

# LLM débarquement-model text image

## CLIP and co

Guillaume Behar et Adrien Zabban

mars 2024

# CLIP: connecter le texte et l'image

(1) Contrastive pre-training

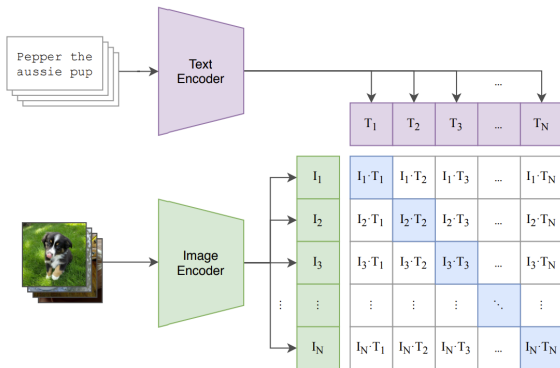


Figure: Modèle CLIP [1] (Open AI)

Proceedings of the 38 th International Conference on Machine Learning, PMLR 139, 2021

# Numpy-like pseudocode

```
# image_encoder - ResNet or Vision Transformer
# text_encoder  - CBOW or Text Transformer
# I[n, h, w, c] - minibatch of aligned images
# T[n, l]       - minibatch of aligned texts
# W_i[d_i, d_e] - learned proj of image to embed
# W_t[d_t, d_e] - learned proj of text to embed
# t            - learned temperature parameter

# extract feature representations of each modality
I_f = image_encoder(I) #[n, d_i]
T_f = text_encoder(T)  #[n, d_t]

# joint multimodal embedding [n, d_e]
I_e = l2_normalize(np.dot(I_f, W_i), axis=1)
T_e = l2_normalize(np.dot(T_f, W_t), axis=1)

# scaled pairwise cosine similarities [n, n]
logits = np.dot(I_e, T_e.T) * np.exp(t)

# symmetric loss function
labels = np.arange(n)
loss_i = cross_entropy_loss(logits, labels, axis=0)
loss_t = cross_entropy_loss(logits, labels, axis=1)
loss   = (loss_i + loss_t)/2
```

# Numpy-like pseudocode2

```

# image_encoder - ResNet or Vision Transformer
# text_encoder - CBOW or Text Transformer
# I[n, h, w, c] - minibatch of aligned images
# T[n, l] - minibatch of aligned texts
# W_i[d_i, d_e] - learned proj of image to embed
# W_t[d_t, d_e] - learned proj of text to embed
# t - learned temperature parameter

# extract feature representations of each modality
I_f = image_encoder(I) #[n, d_i]
T_f = text_encoder(T) #[n, d_t]

# joint multimodal embedding [n, d_e]
I_e = l2_normalize(np.dot(I_f, W_i), axis=1)
T_e = l2_normalize(np.dot(T_f, W_t), axis=1)

# scaled pairwise cosine similarities [n, n]
logits = np.dot(I_e, T_e.T) * np.exp(t)

# symmetric loss function
labels = np.arange(n)
loss_i = cross_entropy_loss(logits, labels, axis=0)
loss_t = cross_entropy_loss(logits, labels, axis=1)
loss = (loss_i + loss_t)/2

```

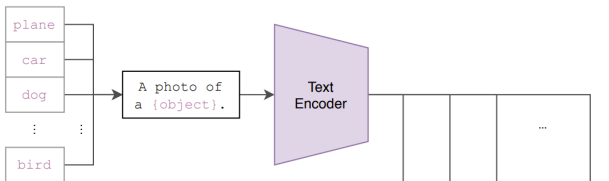
pseudocode

	T <sub>1</sub>	T <sub>2</sub>	T <sub>3</sub>	...	T <sub>N</sub>
I <sub>1</sub>	I <sub>1</sub> ·T <sub>1</sub>	I <sub>1</sub> ·T <sub>2</sub>	I <sub>1</sub> ·T <sub>3</sub>	...	I <sub>1</sub> ·T <sub>N</sub>
I <sub>2</sub>	I <sub>2</sub> ·T <sub>1</sub>	I <sub>2</sub> ·T <sub>2</sub>	I <sub>2</sub> ·T <sub>3</sub>	...	I <sub>2</sub> ·T <sub>N</sub>
I <sub>3</sub>	I <sub>3</sub> ·T <sub>1</sub>	I <sub>3</sub> ·T <sub>2</sub>	I <sub>3</sub> ·T <sub>3</sub>	...	I <sub>3</sub> ·T <sub>N</sub>
⋮	⋮	⋮	⋮	⋮	⋮
I <sub>N</sub>	I <sub>N</sub> ·T <sub>1</sub>	I <sub>N</sub> ·T <sub>2</sub>	I <sub>N</sub> ·T <sub>3</sub>	...	I <sub>N</sub> ·T <sub>N</sub>

CLIP

# CLIP: a zero-shot classifiers

(2) Create dataset classifier from label text



(3) Use for zero-shot prediction

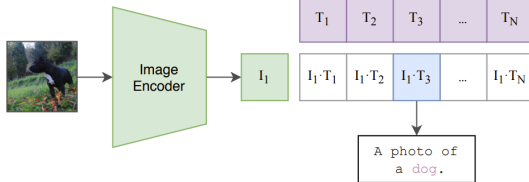
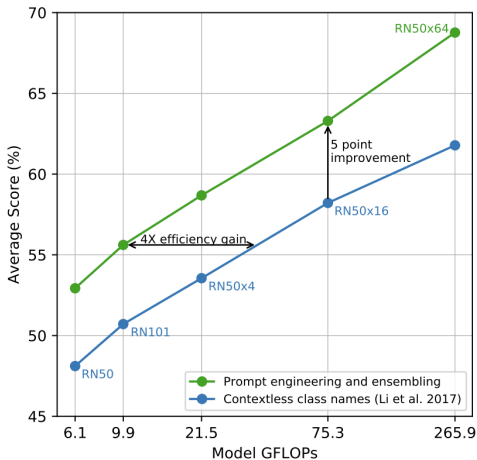


Figure: Clip classifiers

# CLIP: a zero-shot classifiers



**Figure:** Prompt engineering and ensembling improve zeroshot performance.

# CLIP: zero-shot vs few-shot linear probes

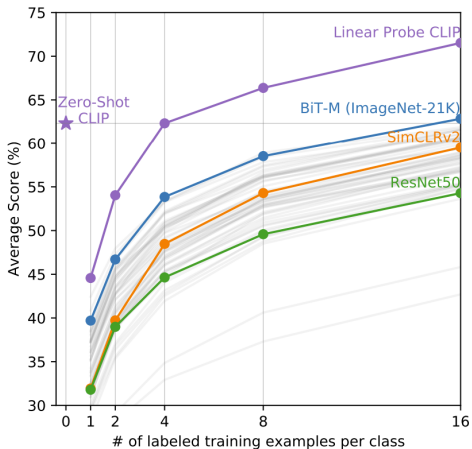


Figure: Zero-shot CLIP outperforms few-shot linear probes

## unCLIP avec DALL-E

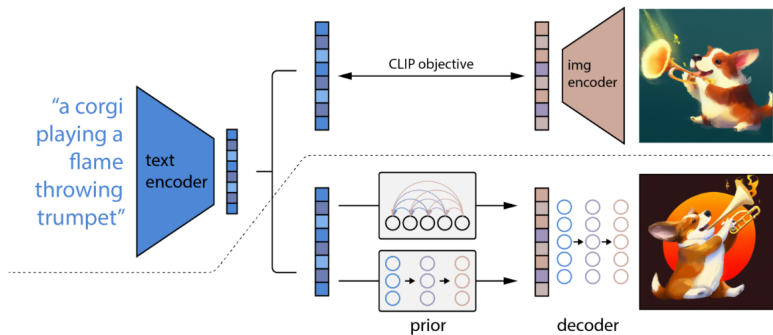


Figure: Modèle DALL-E [2]



# Références

- [1] A. Radford, J. W. Kim, C. Hallacy, *et al.*, “Learning transferable visual models from natural language supervision,” *CoRR*, vol. abs/2103.00020, 2021. arXiv: 2103.00020. [Online]. Available: <https://arxiv.org/abs/2103.00020>.
- [2] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen, *Hierarchical text-conditional image generation with clip latents*, 2022. arXiv: 2204.06125 [cs.CV].