

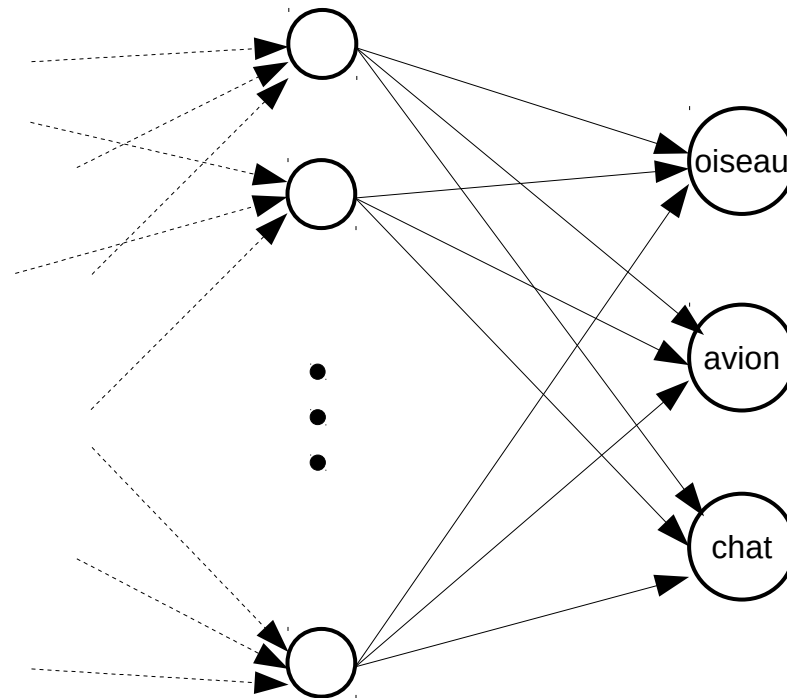
INTRODUCTION AUX RÉSEAUX DE NEURONES

Réseaux de neurones convolutifs
(«CNN» : *Convolutional neural networks*)

Pascal Germain*, 2018

* Merci spécial à [Philippe Giguère](#) pour m'avoir permis de réutiliser une partie de ces transparents.

Mais d'abord : Réseau avec sortie multiclasse

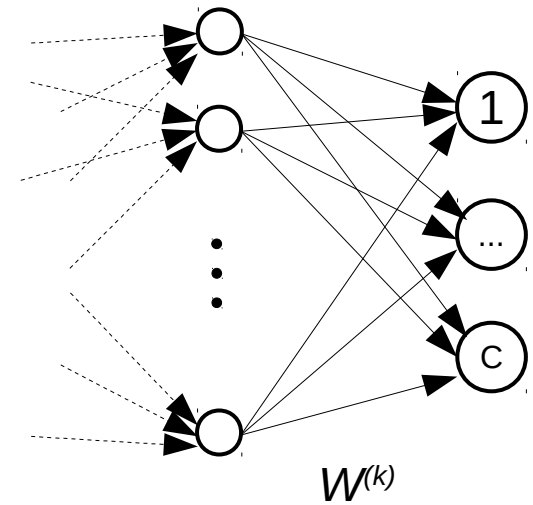


Réseau avec sortie multiclasse

Soit un problème à C classes.

1. Pour un exemple d'apprentissage (x, y) , on représente chaque classe par un entier

$$y \in \{1, \dots, C\}$$



Réseau avec sortie multiclasse

Soit un problème à C classes.

1. Pour un exemple d'apprentissage (x, y) , on représente chaque classe par un entier

$$y \in \{1, \dots, C\}$$

2. On converti y sous la forme d'un vecteur one-hot $\mathbf{y} \in \mathbb{R}^C$, possédant la valeur 1 à l'index correspondant à y , et les valeurs 0 autrement:

$$y = 0 \mapsto \mathbf{y} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad y = 1 \mapsto \mathbf{y} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

3. La couche de sortie sera donnée un vecteur $\hat{\mathbf{y}}$ obtenu en appliquant la fonction d'activation softmax aux valeurs $\mathbf{a} \in \mathbb{R}^C$ propagées par les couches précédentes:

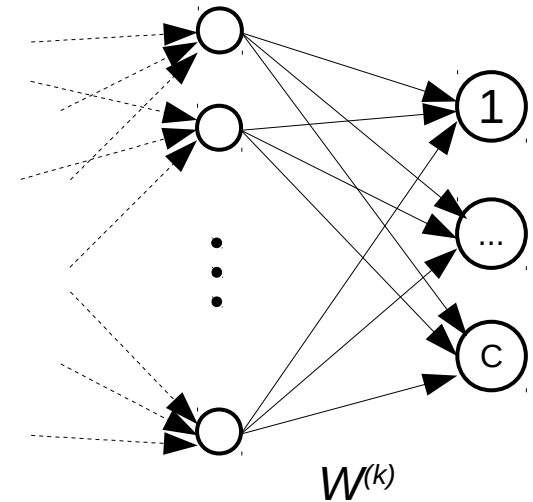
$$\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_C \end{bmatrix} \quad \text{avec } \hat{y}_i = \text{softmax}(a_i) = \frac{e^{a_i}}{\sum_{j=1}^C e^{a_j}}$$

On interprète le vecteur \mathbf{y} comme une distribution de probabilité sur les classes.

4. Lors de l'optimisation du réseau, on minimise la perte du négatif log-vraisemblance:

$$L_{\text{nlv}}(\hat{\mathbf{y}}, y) = -\ln(\hat{y}_y),$$

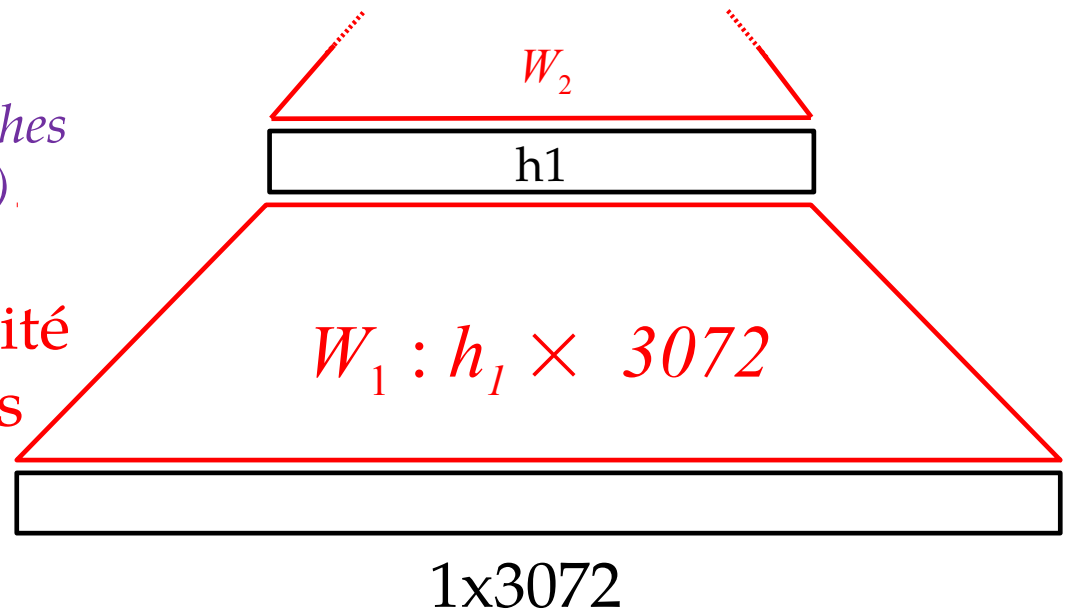
où \hat{y}_y correspond à la sortie d'index y du réseau (la probabilité associée à la classe de l'exemple).



Réseau pleinement connecté
(*fully connected*)

Aussi appelé *perception multi-couches*
(*multi-layers perceptron – MLP*).

grande quantité
de paramètres



Vectorisation
(*flatten*) de l'image

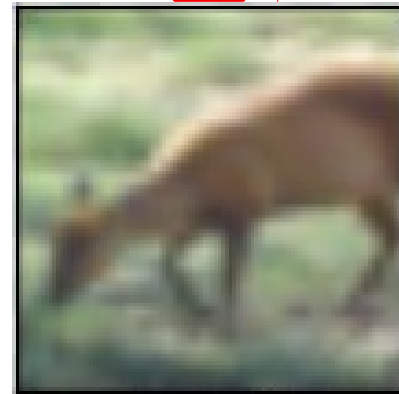
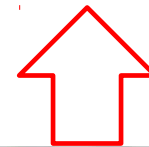


Image $32 \times 32 \times 3$

- **Vectorisation détruit :**
 - relations spatiales
 - canaux de couleurs

La connaissance du problème influence la conception de l'architecture

- Forte corrélation locale dans les valeurs des pixels
 - structure 2D

Convolution *

- Opération mathématique très utilisée :
 - Traitement de signal
 - Probabilités (somme de 2 variables aléatoires)

$$(I \star F)(i, j) = (F \star I)(i, j) = \sum_m \sum_n F(m, n) I(i - m, j - n)$$

- Au sens strict, les réseaux utilisent plutôt la corrélation croisée ★

$$(F \star I)(i, j) = \sum_m \sum_n F(m, n) I(i + m, j + n)$$

Exemple « convolution »

(qui est plus une corrélation croisée, mais bon...)

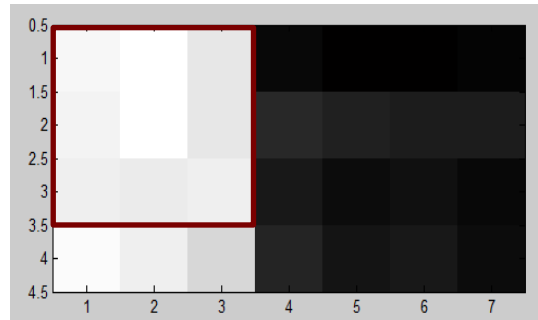
(Appellation de Filtre,
Filter, ou *Kernel*)

$$F_1 = \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix}$$

207	210	195	63	57	56	59
204	212	197	82	76	74	75
202	198	202	72	65	67	63
209	201	187	78	69	71	64

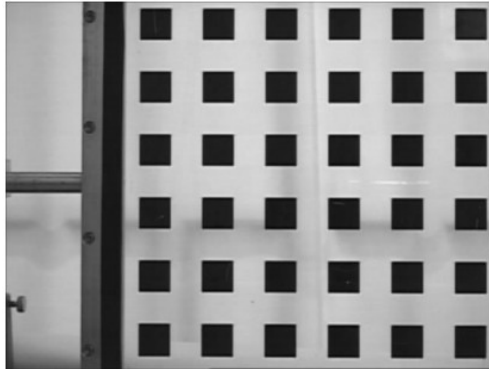
$$(I * F_1)(x, y)$$

	-26	-533	-517	-28		
	-29	-505	-513	-25		

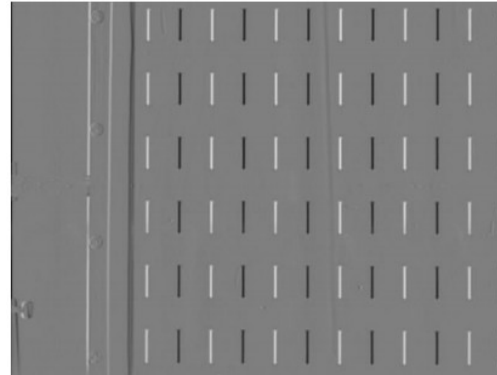


Exemples de filtres calibrés manuellement

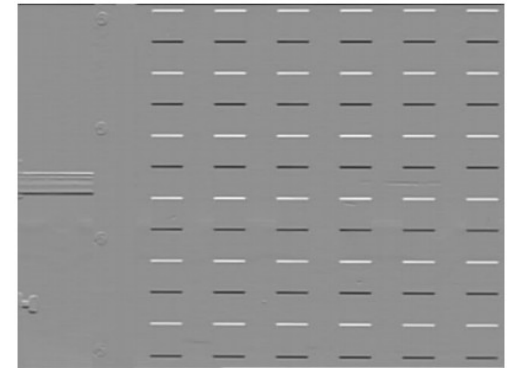
Détection bordure



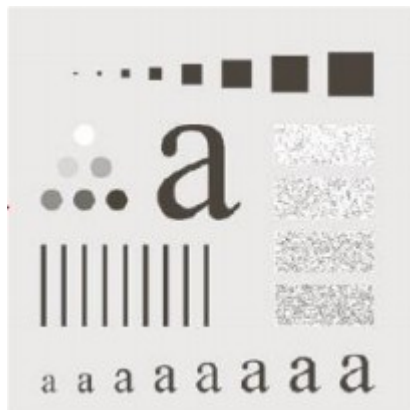
Bordure
verticale $\begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix}$



Bordure
horizontale $\begin{bmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix}$



Flou



$$\begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & 1 & 1 \end{bmatrix}$$

5x5



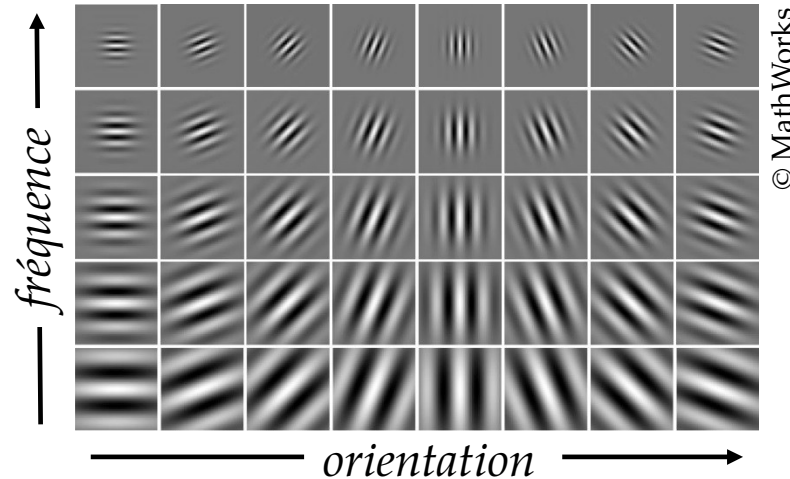
15x15



Filtres

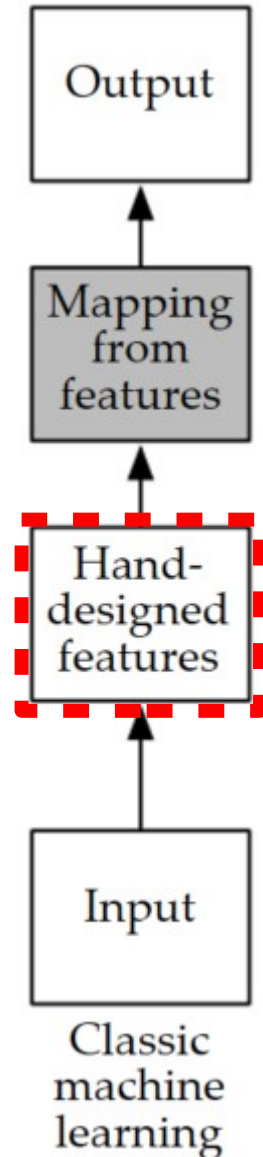
- Vont extraire des *features* de bas niveau

- Filtres *Gabor* :



- Filtres de bordure, ondelettes

- Longtemps été un domaine de recherche
 - que concevriez-vous comme filtre pour MNIST ?
- Comme les filtres CNN sont différentiables, le réseau pourra les modifier à sa convenance
 - les ajuster pour maximiser les performances sur les données d'entraînement

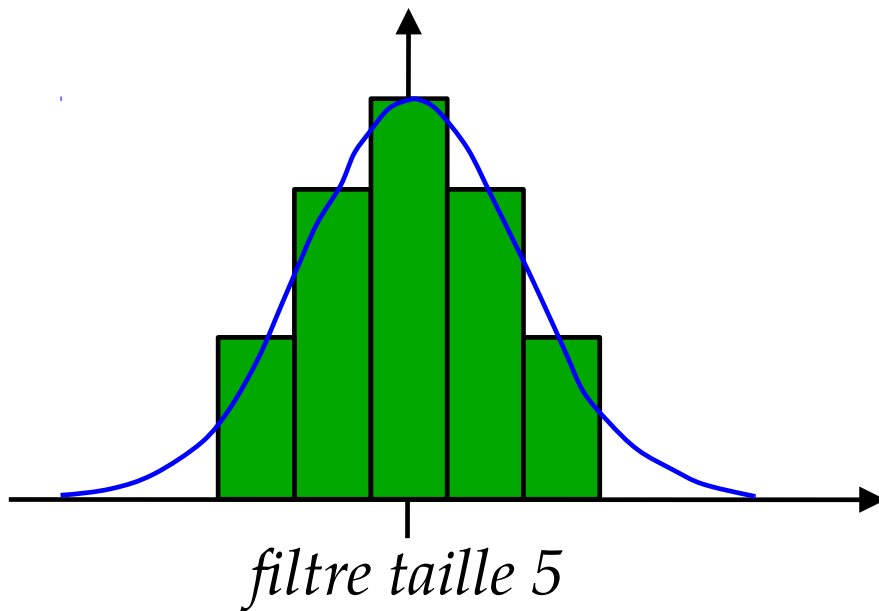


Exemple convolution

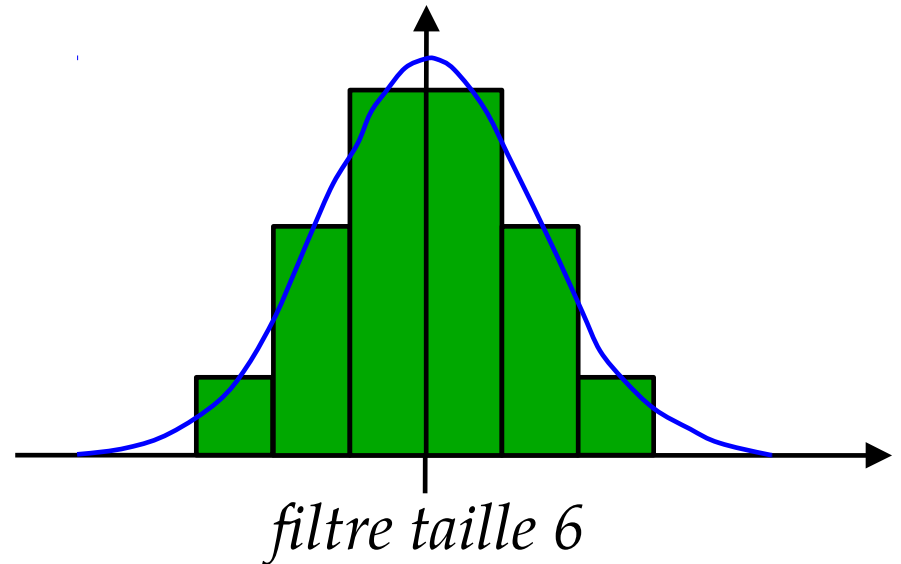
<http://setosa.io/ev/image-kernels/>

Pourquoi filtre taille impaire

- Pas de pixel « milieu » pour filtre taille paire
- Exemple : filtre radialement symétrique



La résultante est imputée à un seul pixel de l'image en sortie

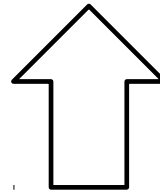


La résultante tombe à cheval entre des pixels de l'image en sortie (aliasing)

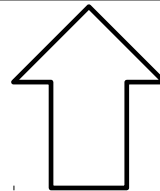
CNN : couche typique

Tenseur (volume 3D)

$(H_{\text{out}} \times W_{\text{out}} \times C_{\text{out}})$



Fonction différentiable,
avec ou sans paramètres



Tenseur (volume 3D)

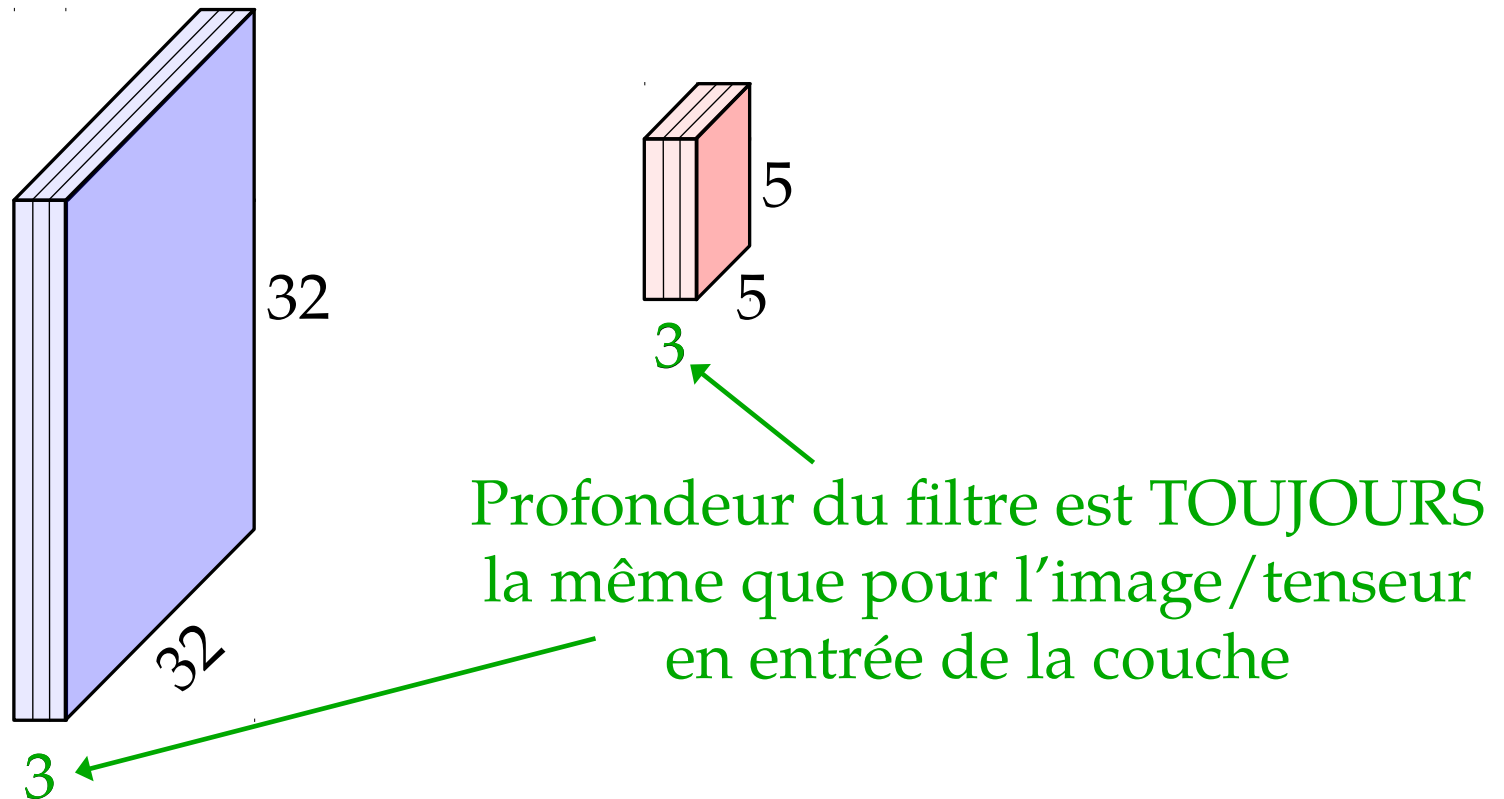
$(H_{\text{in}} \times W_{\text{in}} \times C_{\text{in}})$

Principaux types de couche

- Pleinement connecté (*Fully-Connected*)
⇒ vous connaissez déjà
- Convulsive
- Pooling
 - MaxPooling
 - Average Pooling
 - Stochastic Pooling
 - ...

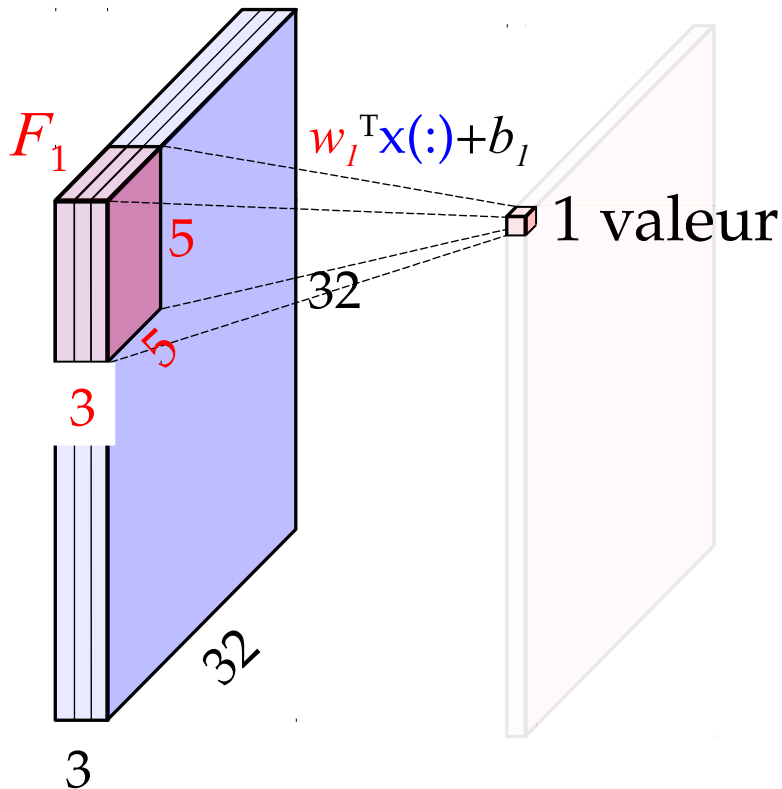
Filtres convolutifs

- Conserver la structure spatiale / couleur / *feature* de l'entrée



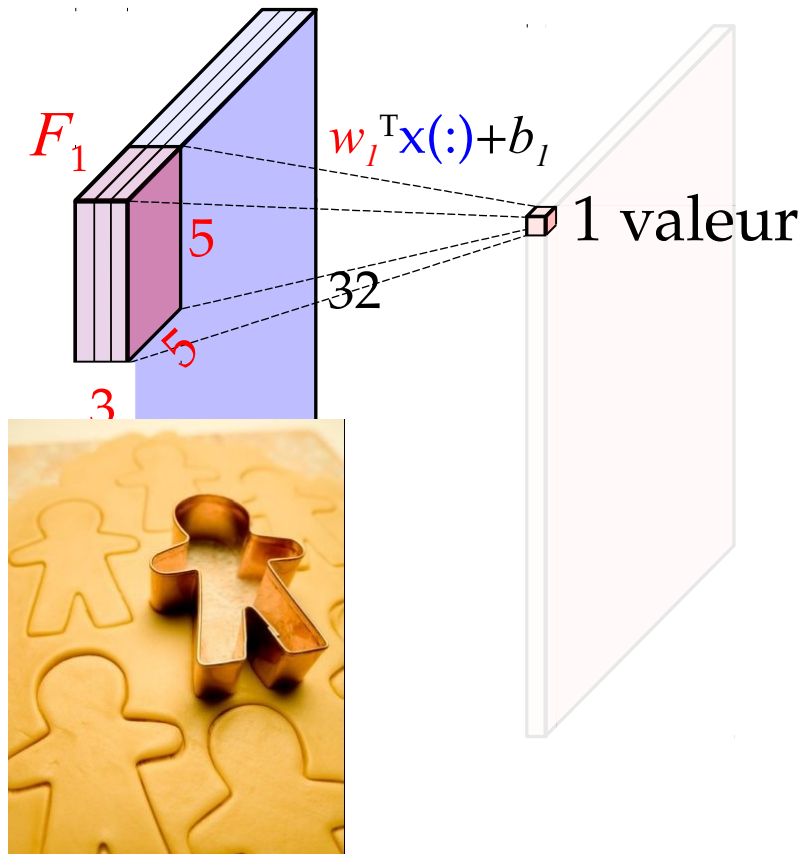
Filtres convolutifs

- Glisser spatialement le filtre F sur l'image, en calculant produit scalaire à chaque endroit de x



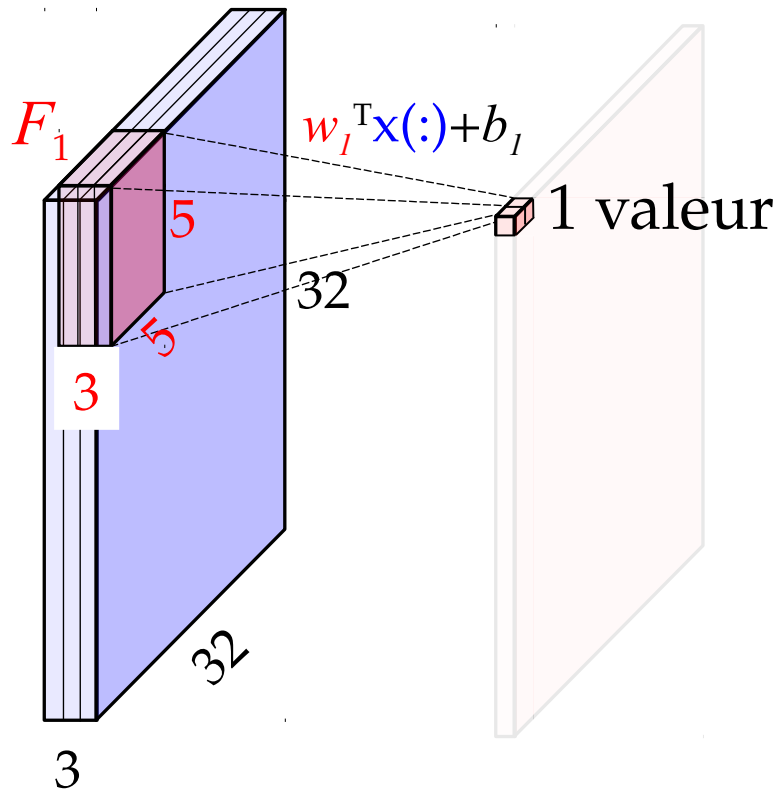
Filtres convolutifs

- Glisser spatialement le filtre F sur l'image, en calculant produit scalaire à chaque endroit de x



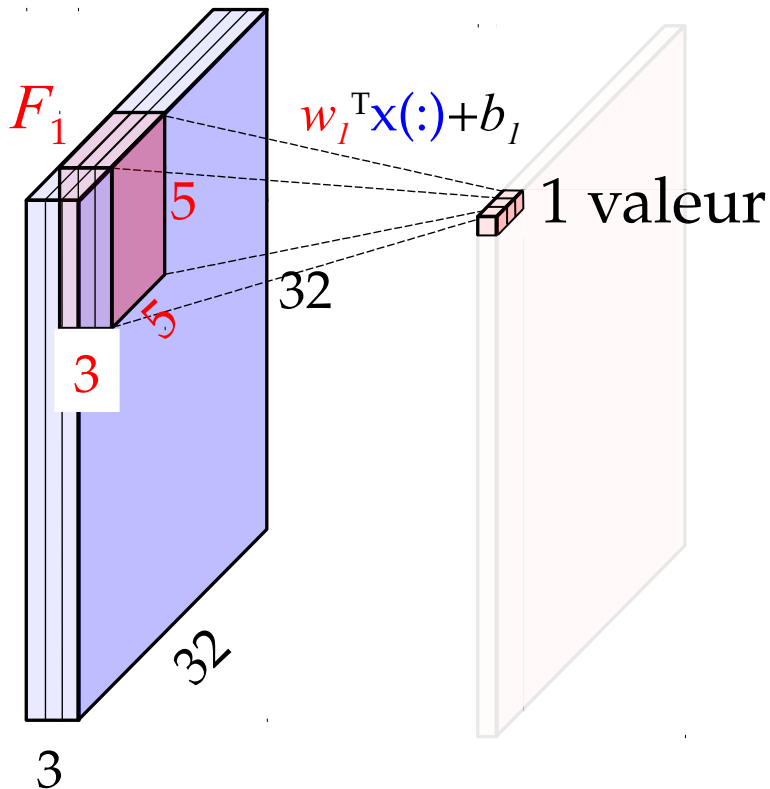
Filtres convolutifs

- Glisser spatialement le filtre F sur l'image, en calculant produit scalaire à chaque endroit de x



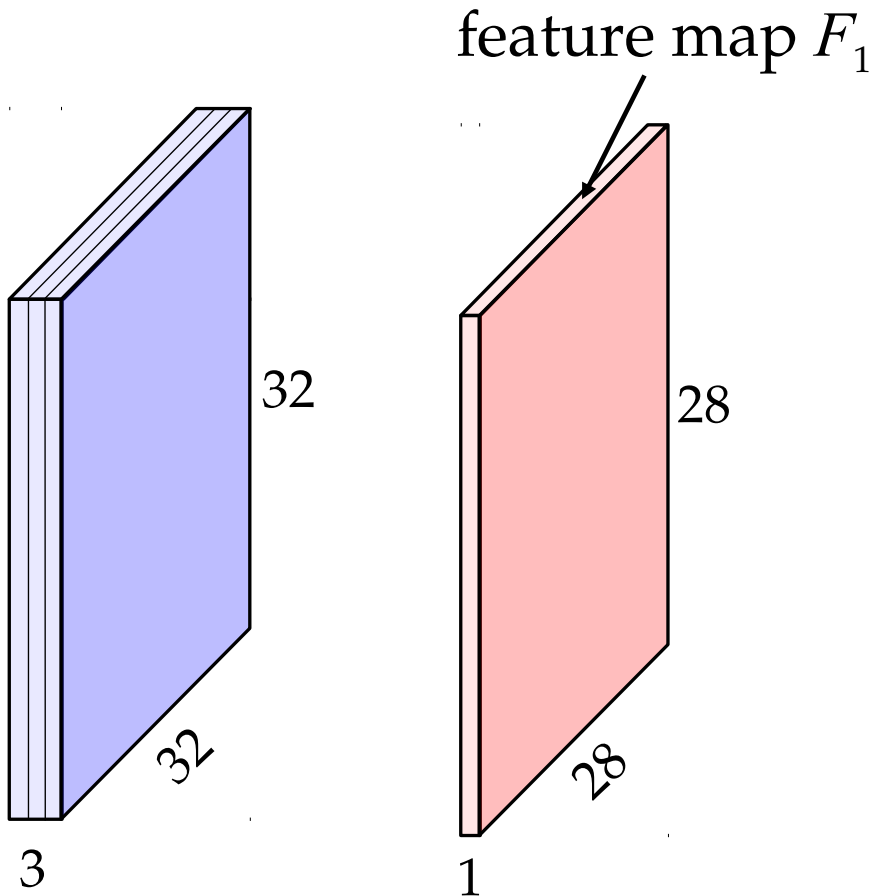
Filtres convolutifs

- Glisser spatialement le filtre F sur l'image, en calculant produit scalaire à chaque endroit de x



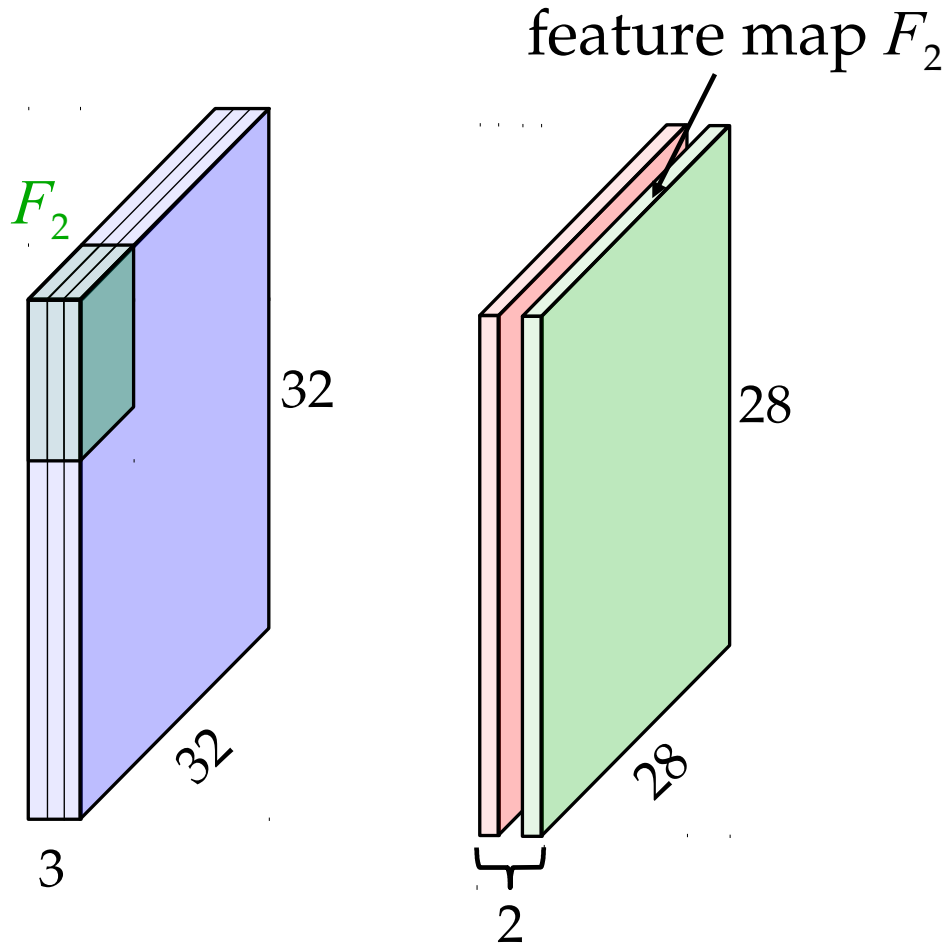
Filtres convolutifs

- Sortie : *feature map*

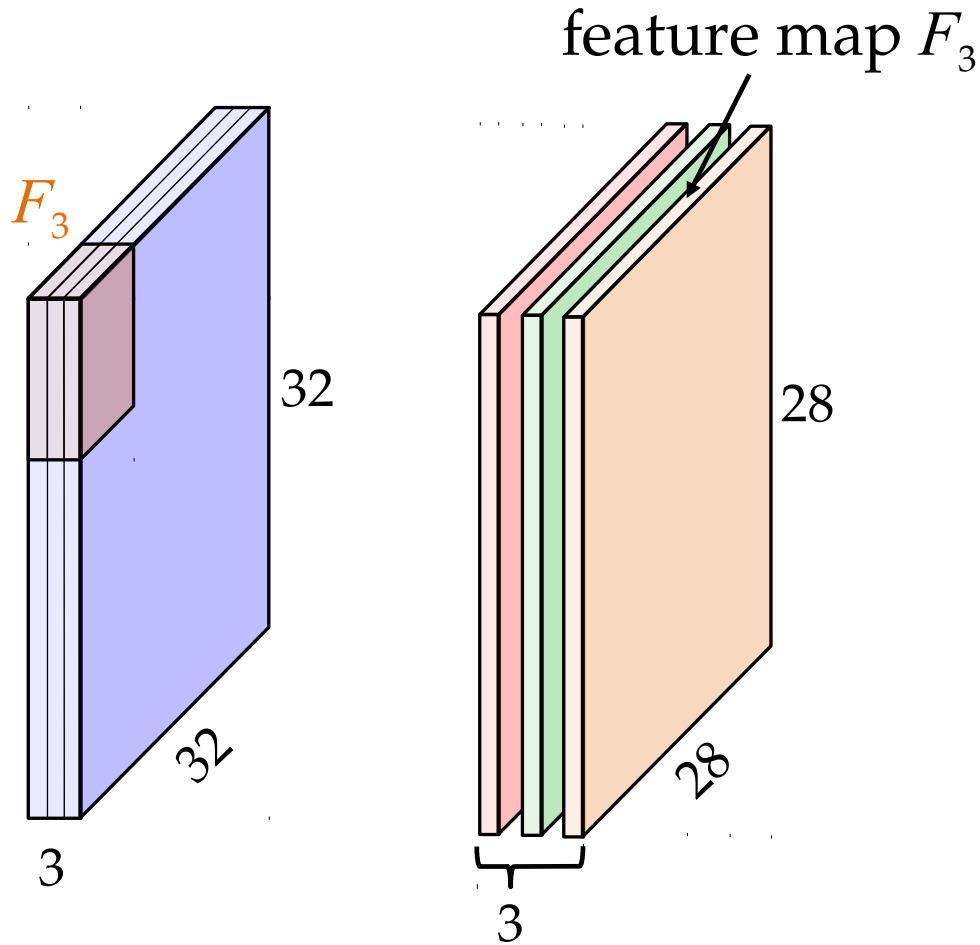


« combien présent est le feature F_1 à cet endroit? »

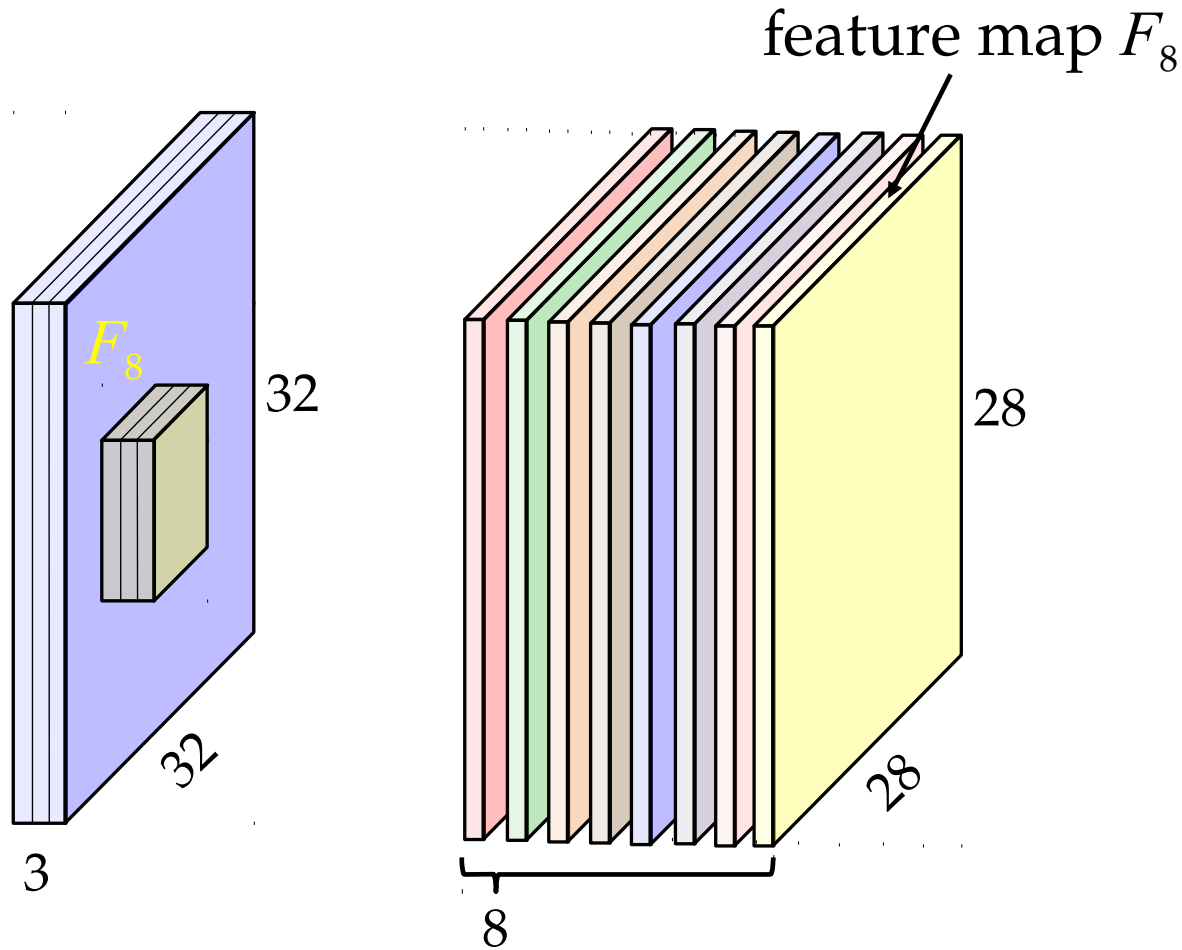
Banque de filtres convolutifs



Banque de filtres convolutifs

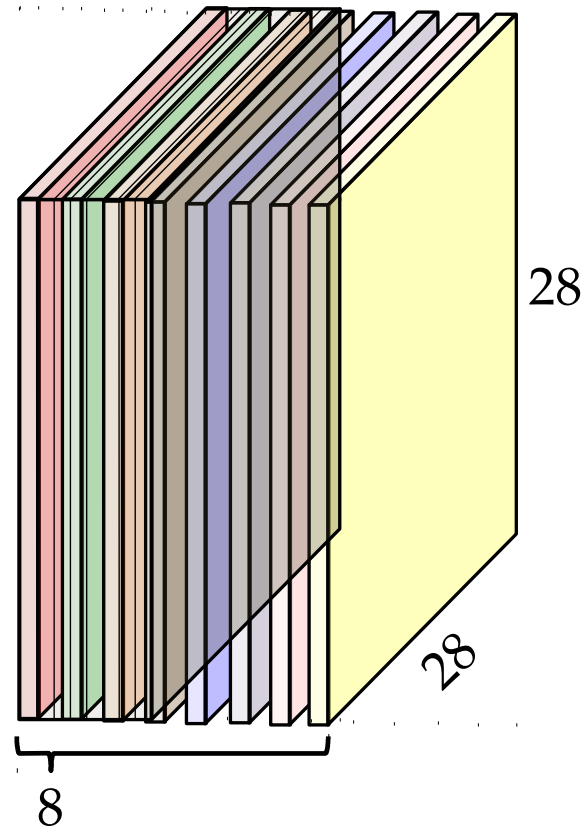


Banque de filtres convolutifs



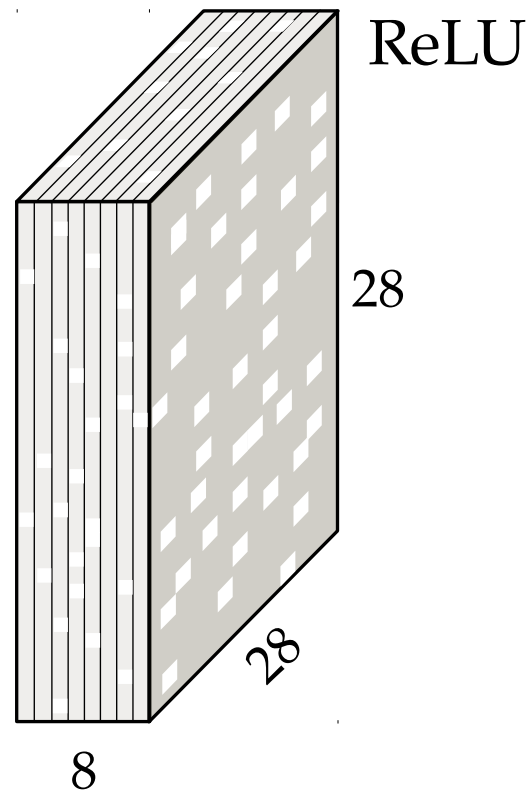
Résultante

Tenseur ordre 3

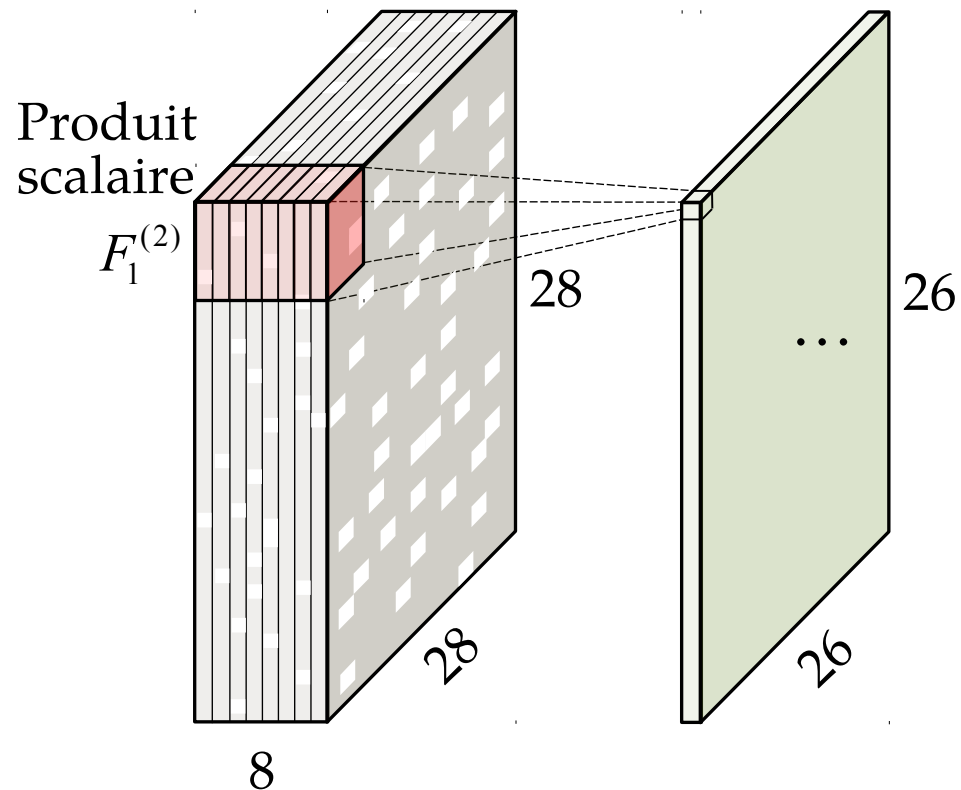
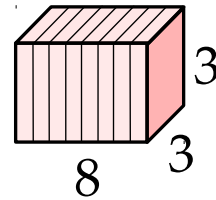


Résultante

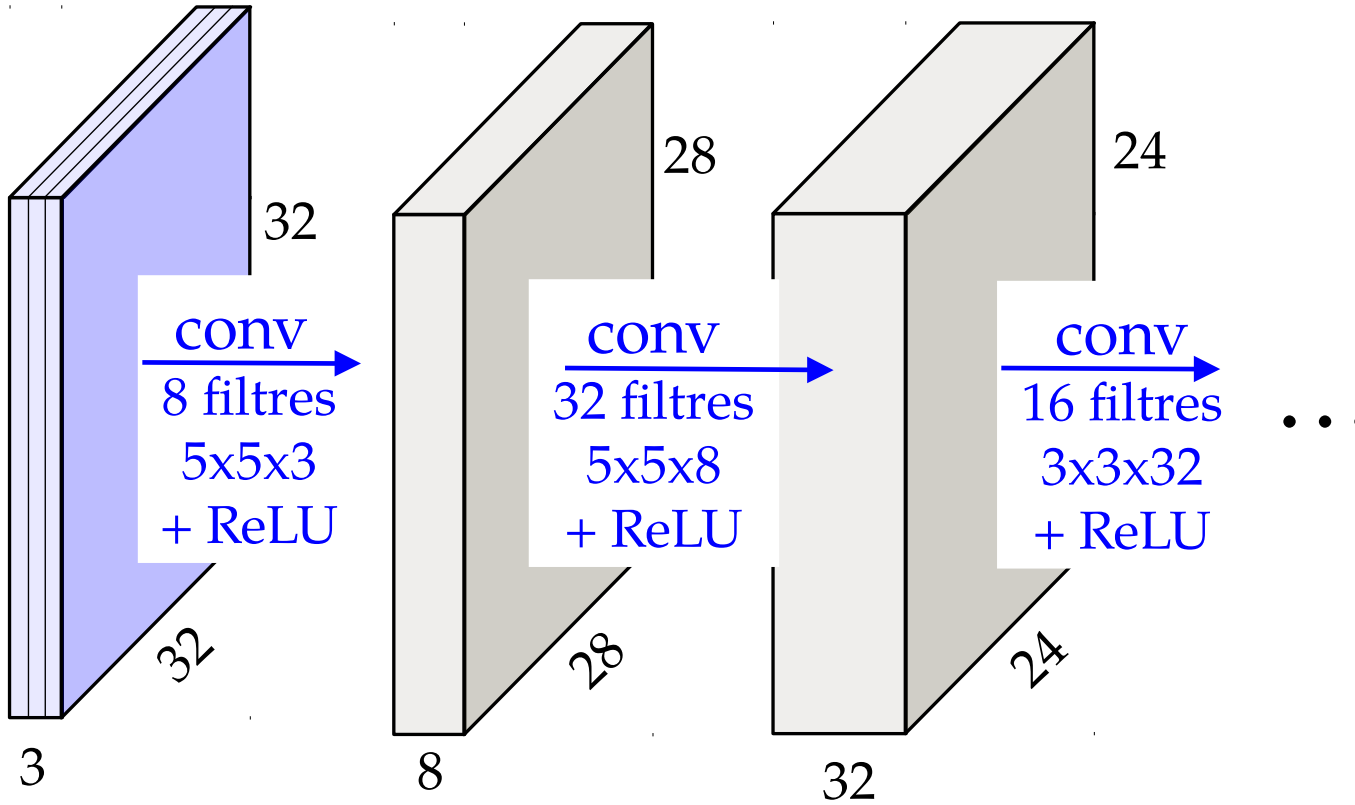
- Applique la non-linéarité sur chaque entrée, individuellement
- Appelé *feature activation map*



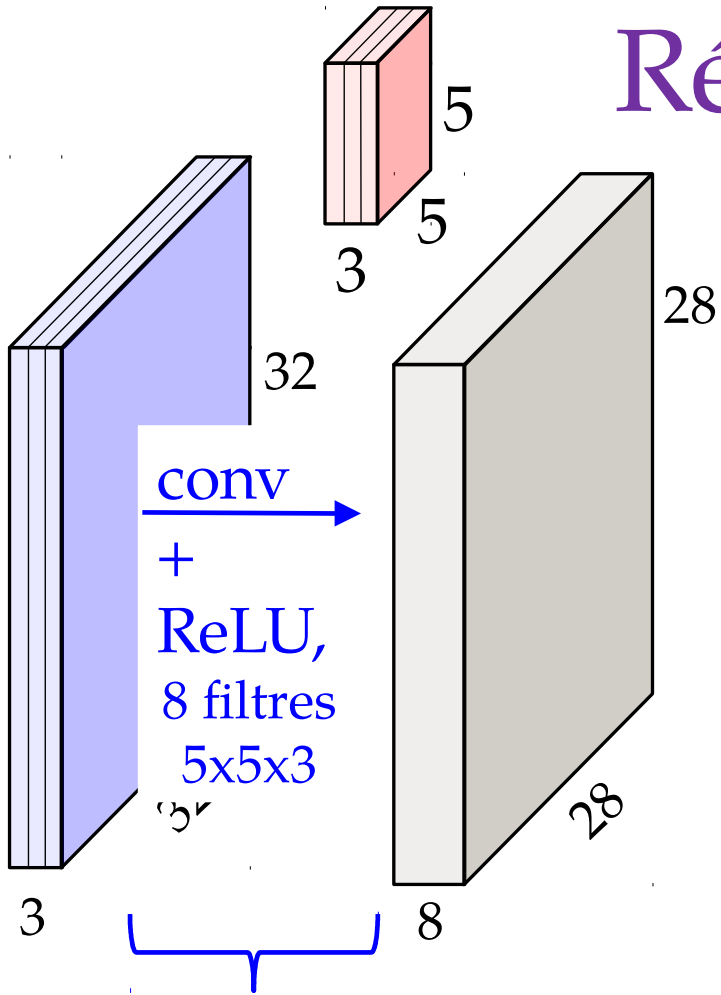
On recommence, avec une banque de filtres différents



Idée de base



Réduction # paramètres



Combien de paramètres pour cette couche?

$$(5 \times 5 \times 3 + 1) \times 8 = 608$$

biais filtres

- Perte de capacité d'expression vs. couche pleinement connectée
 - on ne peut pas établir de lien entre les pixels très éloignés
- Réduction importante du nombre de paramètre

Combien de paramètres pour une couche pleinement connectée ?

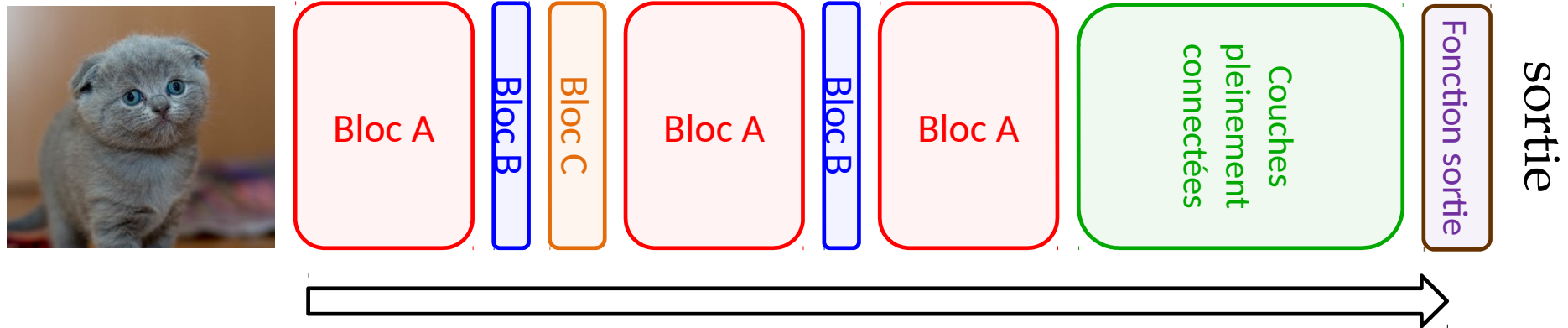
$$28 \times 28 \times 8 = 6272 \text{ neurones}$$

$$32 \times 32 \times 3 + 1 : 3073 \text{ param./neurones}$$

$$\text{Total : } 19,273,856 \text{ paramètres}$$

Approche par bloc

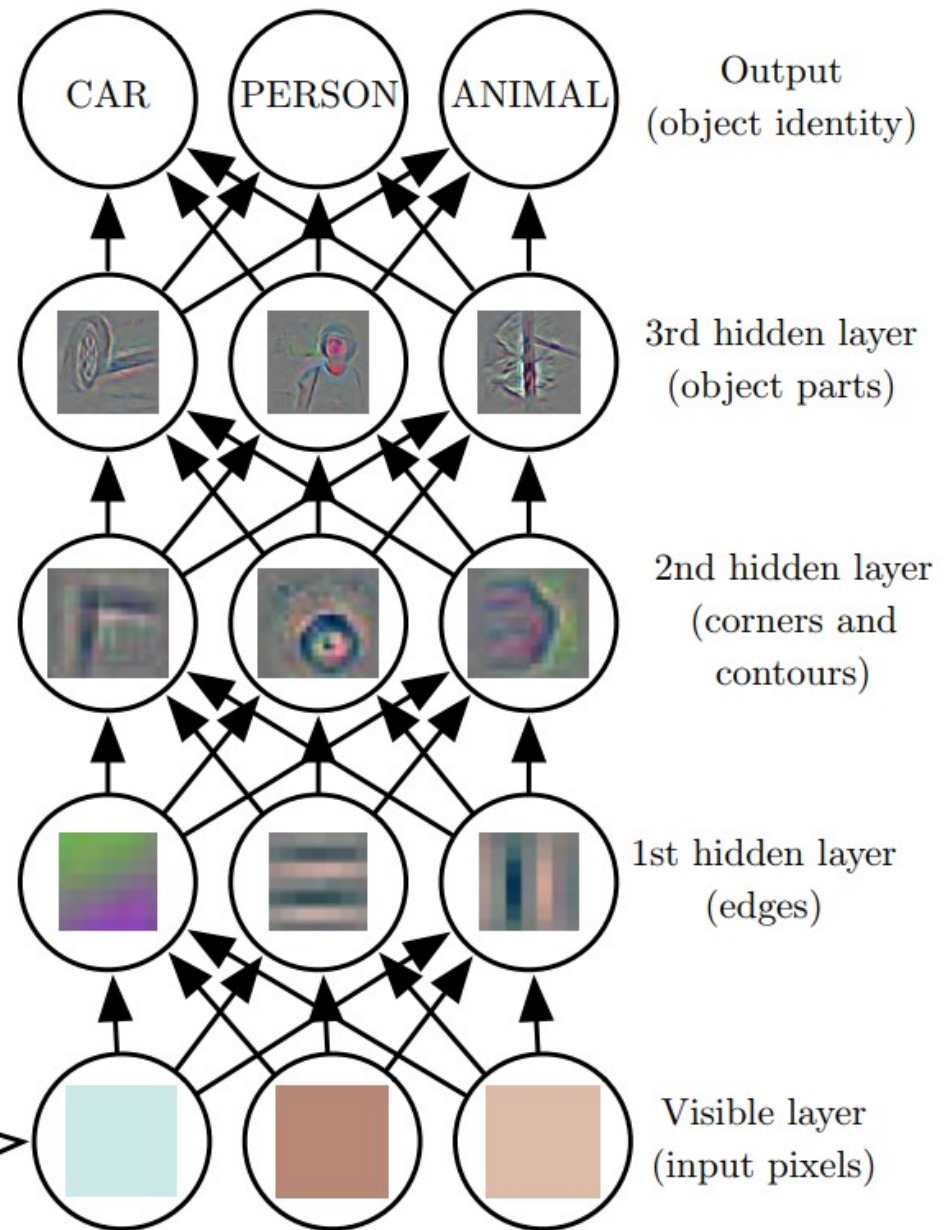
- La plupart des architectures sont organisées par alternance de blocs



- Facile de combiner des blocs de différents types, jouer sur la profondeur, réutiliser des idées
-

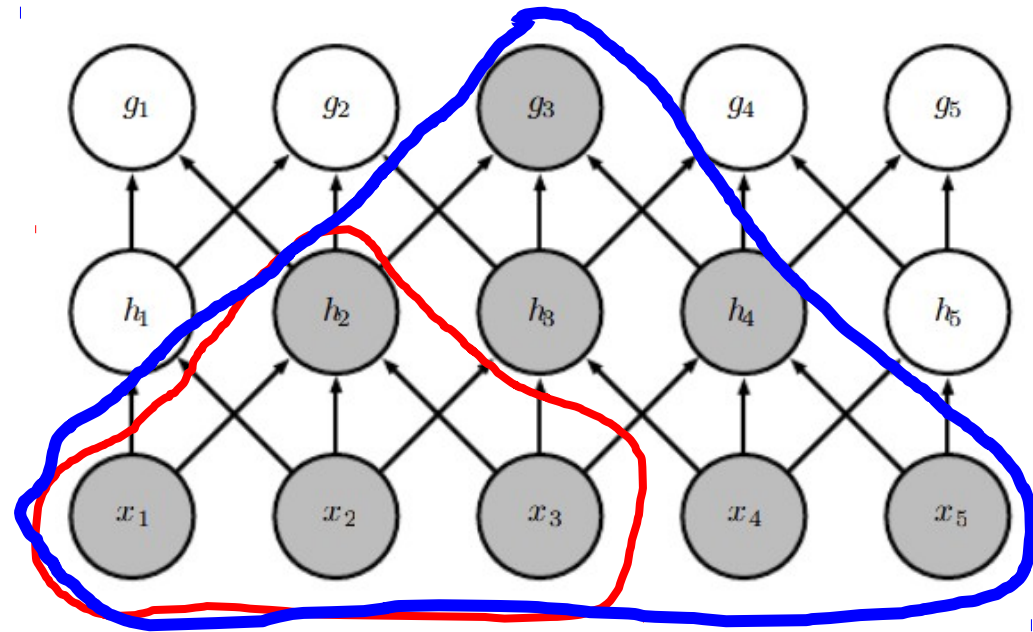
Hiérarchie de filtres

Va établir des liens entre des pixels de plus en plus éloignés

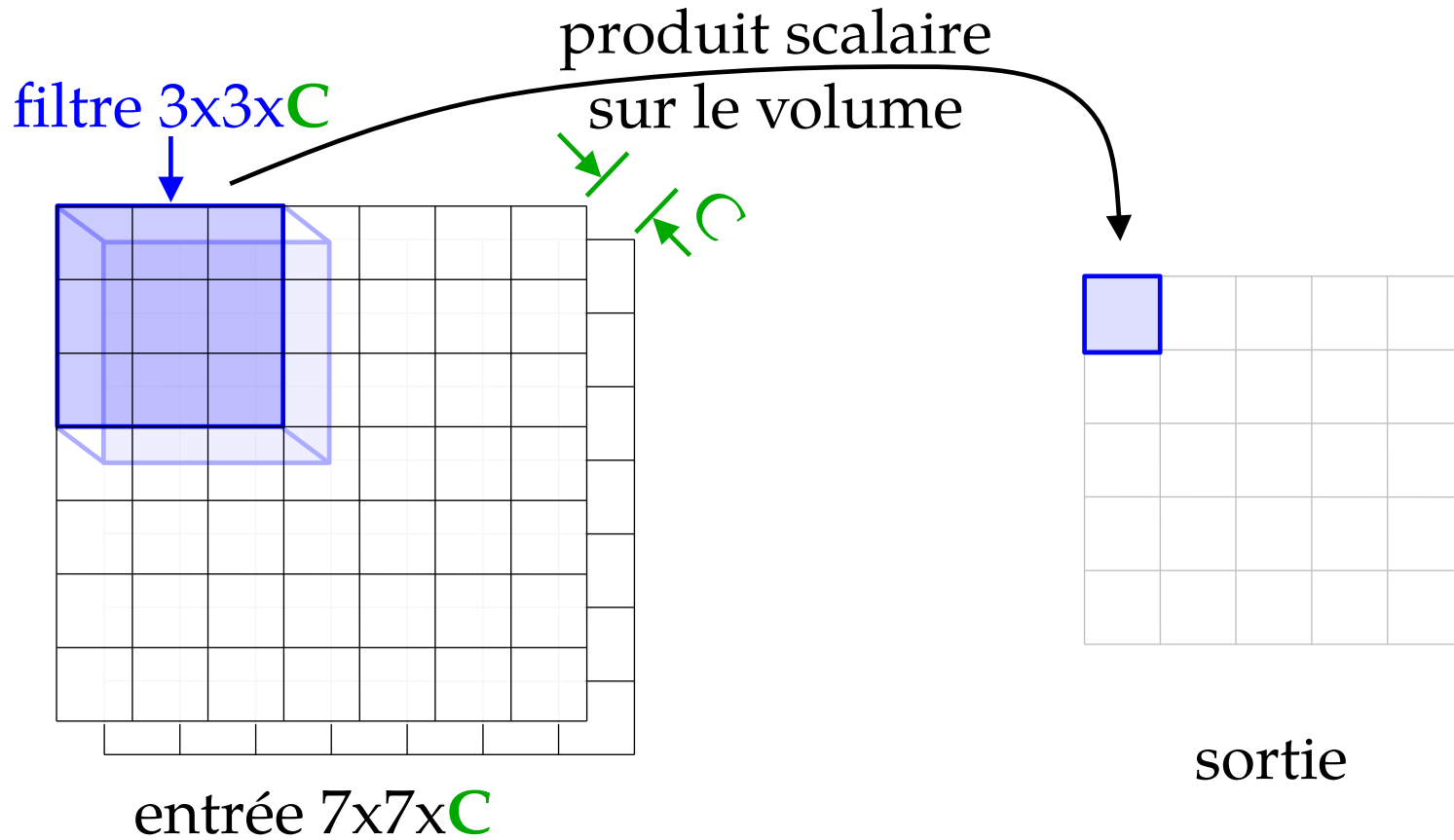


Champ récepteur

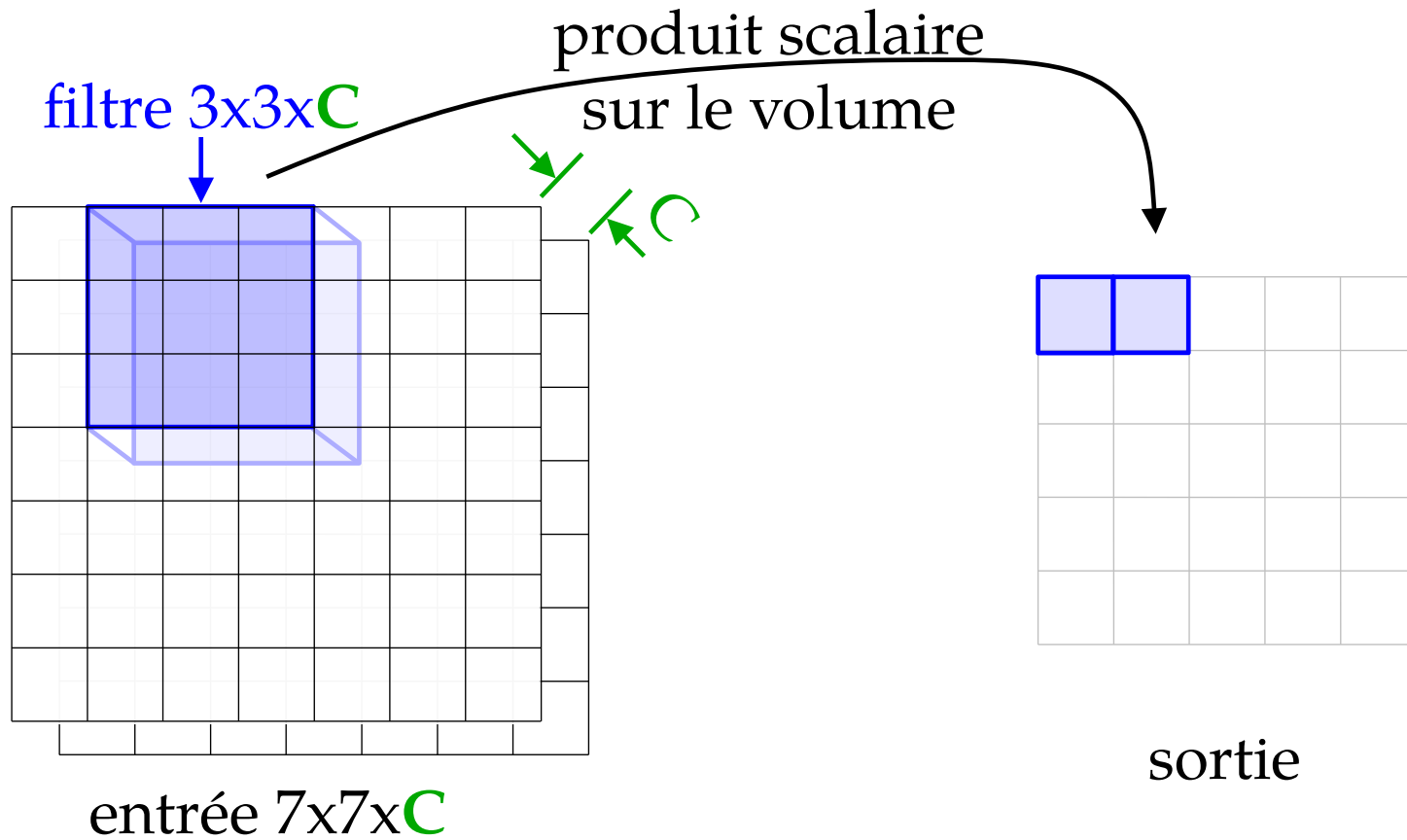
- Augmente avec la profondeur
- Les neurones plus haut peuvent prendre des décisions basées sur des champs plus grands



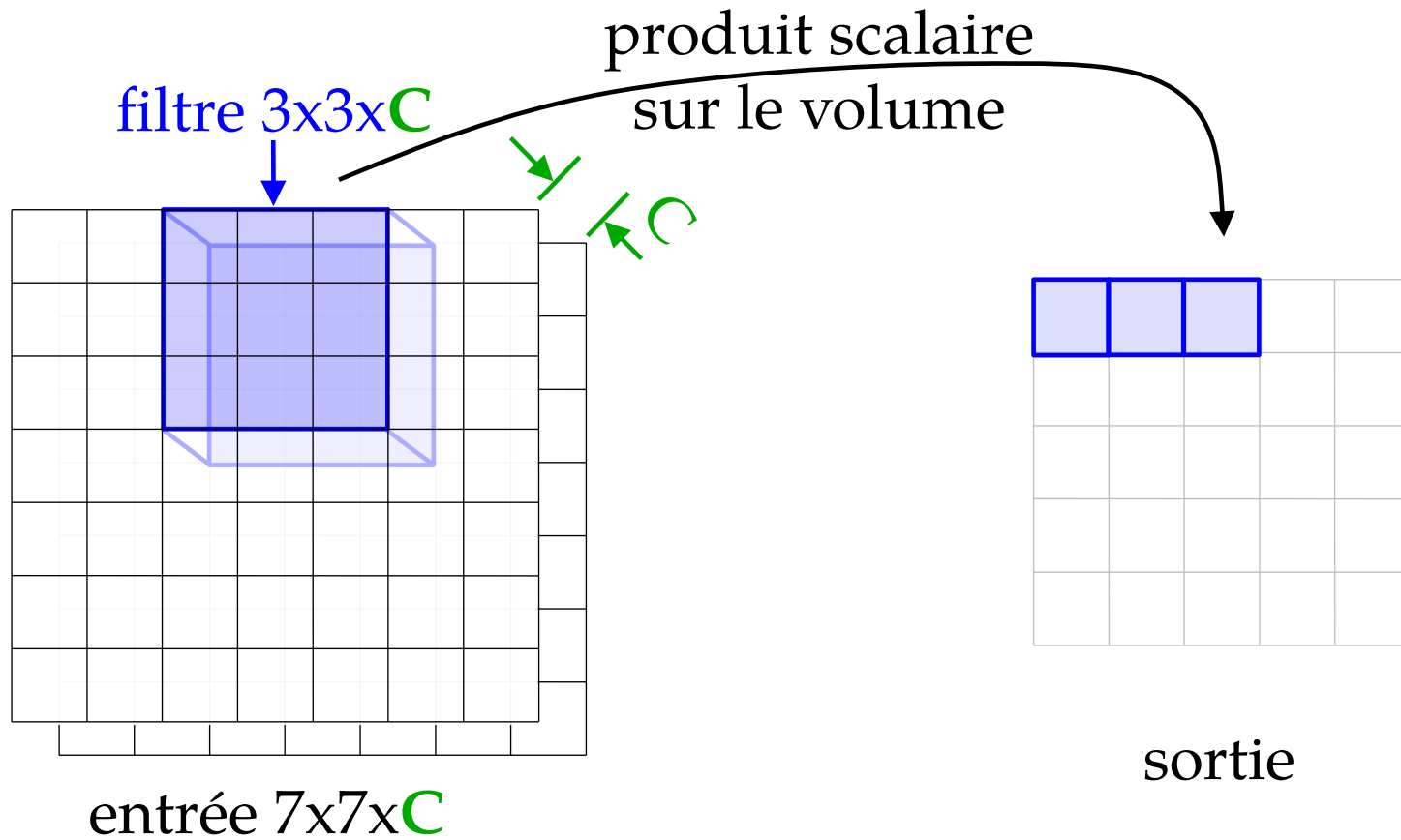
Effet de bord : redimensionnalisation



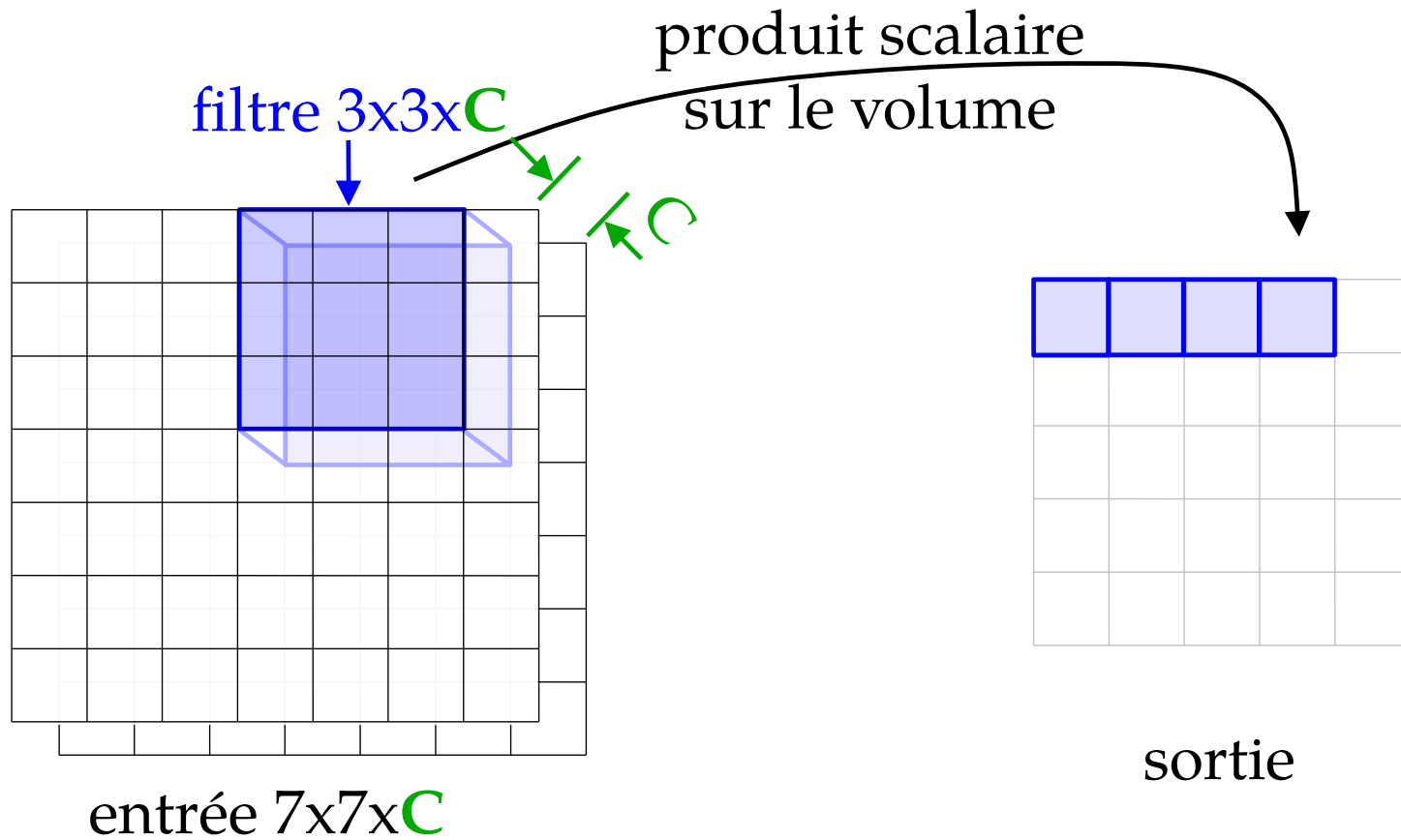
Effet de bord : redimensionnalisation



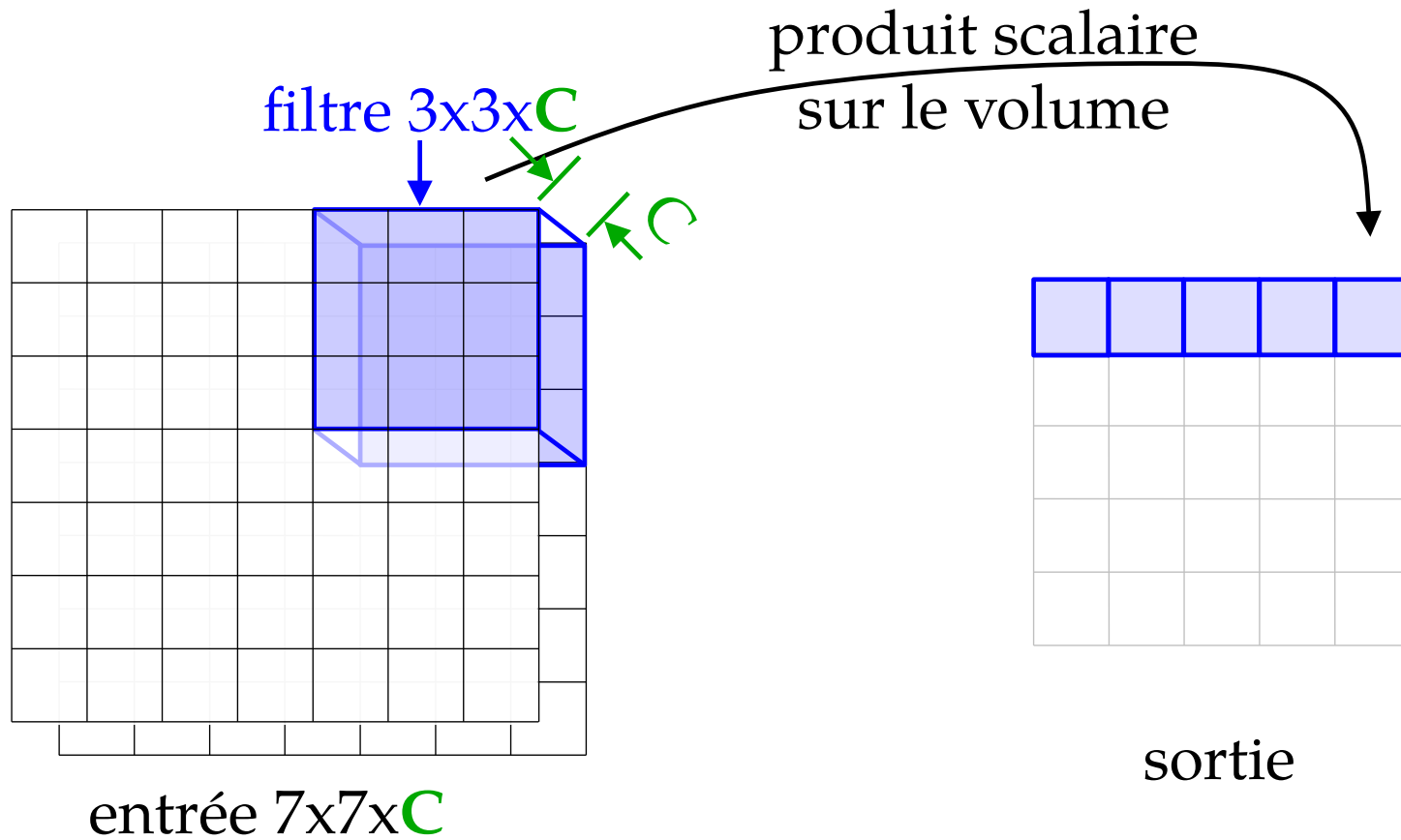
Effet de bord : redimensionnalisation



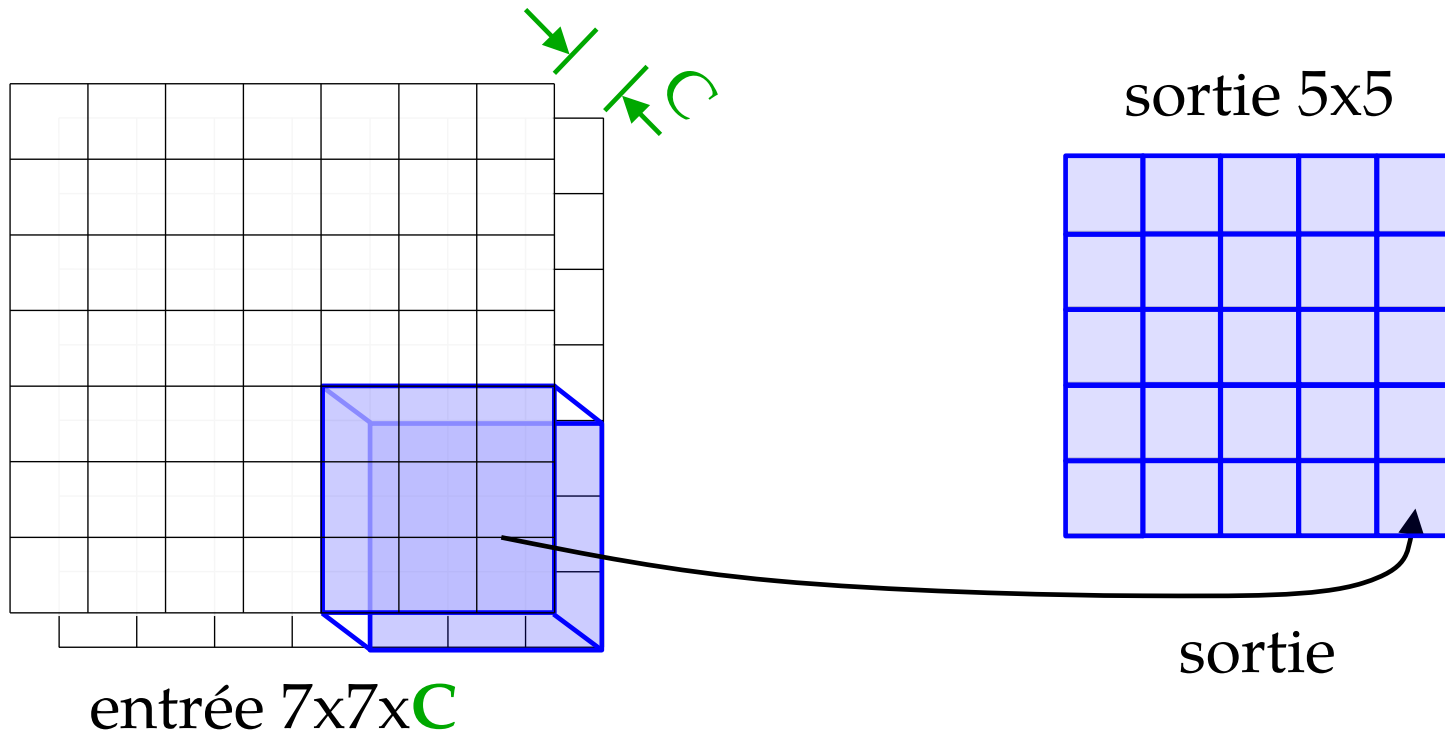
Effet de bord : redimensionnalisation



Effet de bord : redimensionnalisation

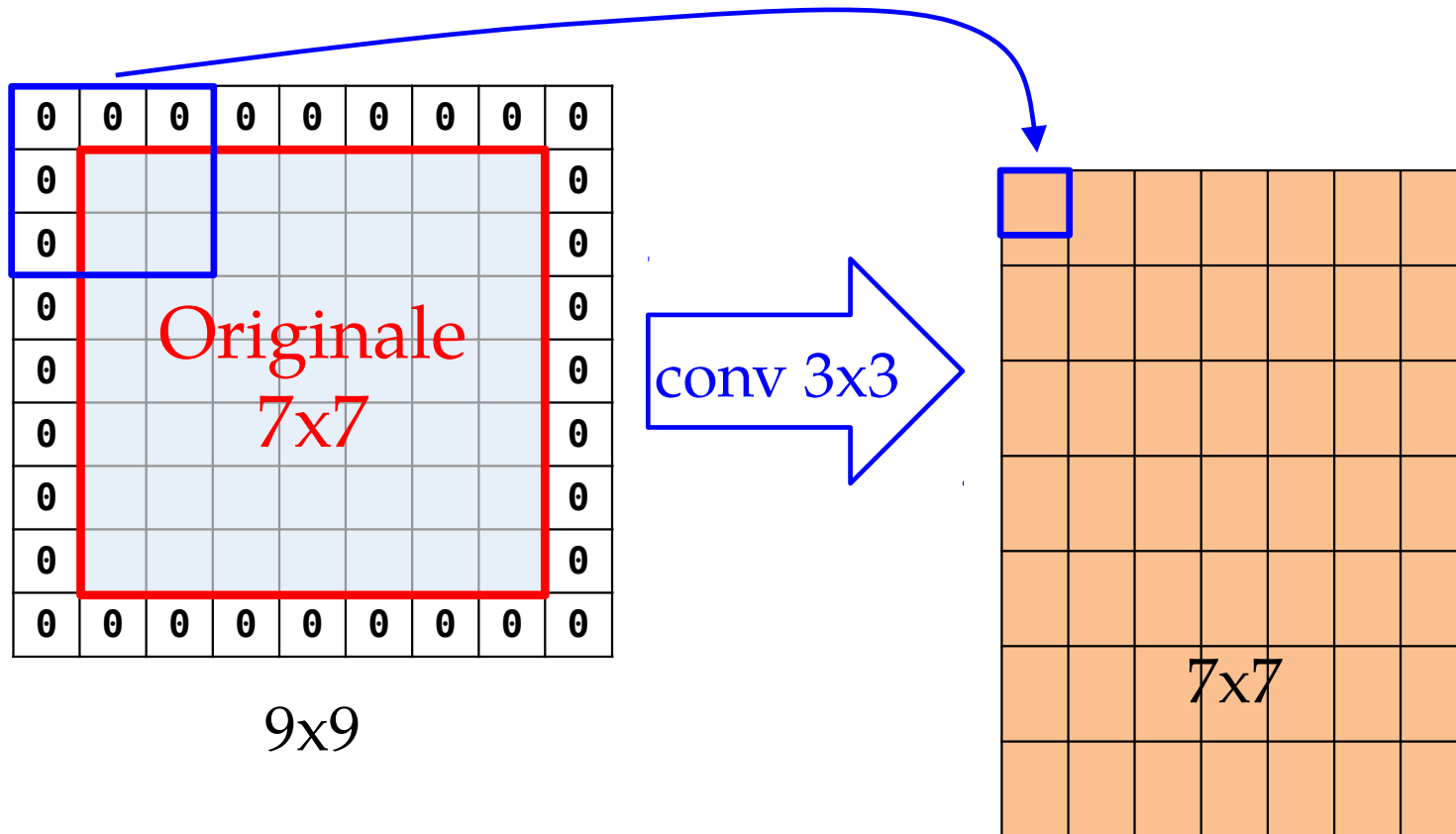


Effet de bord : redimensionnalisation



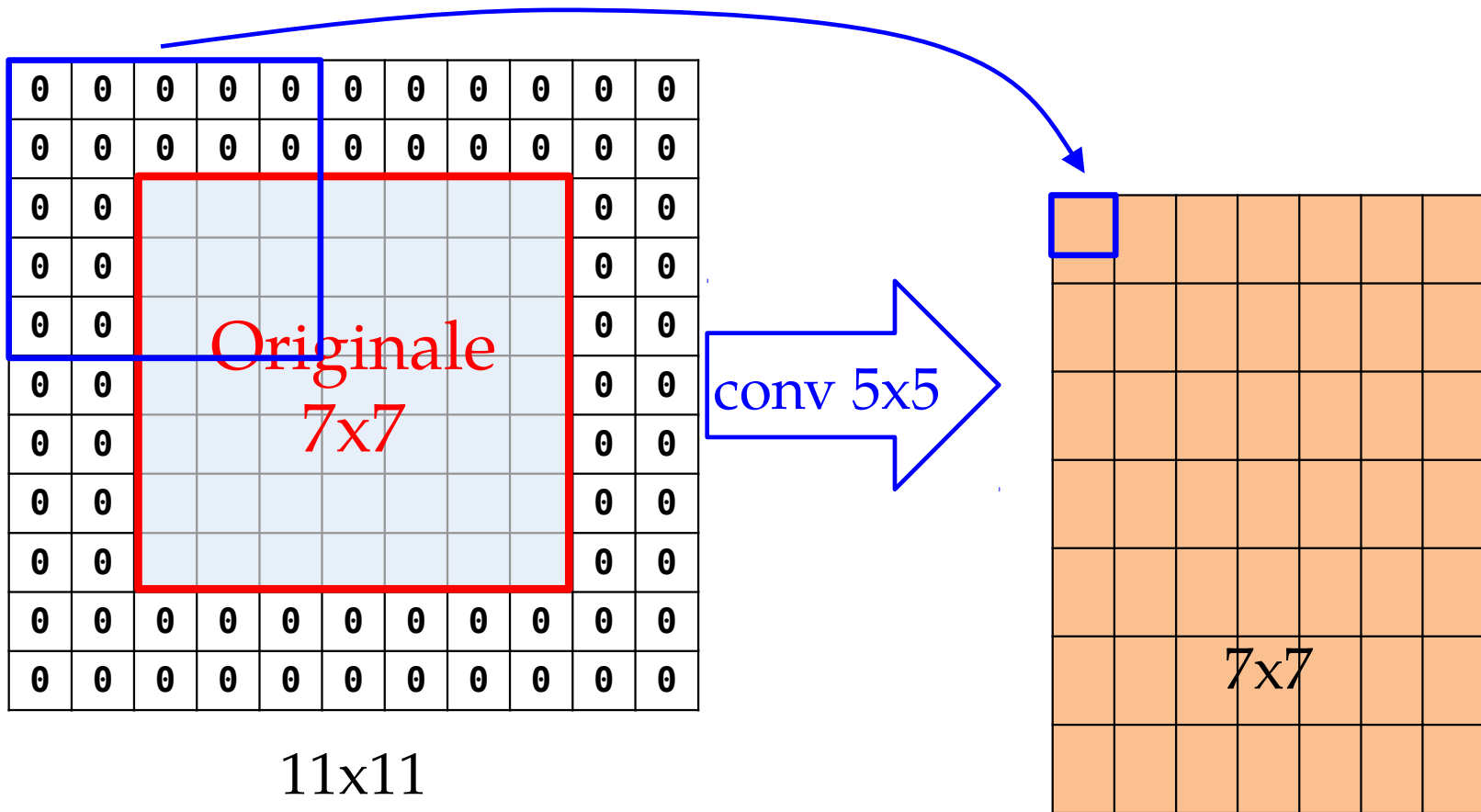
Ajout de zéros : *zero padding*

- Ajoute des entrées à 0 en bordure



(par simplicité, j'ai retiré la 3^{ème} dimension)

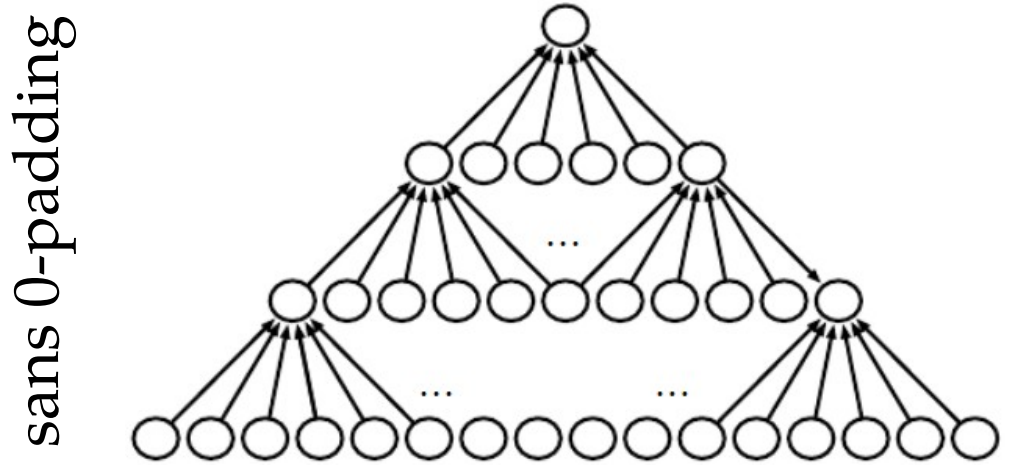
Ajout de zéros : *zero padding*



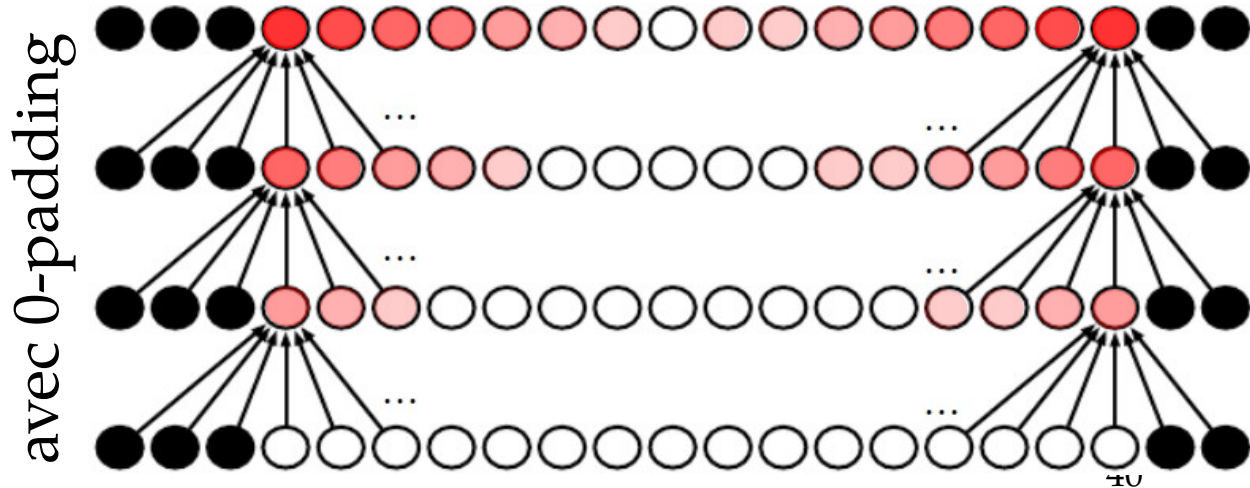
- Pour filtre 3x3, padding de largeur 1
- Pour filtre 5x5, padding de largeur 2
- Pour filtre 7x7, padding de largeur 3

Zero-padding

- Conserve la largeur du pipeline
- Attention aux effets de bord, où les entrées sont moins informatives à cause du 0-padding



rougeur : impact du 0-padding



(ne semble pas trop problématique)

Quelques questions!

- Si j'ai une image de $224 \times 224 \times 3$ en entrée, complétez la taille du filtre : $5 \times 5 \times ?$

Réponse : $5 \times 5 \times 3$

- Si j'ai 32 filtres 5×5 appliqués sur cette image $224 \times 224 \times 3$, sans 0-padding, quelle sont les dimensions du tenseur de sortie?

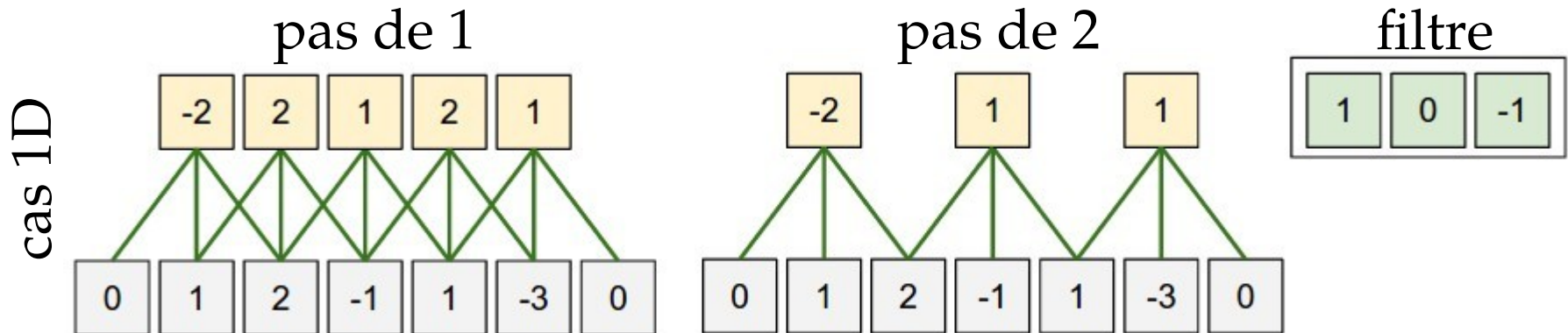
Réponse : $220 \times 220 \times 32$

- Quelle est la largeur du 0-padding pour un tenseur de $64 \times 64 \times 3$, si on applique un filtre 9×9 ?

Réponse : largeur de 4

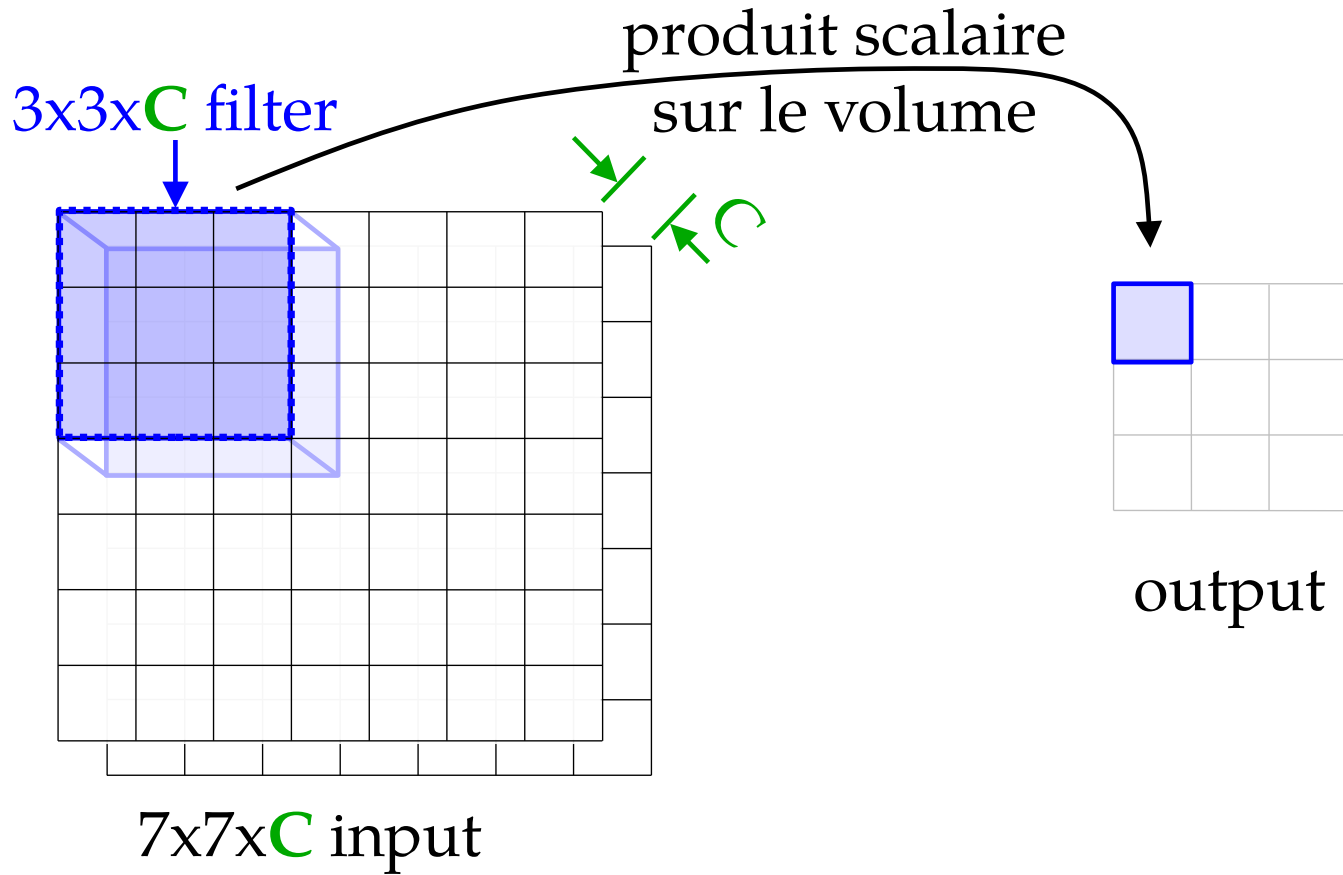
Pas (*stride*)

- **Pas** (*stride*) : saut dans l'application de la convolution

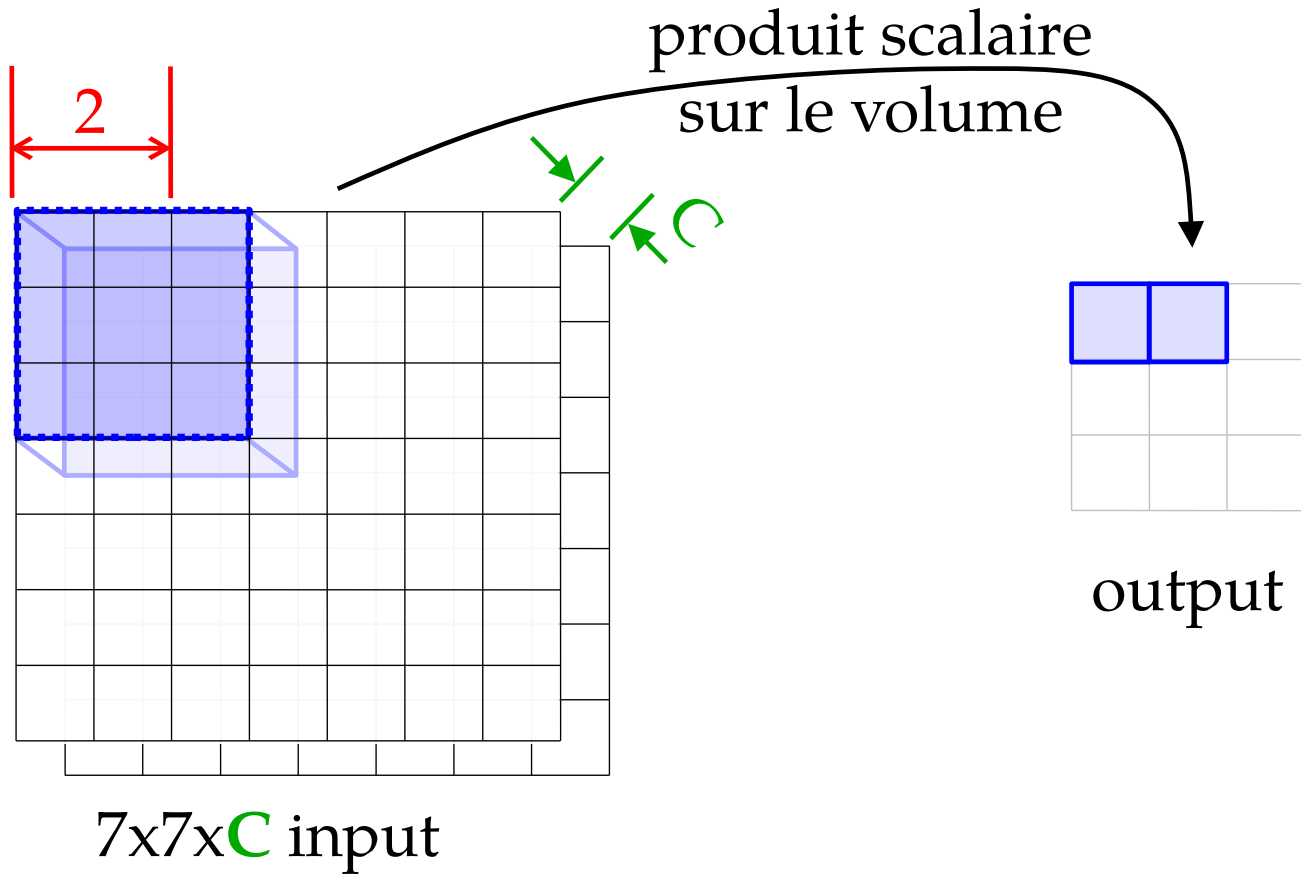


- Le ***pas*** est rarement plus de 3
 - Si plus d'1, réduit grandement la taille de sortie HxW
- Pas toujours possible d'avoir un nombre entier d'application de convolution, si le ***pas*** n'est pas 1
 - Par exemple, entrée 7x7, filtre 3x3, ***pas*** de 3
 - Libraire peut automatiquement faire du 0-padding, couper l'image (*crop*) ou lancer une exception

Pas de 2

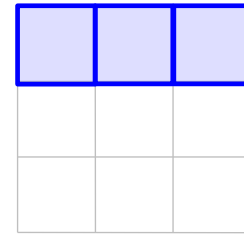
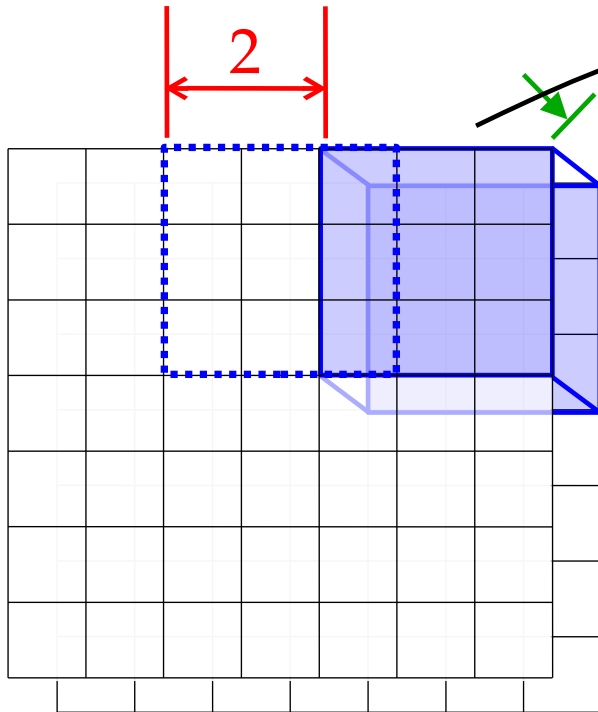


Pas de 2



Pas de 2

produit scalaire
sur le volume



sortie 3x3

7x7x3 input

Tailles entrée et sortie

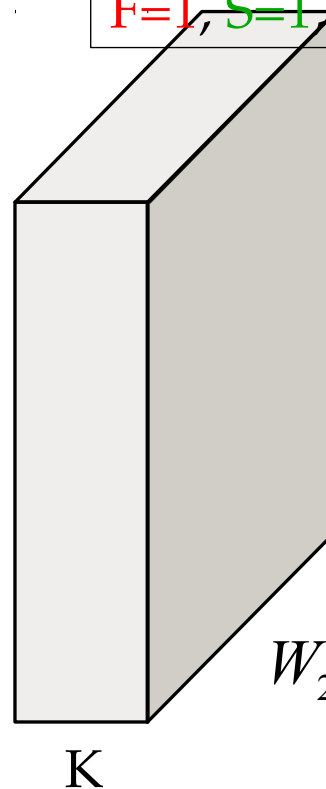
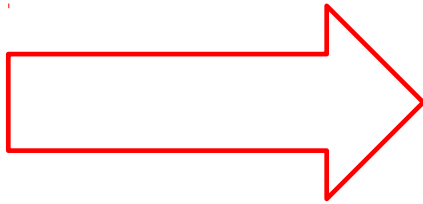
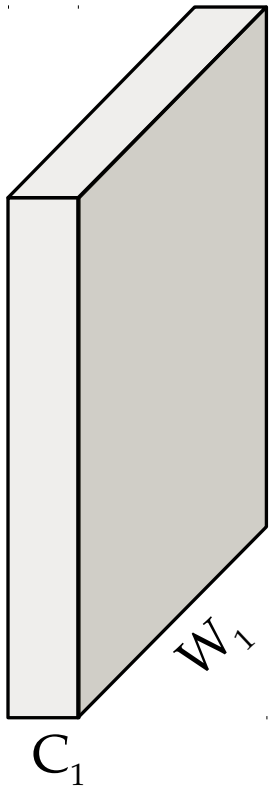
K : nombre de filtres

F : taille du filtre (**F**x**F**)

S : pas (stride)

P : quantité de 0-padding

Valeurs typiques (cs231n)
K : puissances de 2
F=3, S=1, P=1
F=5, S=1, P=2
F=5, S=2, P=autant que nécessaire
F=1, S=1, P=0



$$H_2 = \frac{H_1 - F + 2P}{S} + 1$$

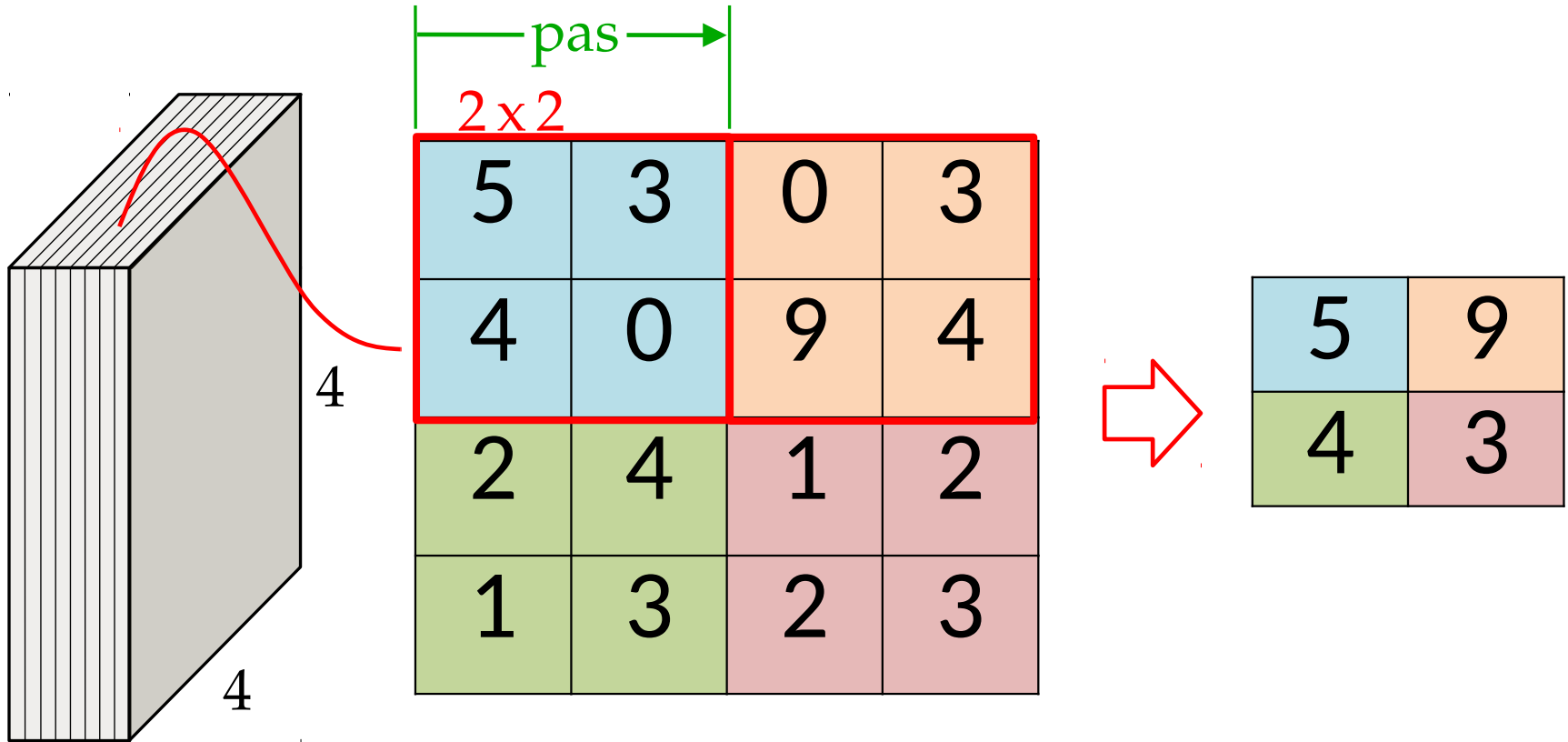
$$W_2 = \frac{W_1 - F + 2P}{S} + 1$$

paramètres : **F**•**F**•**C**₁•**K** poids + **K** biais

Pooling

Max Pooling

- Appliqué pour chaque tranche, indépendamment



On doit spécifier le «pas»

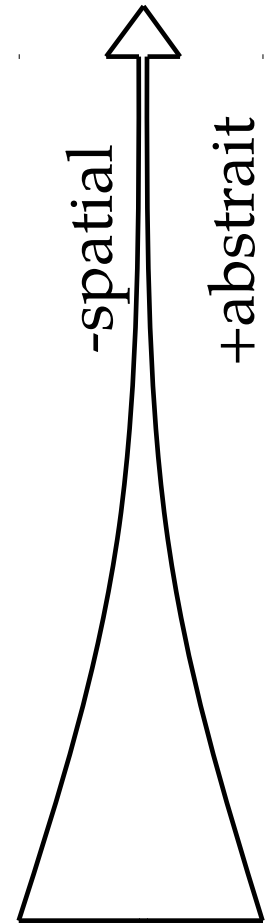
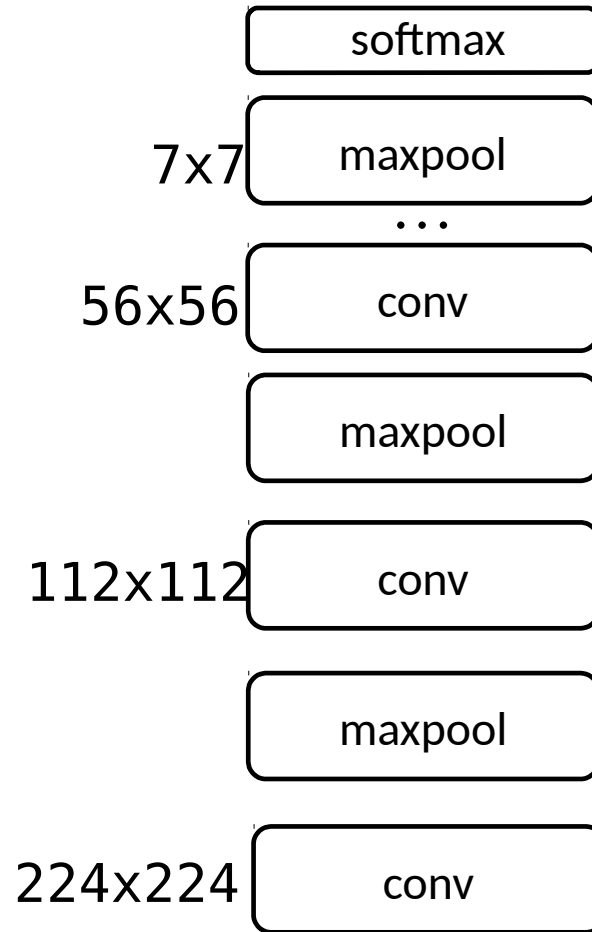
Max Pooling

- Réduit la dimension du feature map
- Souvent, on en profite pour augmenter le nombre de filtre
 - la « quantité » d'information reste similaire
 - augmente la richesse de la représentation / abstraction
- Pourquoi maxpool au lieu de faire une moyenne?
 - maxpool: détecte la présence d'un feature dans une région
 - avgpool: va en partie noyer cette valeur (ou compter le nombre)

Pooling

- Augmente champ réceptif rapidement
- Réduit le nombre de paramètre
- Confère une certaine invariance aux transformations géométriques (rotation, translation, échelle)

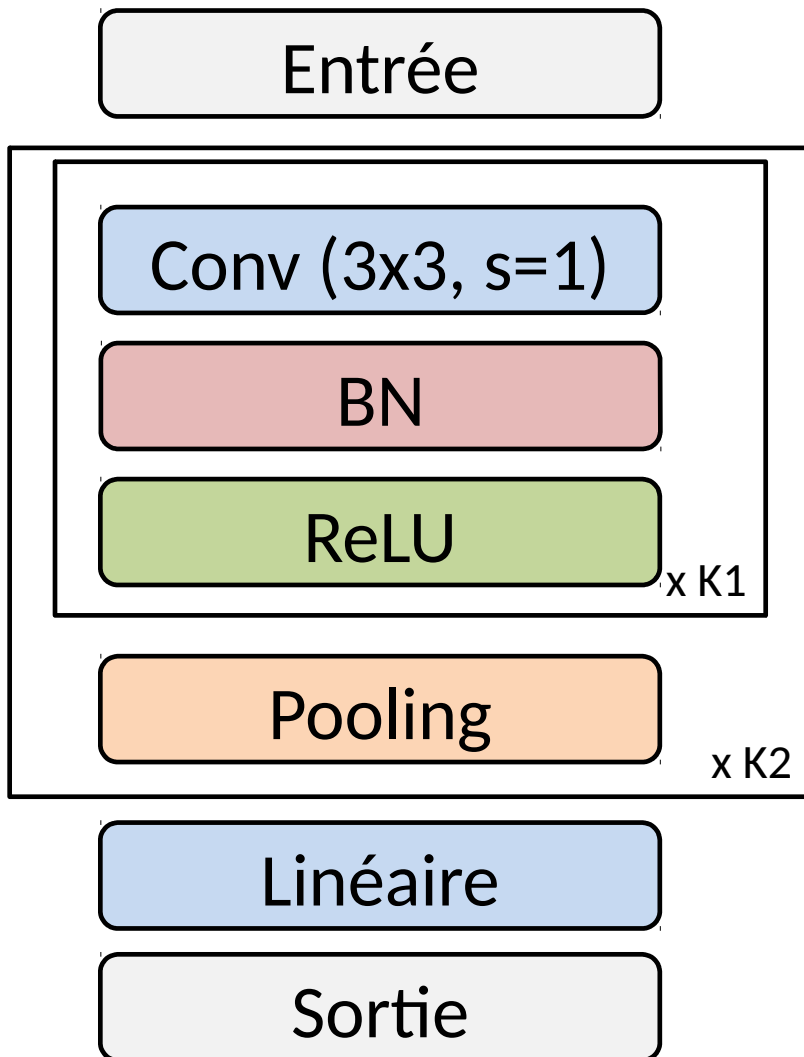
Classe: aucune spatialité



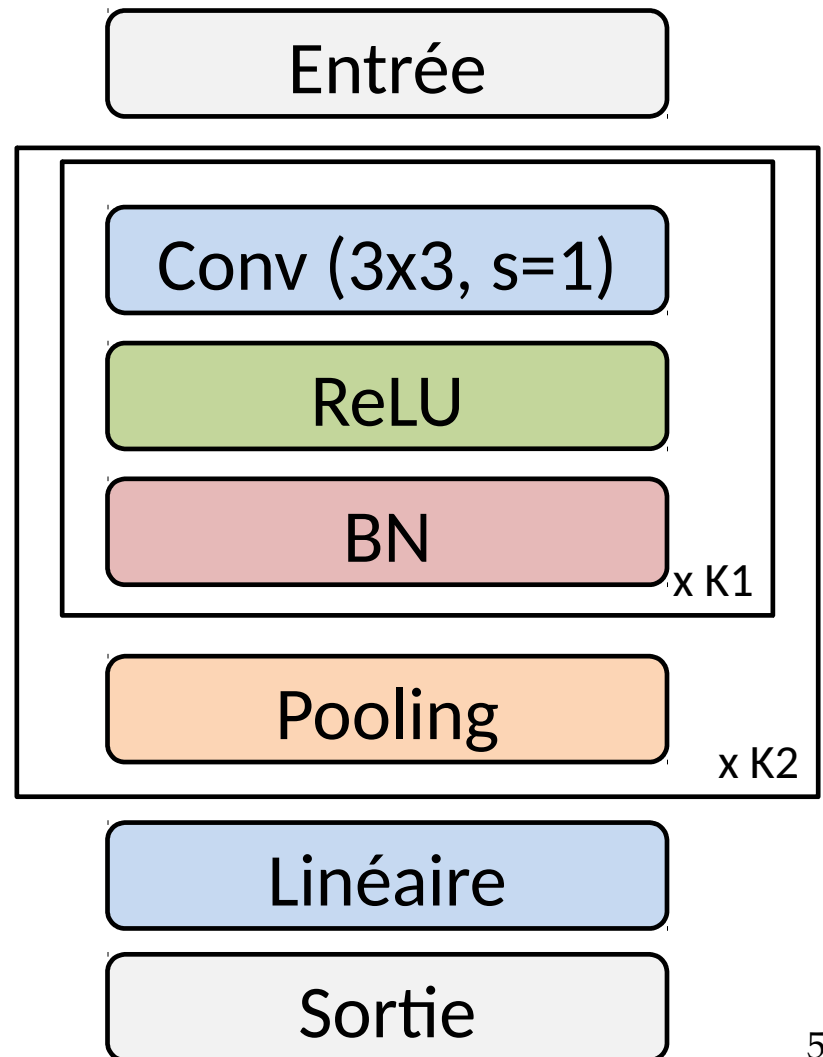
100% spatial

Position de la batch norm

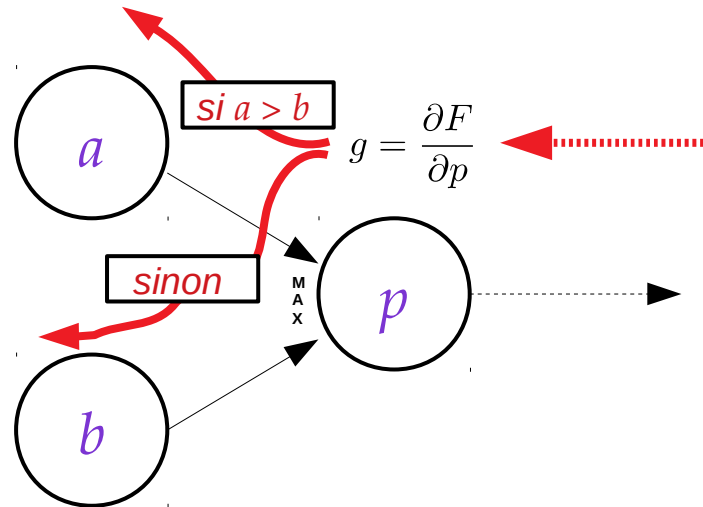
Classique



Nouvelle tendance



Rétropropagation du gradient



$$p = \max(a, b) = \begin{cases} a & \text{si } a > b \\ b & \text{sinon.} \end{cases}$$

$$\frac{\partial p}{\partial a} = \max(a, b) = \begin{cases} 1 & \text{si } a > b \\ 0 & \text{sinon.} \end{cases}$$

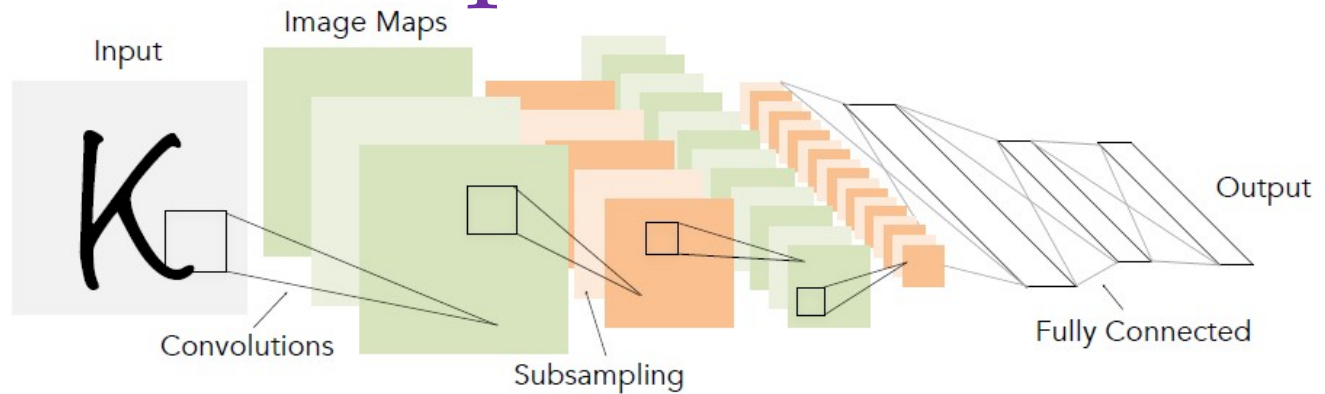
$$\frac{\partial p}{\partial b} = \max(a, b) = \begin{cases} 0 & \text{si } a > b \\ 1 & \text{sinon} \end{cases}$$

Architectures

Ne date pas d'hier

1998

LeCun et al.



of transistors



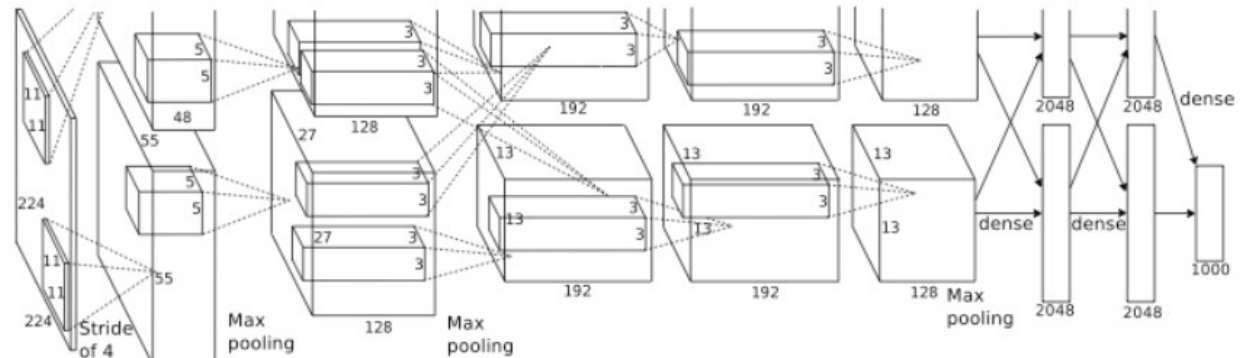
10^6

of pixels used in training

10^7 **NIST**

2012

Krizhevsky et al.



of transistors GPUs



10^9



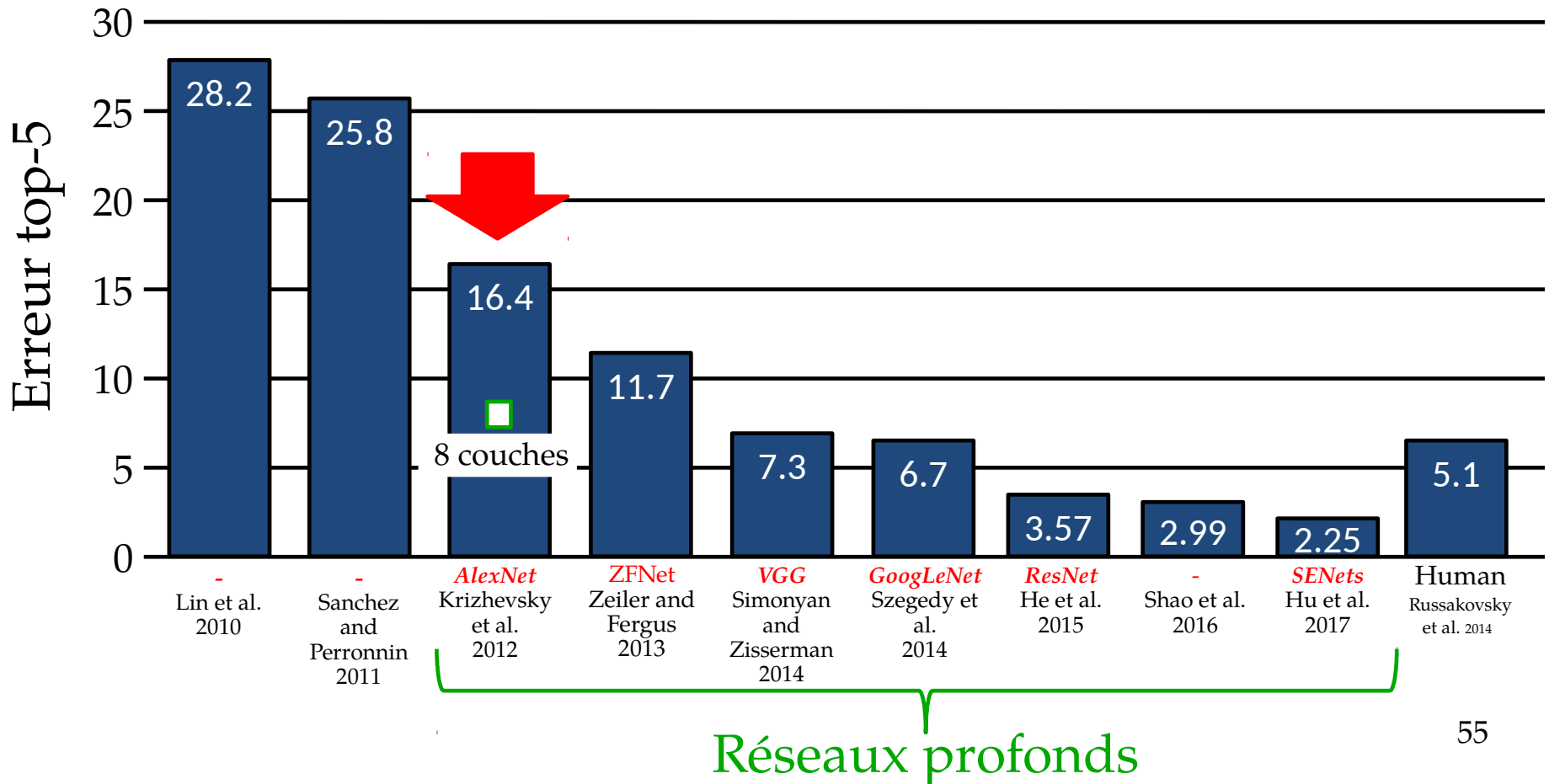
of pixels used in training

10^{14} **IMAGENET**

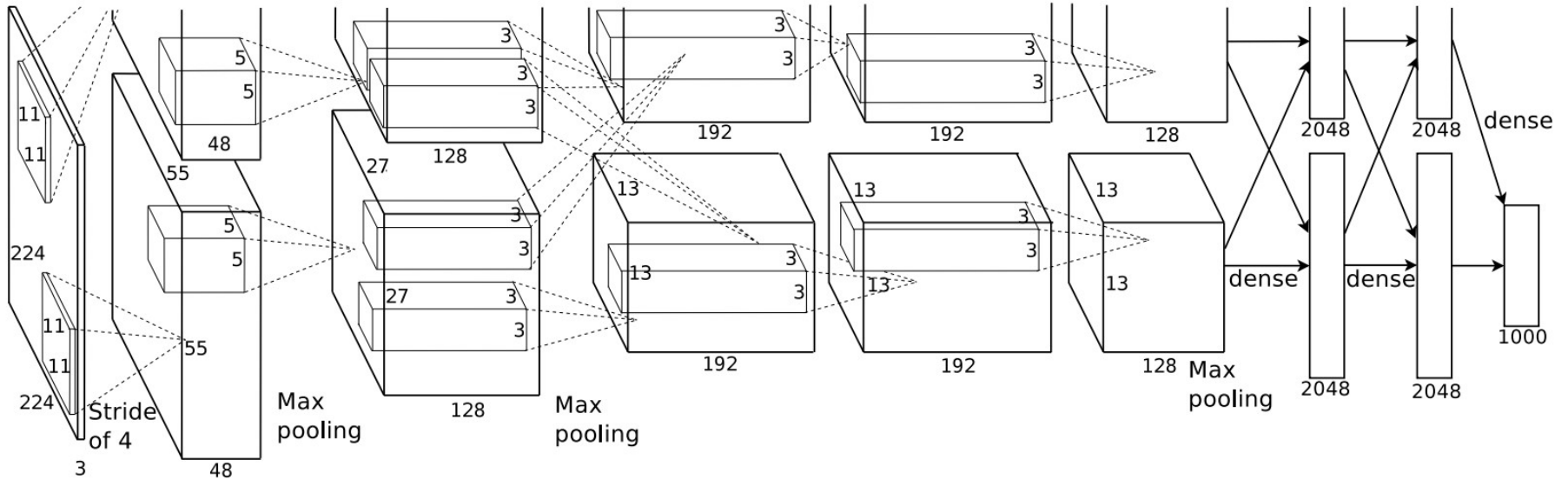
Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

Large Scale Visual Recognition Challenge

- *Image Classification Challenge* :
 - 1,000 classes d'objets
 - 1,431,167 images



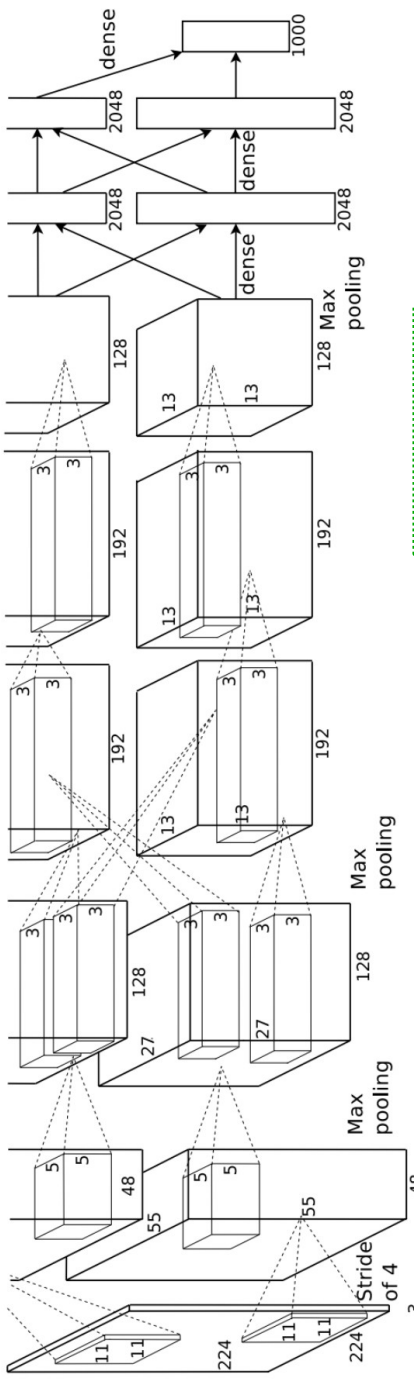
AlexNet



- 8 couches
- 60M paramètres
- Apparition des ReLU
- Dropout de 0.5
- Entraîné sur deux cartes GTX 580 (3 Go) en parallèle

AlexNet

majorité des paramètres



[1000] FC8: 1000 neurons (class scores)

[4096] FC7: 4096 neurons

[4096] FC6: 4096 neurons

Classificateur puissant (besoin dropout)

[6x6x256] MAX POOL3: 3x3 filters at stride 2

[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1

[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1

[13x13x256] NORM2: Normalization layer

[13x13x256] MAX POOL2: 3x3 filters at stride 2

[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2

[27x27x96] NORM1: Normalization layer

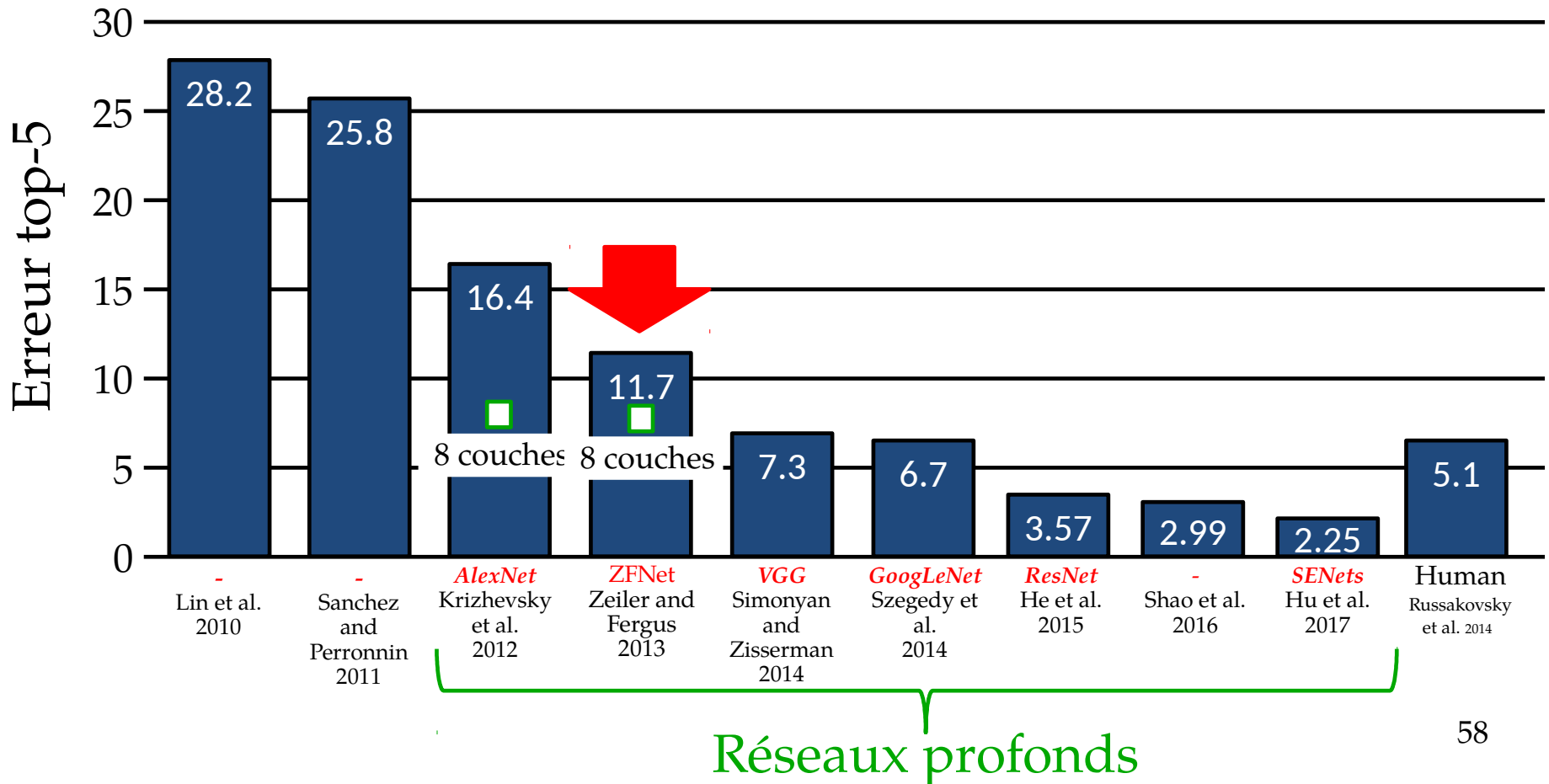
[27x27x96] MAX POOL1: 3x3 filters at stride 2

[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0

[227x227x3] INPUT

réduction rapide

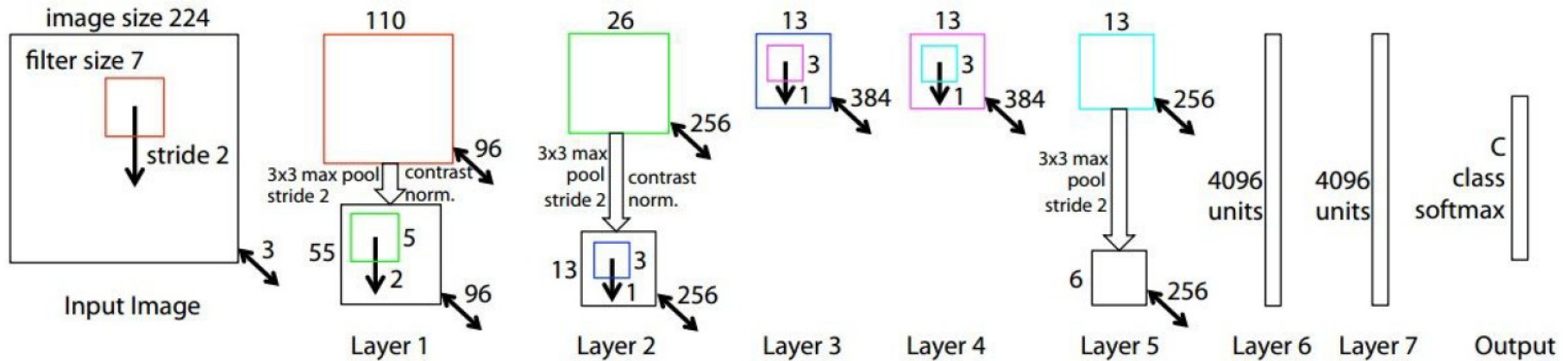
Large Scale Visual Recognition Challenge



ZFNet

ZFNet

[Zeiler and Fergus, 2013]



TODO: remake figure

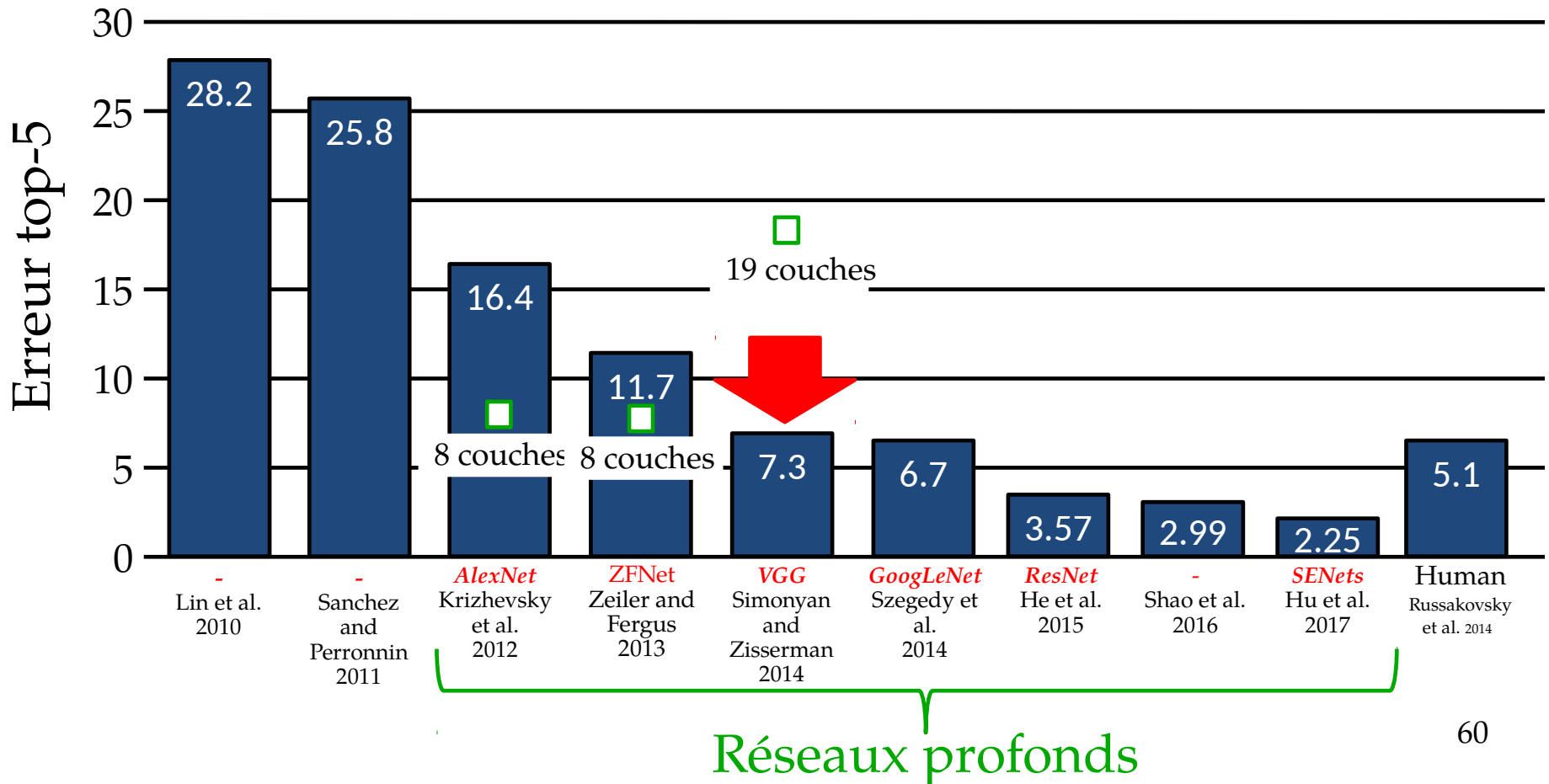
AlexNet but:

CONV1: change from (11x11 stride 4) to (7x7 stride 2)

CONV3,4,5: instead of 384, 384, 256 filters use 512, 1024, 512

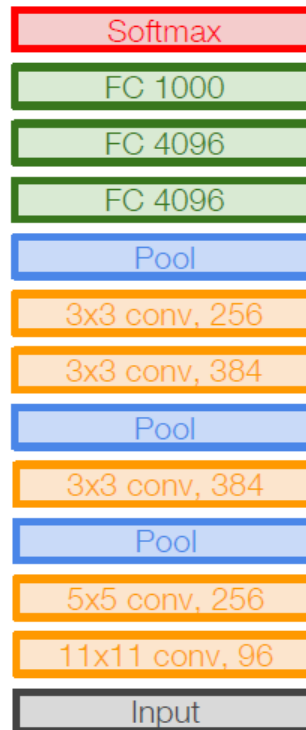
ImageNet top 5 error: 16.4% -> 11.7%

Large Scale Visual Recognition Challenge

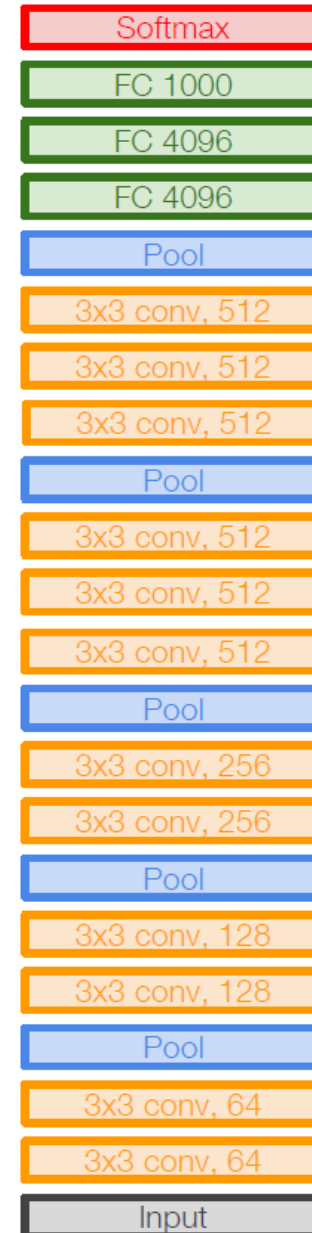


VGGNet

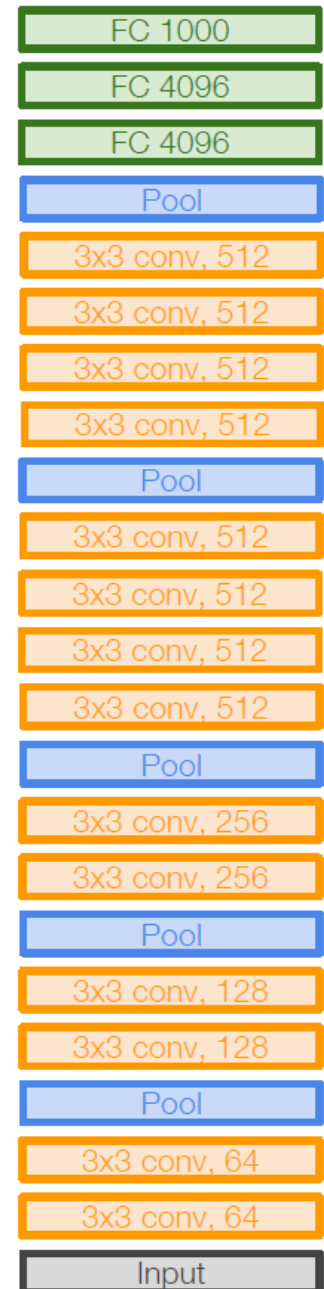
- Toujours 3 couches **fully-connected** comme classificateur
- 16-19 couches
- 138M paramètres
- Que des convolutions 3x3
- Empilement de 3 convolution 3x3 a le même champ récepteur qu'un filtre 7x7
 - Mais plus de non-linéarité (si ReLU)
 - Moins de paramètres : $3(3^2C^2)$ vs. 7^2C^2 , avec C channels en entrée-sortie (économie de 45%)



AlexNet



VGG16



VGG19