

Задачи CV

Шабалин Евгений

Что такое CV

Computer Vision - область компьютерных наук, включающая в себя теорию и технологию создания машин, которые могут анализировать изображения или видео (много изображений подряд).

Задачи CV

С помощью компьютерного зрения люди находят какие-либо паттерны, закономерности и особенности на изображениях.

Выделяют несколько различных разделов.

Computer Vision tasks

Semantic Segmentation



GRASS, CAT,
TREE, SKY

No objects, just pixels

**Classification
+ Localization**



CAT

Single Object

**Object
Detection**



DOG, DOG, CAT

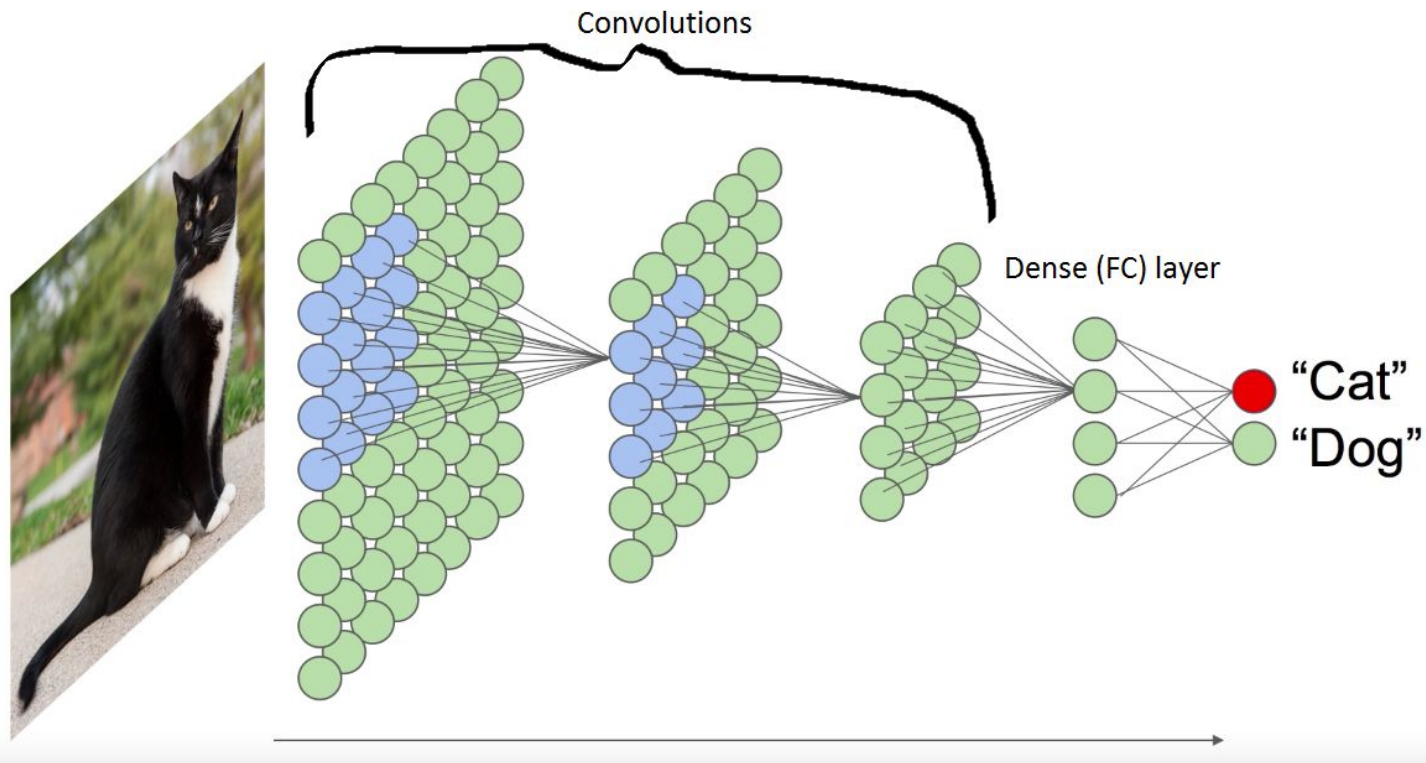
Multiple Object

**Instance
Segmentation**



DOG, DOG, CAT

Классификация



Сегментация

Сегментация — процесс деления изображения на несколько областей (сегментов). Это позволяет упростить изображение, чтобы его было легче анализировать. Можно понимать сегментацию как попиксельную классификацию

Сегментация

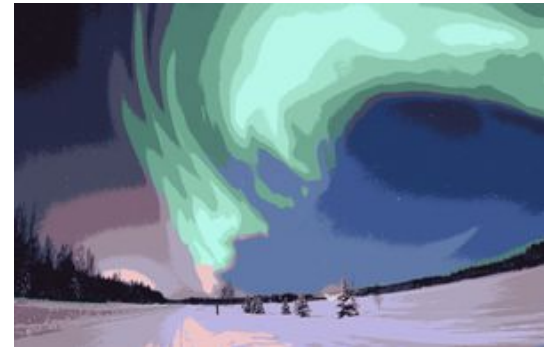
Без нейронных сетей:

- Кластеризация
- Гистограммы
- Выделение границ
- Метод разрастания
- Деление графа

Сегментация

Без нейронных сетей:

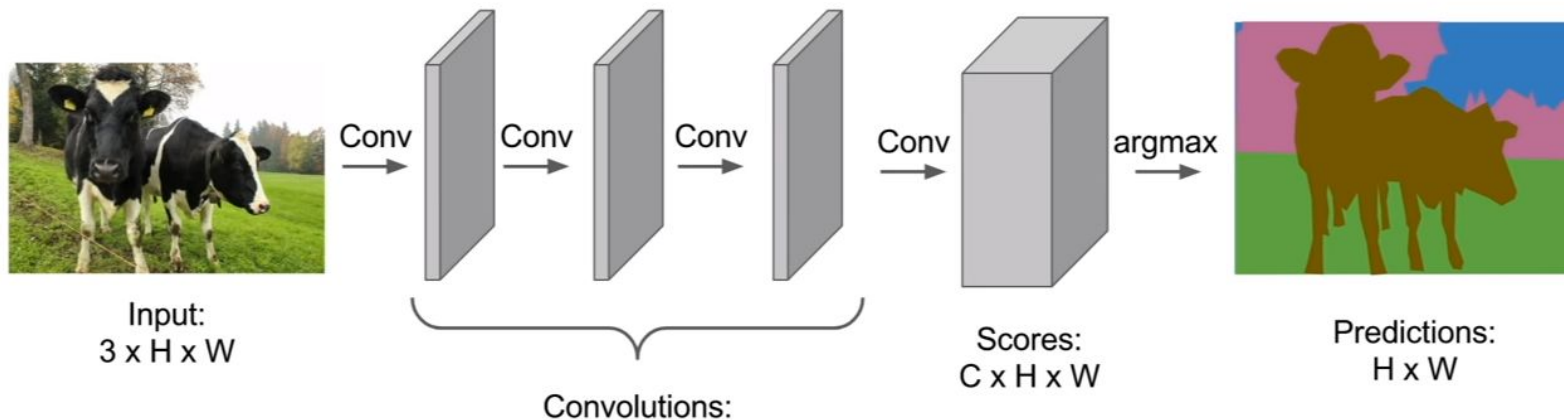
- Кластеризация
- Гистограммы
- Выделение границ
- Метод разрастания
- Деление графа



Сегментация

С использованием нейронок:

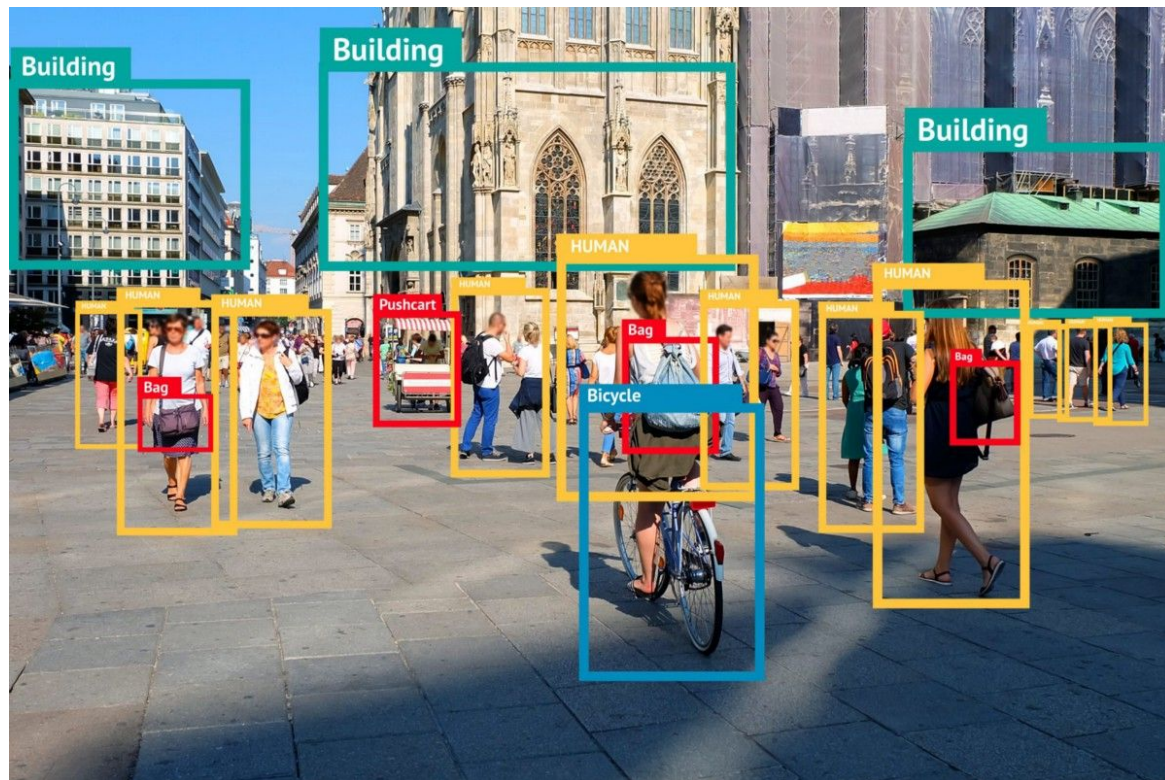
Fully Convolutional Network (FCN)



Детектирование

Задача детектирования заключается в поиске заданных объектов на изображении и их выделении.

Детектирование



Детектирование на видео

Проблема: картинок очень много и обрабатывать отдельно каждую слишком долго.

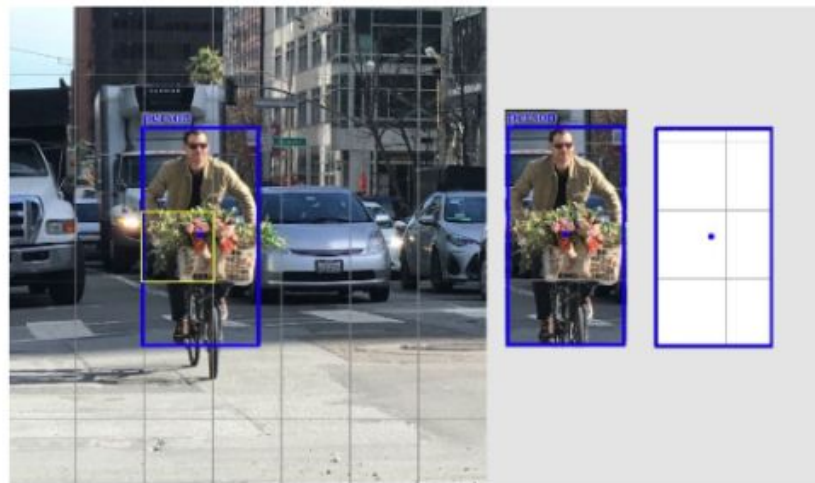
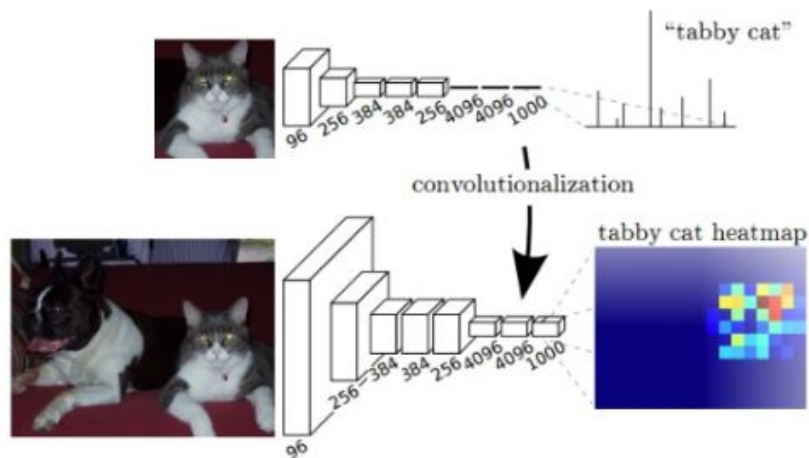
Что делать?

Детектирование на видео

Проблема: картинок очень много и обрабатывать отдельно каждую слишком долго.

Есть 2 варианта: one-shot и two-shot. Первый неточный, второй медленный

One shot detectors



Для каждой ячейки в последнем conv слое предсказываем координаты бокса и класс объекта с центром в ячейке.

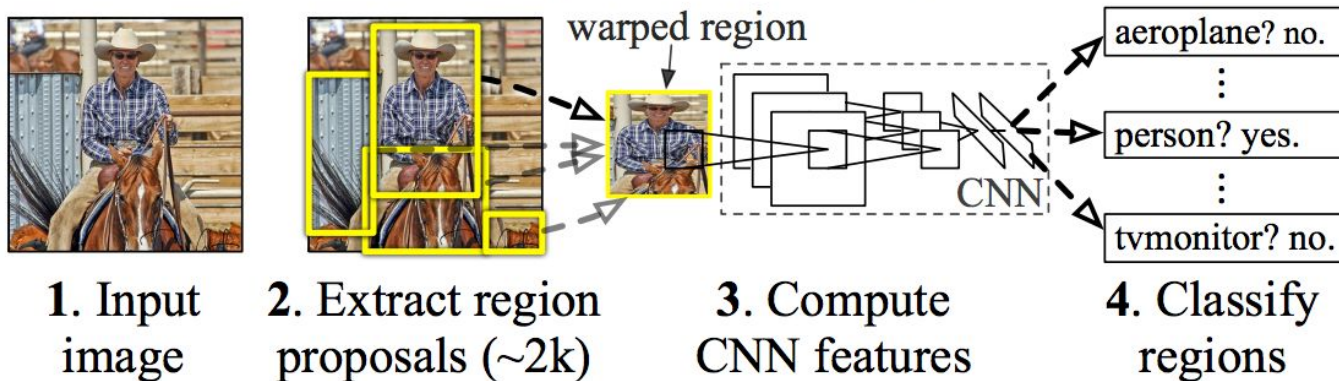
R-CNN

Подход относительно прост: возьмем существующее решение (AlexNet, например) и изменим входные данные. С помощью selective search сгенерируем много различных регионов, в которых может находиться объект и будем смотреть отдельно на них.

R-CNN = Selective Search
+ Classification



R-CNN: *Regions with CNN features*



R-CNN

Хорошо? Не очень

Получились хорошие результаты, но с очень большими затратами по времени (47 секунд на одну картинку для VGG16).

Fast/Faster R-CNN

Fast R-CNN: добавим region of interest (RoI) pooling вместо обычного

Faster R-CNN: вместо определения регионов по изначальному изображению будем это делать после CNN слоев

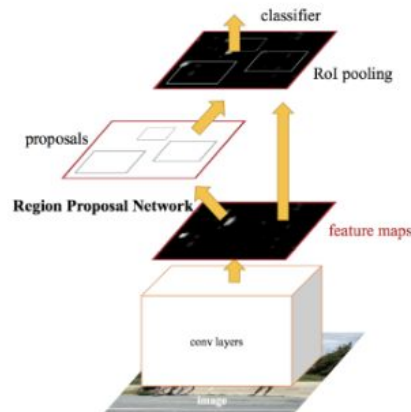
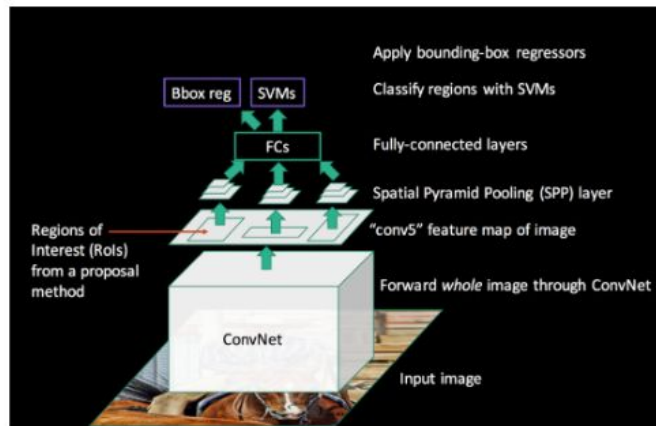
Fast/Faster R-CNN

Fast => Faster

Fast R-CNN => Faster R-CNN

Вычисляем proposals самой сетью.

2 секунды => 0.2 секунды (10 раз быстрее)

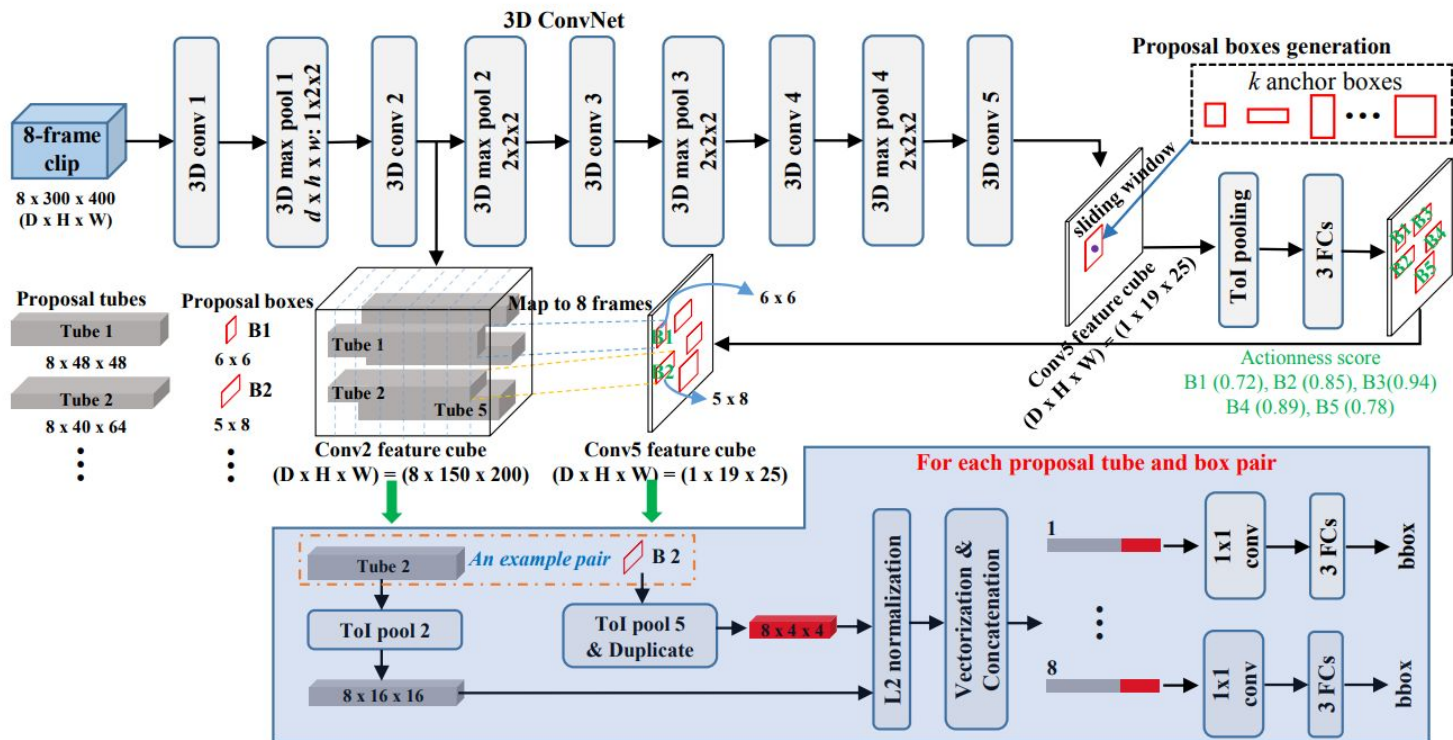


T-CNN

Все предыдущие подходы были разработаны для статичных изображений и никак не использовали факт, что в близких по времени картинках особенности будут расположены в похожих местах.

Взяли все подходы из Faster R-CNN и адаптируем их для клипов по 8 кадров (Tubes).

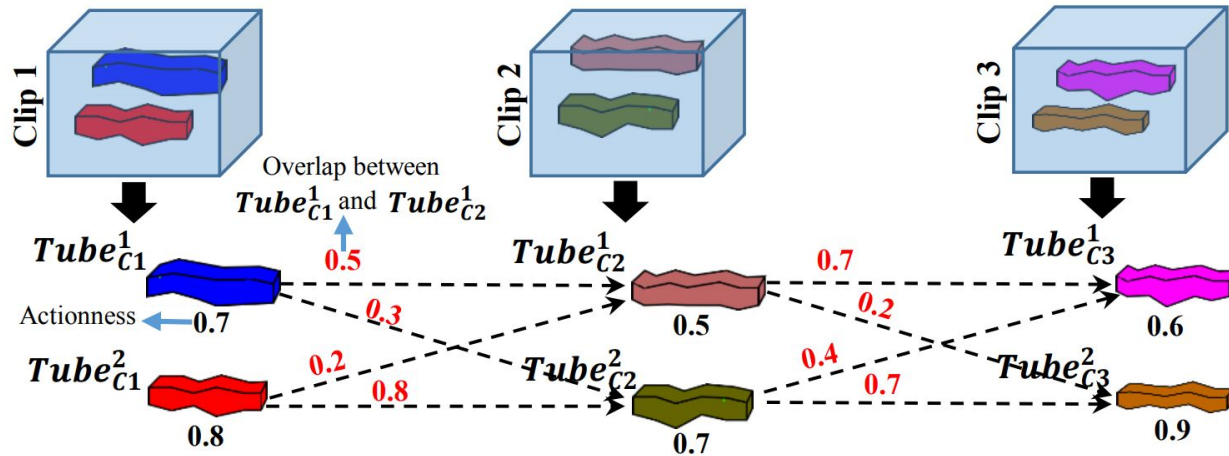
T-CNN



T-CNN

$$S = \frac{1}{m} \sum_{i=1}^m \textit{Actionness}_i + \frac{1}{m-1} \sum_{j=1}^{m-1} \textit{Overlap}_{j,j+1} \quad (2)$$

T-CNN

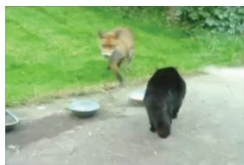


$$\begin{aligned}
 S1 &= Tube_{C1}^1 \rightarrow Tube_{C2}^1 \rightarrow Tube_{C3}^1 = 1.8/3 + 1.2/2 & S2 &= Tube_{C1}^1 \rightarrow Tube_{C2}^1 \rightarrow Tube_{C3}^2 = 2.1/3 + 0.7/2 \\
 S3 &= Tube_{C1}^1 \rightarrow Tube_{C2}^2 \rightarrow Tube_{C3}^1 = 2.0/3 + 0.7/2 & S4 &= Tube_{C1}^1 \rightarrow Tube_{C2}^2 \rightarrow Tube_{C3}^2 = 2.3/3 + 1.0/2 \\
 S5 &= Tube_{C1}^2 \rightarrow Tube_{C2}^1 \rightarrow Tube_{C3}^1 = 1.9/3 + 0.9/2 & S6 &= Tube_{C1}^2 \rightarrow Tube_{C2}^1 \rightarrow Tube_{C3}^2 = 2.2/3 + 0.4/2 \\
 S7 &= Tube_{C1}^2 \rightarrow Tube_{C2}^2 \rightarrow Tube_{C3}^1 = 2.1/3 + 1.2/2 & S8 &= Tube_{C1}^2 \rightarrow Tube_{C2}^2 \rightarrow Tube_{C3}^2 = 2.4/3 + 1.5/2
 \end{aligned}$$

Seq-NMS

Toy example video

t=0



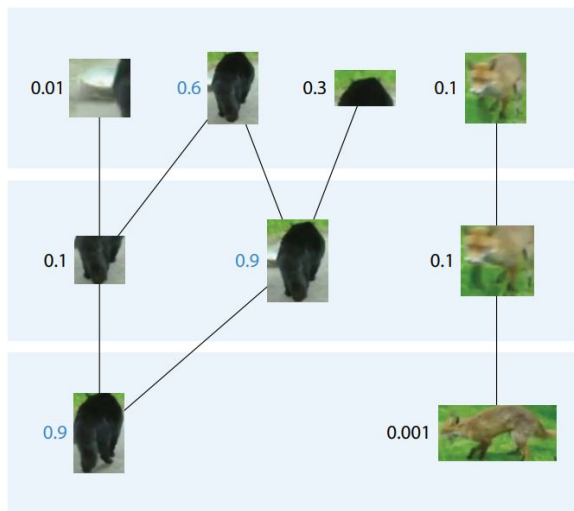
t=1



t=2



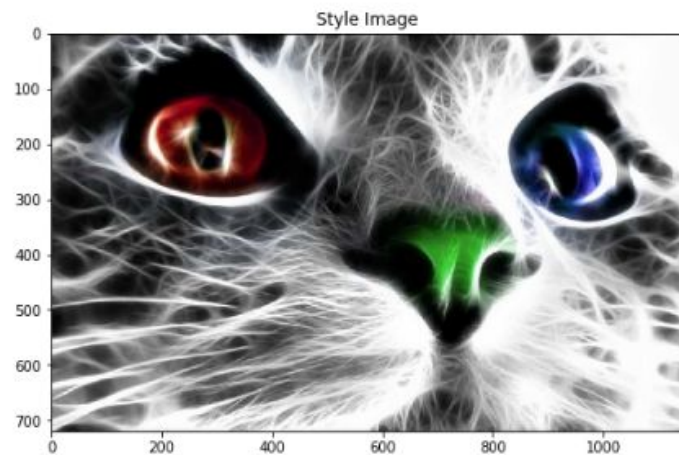
Build inter-frame IoU graph



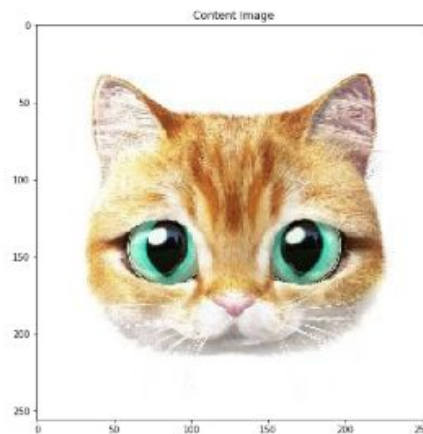
After Seq-NMS rescoring



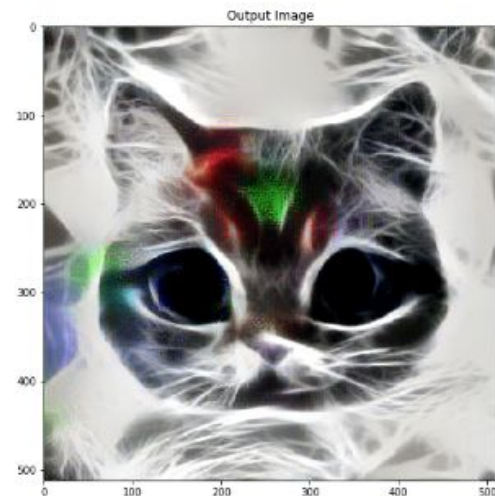
Перенос стиля



+



=



Перенос стиля

Как работает?

Перенос стиля

Начнем с контента

$$\mathcal{L}_{\text{content}}(\vec{p}, \vec{x}, l) = \frac{1}{2} \sum_{i,j} (F_{ij}^l - P_{ij}^l)^2 .$$

\vec{p} - исходная картинка, \vec{x} - то, что мы хотим получить

Перенос стиля

Теперь стиль

$$G_{ij}^l = \sum_k F_{ik}^l F_{jk}^l.$$

Перенос стиля

Теперь стиль

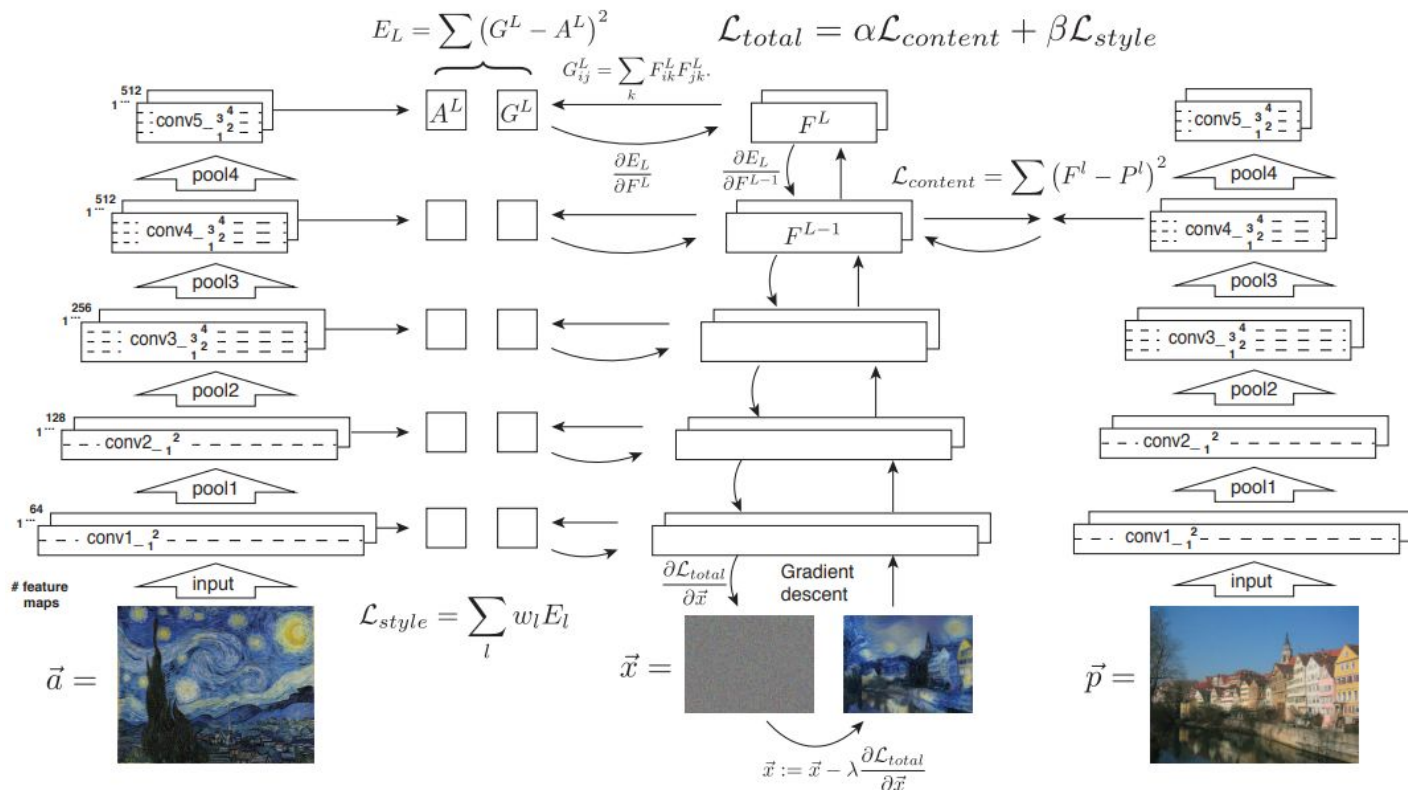
$$G_{ij}^l = \sum_k F_{ik}^l F_{jk}^l.$$

$$E_l = \frac{1}{4N_l^2 M_l^2} \sum_{i,j} (G_{ij}^l - A_{ij}^l)^2$$

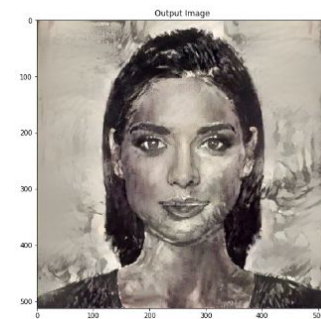
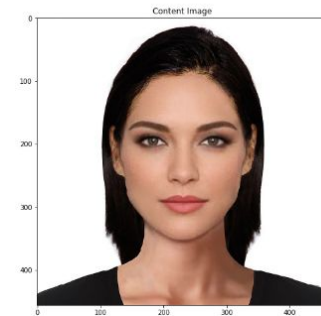
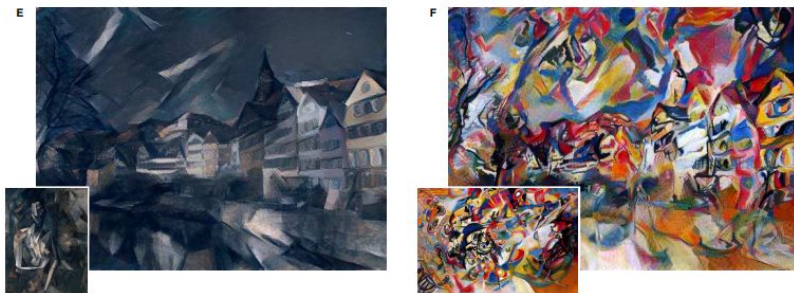
and the total style loss is

$$\mathcal{L}_{\text{style}}(\vec{a}, \vec{x}) = \sum_{l=0}^L w_l E_l,$$

Перенос стиля



Результаты



Что нас ждет

Вопросы

- Чем отличается сегментация от детектирования
- Что такое RoI и как работает
- Как создается представления контента и стиля в задаче переноса стиля и почему по-разному

ИСТОЧНИКИ

- R-CNN <https://arxiv.org/abs/1311.2524>
- Fast R-CNN <https://arxiv.org/abs/1504.08083>
- Faster R-CNN <https://arxiv.org/abs/1506.01497>
- T-CNN <https://arxiv.org/pdf/1703.10664.pdf>
- Style transfer

https://www.cv-foundation.org/openaccess/content_cvpr_2016/papers/Gatys_Image_Style_Transfer_CVPR_2016_paper.pdf