

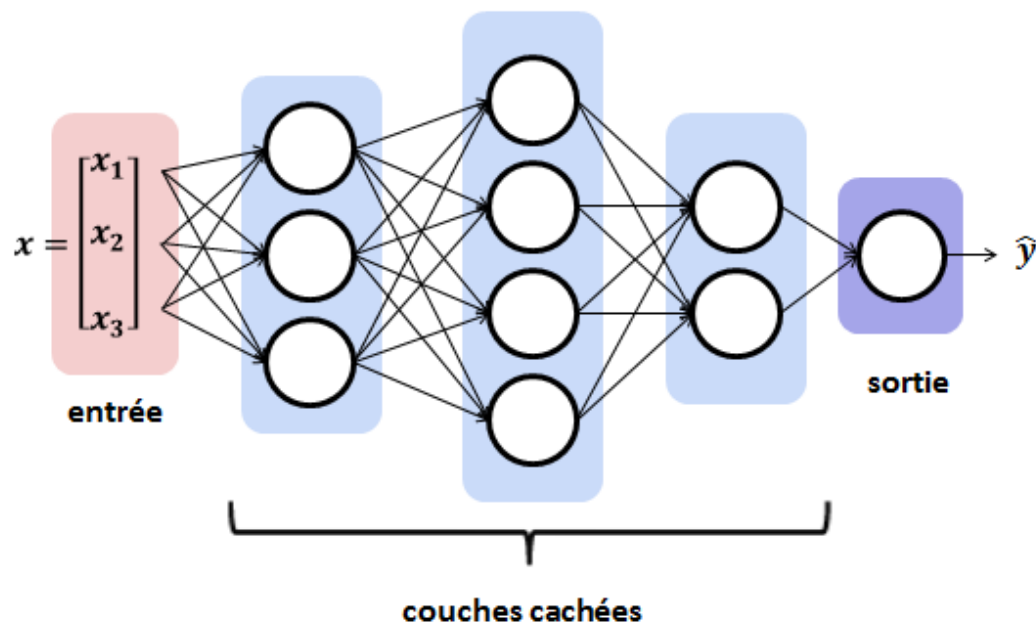
Traitement et Analyse d'Images

Apprentissage profond

Réseaux de neurones convolutionnels

Analyse d'images par réseaux de neurones

► Perceptron multicouches (MLP)



- Vectorisation d'une image en entrée

- Image 256 x 256

$$x^{(i)} \in \mathbb{R}^{[65536 \times 1]}$$

- Réseau à 1 couches cachée avec $n^{[1]} = n_x$

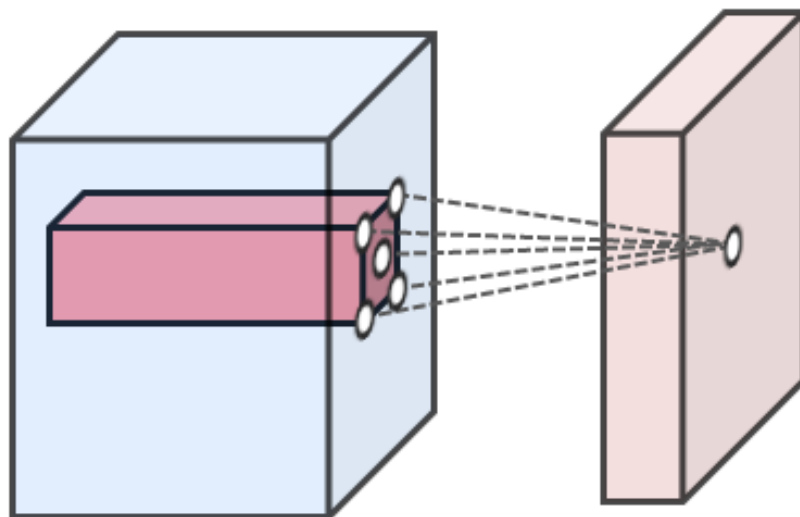
$$\# \text{ param} > 4 M$$

Beaucoup de paramètres pour des réseaux peu profonds

Analyse d'images par réseaux de neurones

► Comment adapter les réseaux de neurones aux images ?

- Exploitation de la notion de *convolution*



- Paramètres partagés
- Connectivité locale

Moins de paramètres pour des architectures plus profondes

Briques fondamentales

Convolutions, cartes de caractéristiques,
champ réceptif (*receptive field*)

Rappel sur la convolution

3 ₀	3 ₁	2 ₂	1	0
0 ₂	0 ₂	1 ₀	3	1
3 ₀	1 ₁	2 ₂	2	3
2	0	0	2	2
2	0	0	0	1

12	12	17
10	17	19
9	6	14

$$\begin{pmatrix} 0 & 1 & 2 \\ 0 & 0 & 1 \\ 0 & 1 & 2 \end{pmatrix} * \begin{pmatrix} 3 & 3 & 2 & 1 & 0 \\ 0 & 0 & 1 & 3 & 1 \\ 3 & 1 & 2 & 2 & 3 \\ 2 & 0 & 0 & 2 & 2 \\ 2 & 0 & 0 & 0 & 1 \end{pmatrix} = \begin{pmatrix} 12 & 12 & 17 \\ 10 & 17 & 19 \\ 9 & 6 & 14 \end{pmatrix}$$

filtre

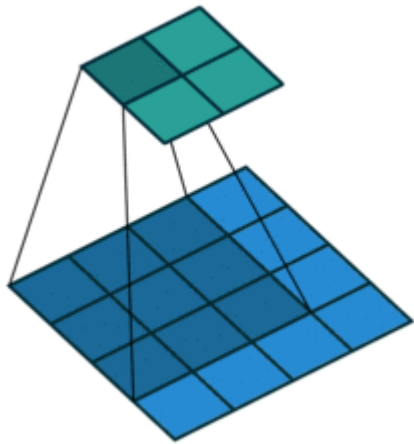
image

sortie

$$G = h * F$$

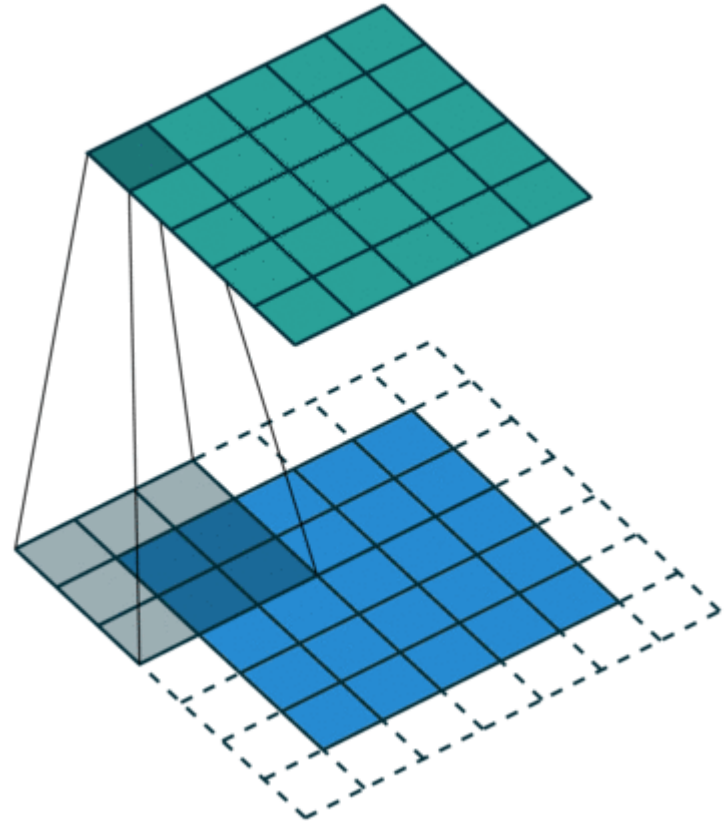
$$G[i, j] = \sum_{u=-k}^k \sum_{v=-k}^k h[u, v] F[i - u, j - v]$$

Remplissage / Padding



Sans padding

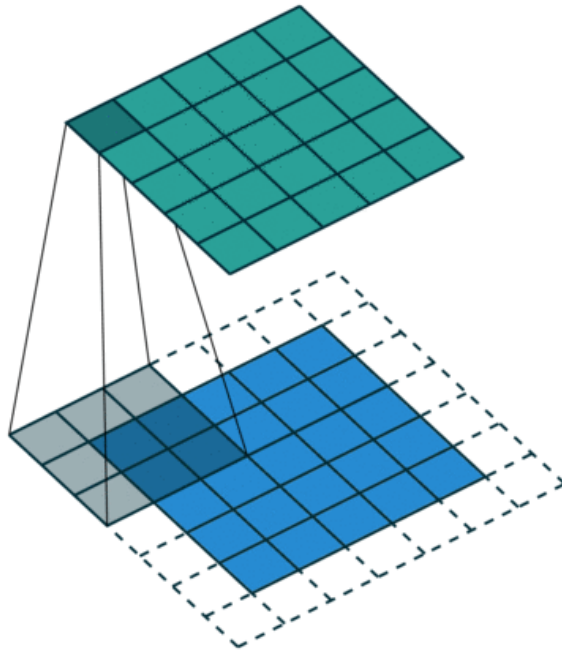
Dimensions spatiales réduites



Avec padding

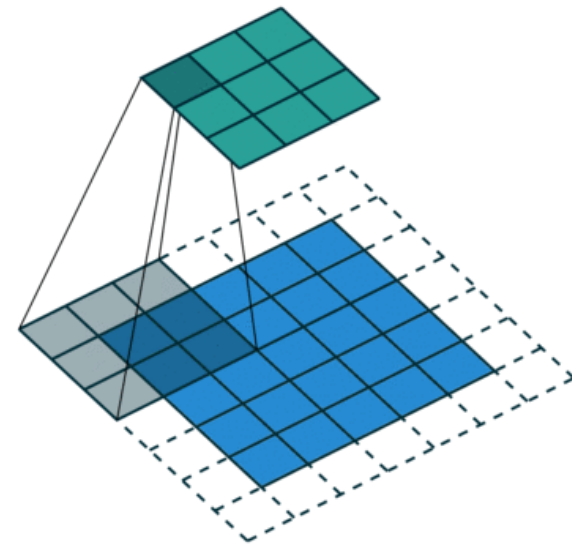
Dimensions spatiales préservées

Pas / Stride



Stride = 1

Dimensions spatiales préservées

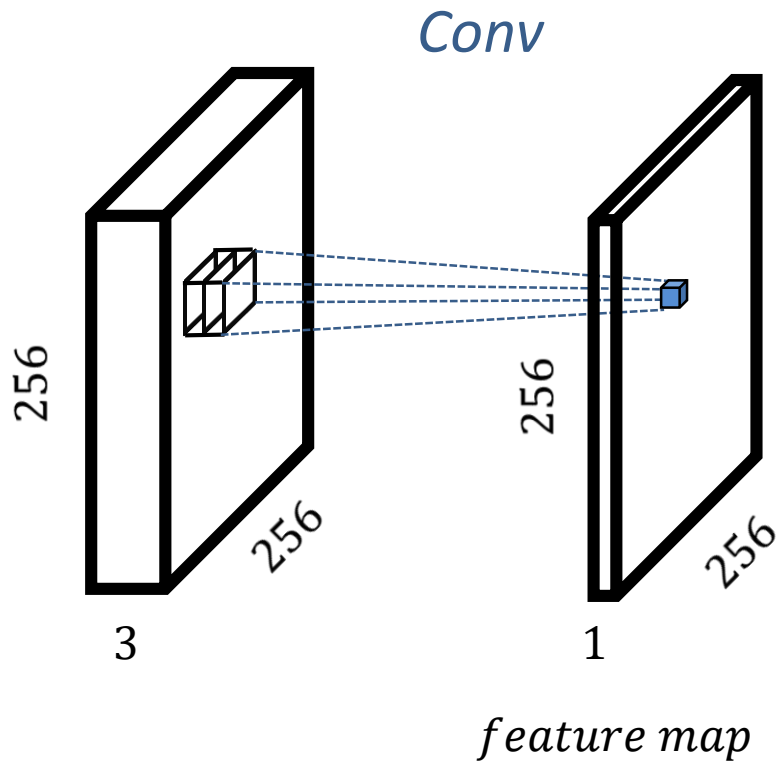


Stride = 2

Dimensions spatiales réduites

Réseaux de neurones convolutionnels

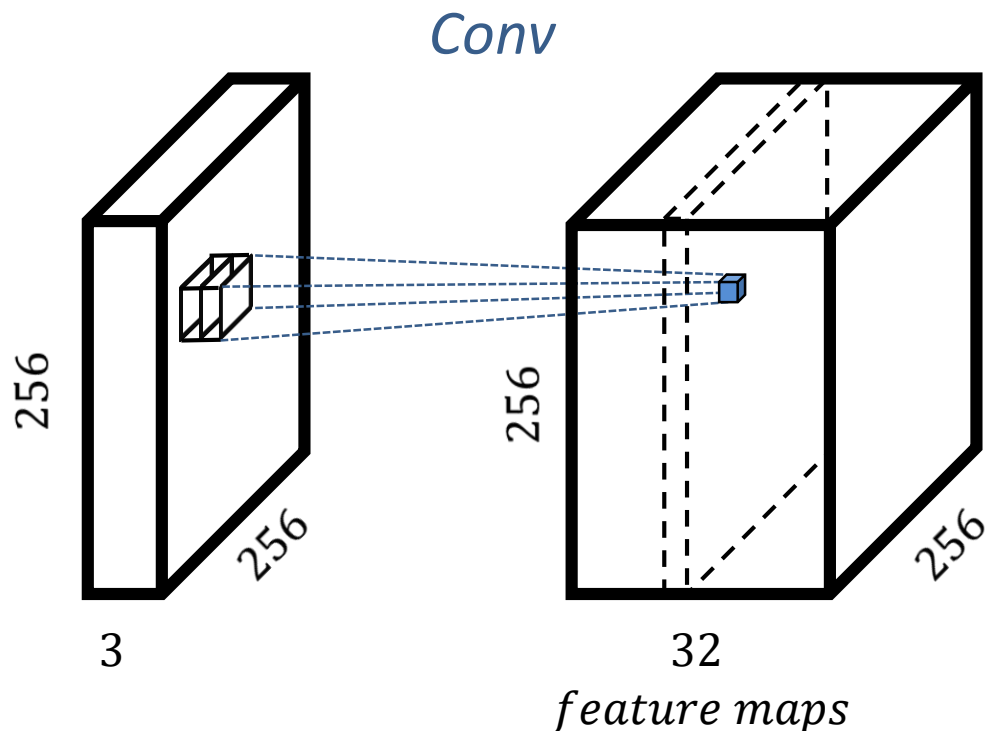
- Paramètres d'apprentissage - *poids de chaque filtre*



- Ex. filtre de taille 3×3
$$\# \text{ param} = 3 \times 3 \times 3 + 1$$
$$= 28$$
- *Conv*
filtrage de l'image d'entrée +
fonction d'activation
feature map / canal

Réseaux de neurones convolutionnels

- ▶ Plusieurs *feature maps* (canaux) par couche



- Ex. filtre de taille 3×3

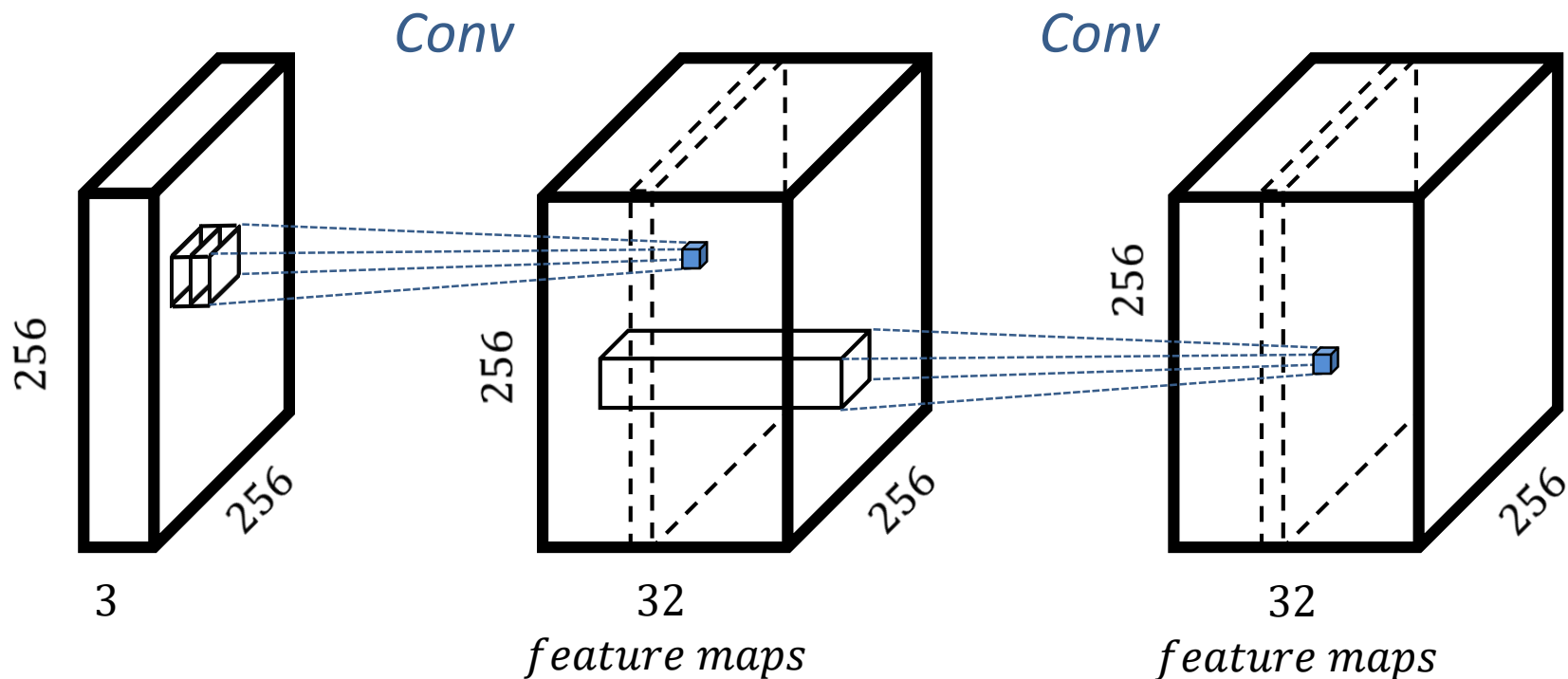
param

$$= 32 \times (3 \times 3 \times 3 + 1)$$

$$= 896$$

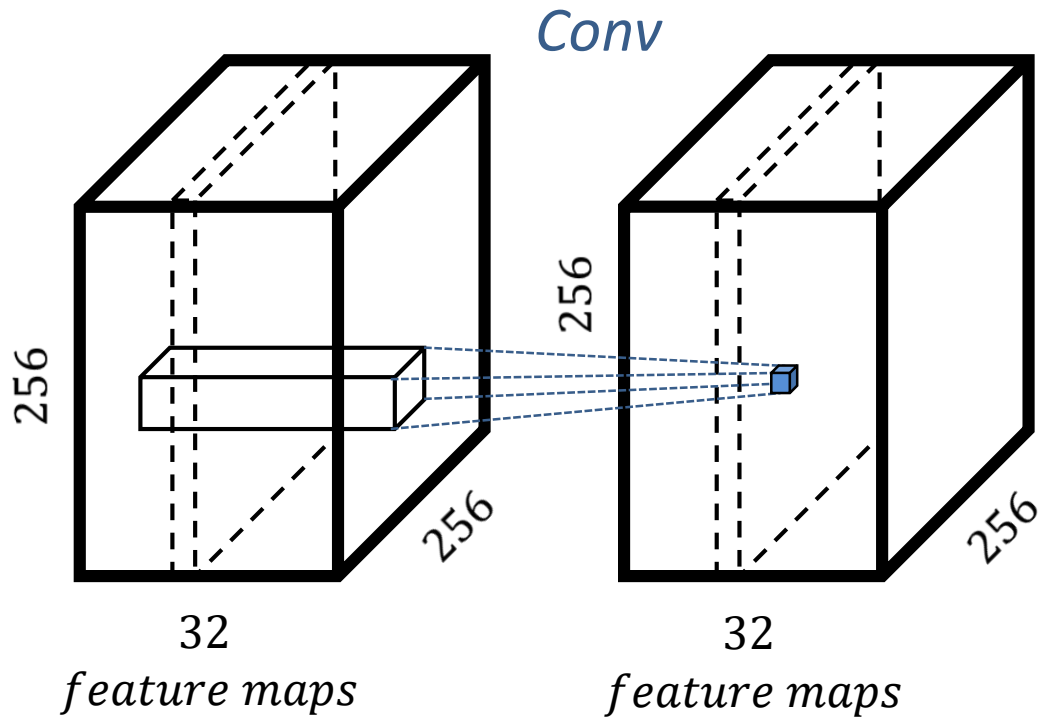
Réseaux de neurones convolutionnels

► Approche multicouche



Réseaux de neurones convolutionnels

► Approche multicouche



- Ex. filtre de taille 3×3

param

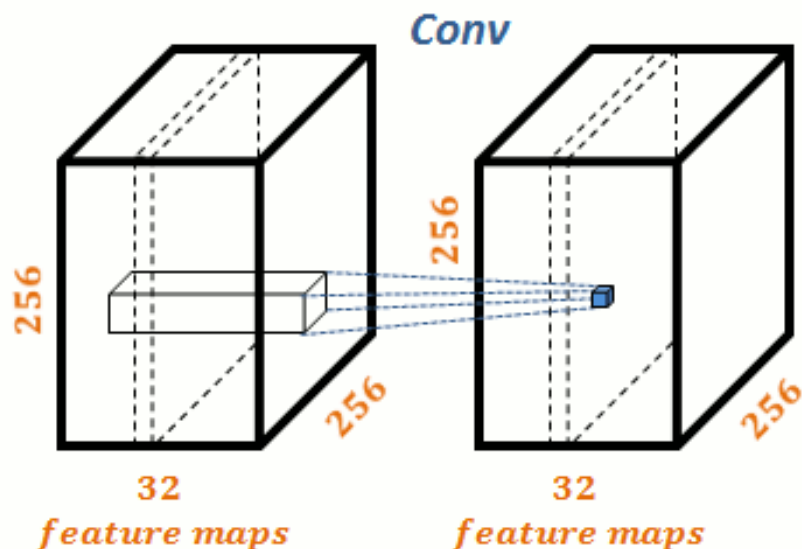
$$= 32 \times (3 \times 3 \times 32 + 1)$$

$$= 9248$$

Réseaux de neurones convolutionnels

► Approche multicouches

- Taille du filtre = nombre de canaux d'entrée
- Nombre de canaux de sortie = nombre de filtres
- Noyau d'une convolution 2D est un filtre 3D
- Noyau d'une convolution 3D est un filtre 4D



Pooling

- Réduction de la résolution spatiale des *canaux (feature maps)*
- Diminution de l'empreinte mémoire / coût calculatoire
- Invariance pour de petites translations, rotations et changements d'échelle
- Appliqué individuellement à chaque canal d'entrée

Max pooling

135	212	189	56
164	201	204	145
30	126	189	156
36	45	38	12

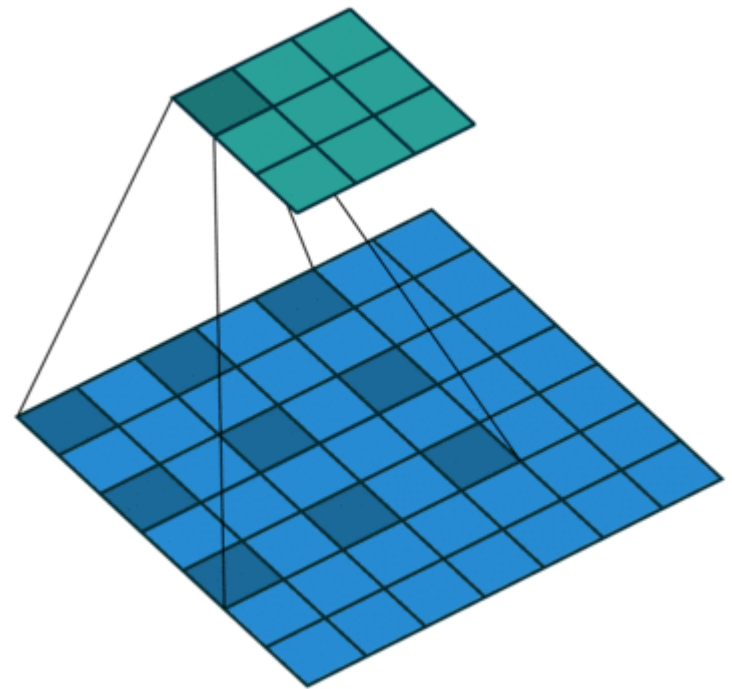


Taille pool
 2×2
(Stride = 2)

212	204
126	189

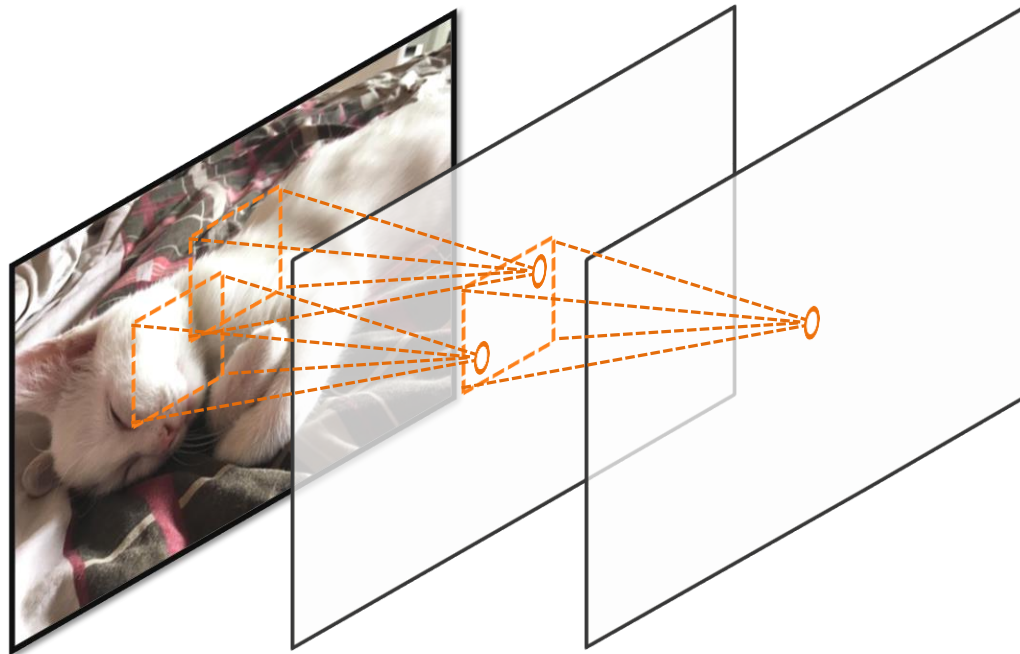
Convolution dilatée (à trous)

- Augmente la taille du filtre en rajoutant des espaces entre les éléments
- Contrôle par le paramètre de dilatation d
- Réduit la dimension spatiale des *feature maps* en sortie



Proche de l'opération de pooling mais ici des paramètres doivent être appris

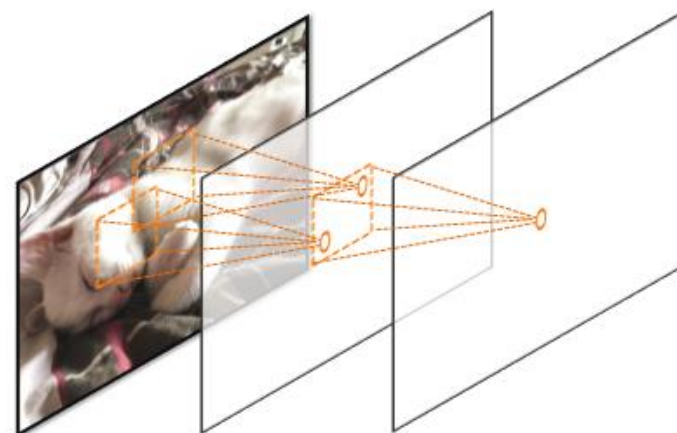
Champ réceptif / Receptive field



Région de l'image d'entrée qui affecte la valeur d'un neurone d'une couche

Champ réceptif / Receptive field

- Un champ réceptif large est nécessaire pour capturer une information de contexte spatial
- Au prix du nombre de paramètres
- Le champ réceptif augmente avec la profondeur d'un réseau

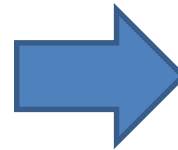


Comment avoir un champ réceptif large sans un nombre trop élevé de paramètres ?

Applications

Classification d'images

Classification d'images



Chien - 10%
Chat - 85%
Cheval - 5%
Chaise - 0%

Prédire une seule classe (ou une distribution de probabilité sur un ensemble de classes) pour une image donnée

Difficultés



Ce que voit un humain

37	49	43	43	63	45	51	56	65	59	28
47	64	68	37	48	56	37	47	61	47	65
56	67	64	39	80	66	31	48	49	33	45
38	49	32	75	48	49	71	35	47	27	62
61	62	33	64	60	49	35	40	70	49	47
52	32	31	56	34	32	34	27	43	36	60
34	77	26	36	46	27	62	76	70	65	27
69	36	49	37	34	41	75	61	69	46	76
31	40	62	30	67	43	54	77	72	72	70
42	69	65	76	73	61	64	34	53	66	67
39	52	55	64	45	78	34	76	60	57	70
55	68	31	56	63	77	78	69	78	35	47
31	47	60	43	42	51	55	57	50	78	59
48	41	32	35	55	39	63	50	29	40	57
31	28	33	60	71	68	28	40	73	76	55
32	63	31	80	58	67	70	67	60	38	41
69	49	33	35	44	66	67	38	45	46	39
42	50	35	40	42	66	32	29	80	30	50
59	59	36	47	50	31	54	68	38	61	38
79	29	43	30	49	63	43	62	61	35	70

Ce que voit un ordinateur

D'autres difficultés

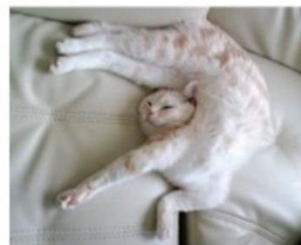
Point de vue d'observation



Variation d'échelle



Déformation



Occlusion



Conditions d'illumination



Texture



Variation intra-class



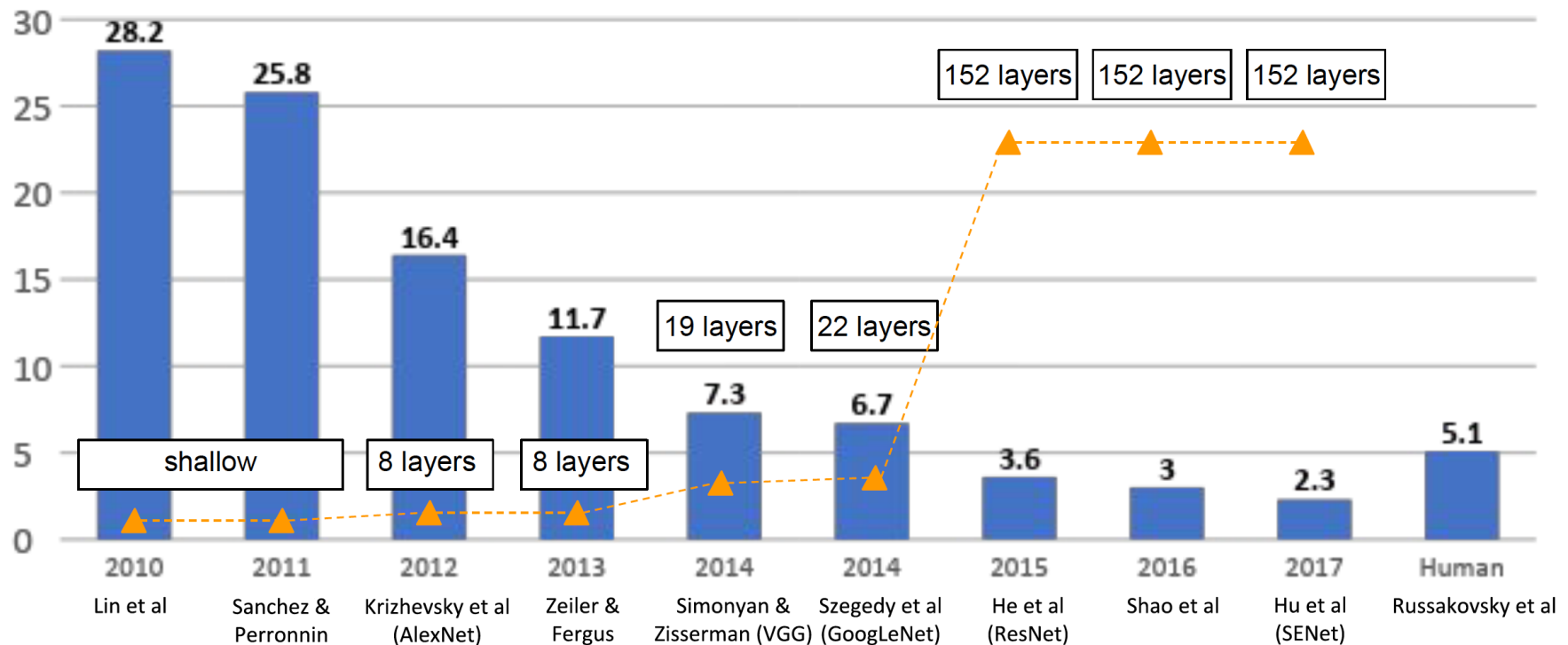
Simple pour un humain, qu'en est-il pour un ordinateur ?



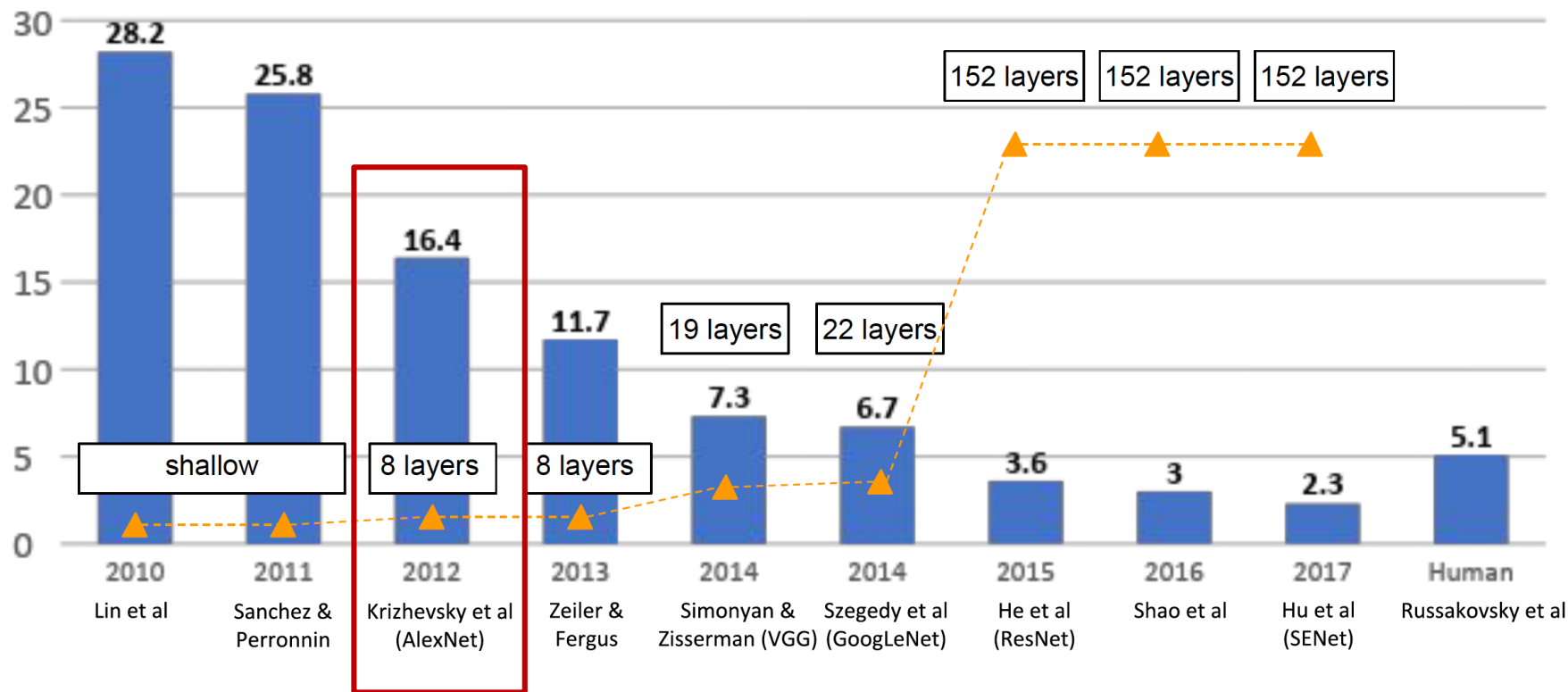
- Challenge en classification d'images (2010 → 2017)
- 1 000 classes d'objet à reconnaître
- 1 431 167 images

ImageNet

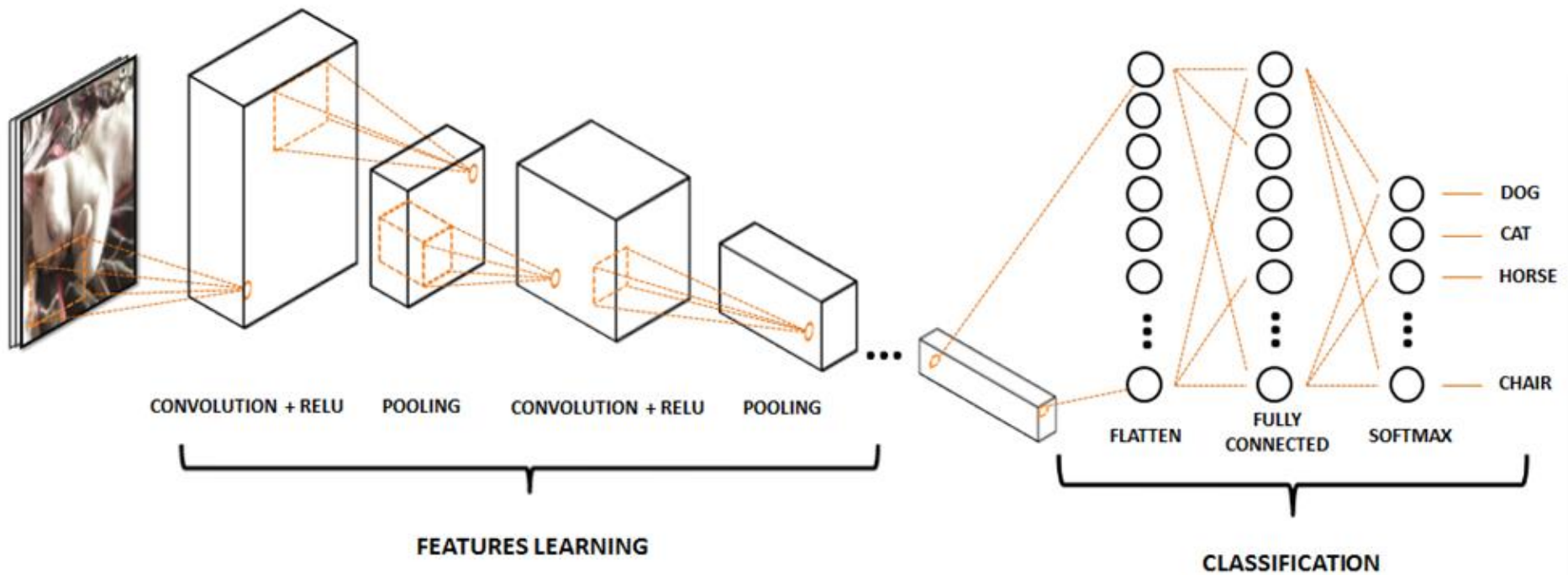
► Classement annuel



AlexNet [Krizhevsky, NIPS, 2012]



AlexNet [Krizhevsky, NIPS, 2012]



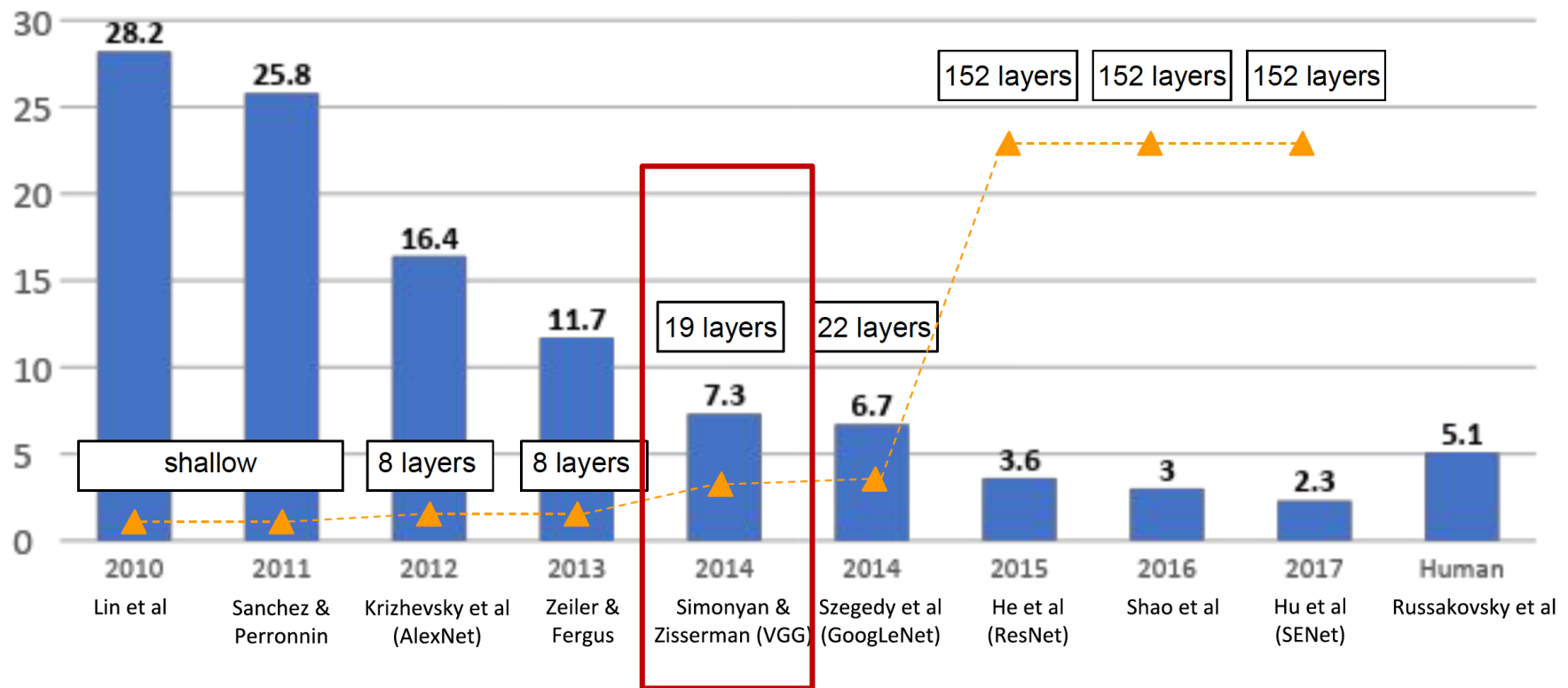
AlexNet [Krizhevsky, NIPS, 2012]

- Premier modèle à bien performer sur la base de données ImageNet
- Exploite des techniques encore utilisées (ReLU, augmentation de données, dropout)
- Utilise des GPU pour l'entraînement

Responsable de la révolution de l'apprentissage profond en vision par ordinateur



VGG [Simonyan and Zisserman, arxiv, 2014]



VGG [Simonyan and Zisserman, arxiv, 2014]

- Architecture plus simple

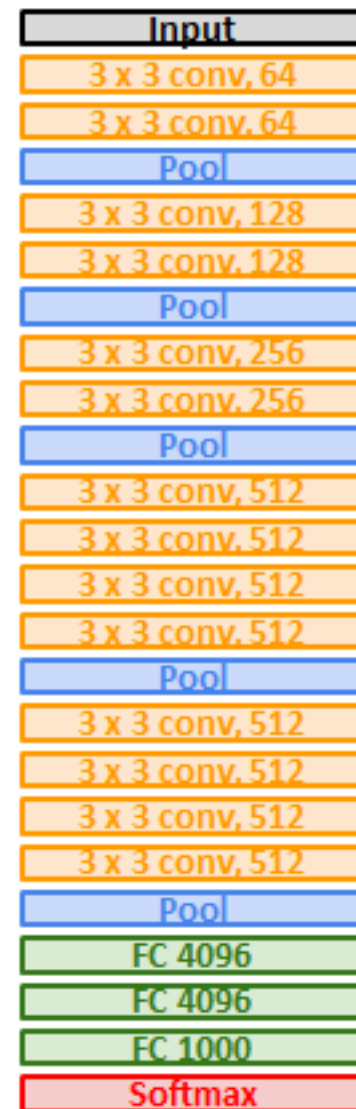
Convolutions 3×3 , ReLU et max pooling 2×2

- Réseau plus profond

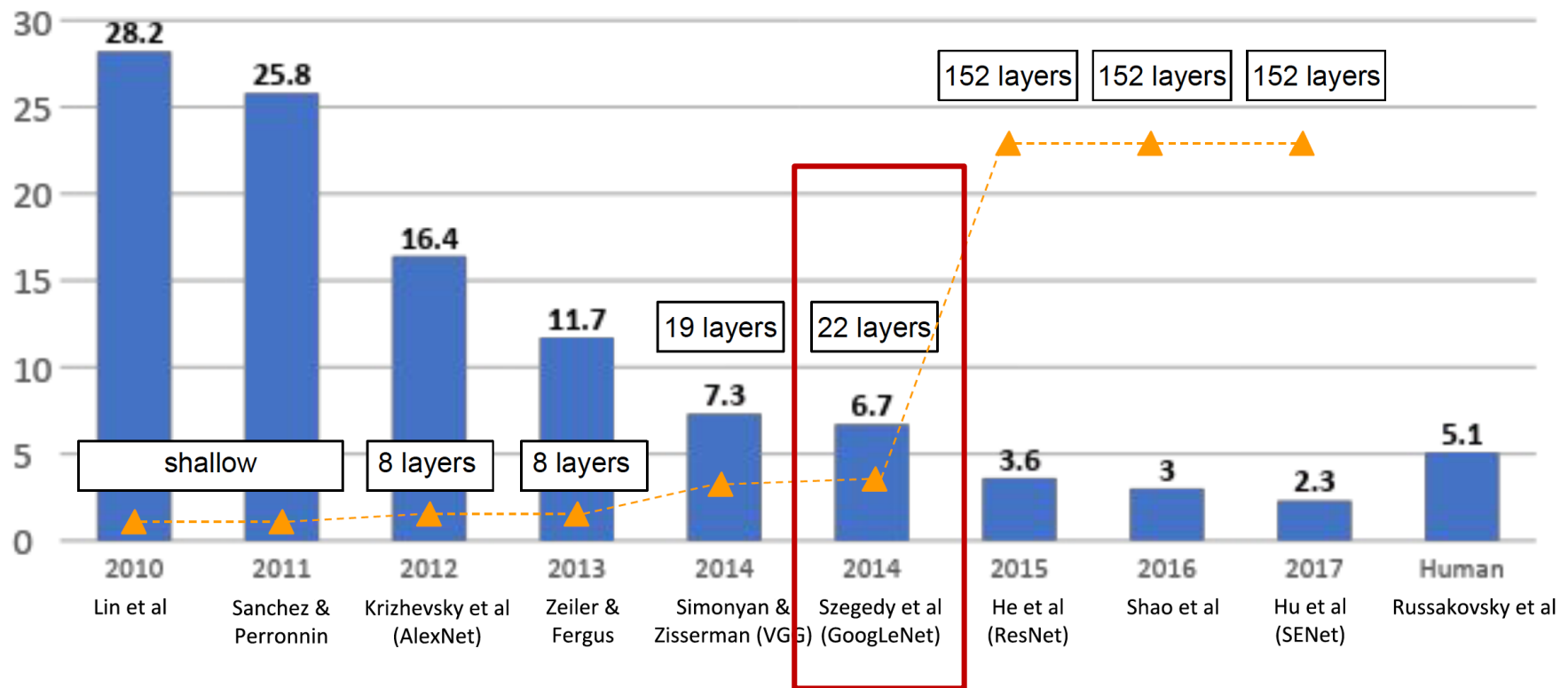
19 couches (8 couches pour AlexNet)

- Idée clé:

Mettre en cascade 2 convolutions 3×3 produit le même champ réceptif qu'un convolution 5×5 mais avec moins de paramètres

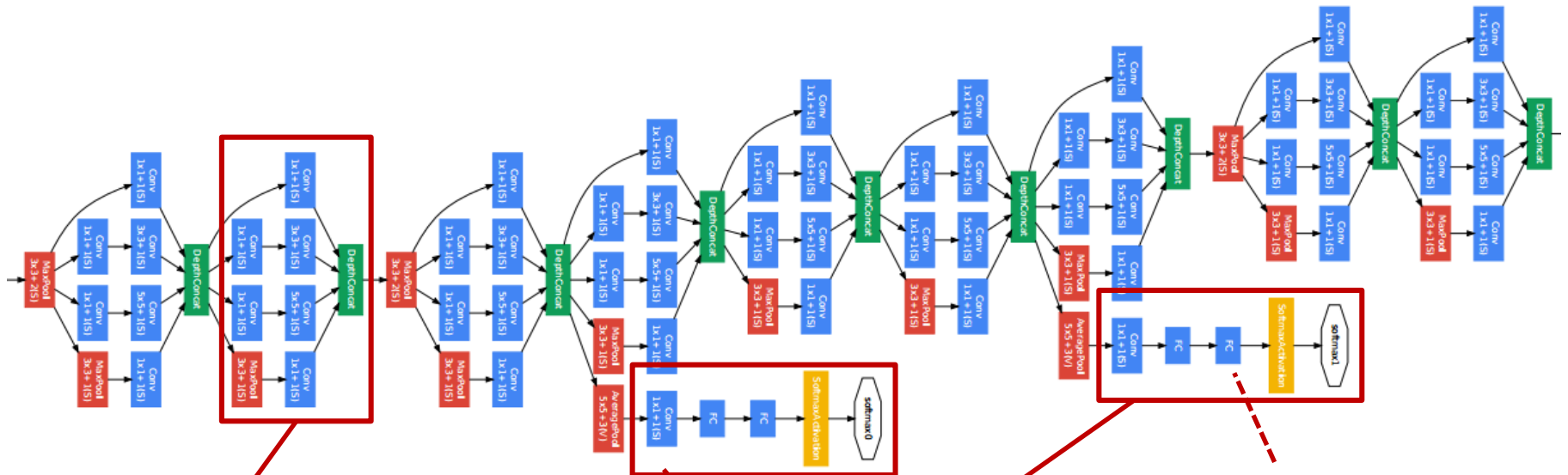


GoogLeNet (Inception v1) [Szegedy, CVPR, 2015]



GoogLeNet (Inception v1) [Szegedy, CVPR, 2015]

► Réseau entièrement repensé pour être très profond



1) Blocs répétitifs
Inception module

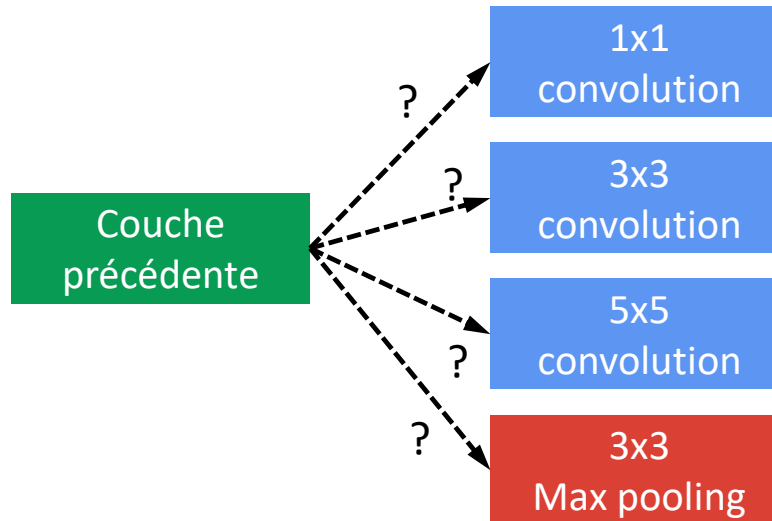
2) Fonctions de pertes intermédiaires
afin d'injecter des gradients dans les
couches intermédiaires

3) Couches FC remplacées
par du average pooling
(moins de paramètres)

Inception module

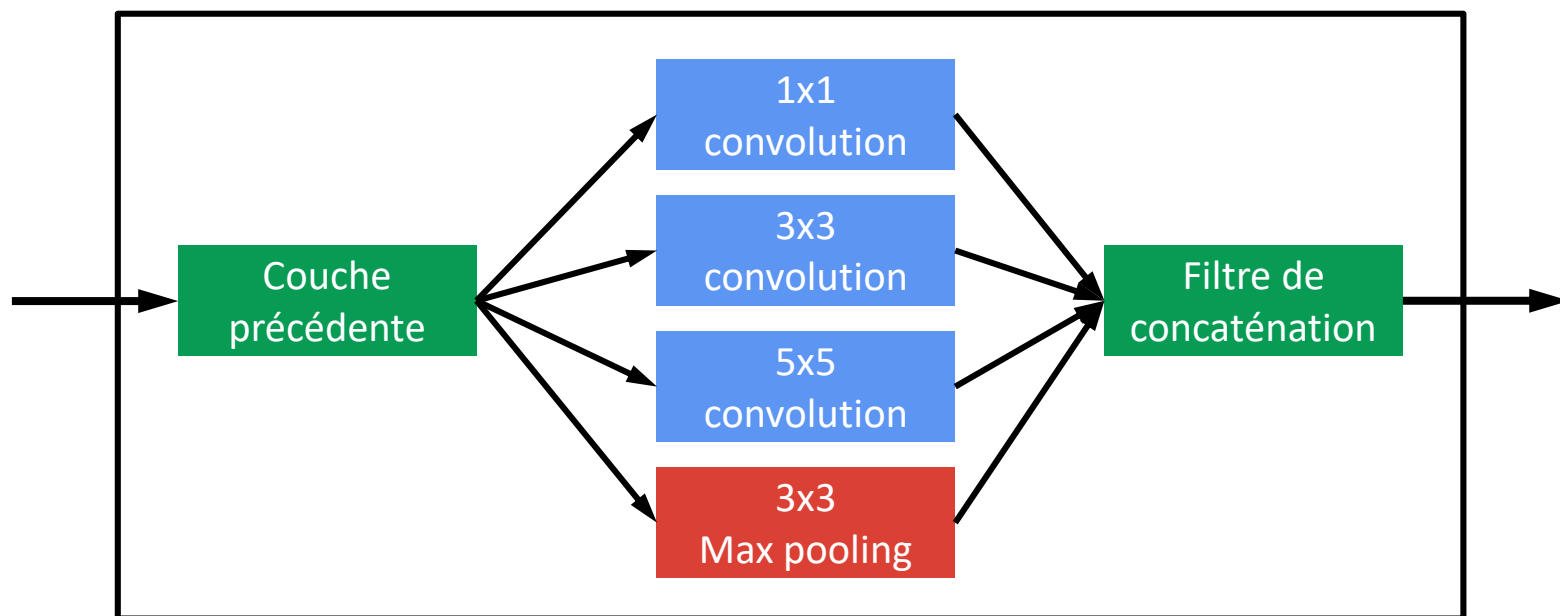
► Choix pour chaque couche

- Convolution ou pooling ?
- Si convolution, quelle taille de filtre ?



Inception module – Idée clé

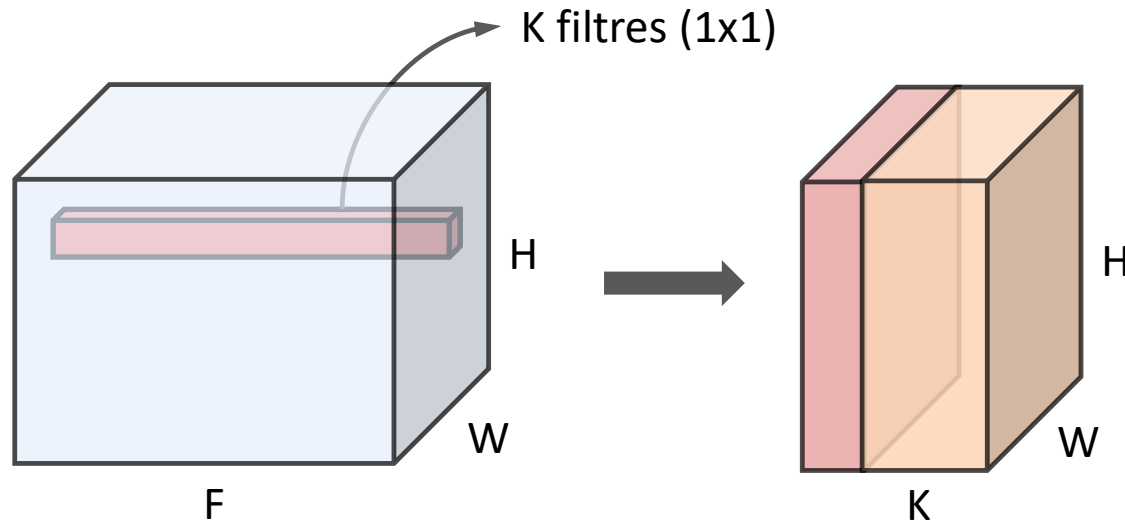
- Calculer toutes les sorties en parallèle
- Concaténation des résultats
- Laisser l'apprentissage choisir !



Problème: cela produit trop de sorties et de paramètres

Inception module – Idée clé

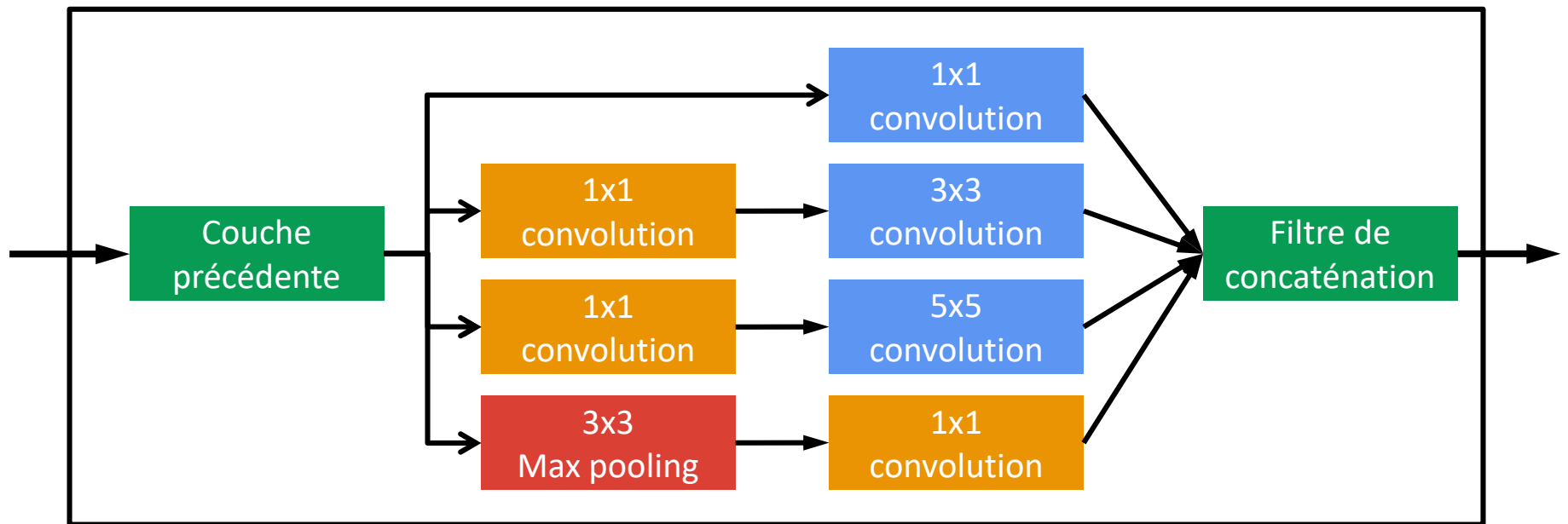
► Convolution 1×1



Réduction de la dimension pour $K < F$

Agit comme une fonction de « feature pooling » que l'on peut apprendre

Inception module – Idée clé

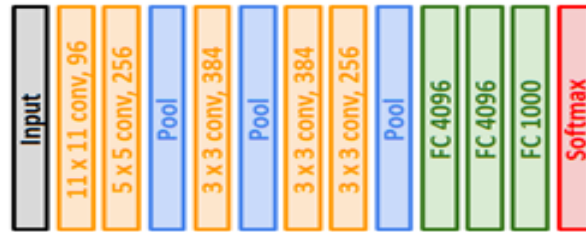


Réduction des dimensions via l'utilisation de couches d'étranglement composées de convolutions 1×1

Efficacité

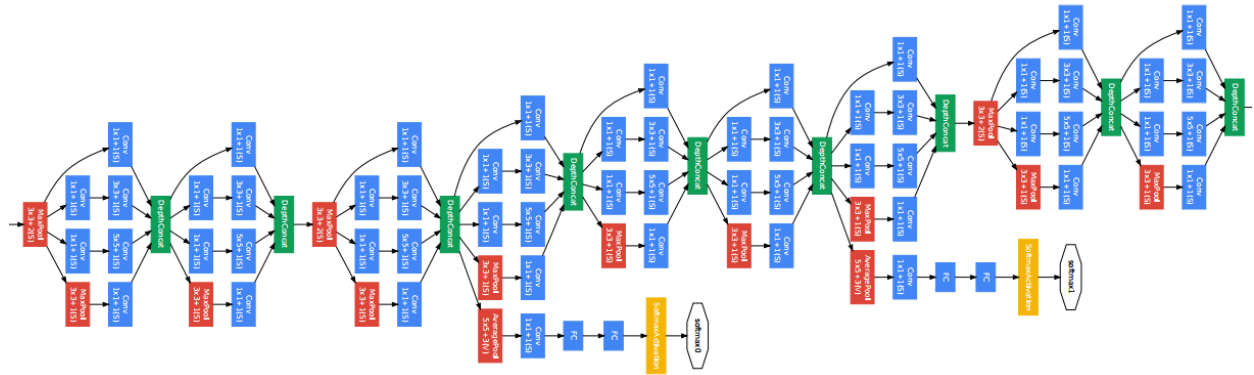
AlexNet

8 couches
param ~ 62M



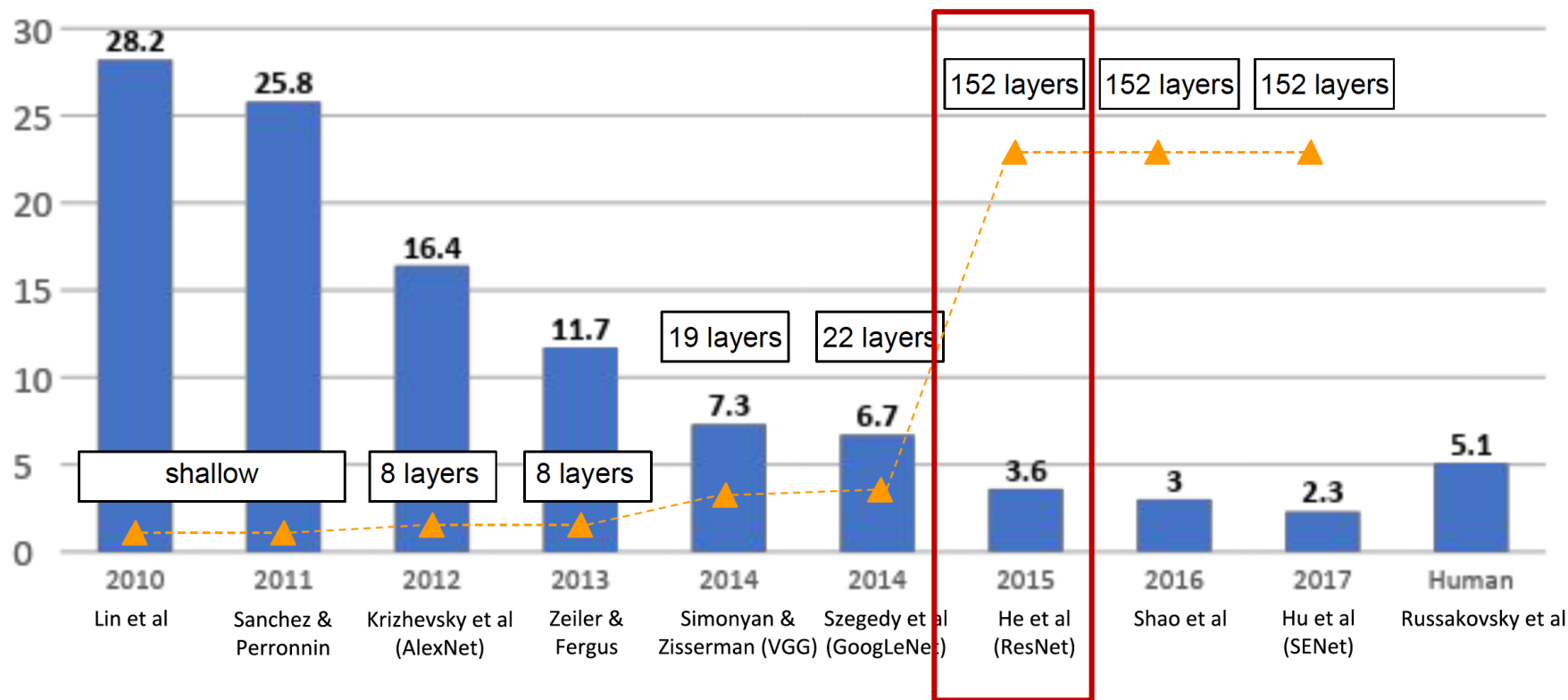
GoogLeNet

22 couches
param ~ 5M



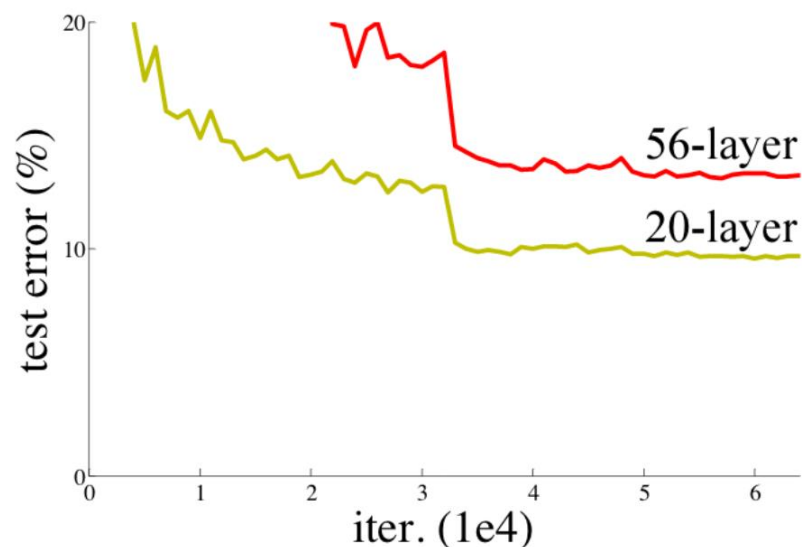
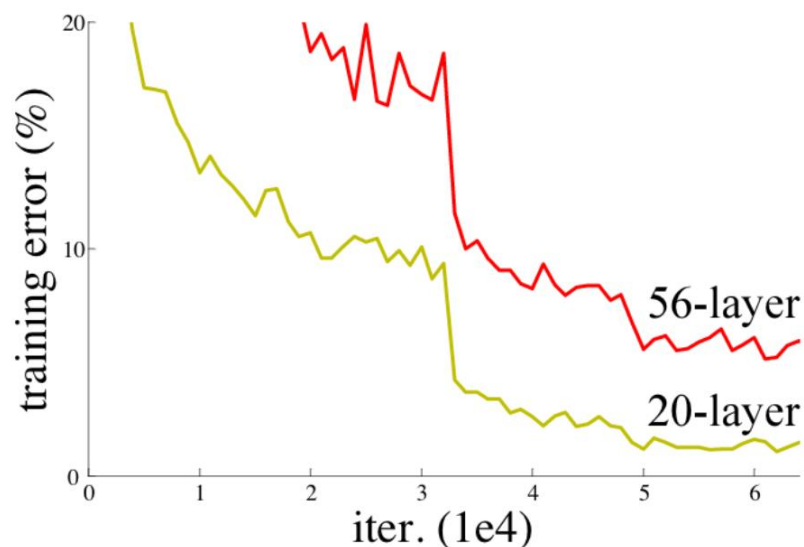
GoogLeNet possède 12x moins de paramètres qu'AlexNet !

ResNet [He, CVPR, 2016]



ResNet [He, CVPR, 2016]

- Que se passe t'il pour un réseau encore plus profond ?



- Erreurs d'entraînement et de test plus grandes !

Problème d'optimisation (*vanishing gradient*)

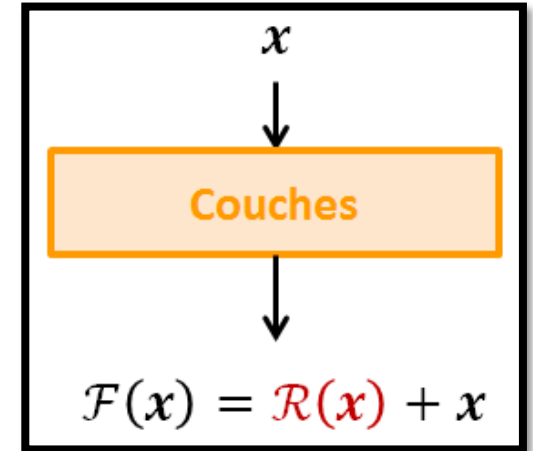
ResNet [He, CVPR, 2016]

► Idée clé

Estimer le résidu plutôt que la transformation

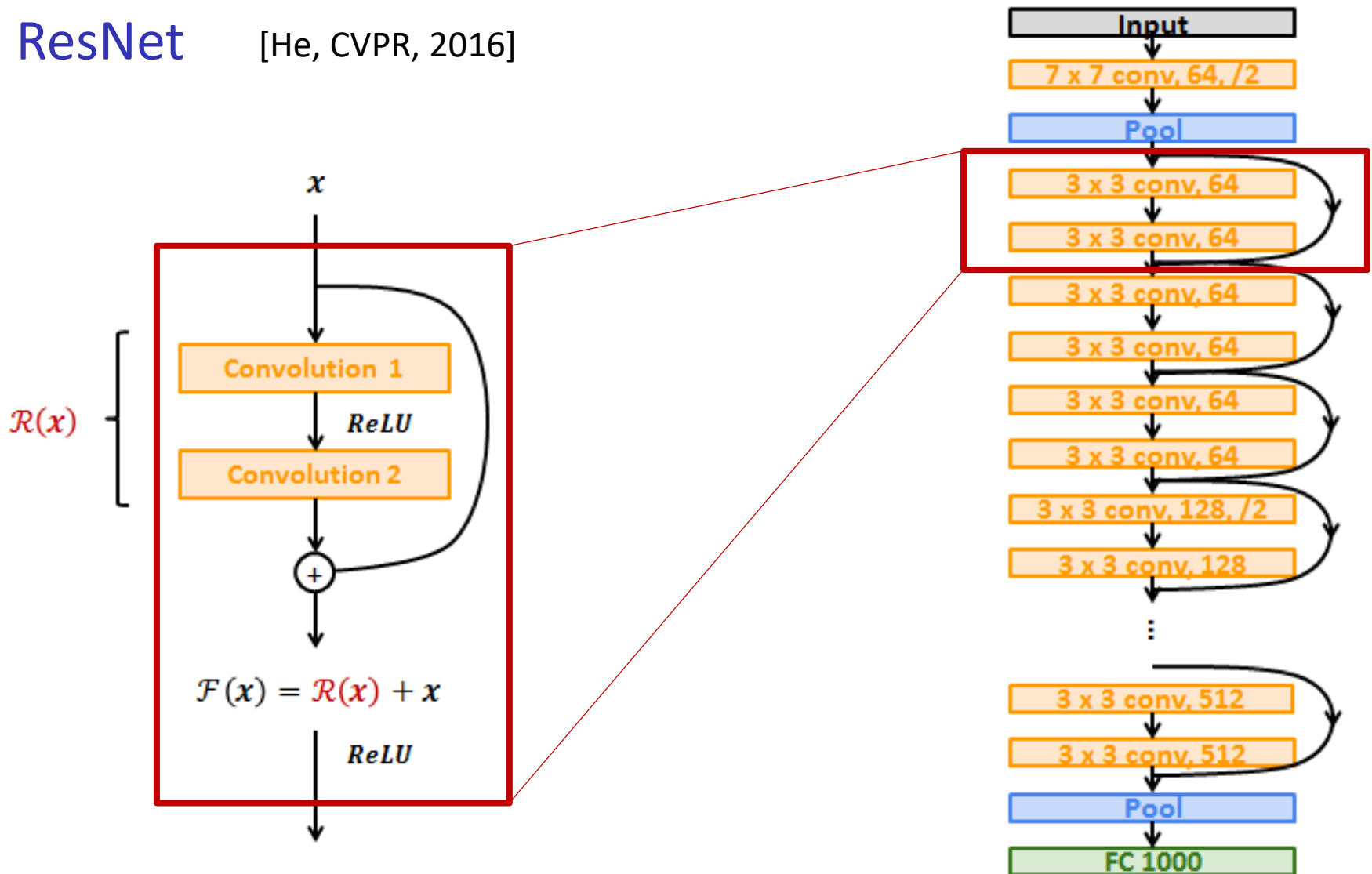
► Avantages

- Modélisation de moins d'information, potentiellement plus simple à apprendre
- Connexions résiduelles préservent le flux de gradient lors de la propagation arrière
- Conception d'architectures très profondes (> 100 couches)



ResNet

[He, CVPR, 2016]

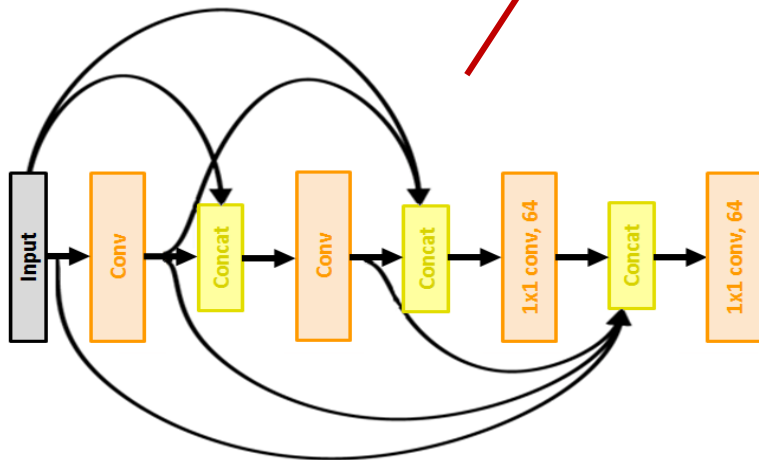
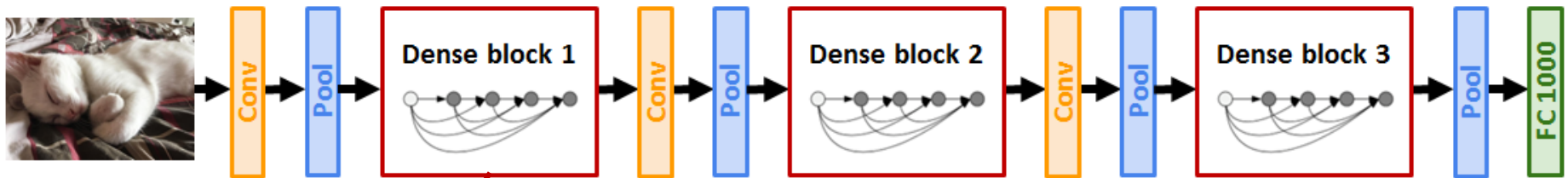


DenseNet (Densely connected)

[Huang, CVPR, 2017]

► Idée clé

Les *features* calculées dans une couche sont concaténées avec les entrées de toutes les autres couches d'un bloc



- Utilisation efficace de *features* multi-échelles
- Propagation du gradient au travers de chaque couche lors de la propagation arrière

Applications

Segmentation sémantique

Segmentation sémantique d'images

- ▶ Prédire la bonne classe pour chaque pixel d'une image



Image d'entrée

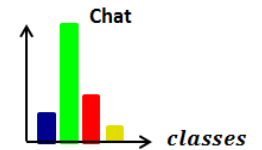
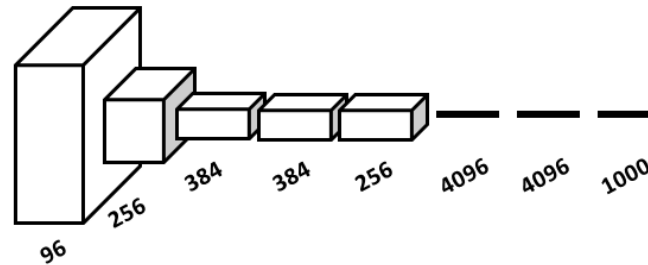


Segmentation

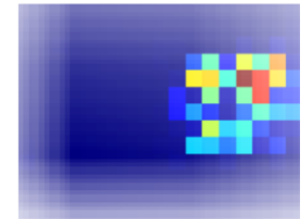
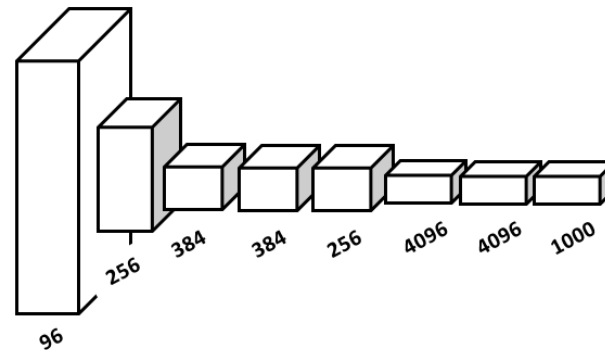
Peut-être vu comme un problème de classification dense et structuré

Fully-CNN: de la classification à la segmentation

Classification
standard CNN



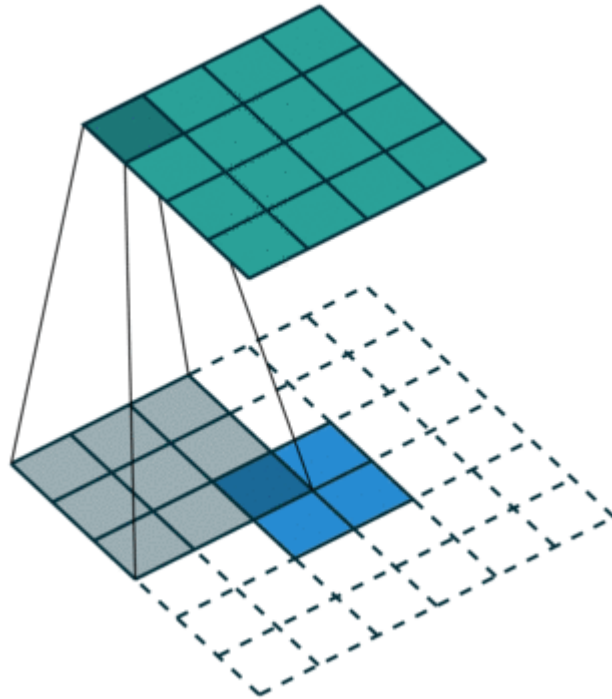
Fully-CNN



- Génération de cartes de segmentation très grossières

Ajout d'opérations de sur-échantillonnage à la fin du
réseau

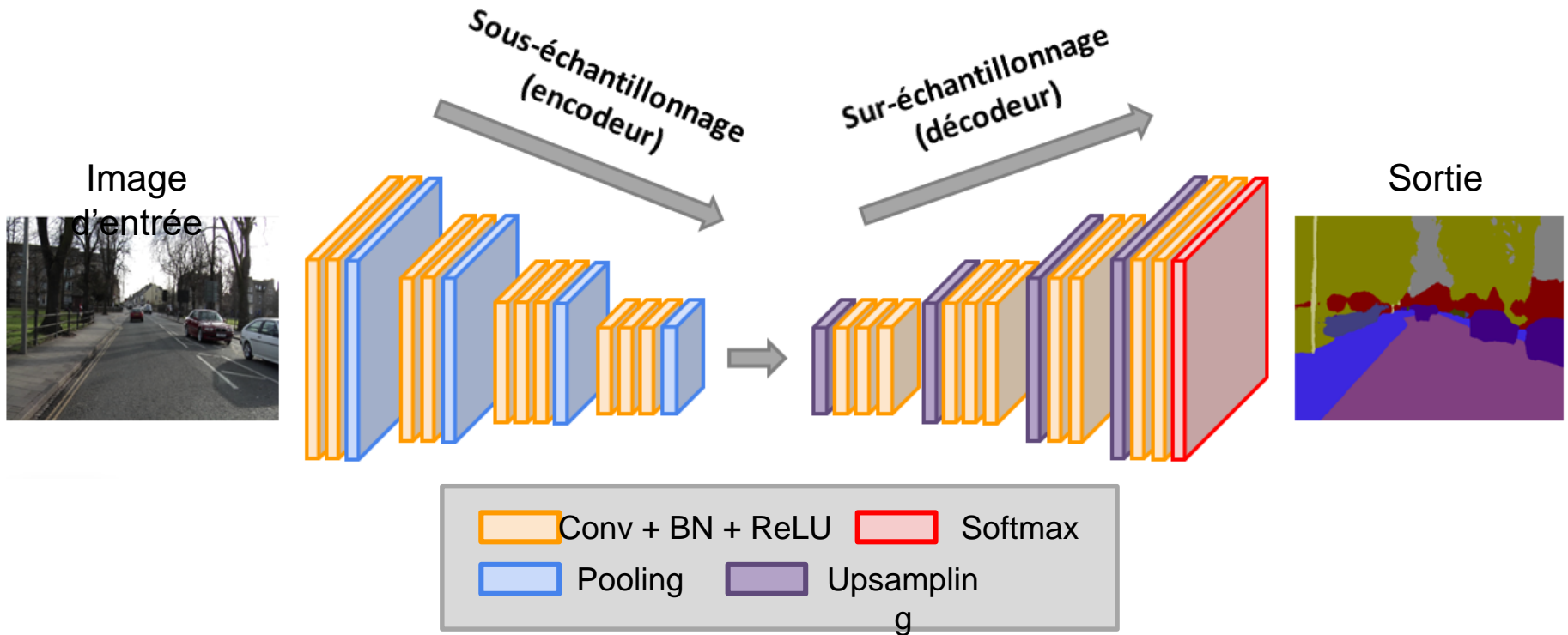
Opération de sur-échantillonnage



Padding 2, stride = 1

Dimensions spatiales augmentées

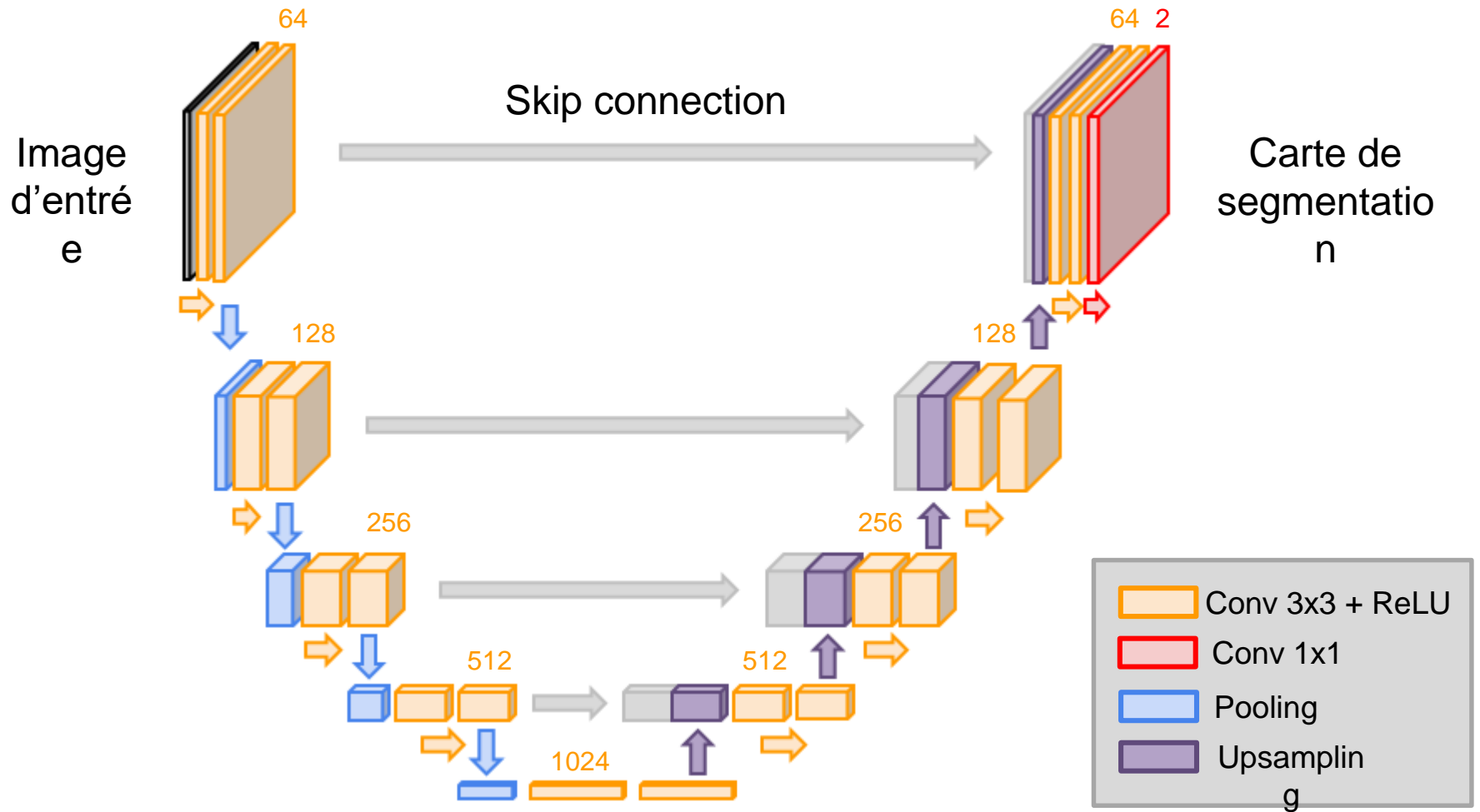
Architectures de type encodeur / decodeur



- Résolution spatiale perdue durant le sous-échantillonnage

Ajout de « skip connections » entre l'encodeur et le décodeur

U-Net



nnU-Net

► Méthode basée sur l'architecture U-Net

- Mise en forme automatique de la base de données
- Planification des patches / topologie de réseaux / batches
- Stratégie efficace de généralisation des réseaux
 - ✓ Augmentation de données en entraînement et en inférence
 - ✓ Nombre d'itération maximum
 - ✓ Décroissance progressive du taux d'apprentissage
 - ✓ Fusion de résultats de plusieurs U-Net

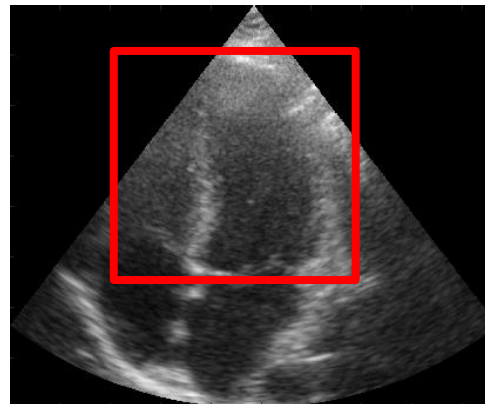
nnU-Net

- ▶ Mise en forme automatique de la base de données
 - Rééchantillonnage de l'ensemble de la base de données suivant la résolution médiane
 - ✓ Permet de travailler sur des images homogènes vis-à-vis des dimensions physiques
 - Normalisation des intensités
 - ✓ Tiens compte uniquement des pixels non nulles.
 - ✓

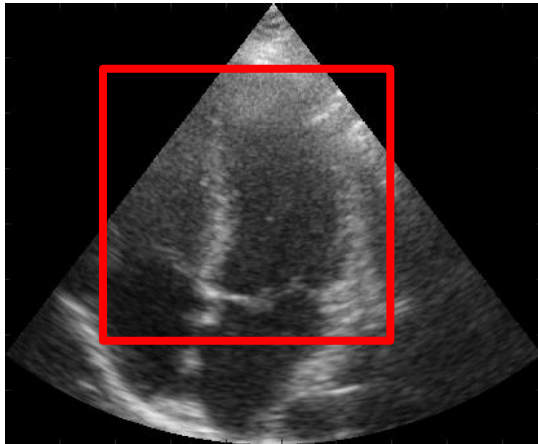
nnU-Net

- ▶ Planification des patches / topologies de réseaux / mini-batch
 - Dimensionnement en tenant compte de la capacité du GPU
 - Données d'entrées: patches de taille maximale
 - ✓ permet d'éviter de redimensionner les images

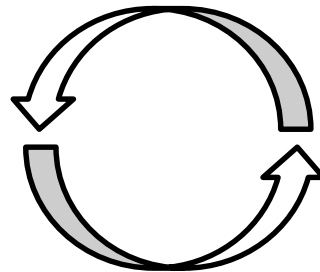
Patch à mettre
en entrée d'un
réseau U-Net



Patches



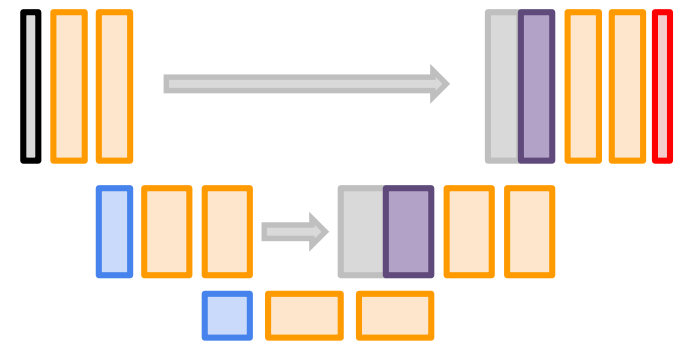
Dimensions max.



Mini-batch

Taille min. = 2

Topologie de réseaux



Profondeur max.

Jusqu'à feature map 2x2

nnU-Net

► Stratégie de généralisation des réseaux

- Sélection aléatoire de patch dans chaque image lors de la création des batch
- 1000 epoch avec 250 itérations par epoch
 - ✓ On ne verra pas toute la bd par epoch => bagging !
- Décroissance progressive du taux d'apprentissage au cours des epochs
 - ✓ Convergence progressive vers un optimum

nnU-Net

► Stratégie de généralisation des réseaux

- Augmentation de données lors de la génération de batch
 - ✓ Rotation, changement d'échelle, effet miroir
 - ✓ Changement d'intensité, floutage gaussien, contraste, luminosité
- Augmentation de données en inférence
 - ✓ Segmentation d'une image et de ses versions miroir
 - ✓ Moyennage des résultats obtenus après effet miroir inverse pour filtrer les zones incertaines

nnU-Net

► Pour la segmentation d'images 3D

- Apprentissages indépendants de plusieurs réseaux
 - ✓ U-Net 2D / patches 2D
 - ✓ U-Net 3D / patches 3D
 - ✓ U-Net 3D / volumes entiers sous résolus
- Approche de type Boosting
 - ✓ Moyennage des résultats de deux architectures sur la base de données de validation au cours de l'apprentissage
 - ✓ La meilleure combinaison de réseaux est utilisée en inférence

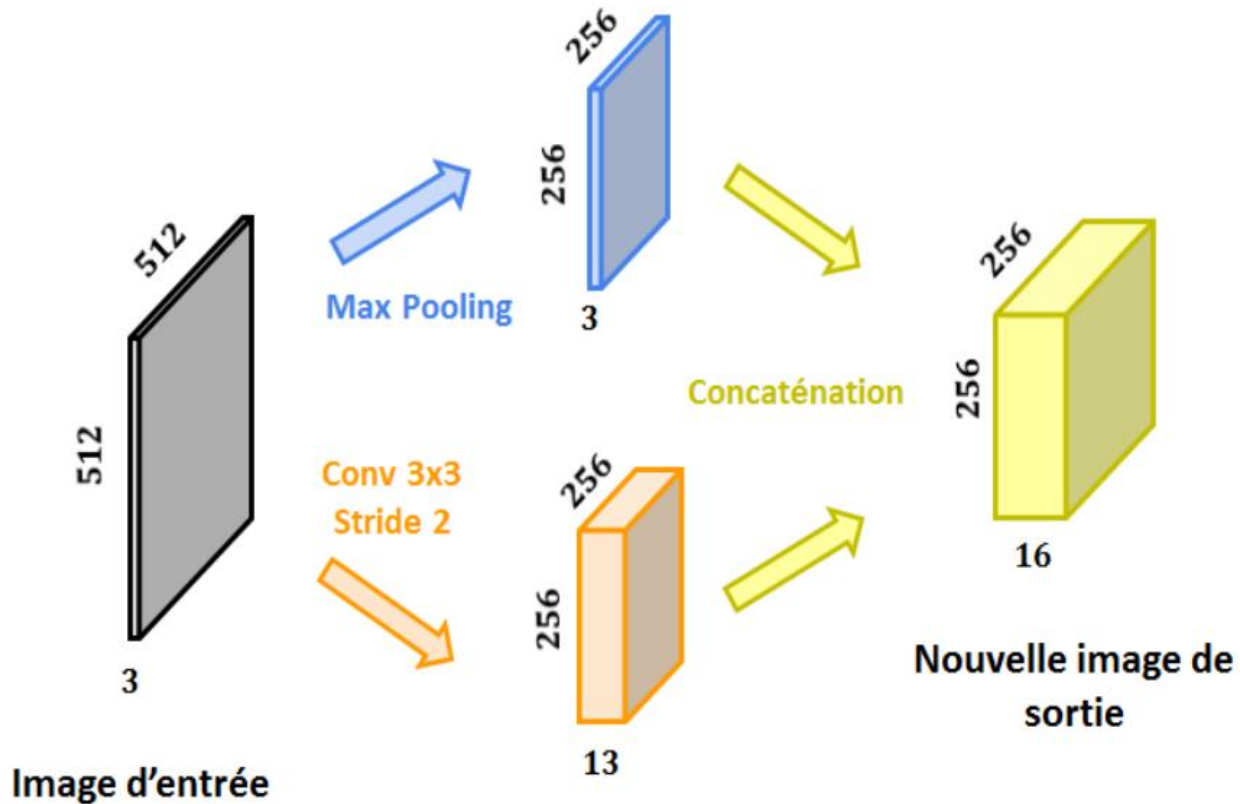
nnU-Net

► Segmentation en inférence

- Utilisation d'une fenêtre glissante
 - ✓ Patches avec un chevauchement de la moitié de la taille du patch
 - ✓ Moyennage des résultats obtenus

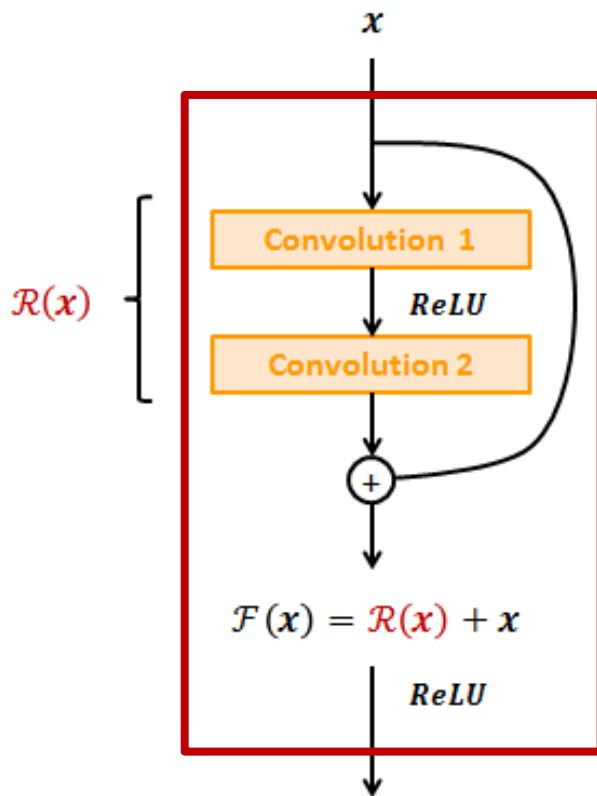
ENet

► Modélisation de l'image d'entrée

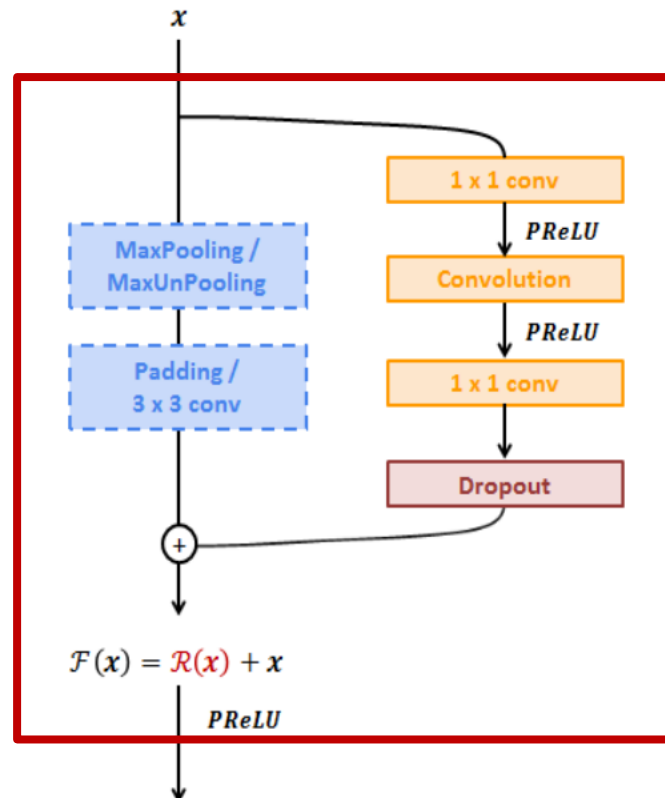


ENet

► Architecture basée sur ResNet



Bloc ResNet



Bloc ENet

*Projection / réduction
de dimensionnalité*

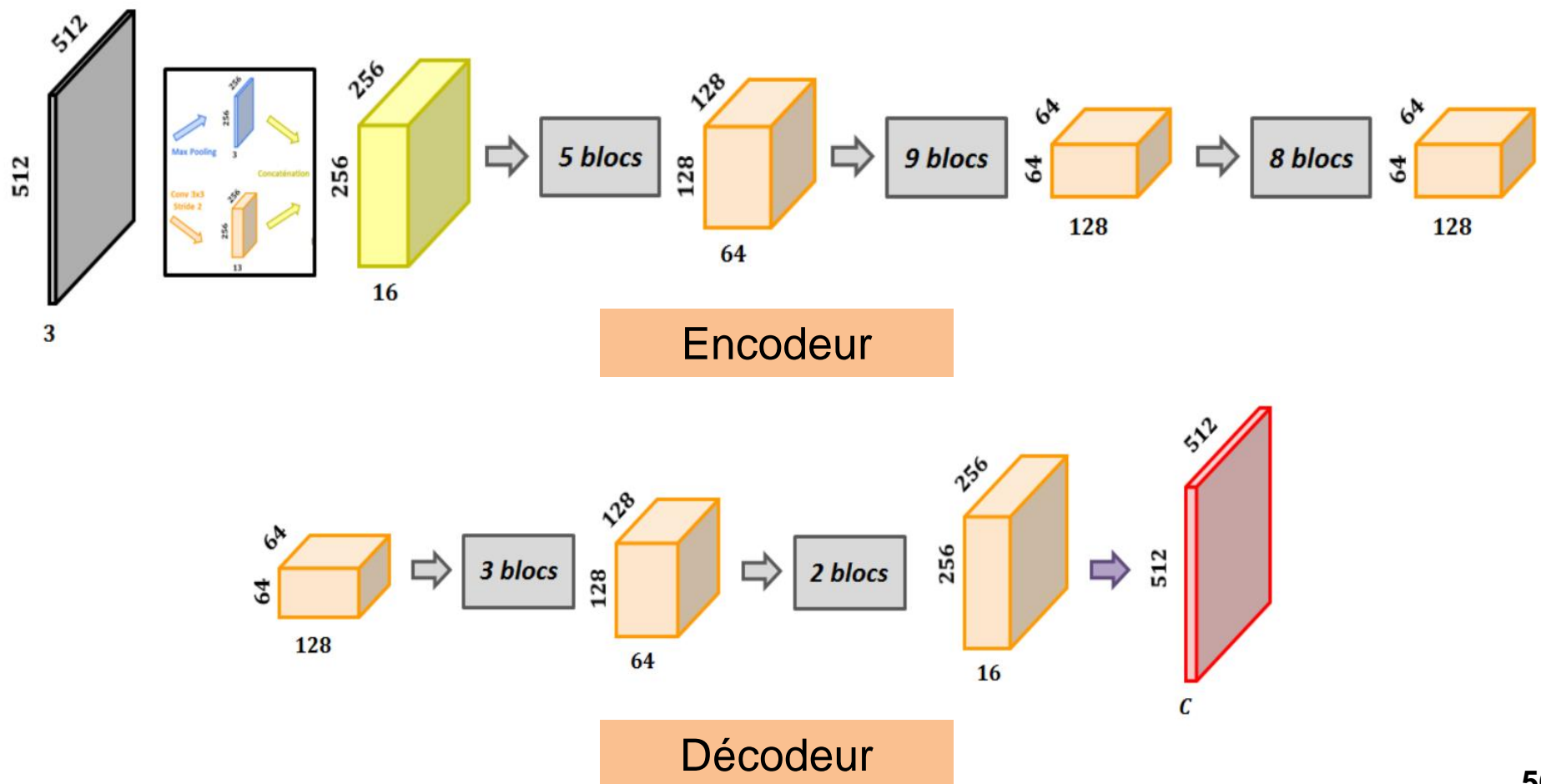
Création d'information

*Expansion de l'info au
travers de canaux*

Régularisation

ENet

► Architecture asymétrique



ENet

► Performances

- Qualité de segmentation équivalente voir meilleure que l'état de l'art en apprentissage profond
- # paramètres: 0.37 M
- Taille du réseau < 6 MB
- Temps d'exécution (NVIDIA TitanX)

640x360 px	7 ms
1280x720 px	21 ms
1920x1080 px	46 ms

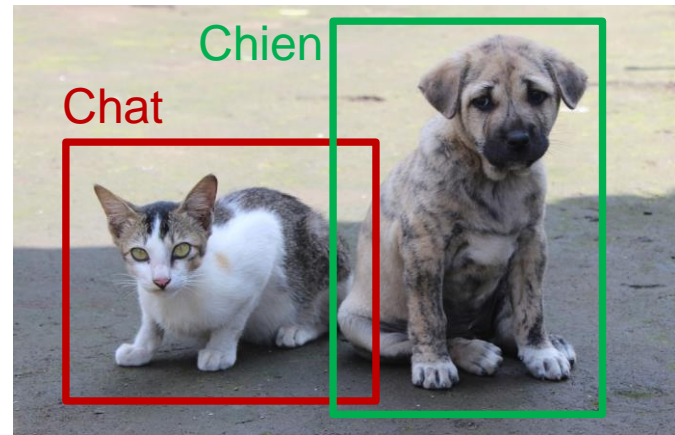
Applications

Détection d'objets

Détection d'objets



Image d'entrée

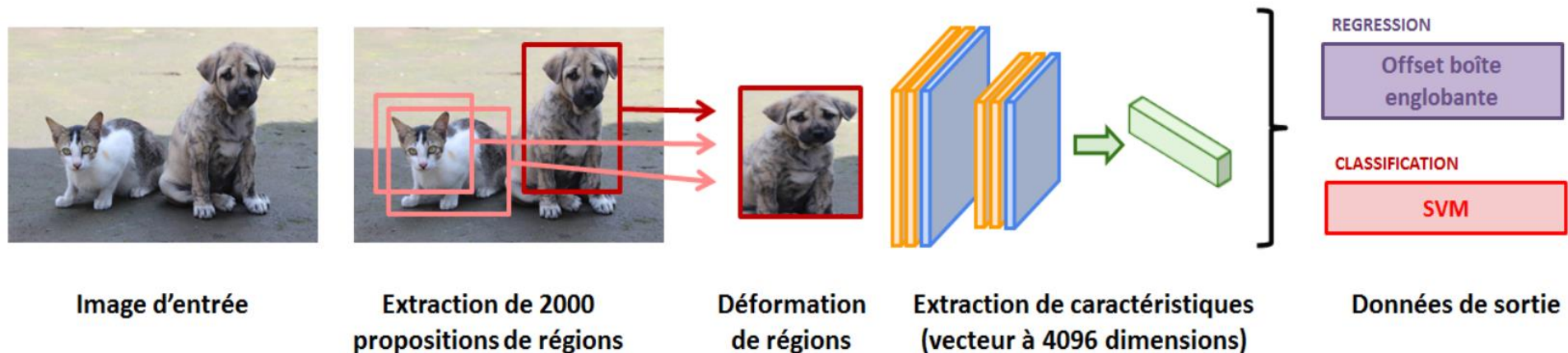


Classes détectées

Trouver les objets/classes présents dans une image ainsi que leur localisation

R-CNN (Region-CNN)

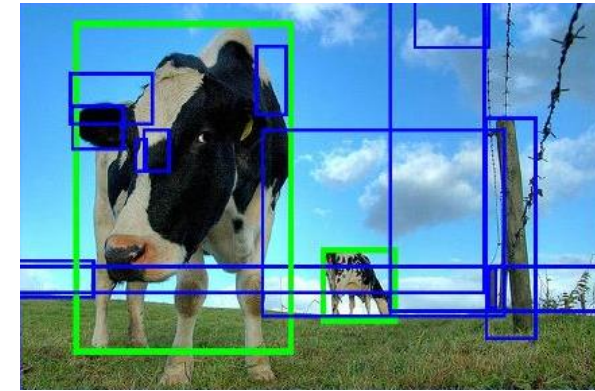
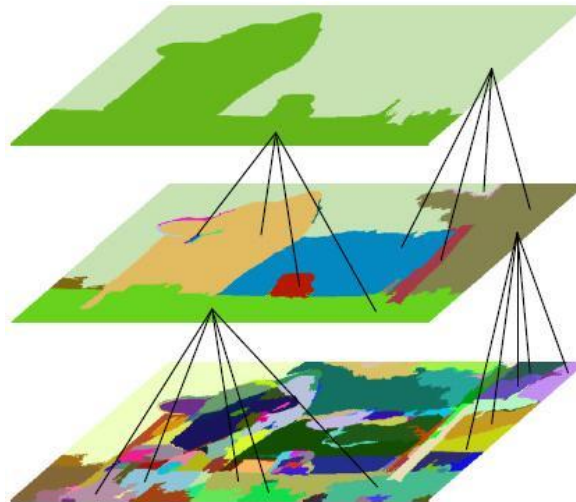
► Phase d'entraînement



Plusieurs étapes d'entraînement (CNN, SVM, régression boîtes englobantes)

R-CNN

► Extraction de régions



- Méthode classique par graphe

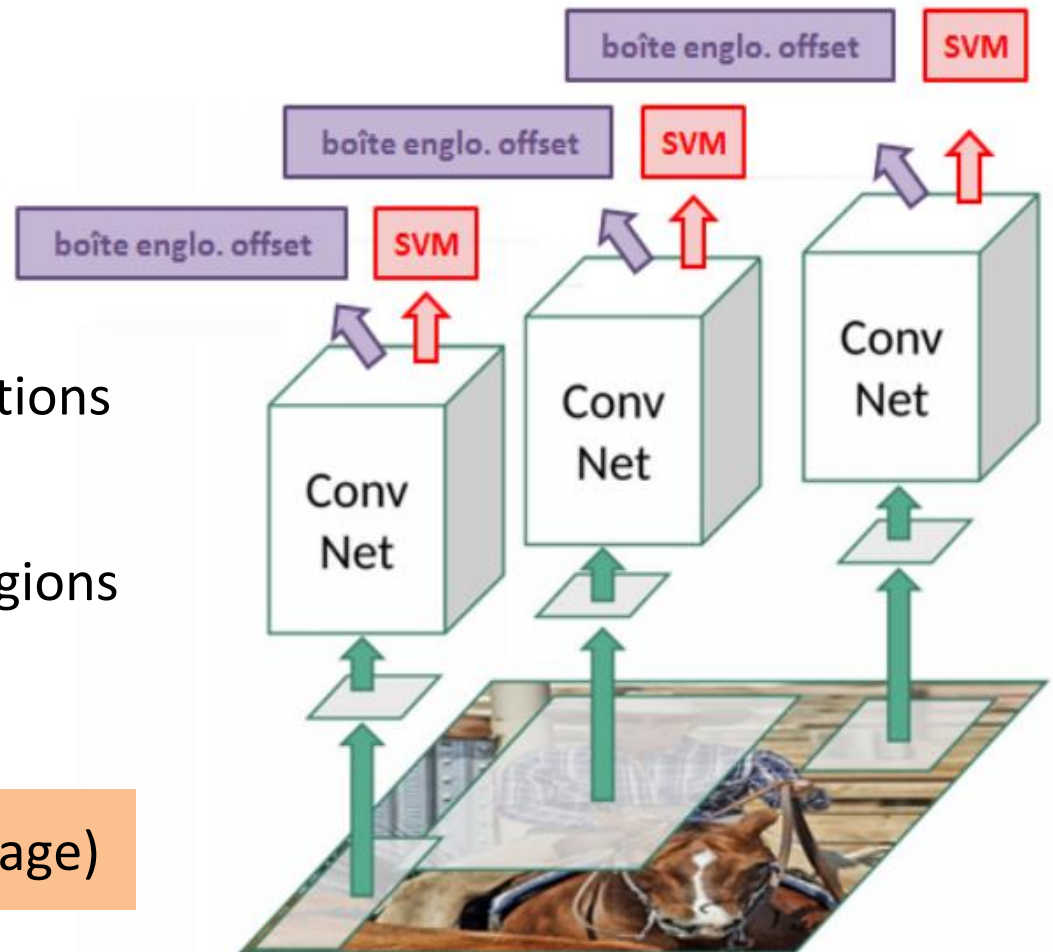
Génération non négligeable de mauvais candidats

R-CNN (Region-CNN)

► Phase d'inférence

- Extraction des propositions de régions
- Inférences sur 2000 régions

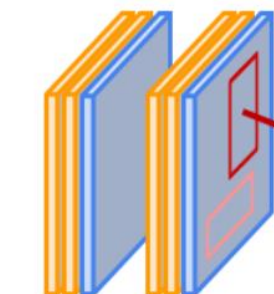
Très lent (~50 s par image)



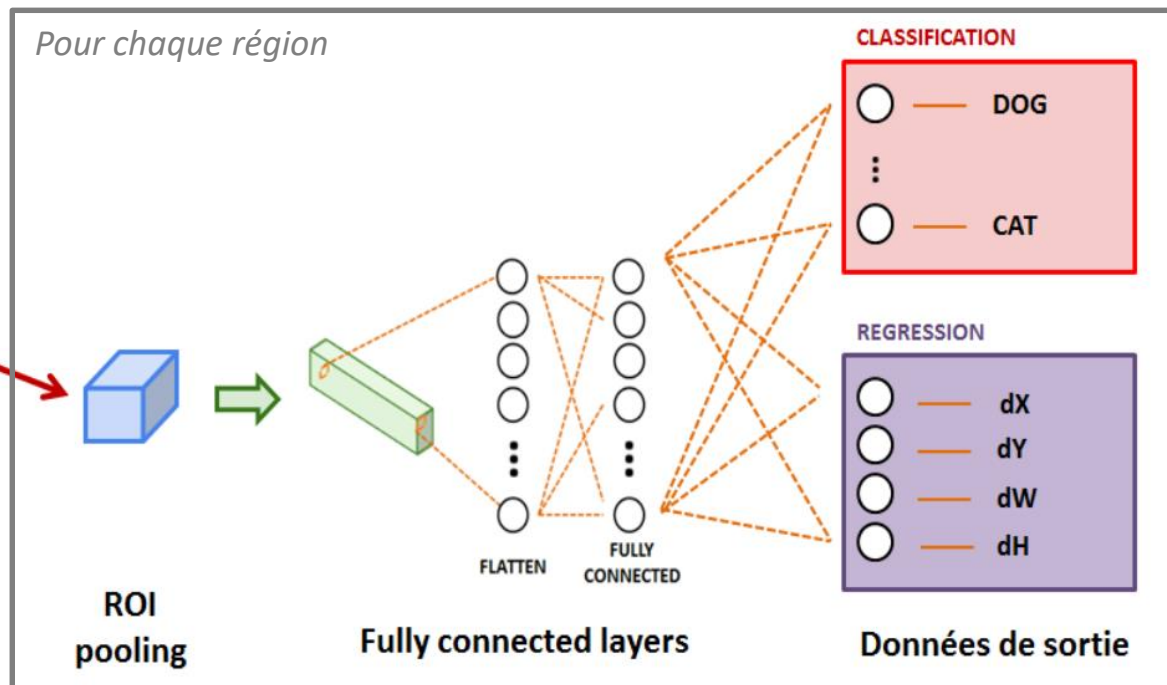
Fast R-CNN



Image d'entrée



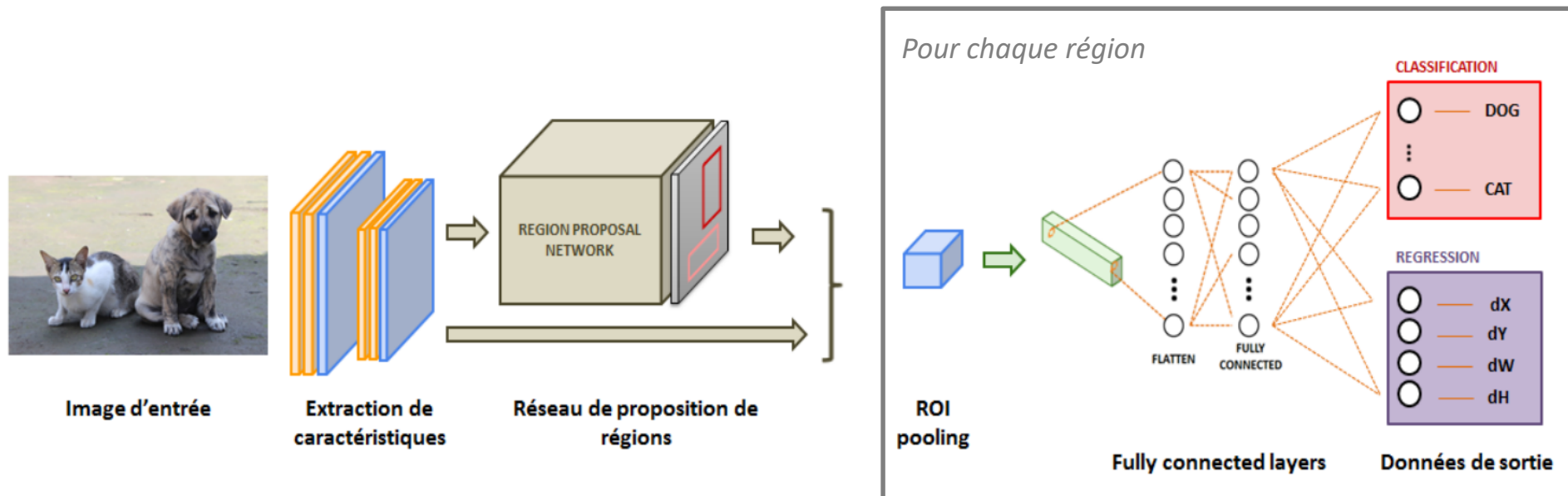
Extraction de caractéristiques



- 20x plus rapide que R-CNN en inférence !

Extraction de proposition de régions reste un point faible de la méthode

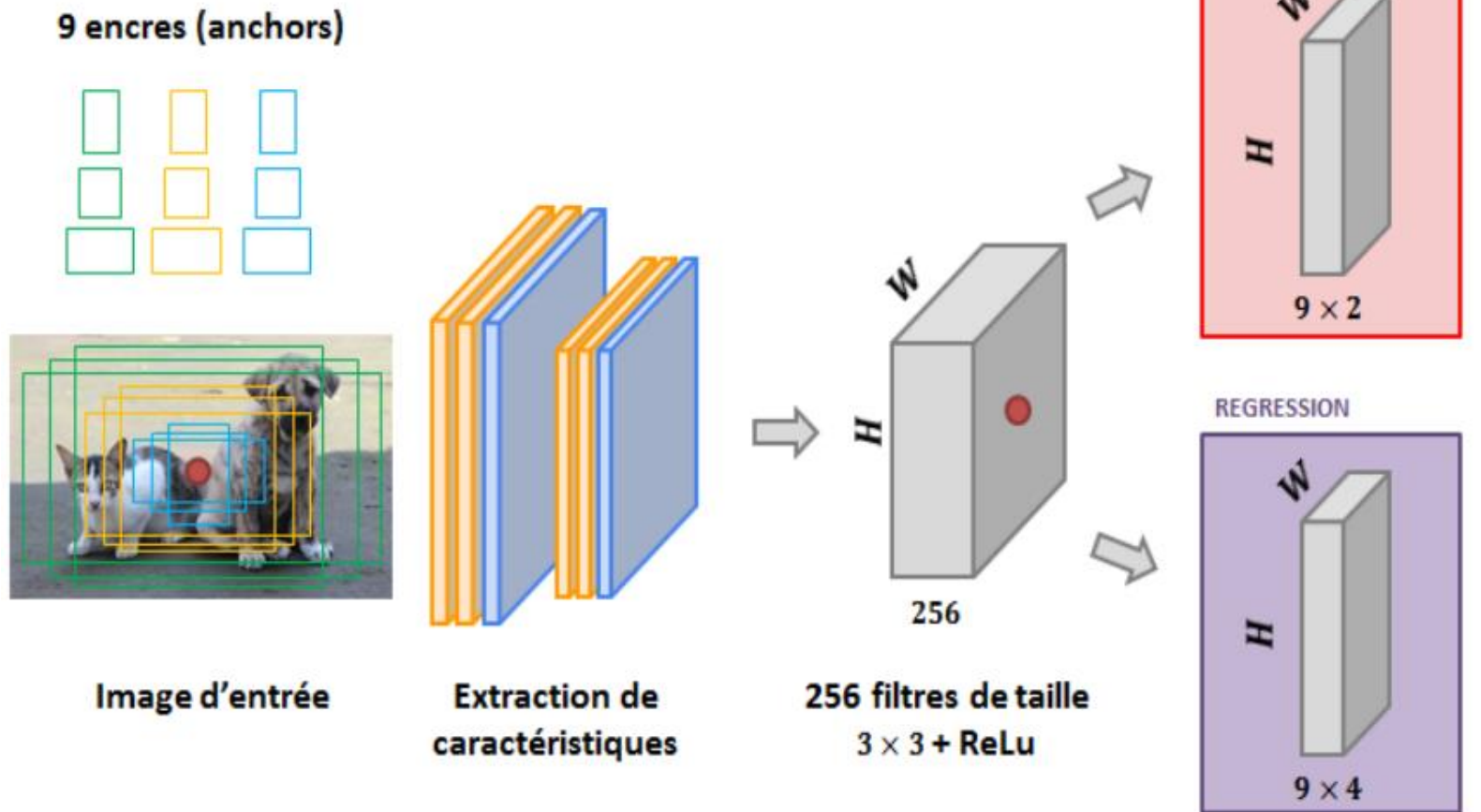
Faster R-CNN



- Réseau de proposition de régions intégré

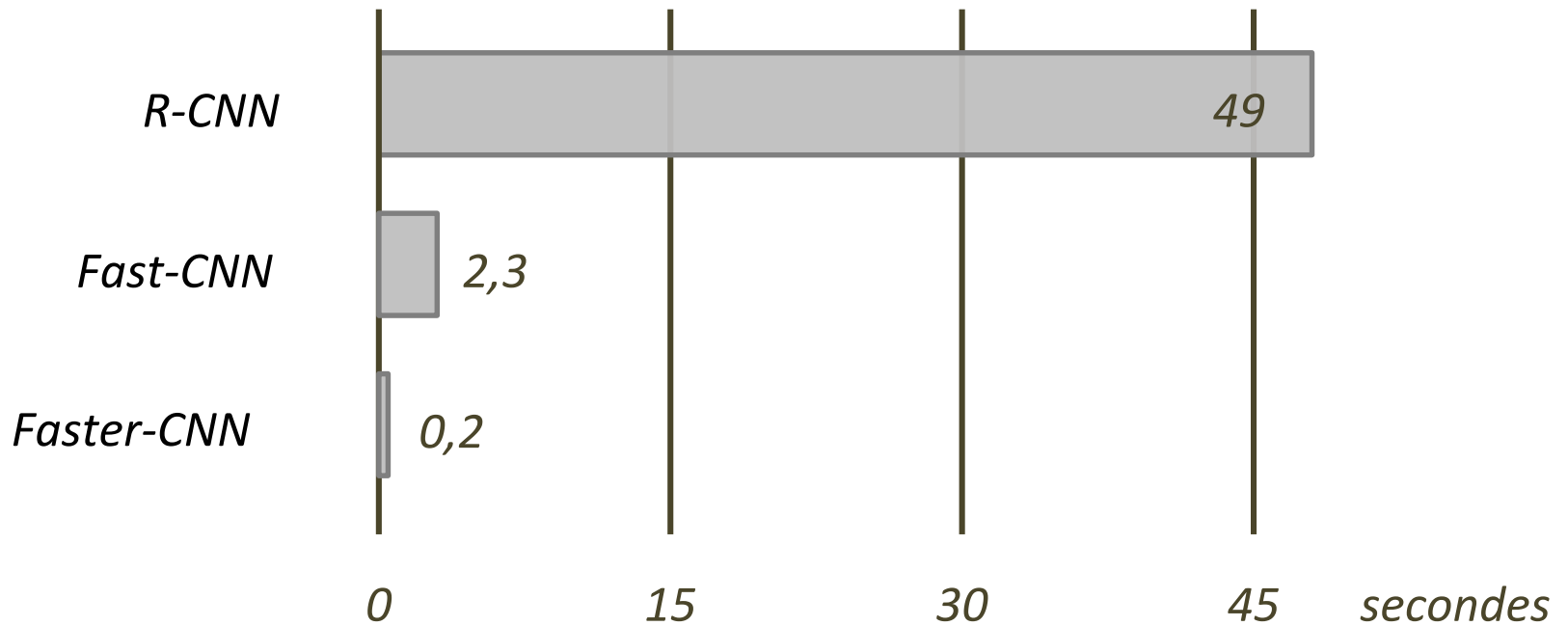
Réseau global entraînable de bout en bout !

Region Proposal Network (RPN)



$W \times H \times 9$ régions détectées

Performance d'exécution (phase d'inférence)



Applications

Segmentation d'instances

Segmentation d'instances



Image d'entrée

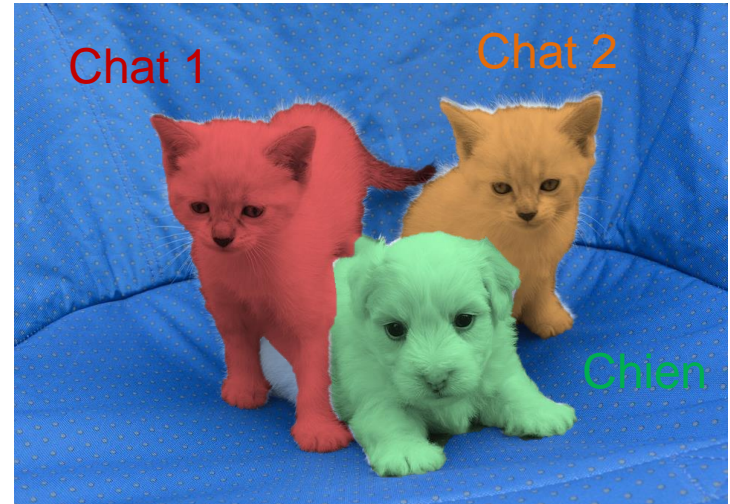
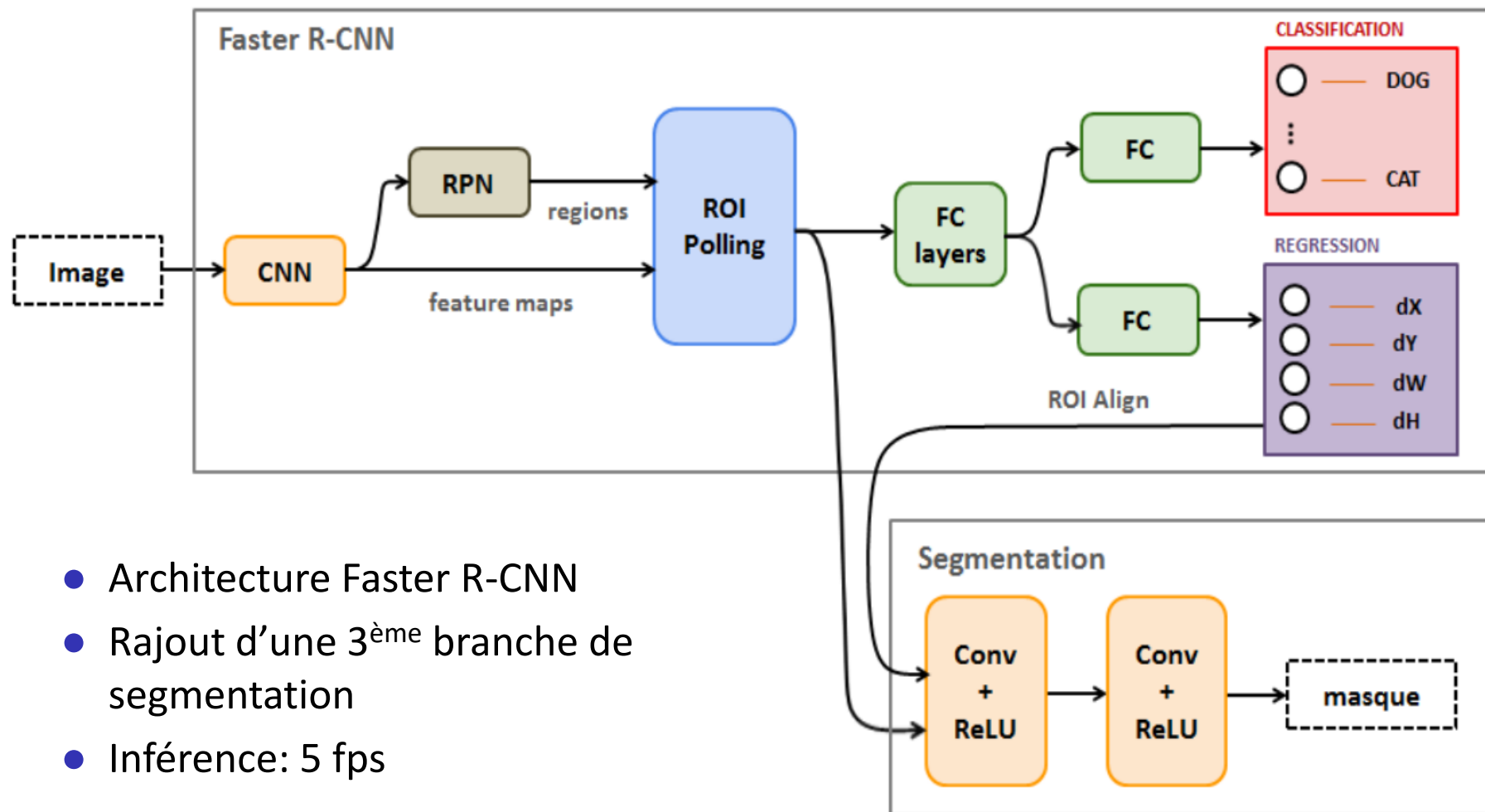


Image segmentée

Détecter et segmenter toutes les instances des objets/classes présents dans l'image

Mask R-CNN



- Architecture Faster R-CNN
- Rajout d'une 3^{ème} branche de segmentation
- Inférence: 5 fps

That's all folks
