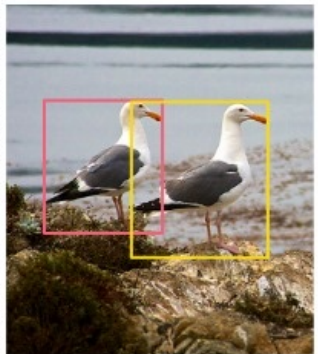
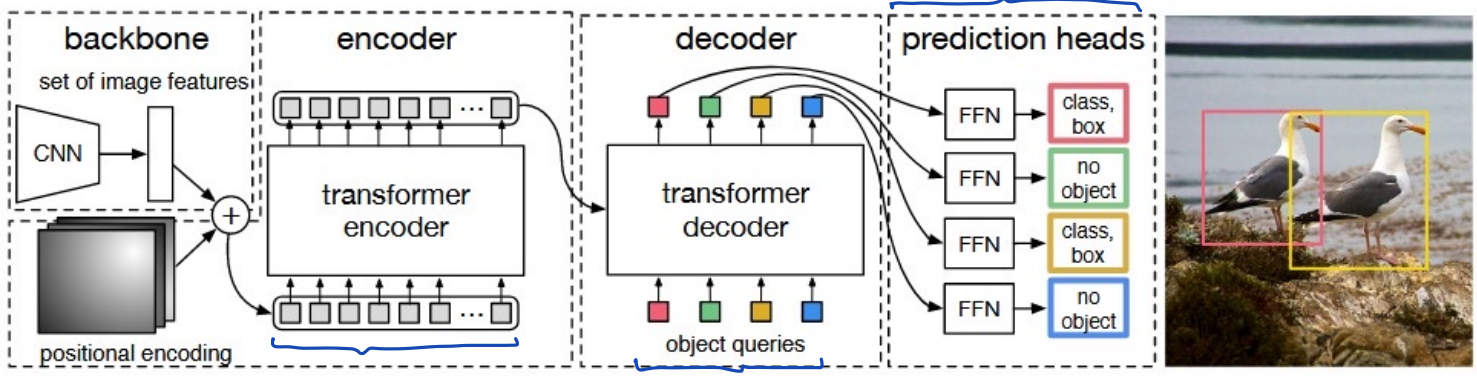


Detr:

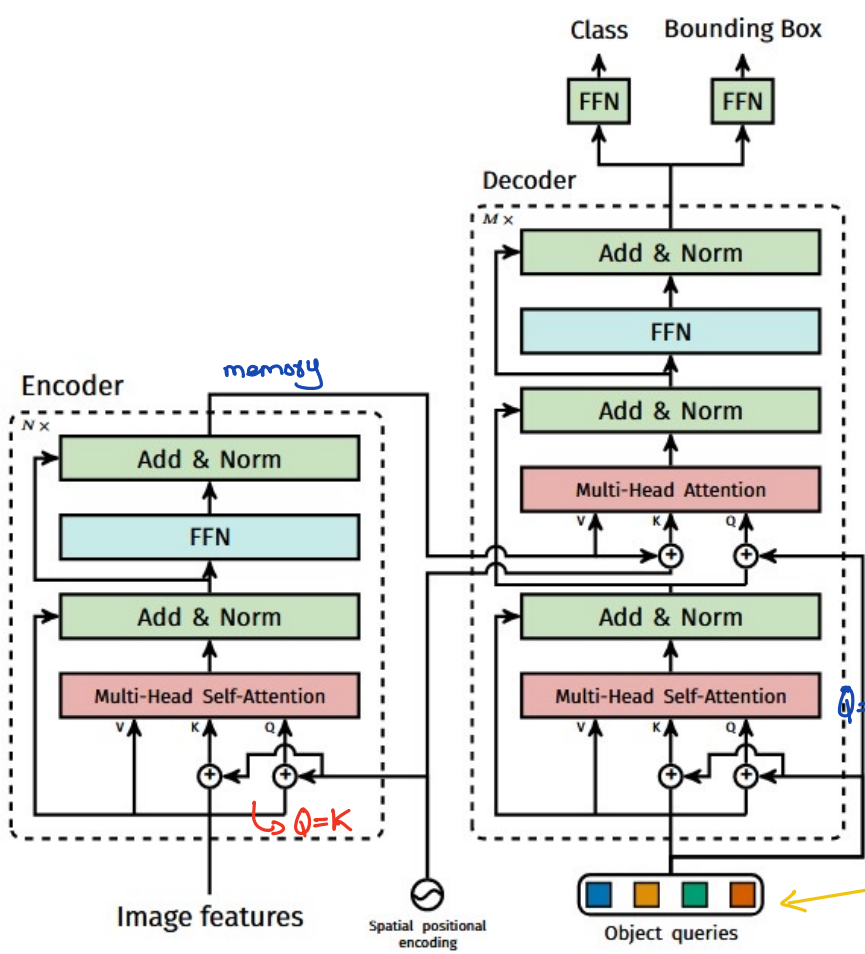
↳ loss with last decoder layer  
 minimize it with grad torch  
 ↳ Distinct bi-partite matching



↳ flatten image with positional encoding

query-embed → learnable  
 tgt → set to 0

↳ Actual loss that's used for backprop is from each decoder layer's output



evolving state (what has been found so far)

set to 0

$$Q = K = \text{query\_embed} + \text{tgt}$$

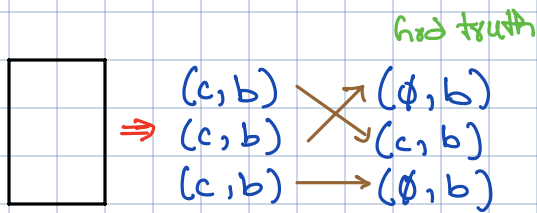
no embedding  
 ↳ reasonable (what to look for)

serve as anchor points to be added at each decoder layer.

Why do this?

- Adding the query embedding at each layer creates a form residual connection.

These help with gradient flow. Without this, the network might forget the initial query information.



match left to right  
where  $L$  is minimum.

$$\hat{\sigma} = \underset{\text{with assignment}}{\operatorname{argmin}} \sum_i^N \mathcal{L}_{\text{match}}(y_i, \hat{y}_{\sigma(i)})$$

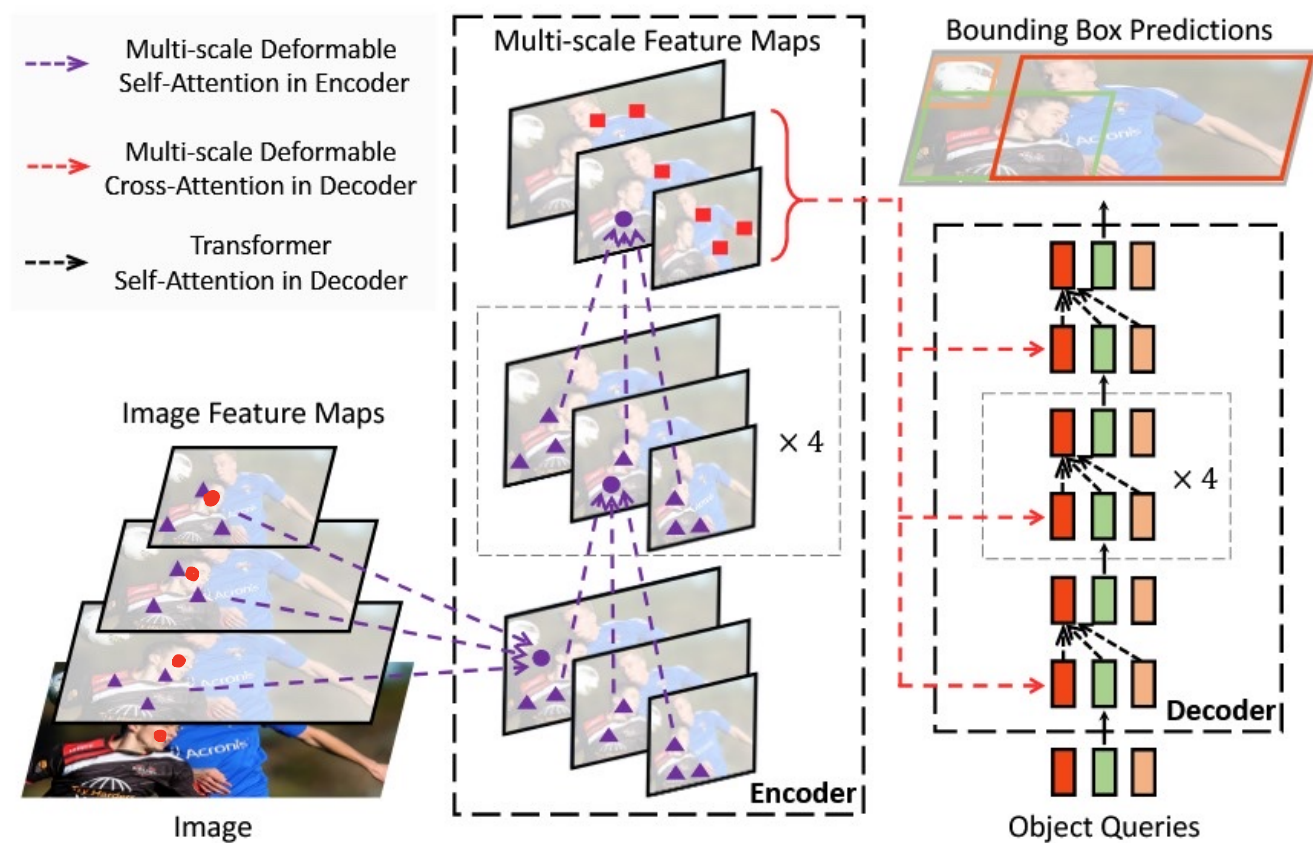
$$\mathcal{L}_{\text{match}} = -1_{\{c_i \neq \phi\}} \hat{p}_{\sigma(i)}(c_i) + 1_{\{c_i \neq \phi\}} \mathcal{L}_{\text{box}}(b_i, \hat{b}_{\sigma(i)})$$

Actual loss for training:

$$\mathcal{L}_{\text{Hungarian}} = \sum_i^N \left[ \log \hat{p}_{\hat{\sigma}(i)}(c_i) + 1_{\{c_i \neq \phi\}} \mathcal{L}_{\text{box}}(b_i, \hat{b}_{\hat{\sigma}(i)}) \right]$$

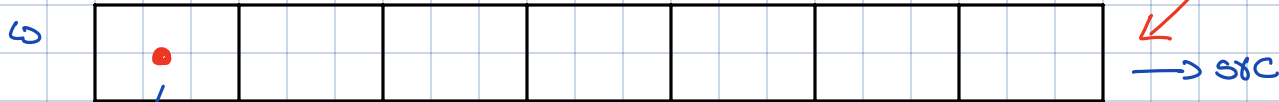
We do this for each decoder layer.

# Deformable Dets:



↳ flatten the four level image into a sequence.

↳ Add level & positional encoding to it  $\rightarrow pos = level + \text{positional encoding}$



↳ A reference point

↳ find it at all levels

↳ Pass this (src + pos) through a Linear layer to get 4 offsets & Attention weights.

↳ Do this for each level

↳ Get one output from each level

↳ Sum them & pass through a linear layer to get final output

↳ Repeat for each attention head

$$\text{MSDeformAttn}(z_q, \hat{p}_q, \{x^l\}_{l=1}^L) = \sum_{m=1}^M W_m \left[ \sum_{l=1}^L \sum_{k=1}^K A_{mlqk} \cdot W'_m x^l (\phi_l(\hat{p}_q) + \Delta p_{mlqk}) \right], \quad (3)$$

For query,

↳ Same process is done

↳ Except the reference point is on the image & found through a linear layer.

For a single scale:

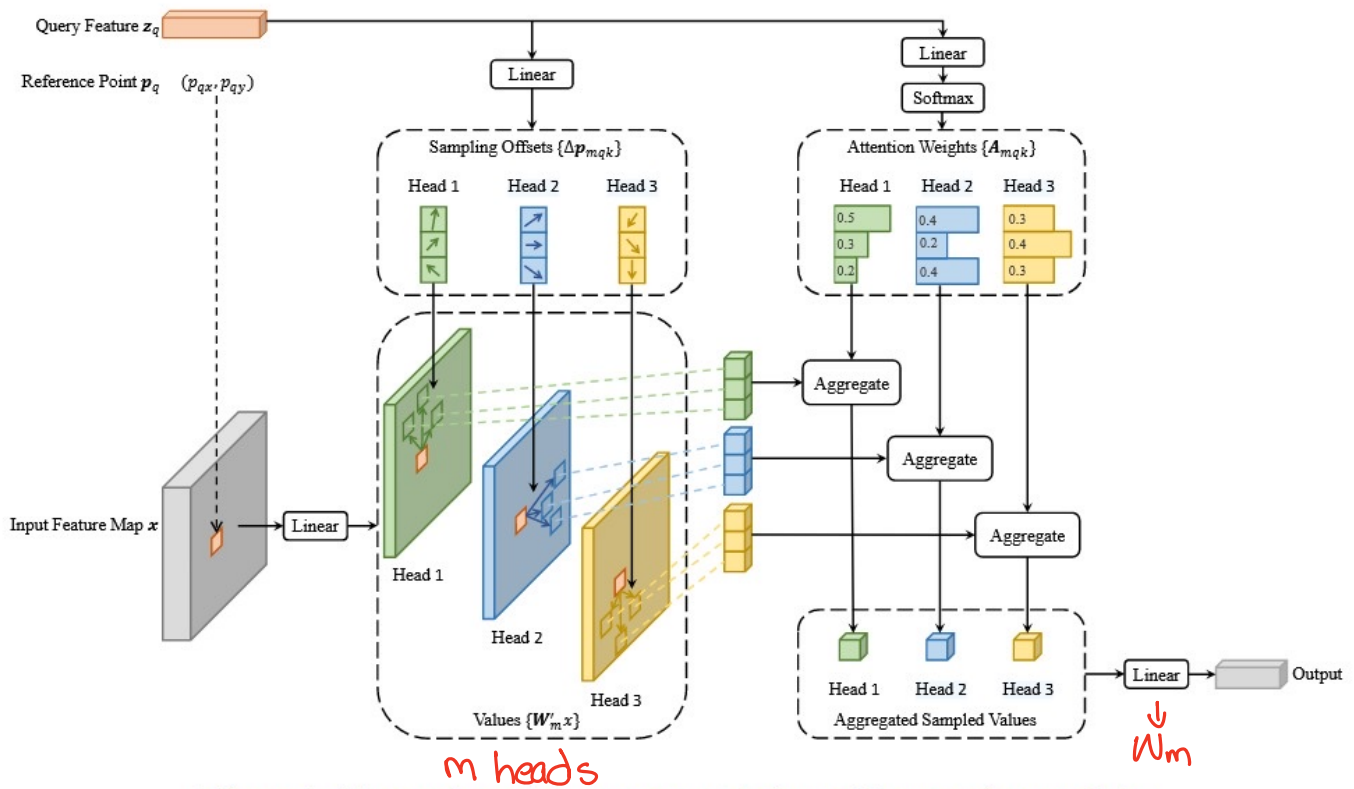


Figure 2: Illustration of the proposed deformable attention module.

## Dab - Detr:

- ↳ The authors recognized the object queries as slowing convergence.
  - ↳ From previous work, they know that cross-attn module is the one slowing the convergence.
  - ↳ Since later steps are same in encoder & de-coder, they conclude the issue is likely from queries.
- ↳ To find the reason, they have two hypothesis

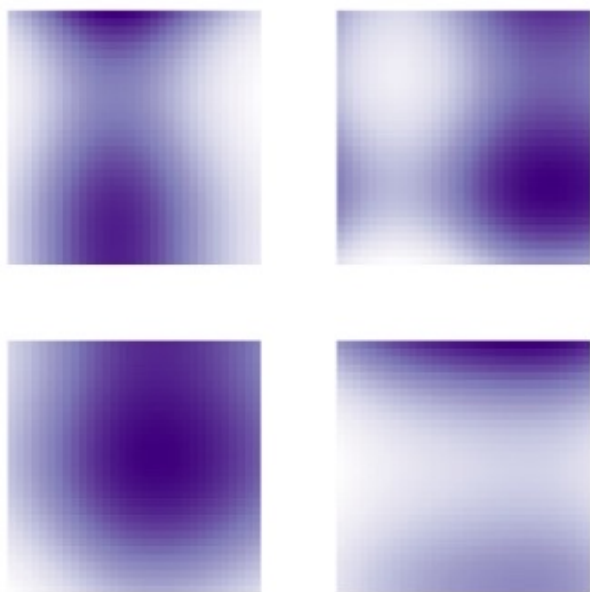
### 1. The Queries are hard to learn

- ↳ So, they fix learned Queries & train other modules & find that the improvement is minimal. So, this can't be it

### 2. The positional information in the queries is not encoded the same way as sinusoidal positional encoding used for image features.

- ↳ While sinusoidal encoding provides a structured & consistent way to represent positions for the image features from the start, the learnable queries have to gradually discover useful positional representations during training, which might slow down convergence.

- ↳ Furthermore, the model learns bad queries still



(a) DETR

- ↳ Computed as learned queries

Positional encoding of image

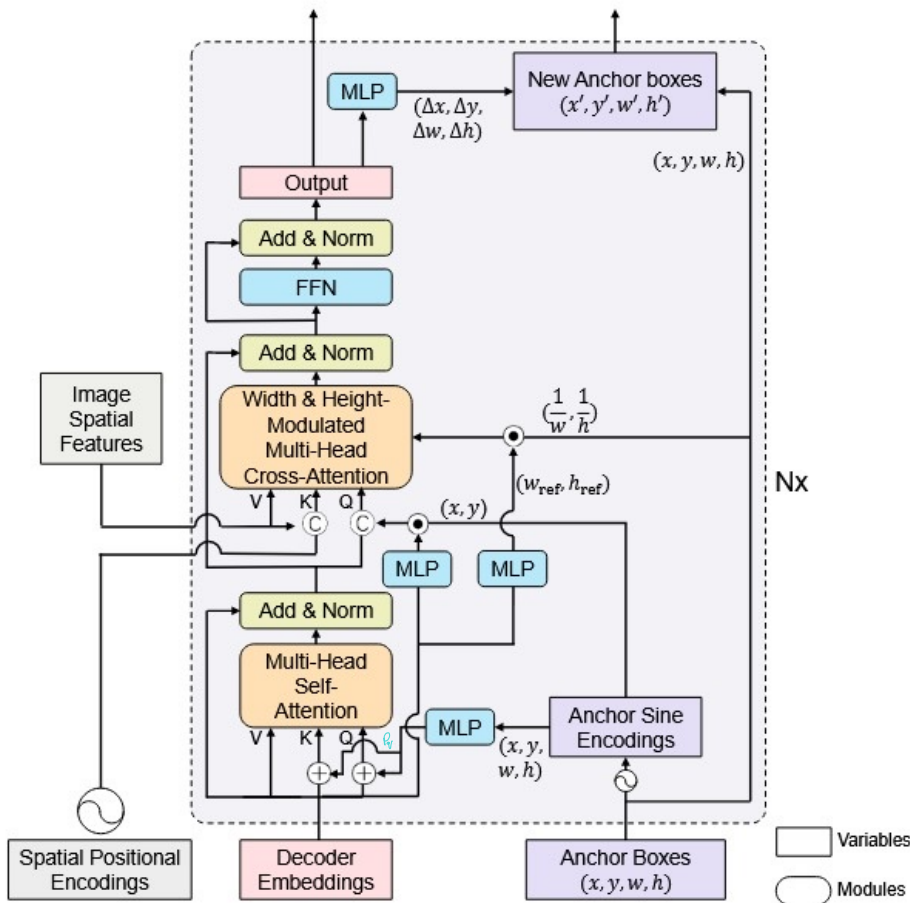
- ↳ All over the place

- ↳ No sense of scale

- ↳ No sense of importance of certain points more than others (Uniform)

- ↳ Query embeddings are needed to compute attention & they are NOT doing a good job at it.





↳ The Anchor boxes are added to output prediction of each decoder layers.

↳ This way anchor boxes are also learned

↳ They serve as priors for the position, size, shape of bbox of obj.

Content queries (Decoder Embeddings)  
 ↳  
 Positional queries (Anchor Boxes)  
 are fed to Decoder as dual queries.

↳ The Anchor boxes are learned directly & derive positional queries from these anchors.

$$\text{Self Attn: } Q_v = C_v + P_v, K = C_v + P_v, V_v = C_v$$

$$\text{Cross-Attn: } Q_v = \text{Cat}[C_v, \text{PE}(x_v, y_v)] \cdot \text{MLP}^{\text{cross}}(C_v)$$

$$K_{x,y} = \text{Cat}(F_{x,y}, \text{PE}(x,y)), V_{x,y} = F_{x,y}$$

↓  
 image feature at position (x,y)

In DETR models, when computing attention between queries and image features, both content and position information get mixed together. This makes it hard to know how much each aspect contributes to the similarity scores.

**What they're doing is:**

1. Separating content and position information by concatenating them rather than adding them
2. This lets the model learn to pay attention to each aspect separately

The MLP(csq) part is addressing a scaling issue. Position information and content information might naturally have different magnitudes or importance. They're using the content information to produce a custom scaling factor for the position information.

**Think of it like this:** Some objects might need precise position information (like small objects), while others might rely more on visual appearance. This approach lets the model adjust how much it cares about position versus content for each query, making it more flexible