

Segment Anything - SAM

Méthodes Avancées en Segmentation d'Images

Strasbourg
23.02.2024

Robert Maria

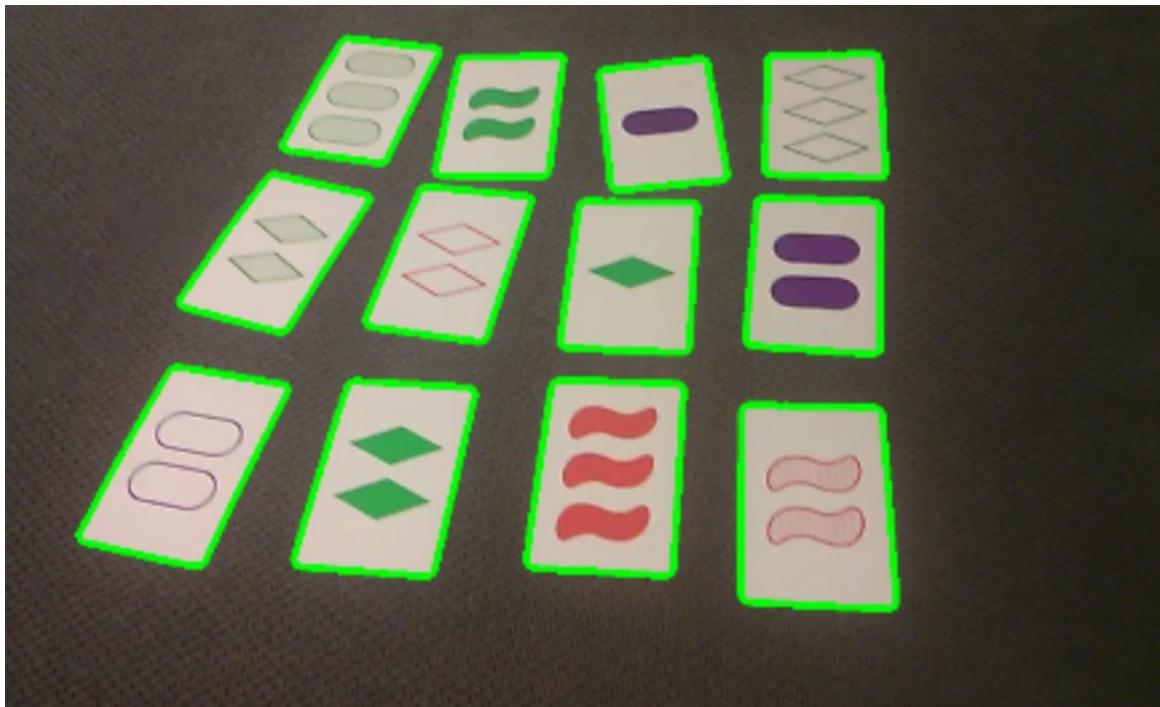
La segmentation d'image



- diviser une image en parties ou segments

La segmentation d'image

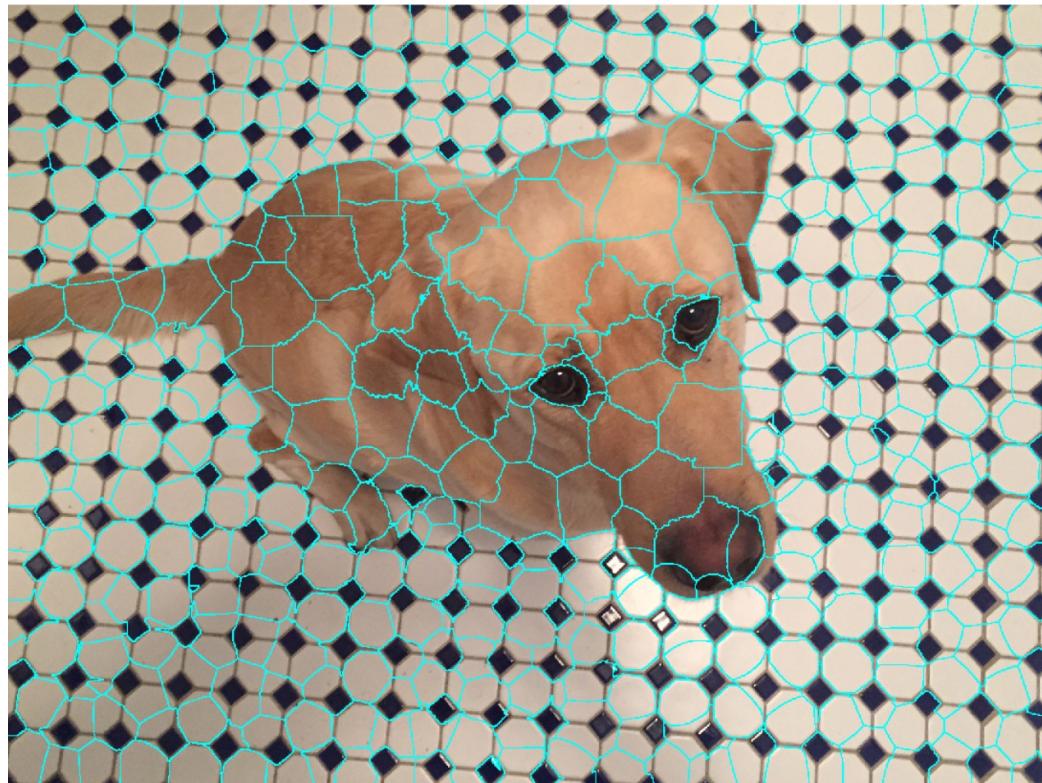
1. Détection de Contours



- étape préliminaire dans de nombreuses tâches de segmentation
- délimiter les frontières avant un traitement plus approfondi

La segmentation d'image

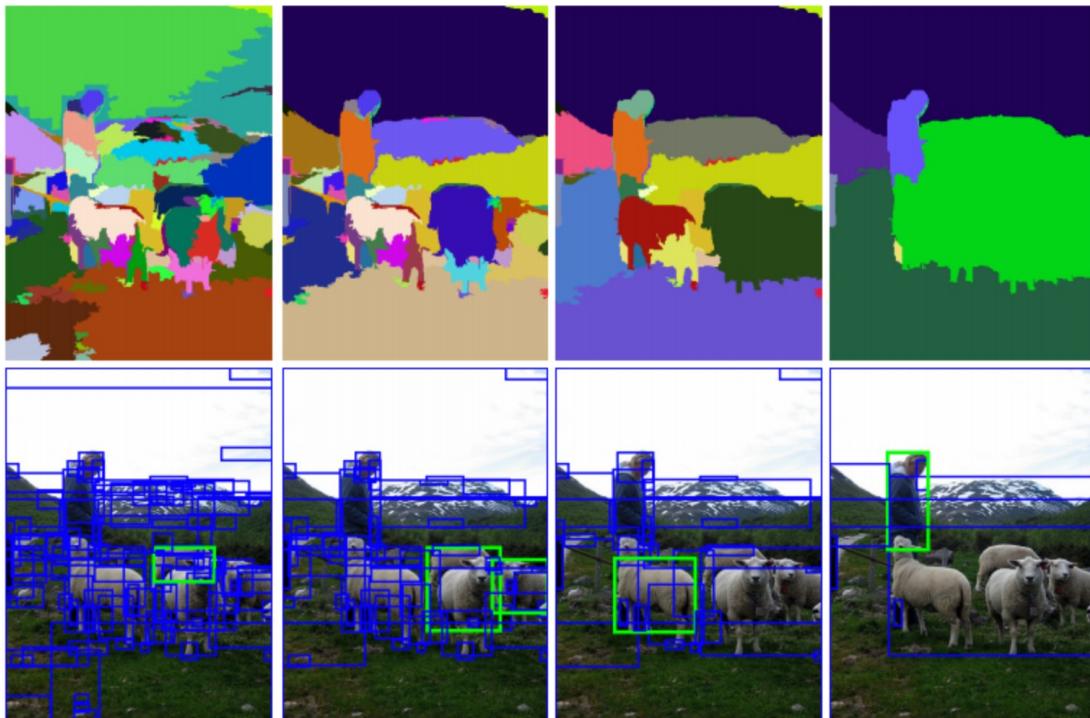
2. La superpixelisation



- regrouper des pixels voisins (ayant des caractéristiques similaires) en unités plus grandes appelées superpixels.

La segmentation d'image

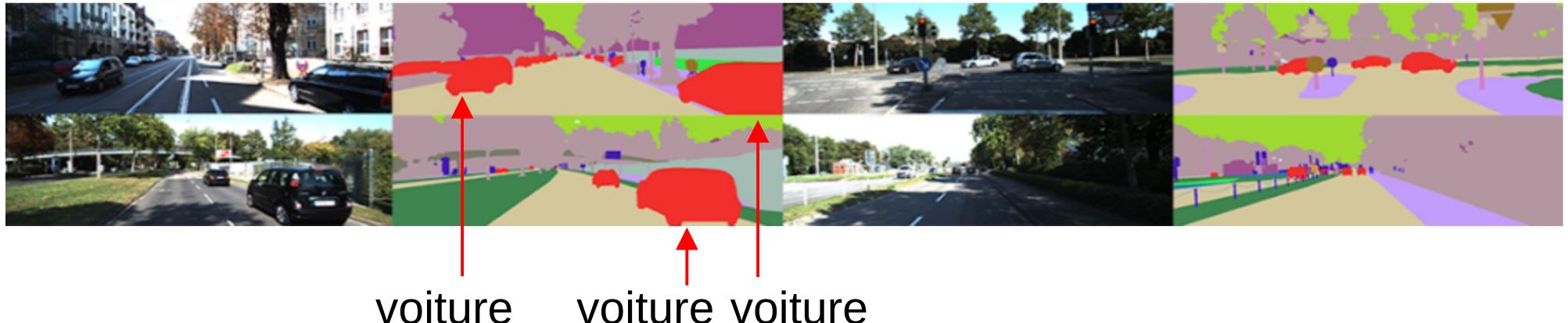
3. La proposition d'objets ("object proposal")



- génération de régions candidates
- Ex: selective search

La segmentation d'image

4. Segmentation Sémantique



chaque pixel de l'image est classé dans une catégorie spécifique

La segmentation d'image

5. Segmentation par Instance



- distinguer entre différentes instances de la même classe

La segmentation d'image

5. Segmentation Panoptique



- Segmentation par instance + Segmentation sémantique

Objectifs - SAM

Zero-shot learning

Classes non vues lors de l'entraînement

Few-shots learning

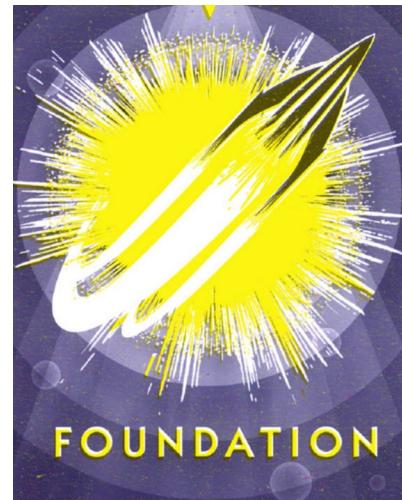
Petit nombre d'exemples (1 ou 2)
pour chaque classe

Flexible prompting

Points, rectangles, texte

Flexibilité à l'ambiguïté

Un nouveau modèle de fondation



source

Modèles de fondation

Succès = $f(\text{tâche}, \text{modèle}, \text{données})$

Modèles de fondation

Succès = $f(tâche, modèle, données)$

tâche = $f(\text{énergie, calcul, mémoire, temps})$

tâche = meilleure précision avec la **charge de ressources**

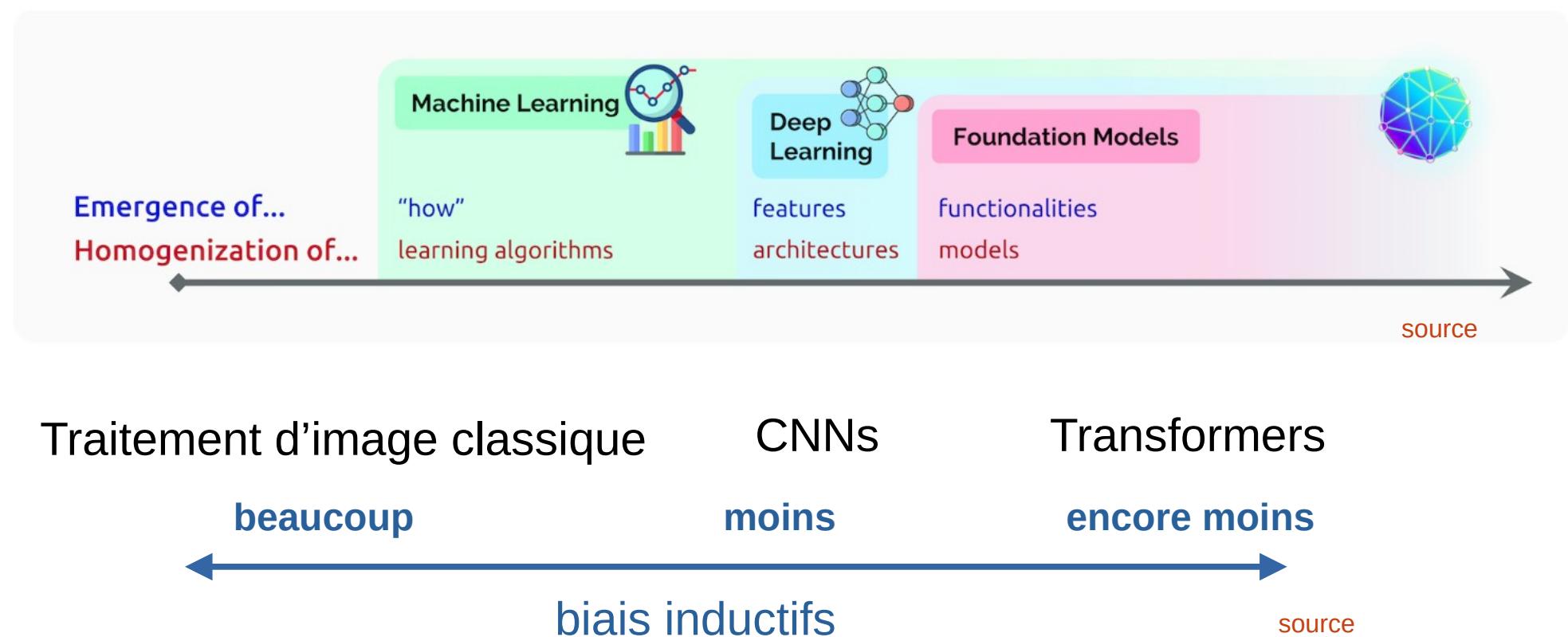


change avec le temps

Modèles de fondation

L'IA connaît un changement de paradigme

émergence de modèles (BERT, DALL-E, GPT-3, CLIP) qui peuvent être adaptés à une vaste gamme de tâches en aval.



Modèles de fondation

L'émergence

le comportement d'un système est induit implicitement plutôt qu'explicitement construit

L'homogénéisation

consolidation des méthodologies pour construire des systèmes à travers un large éventail d'applications

levier important pour de nombreuses tâches

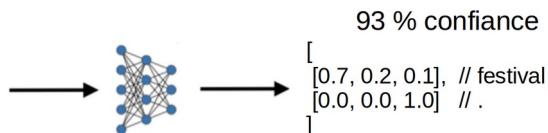
points uniques de défaillance

Modèles de fondation

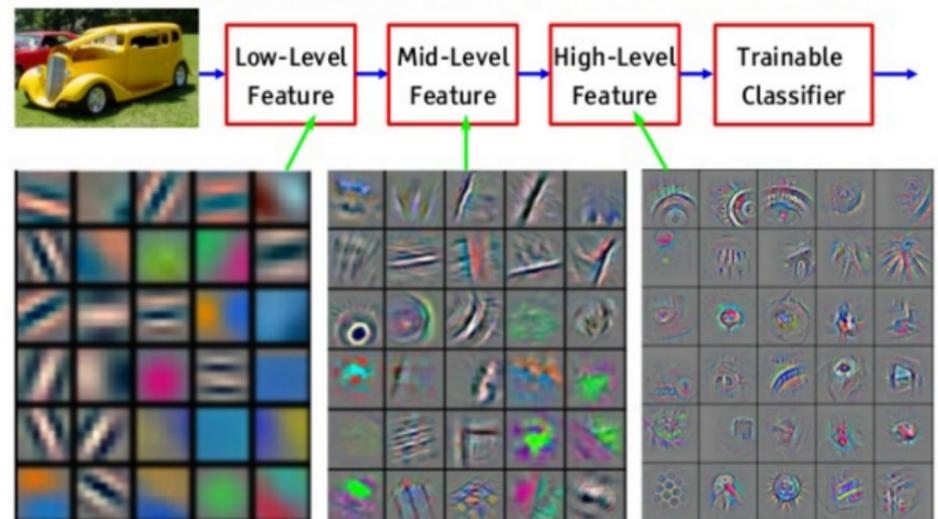
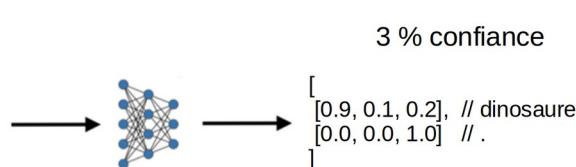
Succès = f(tâche, modèle, données)

“Modèles de fondation”

```
[  
  [0.2, 0.1, 0.7], // Le  
  [0.4, 0.3, 0.8], // château  
  [0.1, 0.5, 0.6], // est  
  [0.6, 0.4, 0.3], // magnifique  
  [0.3, 0.7, 0.2], // ment  
  [0.5, 0.2, 0.9], // illuminé  
  [0.2, 0.6, 0.4], // pour  
  [0.2, 0.1, 0.7], // le  
]
```



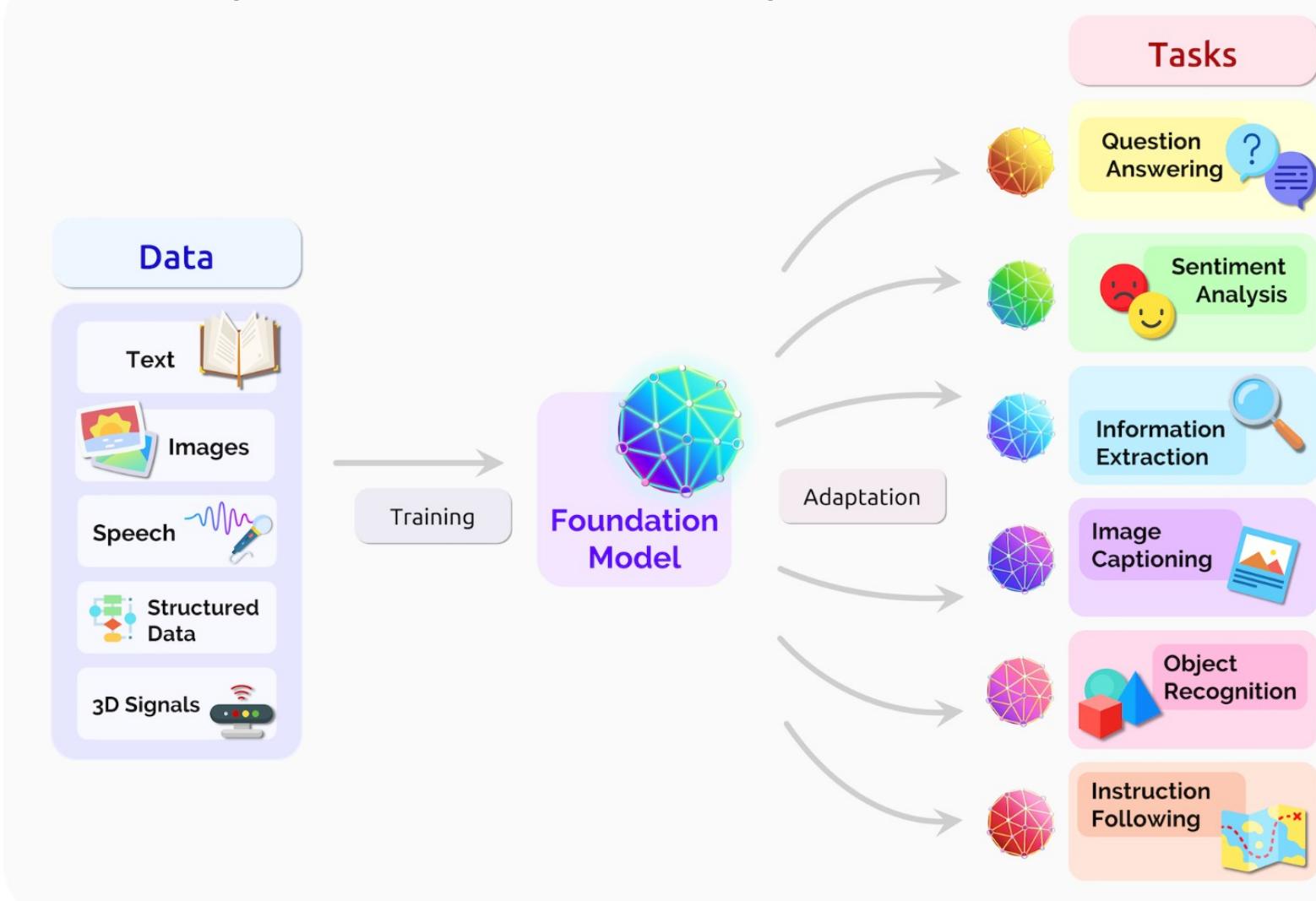
```
[  
  [0.2, 0.1, 0.7], // Le  
  [0.4, 0.3, 0.8], // château  
  [0.1, 0.5, 0.6], // est  
  [0.6, 0.4, 0.3], // magnifique  
  [0.3, 0.7, 0.2], // ment  
  [0.5, 0.2, 0.9], // illuminé  
  [0.2, 0.6, 0.4], // pour  
  [0.2, 0.1, 0.7], // le  
]
```



apprentissage par transfert + échelle
préentraînement + fine-tuning

Modèles de fondation

Succès = f(tâche, modèle, données)



source

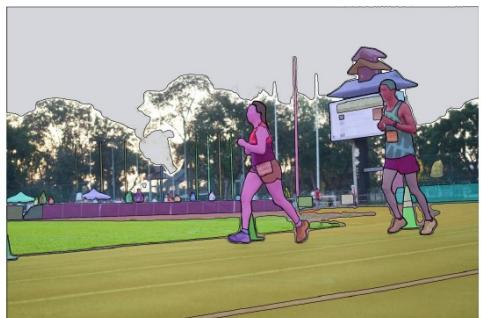
Segment Anything (SAM)

Modèle de fondation dans la segmentation d'images

Segment Anything (SAM)

SA-1B dataset release:

50-100 masks



100-200 masks

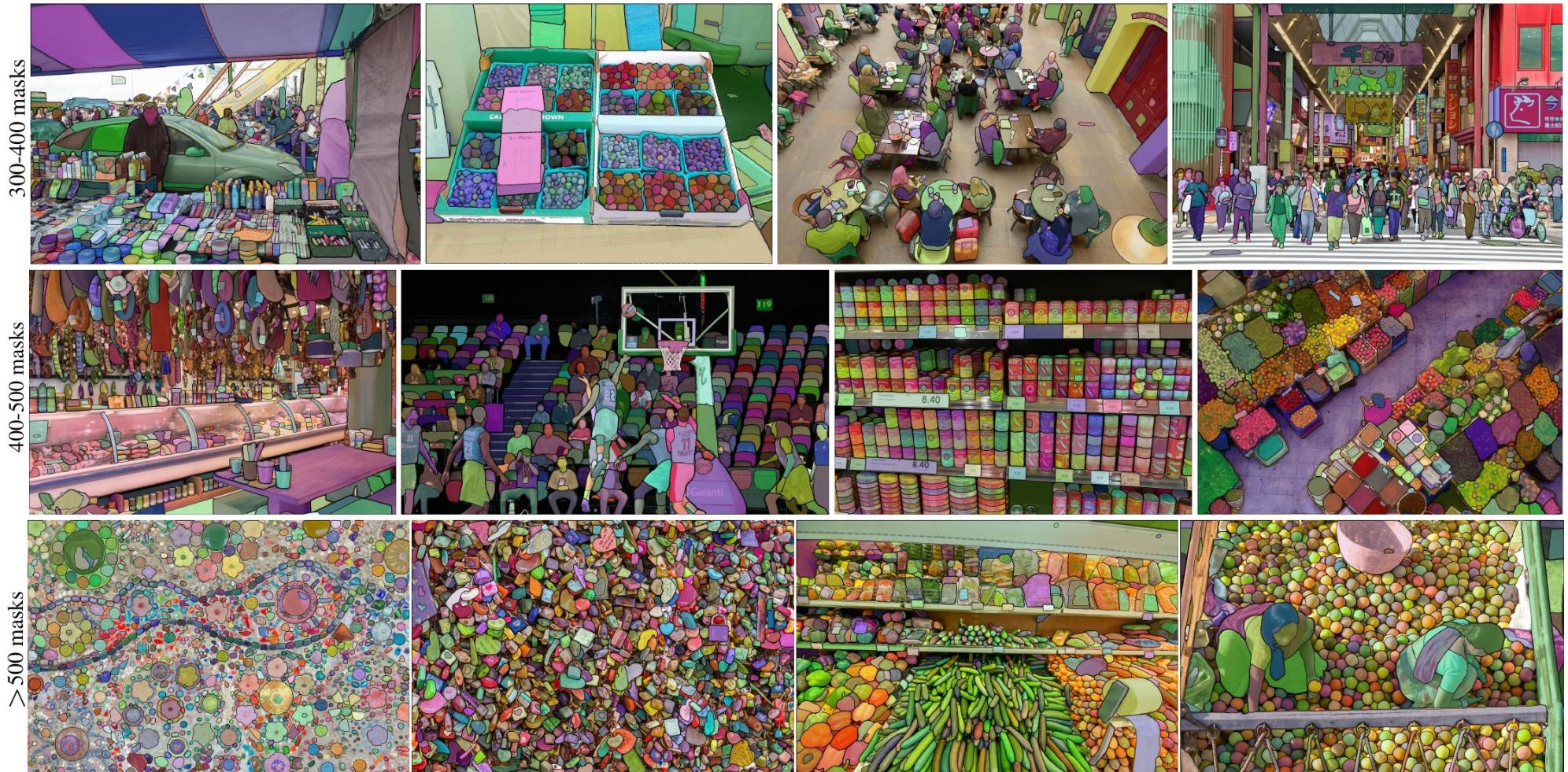


200-300 masks

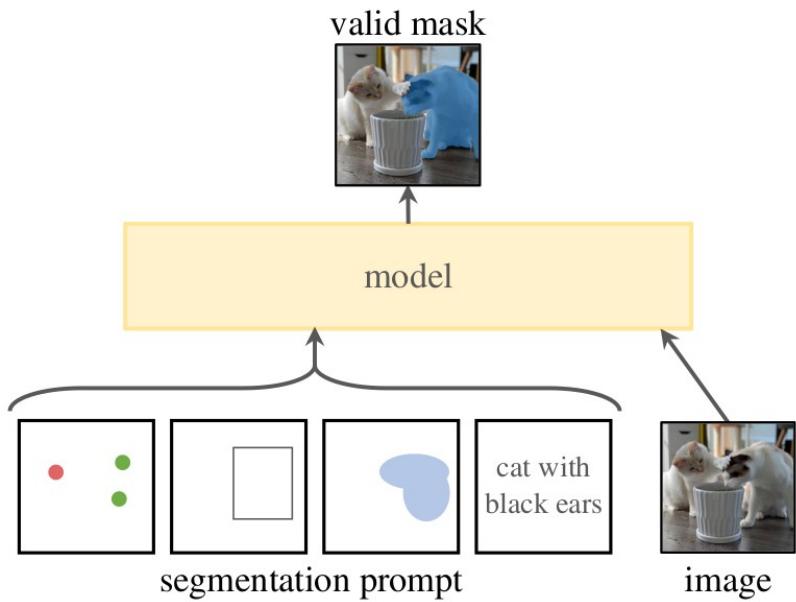


Segment Anything (SAM)

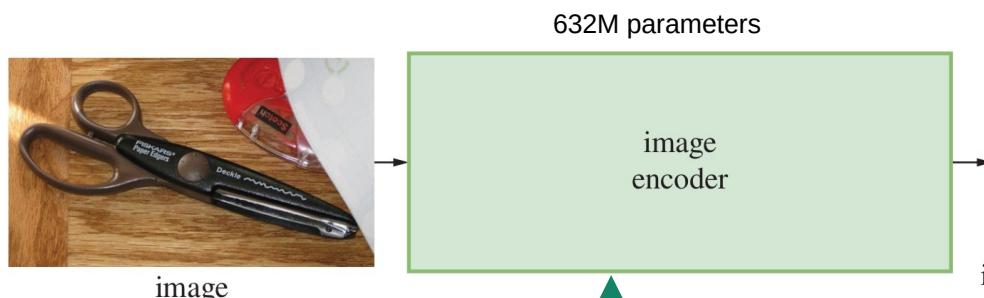
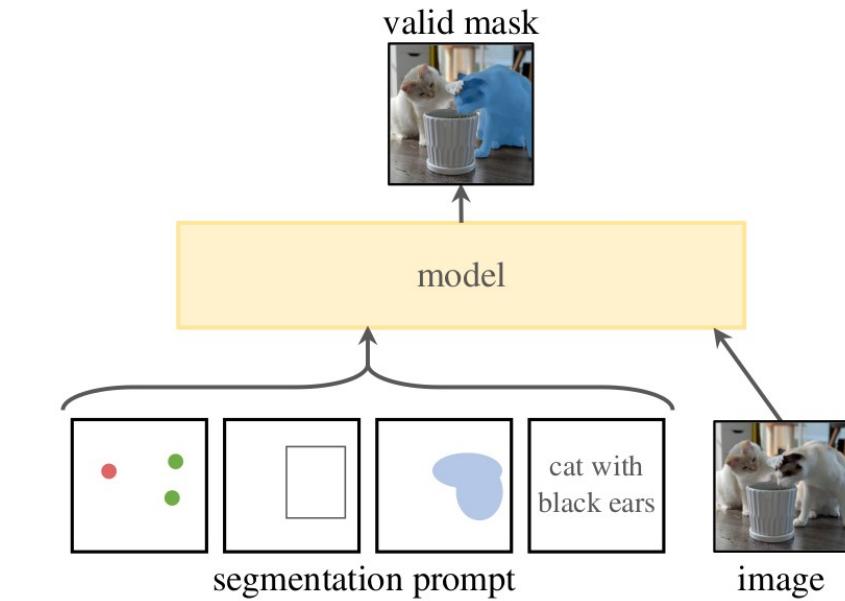
SA-1B dataset release:



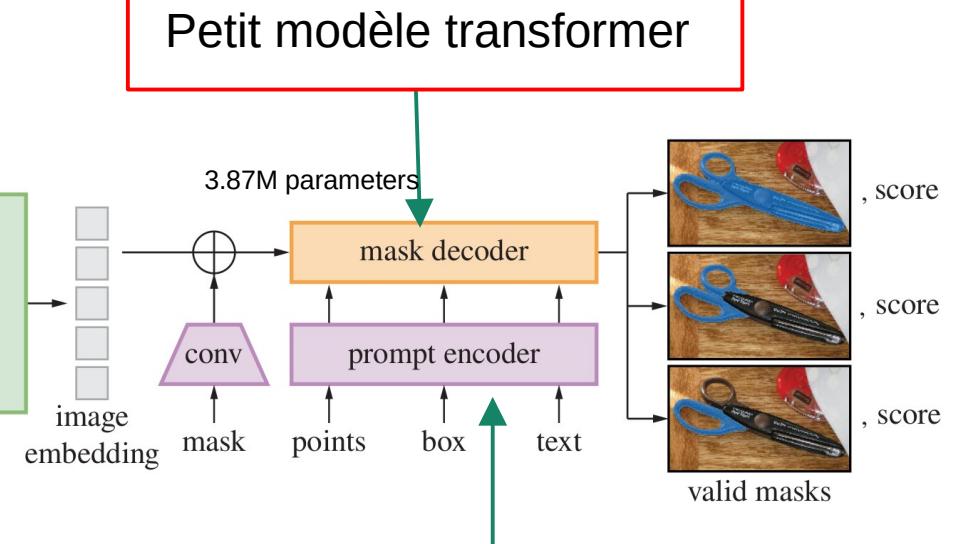
Segment Anything (SAM)



Segment Anything (SAM)

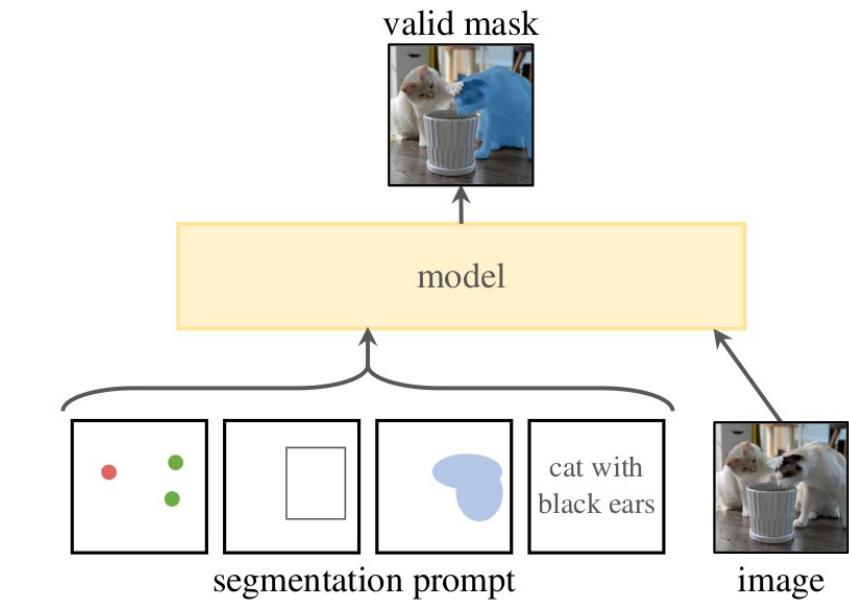


Convertit l'image en une
représentation vectorielle



Convertit le prompt en une
représentation vectorielle

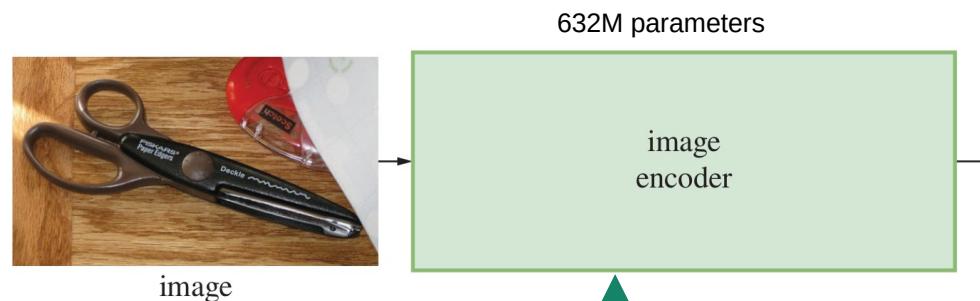
Segment Anything (SAM)



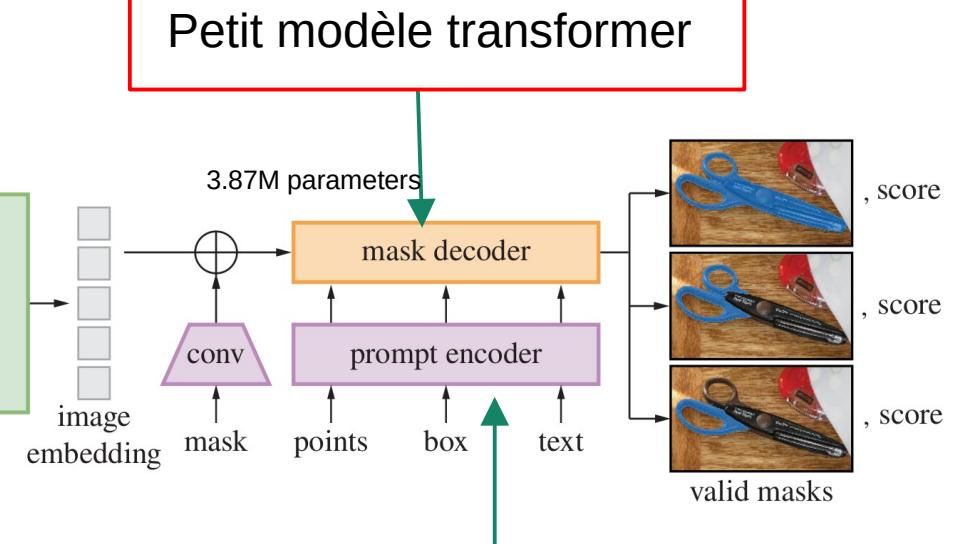
Encodeur d'image = ViT pré-entraîné MAE

ViT = Visual Transformer

MAE = Masked autoencoders



Convertit l'image en une
représentation vectorielle

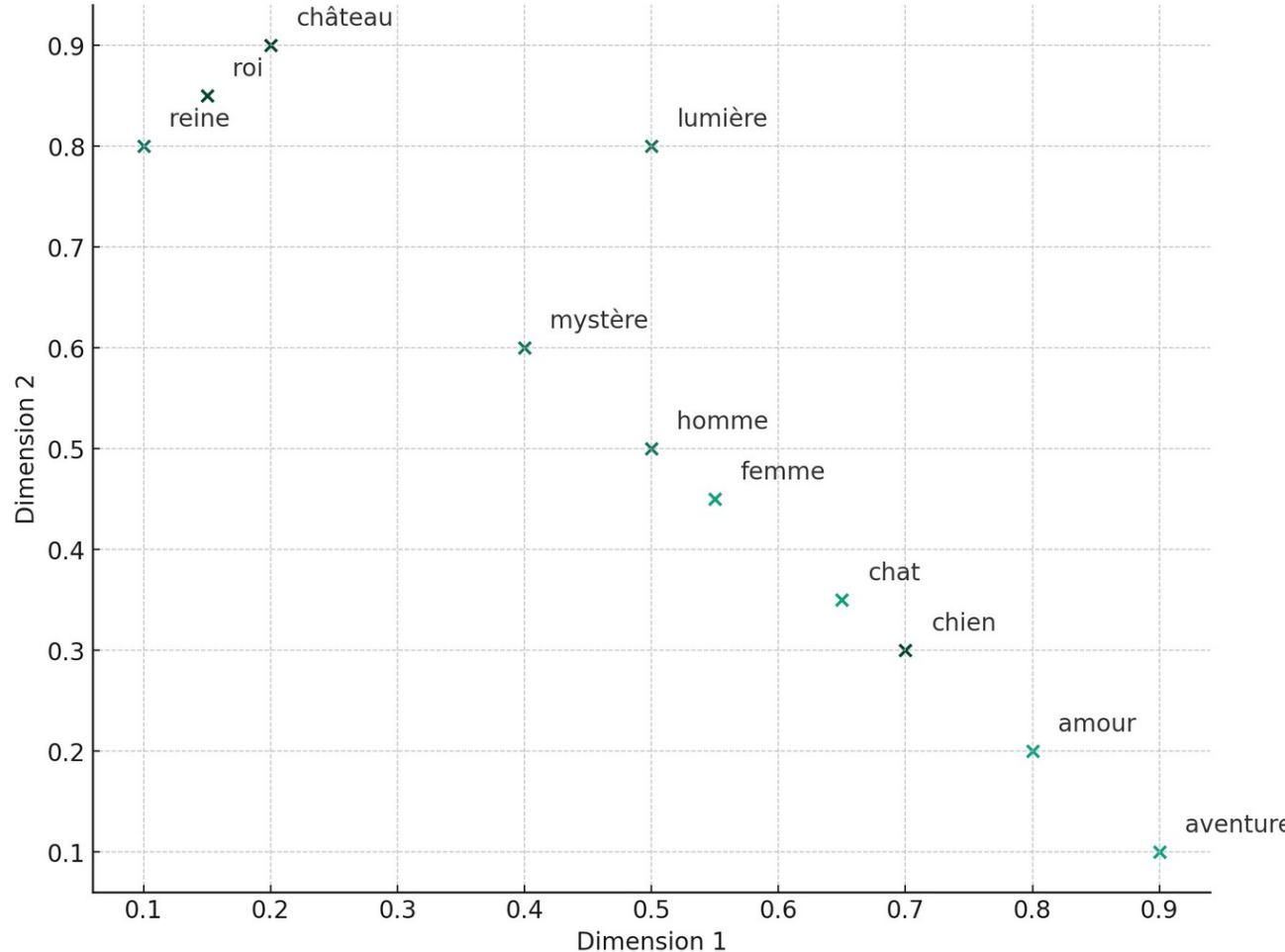


Convertit le prompt en une
représentation vectorielle

Embeddings

Embeddings = représentations vectorielles d'images, de texte

Embeddings



Roi – homme + femme = reine

$$\begin{aligned} \text{roi} &= [0.9, 0.1] \\ \text{homme} &= [0.5, 0.5] \\ \text{femme} &= [0.4, 0.6] \\ \text{reine} &= [0.8, 0.2] \end{aligned}$$

$$\begin{aligned} \text{roi} - \text{homme} + \text{femme} \\ &= [0.9, 0.1] - [0.5, 0.5] + [0.4, 0.6] \\ &= [0.8, 0.2] \end{aligned}$$

Transformers

Transformers

Attention Is All You Need

Ashish Vaswani*

Google Brain

avaswani@google.com

Noam Shazeer*

Google Brain

noam@google.com

Niki Parmar*

Google Research

nikip@google.com

Jakob Uszkoreit*

Google Research

usz@google.com

Llion Jones*

Google Research

llion@google.com

Aidan N. Gomez* [†]

University of Toronto

aidan@cs.toronto.edu

Lukasz Kaiser*

Google Brain

lukaszkaiser@google.com

Illia Polosukhin* [‡]

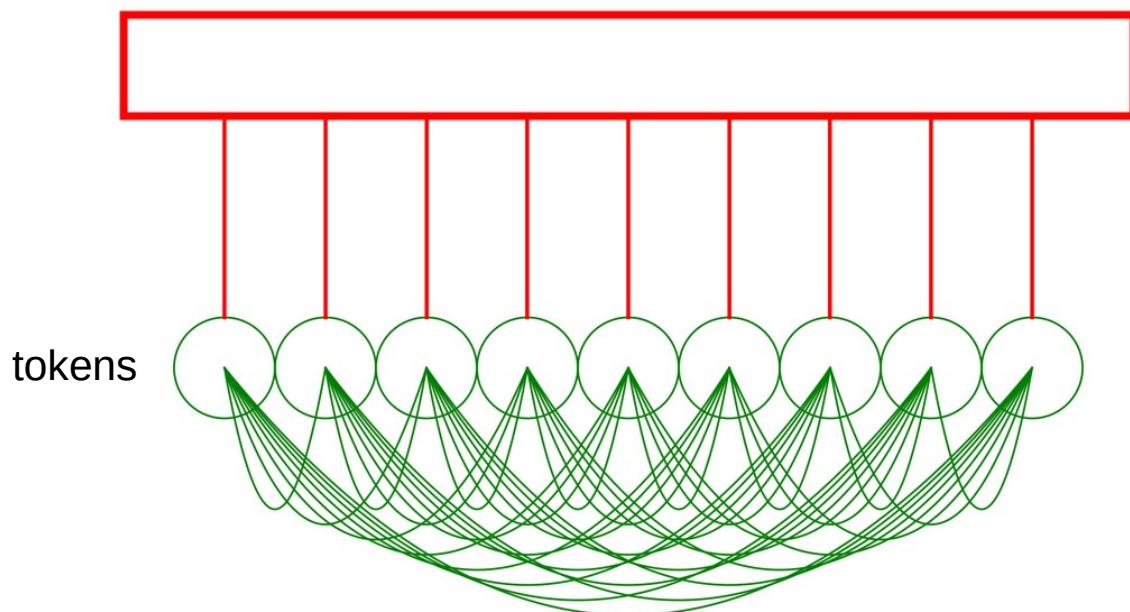
illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

Transformers

Text



Attention = pairwise inner-product

Quadratic operation

Images

0	2	15	0	0	11	10	0	0	0	0	9	9	0	0	0	0
0	0	0	4	60	157	236	255	255	177	95	61	32	0	0	29	
0	10	16	119	238	255	244	245	243	250	349	255	222	103	10	0	
0	14	170	255	255	244	254	255	253	245	255	249	253	251	124	1	
2	98	255	228	255	251	254	211	141	116	122	715	251	238	255	49	
13	217	243	255	155	33	226	52	2	0	10	13	232	255	255	36	
16	229	252	254	49	12	0	0	7	7	0	70	237	252	236	62	
6	141	245	255	217	25	11	9	3	0	115	236	243	255	137	0	
0	87	252	250	248	215	60	0	1	121	252	255	248	144	6	0	
0	13	113	255	255	245	255	182	181	248	252	242	208	36	0	19	
1	0	5	117	251	255	241	255	247	255	241	162	17	0	7	0	
0	0	0	4	58	251	255	246	254	253	255	129	11	0	1	0	
0	0	4	97	255	255	255	248	257	255	244	255	182	10	0	4	
0	22	206	252	246	281	241	100	24	113	255	245	255	194	9	0	
0	111	255	242	255	158	24	0	0	6	39	255	232	230	56	0	
0	218	251	250	137	7	11	0	0	2	62	255	250	125	3		
0	173	255	255	101	9	20	0	13	3	33	182	251	245	61	0	
0	107	251	241	255	230	98	55	19	118	217	248	253	255	52	4	
0	18	146	250	255	247	255	255	255	249	255	240	255	129	0	5	
0	0	23	113	215	255	250	248	255	255	248	245	118	14	12	0	
0	0	6	1	0	52	153	233	256	252	187	37	0	0	4	1	
0	0	5	5	0	0	0	0	0	14	1	0	6	5	0	0	

Visual Transformer (ViT)

Published as a conference paper at ICLR 2021

AN IMAGE IS WORTH 16X16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

**Alexey Dosovitskiy^{*,†}, Lucas Beyer^{*}, Alexander Kolesnikov^{*}, Dirk Weissenborn^{*},
Xiaohua Zhai^{*}, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer,
Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, Neil Houlsby^{*,†}**

^{*}equal technical contribution, [†]equal advising

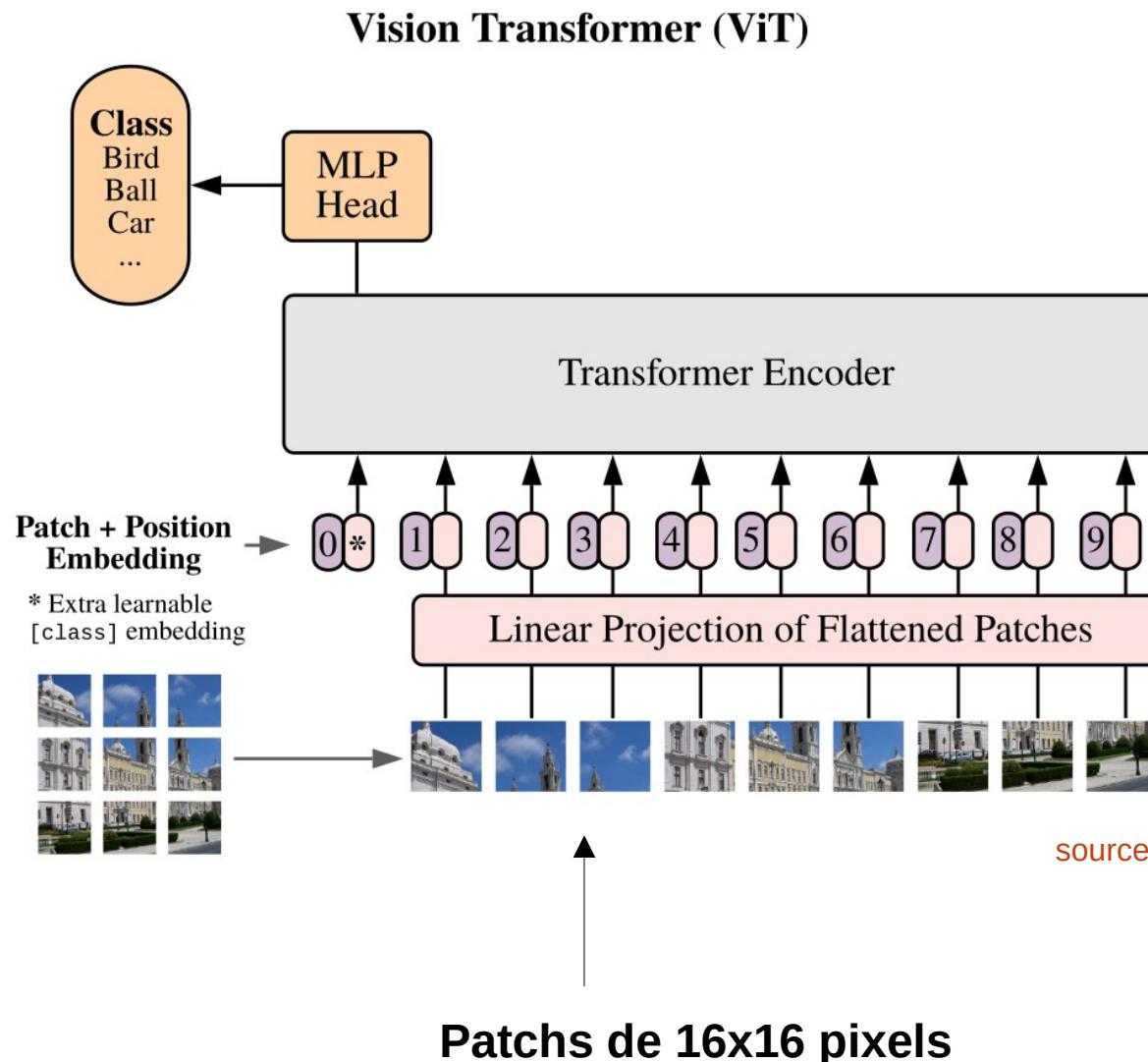
Google Research, Brain Team

{adosovitskiy, neilhoulsby}@google.com

ABSTRACT

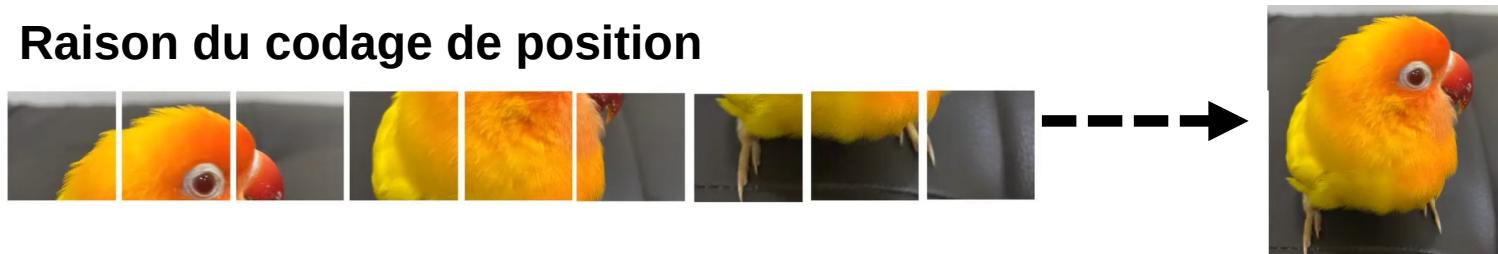
While the Transformer architecture has become the de-facto standard for natural language processing tasks, its applications to computer vision remain limited. In vision, attention is either applied in conjunction with convolutional networks, or used to replace certain components of convolutional networks while keeping their overall structure in place. We show that this reliance on CNNs is not necessary and a pure transformer applied directly to sequences of image patches can perform very well on image classification tasks. When pre-trained on large amounts of data and transferred to multiple mid-sized or small image recognition benchmarks (ImageNet, CIFAR-100, VTAB, etc.), Vision Transformer (ViT) attains excellent results compared to state-of-the-art convolutional networks while requiring substantially fewer computational resources to train.^[1]

Visual Transformer (ViT)

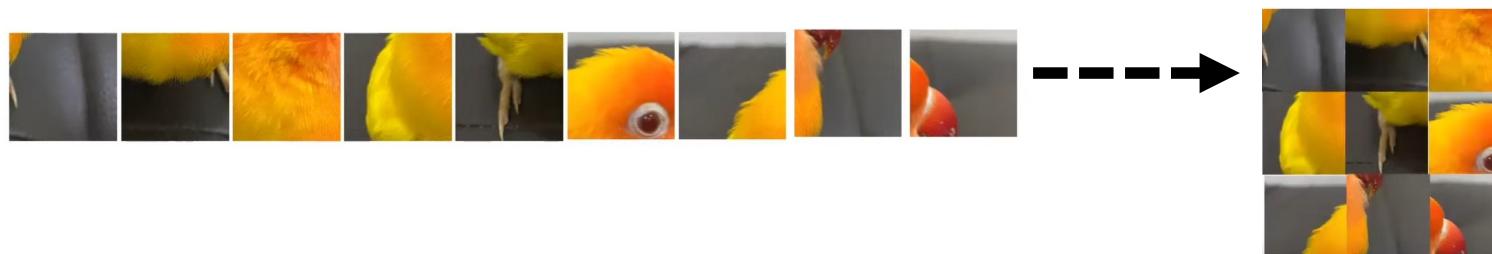


Visual Transformer (ViT)

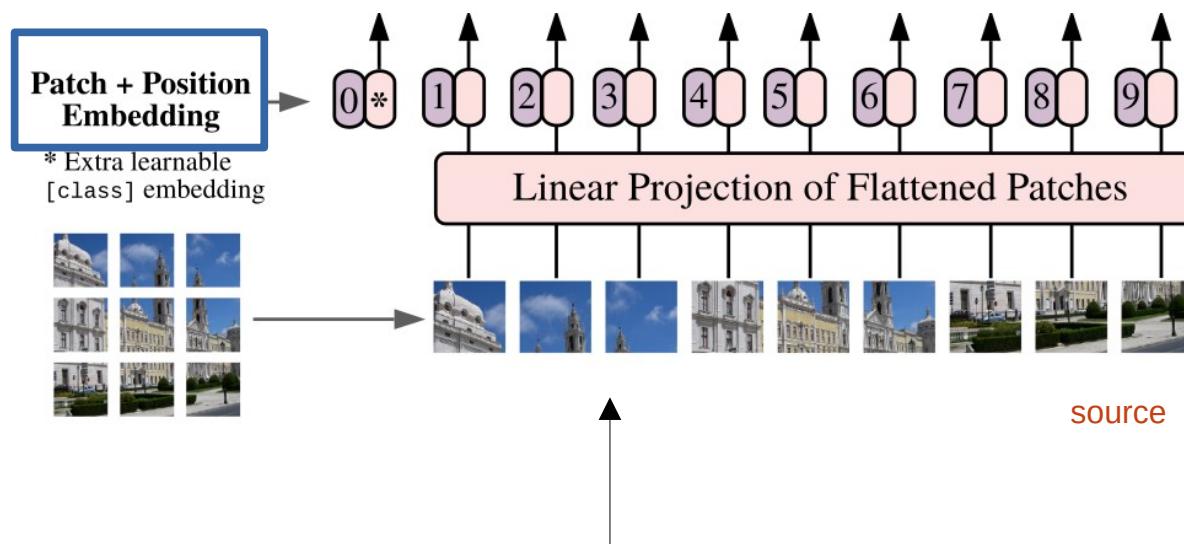
Raison du codage de position



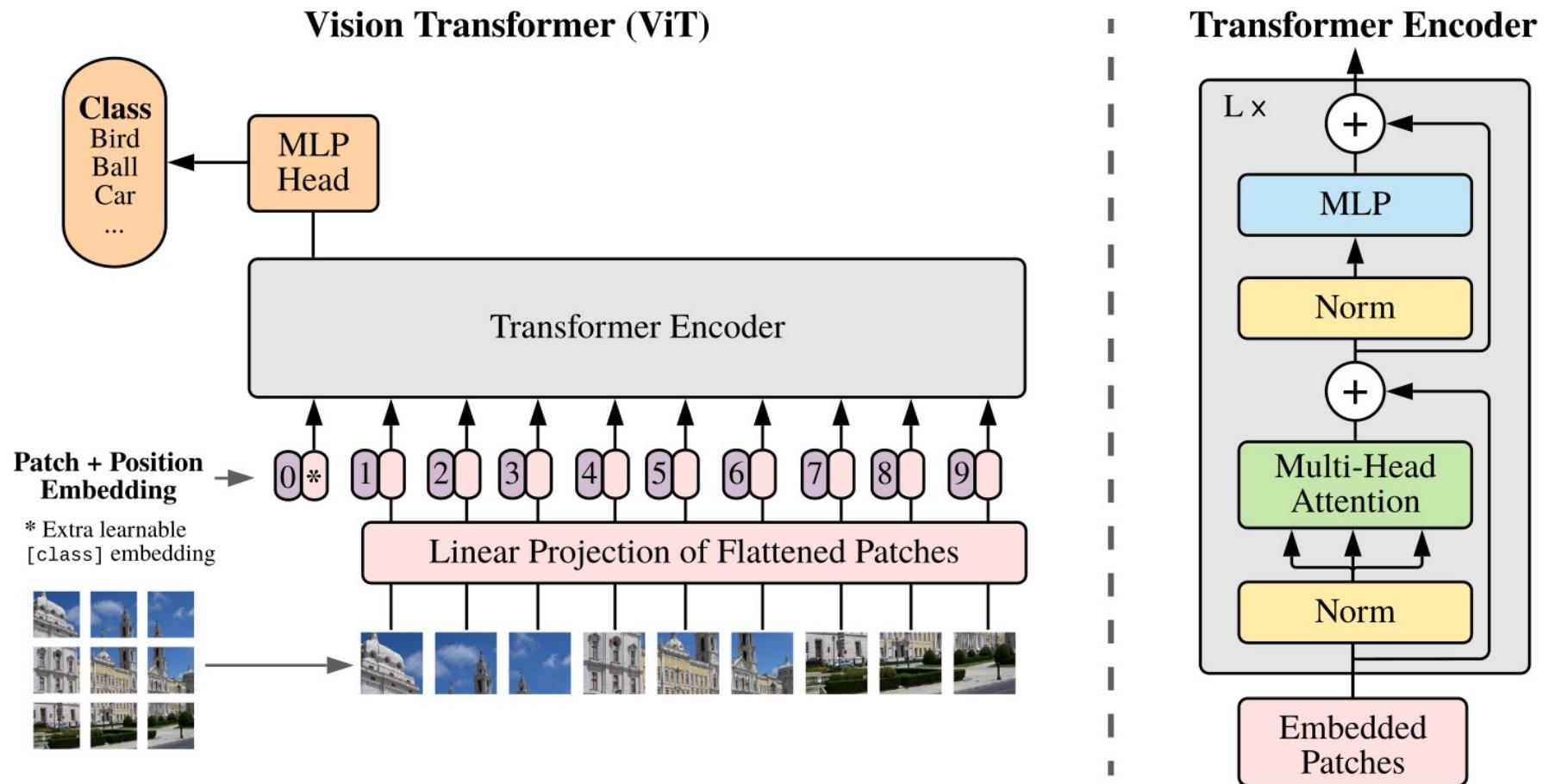
Oiseau



Oiseau ?



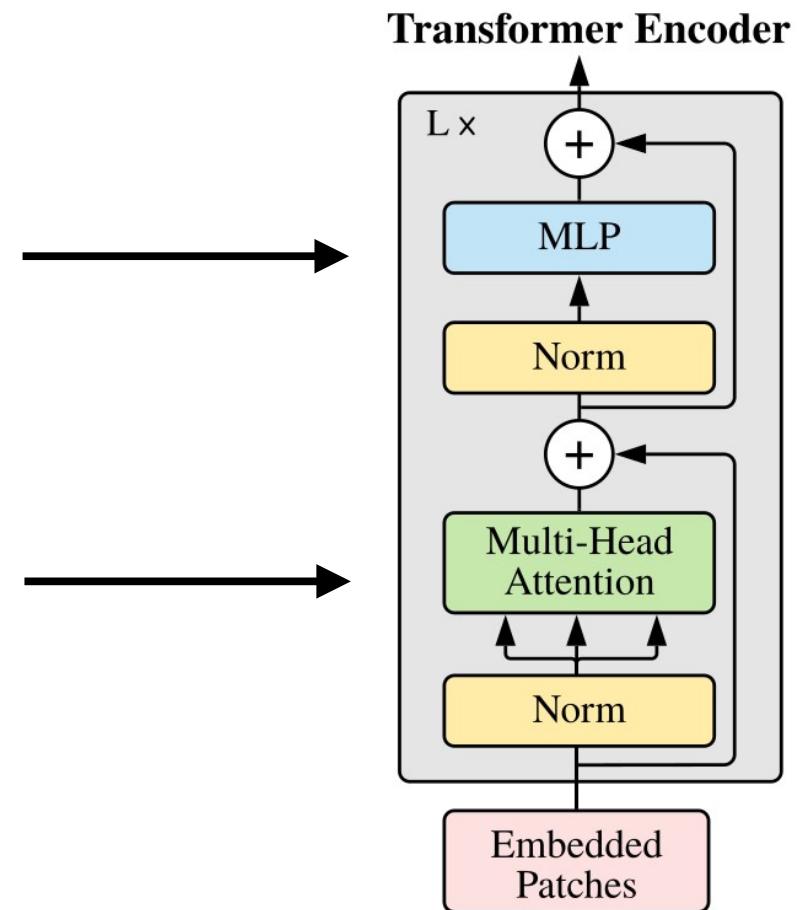
Visual Transformer (ViT)



Visual Transformer (ViT)

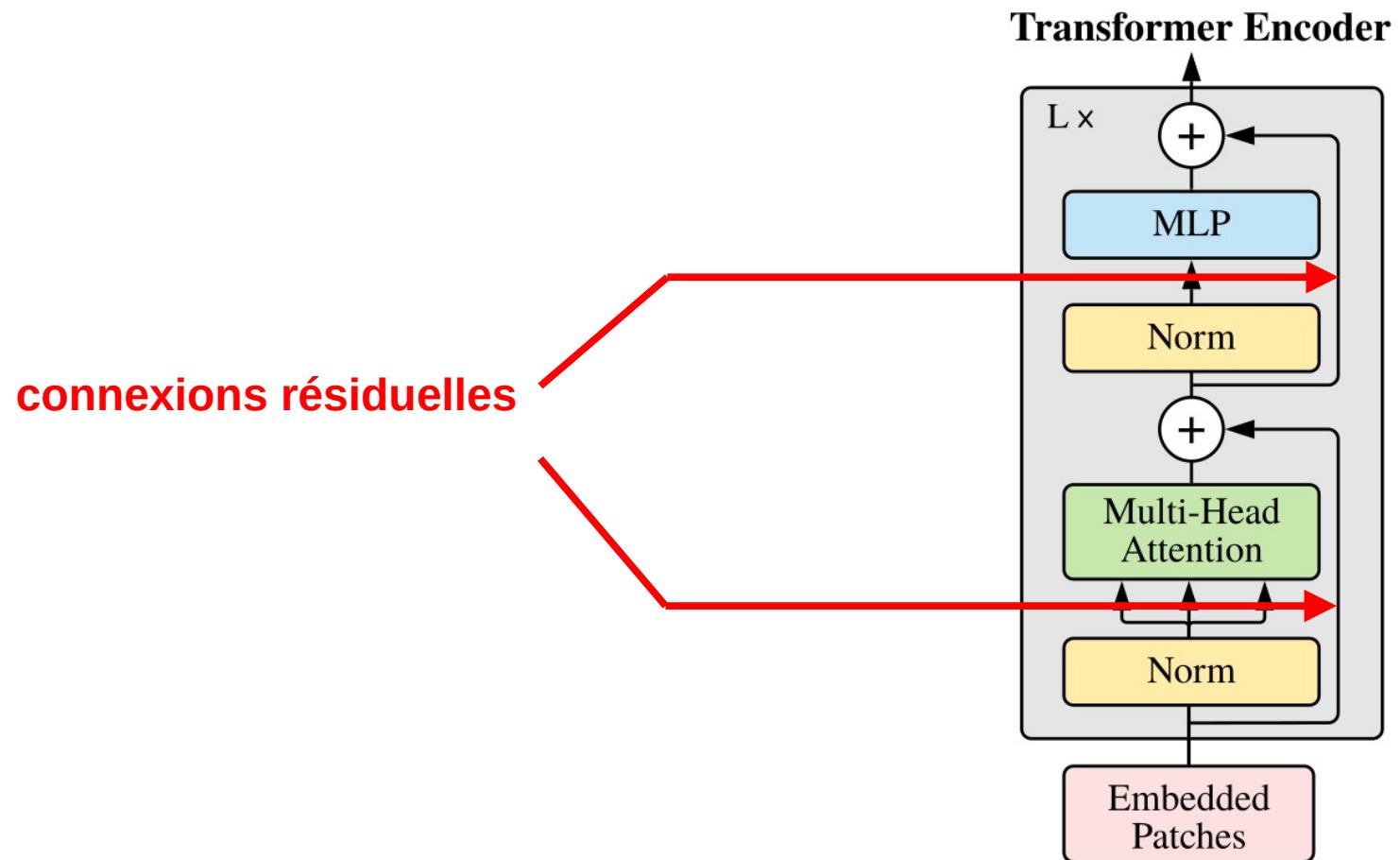
permet aux embeddings de
« penser » de manière
indépendante

permet aux embeddings de
communiquer



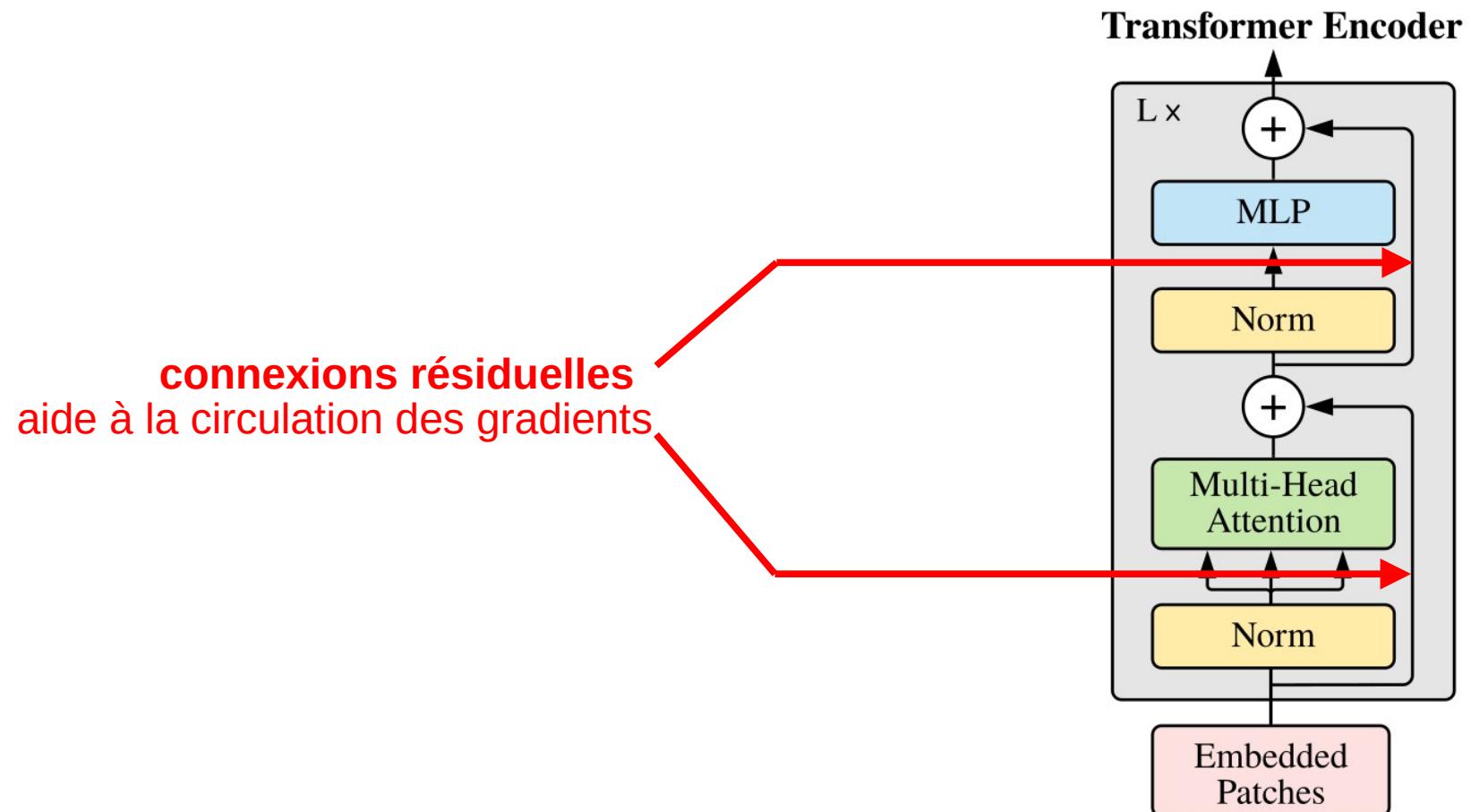
source

Visual Transformer (ViT)



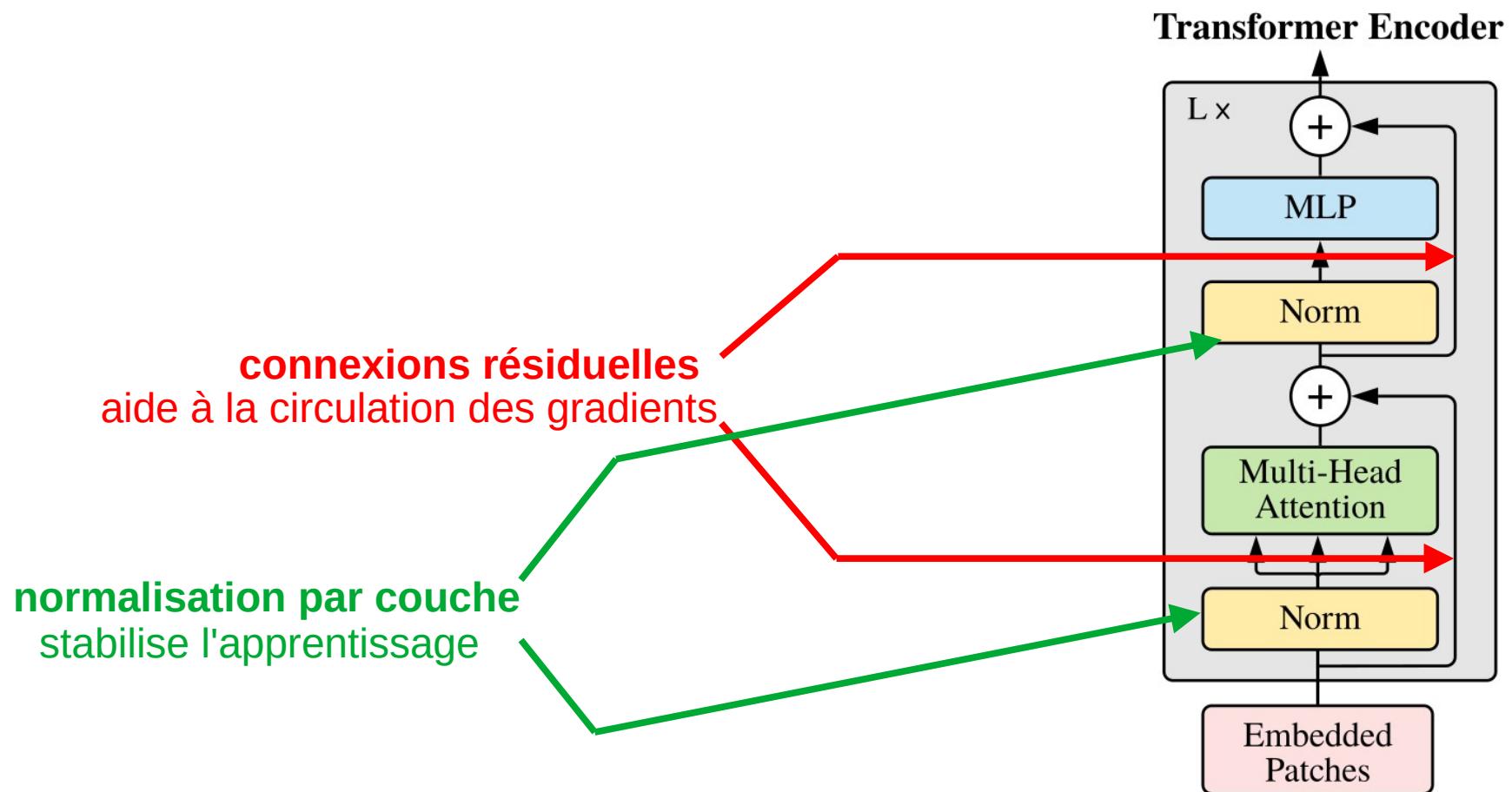
source

Visual Transformer (ViT)



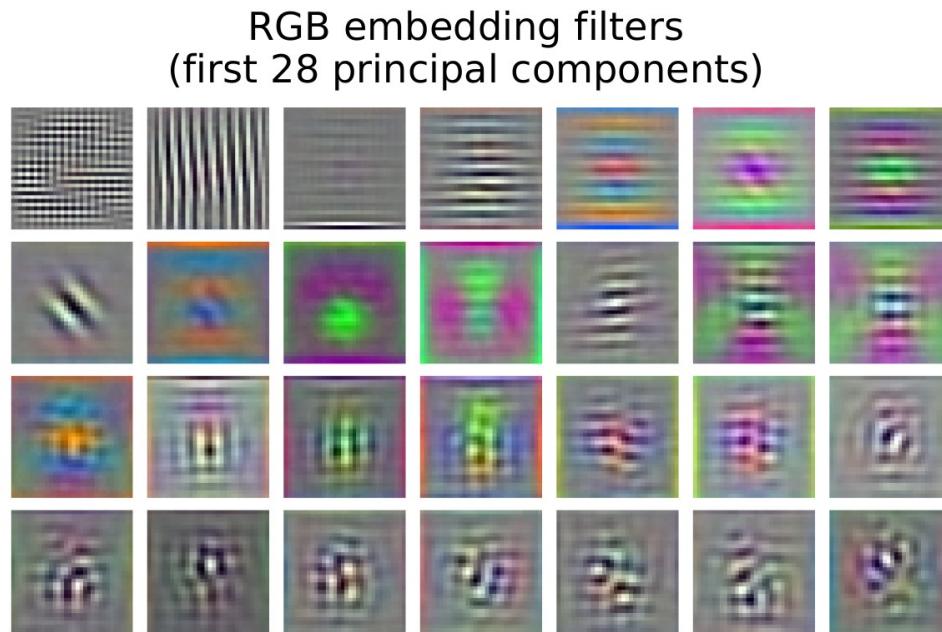
source

Visual Transformer (ViT)

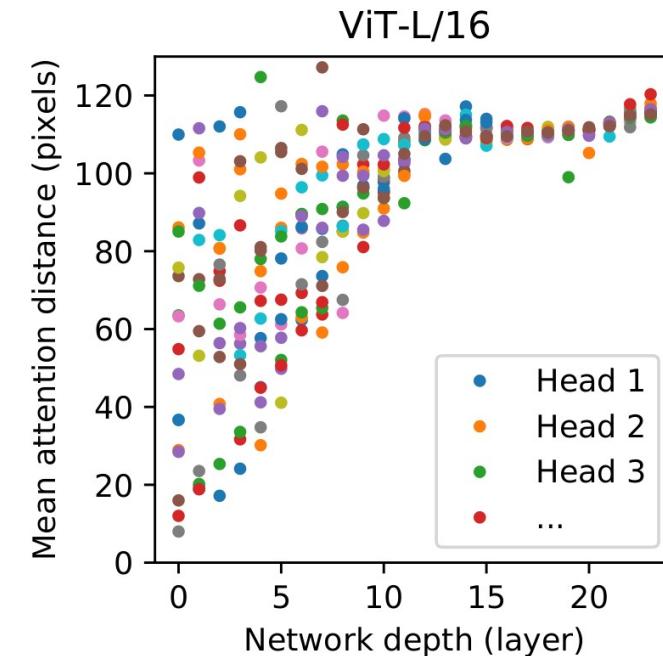
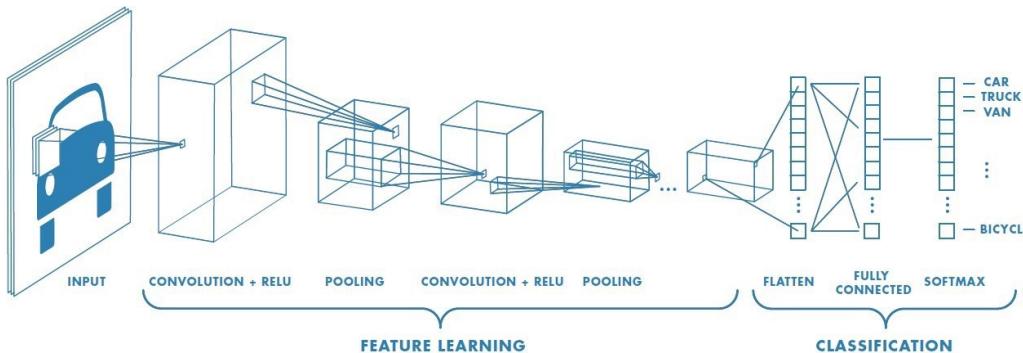


source

Visual Transformer (ViT)



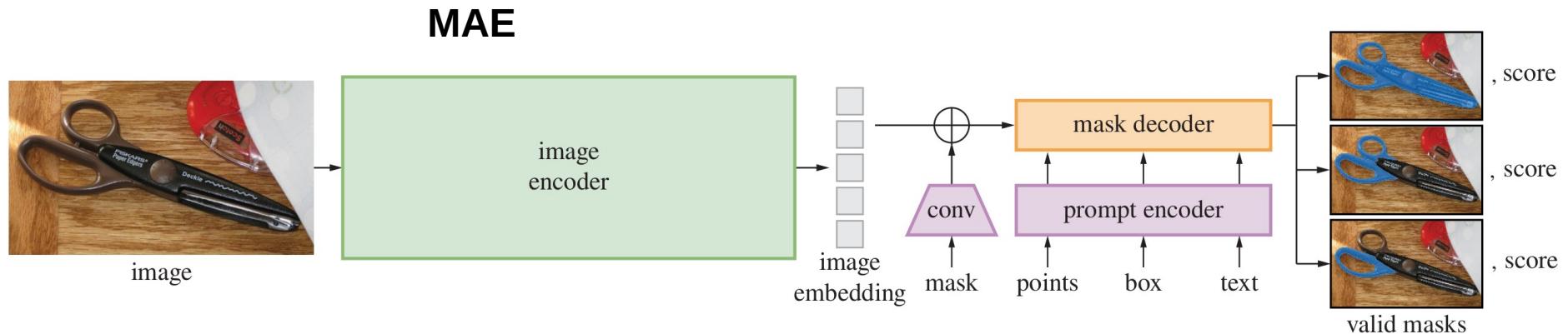
Feature maps ressemblant à des CNNs



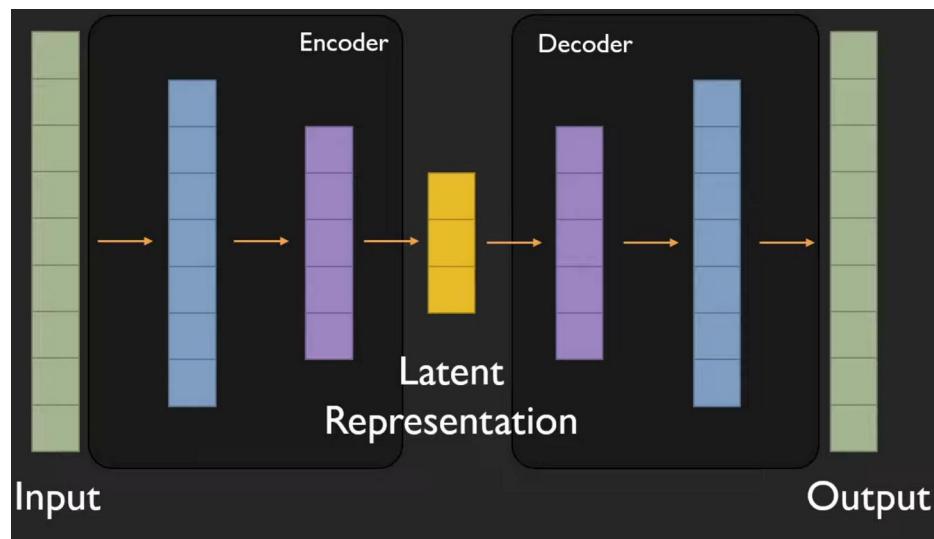
Champs récepteurs globaux dès le début

Les CNN utilisent des représentations hiérarchiques

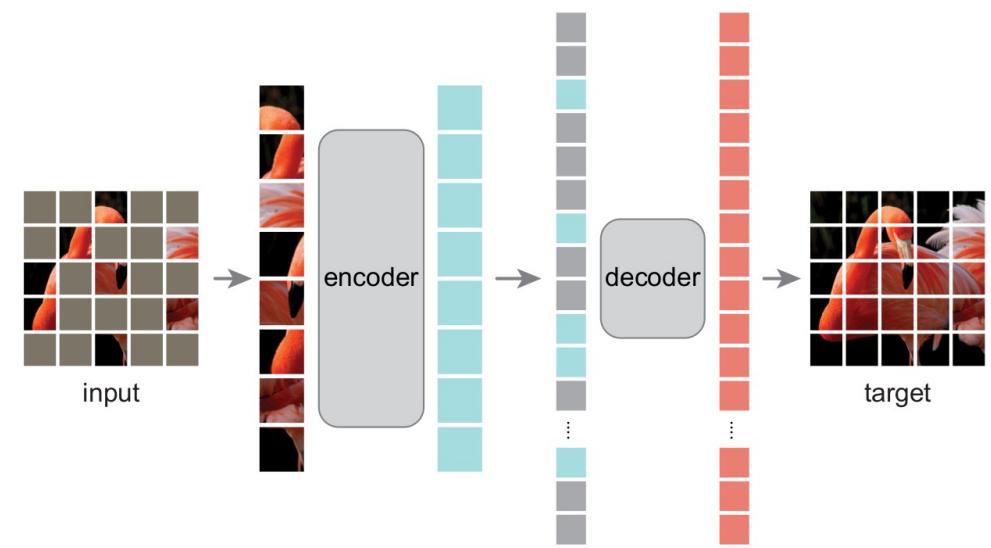
Masked Autoencoder (MAE)



Autoencoder

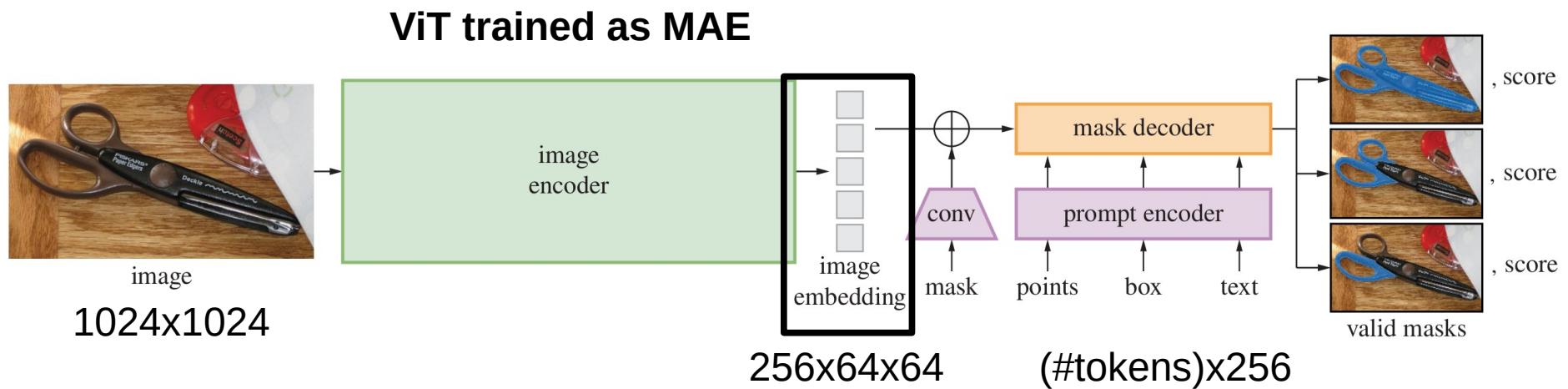


Masked Autoencoder

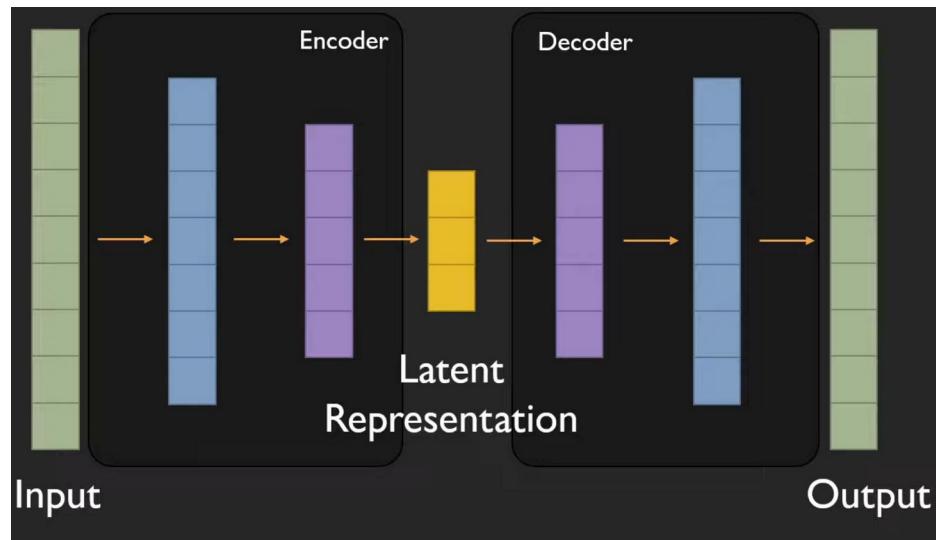


Fonctionne en auto-supervision.

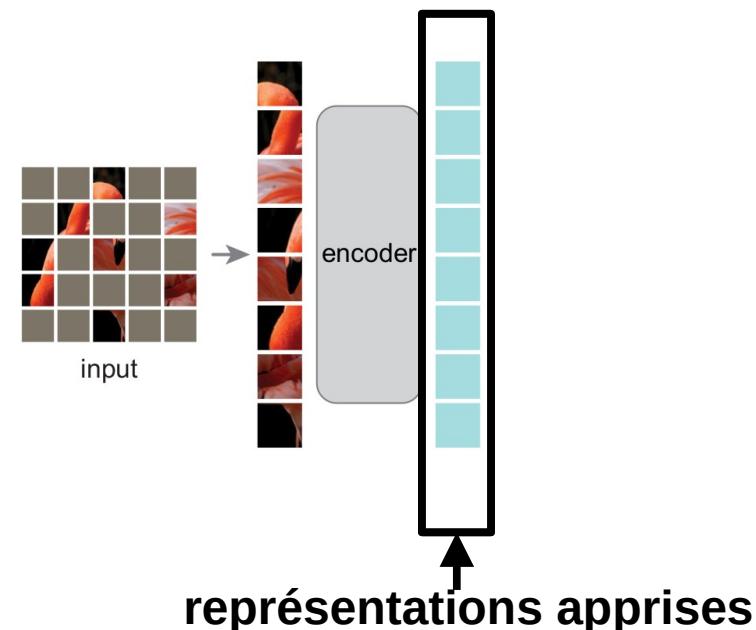
Masked Autoencoder (MAE)



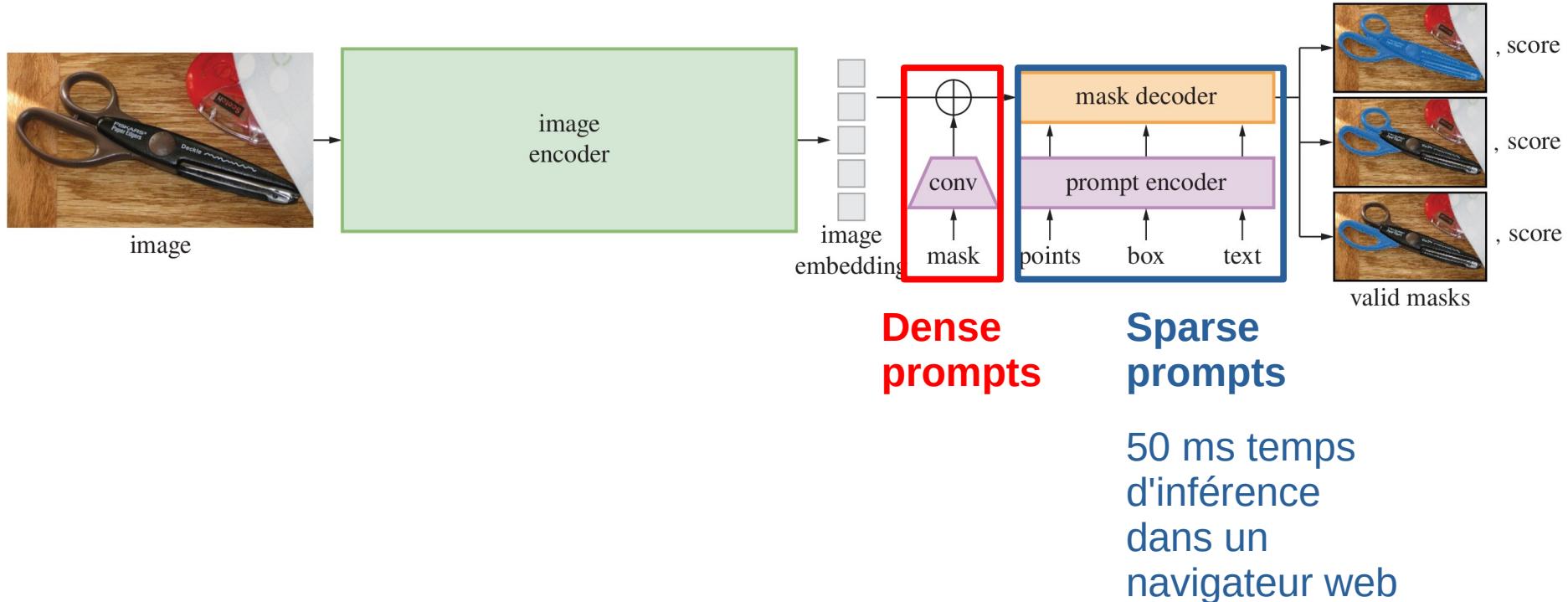
Autoencoder



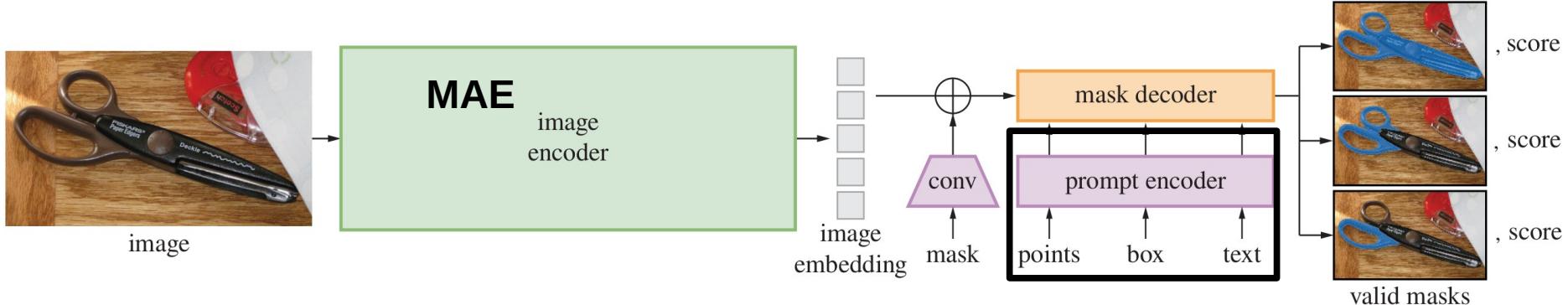
Masked Autoencoder



Segment Anything



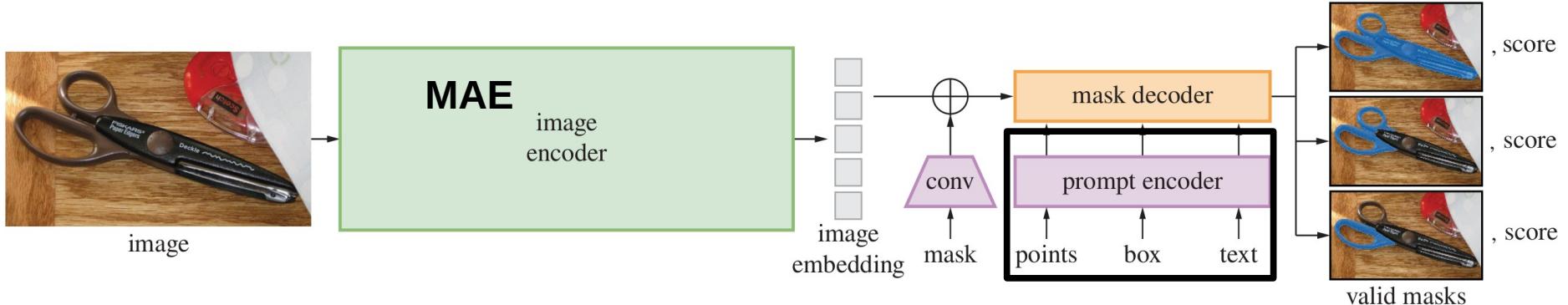
Segment Anything



CLIP = Contrastive Language – Image Pretraining

Positional encodings + CLIP model

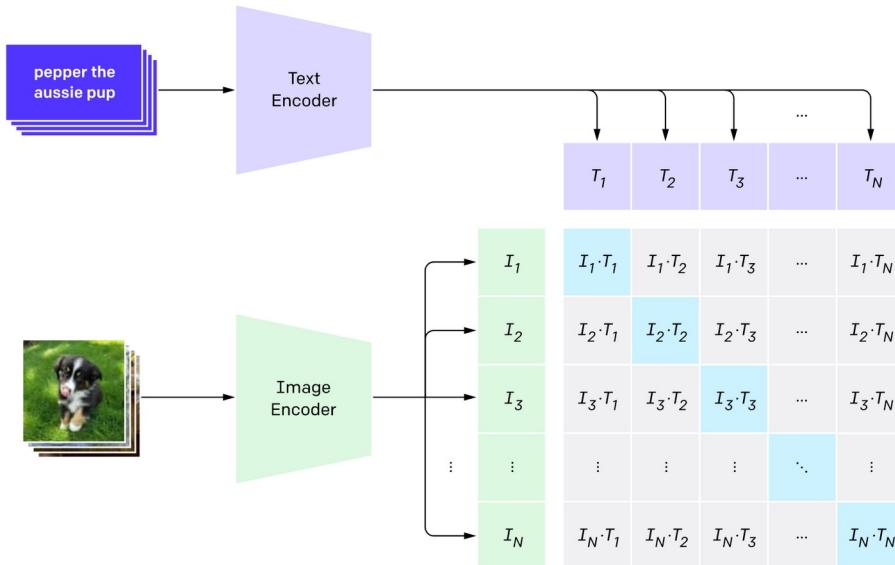
Segment Anything



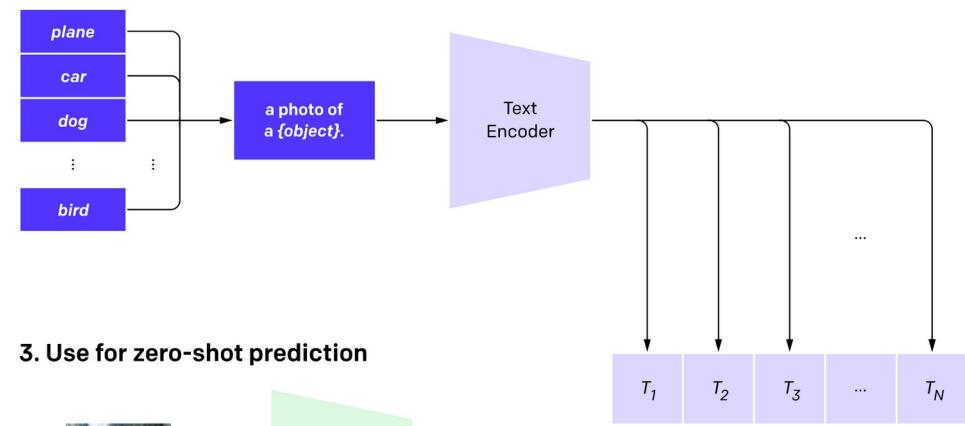
CLIP = Contrastive Language – Image Pretraining

Positional encodings + CLIP model

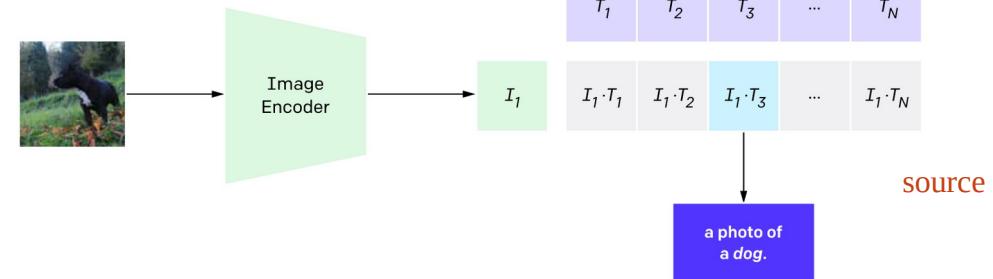
1. Contrastive pre-training



2. Create dataset classifier from label text



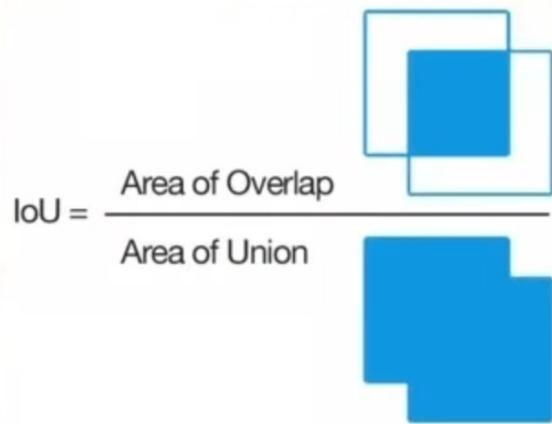
3. Use for zero-shot prediction



source

L'entraînement

Pour chaque objet – produire 3 masques : entier, partie, sous-partie



les masques ayant le plus haut IoU avec les données de référence est utilisé pour calculer la perte (loss)

DICE LOSS

$$\frac{2 \times (\text{Overlapping Pixels})}{\text{Total pixels in Pred \& Target}}$$

A Venn diagram consisting of two overlapping circles, one blue and one orange. The intersection of the two circles is shaded black. The formula above it indicates that the overlapping area is multiplied by 2.

$$\frac{2 \times \text{Intersection Area}}{\text{Union Area}}$$

+ **FOCAL LOSS**

$$FL = -(1 - P_t)^\gamma \log(P_t)$$
$$CE = -\log(P_t)$$

L'entraînement interactif



image d'entraînement

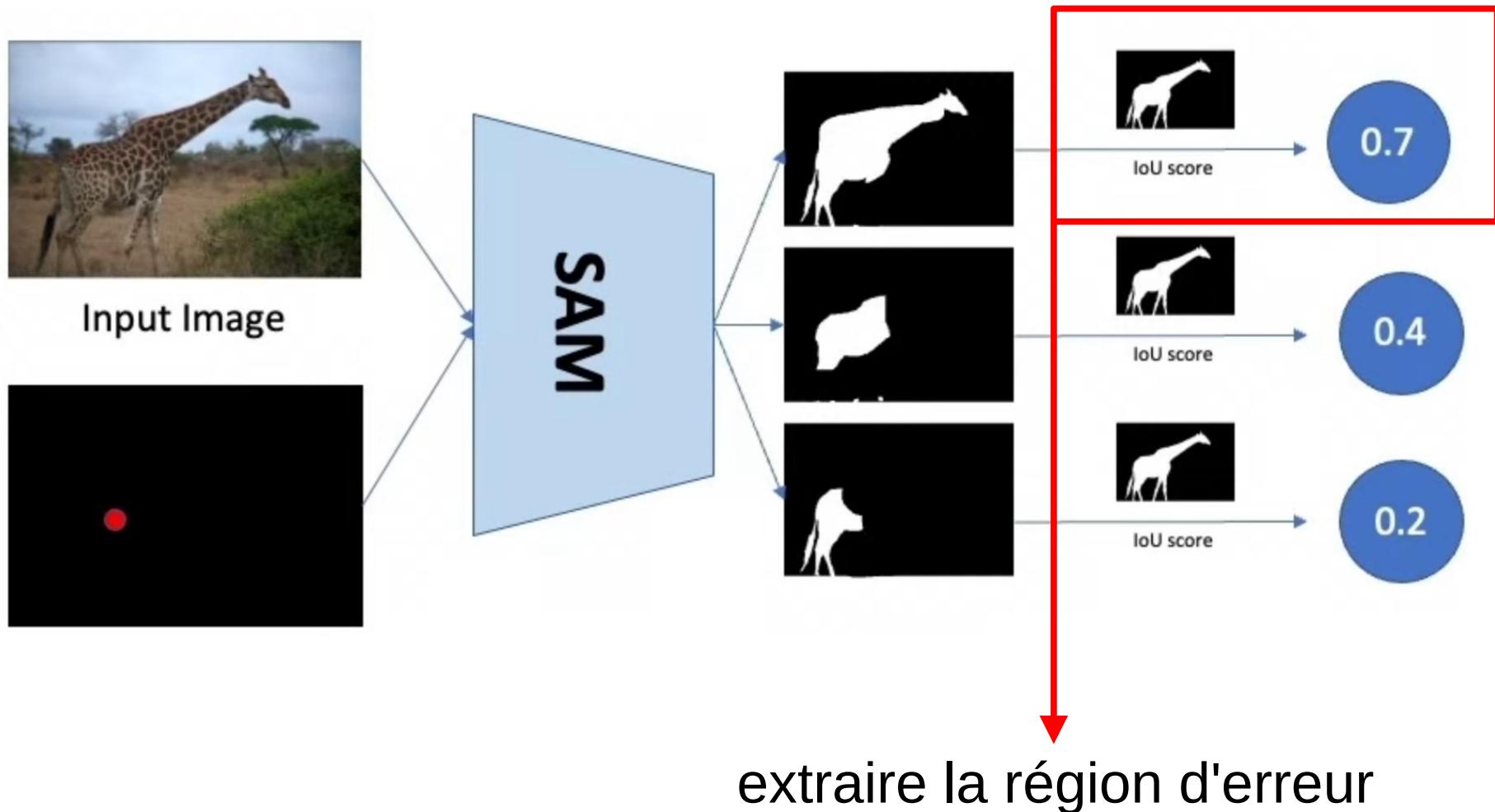


image cible

Prompt aléatoire



L'entraînement interactif



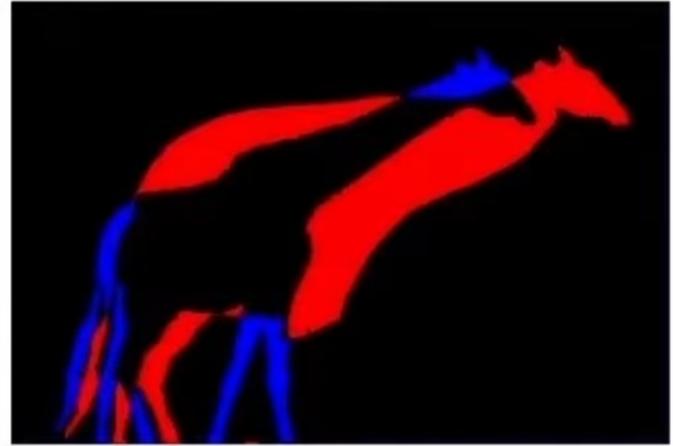
L'entraînement interactif



image cible



prédiction



région d'erreur

la région d'erreur = la différence entre la cible et le masque prédict

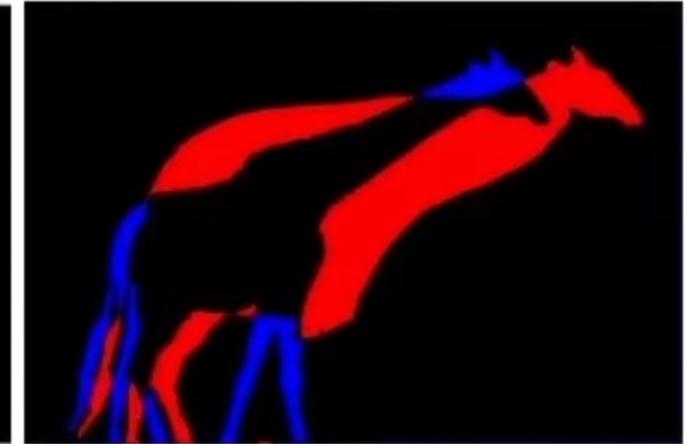
L'entraînement interactif



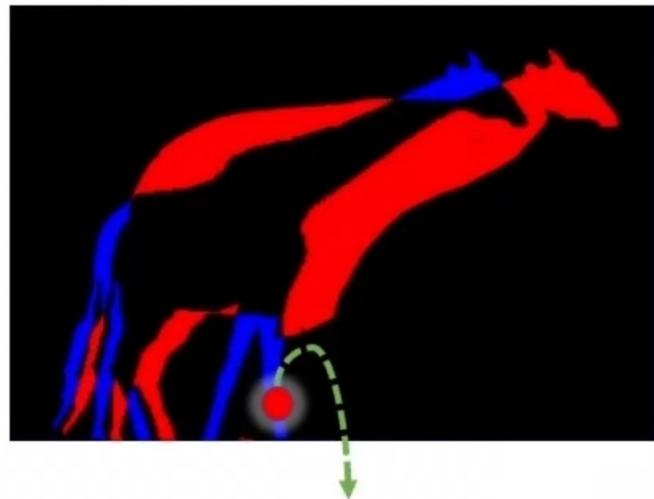
image cible



prédiction

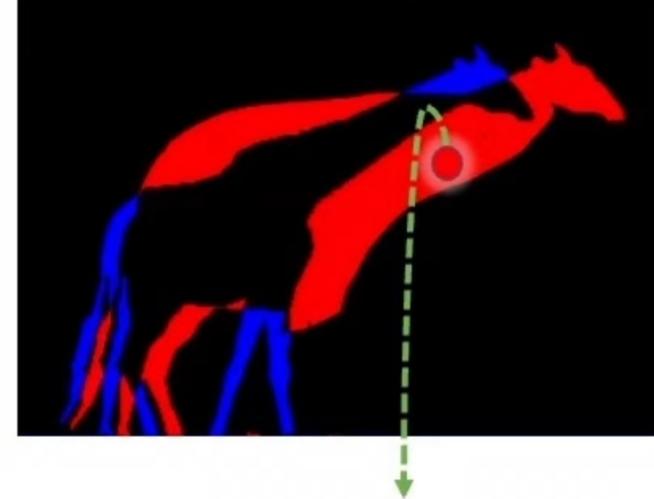


région d'erreur



Faux négatif

Point d'avant-plan
(foreground)



Faux positif

Point d'arrière-plan
(background)

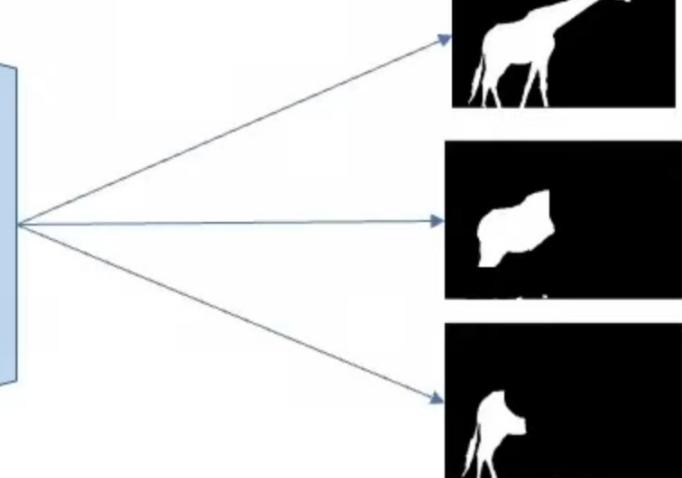
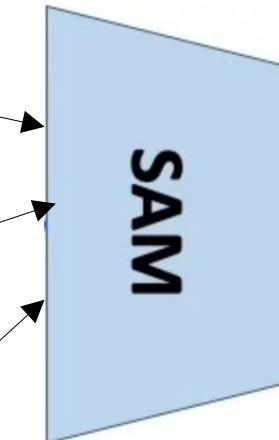
L'entraînement interactif



Sparse prompts



Dense prompt



Segment Anything

DEMO