

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

ALLOUACHE et ELMAM

Sorbonne Université

Introduction

Les Transformers sont appliqués sur plusieurs tâches de NLP, ses applications à la vision restent limitées. Dans cet article, on montre que pour une tâche de classification d'image, on utilise un transformer pur appliqué sur des séquences de patches d'image. Ainsi, le retour aux CNN n'est plus nécessaire. On explore la reconnaissance d'images à des échelles plus grandes que l'ensemble de données standard d'ImageNet. On divise l'image en un ensemble de patches. Les patches sont traités de la même manière que les mots dans une application de NLP.

Points-clés

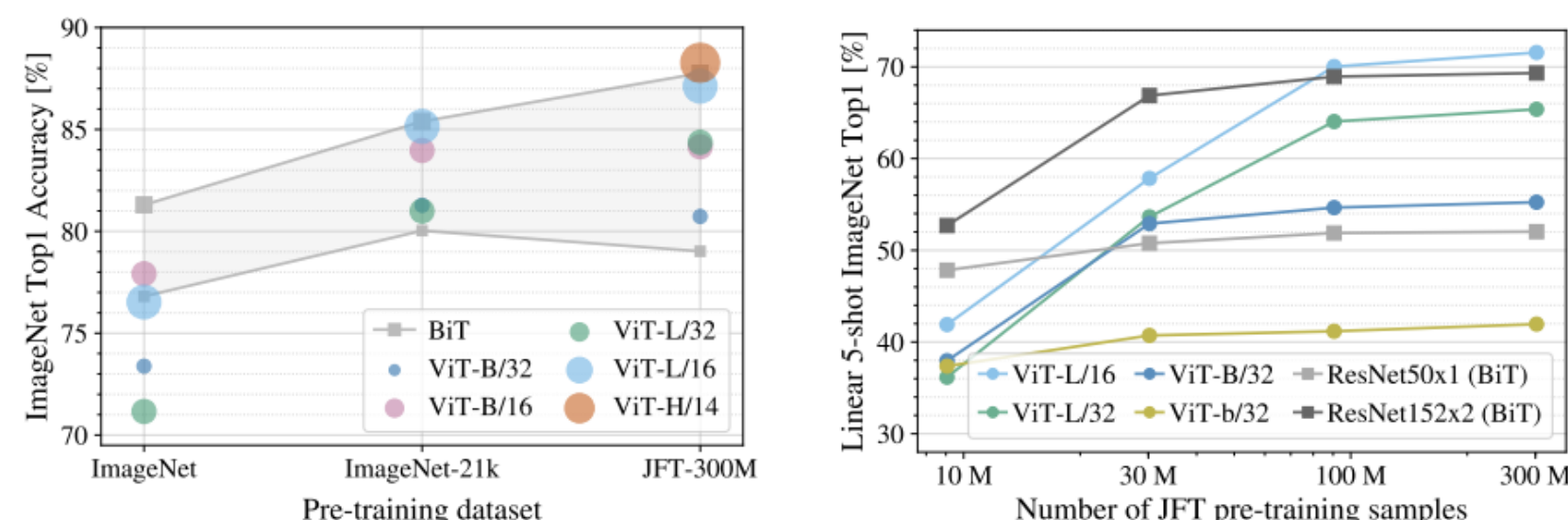
- Application de principe des Transformers pour la classification d'images.
- Représentation d'image sous forme d'une séquence de patches.
- Limitation de la taille des patches.
- La nécessité d'un large dataset pour l'entraînement.
- L'obtention de visualisation intermédiaire plus significative.
- Les architectures hybrides (CNN+Transformers).
- Moins de ressources pour l'entraînement que les modèles CNNs.

Experimentation

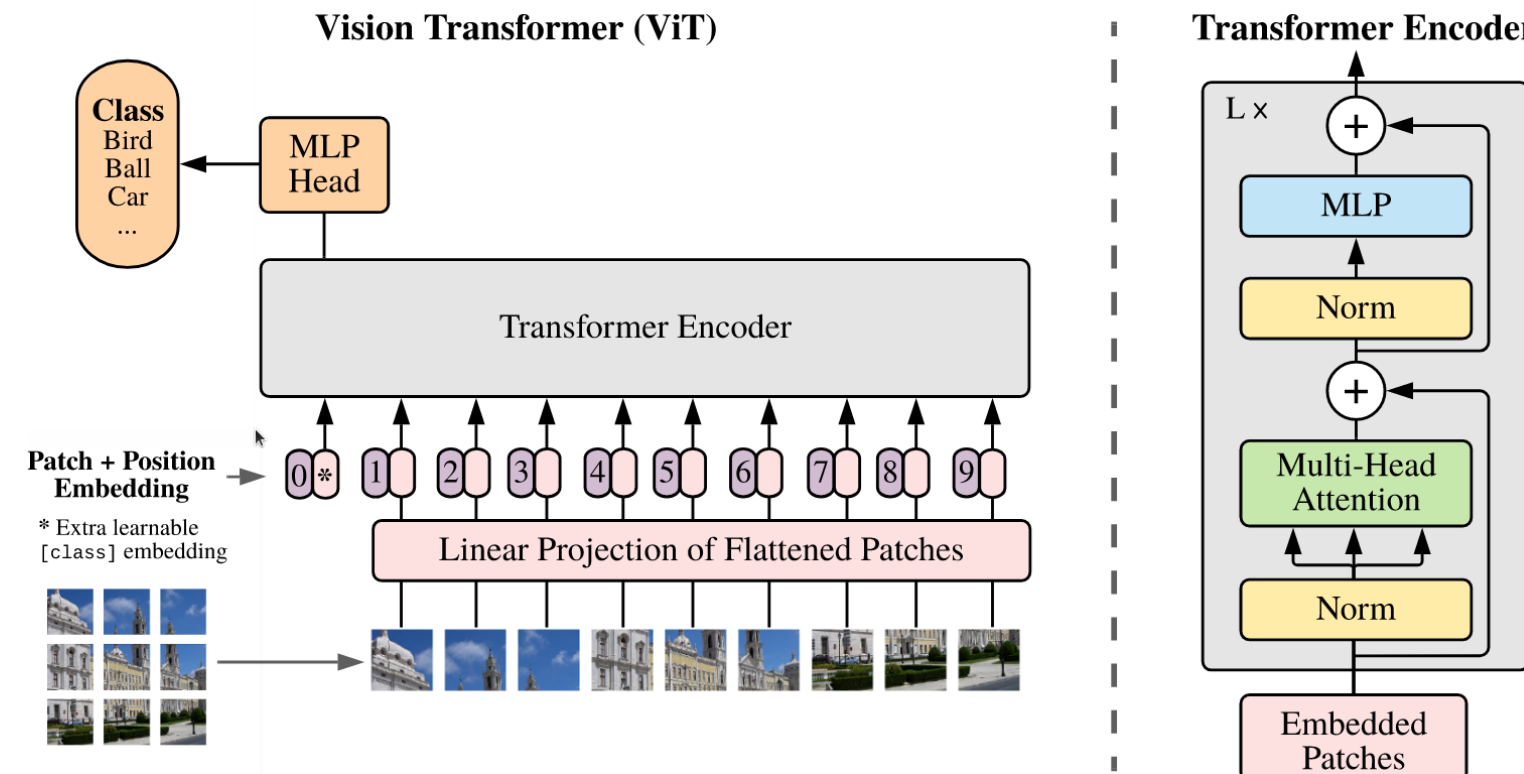
- **Datasets:** Entraînement sur *JFT*, *ImageNet*, ensuite application d'un transfert learning par les auteurs de l'article. Dans notre cas, on a entraîné le modèle sur CIFAR10.
- La configuration ViT utilisée par les auteurs est basée sur celles utilisées pour *BERT* (Devlin et al., 2019).

Model	Layers	Hidden size D	MLP size	Heads	Params
ViT-Base	12	768	3072	12	86M
ViT-Large	24	1024	4096	16	307M
ViT-Huge	32	1280	5120	16	632M

- ViT-L/16 signifie la variante "Large" avec une taille de patch d'entrée de 16×16 .
- La longueur de la séquence du transformer est inversement proportionnelle au carré de la taille du patch, ce qui fait que les modèles avec une taille de patch plus petite sont plus lourds en termes de calcul.
- Preprocessing: Basé sur celui utilisé dans *Kolesnikov et al (2020)*.
- Fonction de loss: label smoothing qui favorise la généralisation.
- Self-Supervision: pour prédire les couleurs des patches.
- Afin d'avoir de meilleures performances sur des petits datasets, on applique des régularisations: Adam avec un `weight_decay = 0.1`, `linear schedulers`.



Vision Transformer ViT



1 Embedding

- **Construction des Patches:** Image $\mathbf{x} \in \mathbb{R}^{H \times W \times C}$ tel que (H, W) est la résolution de l'image et C est le nombre de channels. On aplatit cette image vers une séquence de patches $\mathbf{x}_p \in \mathbb{R}^{N \times (P^2 \cdot C)}$ avec $N = HW/P^2$ et P la taille de patch.
- **Projection linéaire des patches (Patch Embedding):** Avec les paramètres $\mathbf{E} \in \mathbb{R}^{(P^2 \cdot C) \times D}$, on projette linéairement chaque patch \mathbf{x}_p vers une dimension D qui est la dimension de modèle.
- **Le token de classification:** Comme le token `[class]` de *BERT*, on ajoute à la séquence des Embeddings des patch $\mathbf{z}_0^0 = \mathbf{x}_{\text{class}}$ dont l'état à la sortie de l'encodeur \mathbf{Z}_L^0 sera la représentation de l'image d'entrée.
- **Embedding des Position:** $\mathbf{E}_{\text{pos}} \in \mathbb{R}^{N+1 \times D}$ Des Embeddings de position sont ajoutées aux Embeddings des patches pour conserver les informations de position. On utilise une matrice de poids $\mathbf{E}_{\text{pos}} \in \mathbb{R}^{(N+1) \times D}$ pour cet Embedding.

À la fin de l'étape de l'embedding, on obtient l'ensemble \mathbf{z}_0 qui sera l'entrée de l'encodeur:

$$\mathbf{z}_0 = [\mathbf{x}_{\text{class}}; \mathbf{x}_p^1 \mathbf{E}; \mathbf{x}_p^2 \mathbf{E}; \dots; \mathbf{x}_p^N \mathbf{E}] + \mathbf{E}_{\text{pos}}, \quad \mathbf{E} \in \mathbb{R}^{(P^2 \cdot C) \times D}, \mathbf{E}_{\text{pos}} \in \mathbb{R}^{(N+1) \times D}$$

2 Architecture Encodeur

Transformer encodeur consiste à alterner entre les couches Multi-Head-Self-Attention (MSA), des blocs (MLP). *LayerNorm* (LN) est appliquée avant chaque bloc, les connexions résiduelles sont appliquées après chaque bloc. *MLP* contient 2 couches linéaires et une activation non-linéaire *GELU*.

$$\begin{aligned} \mathbf{z}'_\ell &= \text{MSA}(\text{LN}(\mathbf{z}_{\ell-1})) + \mathbf{z}_{\ell-1}, \quad \ell = 1 \dots L \\ \mathbf{z}_\ell &= \text{MLP}(\text{LN}(\mathbf{z}')) + \mathbf{z}'_\ell, \quad \ell = 1 \dots L \\ \mathbf{y} &= \text{LN}(\mathbf{z}_L^0) \end{aligned}$$

2.1 MSA: Multi-Head Self-Attention

$$\begin{aligned} [\mathbf{q}, \mathbf{k}, \mathbf{v}] &= \mathbf{z} \mathbf{U}_{qkv} & \mathbf{U}_{qkv} &\in \mathbb{R}^{D \times 3D_h} \\ A &= \text{softmax}\left(\frac{\mathbf{q} \mathbf{k}^T}{\sqrt{D_h}}\right) & A &\in \mathbb{R}^{N \times N} \\ \text{SA}(\mathbf{z}) &= A \mathbf{v}. \end{aligned}$$

Pour chaque élément d'une séquence $\mathbf{z} \in \mathbb{R}^{N \times D}$ est projeté vers une dimension $3 * D_h$ avec h le nombre de *heads* afin d'obtenir une représentation $[k, q, v]$, ensuite on calcule une somme pondérée sur toutes les valeurs \mathbf{v} dans la séquence. Les poids d'attention A_{ij} sont basés sur la similarité entre les représentations \mathbf{q}^i et \mathbf{k}^j .

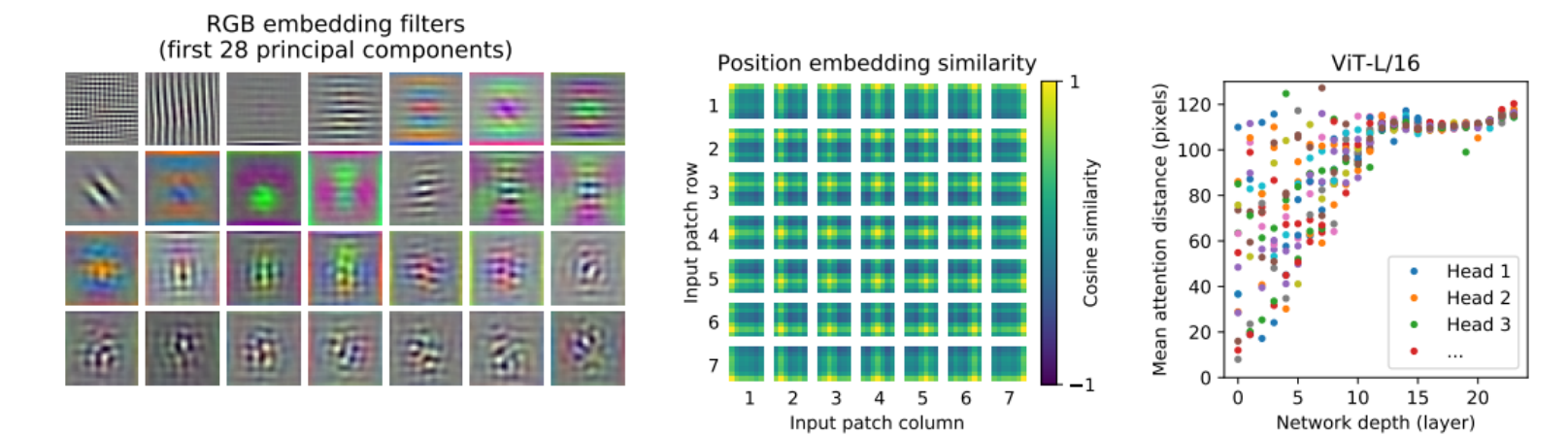
$$\text{MSA}(\mathbf{z}) = [\text{SA}_1(\mathbf{z}); \text{SA}_2(\mathbf{z}); \dots; \text{SA}_k(\mathbf{z})] \mathbf{U}_{msa} \quad \mathbf{U}_{msa} \in \mathbb{R}^{k \cdot D_h \times D}$$

Multihead self-attention (MSA) c'est une extension de *SA* dans laquelle on exécute k self-attention appelée "*heads*", en parallèle, ensuite on projette leur concaténation vers D pour revenir à la dimension de l'encodeur. Pour le modèle *ViT* on choisit $D_h = D/k$.

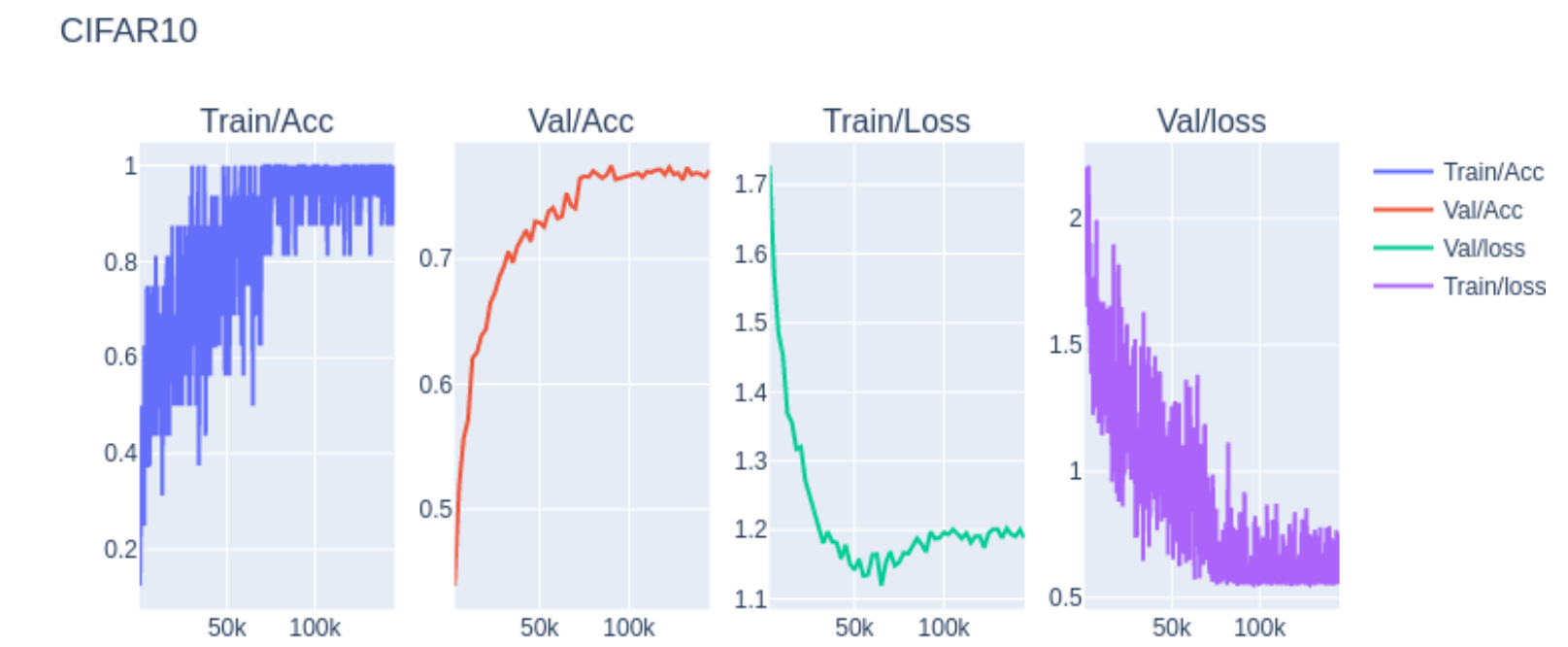
Résultats

Résultats obtenus par les auteurs de l'article:

	Ours-JFT (ViT-H/14)	Ours-JFT (ViT-L/16)	Ours-I21K (ViT-L/16)	BiT-L (ResNet152x4)	Noisy Student (EfficientNet-L2)
ImageNet	88.55 ± 0.04	87.76 ± 0.03	85.30 ± 0.02	87.54 ± 0.02	88.4/88.5*
ImageNet ReaL	90.72 ± 0.05	90.54 ± 0.03	88.62 ± 0.05	90.54	90.55
CIFAR-10	99.50 ± 0.06	99.42 ± 0.03	99.15 ± 0.03	99.37 ± 0.06	—
CIFAR-100	94.55 ± 0.04	93.90 ± 0.05	93.25 ± 0.05	93.51 ± 0.08	—
Oxford-IIT Pets	97.56 ± 0.03	97.32 ± 0.11	94.67 ± 0.15	96.62 ± 0.23	—
Oxford Flowers-102	99.68 ± 0.02	99.74 ± 0.00	99.61 ± 0.02	99.63 ± 0.03	—
VTAB (19 tasks)	77.63 ± 0.23	76.28 ± 0.46	72.72 ± 0.21	76.29 ± 1.70	—
TPUv3-core-days	2.5k	0.68k	0.23k	9.9k	12.3k



Nos résultats:



Conclusion

Nous avons exploré l'application directe de Transformers à la reconnaissance d'images et nous avons constaté que les résultats du modèle correspondent ou dépassent l'état de l'art sur de nombreux ensembles de données de classification d'images.

Avec l'augmentation constante de la taille des ensembles de données, et le développement continu de méthodes non supervisées et semi-supervisées, le développement de nouvelles architectures de vision qui s'entraînent plus efficacement sur ces ensembles de données devient de plus en plus important. Nous pensons que *ViT* est une étape préliminaire vers des architectures génériques et évolutives qui peuvent résoudre de nombreuses tâches de vision, ou même des tâches provenant de nombreux domaines.

Informations supplémentaires

Code source: github.com/YacineAlI/ViT-pytorch-lightning
Article: <https://arxiv.org/abs/2010.11929>