Soutenance de projet 30 janvier 2025



## HeadMind Partners

## Challenge IA

Etude de cas HeadMind Partners - Dior



BOUAITA Rayane - DAVID Erwan - EL ANATI Pierre - FAYNOT Guillaume - TRIER Gabriel

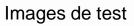


## 01 Problématique

### Contexte









BDD d'images

Catégories produits	%
Bags	30
RTW	23
Accessories	20
Shoes	12
SLG	11
Watches	0,04



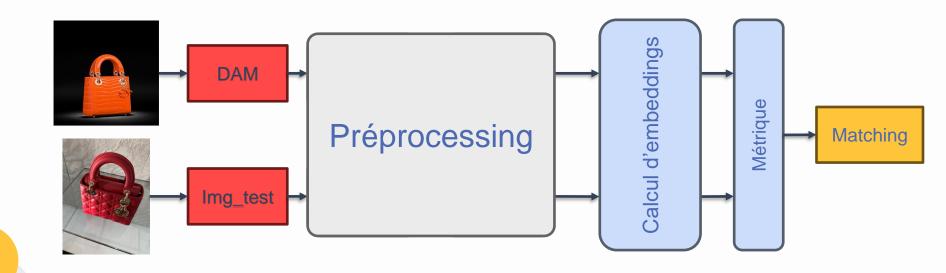




## 02 Pipeline

## **Pipeline**









## 03 Préprocessing

## Préprocessing



Images de test



BDD d'images



### O

#### Problèmes rencontrés

- Taille (HD)
- Dimensions (carré vs rectangle)
- Arrière-plan

### Solutions apportées

- Suppression de l'arrière-plan (rembg)
- Crop autour de la forme isolée (carré)
- Diminution de la résolution (256x 256)





## RemBg

- Module python open source
- Suppression d'arrière plan via modèles profonds



Image initiale

#### Fonctionnement de rembg

- Conversion en .png
- Normalisation par canal (RGB)
- Segmentation avec U²-Net (classification foreground vs background)
- Application du masque binaire au canal alpha















### Annexe U<sup>2</sup>-Net

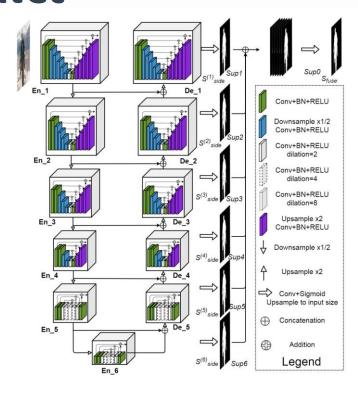




Figure 5. Illustration of our proposed U<sup>2</sup>-Net architecture. The main architecture is a U-Net like Encoder-Decoder, where each stage consists of our newly proposed residual U-block (RSU). For example, **En.1** is based on our RSU block shown in Fig. 2(e). Detailed configuration of RSU block of each stage is given in the last two rows of Table 1.



# 04 Data augmentation

## **Data augmentation**





Flip horizontal

#### Photo mal orientée

Rotation +/- 45°

#### Objet mal éclairé

- Ajustement des couleurs
- Ajustement de la luminosité

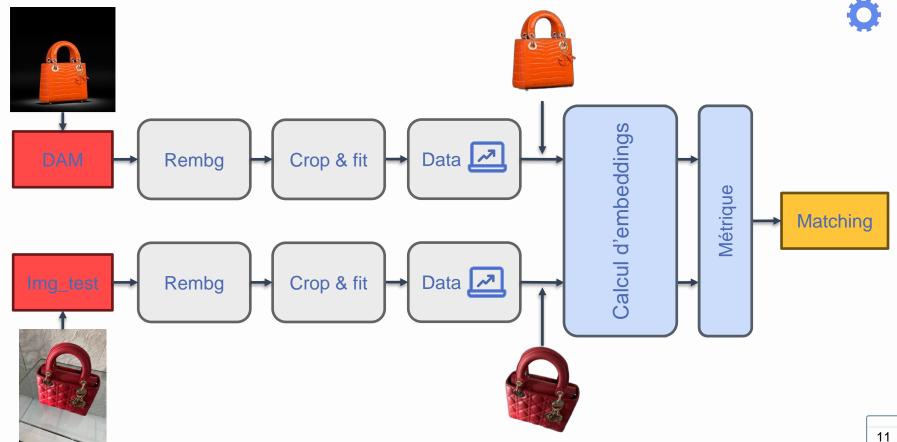
#### Objet mal cadré

Redimensionnement





## Bilan - Pipeline

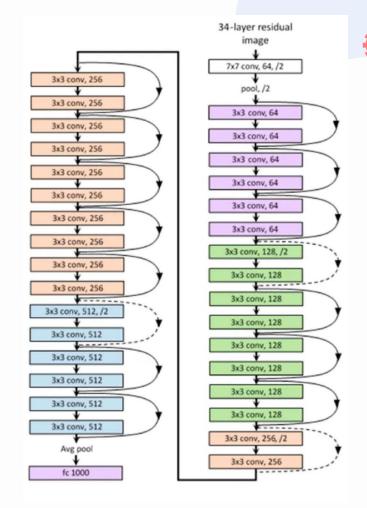




## 05 Modèles

### ResNet

- Architecture introduite en 2015 (Microsoft Research)
- 1ère place ILSVRC → 3,57%
- Skip connections
- Entrainé sur ImageNet
- Utilisé comme extracteur de caractéristiques











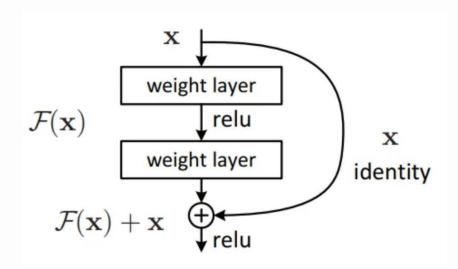
## **Skip Connections**

**Problème: Vanishing Gradient** 



**Solution: Skip Connections** 

- Apprentissage de fonction résiduelle
- Résolution des limites de profondeur







## ResNet50

,=======								
layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer		
conv1	112×112			7×7, 64, stride 2	2			
conv2_x	56×56	3×3 max pool, stride 2						
		$\left[\begin{array}{c} 3\times3,64\\ 3\times3,64 \end{array}\right]\times2$	$\left[\begin{array}{c} 3\times3,64\\ 3\times3,64 \end{array}\right]\times3$	$   \begin{bmatrix}     1 \times 1, 64 \\     3 \times 3, 64 \\     1 \times 1, 256   \end{bmatrix} \times 3 $	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$		
conv3_x	28×28	$\left[\begin{array}{c} 3\times3, 128\\ 3\times3, 128 \end{array}\right] \times 2$	$\left[\begin{array}{c} 3\times3, 128\\ 3\times3, 128 \end{array}\right] \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 8$		
conv4_x	14×14	$\left[\begin{array}{c} 3\times3,256\\ 3\times3,256 \end{array}\right]\times2$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 23$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 36$		
conv5_x	7×7	$\left[\begin{array}{c} 3\times3,512\\ 3\times3,512 \end{array}\right]\times2$	$\left[\begin{array}{c}3\times3,512\\3\times3,512\end{array}\right]\times3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$ \left[\begin{array}{c} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array}\right] \times 3 $		
	1×1	average pool, 1000-d fc, softmax						
FLOPs		$1.8 \times 10^{9}$	$3.6 \times 10^9$	$3.8 \times 10^9$	$7.6 \times 10^9$	11.3×10 <sup>9</sup>		



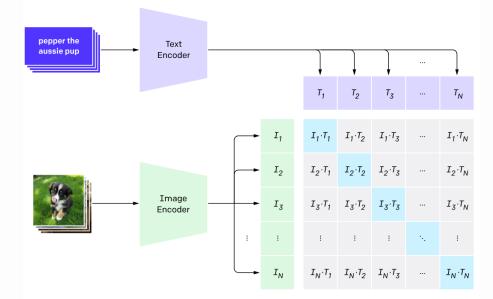






#### Formation auto-supervisée :

- Deux réseaux de neurones distincts : ViT et transfomer
- Générer deux vecteurs de même dimension dans un espace commun
- Rapprocher les paires correspondantes et à éloigner les paires incorrectes
- La perte contrastive InfoNCE









#### **Motivations**

#### Meilleure correspondance sémantique des images

CLIP encode des concepts visuels de manière plus riche qu'un ResNet, qui lui se limite aux caractéristiques discriminatives utiles pour la classification.

#### **Robustesse aux Variations Visuelles**

Entraîné sur des données variées avec un alignement image-texte, ce qui lui permet de mieux capturer les similitudes visuelles sous différents angles.

#### Zero-shot

CLIP peut fonctionner directement sans réentraînement.



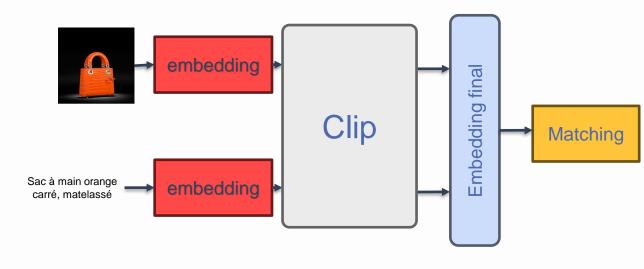




#### **Embedding des images**

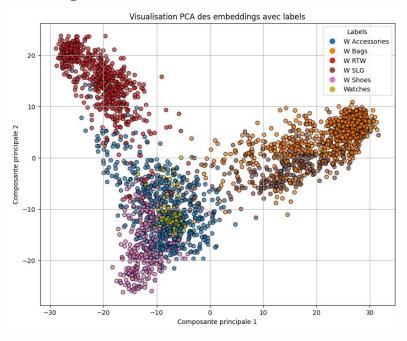
ResNet50	Clip				
Non	Non				
Couleur	Couleur				
Cosine	Cosine				
41,25%	31,25%				
55%	51.25%				
61%	61.25%				

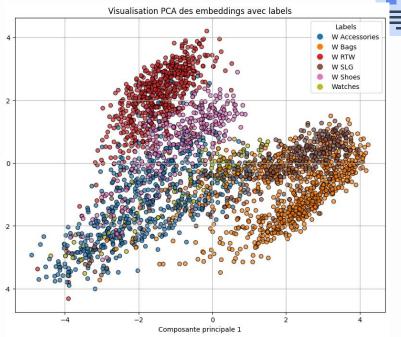
#### Embedding des images et du texte











ResNet50 CLIP



### DINOv2





#### Pré-traitement de l'image

Redimensionnement, conversion en tenseur, normalisation

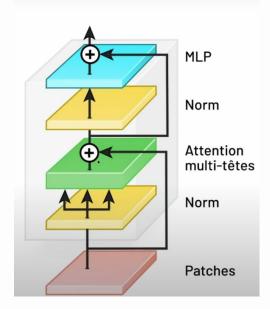
#### **Traitement**

Découpage en patchs, encodage des patchs (partie locale de l'image), Ajout d'un token CLS (résumé global de l'image) Passage des patchs dans un transformer -> multi heads attention

#### **Extraction du token**

Représentation riche et compressée de l'image

#### **Encodeur Transformer**





### DINOv2





#### **Motivations**

#### Par rapport à ResNet

- Plus robuste aux variations intra-classe (ex: variations d'éclairage, d'angles de prise de vue ou de bruit dans les images)
- Contrairement à ResNet (réseau convolutif classique), DINOv2 utilise un Vision Transformer
   (VIT)

#### Par rapport à Clip

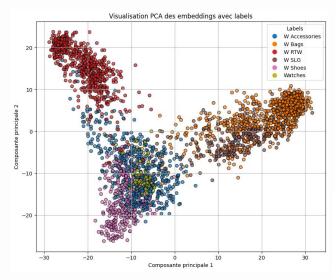
Encoder d'image plus finement optimisé



### **DINOv2**



DinoV2	ResNet50
Non	Non
FlipV+Couleur	FlipH +Couleur
Cosine	Cosine
42,50%	45%
55%	56,25%
60%	65%

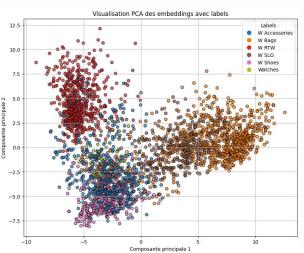


ResNet50





#### DinoV2



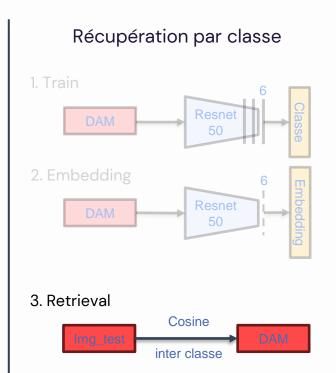


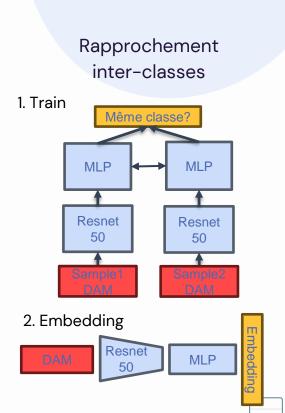


## 06 Fine-Tuning

## Fine-Tuning

## Spécialisation sur le Dataset 1. Train Resnet 2. Embedding Resnet 50 3. Retrieval Cosine







## 07 Métriques

## Métriques



Comment retrouver une image de train à partir d'une image de test ?

Quelle métrique pour trouver un minimum de distance entre les embeddings?

#### Similarité cosinus

$$D(A,B) = \frac{A \cdot B}{\|A\| \|B\|}$$

$$D(A,B) = \sum_{i} |A_i -$$

$$D(A,B) = \sum_{i} |A_i - B_i|$$
  $D(A,B) = \sqrt{\sum_{i} (A_i - B_i)^2}$ 





### Résultats



Embeddings	ResNet50	ResNet50	ResNet50	ResNet50	ResNet50	ResNet50	ResNet50	DinoV2	Clip	ResNet50	ResNet50
Fine tuning	Non	Non	Non	Non	Non	Non	Non	Non	Non	Oui	Oui+classe
Data augmentation	Non	Non	Non	Flip H	Couleur	FlipH +Couleur	FlipV+Couleur	FlipV+Couleur	Couleur	Non	Non
Méthode	Cosine	L1	L2	Cosine	Cosine	Cosine	Cosine	Cosine	Cosine	Cosine	Cosine
Accuracy Top 1	45%	45%	45%	45%	41,25%	45%	37,50	42,50%	31,25%	28%	46%
Accuracy Top 3	58,75%	57,5%	56,25%	60%	55%	56,25%	54%	55%	51.25%	46%	59%
Accuracy Top 5	67.5%	63 75%	66 25%	68 75%	61%	65%	59%	60%	61 25%	53%	67%

#### Constatations

- ResNet50 > DinoV2 > Clip
- FlipH > No Data Augmentation
- Cosine > L1 & Cosine > L2





# O7 Conclusion

# Merci

