

Introduction

Featured Prediction Competition

2018 Data Science Bowl

Find the nuclei in divergent images to advance medical discovery

Booz Allen Hamilton · 739 teams · 8 months ago

DATA SCIENCE BOWL
Passion. Curiosity. Purpose.

Presented by
Booz | Allen | Hamilton & kaggle

细胞核位置检测比赛，精确找出细胞核的位置，用一个mask定位



Spot Nuclei. Speed Cures.

The challenge: Create an algorithm to automate nucleus detection

40%
of all deaths are caused by illnesses like heart disease and cancer¹

75%
of rare diseases affect children²

30%
of affected children with rare diseases die before age 5²

Finding the nucleus helps to...

- locate cells in varied conditions to enable faster cures
- free biologists to focus on solutions
- improve throughput for research and insight
- reduce time-to-market for new drugs— currently 10 years
- increase # of compounds for experiments
- improve health and increase quality of life

Nuclei are distinctive in images and can help researchers locate cells

Nuclei take many shapes across the body's 30 trillion cells

SOURCES:

¹Heart disease and cancer causing 40% of all deaths – World Health Organization (WHO). 2015 – 56.4 million deaths. <http://www.who.int/mediacentre/factsheets/fs310/en/>
2015 – 17.7 deaths from cardiovascular diseases (CVDs). <http://www.who.int/mediacentre/factsheets/fs317/en/>
2015 – 8.8 million deaths from cancer. <http://www.who.int/mediacentre/factsheets/fs297/en/>

²Childhood diseases - 75% of rare diseases affect children; 30% die before age 5 – European Society of Paediatric Oncology (SIOPe). <http://www.siope.eu/SIOPE-EU/English/SIOPE-EU/Advocacy-Activities/Rare-Diseases/page.aspx/148>

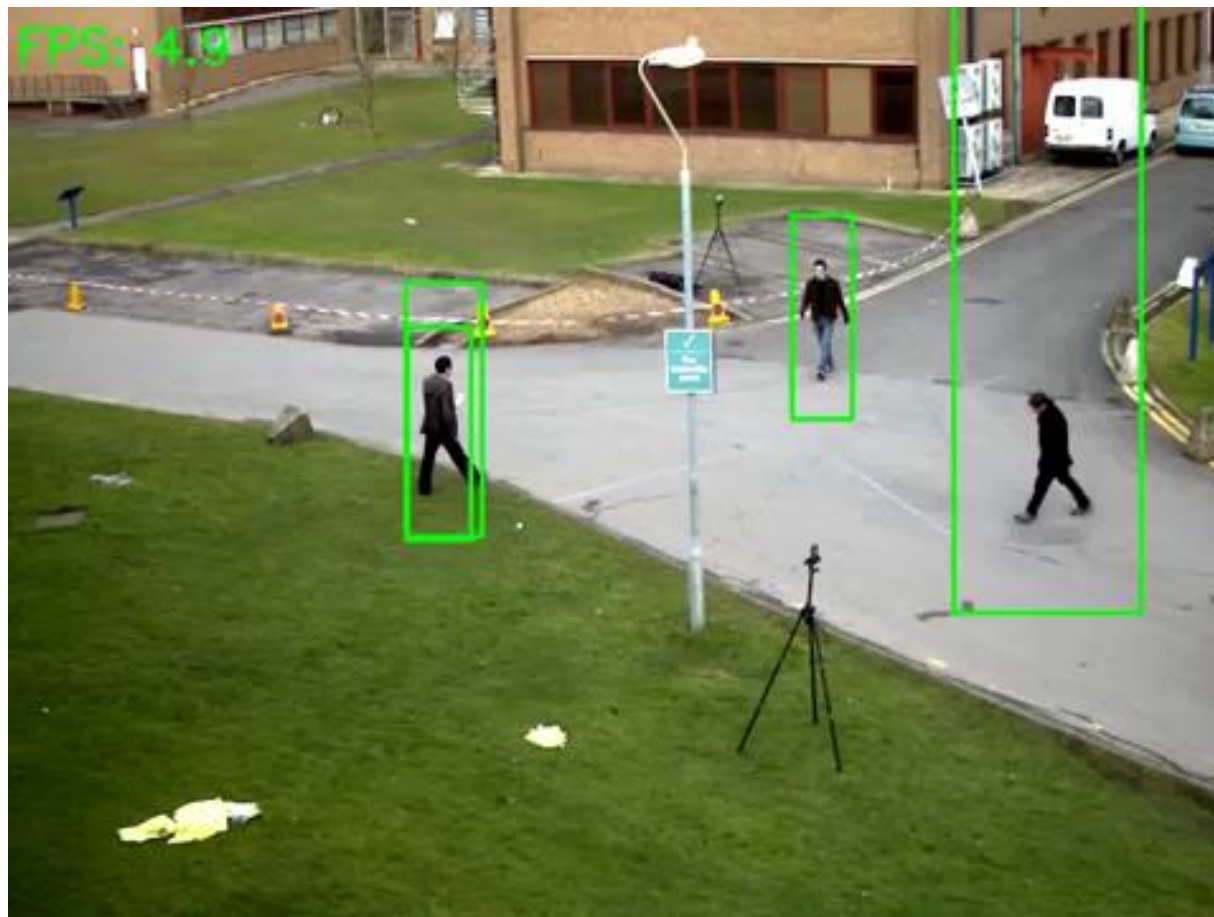
为什么要做图像分割(image segmentation?)

相关应用

- 图像分割
- 之后对提

1. 机器视觉
2. 人脸识别
3. 指纹识别
4. 交通控制系统
5. 在卫星图像中定位物体（道路、森林等）
6. 行人检测
7. 医学影像，包括:
 - （1）肿瘤和其他病理的定位
 - （2）组织体积的测量
 - （3）计算机引导的手术
 - （4）诊断

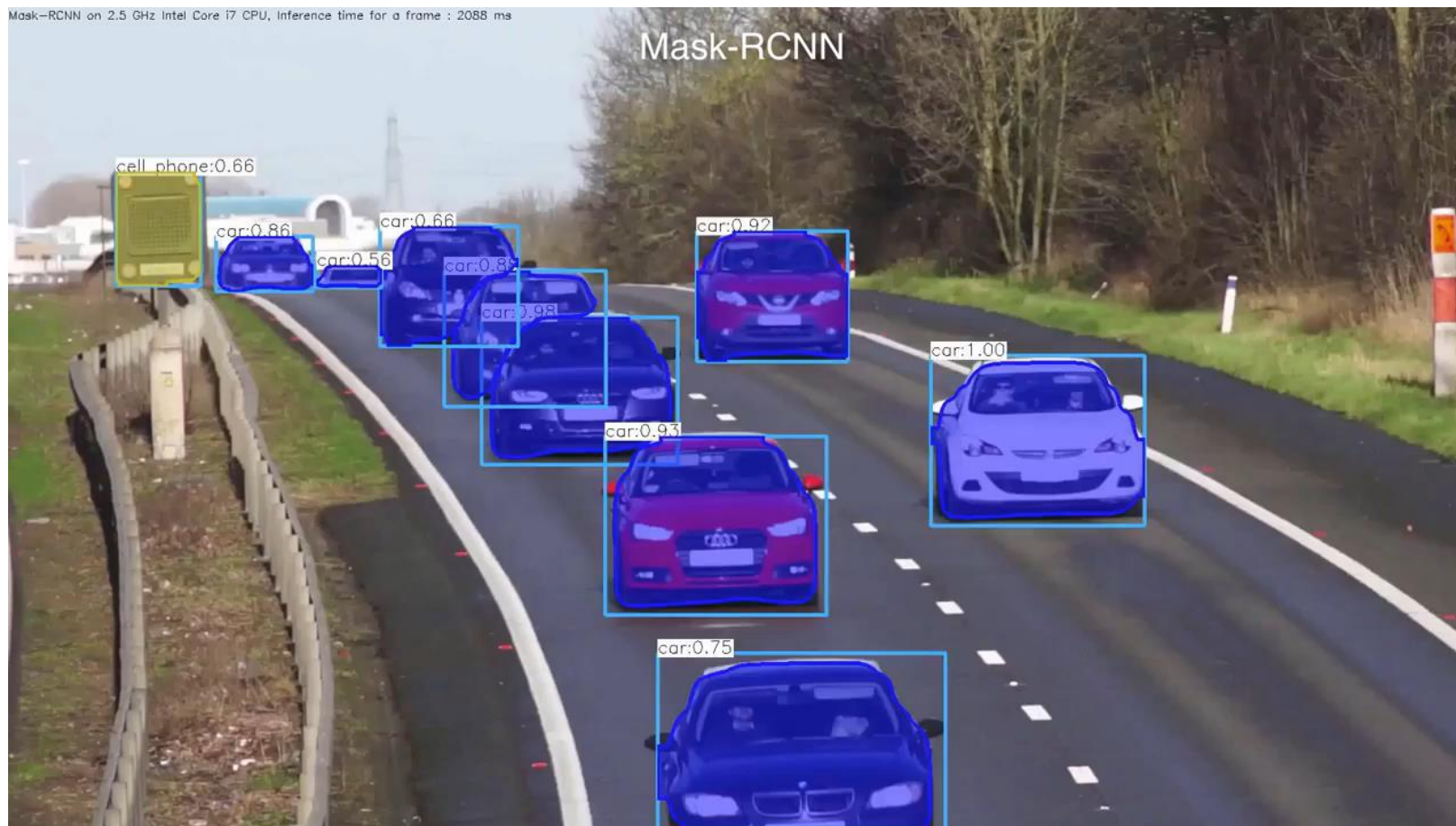
传统图像分割实例 (HOG行人检测)



DL图像分割实例(YOLO)



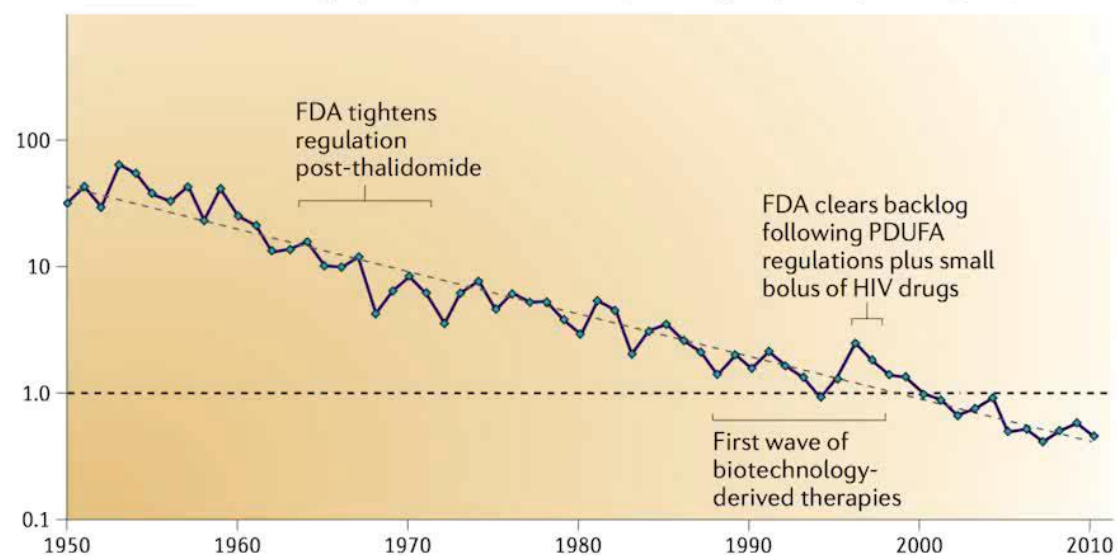
DL图像分割实例(像素级别, Mask RCNN)



医学图像分割的应用

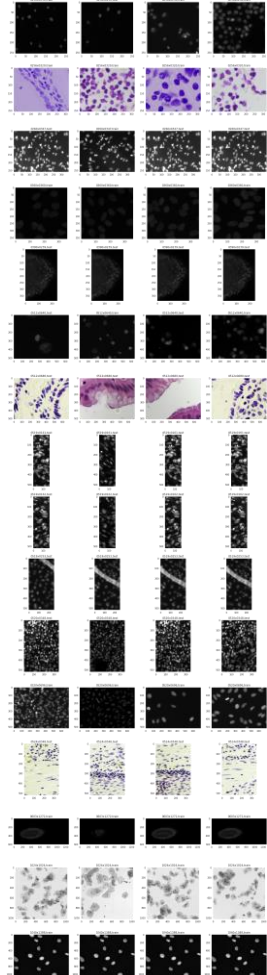
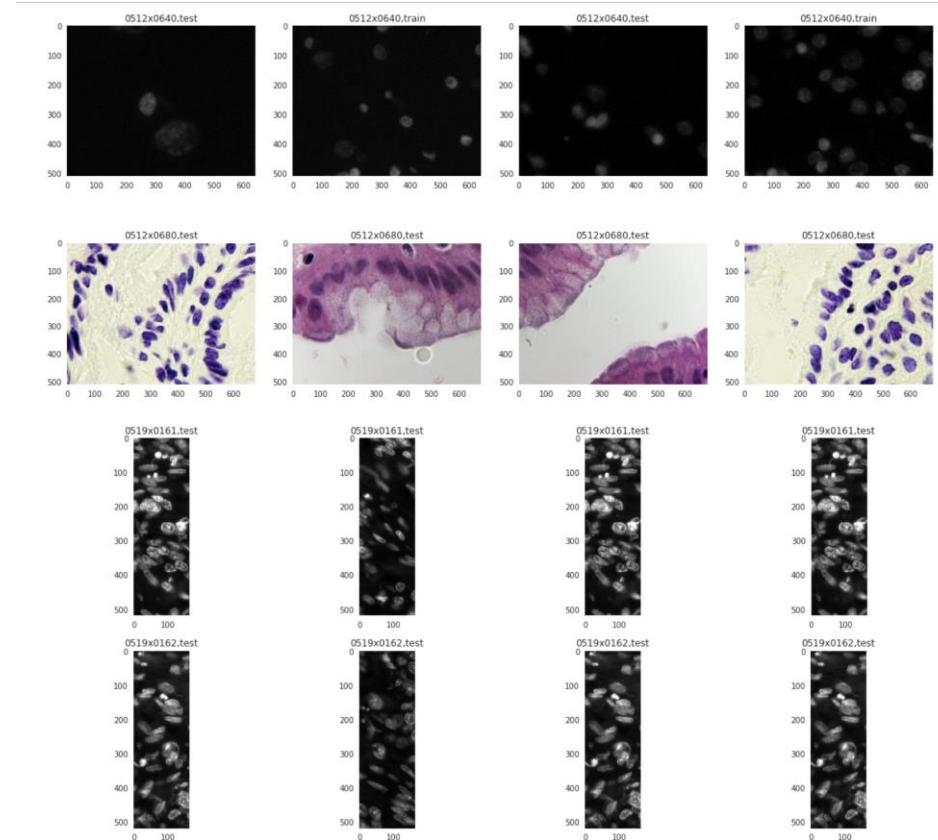
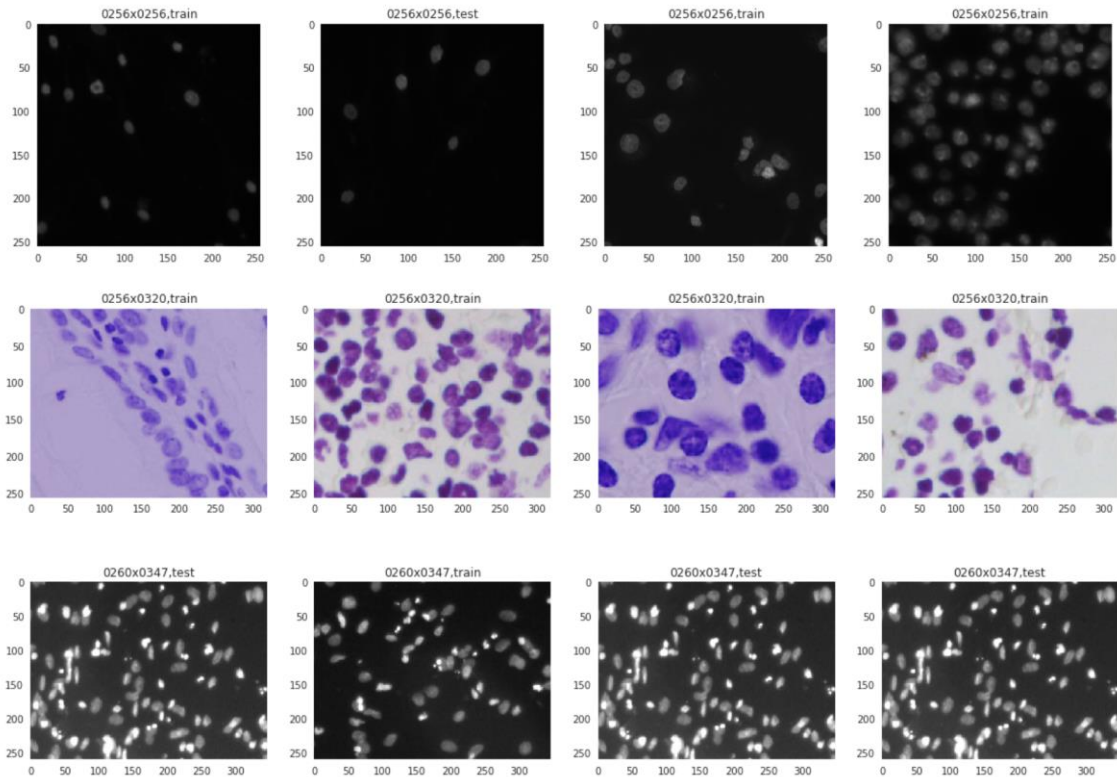
Drug discovery is increasingly difficult:
Eroom's law

Number of drugs per \$1 billion R&D spending (\$USD, inflation-adjusted)



Before Modeling-Exploratory data Analysis

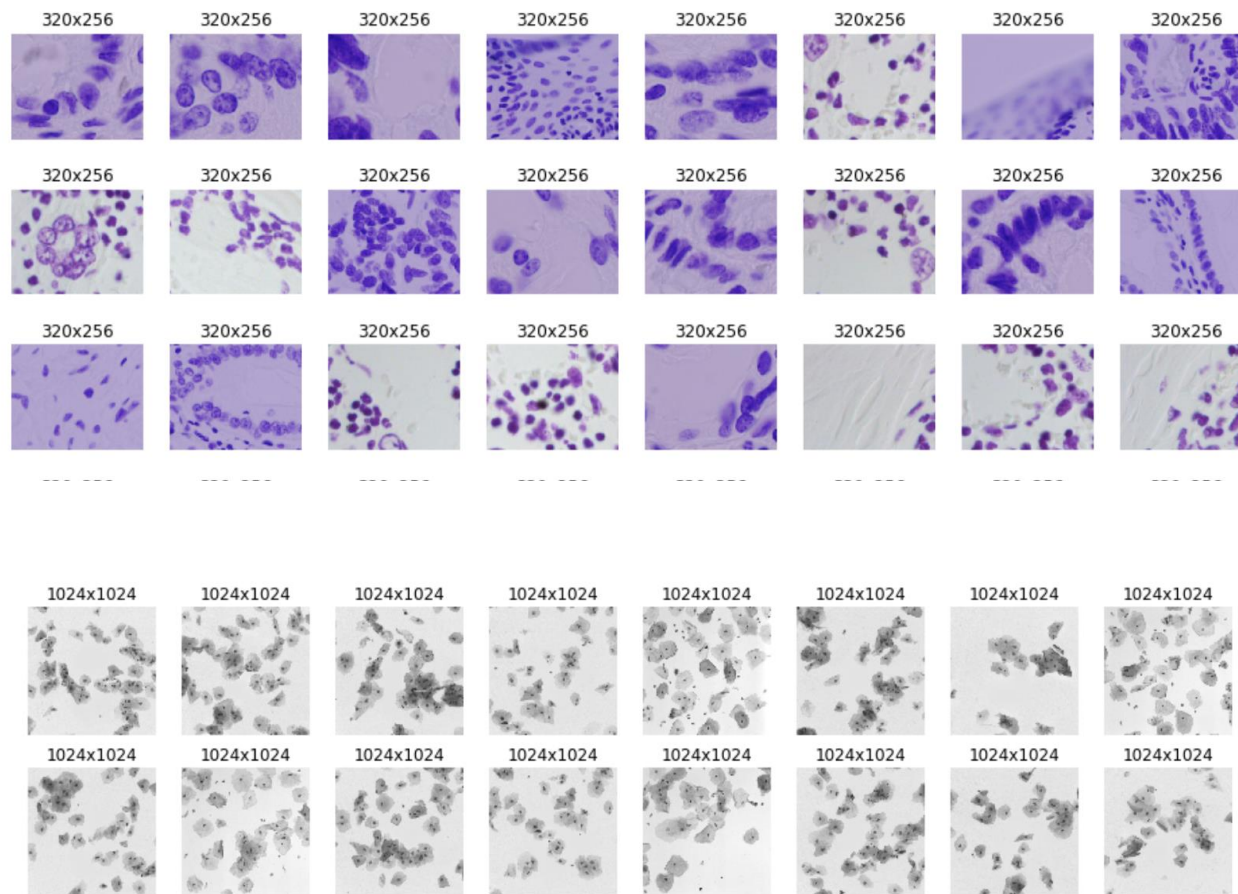
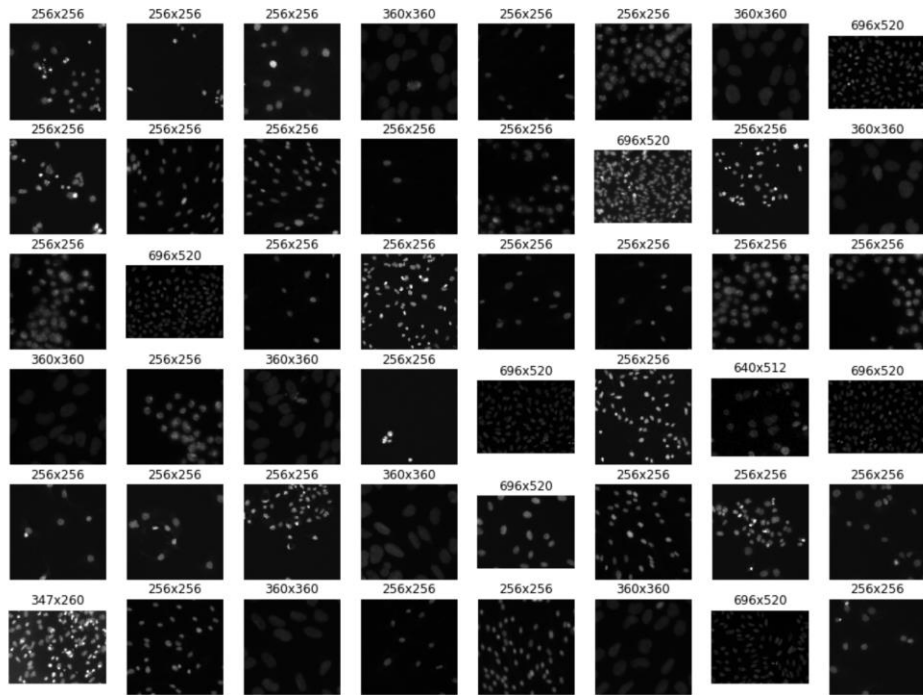
```
sample['label'] = sample[['shape', 'train_or_test']].apply(lambda x: '{}'.format(x[0], x[1]), axis=1)
show_row_col(sample, 4, path_col='path', label_col='label', mode='file')
```



Before Modeling-Exploratory data Analysis

将图片从RGB空间转换至HSV空间，KMean进行聚类

```
plot_images(trainPD[trainPD[HSV_CLUSTER] == 0][IMAGE_ID].values, 6, 8)
```



Before Modeling-Evaluation Metrics

the mean average precision at different intersection over union (IoU) thresholds
在不同交并比阈值下的平均精度

This tells us there are a few different steps to getting the score reported on the leaderboard. For each image...

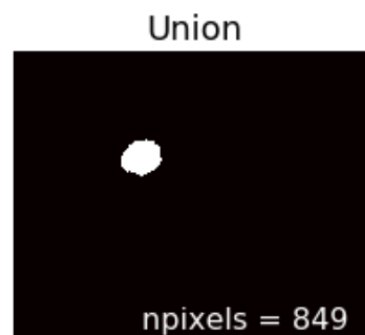
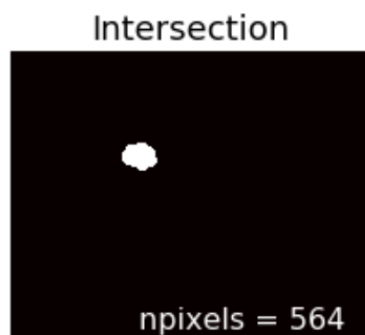
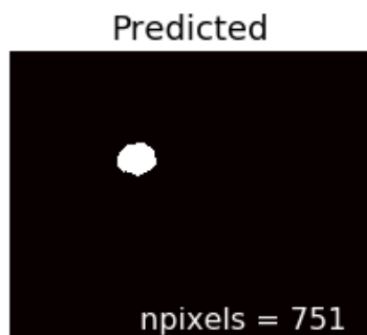
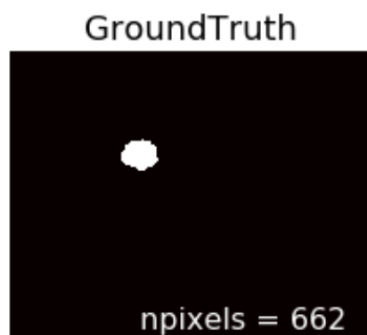
1. For each submitted nuclei "prediction", calculate the Intersection of Union metric with each "ground truth" mask in the image.
2. Calculate whether this mask fits at a range of IoU thresholds.
3. At each threshold, calculate the precision across all your submitted masks.
4. Average the precision across thresholds.

Across the dataset...

1. Calculate the mean of the average precision for each image.

Before Modeling-Evaluation Metrics

IOU 交并比



$$IoU(A, B) = \frac{A \cap B}{A \cup B} = \frac{564}{849} = 0.664$$

Does this IoU hit at each threshold?

| | |
|------|-------|
| 0.50 | True |
| 0.55 | True |
| 0.60 | True |
| 0.65 | True |
| 0.70 | False |
| 0.75 | False |
| 0.80 | False |
| 0.85 | False |
| 0.90 | False |
| 0.95 | False |

Name: GT-P, dtype: bool

Next, we sweep over a range of IoU thresholds to get a vector for each mask comparison. The threshold values range from 0.5 to 0.95 with a step size of 0.05: `(0.5, 0.55, 0.6, 0.65, 0.7, 0.75, 0.8, 0.85, 0.9, 0.95)`.

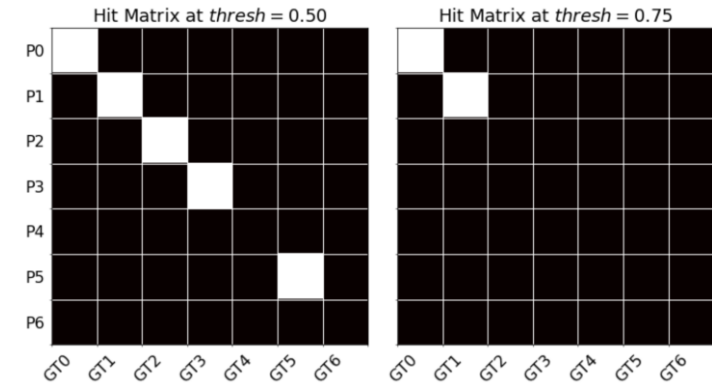
Before Modeling-Evaluation Metrics

$$Precision(t) = \frac{TP(t)}{TP(t) + FP(t) + FN(t)}$$

- TP: a single predicted object matches a ground truth object with an IoU above the threshold
- FP: a predicted object had no associated ground truth object.
- FN: a ground truth object had no associated predicted object.

In the above matrix...

- The number of **true positives** is equal to the number of predictions with a "hit" on a true object.
- The number of **false positives** is equal to the number of predictions that don't hit anything.
- The number of **false negatives** is equal to the number of "ground truth" objects that aren't hit.



Precision values at each threshold:

$t(0.50) = 0.556$

$t(0.55) = 0.400$

$t(0.60) = 0.400$

$t(0.65) = 0.400$

$t(0.70) = 0.167$

$t(0.75) = 0.167$

$t(0.80) = 0.077$

$t(0.85) = 0.077$

$t(0.90) = 0.000$

$t(0.95) = 0.000$

Mean precision for image is: 0.224

At a threshold of 0.50...

TP = 5

FP = 2

FN = 2

p = 0.556

At a threshold of 0.75...

TP = 2

FP = 5

FN = 5

p = 0.167

Code

Preprocess-Resize&Scaling&DataAug

网络需要一个统一维度的输入

```
img = resize(img, (256, 256), mode='constant', preserve_range=True)
```

数据归一加速训练过程

```
In [15]: 1 inputs = Input((input_H, input_W, channel))
          2 c = Lambda(lambda x: x/255)(inputs)
```

通过 平移 旋转 裁剪 镜像 翻转 缩放的方式来增强数据(增加可训练图片)

```
In [23]: 1 datagen = ImageDataGenerator(featurewise_center = False, samplewise_center = False, featurewise_std_normalization = False,
          2                               samplewise_std_normalization = False,
          3                               rotation_range = 90,
          4                               width_shift_range = 0.2,
          5                               height_shift_range = 0.2,
          6                               shear_range = 0.1,
          7                               zoom_range = 0.2,
          8                               horizontal_flip = True,
          9                               vertical_flip = True,
          10                              fill_mode = 'reflect')
```


Preprocess-DataAugValidation&OptTarget

增强数据后训练的效果不一定总是有益的 需要验证有效的增强，通过在训练时分割训练测试集比例

```
results = model.fit(X_train, Y_train, validation_split=0.1)
```

对于输出的mask 每个mask判断该像素为属于细胞核或不属于细胞核是一个二分类问题
所以最后一层选sigmoid function 优化目标为交叉熵代价函数

$$C = -\frac{1}{n} \sum_x [y \ln a + (1 - y) \ln(1 - a)]$$

```
model.compile(optimizer='adam', loss='binary_crossentropy', metrics=[mean_iou])
```

U-Net

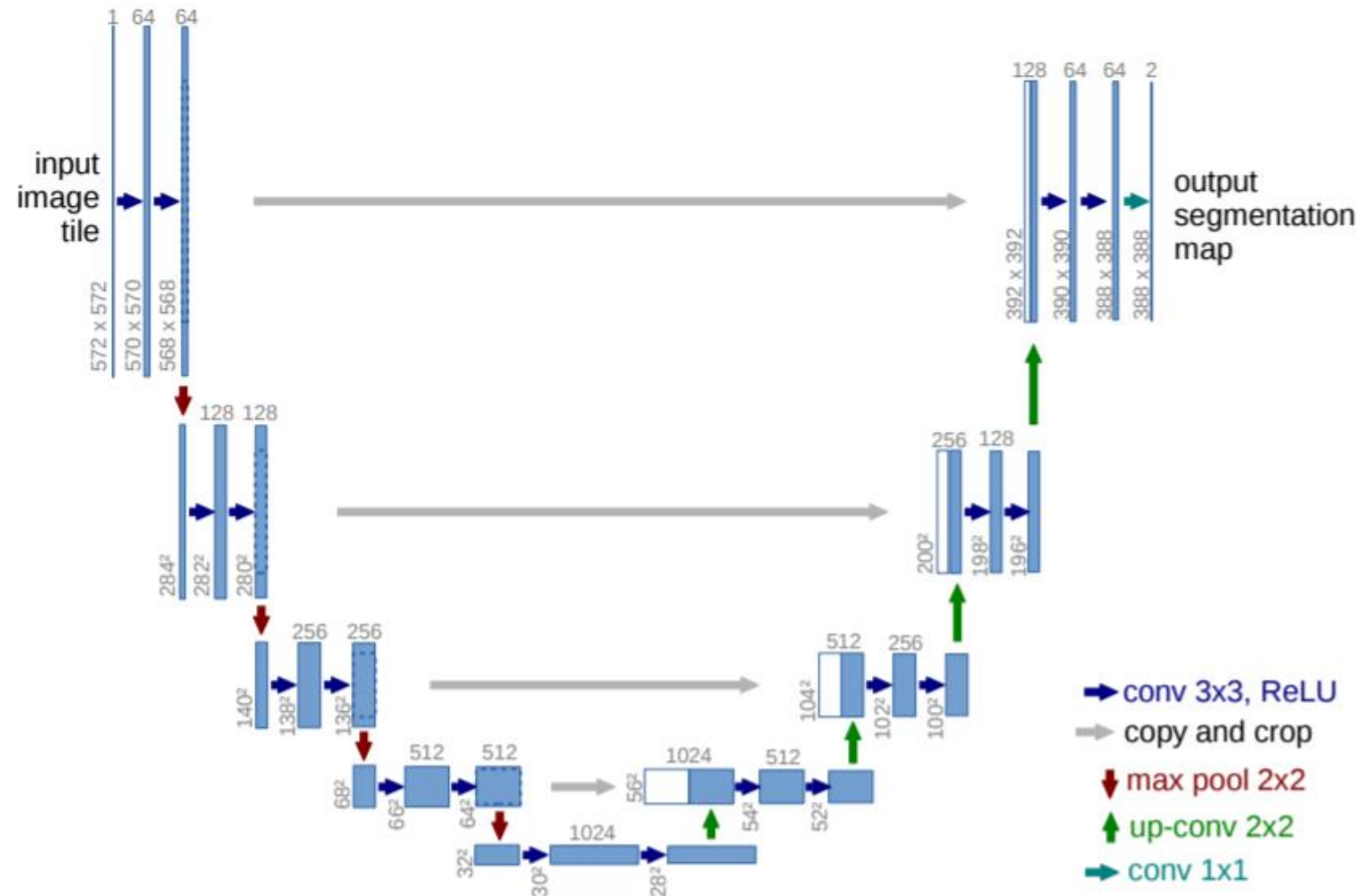
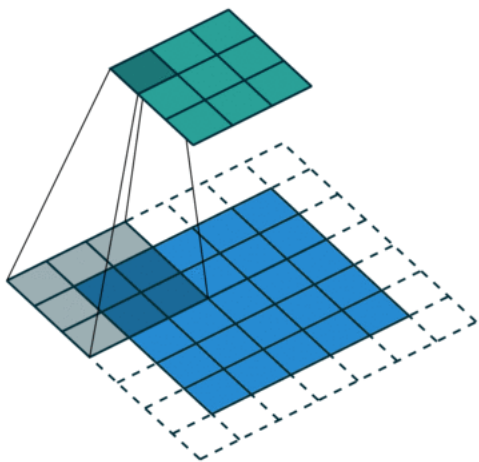


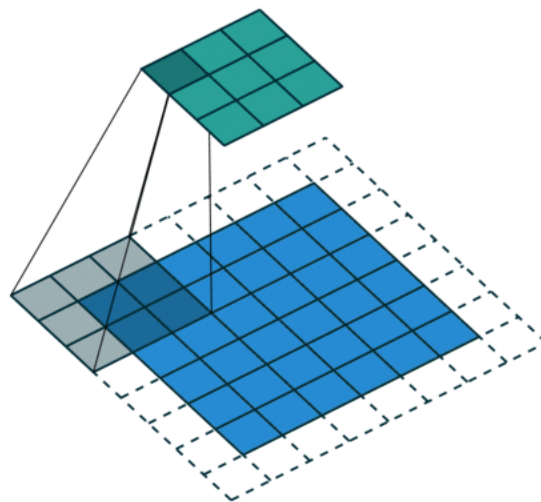
Fig. 1. U-net architecture (example for 32x32 pixels in the lowest resolution). Each blue box corresponds to a multi-channel feature map. The number of channels is denoted on top of the box. The x-y-size is provided at the lower left edge of the box. White boxes represent copied feature maps. The arrows denote the different operations.

Convolution & Deconvolution

下图表示参数为 ($i=5, k=3, s=2, p=1$) 的卷积计算过程, 从计算结果可以看出输出特征的尺寸为 ($o=3$)



下图表示参数为 ($i=6, k=3, s=2, p=1$) 的卷积计算过程, 从计算结果可以看出输出特征的尺寸为 ($o=3$)。

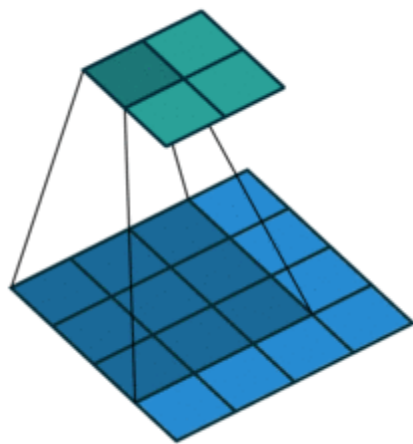
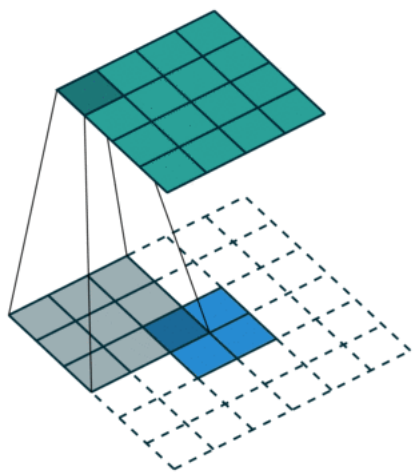


$$o = \left\lfloor \frac{i + 2p - k}{s} \right\rfloor + 1.$$

Convolution & Deconvolution

下图表示的是参数为($i'=2, k'=3, s'=1, p'=2$)的反卷积操作，其对应的卷积操作参数为($i=4, k=3, s=1, p=0$)。我们可以发现对应的卷积和非卷积操作其($k=k', s=s'$)，但是反卷积却多了 $p'=2$ 。通过对比我们可以发现卷积层中左上角的输入只对左上角的输出有贡献，所以反卷积层会出现 $p'=k-p-1=2$ 。通过示意图，我们可以发现，反卷积层的输入输出在 $s=s'=1$ 的情况下关系为：

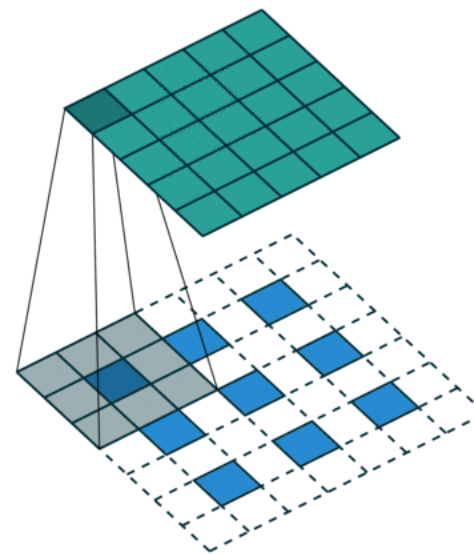
$$o' = i' - k' + 2p' + 1 = i' + (k - 1) - 2p$$



一般的，在实际应用中会采取以下方式进行反卷积

```
from keras.layers.convolutional
import Conv2D, Conv2DTranspose

Conv2DTranspose(1, (3, 3), strides=1)
```



Summary

- 图像检测&图像分割相关技术的介绍
- 图像分割应用：kaggle竞赛-细胞核分割
 - 竞赛前的数据探索
 - 神经网络训练前数据的预处理
 - U-Net介绍及Keras实现
 - 通过预训练的res-50进行迁移学习

Rethinking ImageNet Pre-training

Kaiming He Ross Girshick Piotr Dollár

Facebook AI Research (FAIR)

对于训练集数据较少的情况下，U-Net可以获得比较好的效果。