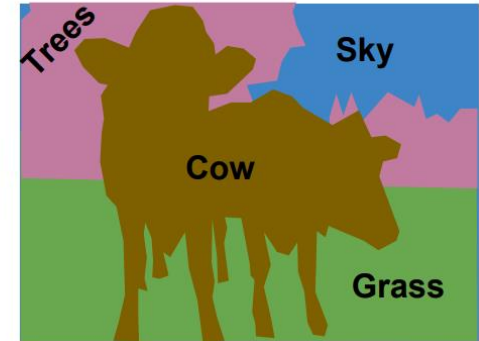
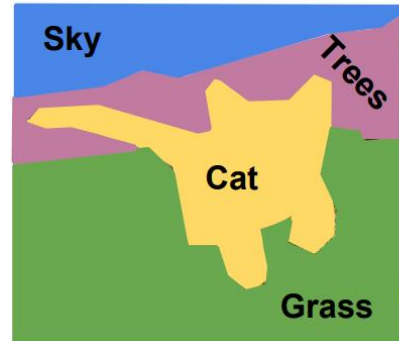


Segmentation

Segmentation sémantique

- Attribuer une classe pour chaque pixel
- Pas tenir compte des différentes instances
 - toutes les vaches, arbres, etc. sont groupés ensemble
- Limite les applications : ne va pas distinguer les voitures, empêchant le *tracking*

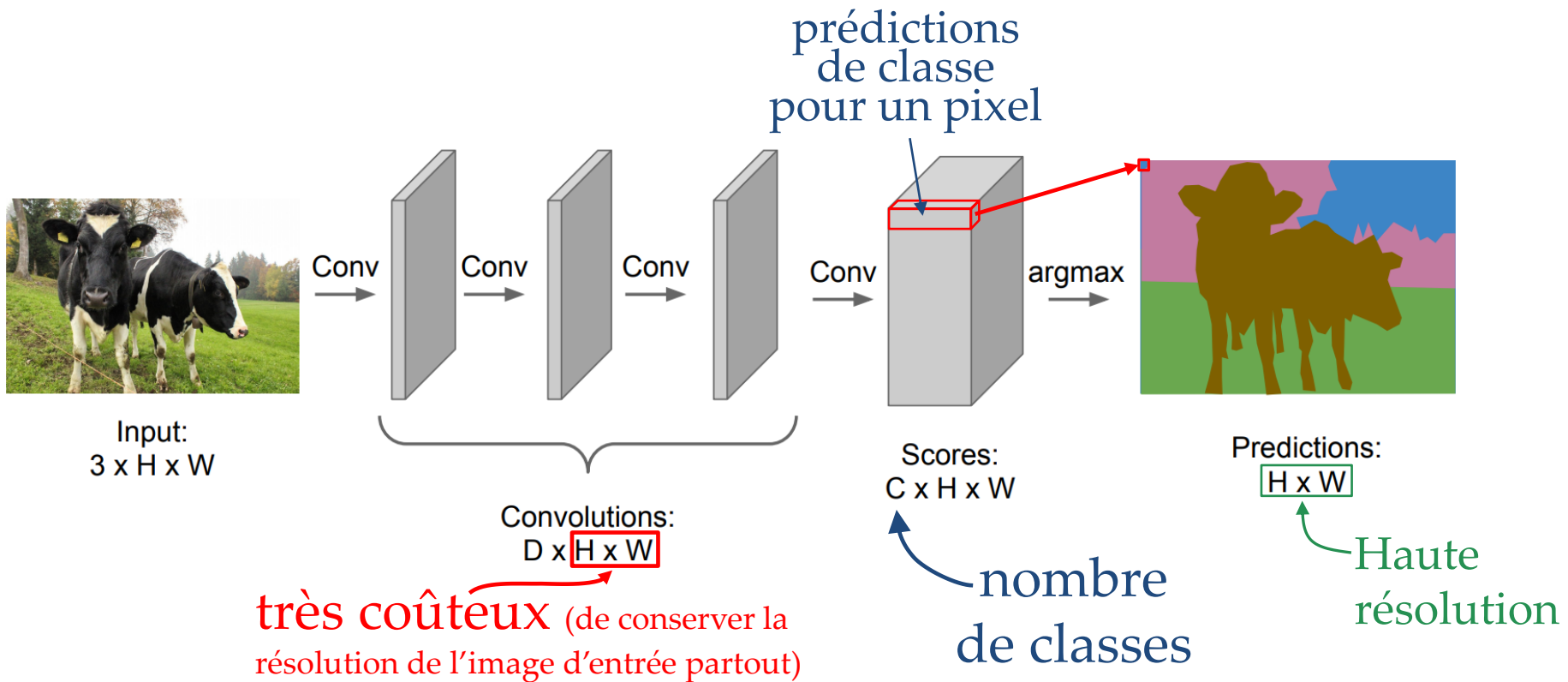


Segmentation sémantique

- S'entraîne en calculant la perte pour chacun des pixels
 - nécessite d'avoir la vérité-terrain pour **TOUS LES PIXELS** (\$\$\$)
- Souvent, fort déséquilibre des classes (**background** dominant)
 - augmenter la perte pour les classes moins représentées
 - *focal loss* peut aider à réduire l'impact

(naïve) Approche Fully-Convolutional

- Dernières couches de convolutions agissent comme des *fully-connected*, pour faire la classification

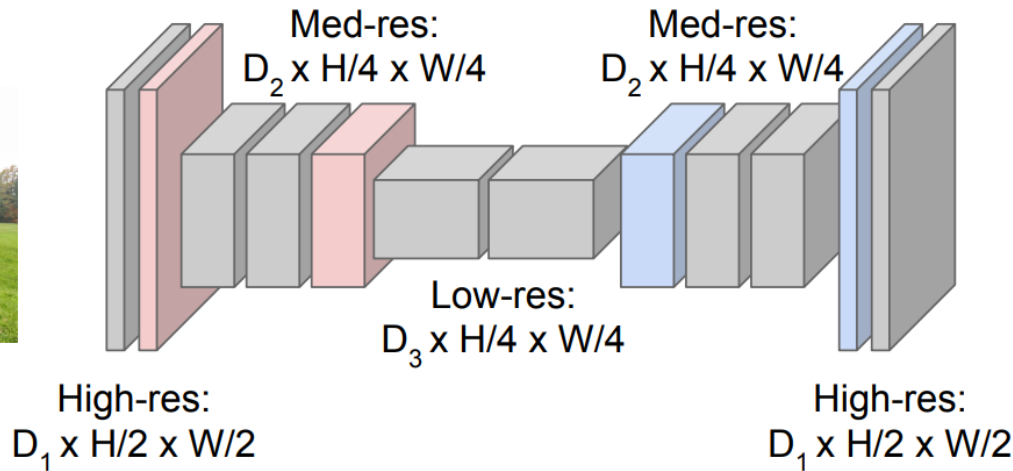


Hourglass : réduire les calculs

downsampling upsampling



Input:
 $3 \times H \times W$

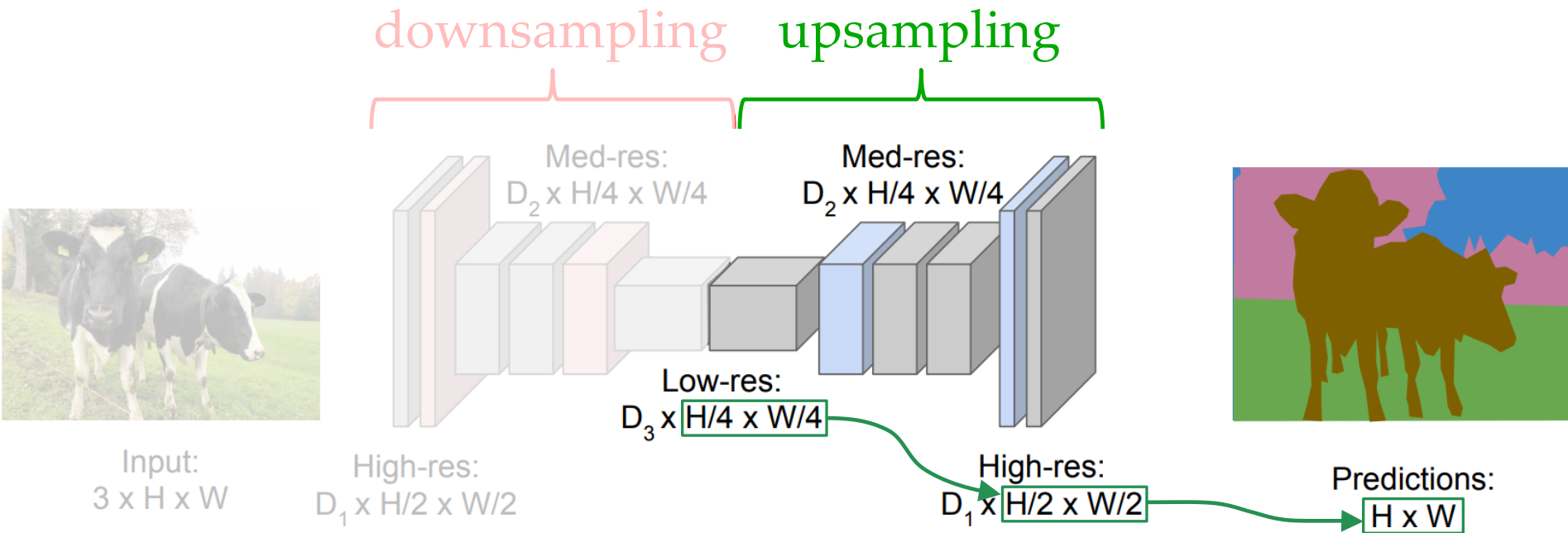


Predictions:
 $H \times W$

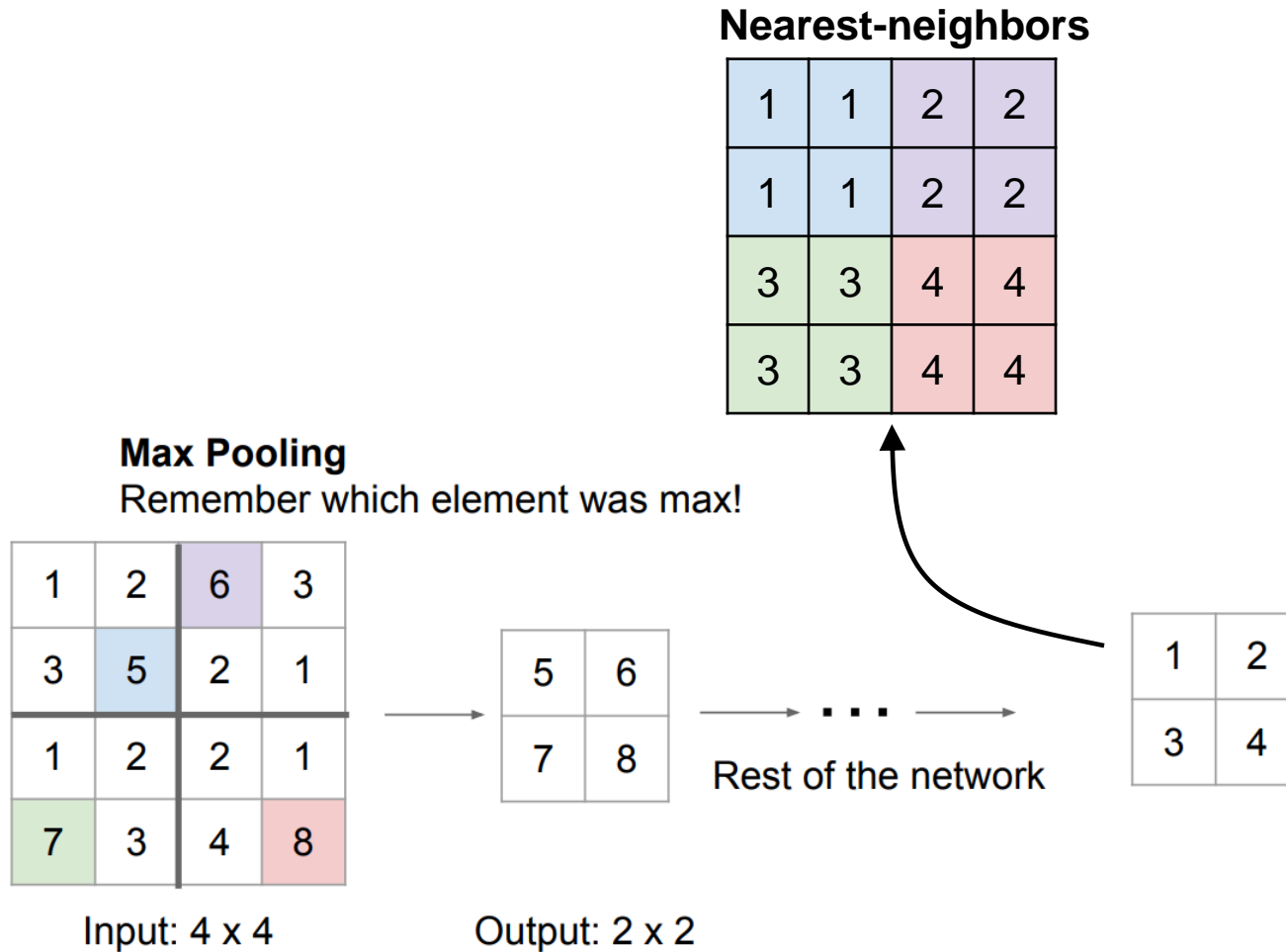
- Perte d'information spatiale fine
 - baisse de résolution H, W
 - *verra plus loin comment corriger cela*

Opération d'upsampling

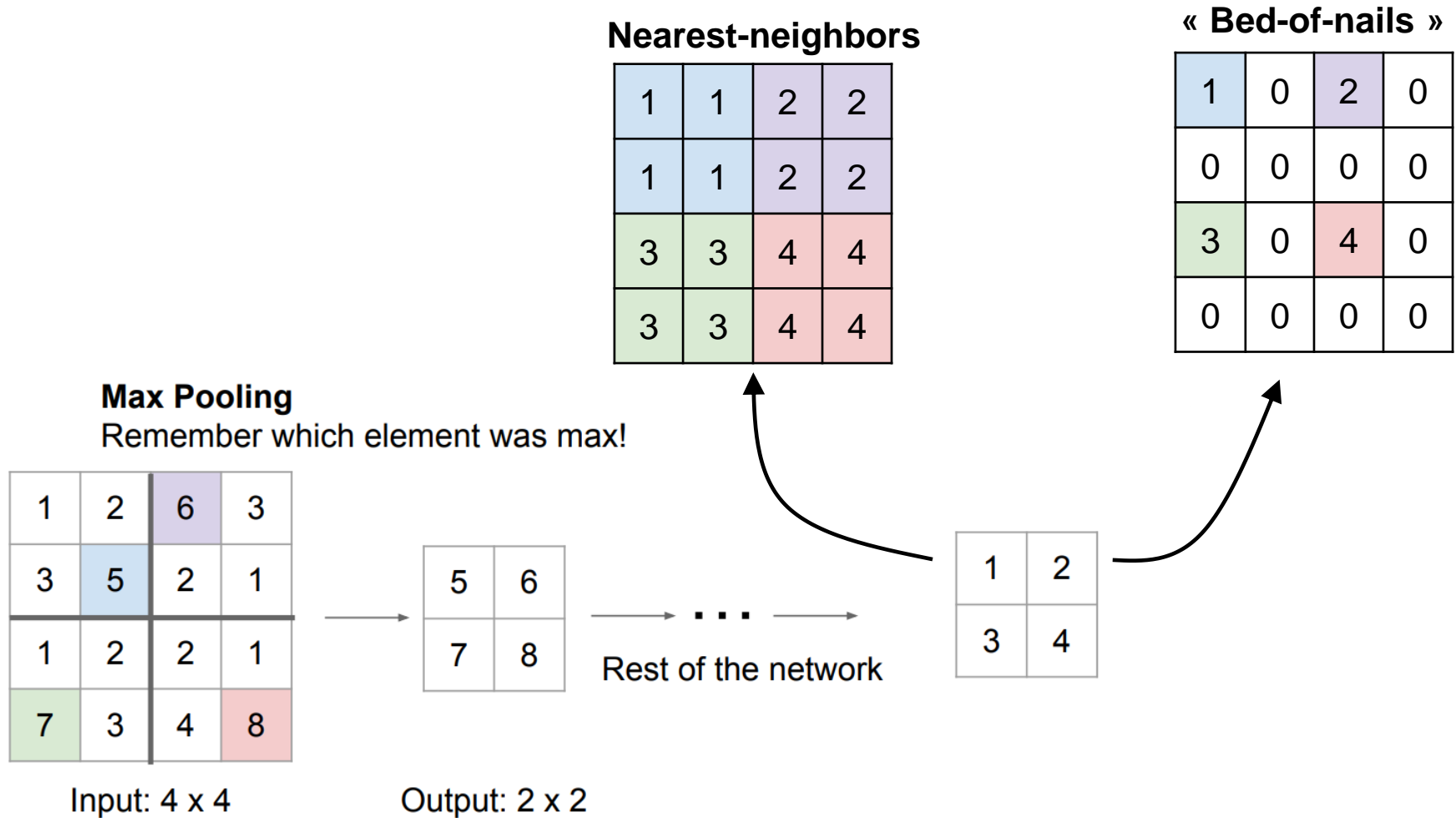
- Pour augmenter la résolution spatiale
- Plusieurs approches possibles



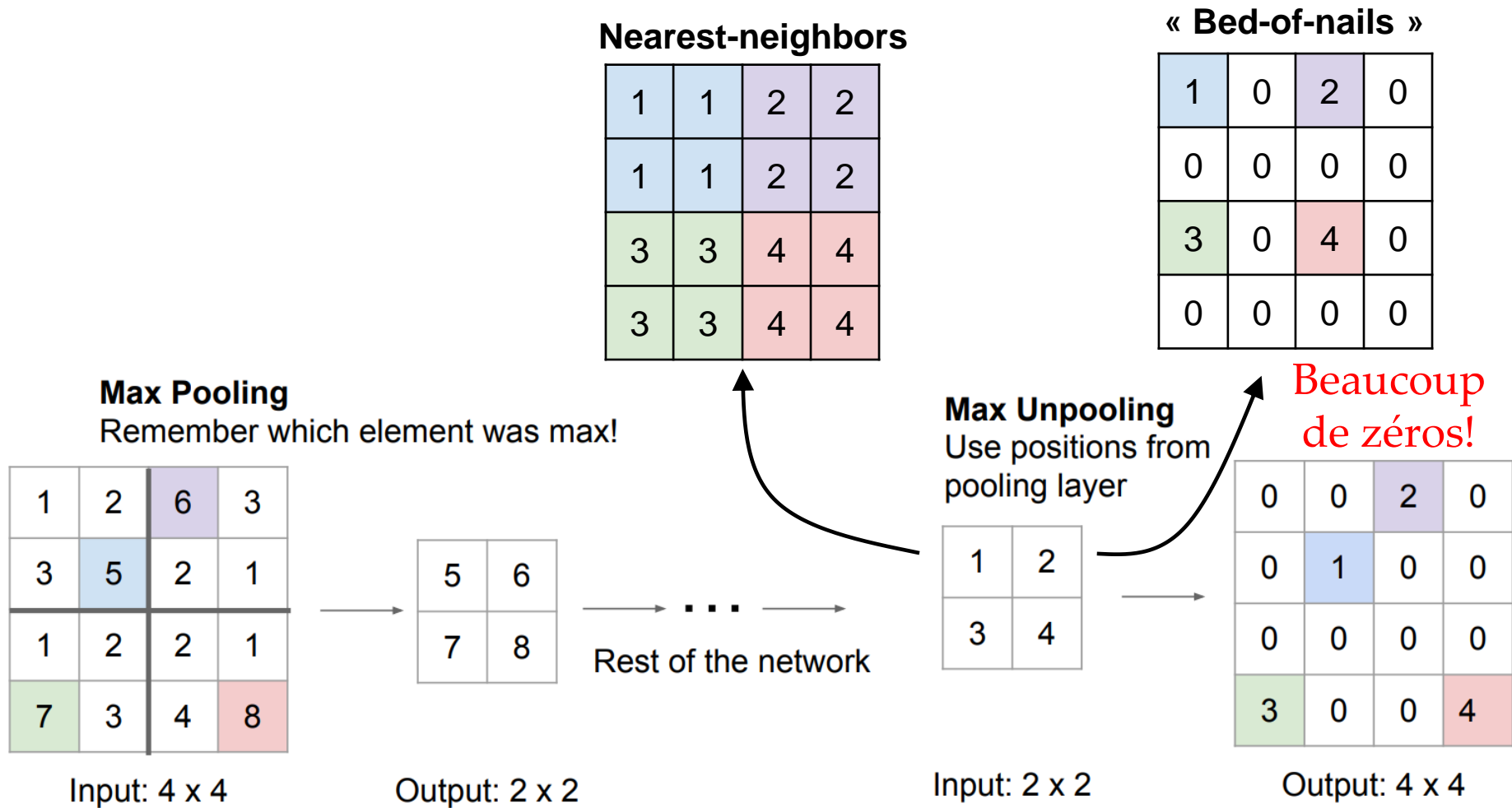
Upsampling : sans paramètres



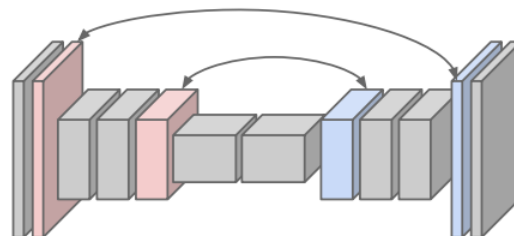
Upsampling : sans paramètres



Upsampling : sans paramètres



Corresponding pairs of
downsampling and
upsampling layers



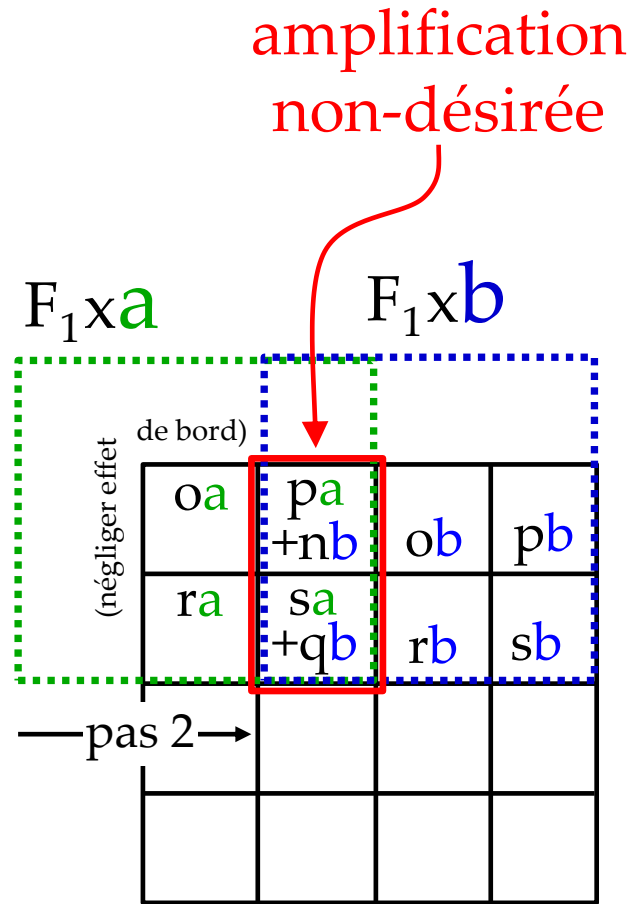
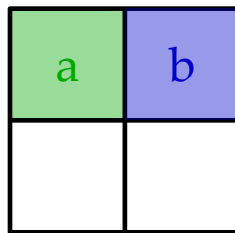
Upsampling avec paramètres : upconvolution

- Des filtres F appris
- Exemple : filtre 3x3, pas de 2
- Peut créer des **artéfacts en forme d'échiquier** dans la sortie
- S'évite avec
 - filtre 2x2 stride 2
 - filtre 4x4 stride 2
 - etc...

F_1

k	l	m
n	o	p
q	r	s

3×3



Upsampling 2×2 → 4×4

Other names:

- Deconvolution (bad)
- Upconvolution
- Fractionally strided convolution
- Backward strided convolution

Upsampling : sub-pix + shuffle

Checkerboard artifact free sub-pixel convolution

A note on sub-pixel convolution, resize convolution and convolution resize

Andrew Aitken*, Christian Ledig*, Lucas Theis*, Jose Caballero, Zehan Wang, Wenzhe Shi*
Twitter, Inc.¹

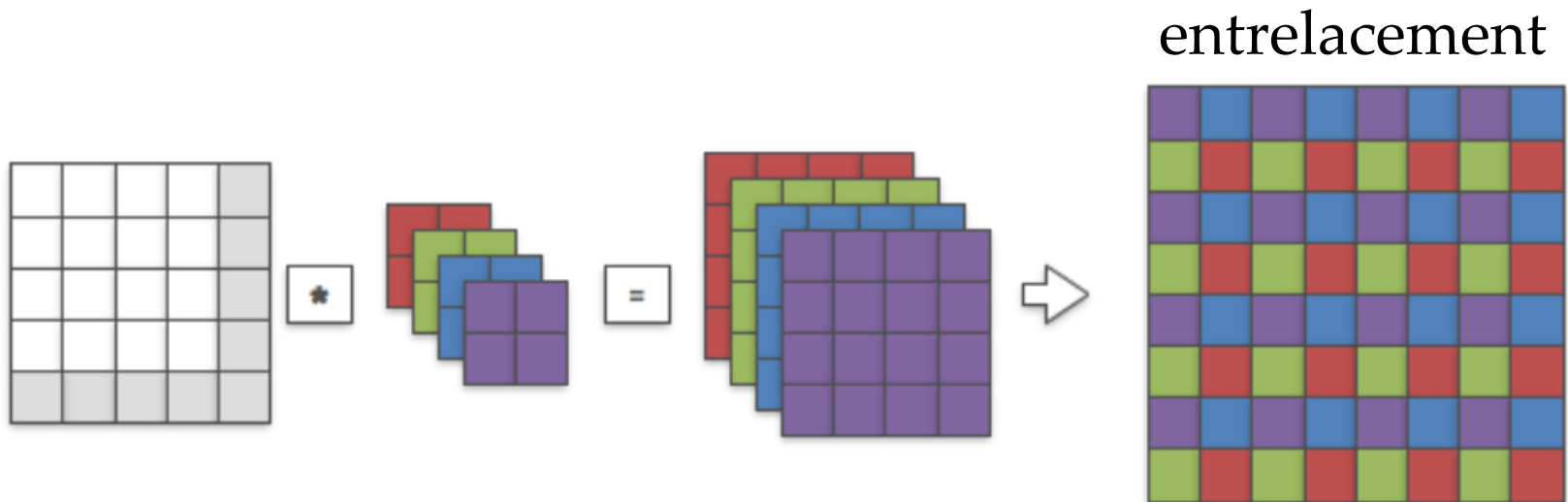


Figure 2: Sub-pixel convolution can be interpreted as convolution + shuffling.

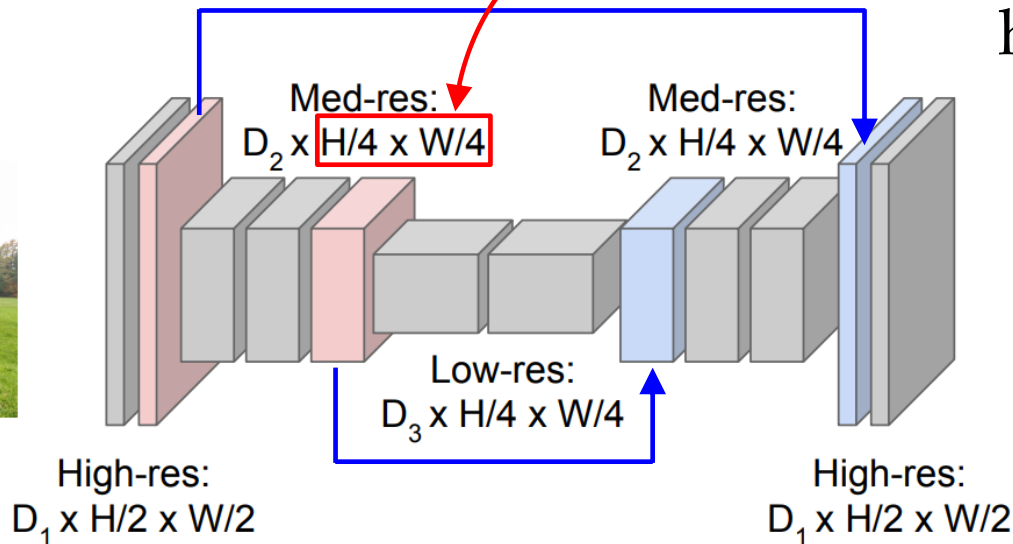
Perte d'information spatiale fine

- Approche *hourglass* fait **perdre l'information** de haute résolution

... nécessaire pour segmentation à haute résolution



Input:
 $3 \times H \times W$



Predictions:
 $H \times W$

- Pour compenser, ajouter des *skips connections* entre les couches de même résolution : U-Net
- Rappel : *skip* n'ajoute pas de paramètres

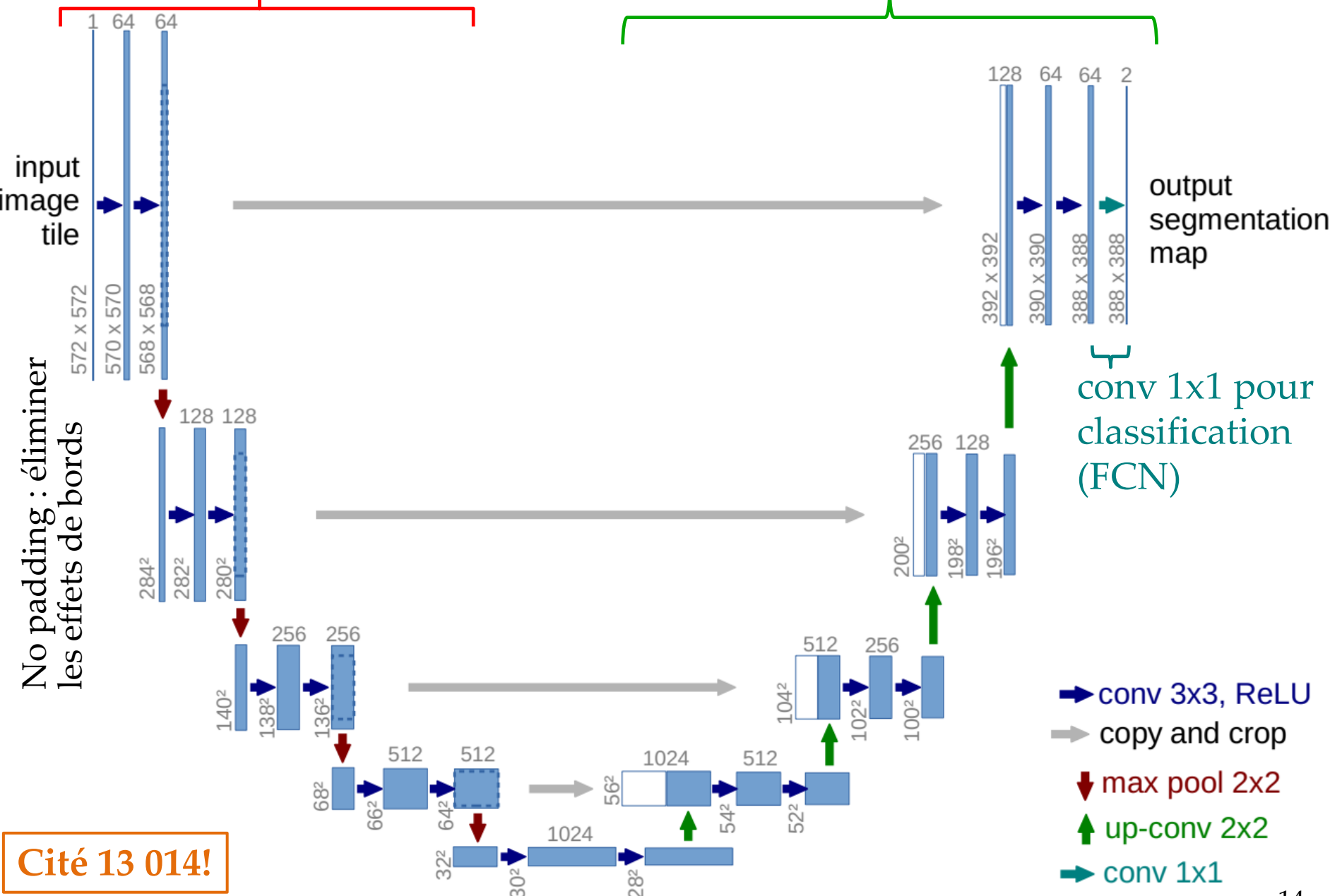
U-net

- Vient du domaine de l'imagerie biomédicale
- Quantité de données plus limitée
- Emploi des déformations élastiques pour faire *data augmentation*
 - rappel : comprendre la physique du problème
 - tissus humains sont mous et déformables
- Perte qui tient compte de la position
 - perte plus élevée pour les erreurs sur background en pourtour des autres classes

downsampling

U-net

upsampling



Cité 13 014!

Segmentation d'instances

Segmentation d'instances

- Hybride entre **segmentation** et **détection**
- Peut le percevoir comme le problème complet de la vision 2D



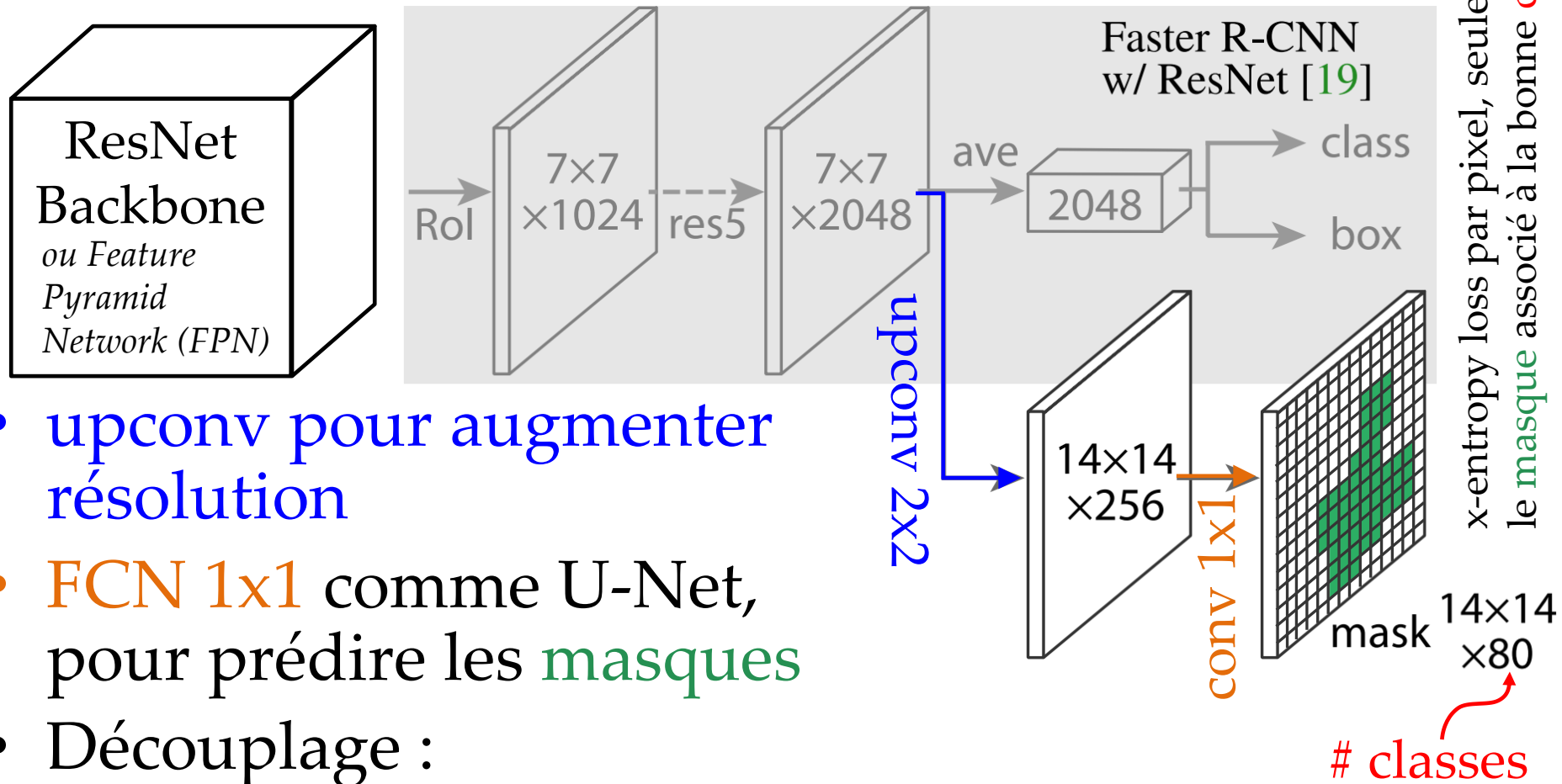
Tiré de cs231n

DOG, DOG, CAT

Mask R-CNN

- Performe de la segmentation d'instances
- Modification de Faster R-CNN
 - Simple ajout d'une tête supplémentaire pour prédire un masque binaire !
- N'ajoute qu'un tout petit surplus de calcul
- RoIPool → RoIAlign
- 200 ms par trame (*frame*)

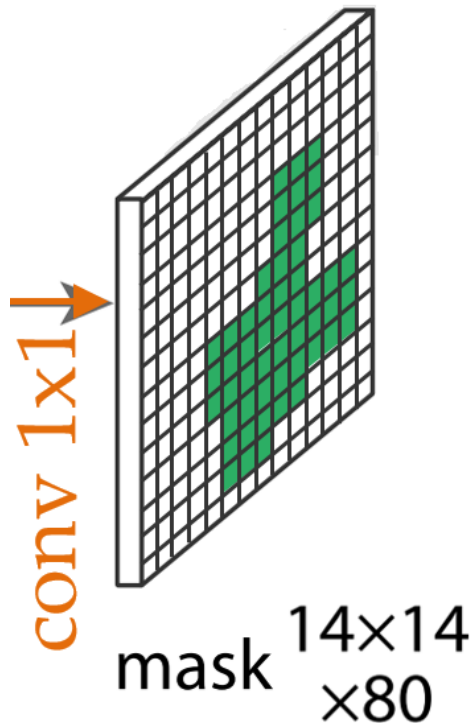
Mask R-CNN : Modification à Faster R-CNN



- **upconv** pour augmenter résolution
- **FCN 1x1** comme U-Net, pour prédire les **masques**
- Découplage :
 - un **masque** par **classe**
 - élément clé des performances

Découplage des masques

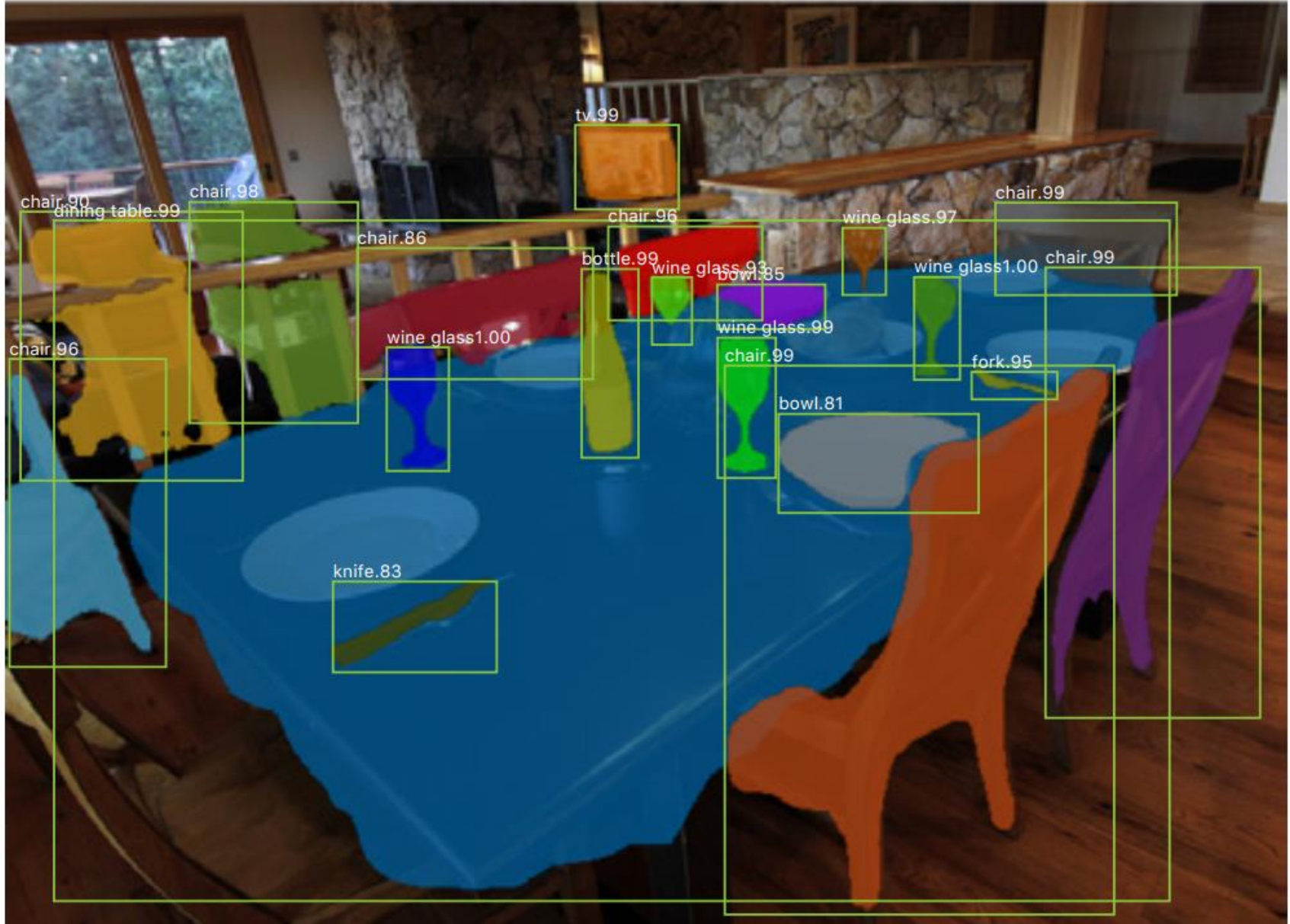
$$\text{softmax} : \hat{y}_i = \frac{\exp(z_i)}{\sum_{j \in \text{groupe}} \exp(z_j)} \quad \left. \vphantom{\sum_{j \in \text{groupe}}} \right\} \text{couplage}$$



	AP	AP ₅₀	AP ₇₅
<i>softmax</i>	24.8	44.1	25.1
<i>sigmoid</i>	30.3	51.2	31.5
	+5.5	+7.1	+6.4

(b) **Multinomial vs. Independent Masks** (ResNet-50-C4): *Decoupling* via per-class binary masks (sigmoid) gives large gains over multinomial masks (softmax).

Mask R-CNN : exemple



Exemple de segmentation

- Données de *Cityscape*



Conclusion

- Réseaux de neurones s'appliquent aussi sur des problèmes de vision numérique au-delà de la classification
- **Détection**
 - doit retrouver les objets d'une même catégorie, sans connaître au préalable leur nombre
 - méthodes basés sur propositions (300-2000)
 - concept d'anchor box
 - les plus rapides : proposition par réseau (RPN)
 - régression des anchor box pour localisation
 - existe aussi des méthodes par grille (YOLO)

Conclusion

- **Segmentation sémantique**
 - cherche à trouver la classe de chaque pixel
 - architecture U-Net
 - hourglass : downsampling et upsampling
 - méthodes d'upsampling (upconv)
 - skip connexion pour préserver l'info. spatiale fine
- **Segmentation d'instance**
 - hybride entre détection + segmentation
 - Mask R-CNN : ajout d'une branche *masque*
 - importance de découpler les masques entre les classes
- Fully-convolutional network (FCN)