

# **Aslam Equation Explanations**





21BAD034-21BAD051.pptx

Presentation



Computer\_Vision\_and\_Deep\_Learning\_to\_Manage\_Safety\_in\_Construction\_Matching\_Images\_of\_Unsafe\_Behavior\_and\_Sem



YOLO-ESCA A High-

Performance\_Safety\_Helmet\_Standard\_Wearing\_Behavior\_Detection\_Model\_Based\_on\_Improved\_YOLOv5.pdf

give me equations explained in the ppt with respect to the Papers that have given with usage



The equations and their usage as presented in the uploaded documents are summarized below:

# **Equations and Key Concepts in the PowerPoint and Related Papers:**

## **EIOU Loss (Efficient Intersection over Union):**

#### **Equation:**

1. EIOU Loss  $L_{EIOU}$  is given by:

$$L_{EIOU} = L_{IOU} + L_{dis} + L_{asn}$$

where:

$$egin{array}{l} \circ \ L_{IOU} = 1 - ext{IOU} + rac{
ho^2(b,b_{gt})}{w_c^2 + h_c^2} \ \circ \ L_{dis} = rac{
ho^2(w,w_{gt})}{w_c^2} \ \circ \ L_{asp} = rac{
ho^2(h,h_{gt})}{h_c^2} \end{array}$$

$$\circ L_{dis} = \frac{\rho^2(w, w_{gt})}{w^2}$$

$$\circ~L_{asp}=rac{
ho^2(h,h_{gt})}{h^2}$$

This loss function improves bounding box prediction in YOLO-ESCA by considering overlap, distance, and aspect ratio between the predicted and ground truth boxes.

Usage: Used in YOLO-ESCA to replace CIOU for better performance on small and dense objects (e.g., helmets).

#### **Attention Mechanisms (CBAM):**

**Equation:** CBAM combines:

• Channel Attention Module (CAM):

$$M_c = \sigma(W_1(W_0(F_{ava})) + W_1(W_0(F_{max})))$$

• Spatial Attention Module (SAM):

$$M_s = \sigma(Conv([F_{avg}; F_{max}]))$$

**Usage**: Enhances feature focus on small objects like helmet straps. CBAM is integrated into YOLO-ESCA's backbone for better target identification.

## **Stacked Cross Attention Network (SCAN):**

## **Equations:**

1. Cosine Similarity:

$$s_{ij} = rac{v_i^T e_j}{\|v_i\| \|e_j\|}$$

2. Attention Weights:

$$a_{ij} = rac{\exp(\lambda s_{ij})}{\sum_{k=1}^n \exp(\lambda s_{ik})}$$

3. Feature Similarity:

$$R(v_i, a_t) = rac{v_i^T a_t}{\|v_i\| \|a_t\|}$$

**Usage**: Used for matching image regions with safety rule features in semantic similarity computations for detecting unsafe behaviors on construction sites.

#### **Triplet Loss:**

**Equation**:

$$L_{triplet} = \max(0, lpha + d(a, p) - d(a, n))$$

Where:

- d(a, p): Distance between the anchor and positive example.
- d(a, n): Distance between the anchor and negative example.
- $\alpha$ : Margin.

**Usage**: Optimizes SCAN by ensuring positive image-rule pairs have higher similarity than negative pairs.

#### **Precision and Recall:**

## **Equations**:

1. Precision:

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

#### 2. Recall:

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

**Usage**: Used as evaluation metrics for models detecting unsafe behaviors and helmet usage.

Let me know if you need any specific section elaborated or equations implemented!



explain me every single equation in the powerpoint and how the input goes and comes out of the equation should be given with example



Here's a detailed explanation of the equations from the PowerPoint, along with their inputs, processing, and outputs. I will describe each equation step-by-step with examples to make the concepts clear.

# 1. EIOU Loss (Efficient Intersection Over Union)

## **Equation:**

$$L_{EIOU} = L_{IOU} + L_{dis} + L_{asp}$$

Where:

$$egin{aligned} ullet & L_{IOU} = 1 - IOU + rac{
ho^2(b,b_{gt})}{w_c^2 + h_c^2} \ ullet & L_{dis} = rac{
ho^2(w,w_{gt})}{w_c^2} \ ullet & L_{asp} = rac{
ho^2(h,h_{gt})}{h_c^2} \end{aligned}$$

$$ullet \ L_{dis} = rac{
ho^2(w,w_{gt})}{w_s^2}$$

$$ullet \; L_{asp} = rac{
ho^2(h,h_{gt})}{h_c^2}$$

### **Explanation:**

# • Inputs:

- **Predicted Box** (b): Defined by its center coordinates, width (w), and height (h).
- $\circ$  **Ground Truth Box** ( $b_{qt}$ ): Actual bounding box with center coordinates, width ( $w_{qt}$ ), and height
- $\circ$  **Diagonal of enclosing box** ( $w_c, h_c$ ): The dimensions of the smallest box enclosing b and  $b_{at}$ .

## • Outputs:

A scalar value representing the loss.

#### **Process:**

#### 1. Calculate IOU:

$$IOU = rac{ ext{Area of Overlap}}{ ext{Area of Union}}$$

Example: If the overlap area is 25 and the union area is 100, IOU=0.25.



- 2. **Center Distance** ( $ho^2(b,b_{qt})$ ): Compute the squared Euclidean distance between the centers of band  $b_{at}$ .
- 3. **Aspect Loss** ( $L_{asp}$ ) and **Distance Loss** ( $L_{dis}$ ): Measure the difference in width and height.

## **Example:**

- Predicted box: Center (10, 10), w = 4, h = 6.
- Ground truth box: Center (12,12),  $w_{at}=5$ ,  $h_{at}=6$ .
- Enclosing box: Width = 6, Height = 8.

## Compute each term:

- $ullet \ L_{IOU} = 1 0.25 + rac{((10-12)^2 + (\overline{10-12})^2)}{6^2 + 8^2} = 1.25.$
- $ullet L_{dis} = rac{(4-5)^2}{6^2} = 0.027. \ ullet L_{asp} = rac{(6-6)^2}{8^2} = 0. \$

Final Loss:  $L_{EIOU} = 1.25 + 0.027 + 0 = 1.277$ .

# 2. CBAM (Convolutional Block Attention Module)

### **Equation:**

1. Channel Attention:

$$M_c = \sigma(W_1(W_0(F_{avg})) + W_1(W_0(F_{max})))$$

2. Spatial Attention:

$$M_s = \sigma(Conv([F_{avq}; F_{max}]))$$

# **Explanation:**

- Inputs:
  - $\circ~$  Feature map F of size C imes H imes W , where C is the number of channels, H is height, and Wis width.
- Outputs:
  - Enhanced feature map with spatial and channel attention.

#### **Process:**

- 1. Channel Attention:
  - $\circ~$  Compute global average pooling ( $F_{avq}$ ) and max pooling ( $F_{max}$ ) across spatial dimensions.
  - Pass these through a shared fully connected network.
  - $\circ$  Combine and apply sigmoid activation to get channel weights ( $M_c$ ).
- 2. Spatial Attention:
  - $\circ$  Pool along the channel axis to generate  $F_{avg}$  and  $F_{max}$ .
  - Concatenate and convolve the result.

- $\circ$  Apply sigmoid to produce spatial weights ( $M_s$ ).
- 3. Multiply  $M_c$  and  $M_s$  with F to enhance the attention-weighted feature map.

# **Example:**

Feature map F=3 imes 5 imes 5 (3 channels, 5x5 grid):

- $F_{avg}$  = [Average pixel value for each channel].
- After  $M_c$  computation, weights = [1.2, 0.8, 1.0].
- Multiply these weights with the respective channels.
- ullet Similarly, compute  $M_s$  and refine the spatial features.

# 3. SCAN (Stacked Cross Attention Network)

## **Equations:**

1. Cosine Similarity:

$$s_{ij} = rac{v_i^T e_j}{\|v_i\| \|e_j\|}$$

2. Attention Weight:

$$a_{ij} = rac{\exp(\lambda s_{ij})}{\sum_{k=1}^n \exp(\lambda s_{ik})}$$

3. Feature Fusion:

$$R(v_i, a_t) = rac{v_i^T a_t}{\|v_i\| \|a_t\|}$$

## **Explanation:**

- Inputs:
  - $\circ$  Image region features  $v_i$ .
  - Safety rule text features  $e_i$ .
- Outputs:
  - $\circ$  Similarity score  $s_{ij}$ .

#### **Process:**

- 1. Compute cosine similarity between each image region and each word.
- 2. Normalize scores to obtain attention weights.
- 3. Fuse features using weighted sums to compute the final similarity.

#### **Example:**

- ullet Image has 3 regions:  $v_1 = [1,0,2]$ ,  $v_2 = [0,1,3]$ ,  $v_3 = [2,2,0]$ .
- ullet Text has 2 features:  $e_1=[1,1,1]$  ,  $e_2=[0,1,0]$  .

• Cosine similarity  $s_{ij}$  is computed for each pair. For  $v_1$  and  $e_1$ :

$$s_{11} = rac{1 imes 1 + 0 imes 1 + 2 imes 1}{\sqrt{1^2 + 0^2 + 2^2} imes \sqrt{1^2 + 1^2 + 1^2}} = 0.63$$

Normalize and fuse features for final similarity.

# 4. Triplet Loss

### **Equation:**

$$L_{triplet} = \max(0, \alpha + d(a, p) - d(a, n))$$

### **Explanation:**

- Inputs:
  - $\circ$  Anchor a, positive p, and negative n embeddings.
  - Distance function d (e.g., Euclidean distance).
  - $\circ$  Margin  $\alpha$ .
- Outputs:
  - Scalar loss value.

#### **Process:**

- 1. Calculate distances: d(a, p) and d(a, n).
- 2. Enforce  $d(a,p) + \alpha < d(a,n)$ .

## **Example:**

- Anchor: a=[1,0], Positive: p=[2,1], Negative: n=[0,3].
- Compute distances:

$$d(a,p) = \sqrt{(1-2)^2 + (0-1)^2} = 1.41$$

$$d(a,n) = \sqrt{(1-0)^2 + (0-3)^2} = 3.16$$

• Loss:

$$L_{triplet} = \max(0, 0.2 + 1.41 - 3.16) = 0$$

Let me know if you'd like further clarifications or more examples!



what is the difference between IOU and EIOU and give its advantages over it and disadvantages



## Difference Between IOU and EIOU

## 1. IOU (Intersection Over Union)

IOU measures the overlap between the predicted bounding box and the ground truth bounding box. It is a ratio defined as:

$$IOU = \frac{Area \text{ of Overlap}}{Area \text{ of Union}}$$

- Inputs:
  - $\circ$  Two bounding boxes: Predicted box A and Ground Truth box B.
- Outputs:
  - A value between 0 and 1, where 1 indicates perfect overlap.
- Key Points:
  - Measures overlap but does not consider the distance between the box centers.
  - Fails when the boxes do not intersect (IOU is 0, even if the boxes are close).

#### 2. EIOU (Efficient Intersection Over Union)

EIOU extends IOU by considering:

- 1. **Center Distance**: Penalizes the distance between the centers of the predicted and ground truth boxes.
- 2. Aspect Ratio Difference: Penalizes differences in width and height.

EIOU is defined as:

$$L_{EIOU} = L_{IOU} + L_{dis} + L_{asp}$$

Where:

- ullet  $L_{IOU}=1-\mathrm{IOU}+rac{
  ho^2(b,b_{gt})}{w_c^2+h_c^2}$
- $L_{dis}$  accounts for the width difference.
- ullet  $L_{asp}$  accounts for the height difference.

# **Advantages of EIOU Over IOU**

#### 1. Better Localization:

- **EIOU** accounts for the distance between the centers of the predicted and ground truth boxes, improving convergence during training.
- **IOU** does not consider center distance, which leads to poor predictions when the boxes are close but do not overlap.

#### 2. Aspect Ratio Sensitivity:

- **EIOU** penalizes differences in width and height, leading to better alignment between the predicted box and the ground truth.
- **IOU** ignores shape differences, only focusing on the area of overlap.



## 3. Improved Optimization:

- **EIOU** provides a smoother loss gradient, enabling faster and more stable convergence during training.
- **IOU** can have zero gradients when boxes do not overlap, causing optimization to stall.

#### 4. Handles Non-Overlapping Boxes:

• **EIOU** evaluates distance and aspect ratio even when IOU is zero, ensuring continuous optimization.

# **Disadvantages of EIOU**

## 1. Increased Complexity:

• **EIOU** involves additional computations for center distances and aspect ratios, making it computationally heavier than IOU.

#### 2. Not Suitable for Simple Tasks:

• For tasks where bounding boxes are usually well-aligned, the additional complexity of EIOU may not provide significant benefits.

#### 3. Hyperparameter Sensitivity:

• **EIOU** relies on properly scaled terms for center distance and aspect ratio, requiring careful tuning for best performance.

#### 4. Reduced Interpretability:

• While IOU directly measures overlap (easy to interpret), EIOU combines multiple terms, making it less intuitive.

## **Comparison Table**

Aspect	IOU	EIOU
Overlap Measurement	Based on area overlap only	Incorporates overlap, center distance, and aspect ratio
Center Distance	Ignored	Penalized
Aspect Ratio	Ignored	Penalized
Optimization	Zero gradients when no overlap	Continuous gradient updates
Complexity	Simple	Computationally heavier
Performance	Slower convergence	Faster and smoother convergence

## **Use Case Recommendations**

• Use **IOU** when:



- Simplicity and speed are more important.
- Bounding boxes have a high likelihood of overlapping.
- Use **EIOU** when:
  - Precise bounding box alignment is critical (e.g., small objects, dense environments).
  - Faster convergence during training is desired.

Let me know if you'd like more examples or further clarifications!



explain the center distance aspect in the above generated



# **Understanding the Center Distance Aspect in EIOU**

In **EIOU** (Efficient Intersection Over Union), the center distance aspect penalizes the Euclidean distance between the centers of the predicted bounding box and the ground truth bounding box. This ensures that even if the two boxes do not overlap, the loss function encourages the predicted box to move closer to the ground truth box.

#### **How It Works**

#### 1. Formula for Center Distance:

The center distance term is included in the **EIOU Loss** as part of  $L_{IOU}$ :

$$L_{IOU} = 1 - ext{IOU} + rac{
ho^2(b,b_{gt})}{w_c^2 + h_c^2}$$

Where:

- ullet  $b=(x_p,y_p)$ : Center of the predicted box.
- $b_{qt}=(x_{qt},y_{qt})$ : Center of the ground truth box.
- ullet  $ho(b,b_{gt})=\sqrt{(x_p-x_{gt})^2+(y_p-y_{gt})^2}$ : Euclidean distance between the two centers.
- ullet  $w_c,h_c$ : Width and height of the smallest enclosing box that contains both the predicted and ground truth boxes.

#### 2. Intuition:

- ullet IOU limitation: If the two boxes do not overlap, IOU=0, which halts learning because the gradient becomes zero.
- **EIOU improvement**: By introducing the center distance, the model receives feedback (gradient) even when the boxes do not overlap. This enables the model to adjust the predicted box toward the ground truth.

# Step-by-Step Example

## **Inputs:**

- 1. Predicted bounding box b: Center at (10,10), width w=4, height h=6.
- 2. Ground truth bounding box  $b_{gt}$ : Center at (12,12), width  $w_{gt}=5$ , height  $h_{gt}=7$ .

#### **Process:**

- 1. Compute the center coordinates:
  - Predicted center:  $(x_n, y_n) = (10, 10)$ .
  - $\circ$  Ground truth center:  $(x_{at}, y_{at}) = (12, 12)$ .
- 2. Calculate Euclidean distance ( $\rho(b,b_{at})$ ):

$$ho(b,b_{gt}) = \sqrt{(x_p - x_{gt})^2 + (y_p - y_{gt})^2}$$

Substituting:

$$ho(b,b_{gt})=\sqrt{(10-12)^2+(10-12)^2}=\sqrt{4+4}=\sqrt{8}pprox 2.83$$

- 3. Compute the diagonal of the enclosing box:
  - $\circ$  Enclosing box width:  $w_c = \max(x_p + w/2, x_{gt} + w_{gt}/2) \min(x_p w/2, x_{gt} w_{gt}/2).$
  - $\circ$  Enclosing box height:  $h_c = \max(y_p + h/2, y_{gt} + h_{gt}/2) \min(y_p h/2, y_{gt} h_{gt}/2).$  Substituting:
  - $\circ \ w_c = \max(12, 13.5) \min(8, 9.5) = 13.5 8 = 5.5.$
  - $\circ \ \ h_c = \max(13, 15.5) \min(7, 8.5) = 15.5 7 = 8.5.$
- 4. Normalize center distance:

$$rac{
ho^2(b,b_{gt})}{w_c^2+h_c^2} = rac{(2.83)^2}{(5.5)^2+(8.5)^2} = rac{8}{30.25+72.25} = rac{8}{102.5} pprox 0.078$$

# **Output:**

 $\bullet$  The center distance term contributes 0.078 to the overall EIOU loss. Even without overlap, this term encourages the predicted box to move closer to the ground truth box.

# **Key Advantages of Including Center Distance:**

- 1. Guidance Without Overlap:
  - Unlike IOU, EIOU provides non-zero gradients even when there is no overlap between the boxes, allowing the model to adjust predictions.
- 2. Faster Convergence:
  - The center distance accelerates learning by directly penalizing spatial misalignment, guiding the bounding box toward the correct position.
- 3. Enhanced Localization:
  - By explicitly considering the center alignment, EIOU results in more accurately localized bounding boxes, especially for non-overlapping or partially overlapping cases.

If you'd like further elaboration or visualization, let me know!



how triplet loss is used in that how it minimizes the error in the code as well as given research paper



# **Understanding and Application of Triplet Loss**

#### **Definition:**

Triplet loss is designed to ensure that an anchor image is closer to a positive (similar) image than it is to a negative (dissimilar) image by a margin. The formula is:

$$L_{triplet} = \max(0, lpha + d(a, p) - d(a, n))$$

#### Where:

- a: Anchor embedding.
- p: Positive embedding (similar to the anchor).
- n: Negative embedding (dissimilar to the anchor).
- ullet d(x,y): Distance function, typically Euclidean distance or cosine similarity.
- $\alpha$ : Margin, a hyperparameter that defines the minimum required separation between positive and negative pairs.

# Role in the Research Paper

In the research paper, triplet loss is used in the **Stacked Cross Attention Network (SCAN)** to optimize the alignment of unsafe behavior images with semantic safety rules. The goal is to:

- 1. **Bring Similar Pairs Closer**: Ensure the image features and their corresponding safety rule features are semantically close in the embedding space.
- 2. **Push Dissimilar Pairs Apart**: Increase the separation between mismatched image-rule pairs.

# **How It Works in Reducing Error**

- 1. Anchor, Positive, and Negative Selection:
  - **Anchor**: A feature vector for an image region (e.g., unsafe behavior).
  - **Positive**: A feature vector for the correct safety rule associated with the anchor.
  - **Negative**: A feature vector for an unrelated safety rule.
- 2. Minimizing the Loss:
  - The loss is minimized when:

$$d(a,p) + \alpha < d(a,n)$$

This ensures the correct image-rule pairs are closer than the incorrect pairs by at least a margin  $\alpha$ .

### 3. Impact on the Embedding Space:

 Triplet loss rearranges the embedding space such that clusters of similar pairs (image and rule) form compact regions, while dissimilar pairs are pushed to separate regions. This improves the model's ability to distinguish and retrieve the correct safety rules for given images.

## **How It is Implemented in Code**

Here's a pseudo-code representation of how triplet loss is typically implemented:

```
python
 import torch
 import torch.nn.functional as F
 def triplet loss(anchor, positive, negative, margin=0.2):
     Computes triplet loss.
     Args:
         anchor: Tensor of anchor embeddings (N, D)
         positive: Tensor of positive embeddings (N, D)
         negative: Tensor of negative embeddings (N, D)
         margin: Minimum distance margin between positive and negative pairs.
     Returns:
     Scalar loss value.
     # Compute pairwise distances
     d_ap = F.pairwise_distance(anchor, positive) # Distance between anchor and positive
     d_an = F.pairwise_distance(anchor, negative) # Distance between anchor and negative
     # Compute triplet loss
     loss = torch.relu(margin + d ap - d an) # Apply margin and ReLU for hinge loss
     return loss.mean()
```

# **Code Execution Example**

#### 1. Input:

- $\circ$  Anchor embedding: [0.5, 0.2]
- $\circ$  Positive embedding: [0.6, 0.3]
- $\circ$  Negative embedding: [1.0, 1.0]
- $\circ$  Margin: 0.2

#### 2. Processing:

- $\circ$  Compute d(a,p): Euclidean distance between anchor and positive.
- $\circ$  Compute d(a, n): Euclidean distance between anchor and negative.
- $\circ$  Compute loss using the formula  $\max(0, \alpha + d(a, p) d(a, n))$ .

#### 3. Output:

• A scalar value representing the loss, which is then used for backpropagation.

## Minimizing Error in the Research Context

In the research, triplet loss ensures:

## 1. Correct Matching:

 $\circ$  Unsafe behavior images (e.g., "working at height without harness") are matched to the correct safety rule (e.g., "must wear a harness") by reducing d(a,p).

#### 2. Error Minimization:

 $\circ$  Incorrect associations (e.g., matching "working at height without harness" to an unrelated rule like "wear a helmet") are penalized by increasing d(a,n).

# 3. Better Model Training:

 By learning to separate positives and negatives effectively, the model achieves higher recall and precision in retrieving the correct safety rule for an image.

# **Advantages in SCAN and Similar Applications**

## 1. Hard Negatives:

 Triplet loss focuses on "hard negatives" (cases where the anchor and negative are very similar but incorrect), making the model robust to challenging cases.

## 2. Semantic Embedding Space:

• Aligns the embedding space to capture semantic similarities between images and safety rules, improving accuracy in unsafe behavior detection.

Let me know if you'd like me to dive deeper into the implementation or mathematical aspects!



### Uploaded image

 $S^{\prime}$  \_LSE (I,T)=log  $\left(\sum 1_{i=1}^{i=1} k \right)$   $\left(\sum 2_{i=1}^{i=1} k \right$ 

explain this equation from pdf-1 given with respect to architecture



The provided equation is the **Log-Sum-Exp (LSE) pooling function** from the **Stacked Cross Attention Network (SCAN)** architecture as described in the paper. Let's break it down and explain it in detail with respect to the architecture:

# **Equation:**

$$S_{LSE}'(I,T) = \log \left( \sum_{i=1}^k \exp(\lambda_2 R(v_i, a_t^i)) 
ight)$$

# **Components:**

- 1. I: The image input, represented by features of k regions.
- 2. T: The safety rule text input, represented by features of n words.
- 3.  $v_i$ : The feature vector of the i-th region in the image.



- 4.  $a_t^i$ : The attended text feature vector for the *i*-th region (determined by text attention mechanism).
- 5.  $R(v_i, a_t^i)$ : The similarity score between the image feature vector  $v_i$  and the attended text feature vector  $a_t^i$ . It is calculated using the cosine similarity:

$$R(v_i, a_t^i) = rac{v_i^ op a_t^i}{\|v_i\| \|a_t^i\|}$$

- 6.  $\lambda_2$ : A scaling factor that amplifies or suppresses the importance of highly relevant pairs.
- 7.  $\log$ : The logarithmic function ensures numerical stability while emphasizing the most relevant pairs.
- 8.  $\sum_{i=1}^{k}$ : Aggregates the similarities across all k image regions.

# **Explanation with Respect to the Architecture:**

## 1. Image Features Representation:

- The input image is passed through a feature extractor (e.g., Faster R-CNN).
- $\circ$  The output is a set of k region-level feature vectors, each representing a region of the image (  $v_1,v_2,\ldots,v_k$ ).

## 2. Text Features Representation:

- $\circ$  The safety rule is tokenized into n words and passed through a Bi-GRU (bidirectional gated recurrent unit) network.
- $\circ$  The Bi-GRU produces word-level embeddings, which are combined using attention to create attended text vectors  $(a_t^1, a_t^2, \dots, a_t^k)$  for each image region.

#### 3. Stacked Cross Attention:

- $\circ$  The cross-attention mechanism computes the relevance between each image region feature (  $v_i$ ) and its corresponding attended text feature ( $a_t^i$ ).
- $\circ$  Cosine similarity ( $R(v_i,a_t^i)$ ) quantifies the alignment between the image region and the safety rule.

## 4. Pooling Function:

Instead of averaging or taking the maximum similarity, LSE pooling is used. It amplifies the
most significant similarity values, balancing focus on highly relevant pairs while still
considering weaker matches.

#### 5. Final Score:

- The logarithmic function applies scaling and ensures stability.
- $\circ$  The result,  $S'_{LSE}(I,T)$ , is a scalar score representing the overall semantic alignment between the image and the text.

# **Example Workflow:**

## 1. Input:

- $\circ$  Image with k=4 regions.
- $\circ$  Text with n=5 words.

## 2. Region Features ( $v_1, v_2, \ldots, v_4$ ):

- Extracted from the image.
- 3. Text Features ( $a_t^1, a_t^2, \ldots, a_t^4$ ):



- Attended features computed for each image region.
- 4. Similarity Computation ( $R(v_i, a_t^i)$ ):
  - Region  $v_1$  and attended text  $a_t^1$ : Similarity = 0.85.
  - $\circ$  Region  $v_2$  and attended text  $a_t^2$ : Similarity = 0.72.
  - $\circ$  Region  $v_3$  and attended text  $a_t^3$ : Similarity = 0.65.
  - $\circ$  Region  $v_4$  and attended text  $a_t^4$ : Similarity = 0.90.
- 5. Log-Sum-Exp Pooling:

$$S_{LSE}'(I,T) = \log\left(\exp(\lambda_2\cdot 0.85) + \exp(\lambda_2\cdot 0.72) + \exp(\lambda_2\cdot 0.65) + \exp(\lambda_2\cdot 0.90)\right)$$

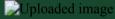
- 6. Final Score:
  - The result is a single value indicating how well the image and text match.

# **Advantages of LSE Pooling in SCAN:**

- 1. Smooth Aggregation:
  - Unlike max-pooling, LSE considers all regions but emphasizes the most relevant matches.
- 2. Numerical Stability:
  - The logarithmic function prevents overflow during computation with high similarity scores.
- 3. Focus on High Relevance:
  - The exponential weighting ensures that highly relevant image-text pairs dominate the final score.

Let me know if you'd like additional details or a practical code example!





Soft-Non-Maximum Suppression explain this with respect to the pdf-2 given



The equation in the image explains **Soft Non-Maximum Suppression (Soft-NMS)**, a technique used in object detection to handle overlapping bounding boxes more effectively. Below is the explanation of Soft-NMS in the context of **PDF-2 (YOLO-ESCA)**:

# **Soft-NMS: Definition and Purpose**

- In traditional Non-Maximum Suppression (NMS), overlapping bounding boxes are suppressed (set to zero) if their Intersection over Union (IOU) exceeds a threshold. However, this hard thresholding can lead to missed detections, especially in dense object scenarios or when objects overlap significantly.
- **Soft-NMS** modifies this by reducing the scores of overlapping boxes instead of removing them outright. This allows the detector to retain multiple overlapping boxes with reduced confidence scores, which can help in handling occluded or densely packed objects.

## **Equation:**

$$s_i = s_i \cdot e^{-rac{ ext{IOU}(M,b_i)^2}{\sigma}}$$

Where:

- $s_i$ : Confidence score of the i-th bounding box.
- $\mathrm{IOU}(M,b_i)$ : Intersection over Union between the box  $b_i$  and the box M with the highest confidence.
- $\sigma$ : A hyperparameter controlling the decay rate of the confidence score.

# **Explanation in Context of YOLO-ESCA (PDF-2):**

#### 1. Problem with NMS:

- o In construction site scenarios (e.g., helmet detection), objects such as helmets or people are often occluded or closely packed, leading to high IOU values between bounding boxes.
- Traditional NMS may suppress correct detections, leading to a higher missed detection rate.

#### 2. Soft-NMS Solution:

- Instead of outright removing bounding boxes with high IOU, Soft-NMS decays their confidence scores based on their overlap with the highest confidence box.
- This ensures that occluded objects or densely packed regions are still detected but with lower confidence.

#### **How It Works:**

#### 1. Input:

- A set of bounding boxes with their confidence scores and IOU values.
- Example:
  - Box 1: Confidence = 0.9, IOU = 0.8.
  - Box 2: Confidence = 0.7, IOU = 0.6.

### 2. Process:

- $\circ$  Identify the bounding box M with the highest confidence (e.g., Box 1).
- $\circ$  For every other box  $b_i$ , calculate the decay factor based on its IOU with M:

$$ext{Decay Factor} = e^{-rac{ ext{IOU}(M,b_i)^2}{\sigma}}$$

 $\circ$  Update the confidence score of  $b_i$ :

$$s_i = s_i \cdot \text{Decay Factor}$$

#### 3. Output:

Adjusted confidence scores for all bounding boxes.

# **Advantages in YOLO-ESCA:**

#### 1. Improved Recall:

 Soft-NMS reduces the likelihood of missing detections in dense or occluded environments, such as helmets on workers in crowded construction sites.



# 2. Handles Overlapping Targets:

• Essential for detecting small targets like helmet chin straps or overlapping objects like workers standing close to each other.

## 3. Flexibility:

 $\circ$  The hyperparameter  $\sigma$  allows fine-tuning the suppression behavior, balancing precision and recall.

# **Example Workflow:**

### 1. Input Boxes:

- Box A: Confidence = 0.9, IOU = 0.7 with the highest confidence box.
- Box B: Confidence = 0.8, IOU = 0.5 with the highest confidence box.

## 2. **Soft-NMS Decay**:

For Box A:

Decay Factor 
$$=e^{-\frac{0.7^2}{0.5}}=e^{-0.98}pprox 0.375$$

New confidence:  $0.9 \cdot 0.375 = 0.337$ .

For Box B:

Decay Factor 
$$= e^{-\frac{0.5^2}{0.5}} = e^{-0.5} \approx 0.606$$

New confidence:  $0.8 \cdot 0.606 = 0.485$ .

#### 3. Output:

- Adjusted confidence scores:
  - Box A: Confidence = 0.337.
  - Box B: Confidence = 0.485.

# **Key Takeaways for YOLO-ESCA:**

- **Soft-NMS ensures better performance** in complex construction site environments by reducing missed detections and improving localization of overlapping objects.
- By retaining overlapping boxes with decayed confidence scores, it improves detection robustness in scenarios with small, dense, or occluded objects.

Let me know if you'd like further clarification or examples!



Uploaded image

**Uploaded** image

Explain this feature from the pdf-2



These diagrams illustrate the **Convolutional Block Attention Module (CBAM)**, which is a lightweight and effective attention mechanism used to refine feature maps in neural networks by sequentially

applying **channel attention** and **spatial attention**. Below is a detailed explanation of each module with respect to its context in **PDF-2 (YOLO-ESCA)**.

#### **Overview of CBAM**

CBAM applies attention mechanisms along two dimensions:

- 1. Channel Attention Module (CAM): Focuses on the importance of each feature channel.
- 2. **Spatial Attention Module (SAM)**: Focuses on the importance of spatial locations within the feature map.

The goal is to refine the input features by selectively emphasizing relevant information in both the channel and spatial dimensions.

# 1. Channel Attention Module (CAM)

**Diagram (b)** explains the structure of the **Channel Attention Module**:

#### **Process:**

## 1. Input:

 $\circ \:$  A feature map F of shape (C,H,W), where C is the number of channels, and H imes W is the spatial resolution.

# 2. Pooling:

• Apply **global average pooling** and **global max pooling** across the spatial dimensions to generate two channel-wise descriptors:

$$F_{avg} \in \mathbb{R}^{C imes 1 imes 1}, \quad F_{max} \in \mathbb{R}^{C imes 1 imes 1}$$

#### 3. Shared MLP:

 $\circ$  Pass  $F_{avg}$  and  $F_{max}$  through a shared multi-layer perceptron (MLP) with a reduction ratio r to reduce the dimensionality and then restore it.

$$M_C = \sigma(\mathrm{MLP}(F_{ava}) + \mathrm{MLP}(F_{max}))$$

 $\circ~$  Here,  $\sigma$  is the sigmoid function that scales the output to [0,1].

#### 4. Element-wise Multiplication:

 $\circ~$  Multiply the input feature map F with the attention weights  $M_C$  (broadcasted to (C,H,W) ):

$$F' = F \cdot M_C$$

#### **Purpose:**

 Helps the network focus on important feature channels, such as those associated with small objects like helmet straps in the YOLO-ESCA context.

# 2. Spatial Attention Module (SAM)

# **Diagram (c)** explains the structure of the **Spatial Attention Module**:

#### **Process:**

#### 1. Input:

 $\circ$  The output from the channel attention module,  $F' \in \mathbb{R}^{C imes H imes W}$  .

#### 2. Pooling:

 Apply max pooling and average pooling along the channel dimension to generate two spatial maps:

$$F_{avg} \in \mathbb{R}^{1 imes H imes W}, \quad F_{max} \in \mathbb{R}^{1 imes H imes W}$$

#### 3. Concatenation:

• Concatenate the pooled outputs along the channel dimension:

$$F_{concat} = [F_{avg}; F_{max}] \in \mathbb{R}^{2 imes H imes W}$$

#### 4. Convolution:

 $\circ$  Pass the concatenated map through a 3 imes 3 convolution to produce a spatial attention map  $M_S$ :

$$M_S = \sigma(\operatorname{Conv}(F_{concat}))$$

 $\circ$  Here,  $\sigma$  is the sigmoid function.

#### 5. Element-wise Multiplication:

 $\circ$  Multiply the spatial attention map  $M_S$  with F' (broadcasted to (C,H,W)):

$$F'' = F' \cdot M_S$$

#### **Purpose:**

• Highlights **important spatial regions**, such as specific areas in an image where a helmet or unsafe behavior is detected.

## 3. Combined Attention Mechanism

**Diagram (a)** shows the integration of CAM and SAM:

- 1. The **input feature map** F is first processed by the **Channel Attention Module (CAM)** to produce F'.
- 2. The refined feature map F' is then processed by the **Spatial Attention Module (SAM)** to produce the final output F''.

# **Advantages in YOLO-ESCA:**

#### 1. Small Object Detection:

 CBAM enhances attention on small, dense objects (e.g., helmets, chin straps) that might be overlooked in standard feature extraction.



#### 2. Reduced Model Complexity:

o CBAM is lightweight and introduces negligible computational overhead, making it suitable for real-time applications like UAV-based helmet detection.

# 3. Improved Detection Accuracy:

• By refining features along both the channel and spatial dimensions, CBAM helps the YOLO-ESCA model better differentiate between safety-compliant and non-compliant behaviors.

# **Application Example in YOLO-ESCA:**

- **Input**: A feature map representing an image region with detected helmets.
- CAM: Highlights channels responsible for helmet-related features.
- SAM: Focuses on spatial areas where helmets or safety violations are present.
- **Output**: A refined feature map that enhances the model's ability to detect helmets and identify compliance.

Let me know if you'd like further elaboration or code examples!