

# Rapport TP1 Multimodalité

Hamza Omari<sup>1\*</sup>

<sup>1</sup>\*UFR Mathématiques et Informatique, Université Paris-Cité, Paris,  
75006, France.

Corresponding author(s). E-mail(s): [hamza.omari@etu.u-paris.fr](mailto:hamza.omari@etu.u-paris.fr);

## 1- familiarisation avec le notebook

Le notebook montre l'utilisation des modèles CLIP de la bibliothèque clip. Deux familles de modèles sont présentes : ResNet et ViT. Un jeu de données composé de paires (image, description) est utilisé. En calculant la similarité, on peut évaluer les performances des espaces d'embeddings et réaliser un fine-tuning sur le dataset ROCO. Le notebook présente ensuite l'utilisation du modèle après l'entraînement, son évaluation, ainsi qu'un exemple d'évaluation sur une tâche de classification.

### Evaluations des modèles clips

ils existent deux familles de modèles présentées: ResNet et ViT RN50, RN101, RN50x4, RN50x16, RN50x64  
ViT-B/32, ViT-B/16, ViT-L/14, ViT-L/14@336px

Pour une évaluation quantitative, j'ai d'abord opté pour une méthode simple : calculer les similarités image–texte, les visualiser sous forme de matrice, puis analyser la **trace** de cette matrice. La trace indique directement à quel point le modèle parvient à maximiser la similarité entre chaque image et son texte associé.

Cependant, en visualisant les matrices de similarités, j'ai constaté que certains modèles affichent une **bonne valeur de trace**, mais restent globalement **mauvais** : ils n'arrivent pas à séparer correctement les images et les textes non associés, ce qui se traduit par des valeurs off-diagonales élevées. Comme dans le cas du modèle RN101, voir la figure 1 Ainsi, la trace seule est insuffisante pour mesurer la qualité globale du modèle.

Pour cela, l'évaluation que j'ai retenue consiste à prendre en compte **à la fois** la diagonale (paires correctes) et les éléments hors-diagonale (paires incorrectes), selon la formule suivante :

$$\text{score} = \frac{\text{trace}(S)}{N} - \frac{\sum S - \text{trace}(S)}{N^2 - N} . \quad (1)$$

moyenne des similarités correctes      moyenne des similarités incorrectes

où  $S$  désigne la matrice de similarité et  $N$  le nombre d'images (et donc de textes associés).

Ce score n'augmente que lorsque :

- les similarités correctes (diagonale) augmentent,
- et les similarités incorrectes (hors-diagonale) diminuent.

Il constitue donc une mesure beaucoup plus fiable et robuste que la trace seule pour évaluer la performance globale d'un modèle CLIP.

**les matrices de similarités calculées, pour les autres différents modèles se retrouve en bas de ce document.**

Nous classons ensuite les modèles par ordre décroissant de performance, selon la métrique définie par l'équation 1. Les résultats obtenus sont les suivants



**Fig. 1:** Matrice de similarité pour le modèle RN101

Model		Score
<b>ViT-L/14@336px</b>		<b>0.094371</b>
<b>ViT-L/14</b>		<b>0.092693</b>
<b>RN50x64</b>		<b>0.091421</b>
<b>RN50x16</b>		<b>0.081308</b>
<b>RN50x4</b>		<b>0.061293</b>
<b>ViT-B/16</b>		<b>0.056187</b>
<b>ViT-B/32</b>		<b>0.054636</b>
<b>RN101</b>		<b>0.054573</b>
<b>RN50</b>		<b>0.052011</b>

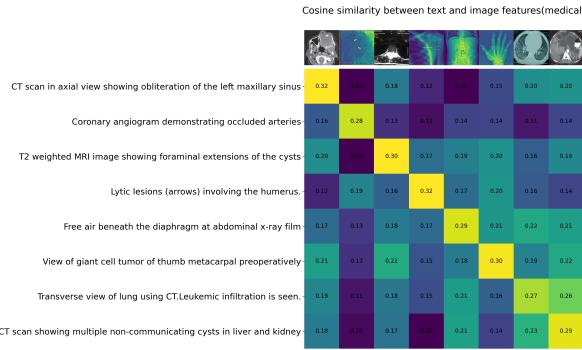
**Fig. 2:** Evaluations

## Performance sur d'autres images médicales

Pour cette tâche, nous sélectionnons des images différents et nous utilisons ChatGPT pour les annoter (c'est-à-dire les décrire). Ensuite, de la même manière que précédemment, nous calculons la matrice de similarité. voir figure 3

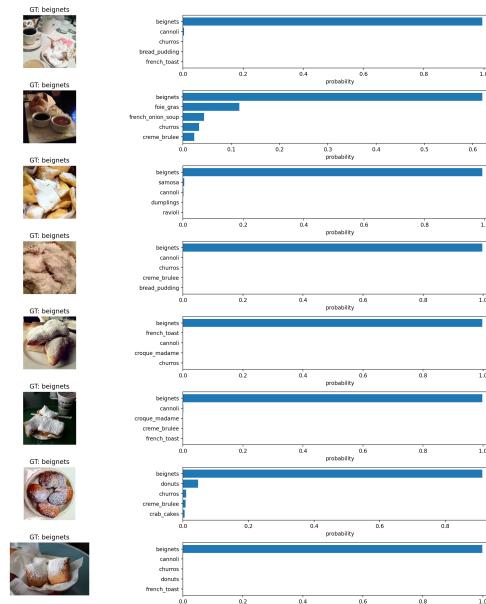
## Évaluation sur une tâche de classification

Pour cette partie, nous utilisons un dataset de Hugging Face, *Food101*. Nous sélectionnons quelques images et définissons des descriptions textuelles correspondant aux différentes classes. Ensuite, nous calculons la similarité entre chaque image et l'ensemble des descriptions des classes choisies. Les logits obtenus sont ensuite



**Fig. 3:** Evaluations sur d'autres images

passés dans une fonction *softmax*, ce qui nous permet de visualiser la distribution des probabilités et d'interpréter les prédictions du modèle.



**Fig. 4:** Evaluations sur la classif

On voit bien que le modèle arrive à bien associer l'image à sa vraie classe.

## Code source

: Continuation sur le même notebook. [GitHub](#)

Cosine similarity for RN50



Cosine similarity for RN50x4

								
a cup of coffee on a saucer	0.42	0.27	0.23	0.22	0.30	0.22	0.24	0.22
a black-and-white silhouette of a horse	0.26	0.45	0.23	0.25	0.26	0.30	0.29	0.24
a rocket standing on a launchpad	0.28	0.29	0.40	0.27	0.31	0.30	0.28	0.30
a red motorcycle standing in a garage	0.26	0.25	0.22	0.40	0.28	0.22	0.24	0.24
a facial photo of a tabby cat	0.27	0.26	0.21	0.21	0.41	0.23	0.25	0.20
a person looking at a camera on a tripod	0.29	0.31	0.30	0.25	0.34	0.40	0.31	0.30
a page of text about segmentation	0.31	0.30	0.23	0.24	0.33	0.28	0.47	0.25
a portrait of an astronaut with the American flag	0.23	0.23	0.29	0.23	0.29	0.26	0.26	0.37

Cosine similarity for RN50x16

								
a cup of coffee on a saucer	0.30	0.17	0.12	0.16	0.18	0.16	0.14	0.14
a black-and-white silhouette of a horse	0.14	0.36	0.12	0.15	0.13	0.15	0.20	0.10
a rocket standing on a launchpad	0.16	0.20	0.28	0.20	0.18	0.19	0.17	0.21
a red motorcycle standing in a garage	0.16	0.18	0.13	0.33	0.16	0.11	0.14	0.13
a facial photo of a tabby cat	0.16	0.18	0.14	0.15	0.29	0.12	0.13	0.12
a person looking at a camera on a tripod	0.18	0.20	0.19	0.15	0.24	0.29	0.21	0.17
a page of text about segmentation	0.20	0.24	0.14	0.14	0.24	0.17	0.37	0.14
a portrait of an astronaut with the American flag	0.10	0.11	0.19	0.11	0.14	0.15	0.16	0.27

Cosine similarity for RN50x64



Cosine similarity for ViT-B/16

								
a cup of coffee on a saucer	0.29	0.13	0.11	0.11	0.19	0.14	0.15	0.14
a black-and-white silhouette of a horse	0.15	0.36	0.15	0.14	0.16	0.19	0.17	0.13
a rocket standing on a launchpad	0.14	0.16	0.28	0.14	0.16	0.21	0.14	0.20
a red motorcycle standing in a garage	0.14	0.16	0.14	0.31	0.16	0.13	0.11	0.15
a facial photo of a tabby cat	0.17	0.17	0.14	0.15	0.32	0.14	0.15	0.16
a person looking at a camera on a tripod	0.15	0.18	0.18	0.16	0.23	0.29	0.18	0.18
a page of text about segmentation	0.18	0.23	0.15	0.13	0.21	0.18	0.35	0.16
a portrait of an astronaut with the American flag	0.15	0.16	0.21	0.13	0.16	0.18	0.13	0.28

Cosine similarity for ViT-B/32

								
a cup of coffee on a saucer	0.29	0.15	0.12	0.12	0.18	0.17	0.14	0.15
a black-and-white silhouette of a horse	0.15	0.35	0.15	0.17	0.15	0.21	0.17	0.11
a rocket standing on a launchpad	0.14	0.17	0.30	0.16	0.18	0.20	0.17	0.19
a red motorcycle standing in a garage	0.13	0.16	0.16	0.32	0.15	0.12	0.14	0.15
a facial photo of a tabby cat	0.17	0.15	0.12	0.12	0.31	0.16	0.12	0.12
a person looking at a camera on a tripod	0.14	0.20	0.21	0.16	0.21	0.30	0.19	0.19
a page of text about segmentation	0.20	0.20	0.16	0.16	0.20	0.20	0.35	0.15
a portrait of an astronaut with the American flag	0.15	0.16	0.22	0.15	0.17	0.17	0.13	0.28

Cosine similarity for ViT-L/14



Cosine similarity for ViT-L/14@336px

