

Classificazione della

# Composizione Fotografica

delle immagini

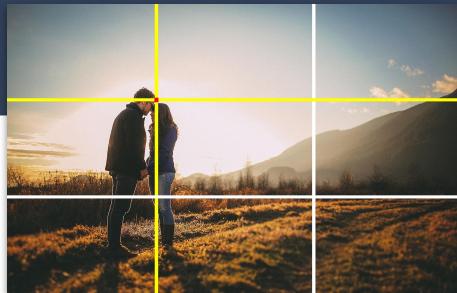
**Relatore:** *Dott. Luigi Celona*

**Correlatori:** *Prof. Raimondo Schettini, Prof. Gianluigi Ciocca*

**Candidato:** *Alessandro Locatelli, 879362*

# Composizione Fotografica

## Posizione dei soggetti



Regola dei terzi (RoT)



Centrato

## Leading Lines



Orizzontale



Incidenti

## Forme



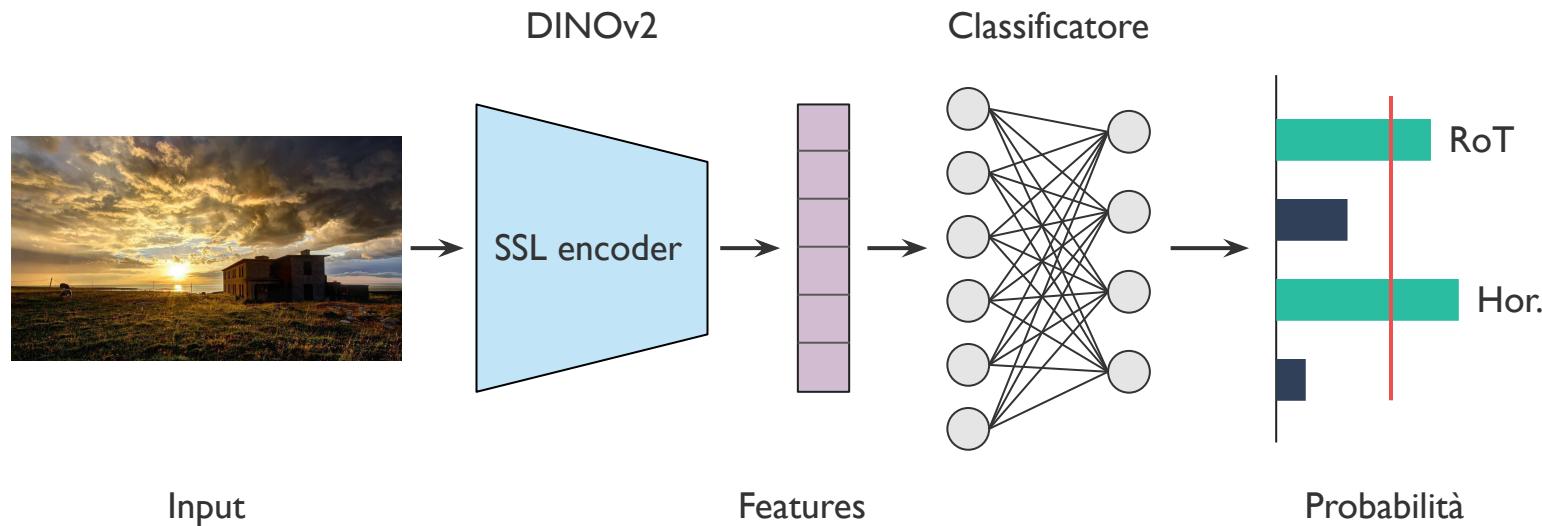
Triangolo



Pattern

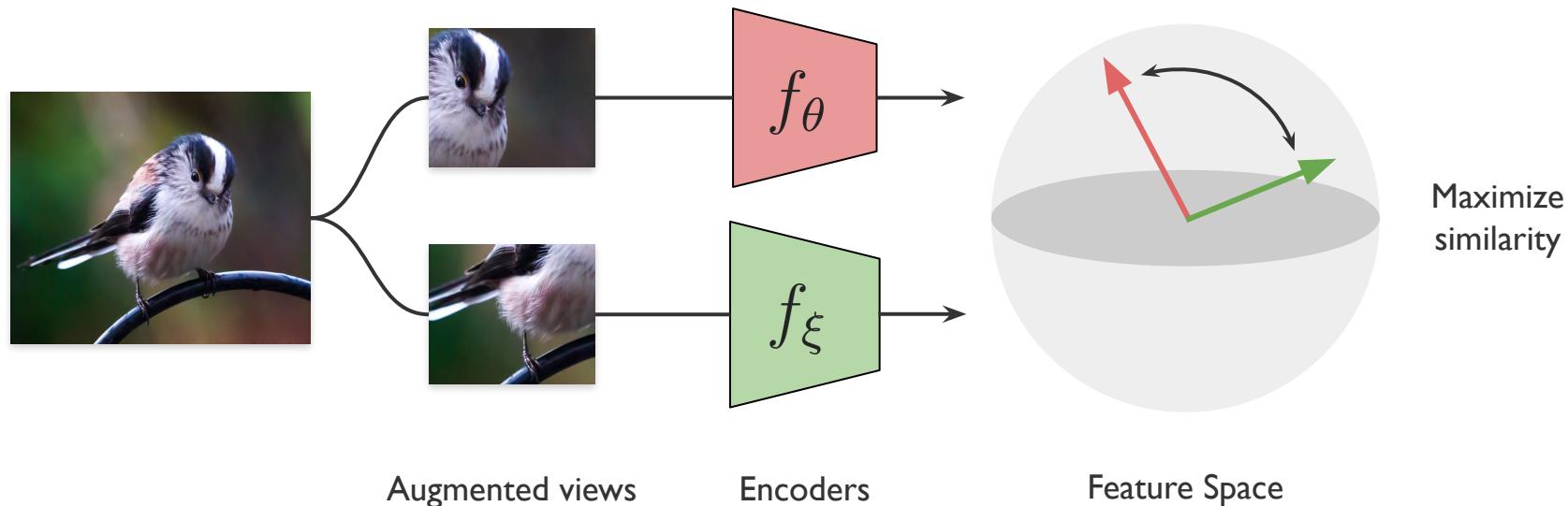
# Obiettivo

Valutare le capacità di generalizzazione di **DINOv2** nella classificazione della composizione



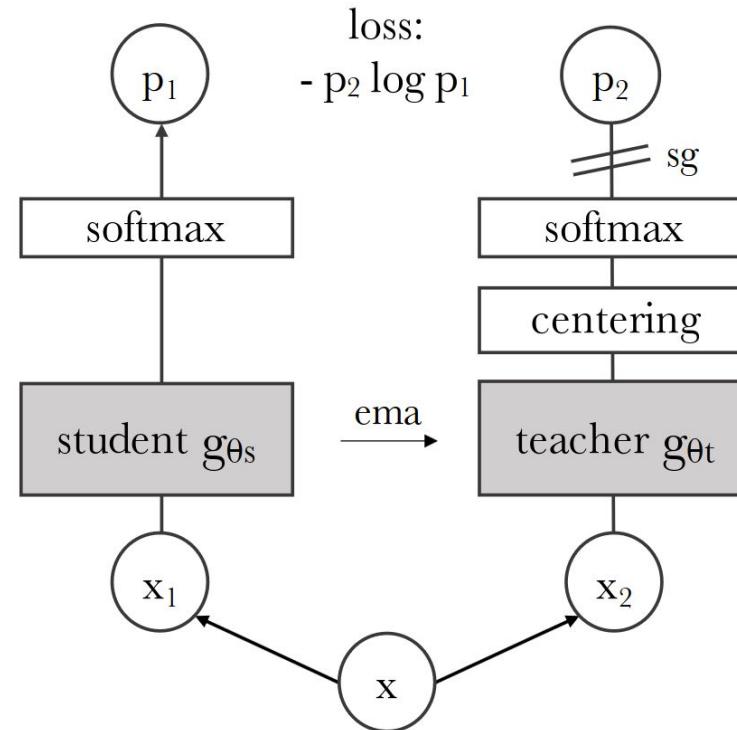
# Self-Supervised Learning (SSL)

Produrre rappresentazioni descrittive e intelligibili di dati non etichettati

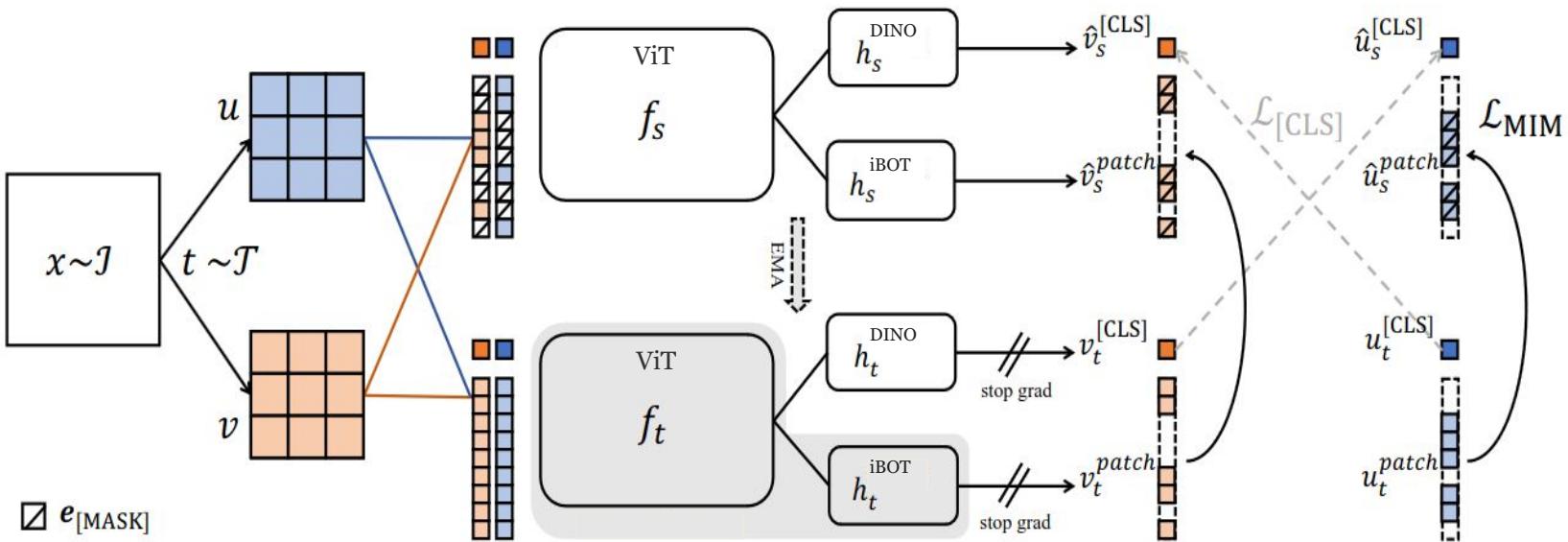


# DINO

Self-distillation with NO labels

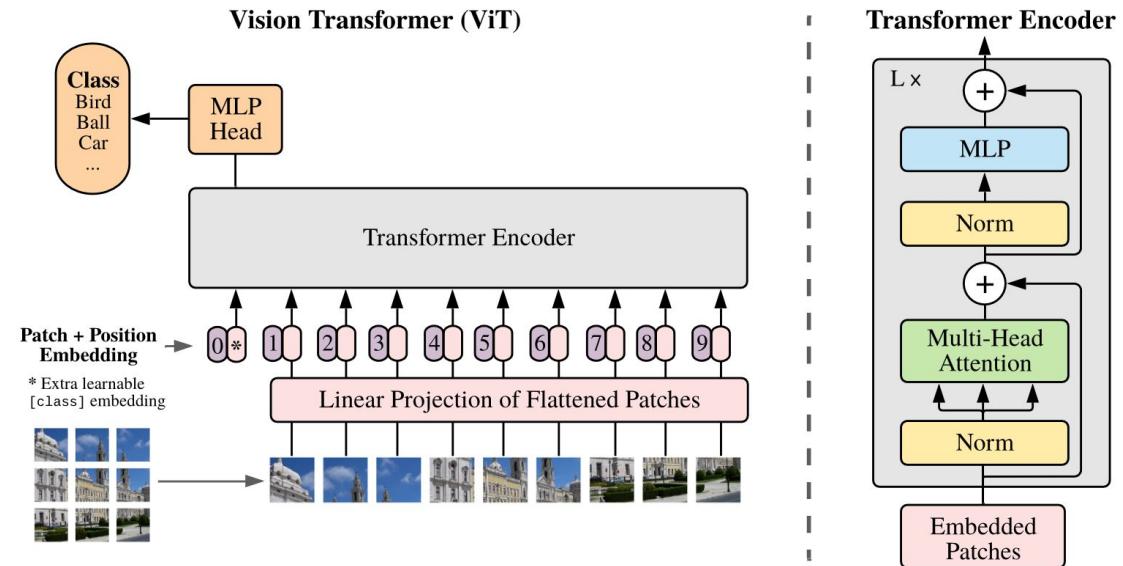


# DINOv2



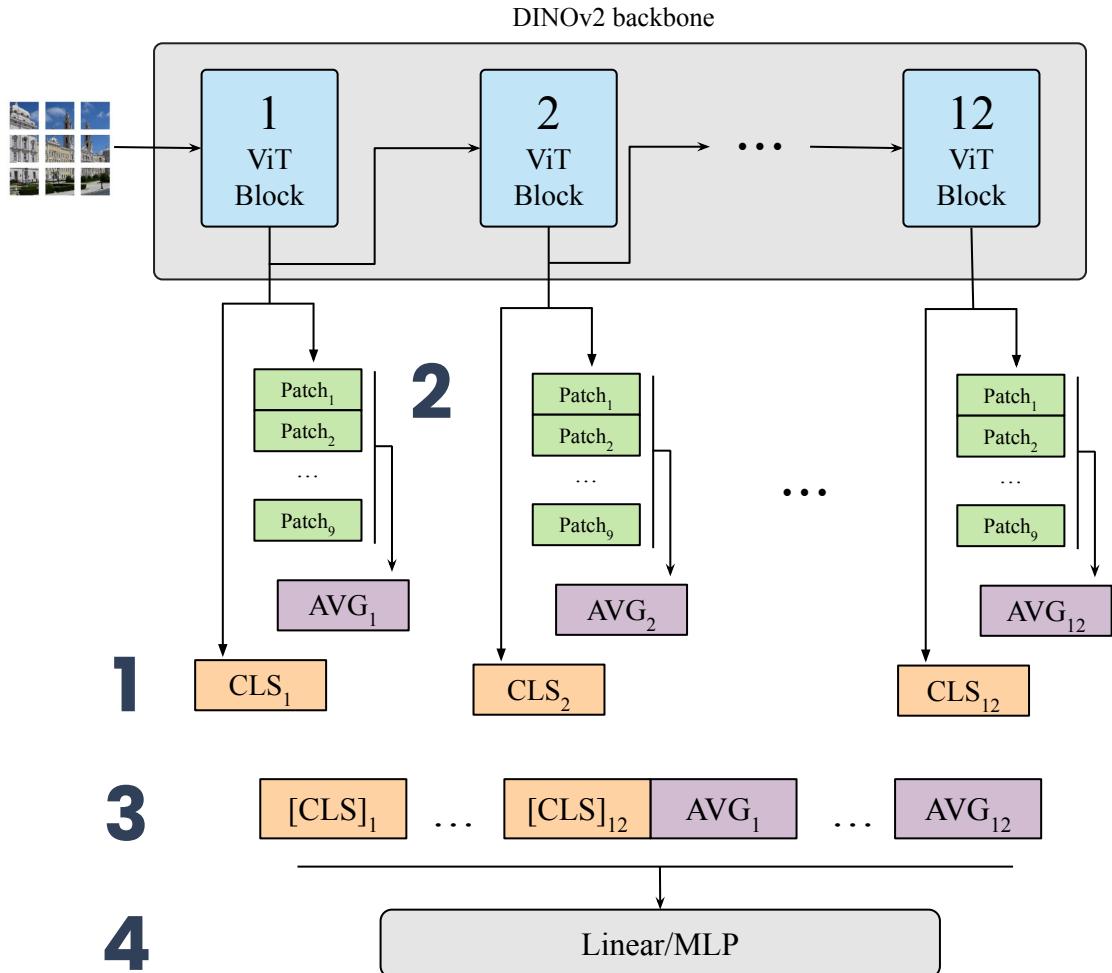
# Vision Transformers

- ▶ Dividere l'immagine in patches
- ▶ Aggiungere token [CLS]
- ▶ Encoder produce vettori di features per i patch



# Metodo

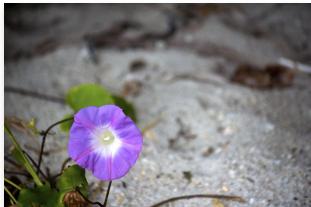
1. Estrazione [CLS]
2. Mediare features dei patch
3. Concatenazione
4. Classificatore



# KU-PCP

Dataset standard nella composizione fotografica

- ▶ 9 classi, multi-label
- ▶ 4.244 immagini etichettate
- ▶ ~20% doppia o tripla label



(a) *RoT*



(b) *Center*



(c) *Symmetric*



(d) *Diagonal*



(e) *Horizontal*



(f) *Curved*



(g) *Vertical*



(h) *Triangle*



(i) *Pattern*

# LODB

Nato per correggere i difetti di KU-PCP

- ▶ 17 classi, multi-label
- ▶ 6.029 immagini etichettate
- ▶ ~200 immagini con doppia label



(a) *O2DiaR.*



(b) *RoT.R*



(c) *O3Tri.*



(d) *O3Li.*



(f) *O2DiaL.*



(g) *RoT.L*



(h) *O2Hor.*



(i) *Oline*

# Risultati

## I. Linear Layer

## 2. MLP (1 hidden, ReLU)

BCE (binary cross-entropy)

Batch size 1024

Encoder ViT-S/14

K-fold cross-validation (5)

NVIDIA Tesla K80

1  
2

Le features di DINOv2 non sono adatte al task

Dataset	Epoche	lr	Acc. (Any)	Acc. (All)	MP	mAP	MF1
KUPCP	300	0.003	0.83	0.63	0.84	0.81	0.79
LODB	200	0.003	0.56	0.52	0.66	0.65	0.49
KUPCP	200	0.005	0.84	0.54	0.83	0.76	0.77
LODB	300	0.002	0.51	0.50	0.68	0.67	0.44

4 [CLS] + 4 AVG

Confronto con altri metodi in letteratura su KU-PCP:

Metodo	Acc. (Any) %
FT_CNN [4]	88.8
RSTN [6]	90.9
DINOv2 + linear	82.9

# Per classe

Scarsa varianza semantica  
intraclasse in KU-PCP  
influenza i risultati

Classe	Precision	Recall	F1-score
<i>RoT</i>	0.6866	0.7480	0.7160
<i>Vertical</i>	0.8571	0.6364	0.7304
<i>Horizontal</i>	0.8390	0.8190	0.8289
<i>Diagonal</i>	0.7372	0.5344	0.6196
<i>Curved</i>	0.7722	0.6630	0.7135
<i>Triangle</i>	0.8487	0.7544	0.7988
<i>Center</i>	0.9135	0.7814	0.8423
<i>Symmetric</i>	0.9770	0.8586	0.9140
<i>Pattern</i>	0.9683	0.9683	0.9683

Metriche per classe in KU-PCP, linear layer

# Varianza semantica intraclasse

L'ipotesi è che l'informazione semantica influenzi le predizioni del classificatore



*Symmetric:* scarsa varianza, metriche migliori

*RoT:* alta varianza, metriche peggiori

# Per classe

**Sbilanciamento** di LODB  
peggiora ulteriormente le  
prestazioni

Classe	Precision	Recall	F1-score
<i>DiaL.</i>	0.7333	0.3438	0.4681
<i>O2DiaR.</i>	0.2381	0.0847	0.1250
<i>O2Hor.</i>	0.4545	0.1923	0.2703
<i>Hor.</i>	0.9630	0.8254	0.8889
<i>Cent.</i>	0.8129	0.7079	0.7568
<i>Pat.</i>	0.9091	0.7843	0.8421
<i>Radi.</i>	0.5294	0.4286	0.4737
<i>RoT.R</i>	0.8333	0.5263	0.6452
<i>Ver.</i>	0.9615	0.8242	0.8876
<i>O3Li.</i>	0.4444	0.1143	0.1818
<i>RoT.L</i>	0.7333	0.2500	0.3729
<i>DiaR.</i>	0.7895	0.3261	0.4615
<i>DiaX.</i>	0.9308	0.8605	0.8943
<i>O3Tri.</i>	0.3333	0.0417	0.0741
<i>O2DiaL.</i>	0.1429	0.0159	0.0286
<i>Oline</i>	0.6250	0.2174	0.3226
<i>Tri.</i>	0.8182	0.7079	0.7590

Metriche per classe su LODB con linear layer

Nome	Descrizione	Numero
<i>Cent.</i>	Soggetto centrato nell'immagine	947
<i>RoT.L</i>	Regola dei terzi, soggetto allineato con la griglia a sinistra	214
<i>RoT.R</i>	Regola dei terzi, soggetto allineato con la griglia a destra	275
<i>O2Dia.L</i>	Due oggetti disposti diagonalmente, dall'angolo in alto a sinistra	301
<i>O2Dia.R</i>	Due oggetti disposti diagonalmente, dall'angolo in alto a destra	293
<i>O2Hor.</i>	Due oggetti posizionati su una linea orizzontale	386
<i>O3Li.</i>	Tre oggetti posizionati sulla stessa linea	157
<i>O3Tri.</i>	Composizione triangolare formata da 3 oggetti	115
<i>Oline.</i>	Più di tre oggetti posizionati sulla stessa linea	110
<i>Pat.</i>	L'immagine contiene dei patterns	255
<i>DiaL.</i>	Composizione a diagonale, dall'angolo in alto a sinistra	446
<i>DiaR.</i>	Composizione a diagonale, dall'angolo in alto a destra	236
<i>Hor.</i>	Struttura orizzontale	574
<i>Tri.</i>	Struttura triangolare	451
<i>Ver.</i>	Struttura verticale	455
<i>Radi.</i>	Struttura radiale	187
<i>DiaX.</i>	Struttura comprendente entrambe le diagonali	861

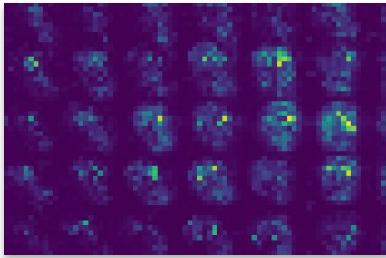
Composizione del dataset LODB

# Sviluppi Futuri

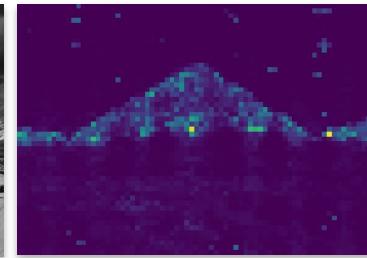
Usare una CNN sulle mappe di self-attention del ViT



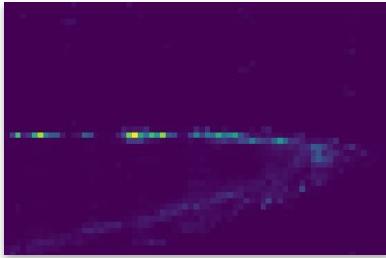
(a) *Pattern*



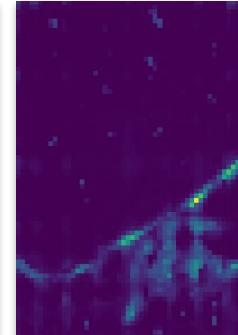
(b) *Triangle*



(c) *Curved, Horizontal*



(d) *RoT, Diagonal*

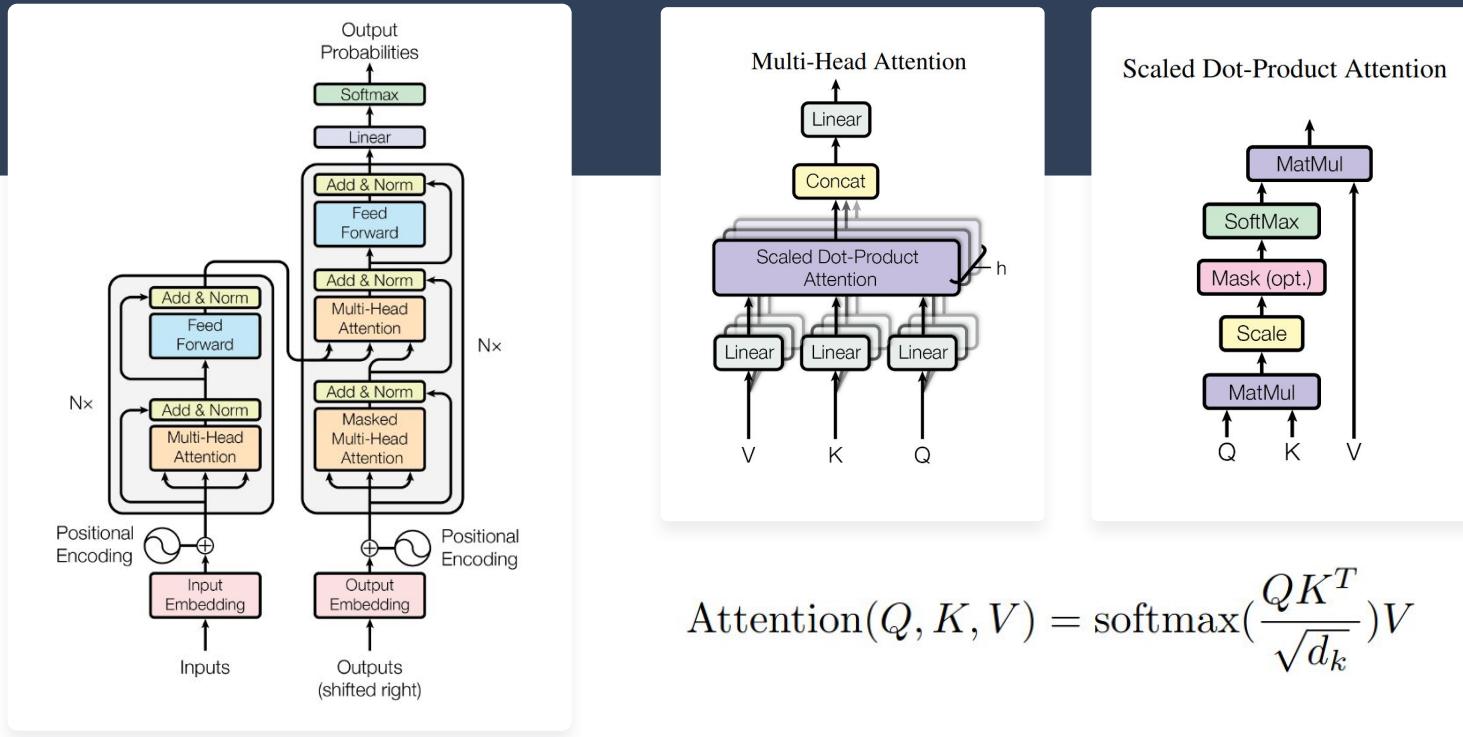


# Riferimenti

- [1] Mathilde Caron et al.“Emerging properties in self-supervised vision transformers”. In: *Proceedings of the IEEE/CVF international conference on computer vision*. 2021, pp. 9650–966
- [2] Maxime Oquab et al.“DINOv2: Learning robust visual features without supervision”. In: *arXiv preprint arXiv:2304.07193* (2023)
- [3] Alexey Dosovitskiy et al.“An image is worth 16x16 words:Transformers for image recognition at scale”. In: *arXiv preprint arXiv:2010.11929* (2020)
- [4] Jun-Tae Lee et al.“Photographic Composition Classification and Dominant Geometric Element Detection for Outdoor Scenes”. In: *Journal of Visual Communication and Image Representation* 55 (mag. 2018)
- [5] Zhaoran Zhao et al.“Self-supervised Photographic Image Layout Representation Learning”. In: *arXiv preprint arXiv:2403.03740* (2024)
- [6] Yaoting Wang et al.“Spatial-invariant convolutional neural network for photographic composition prediction and automatic correction”. In: *Journal of Visual Communication and Image Representation* 90 (2023), p. 103751.

# Extra

# Transformer



# LOSS

$$\text{BCEwithLogits}(x, y) = BCE(\sigma(x), y)$$

$y$  = classificazione reale (0, 1)

$x$  = logit

Sigmoide:

$$p = \sigma(x) = \frac{1}{1 + e^{-x}}$$

Binary cross-entropy:

$$BCE(p, y) = -(y \log(p) + (1 - y) \log(1 - p))$$

# LOSS (MLP)

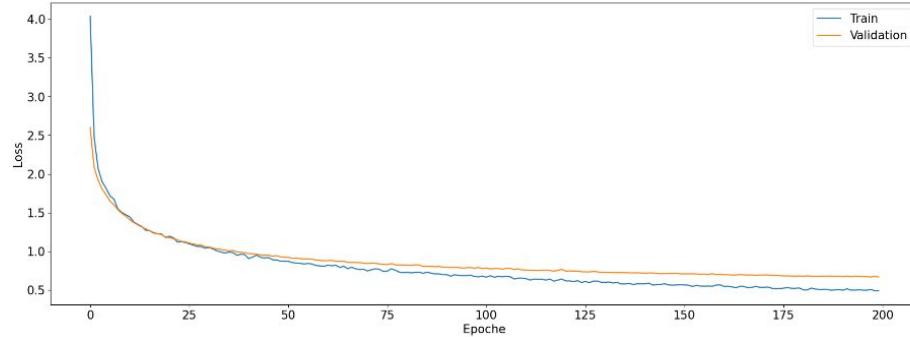
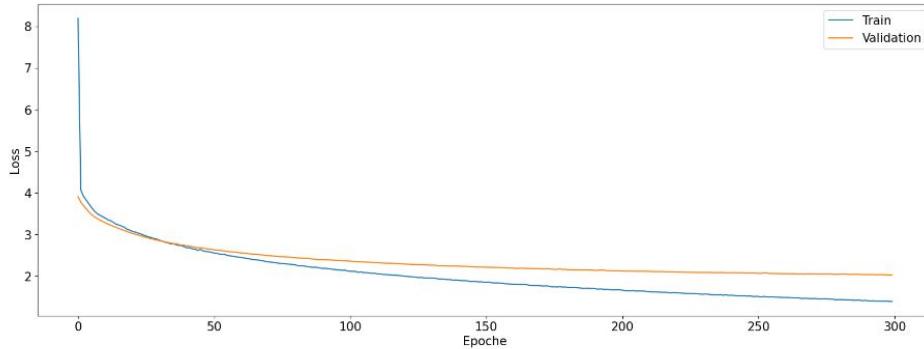


Figura 6.3: Grafico della loss in validation (quinto fold) su KU-PCP e MLP per la configurazione 4 token CLS e 4 AVG patch, 200 epochi, learning rate 0.005, che viene usata in fase di test. Non ci sono segnali di overfitting.

# Loss (Linear)

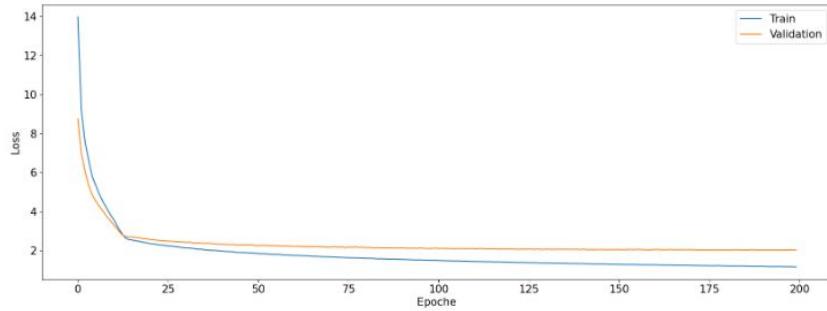


Figura 6.1: Grafico della loss in validation (quinto fold) su LODB e linear layer per la configurazione 4 token CLS e 4 AVG patch, 200 epochi, learning rate 0.003, che viene usata in fase di test. Non ci sono segnali di overfitting.

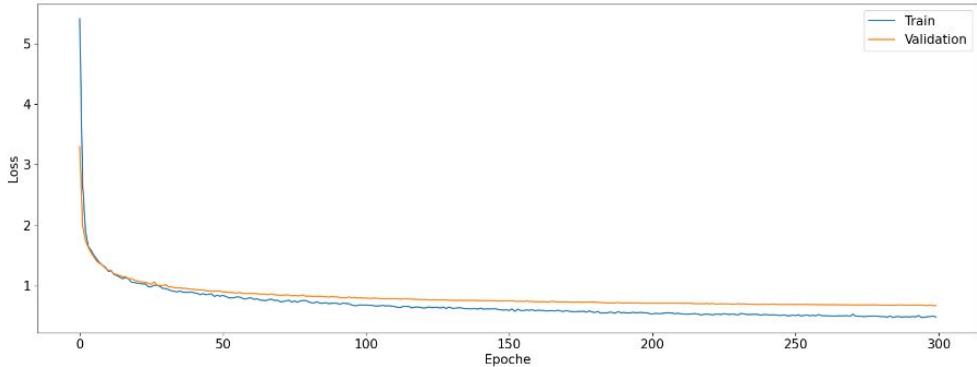


Figura 6.2: Grafico della loss in validation (quinto fold) su KU-PCP e linear layer per la configurazione con 4 token CLS e 4 AVG patch, 300 epochi e learning rate 0.003, che viene usata in fase di test. La discesa della loss non mostra segni di overfitting.

# Metriche

$$Precision = \frac{TP}{TP + FP}$$

$$F1 = 2 \cdot \frac{precision \cdot recall}{precision + recall}$$

$$Recall = \frac{TP}{TP + FN}$$

$$AP = \sum_n (R_n - R_{n-1}) P_n$$

R = recall, P = precision, n-esimo threshold

# Matrici di confusione

		$\sim DiaX$	$DiaX$			$\sim O3Tri$	$O3Tri$
$\sim DiaX$	$\sim DiaX$	805	<b>13</b>	$\sim O3Tri$	937	<b>0</b>	
$DiaX$	$DiaX$	19	<b>111</b>	$O3Tri$	11	<b>0</b>	

Tabella 6.3: Matrici di confusione del quinto fold per  $DiaX$  e  $O3Tri$  nella configurazione 8 CLS e lr = 0.002 e nel fold 5.

# Features

[CLS] e AVG usati singolarmente

Features	lr	Acc. (Any)	Acc. (All)	MP	mAP
CLS <sub>12</sub>	0.05	<b>0.52</b>	0.42	0.46	0.49
CLS <sub>9</sub> ... CLS <sub>12</sub>	0.006	0.45	0.44	0.65	0.58
CLS <sub>5</sub> ... CLS <sub>12</sub>	0.002	0.38	0.37	0.58	0.56
CLS <sub>1</sub> ... CLS <sub>12</sub>	0.0005	0.20	0.19	0.33	0.47

Tabella 6.1: Un confronto fra le prestazioni del layer lineare utilizzando diverse combinazioni di layer di token [CLS]

Features	lr	Acc. (Any)	Acc. (All)	MP	mAP
AVG <sub>12</sub>	0.06	<b>0.52</b>	<b>0.49</b>	<b>0.61</b>	<b>0.59</b>
AVG <sub>9</sub> ... AVG <sub>12</sub>	0.005	0.43	0.42	0.60	0.59
AVG <sub>5</sub> ... AVG <sub>12</sub>	0.002	0.41	0.41	0.55	0.56
AVG <sub>1</sub> ... AVG <sub>12</sub>	0.001	0.35	0.34	0.51	0.54

Tabella 6.4: Un confronto fra le prestazioni del layer lineare utilizzando diverse combinazioni di layer di AVG patch

# Features

[CLS] e AVG usati in coppia

Features	lr	Acc. (Any)	Acc. (All)	MP	mAP
CLS <sub>12</sub> AVG <sub>12</sub>	0.01	0.49	0.46	0.58	0.56
CLS <sub>9</sub> ... CLS <sub>12</sub> AVG <sub>9</sub> ... AVG <sub>12</sub>	0.005	<b>0.51</b>	<b>0.49</b>	<b>0.64</b>	<b>0.63</b>
CLS <sub>5</sub> ... CLS <sub>12</sub> AVG <sub>5</sub> ... AVG <sub>12</sub>	0.001	0.42	0.41	0.63	0.60
CLS <sub>1</sub> ... CLS <sub>4</sub> CLS <sub>9</sub> ... CLS <sub>12</sub> AVG <sub>1</sub> ... AVG <sub>4</sub> AVG <sub>9</sub> ... AVG <sub>12</sub>	0.001	0.39	0.38	0.56	0.58

Tabella 6.5: Un confronto fra le prestazioni del layer lineare utilizzando diverse combinazioni di layer di token [CLS] e AVG patch.

# Predizioni

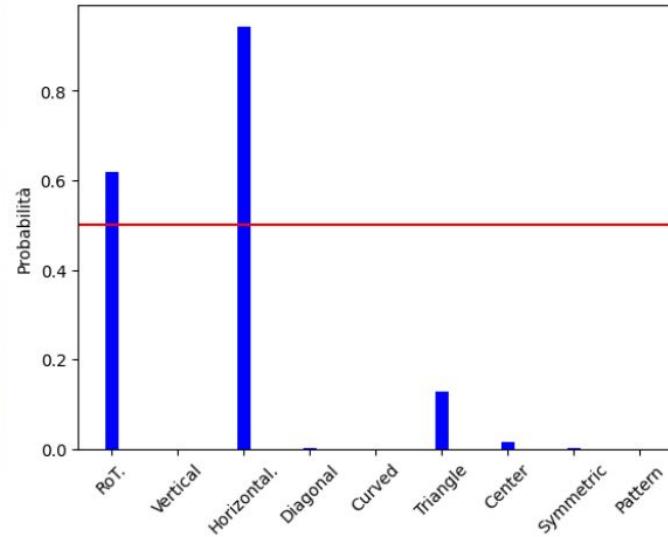


Figura 6.5: Esempio di predizione corretta della rete con linear classifier su KU-PCP. Le classi di ground truth sono *RoT.* e *Horizontal.*. In rosso è segnata la soglia di threshold (posta a 0.5) utilizzata per il calcolo delle metriche.

# Predizioni

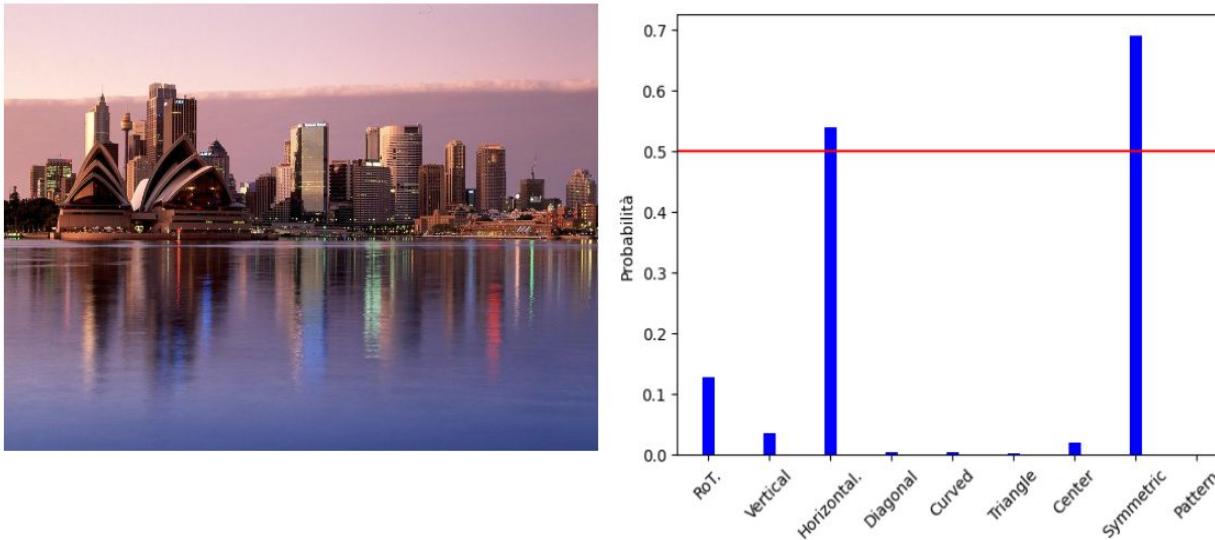


Figura 6.6: Predizione con classificatore lineare su KU-PCP. Le classi di ground truth sono *Vertical* e *Symmetric*, mentre la nostra rete predice *Horizontal* e *Symmetric*.

# Predizioni

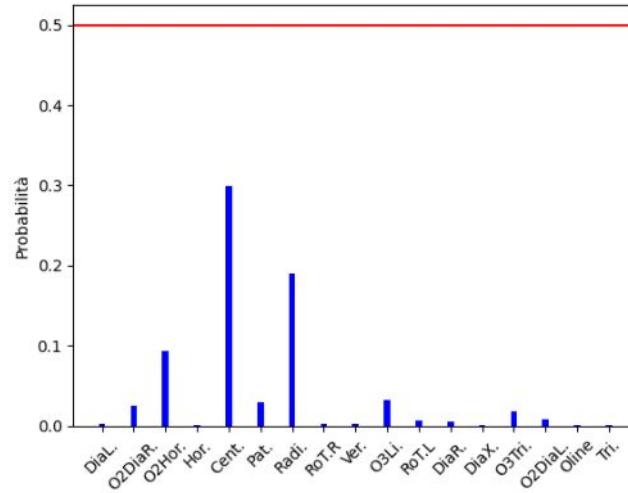


Figura 6.8: Predizione con classificatore lineare su LODB, ground truth *O3Tri..*

# Predizioni

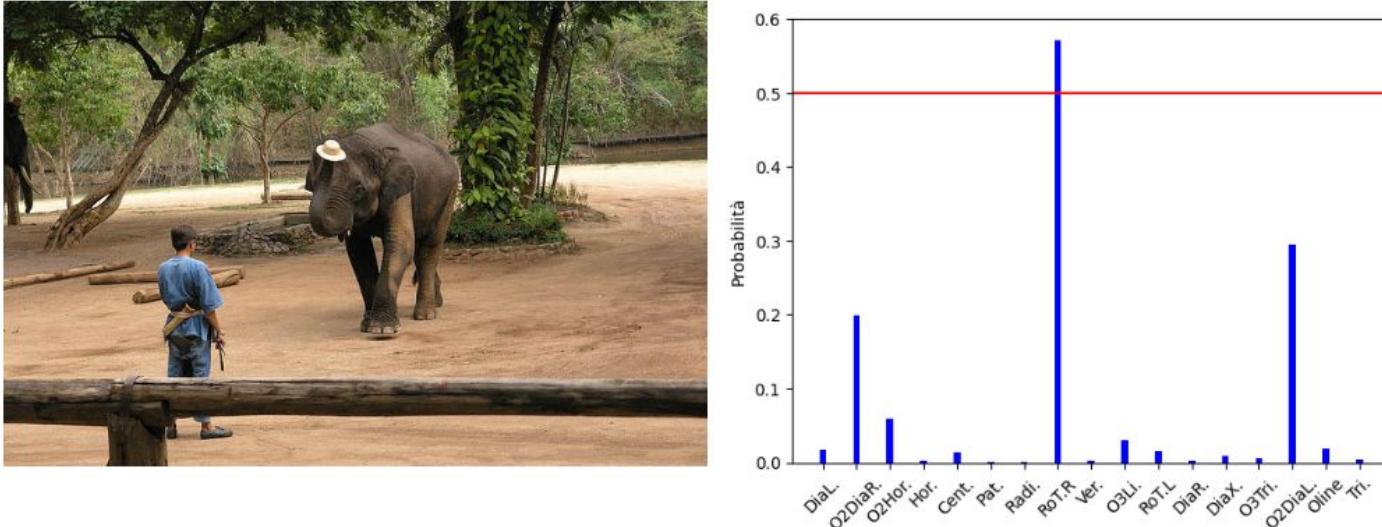


Figura 6.9: Predizione con classificatore lineare su LODB, ground truth  $O2DiaR$ .

# Predizioni

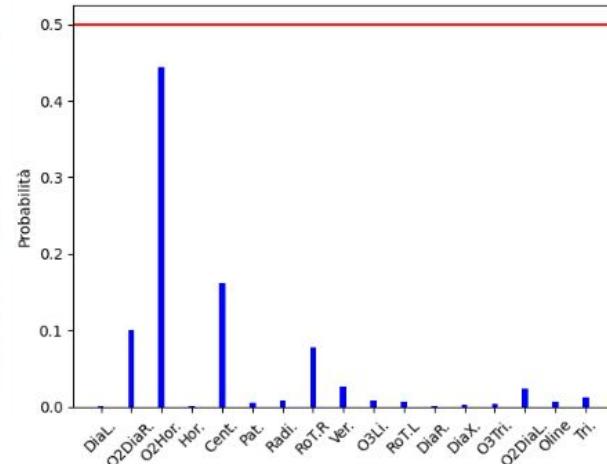


Figura 6.10: Predizione con classificatore lineare su LODB, ground truth *O2Hor..*. Anche su immagini la cui composizione risulta essere piuttosto chiara all'occhio umano, il modello fatica a riconoscerla e spalma le sue predizioni su diverse classi, facendo fatica a prendere una decisione.

# Misclassification



(a) Vertical



(b) O2DiaL.



(c) O3Tri.

Figura 7.4: Esempi di misclassification/forzatura della composizione. Sono indicate le etichette che LODB assegna a ciascuna.

# Augmentations



(a) Original



(b) Crop and resize



(c) Crop, resize (and flip)



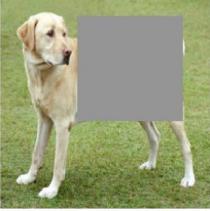
(d) Color distort. (drop)



(e) Color distort. (jitter)



(f) Rotate  $\{90^\circ, 180^\circ, 270^\circ\}$



(g) Cutout



(h) Gaussian noise



(i) Gaussian blur



(j) Sobel filtering

# Momentum Encoder

Aggiornamento dei pesi della rete Teacher con Exponential Moving Average dei pesi dello Student

$$\xi \leftarrow \tau \xi + (1 - \tau) \theta$$

- ▶  $\xi$  pesi del Teacher
- ▶  $\theta$  pesi dello Student
- ▶  $\tau$  momentum coefficient in  $[0, 1]$

inizia a 0.996, aumenta a 1 durante training

Stop Gradient (SG) evita che il Teacher venga aggiornato tramite backpropagation