_j \exp(z_i^T \cdot z_j^I / \tau)}$$

Contrastive Objectives (2)

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs

Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

Image-Text-Label Contrastive Learning (UniCL) [Khosla et al., 2021]

$$\mathcal{L}_{\text{label}}^{I \rightarrow T} = - \sum_i \frac{1}{|P(i)|} \sum_{k \in P(i)} \log \frac{\exp(z_i^I \cdot z_k^T / \tau)}{\sum_j \exp(z_i^I \cdot z_j^T / \tau)}$$

where τ is a temperature hyperparameter, $P(i)$ is the set of positive labels for sample i

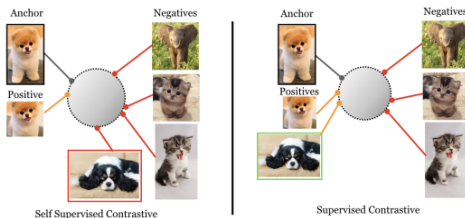


Figure: Supervised vs. self-supervised contrastive losses [Khosla et al., 2021]

Generative Objectives (1)

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

Masked Image Modeling (MIM)

$$\mathcal{L}_{\text{MIM}} = -\frac{1}{B} \sum_i \log f_{\theta}(x_i^I | \hat{x}_i^I)$$

Explanation

- Input image x_i^I is split into patches, some patches are masked: \hat{x}_i^I
- The model reconstructs the masked parts with the log-likelihood loss of correct reconstruction



Figure: Masked Image Modeling (MIM) intuition [Khosla et al., 2021]

Generative Objectives (2)

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

Masked Language Modeling (MLM)

$$\mathcal{L}_{\text{MLM}} = -\frac{1}{B} \sum_i \log f_{\phi}(x_i^T | \hat{x}_i^T)$$

Masked Cross-Modal Modeling (MCM)

$$\mathcal{L}_{\text{MCM}} = -\frac{1}{B} \sum_i \left[\log f_{\theta}(x_i^I | \hat{x}_i^I, \hat{x}_i^T) + \log f_{\phi}(x_i^T | \hat{x}_i^I, \hat{x}_i^T) \right]$$

Image-to-Text Generation

$$\mathcal{L}_{\text{ITG}} = -\sum_{l=1}^L \log f_{\theta}(x_l^T | x_{<l}^T, z^I)$$

Alignment Objectives

CV-2025

A.Korbaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

Image-Text Matching (ITM)

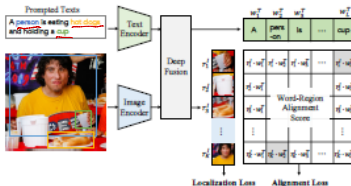
misma *hyper*

$$\mathcal{L}_{\text{ITM}} = p \cdot \log S(z^I, z^T) + (1 - p) \cdot \log(1 - S(z^I, z^T))$$

where $p \in 0, 1$ is label indicating match (1) or mismatch (0).

Region-Word Matching (RWM)

$$\mathcal{L}_{\text{RW}} = p \cdot \log S_r(r^I, w^T) + (1 - p) \cdot \log(1 - S_r(r^I, w^T))$$



Putting It Together

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

Combined Loss

$$\mathcal{L}_{\text{total}} = \lambda_1 \mathcal{L}_{\text{contrastive}} + \lambda_2 \mathcal{L}_{\text{generative}} + \lambda_3 \mathcal{L}_{\text{alignment}}$$

- Loss terms and weights depend on the training goal
- Some models (e.g., CLIP) use only contrastive; others use hybrid objectives (e.g., BLIP)

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

Section 4. Transfer Learning of Large Models

Strategies for LMs

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs

Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

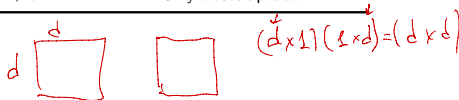
Other Types
of Models

Comparison of Parameter-Efficient Methods

Method	Added Params	Modifies Forward Pass	Key Advantage
Full FT	100%	$h = (W_0 + \Delta W)x$	Highest accuracy
LoRA	$\sim 0.1\%$	$h = W_0x + BAx$	Balance of efficiency/performance
Adapter	$\sim 1\%$	$h = W_0x + W_2(\sigma(W_1x))$	Modular
Prefix Tuning	$\sim 0.5\%$	$[P; x] \rightarrow \text{Attention}$	No backbone changes
BitFit	$\sim 0.01\%$	$h = W_0x + b$	Only biases updated

Mathematical Forms

- **Adapter:** $W_2 \in \mathbb{R}^{d \times r}, W_1 \in \mathbb{R}^{r \times d}$
- **Prefix Tuning:** $P \in \mathbb{R}^{l \times d}$ (prepended tokens)
- **BitFit:** $b \in \mathbb{R}^d$ (bias terms only)



When to Use LoRA?

- Need high parameter efficiency ($r \leq 64$)
- Preserve original model architecture
- Balance between compute and accuracy

LoRA: Low-Rank Adaptation [Hu et al., 2021]

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

Key Mathematical Formulation

For a pretrained weight matrix $W_0 \in \mathbb{R}^{d \times k}$:

$$W = W_0 + \underbrace{BA}_{\text{Low-rank update}} \quad \begin{cases} B \in \mathbb{R}^{d \times r} \\ A \in \mathbb{R}^{r \times k} \\ r \ll \min(d, k) \end{cases}$$

Example: 1024x1024 Layer with Rank=8

- Original params: $1024 \times 1024 = 1,048,576$
- LoRA params: $8 \times (1024 + 1024) = 16,384$
- Reduction: $\frac{16,384}{1,048,576} \approx \underline{1.56\%}$

LoRA in Action

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

Forward Pass Computation

For input $x \in \mathbb{R}^k$:

$$h = W_0 x + \underbrace{B(Ax)}_{\text{Rank-constrained update}}$$

Gradient Flow

- Frozen weights: $\nabla_{W_0} \mathcal{L} = 0$
- Adaptor gradients:

$$\nabla_B \mathcal{L} = (\nabla_h \mathcal{L}) x^\top A^\top \quad \nabla_A \mathcal{L} = B^\top (\nabla_h \mathcal{L}) x^\top$$

Why This Works

- Preserves pretrained knowledge (W_0 fixed)
- Efficient training (only update B, A)
- Low-rank bottleneck prevents overfitting

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Multi-

Section 5. Prospects

Paper reading

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

Vision-Language Models for Vision Tasks: A Survey

Most visual recognition studies rely heavily on crowd-labelled data in deep neural networks (DNNs) training, and they usually train a DNN for each single visual recognition task, leading to a laborious and time-consuming visual recognition paradigm. To address the two challenges, Vision-Language Models (VLMs) have been intensively investigated recently, which learns rich vision-language correlation from web-scale image-text pairs that are almost infinitely available on the Internet and enables zero-shot predictions on various visual recognition tasks with a single VLM. This paper provides a systematic review of visual language models for various visual recognition tasks, including: (1) the background that introduces the development of visual recognition paradigms; (2) the foundations of VLM that summarize the widely-adopted network architectures, pre-training objectives, and downstream tasks; (3) the widely-adopted datasets in VLM pre-training and evaluations; (4) the review and categorization of existing VLM pre-training methods, VLM transfer learning methods, and VLM knowledge distillation methods; (5) the benchmarking, analysis and discussion of the reviewed methods; (6) several research challenges and potential research directions that could be pursued in the future VLM studies for visual recognition. [Zhang et al., 2024].

Future Directions for VLMs [Zhang et al., 2024]

CV-2025

A.Korbaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

Challenges in Pretraining

- **Fine-grained vision-language correlation:** needed for zero-shot dense prediction
- **Unified vision-language modeling:** one transformer for both image and text •
- **Multilingual pretraining:** remove language/cultural bias
- **Data-efficient learning:** less data, more supervision between pairs
- **Leveraging LLMs:** augment image-text pairs using rich linguistic knowledge

Challenges in Transfer Learning

- **Unsupervised transfer:** reduce overfitting to few-shot labels
- **Visual prompts/adapters:** complement text prompting
- **Test-time adaptation:** update prompts dynamically at inference
- **Transfer with LLMs:** generate task-specific prompts automatically

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

Section 6. VLM + Control System = VLAM

CV-2025

Recap: CLIP architecture and loss

DNNs Training Paradigms

VLM's Objectives (losses)

Transfer Learning of Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types of Multi

Unifies vision, language, and action via tokenization: $\text{Action} = \text{Decode}\left(\text{Transformer}\left(\underbrace{\mathcal{V}(I)}_{\text{Visual Tokens}}; \underbrace{\mathcal{T}(C)}_{\text{Language Tokens}}\right)\right)$

where: \mathcal{V} : vision tokenizer (ViT + Action Quantizer); \mathcal{T} : language tokenizer (PaLI-style); actions discretized as $\langle \text{cmd}, x, y, z, \theta \rangle$ tokens.

Key Innovations

- **Action Chunking:** Predicts action sequences autoregressively
- **Cross-Modal Attention:** Attention(Q_{action} , $K_{\text{vision} + \text{lang}}$, $V_{\text{vision} + \text{lang}}$)
- **Chain-of-Thought:** "Plan \rightarrow Verify \rightarrow Execute" token prediction

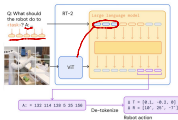


Figure: RT-2's unified architecture [Brohan et al., 2023]

Hands-on Coding with CLIP models (again)

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

CLIP + CLIPSeg = prerequisite for Action

- Get an Image (scene) + text (instruction from a human to a robot)
- Define a set of discrete robot skills (actions) and scene objects, and distractors
- use CLIPSeg for object segmentation and position detection
- use CLIP for skill prediction (VLAM Concept)

The code is available via the [link #1](#).

Blog reading

CV-2025

A.Korinaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

Vision Language Action Models (VLA) Overview: LeRobot Policies Demo

The advent of Generative AI, has fundamentally transformed robotic intelligence, enabling significant strides in how advanced humanoid robots perceive, reason and act in the physical world. This huge progress is primarily attributed in terms of decision making, thanks to LLM and VLMs generalization due to their large scale pre-training. Instead of relying on traditional complex policies which has to be carefully handcrafted for individual low level tasks for fine grained actions, VLA allows robotic control combining vision and language knowledge for better reasoning.

Paper reading

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Models

RT-2: Vision-Language-Action Models Transfer Web Knowledge to Robotic Control

We study how vision-language models trained on Internet-scale data can be incorporated directly into end-to-end robotic control to boost generalization and enable emergent semantic reasoning. Our goal is to enable a single end-to-end trained model to both learn to map robot observations to actions and enjoy the benefits of large-scale pretraining on language and vision-language data from the web. To this end, we propose to co-fine-tune state-of-the-art vision-language models on both robotic trajectory data and Internet-scale vision-language tasks, such as visual question answering. In contrast to other approaches, we propose a simple, general recipe to achieve this goal: in order to fit both natural language responses and robotic actions into the same format, we express the actions as text tokens and incorporate them directly into the training set of the model in the same way as natural language tokens. We refer to such category of models as vision-language-action models (VLA) and instantiate an example of such a model, which we call RT-2. Our extensive evaluation (6k evaluation trials) shows that our approach leads to performant robotic policies and enables RT-2 to obtain a range of emergent capabilities from Internet-scale training. This includes significantly improved generalization to novel objects, the ability to interpret commands not present in the robot training data (such as placing an object onto a particular number or icon), and the ability to perform rudimentary reasoning in response to user commands (such as picking up the smallest or largest object, or the one closest to another object). We further show that incorporating chain of thought reasoning allows RT-2 to perform multi-stage semantic reasoning, for example figuring out which object to pick up for use as an improvised hammer (a rock), or which type of drink is best suited for someone who is tired (an energy drink) [Brohan et al., 2023].

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Multi

Section 7. Other Types of Multimodalities

Other Types of Multimodalities

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Multi

The Five Senses Analogy

Sense	Data Modality	ML Example
Vision	Images/Video	CNNs, ViTs
Auditory	Audio/Waveforms	Spectrogram Transformers
Tactile	Pressure/Texture	Tactile Sensors in Robotics
Olfactory	Chemical Sensors	e-Nose Gas Detection
Gustatory	Molecular Data	Flavor Prediction Models

Emerging Sensor Fusion

- LiDAR+RGB: Autonomous vehicles
- IMU+Vision: Human pose estimation
- Spectrograms+Text: Audio captioning

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Multi-

Section 8. Meta-Learning

Multimodal Optimization Challenges

CV-2025

A.Kornaev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Multi-

Physics-Informed Neural Networks (PINNs)

$$\mathcal{L}_{\text{PINN}} = \underbrace{\lambda_d \|u_\theta(x_i) - u_i\|^2}_{\text{Data Loss}} + \underbrace{\lambda_p \|\mathcal{N}[u_\theta](x_j)\|^2}_{\text{Physics Loss}} + \underbrace{\lambda_r \|\theta\|^2}_{\text{Regularization}}$$

- Multi-objective: Data fitting + PDE residuals [Raissi et al., 2017]
- Loss landscape modality gaps cause training instabilities

Multi-Task Tradeoffs

- Pareto optimality in joint losses
- Gradient conflict quantification:

$$\cos(\nabla_\theta \mathcal{L}_i, \nabla_\theta \mathcal{L}_j) < 0$$

- Solution: Uncertainty weighting [Kendall et al., 2018]

Bibliography

CV-2025

A.Kor-naev,
K.Yakovlev

Recap: CLIP
architecture
and loss

DNNs
Training
Paradigms

VLM's
Objectives
(losses)

Transfer
Learning of
Large Models

Prospects

VLM +
Control
System =
VLAM

Other Types
of Multi-

Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, Pete Florence, Chuyuan Fu, Montse Gonzalez Arenas, Keerthana Gopalakrishnan, Kehang Han, Karol Hausman, Alexander Herzog, Jasmine Hsu, Brian Ichter, Alex Irpan, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Lisa Lee, Tsang-Wei Edward Lee, Sergey Levine, Yao Lu, Henryk Michalewski, Igor Mordatch, Karl Pertsch, Kanishka Rao, Krista Reymann, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Pierre Sermanet, Jaspiar Singh, Anikait Singh, Radu Soricut, Huong Tran, Vincent Vanhoucke, Quan Vuong, Ayzaan Wahid, Stefan Welker, Paul Wohlhart, Jialin Wu, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-2: Vision-language-action models transfer web knowledge to robotic control, 2023. URL <https://arxiv.org/abs/2307.15818>.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. URL <https://arxiv.org/abs/2106.09685>.

Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In **Proceedings of the IEEE conference on computer vision and pattern recognition**, pages 7482–7491, 2018.

Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning, 2021. URL <https://arxiv.org/abs/2004.11362>.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. URL <https://arxiv.org/abs/2103.00020>.

Maziar Raissi, Paris Perdikaris, and George Em Karniadakis. Physics informed deep learning (part i): Data-driven solutions of nonlinear partial differential equations, 2017. URL <https://arxiv.org/abs/1711.10561>.

Jingyi Zhang, Jiaxing Huang, Sheng Jin, and Shijian Lu. Vision-language models for vision tasks: A survey, 2024. URL <https://arxiv.org/abs/2304.00685>.