

Learning with CNNs

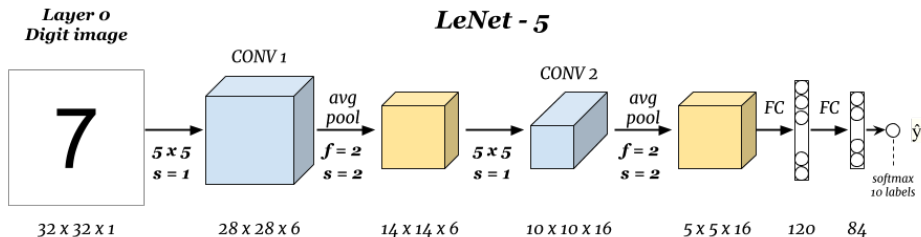
Tianxiang (Adam) Gao

February 13, 2025

Outline

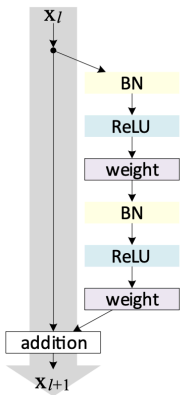
- 1 Classic CNNs: Inception, MobileNets
- 2 Practical Advice for Using CNNs
- 3 Object Detection
- 4 Face Recognition
- 5 Neural Style Transfer

Recap: Convolutional Neural Networks (CNNs)



- **Convolution operation:** It slides a small **filter** over the input image, performing a **locally linear transformation** to produce a feature map that detects patterns.
- **Padding and stride:** Methods for controlling feature map size, preserving spatial dimensions, and improving *computational efficiency*.
- **Convolution over volumes and with multiple filters:** The input can be multi-channel, and so as the output with the use of multiple filters.
- **Weight Sharing and Sparsity:** Neurons in CNNs **share** weights across locations, with each output relying on a **small, localized** input region.
- **Hierarchical Feature Detection:** Early layers capture basic features (like edges), which later layers combine into higher-level features.

Recap: Stabilizing CNN Training



- **Pooling**: Reduces the spatial dimensions of feature maps by downsampling, typically using max or average pooling.
- **Batch Normalization**: Normalizes pre-activation values at hidden layers, mitigating internal covariate shift and accelerating training.
- **Skip Connections**: Create shortcuts by directly adding input to output, stabilizing information flow in deep neural networks.
- **Classic CNNs**: Spatial dimensions shrink while the number of channels increases as depth grows. Overparameterized and deeper CNNs are generally preferred.

Recap: Semantic Segmentation and U-Net



Input

segmented

- 1: Person
- 2: Purse
- 3: Plants/Grass
- 4: Sidewalk
- 5: Building/Structures



Semantic Labels

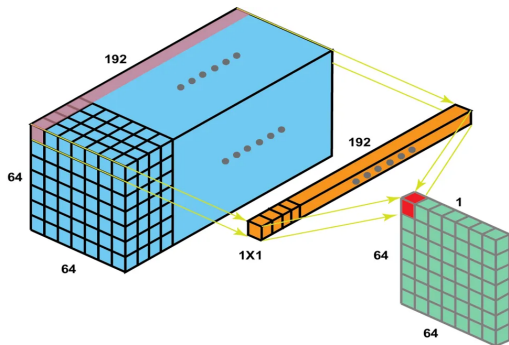
- Assigns a class label to **each pixel** in the image
- The output is an **image** with the same spatial dimensions as the input
- The encoder (typically a CNN) extracts **hierarchical** feature representations
- The decoder uses **transposed convolution** (deconvolution) to restore spatial resolution
- Skip connections **combines** low-level spatial details with high-level semantic features

Outline

- 1 Classic CNNs: Inception, MobileNets
- 2 Practical Advice for Using CNNs
- 3 Object Detection
- 4 Face Recognition
- 5 Neural Style Transfer

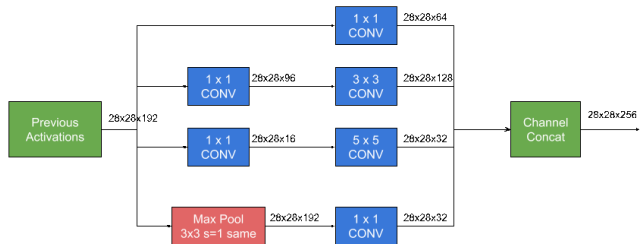
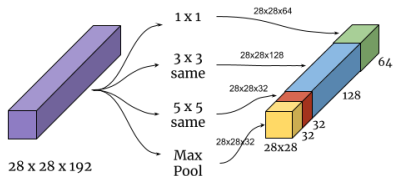
1 × 1 Convolutions

$$\underbrace{\begin{bmatrix} 3 & 0 & 1 \\ 1 & 5 & 8 \\ 2 & 7 & 2 \end{bmatrix}}_{\text{input image } 3 \times 3} * \underbrace{\begin{bmatrix} 2 \end{bmatrix}}_{\text{filter } 1 \times 1} = \underbrace{\begin{bmatrix} 6 & 0 & 2 \\ 2 & 10 & 16 \\ 4 & 14 & 4 \end{bmatrix}}_{\text{feature map } 3 \times 3}$$



- Recall that each channel in the input data serves as an **indicator map** for certain patterns detected by the filter.
- Each pixel, with multiple channels, represents a **set** of indicators that capture the distribution of **multiple patterns** within a local region of the original image.
- Applying a 1 × 1 convolution enables **cross-channel interaction** to detect complex patterns that integrate multiple underlying features.
- This approach is useful for combining lower-level features into **higher-level** representations.

Inception Networks



- Inception networks allow the model to freely select the best filter sizes within a layer, rather than manually choosing them, by providing multiple filter options.
- This approach can be computationally expensive (e.g., a 5x5 filter requires about 120 million operations).

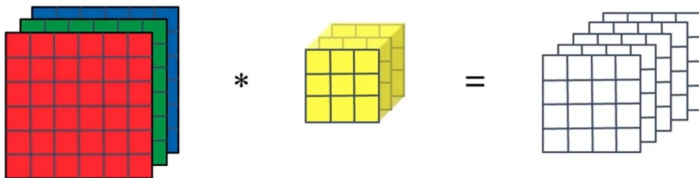
$$(28 \times 28 \times 32) \times (5 \times 5 \times 192) \approx 120 \text{ million}$$

- Applying 1 × 1 convolutions reduces computation by approximately 90%.

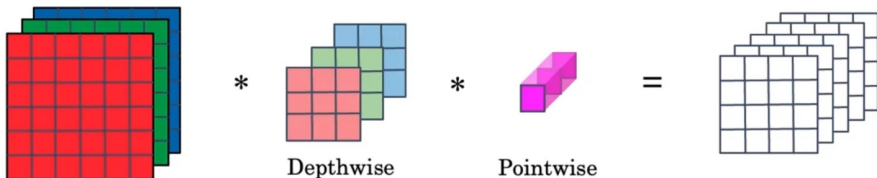
$$(28 \times 28 \times 16) \times (1 \times 1 \times 192) + (28 \times 28 \times 32) \times (5 \times 5 \times 16) \approx 2 \text{ million} + 10 \text{ million}$$

MobileNets

Normal Convolution



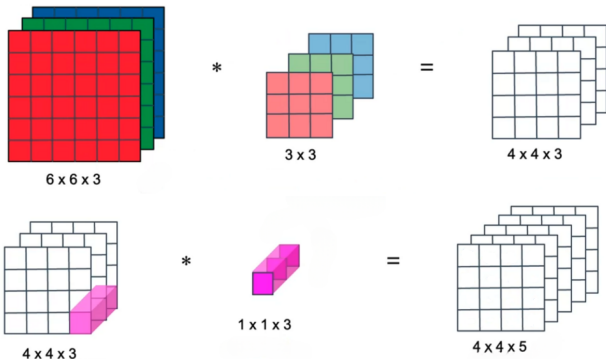
Depthwise Separable Convolution



The computational cost of a normal convolution is:

$$\underbrace{2,160}_{\text{Computational cost}} = \underbrace{3 \times 3 \times 3}_{\# \text{ params}} \times \underbrace{4 \times 4}_{\# \text{ locations}} \times \underbrace{5}_{\# \text{ filters}}$$

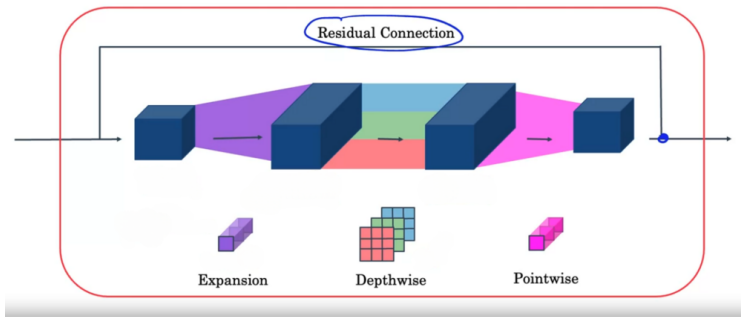
Depthwise Separable Convolution



- Each channel has an **independent** filter for the convolution operation.
- The result is then processed using a pointwise or 1×1 convolution.
- Computational cost ratio: $\frac{1}{C_{out}} + \frac{1}{f^2}$

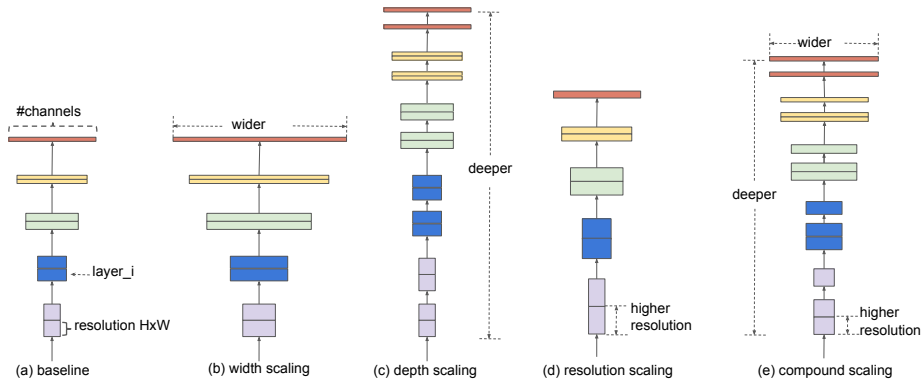
$$672 = \underbrace{3 \times 3}_{\# \text{ params}} \times \underbrace{4 \times 4}_{\# \text{ locations}} \times \underbrace{3}_{\# \text{ filters}} + \underbrace{3}_{\# \text{ params}} \times \underbrace{4 \times 4}_{\# \text{ locations}} \times \underbrace{5}_{\# \text{ filters}}$$

MobileNetV2: Inverted Residual Blocks



- The **expansion convolution** allows CNNs to learn richer functions by increasing the number of channels before depthwise convolution
- When deploying on resource-constrained devices, a **pointwise convolution** is used to project back to a smaller number of channels, reducing memory usage.
- VGG-16 \approx 528 MB, ResNet-152 \approx 230 MB, Inception \approx 85 MB, but MobileNet \approx **5-16 MB**.

EfficientNet: Compound Scaling



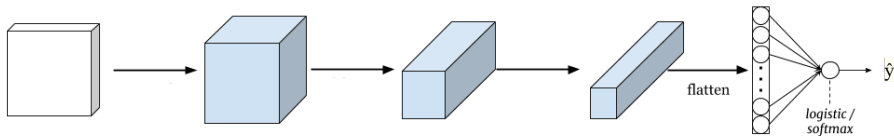
- **Depth (d)**: A deeper CNN can capture **richer hierarchical features** and more complex patterns.
- **Width (w)**: Wider networks can **capture more fine-grained features** and **easier to train**.
- **Resolution (r)**: Higher resolution inputs help CNNs **capture more fine-grained spatial details**.
- **Compound Scaling**: EfficientNet scales all three dimensions **together**: $d = \alpha^\phi$, $w = \beta^{2\phi}$, $r = \gamma^{2\phi}$ such that $\alpha\beta^2\gamma^2 \approx 2$.

Outline

- 1 Classic CNNs: Inception, MobileNets
- 2 Practical Advice for Using CNNs
- 3 Object Detection
- 4 Face Recognition
- 5 Neural Style Transfer

Transfer Learning

- **Definition:** Transfer learning leverages a **pre-trained** model on a large dataset to solve a new, related task on a smaller dataset.
- **Motivation:** Training deep networks from scratch demands significant data and computational resources. The initial layers in CNNs capture fundamental features that are transferable across various tasks, allowing the higher layers to focus on learning task-specific features.



- **Typical Workflow:**
 - ① **Pre-training:** A model is trained on a large dataset (e.g., ImageNet).
 - ② **Freezing Layers:** Early layers are frozen to retain their learned features.
 - ③ **Fine-tuning:** With a lower learning rate, higher layers are retrained on the target dataset, allowing for adaptation to new, task-specific features.
- Frozen early layers can be considered as **fixed** feature extractors, allowing activations to be precomputed, which reduces computation and speeds up training.
- For **larger** target datasets, freeze fewer layers to enable more flexibility in learning higher-level task-specific features during fine-tuning.

Data Augmentation

- **Definition:** Data augmentation involves generating new training samples by applying various transformations to existing data, helping improve model generalization.
- **Motivation:** Increases the diversity of the training dataset, which can reduce overfitting and improve the model's performance on unseen data.
- **Flips, rotations, and scaling:**



- **Random cropping:**



- **Color Adjustments:** Changes to colors or brightness



Mixup



- **Mixup** creates new training samples by blending pairs of images and their corresponding labels:

$$\tilde{\mathbf{x}} = (1 - \lambda)\mathbf{x}_1 + \lambda\mathbf{x}_2, \quad \text{and} \quad \tilde{\mathbf{y}} = (1 - \lambda)\mathbf{y}_1 + \lambda\mathbf{y}_2,$$

where $(\mathbf{x}_i, \mathbf{y}_i)$ are original training samples, and λ is the mixing factor drawn from a Beta distribution.

- This technique helps the model generalize by learning smoother transitions between classes.

Outline

- 1 Classic CNNs: Inception, MobileNets
- 2 Practical Advice for Using CNNs
- 3 Object Detection**
- 4 Face Recognition
- 5 Neural Style Transfer

Sliding Window Detection

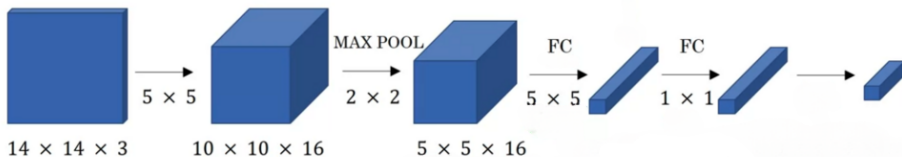
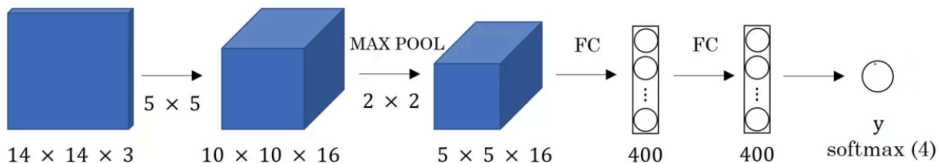
- Use a **pre-trained classifier** to identify cars in images.
- **Scan** the entire image with a sliding window, classifying each cropped section.



- **Problem:** This approach is computationally intensive due to the number of windows.

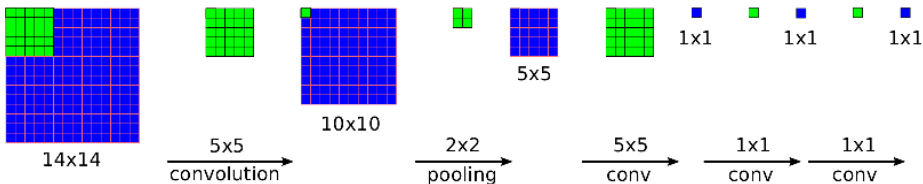
Convolutional Layers as Fully Connected Layers

Fully connected layers can be implemented by convolutional layers:

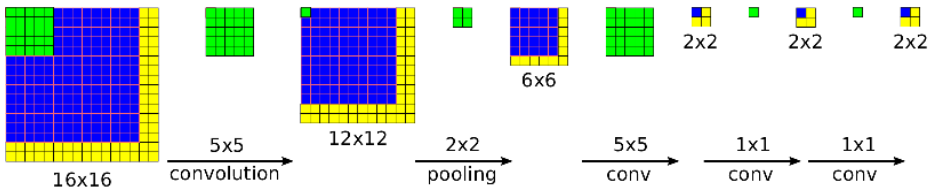


Convolution Implementation of Sliding Windows

ConvNets:

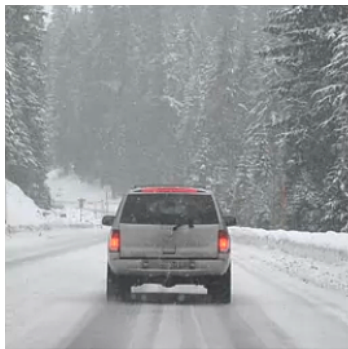


Detection:



- During **training**, a ConvNet produces only a **single** spatial output
- When applied at **test time** over a **larger** image, it produces a spatial output **map**.
- However, with a fixed sliding window size, the bounding boxes may lack precision.

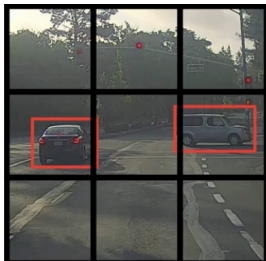
Classification with Object Localization



- **Input:** An image containing a **single object**
- **Output:** Class probability and bounding box $\mathbf{y} = (p_c, b_x, b_y, b_w, b_h)$
- $p_c \in [0, 1]$: Object presence confidence score; (b_x, b_y) : object center; b_w, b_h : width and height
- Use **squared loss** for bounding box predictions
- **Multi-classification:** the output becomes $\mathbf{y} = (p_c, b_x, b_y, b_w, b_h, c_1, c_2, \dots, c_N)$

YOLO Algorithm

YOLO Algorithm Steps: Divide image into $S \times S$ grid cells



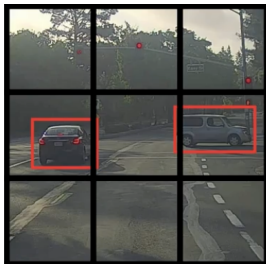
- Each cell predicts **class probabilities** and **bounding boxes**

$$y_s = (p_c, b_x, b_y, b_w, b_h, c_1, \dots, c_N), \quad \forall s \in [S]$$

- In each cell, object center (b_x, b_y) uses local coordinates in $[0, 1] \times [0, 1]$
- b_w and b_h can be greater than 1 (relative to image size)

YOLO Algorithm: IoU and Non-Maximum Suppression (NMS)

Problem: Multiple bounding boxes may be predicted for the same object.



- Choose the bounding box with the highest confidence score p_c as the best detection.
- Remove all other highly overlapping boxes with **Intersection over Union (IoU)** > 0.5 compared to the best box.

$$\text{IoU}(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

- Keep low-overlap boxes as they likely correspond to different objects (e.g., other cars).

Outline

- 1 Classic CNNs: Inception, MobileNets
- 2 Practical Advice for Using CNNs
- 3 Object Detection
- 4 Face Recognition**
- 5 Neural Style Transfer

Verification vs. Recognition

Verification

- **Input:** An image
- **Output:** Confirms if the image matches the claimed person

Recognition

- **Database:** Contains K known identities
- **Input:** An image
- **Output:** Returns the person's ID if in the database; otherwise, not recognized

One-Shot Learning and Similarity Function

Problem: Traditional deep learning models require large amounts of labeled data, but face recognition needs to identify a person with only a few examples.

One-Shot Learning

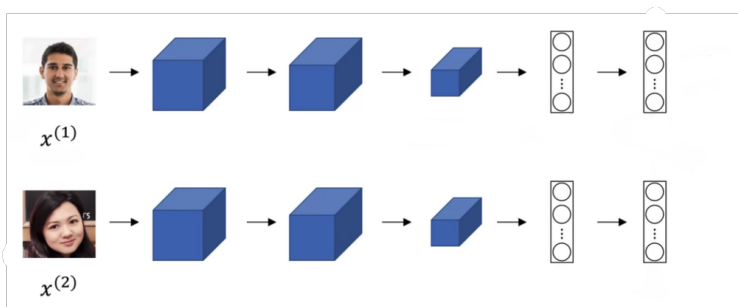
- Learns to recognize a person from a **single example** rather than relying on large datasets.
- Challenges with new identities that were not seen during training.
- Solution: Instead of classification, use a similarity function to compare images

Learning a Similarity Function

- Define a function $d(\mathbf{x}_1, \mathbf{x}_2)$ that measures the **difference** between two images
- If $d(\mathbf{x}_1, \mathbf{x}_2) \leq \text{threshold}$, classify them as the **same** person and return the identity; otherwise, reject.

Siamese Network

Definition: A Siamese network is a neural architecture consisting of **two identical** sub-networks that share weights and parameters.



- The ConvNet $f_{\theta}(x)$ represents the **embedding** (feature representation) of the input image x .
- The similarity between two images is measured using the distance between their **embeddings**:

$$d(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \|\mathbf{f}_{\theta}(\mathbf{x}^{(i)}) - \mathbf{f}_{\theta}(\mathbf{x}^{(j)})\|$$

- The model learns parameters θ such that:
 - If $\mathbf{x}^{(i)}$ and $\mathbf{x}^{(j)}$ belong to the same person, $d(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$ is small.
 - If they belong to different people, $d(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$ is large.

Triplet Loss



Anchor



Positive



Negative

- Training uses image triplets: an **anchor** a , a **positive** p , and a **negative** n .
- Objective:

$$d(\mathbf{a}, \mathbf{p}) + \alpha \leq d(\mathbf{a}, \mathbf{n})$$

where $\alpha > 0$ is the **margin**, preventing trivial solutions.

- Triplet loss:

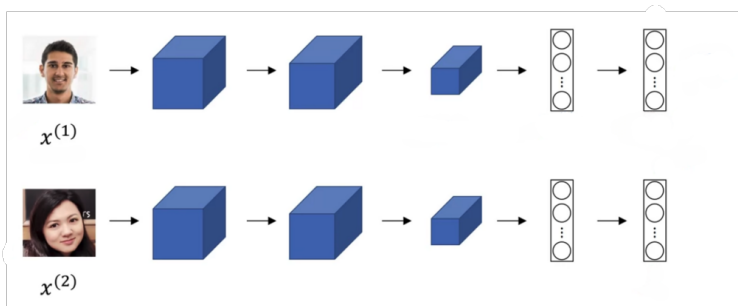
$$\ell(\mathbf{a}, \mathbf{p}, \mathbf{n}) = [d(\mathbf{a}, \mathbf{p}) + \alpha - d(\mathbf{a}, \mathbf{n})]_+,$$

where $[x]_+ = \max\{x, 0\}$.

$$\mathcal{L}(\theta) = \sum_{(\mathbf{a}, \mathbf{p}, \mathbf{n})} \ell(\mathbf{a}, \mathbf{p}, \mathbf{n})$$

- Requires multiple images per identity for effective learning.

Face Recognition via Binary Classification



- The ConvNet $f_{\theta}(x)$ extracts a feature **embedding** from the input image x .
- The similarity between two images is measured as the distance between their embeddings:

$$d(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \|\mathbf{f}_{\theta}(\mathbf{x}^{(i)}) - \mathbf{f}_{\theta}(\mathbf{x}^{(j)})\|$$

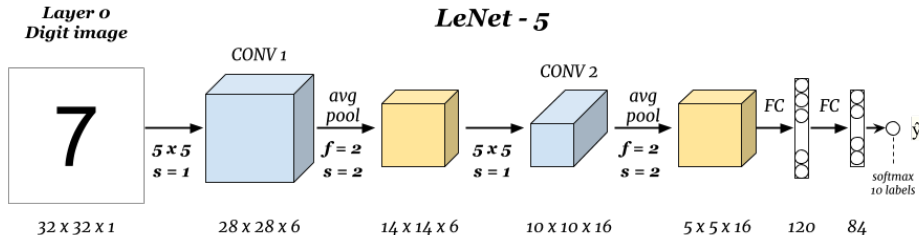
- The distance $d(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$ is passed through a **logistic activation** $\sigma(\cdot)$ to classify whether they belong to the same person:

$$\hat{y} = \sigma(d(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}))$$

Outline

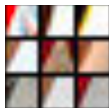
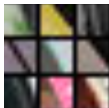
- 1 Classic CNNs: Inception, MobileNets
- 2 Practical Advice for Using CNNs
- 3 Object Detection
- 4 Face Recognition
- 5 Neural Style Transfer**

Understanding CNNs Through Visualization



How to Visualize CNN Filters:

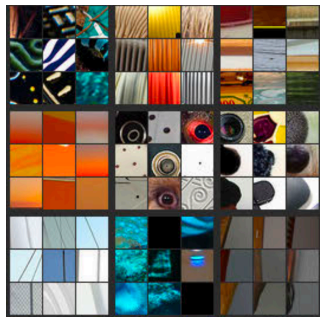
- Process many input images through a **well-trained CNN**.
- Select a **specific convolutional filter** in a layer.
- Find **9 different images** where this filter had the **highest activations** in the feature maps.
- Use DeconvNet to **reconstruct** the input regions responsible for these activations.
- This produces **9 reconstructed patches**, revealing what pattern this filter detects.



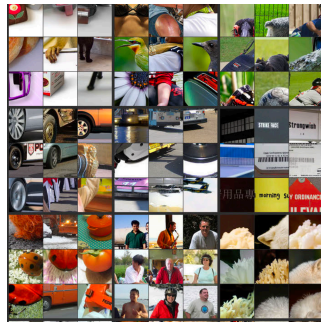
Visualizing Deep Layers



Layer 1



Layer 2



Layer 3

Neural Style Transfer



Content



Style



Transferred

- Randomly initialize an initial image g
- Optimize g using gradient descent to minimize the total loss:

$$\mathcal{L}(g) = \alpha \mathcal{L}_{\text{content}}(g, c) + \beta \mathcal{L}_{\text{style}}(g, s),$$

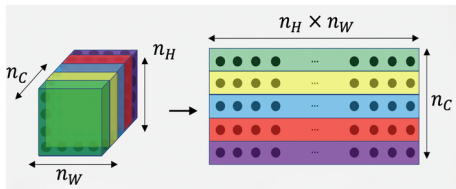
where

$$\mathcal{L}_{\text{content}}(g, c) = \sum_{\ell} \|\mathbf{a}^{\ell}(g) - \mathbf{a}^{\ell}(c)\|^2$$

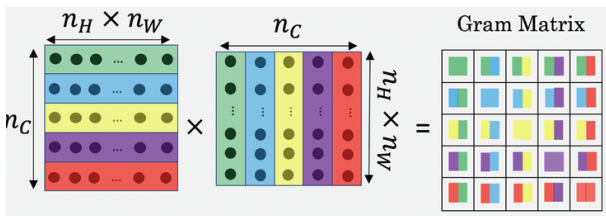
Here, \mathbf{a}^{ℓ} represents the activations of layer ℓ in a pre-trained CNN, e.g., VGG.

Style Matrix

- The **style** is defined as how correlated the activations are across different **channels**.



- The style matrix $\mathbf{G}^\ell \in \mathbb{R}^{n_c \times n_c}$ is defined by: $G_{k\bar{k}}^\ell := \langle \mathbf{a}^k, \mathbf{a}^{\bar{k}} \rangle, \forall k, \bar{k} \in [n_c]$.



- Then the style cost is

$$\mathcal{L}_{\text{style}}(\mathbf{g}, \mathbf{s}) = \sum_{\ell} \frac{1}{n \times n \times n_c} \|\mathbf{G}^\ell(\mathbf{g}) - \mathbf{G}^\ell(\mathbf{s})\|^2$$