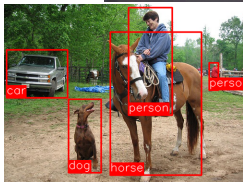
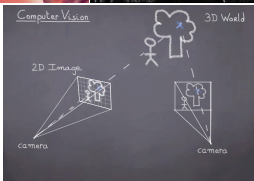
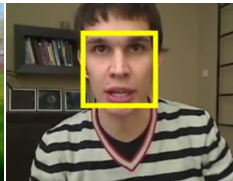
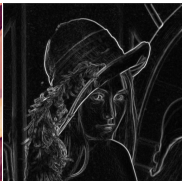
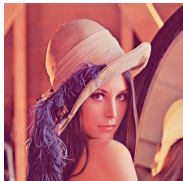


# Computer vision and image processing introduction

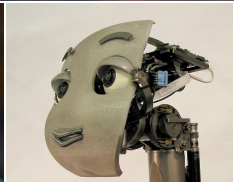
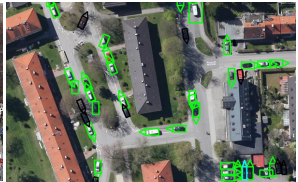
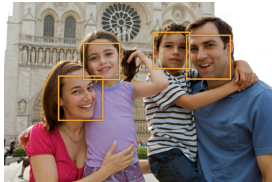
Ronan Sifre

Credits to Yannis Avrithis <https://sif-dlv.github.io/>

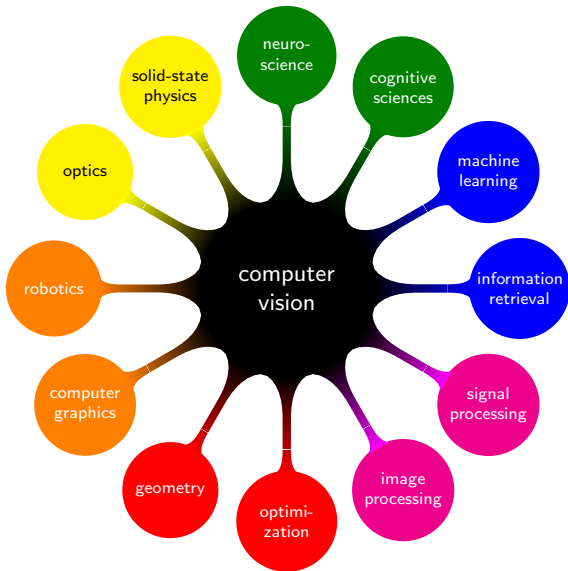
# computer vision in images



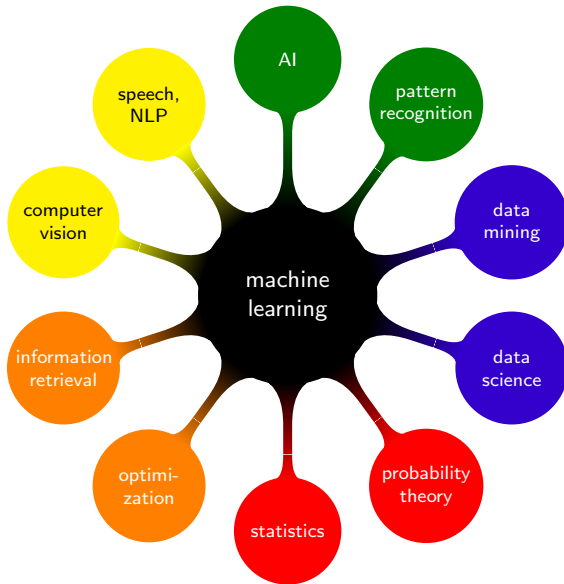
# computer vision in images



# computer vision—related fields



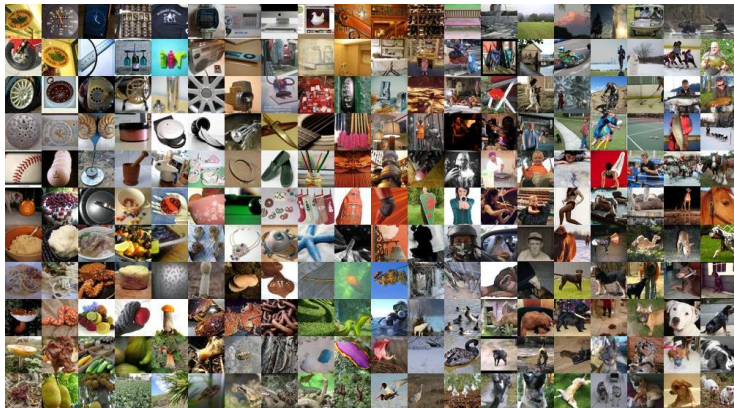
# machine learning—related fields



**modern deep learning**

# ImageNet

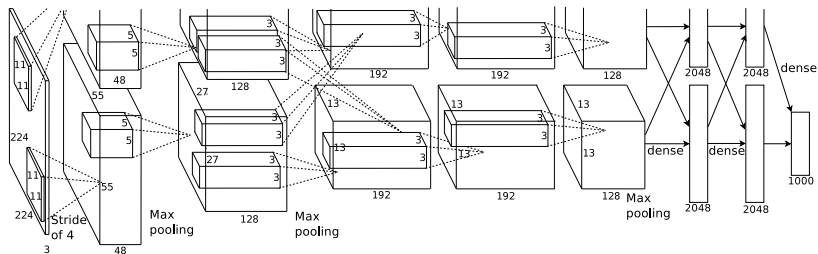
[Russakovsky et al. 2014]



- 22k classes, 15M samples
- ImageNet Large-Scale Visual Recognition Challenge (ILSVRC): 1000 classes, 1.2M training images, 50k validation images, 150k test images

# AlexNet

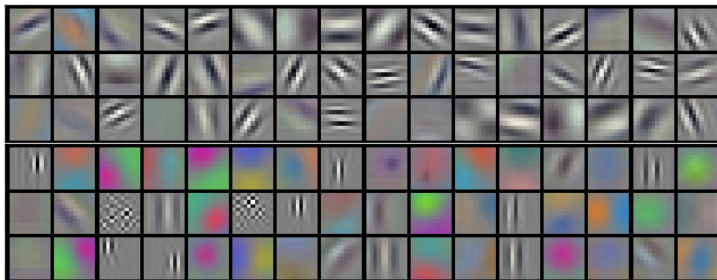
[Krizhevsky et al. 2012]



- implementation on two GPUs; connectivity between the two subnetworks is limited
- ReLU, data augmentation, local response normalization, dropout
- outperformed all previous models on ILSVRC by 10%

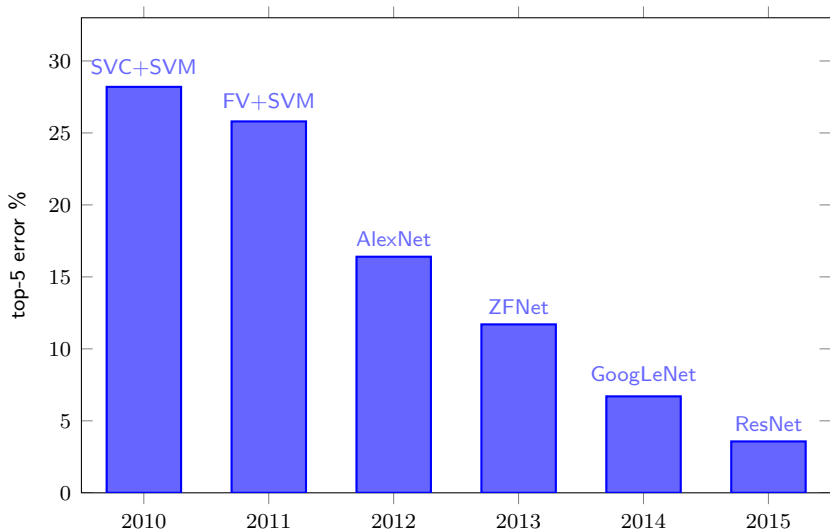


## learned layer 1 kernels



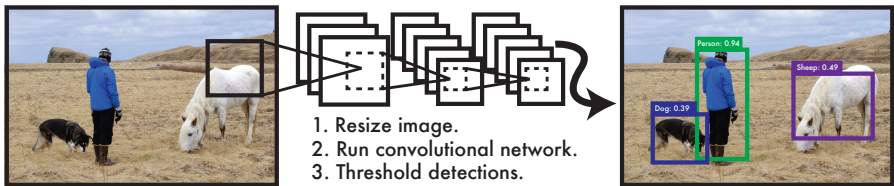
- 96 kernels of size  $11 \times 11 \times 3$
- top: 48 GPU 1 kernels; bottom: 48 GPU 2 kernels

# ImageNet classification performance



# object detection

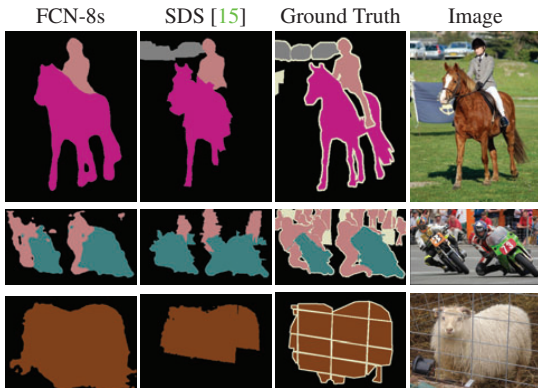
[Redmon et al. 2016]



- learn to detect objects as a single classification and regression task, without scanning the image or detecting candidate regions
- first object detector to operate at 45fps

# semantic segmentation

[Long et al. 2015]



- learn to upsample
- apply to pixel-dense prediction tasks

# instance segmentation and pose estimation

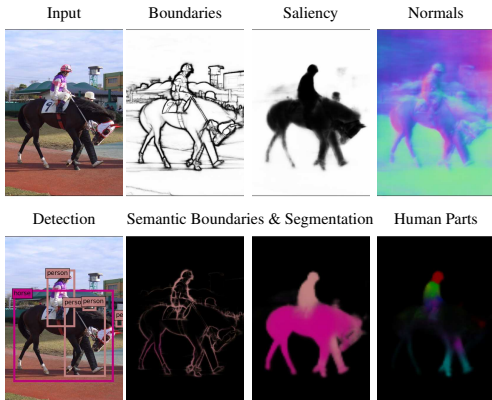
[He et al. 2017]



- semantic segmentation per detected region
- pose estimation as regression

# multi-task learning

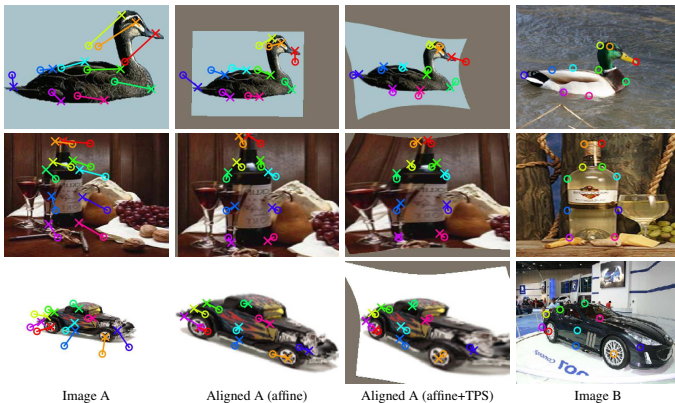
[Kokkinos 2017]



- learn several vision tasks with a joint network architecture including task-specific skip layers

# geometric matching

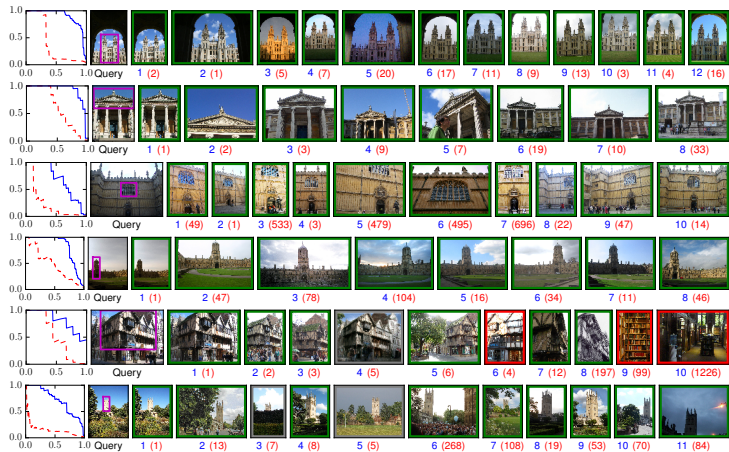
[Rocco et al. 2017]



- mimic the standard steps of feature extraction, matching and simultaneous inlier detection and model parameter estimation
- still trainable end-to-end

# image retrieval

[Gordo et al. 2016]

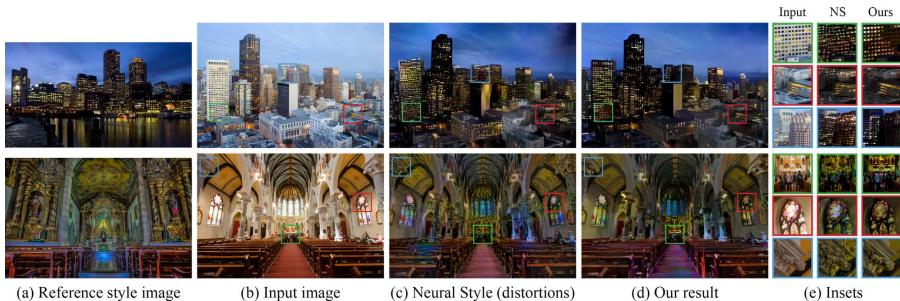


- learn to match
- apply as generic feature extractor



# photorealistic style transfer

[Luan et al. 2017]



- generate same scene as input image
- transfer style from reference image
- photorealism regularization

# image captioning

[Vinyals et al. 2017]

A person riding a motorcycle on a dirt road.



Two dogs play in the grass.



A skateboarder does a trick on a ramp.



A dog is jumping to catch a frisbee.



A group of young people playing a game of frisbee.



Two hockey players are fighting over the puck.



A little girl in a pink hat is blowing bubbles.



A refrigerator filled with lots of food and drinks.



A herd of elephants walking across a dry grass field.



A close up of a cat laying on a couch.



A red motorcycle parked on the side of the road.



A yellow school bus parked in a parking lot.



Describes without errors

Describes with minor errors

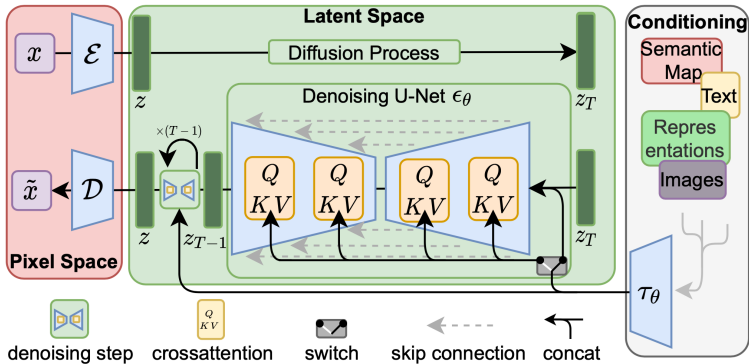
Somewhat related to the image

Unrelated to the image

- image description by deep CNN
- language generation by RNN

# Generative models

GAN, Diffusion, VAE, MAE, DAE.



# Self-supervised models

RotNet, Deep Cluster, BYOL, DINO, iBOT

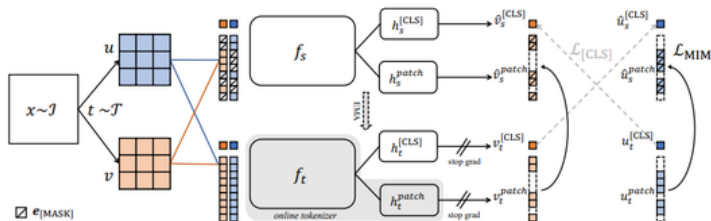


Figure 3: **Overview of iBOT framework, performing masked image modeling with an *online tokenizer*.** Given two views  $u$  and  $v$  of an image  $x$ , each view is passed through a teacher network  $h_t \circ f_t$  and a student network  $h_s \circ f_s$ . iBOT minimizes two losses. The first loss  $\mathcal{L}_{[CLS]}$  is self-distillation between cross-view [CLS] tokens. The second loss  $\mathcal{L}_{MIM}$  is self-distillation between in-view patch tokens, with some tokens masked and replaced by  $e_{[MASK]}$  for the student network. The objective is to reconstruct the masked tokens with the teacher networks' outputs as supervision.