

# 深層學習入門

## 2. 深層學習

# 深層学習

---

ニューラルネットワーク

畳み込みニューラルネットワーク

物体分類

物体検出・セグメンテーション

敵対的生成ネットワーク

# 深層学習

---

## ニューラルネットワーク

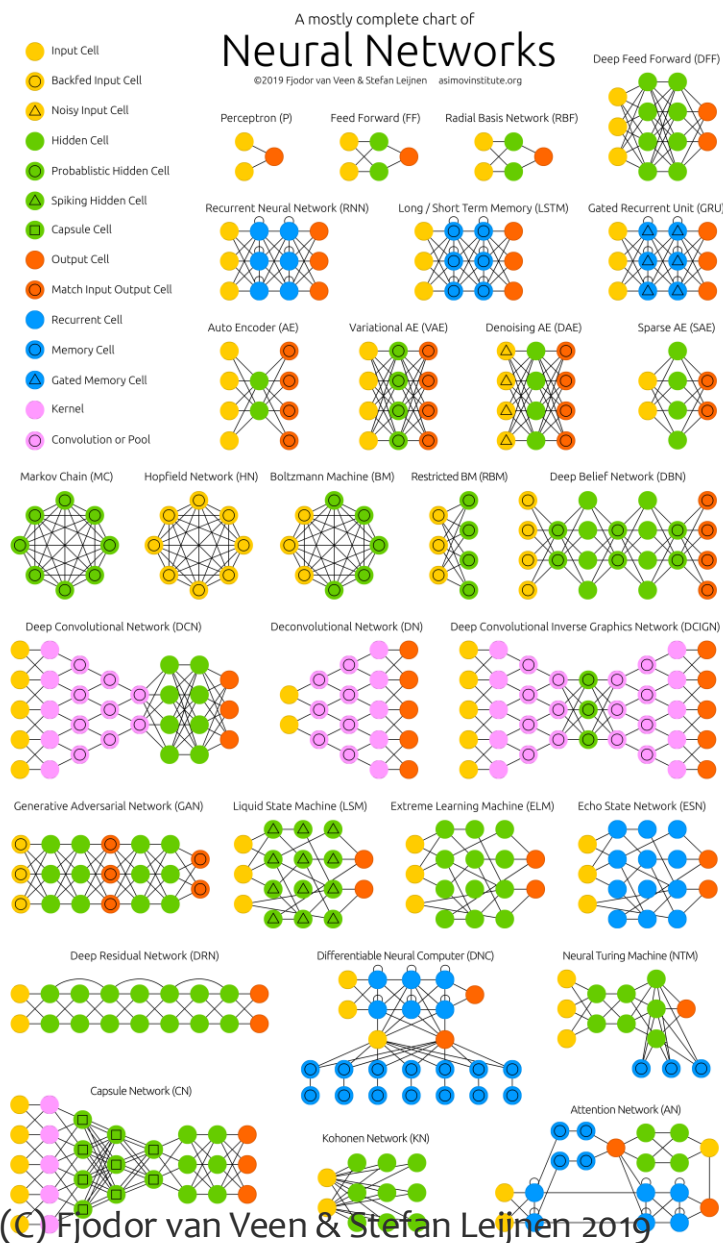
畳み込みニューラルネットワーク

物体分類

物体検出・セグメンテーション

敵対的生成ネットワーク

# ニューラルネットワーク



- 多層化パーセプトロン

入力から出力まですべての層で完全に結合されているニューラルネットワーク。分類や回帰に用いられる。

- 畳み込みニューラルネットワーク (CNN)

畳み込み演算 (移動平均) により画像の特徴を抽出する層を加えたニューラルネットワーク。画像分類や物体検出に広く使われている。

- 再帰型ニューラルネットワーク (RNN)

過去の情報を記録するパラメータを加えたニューラルネットワーク。時系列データを利用した予測などに使われる。LSTM や GRU などの改良型がある。

- 敵対的生成ネットワーク (GAN)

敵対的生成ネットワークは、2つのニューラルネットワーク (生成器と識別器) を対立的に学習させ、その後、生成器だけを利用してデータ生成させる。

# 深層学習

---

ニューラルネットワーク

**畳み込みニューラルネットワーク**

**物体分類**

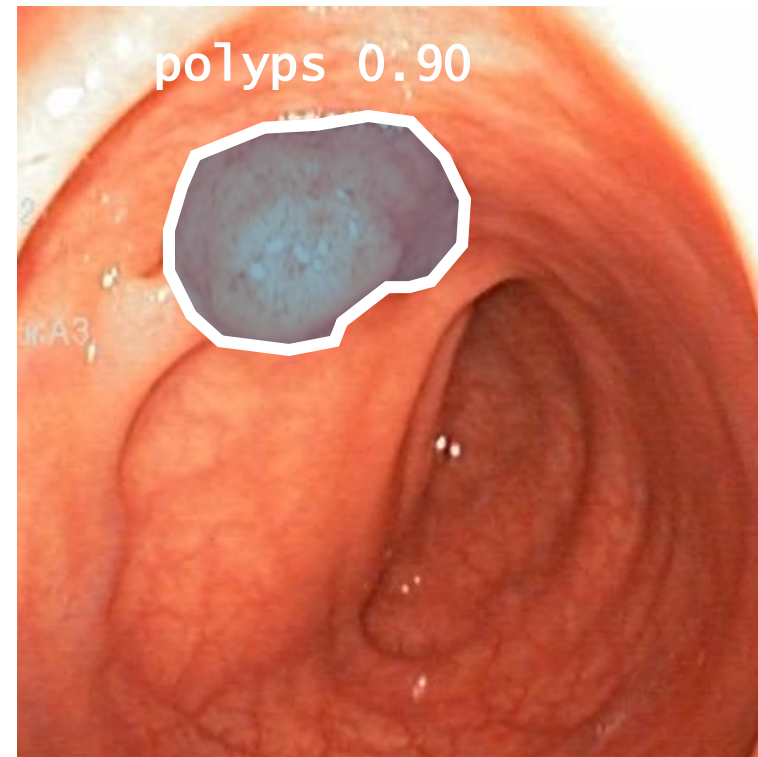
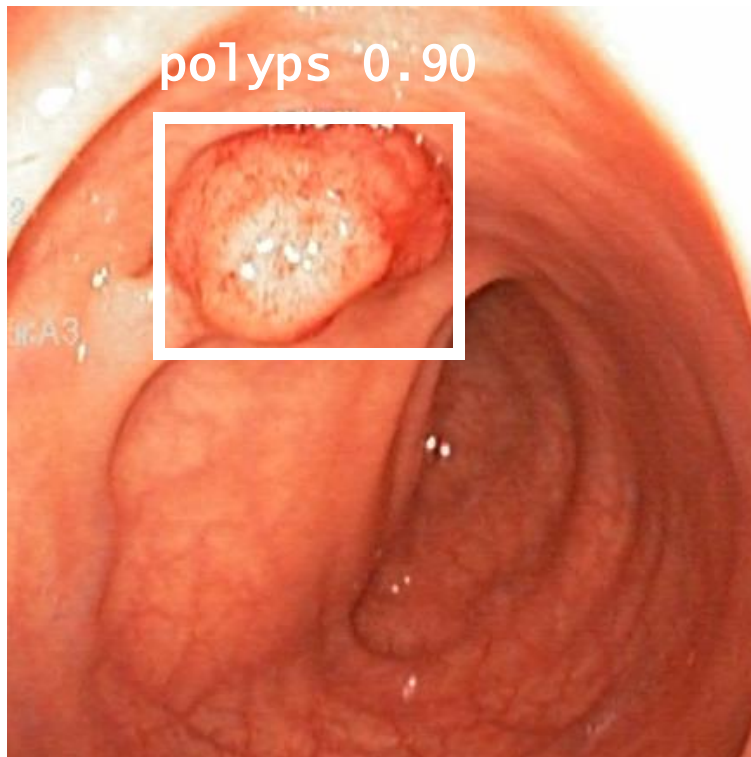
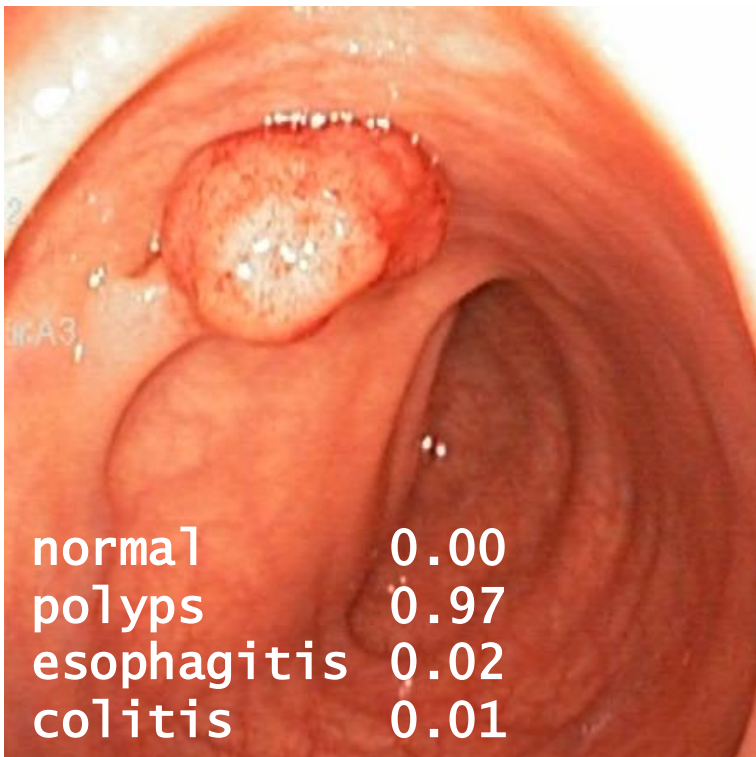
物体検出・セグメンテーション

敵対的生成ネットワーク

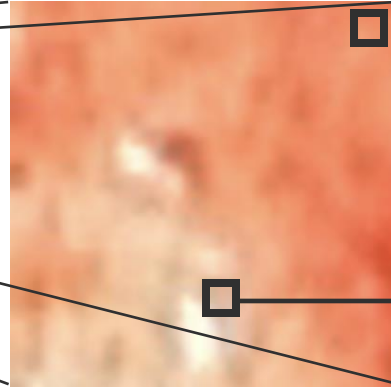
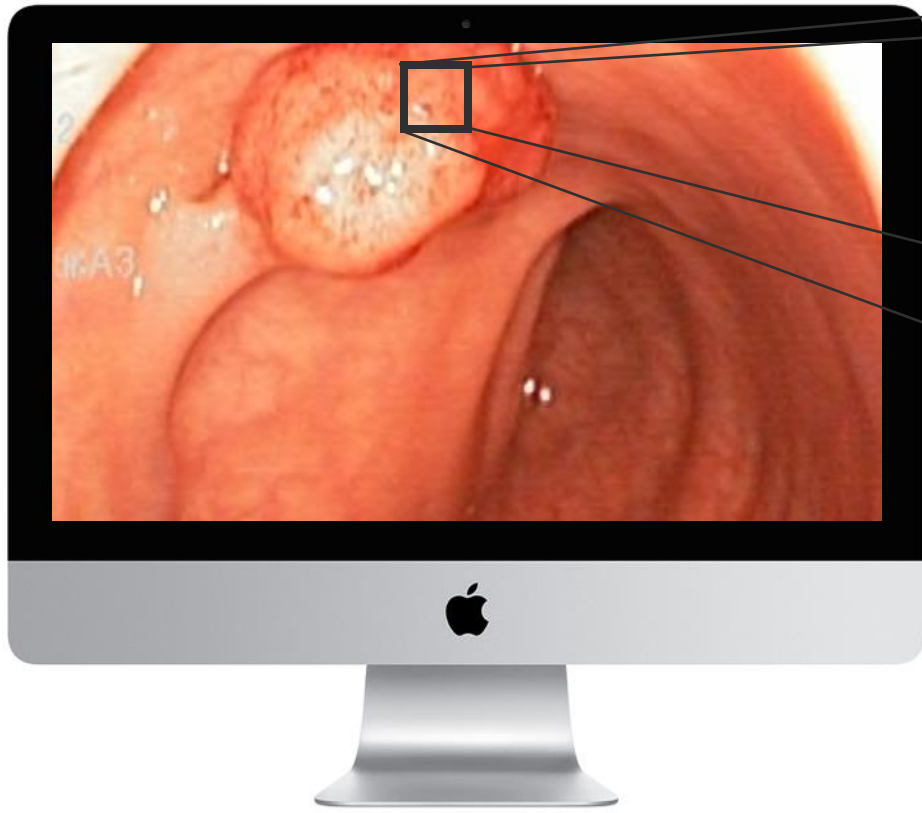
## 物体分類

## 物体検出

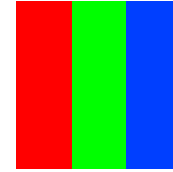
## セグメンテーション



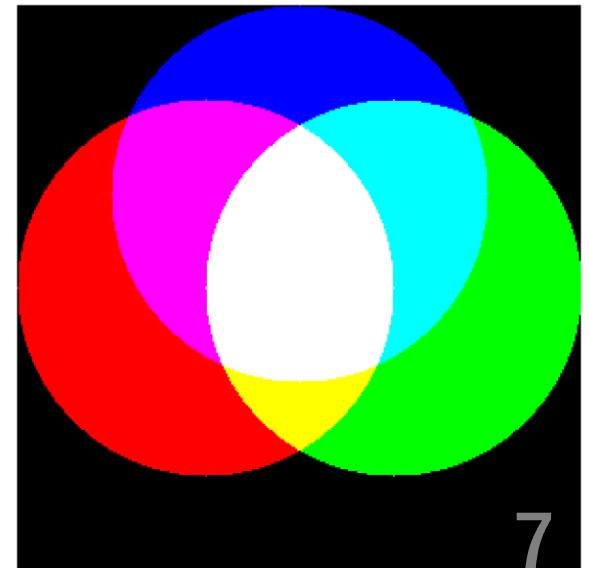
# 物体認識



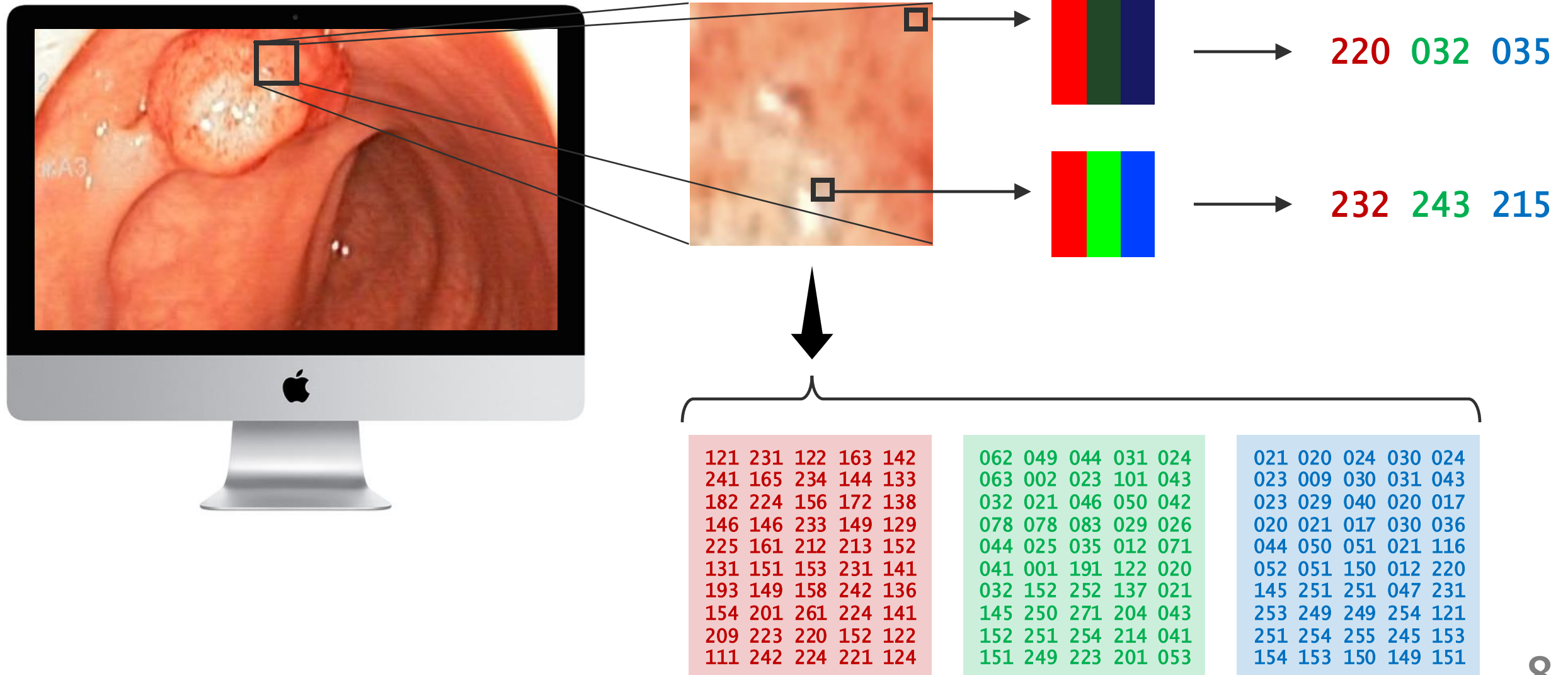
220 032 035



232 243 215

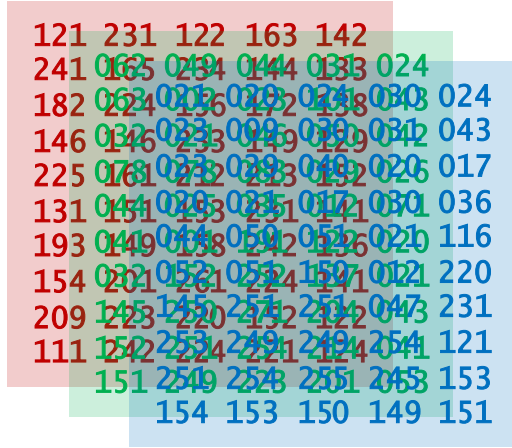
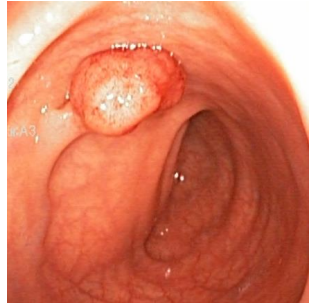


# 物体認識

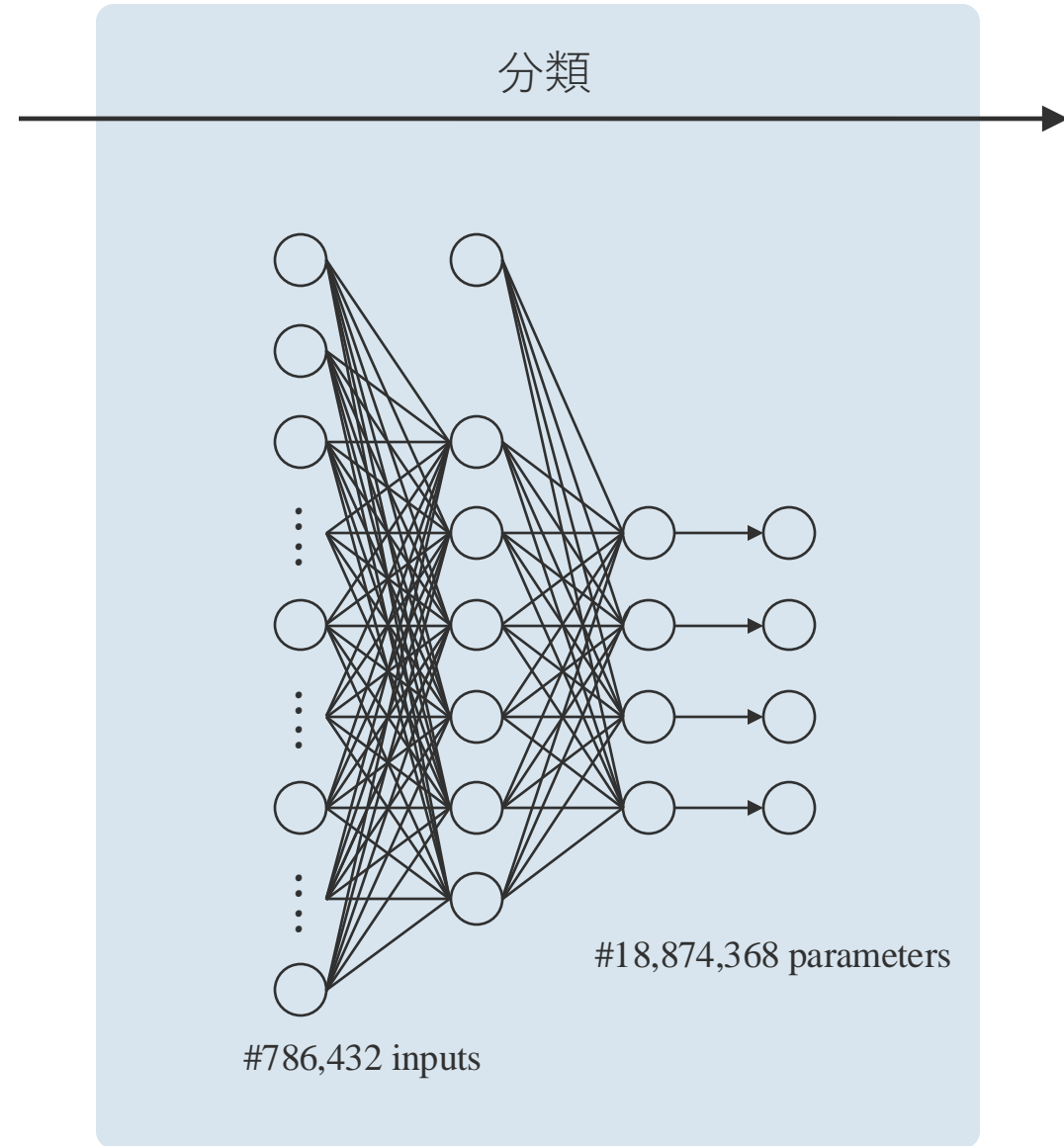




# 物体分類

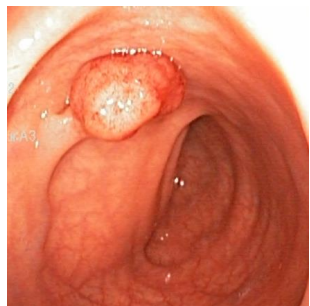


$$512 \times 512 \times 3 = 786,432$$



normal	0.00
polyps	0.97
esophagitis	0.02
colitis	0.01

# 物体分類



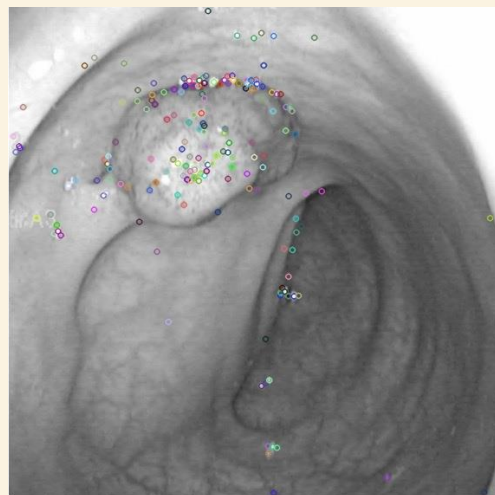
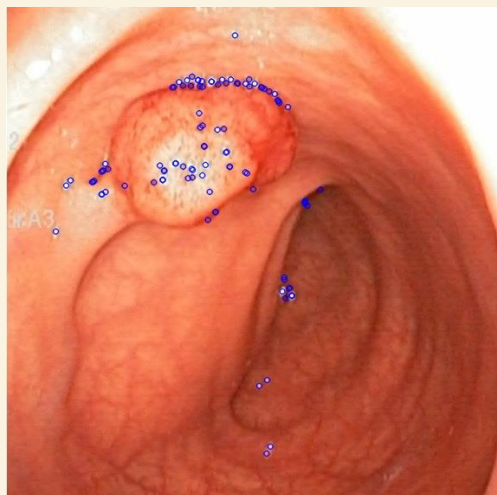
121 231 122 163 142  
 241 165 049 044 033 024  
 182 062 002 028 021 030 024  
 146 032 023 008 030 031 043  
 225 078 073 028 049 026 017  
 131 044 028 031 017 030 036  
 131 151 233 231 141 071  
 193 040 055 050 052 021 116  
 154 201 052 051 150 012 220  
 209 145 150 254 254 043 231  
 111 223 220 152 122 043  
 111 252 253 249 249 254 121  
 151 249 253 251 248 153  
 154 153 150 149 151

$512 \times 512 \times 3 = 786,432$

特徴抽出

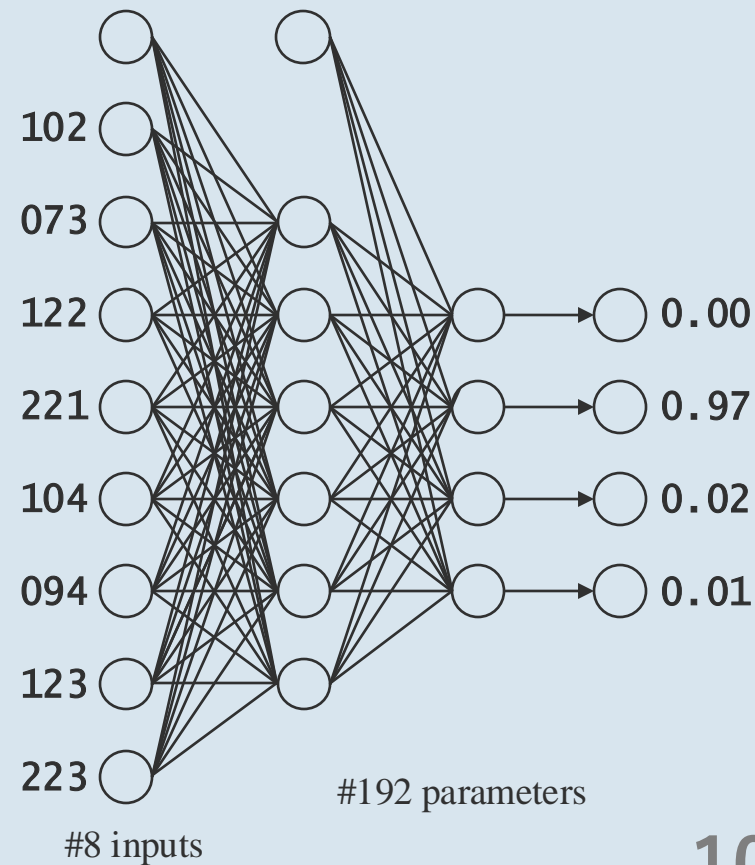
102  
 073  
 122  
 221  
 104  
 094  
 123  
 223

SIFT  
 SURF  
 AKAZE

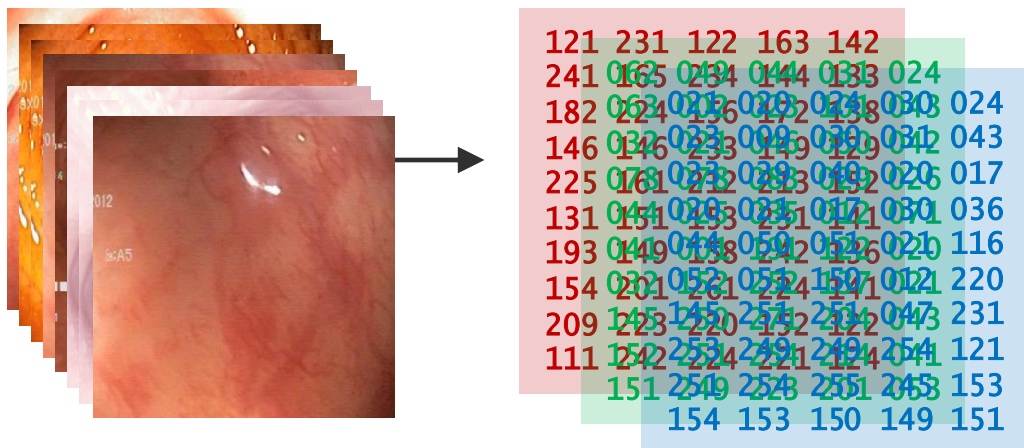


分類

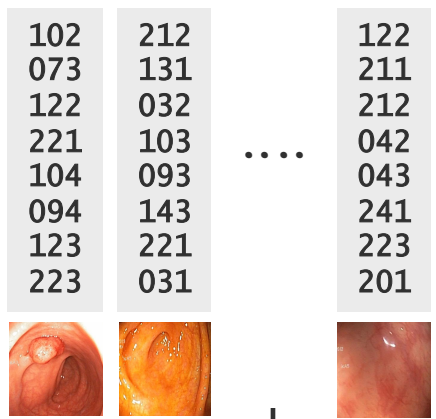
normal 0.00  
 polyps 0.97  
 esophagitis 0.02  
 colitis 0.01



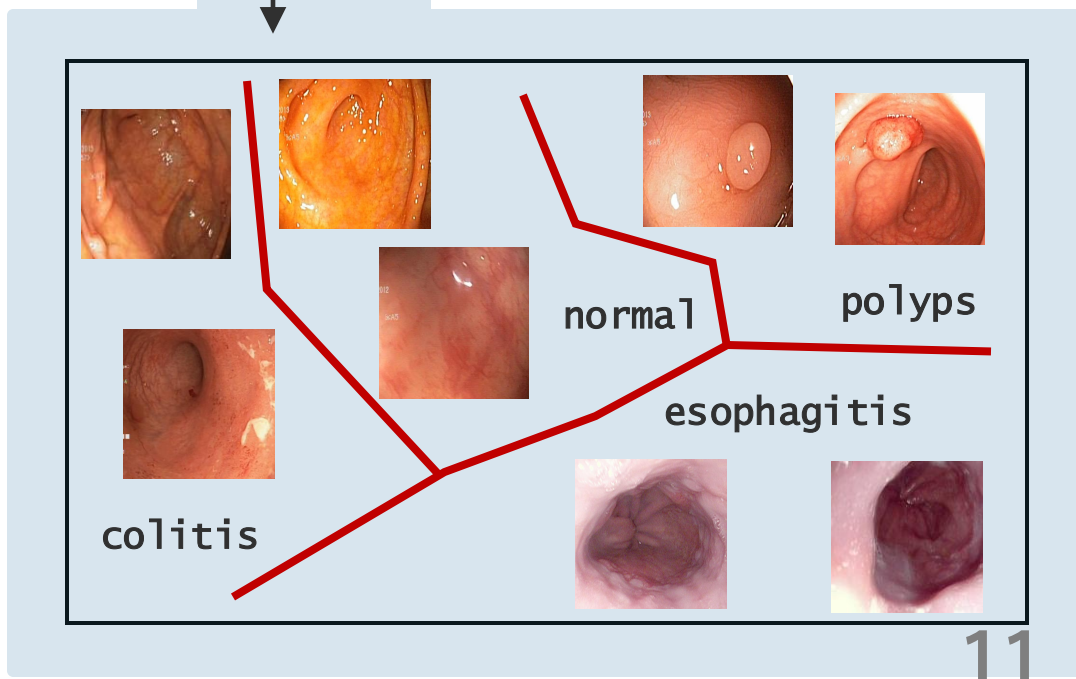
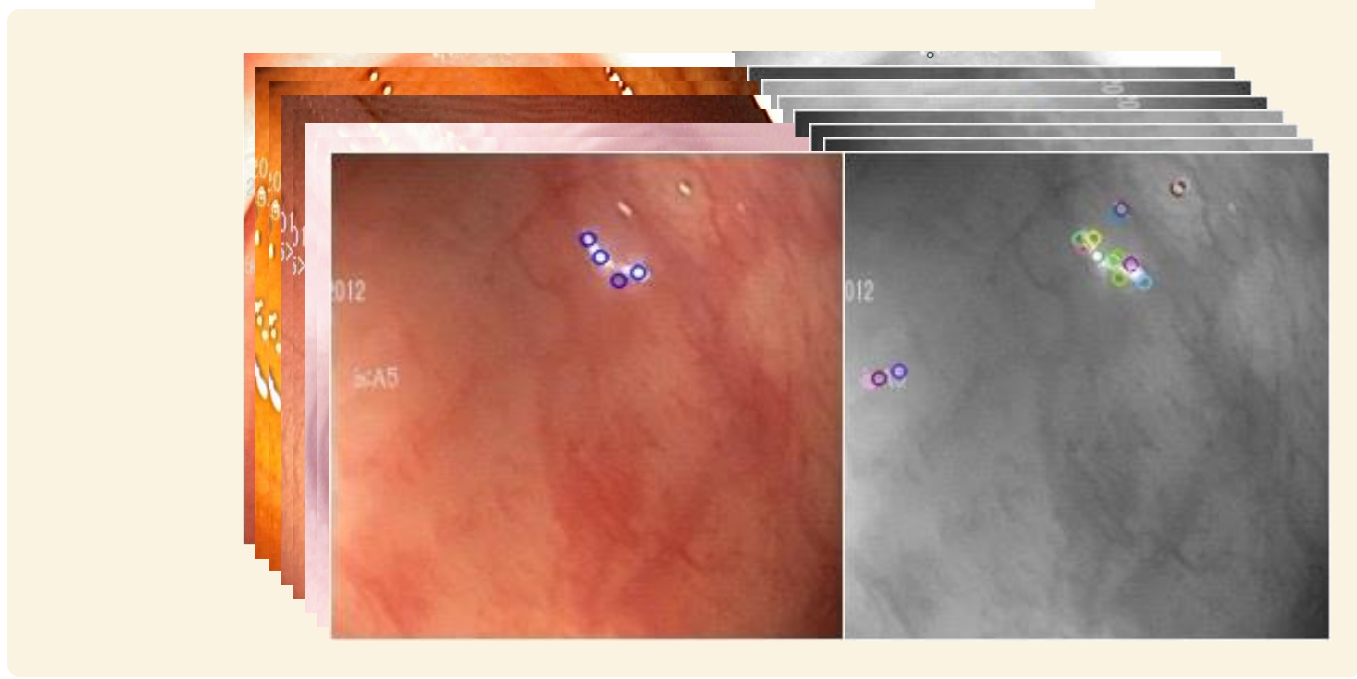
# 物体分類 (學習)



特徴抽出



学習



# 物体分類 (推論)



128 213 154 214 211  
 251 023 049 044 033 024  
 224 052 002 028 021 030 032  
 186 012 023 008 030 031 075  
 215 048 073 028 049 026 032  
 241 134 028 081 012 030 078  
 215 149 058 242 026 020  
 223 142 052 051 150 012 198  
 228 225 220 152 122 043  
 227 242 253 249 244 254 132  
 211 249 228 251 248 186  
 021 133 111 034 213

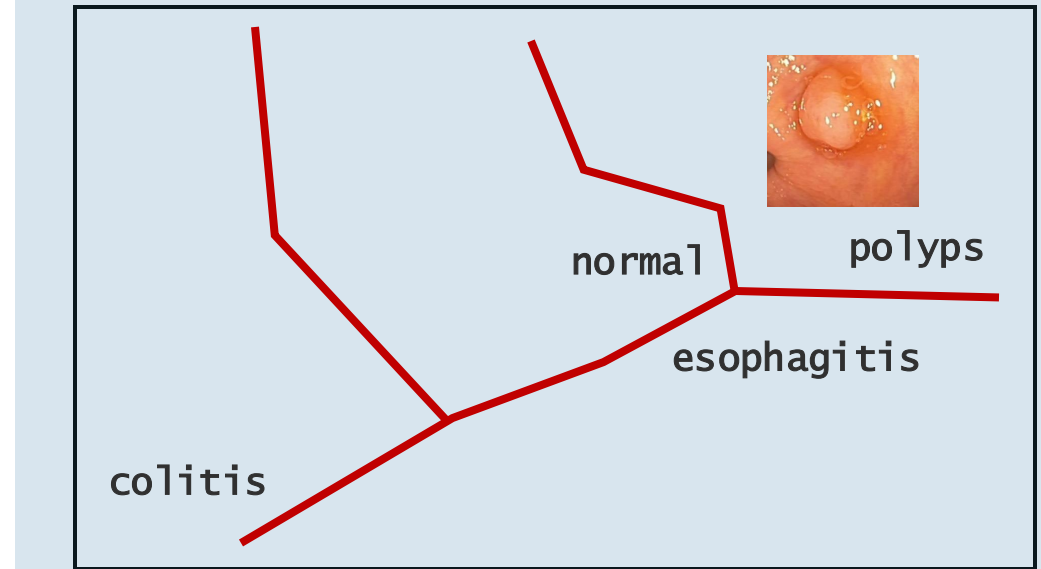
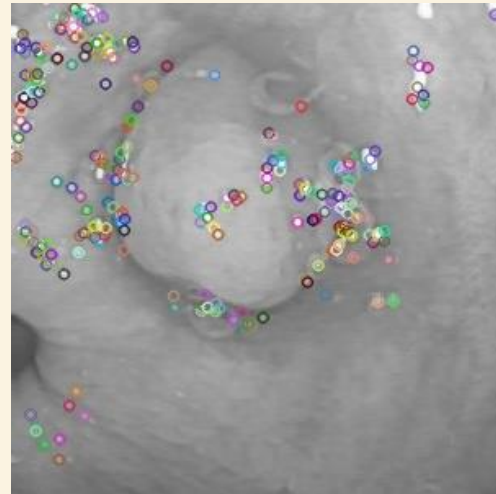
特徴抽出

220  
 023  
 019  
 233  
 098  
 212  
 023  
 167



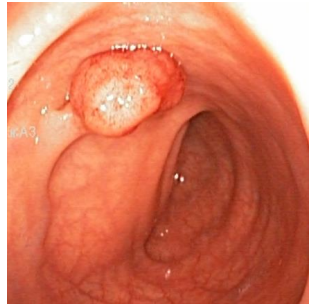
分類

normal	0.01
polyps	0.99
esophagitis	0.00
colitis	0.00





# 物体分類



128 213 154 214 211  
251 025 049 044 033 024  
224 052 002 028 021 030 032  
186 012 023 008 080 031 075  
215 048 073 028 049 026 032  
241 134 028 081 017 030 078  
215 149 058 050 052 021 223  
223 142 052 051 150 012 198  
228 213 150 251 254 043 193  
227 242 253 249 244 254 132  
211 249 228 251 248 186  
021 133 111 034 213

特徴抽出

102  
073  
122  
221  
104  
094  
123  
223

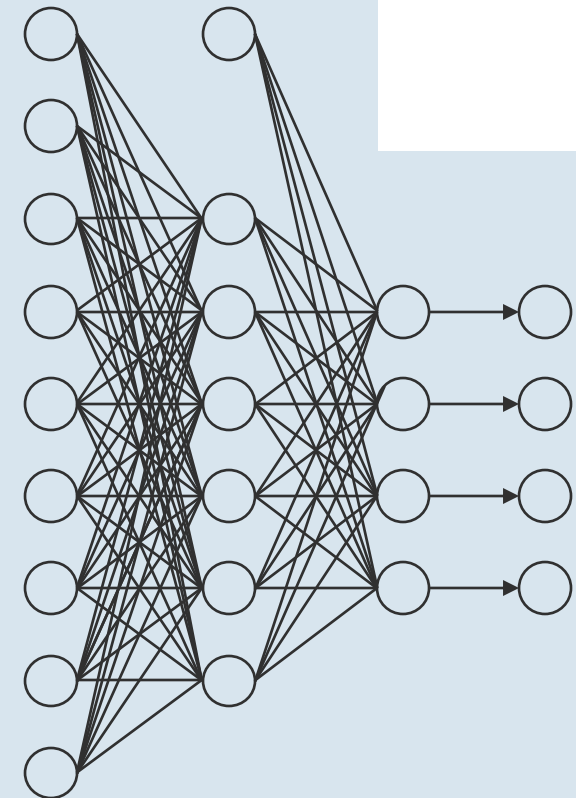


分類

normal	0.01
polyps	0.99
esophagitis	0.00
colitis	0.00

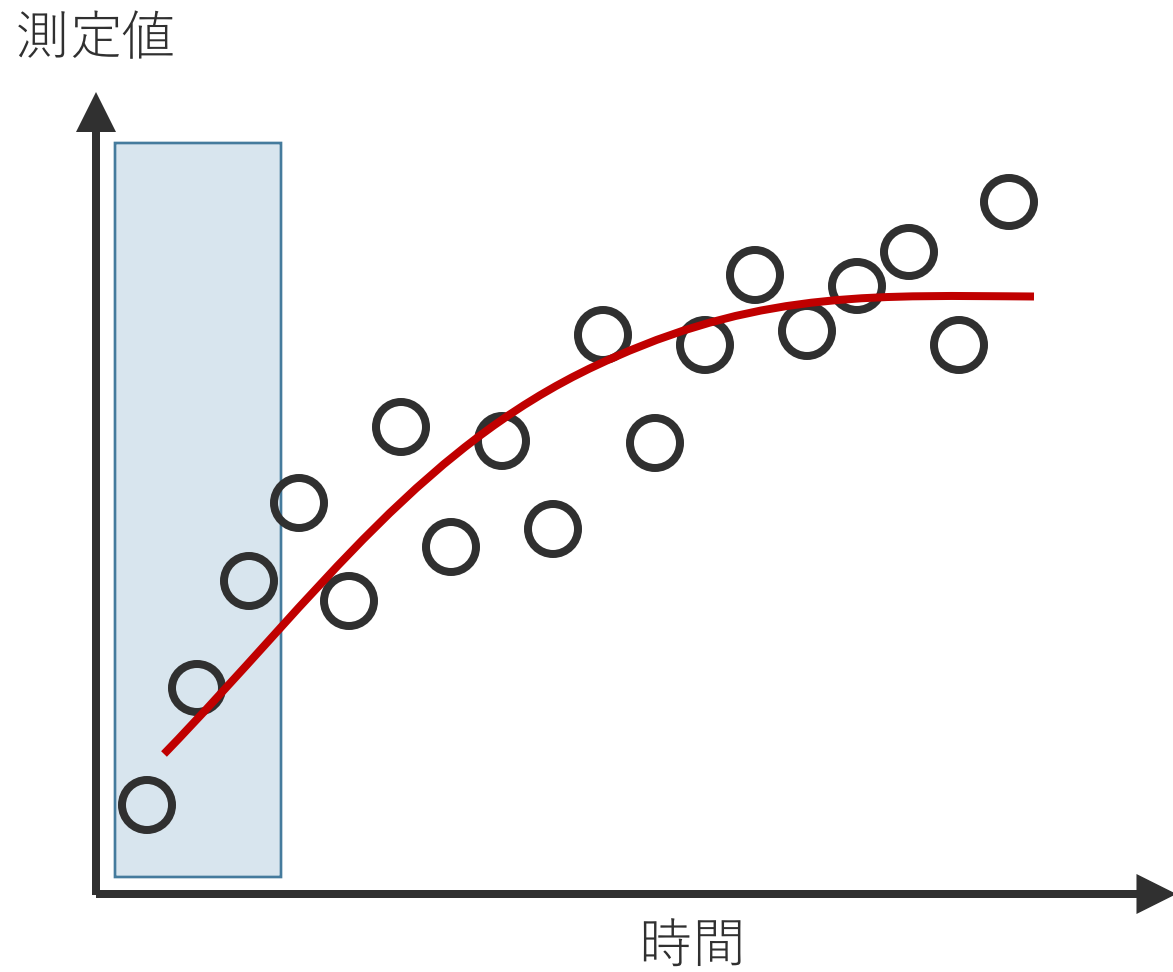
特徴抽出を、ニューラルネットワークのように自動的に学習する仕組みは実現できないだろうか？

- 畳み込み演算
- プーリング演算

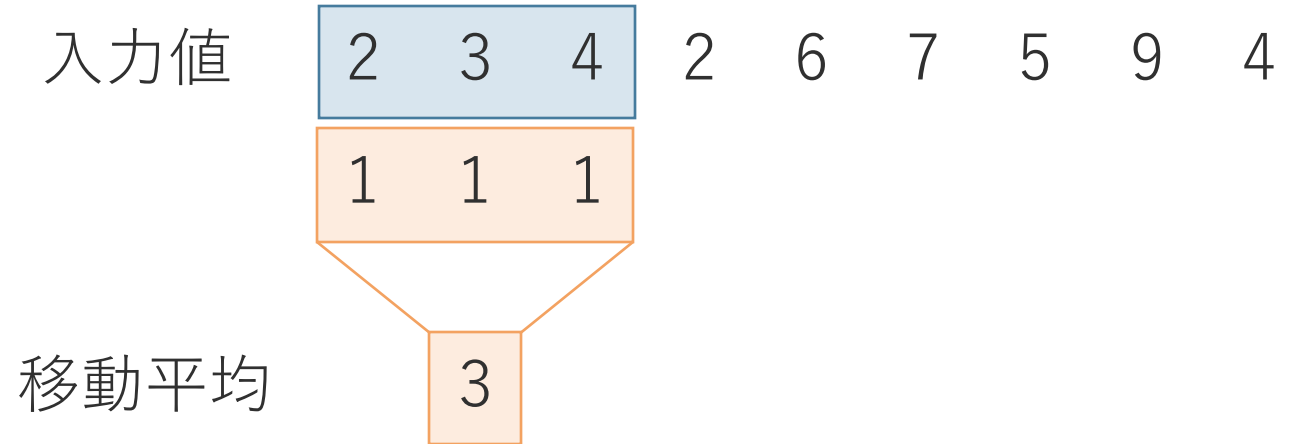
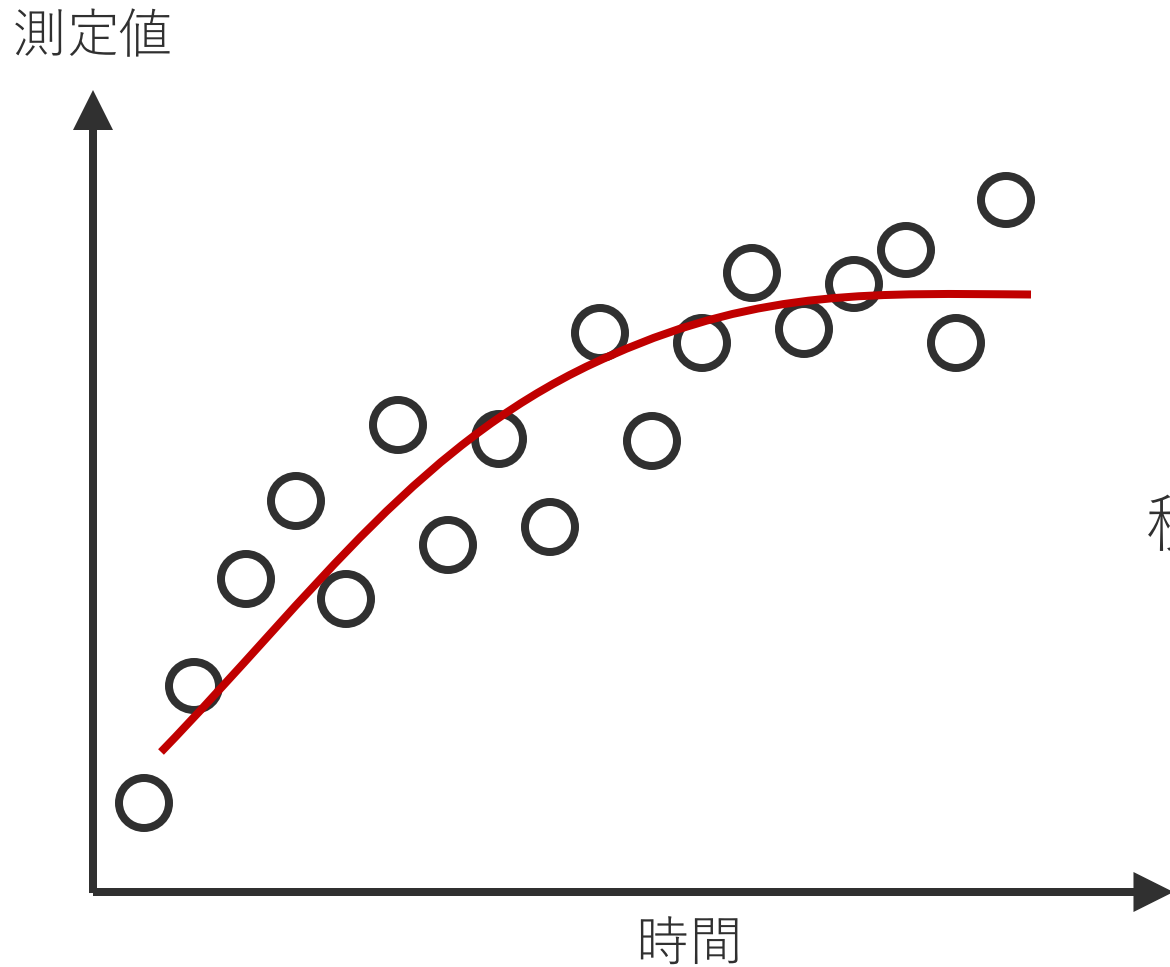


# 移動平均

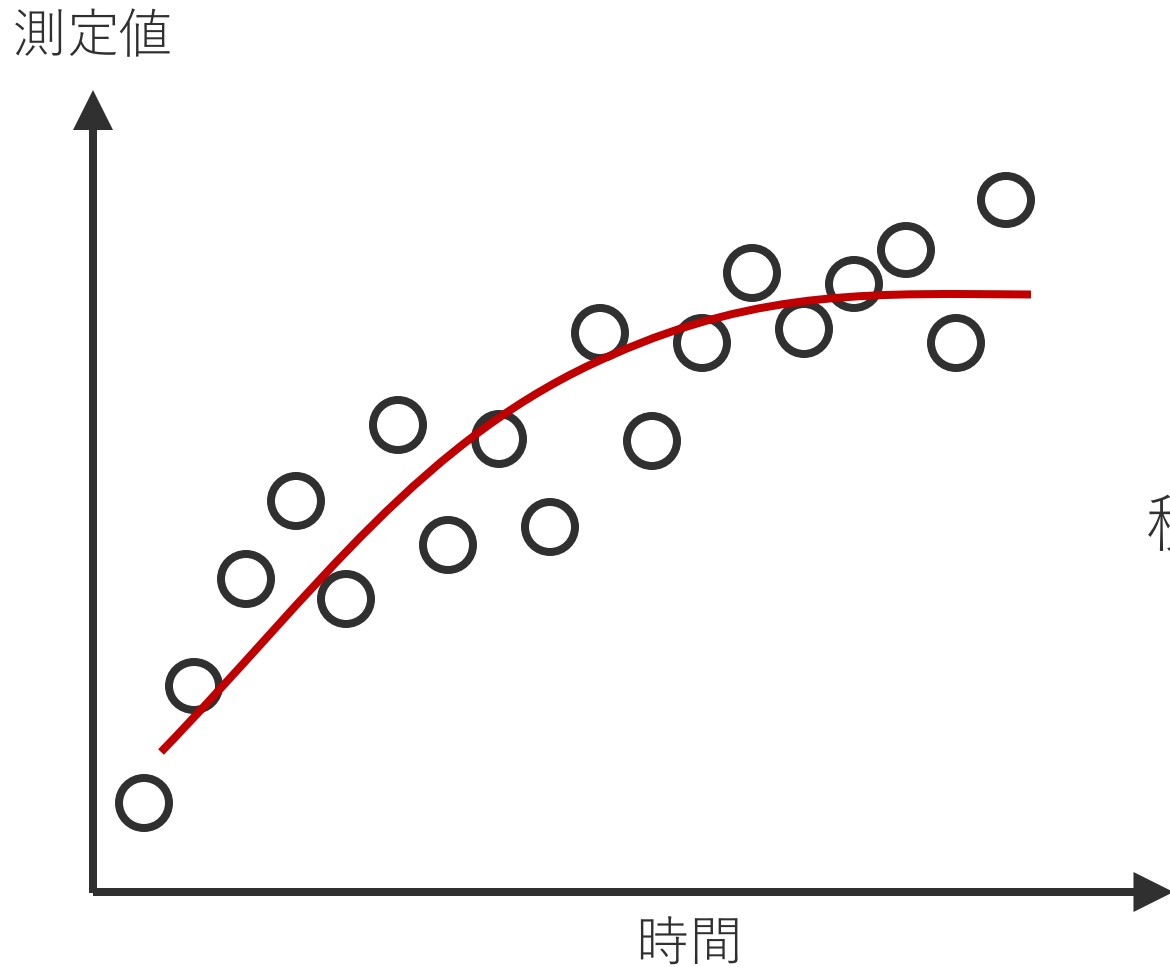
---



# 移動平均



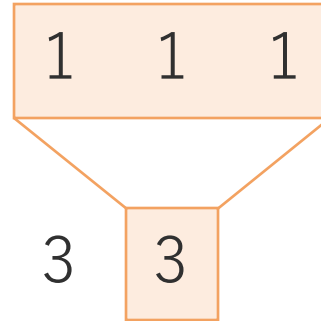
# 移動平均



入力値

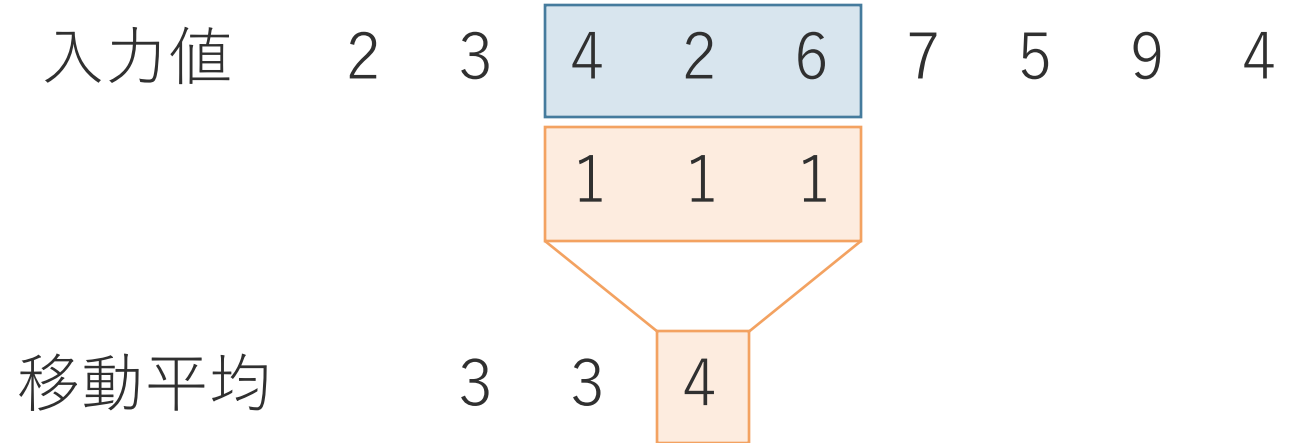
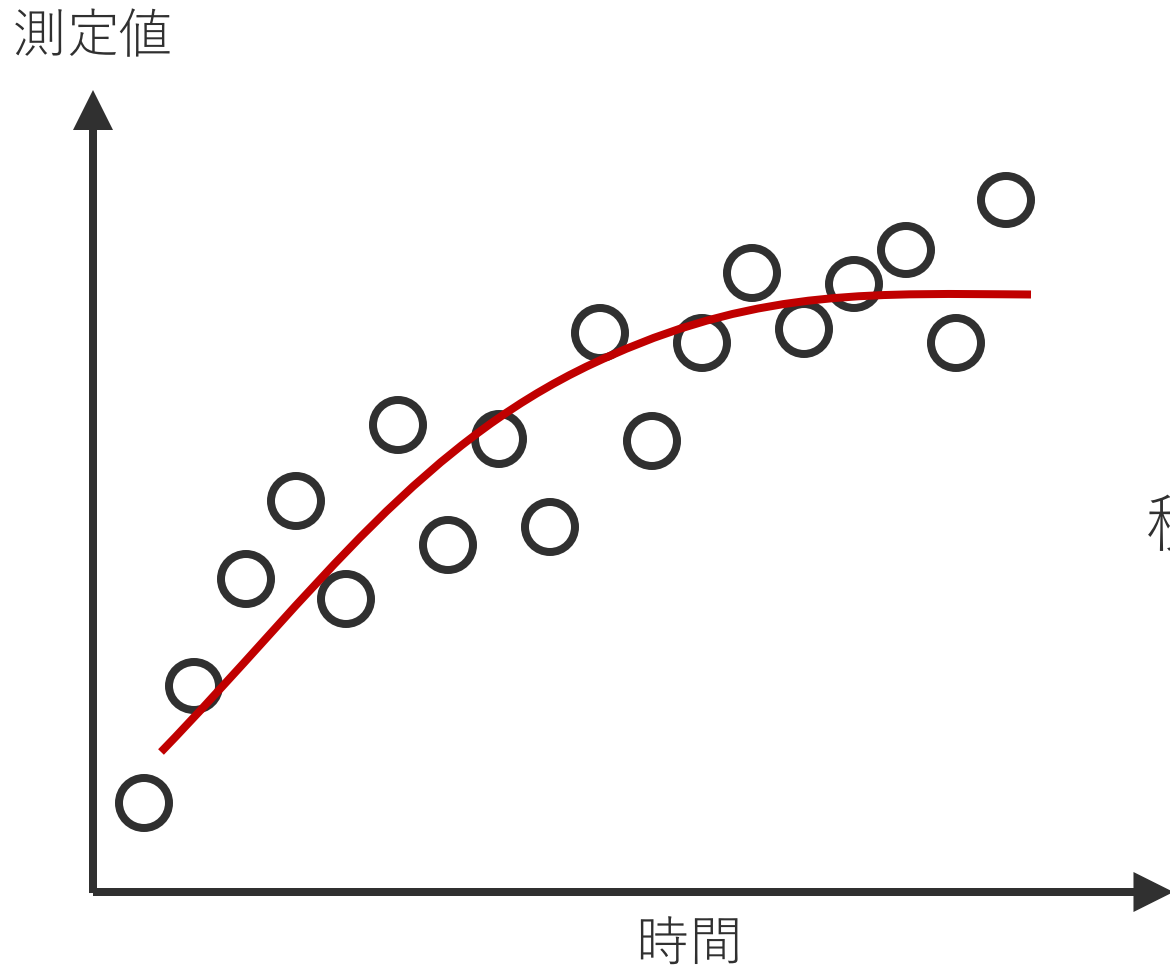
2 3 4 2 6 7 5 9 4

移動平均

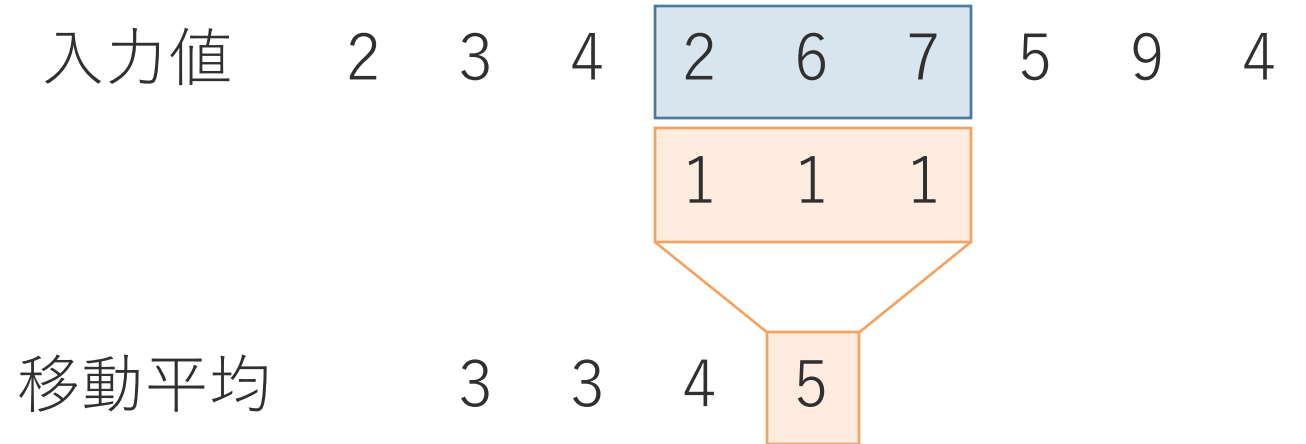
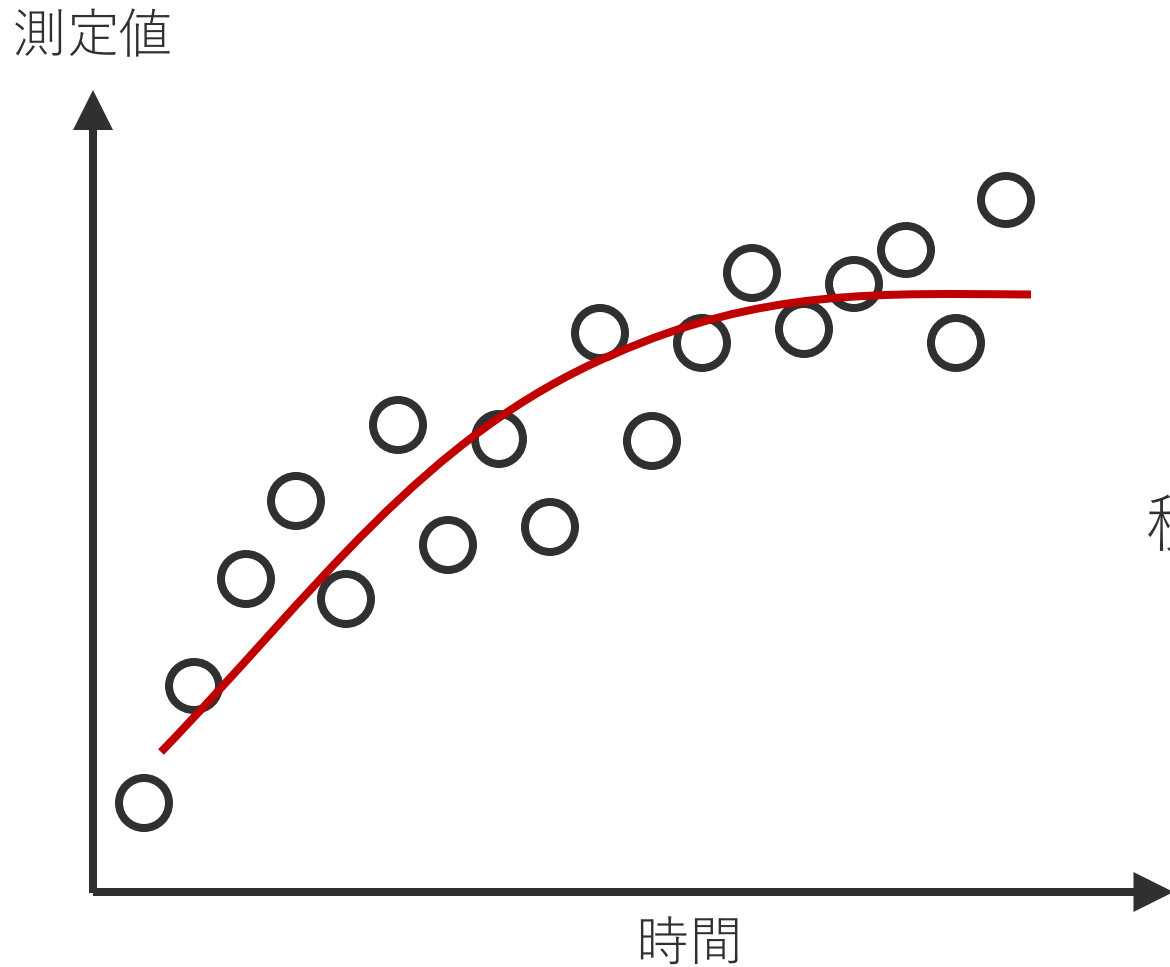




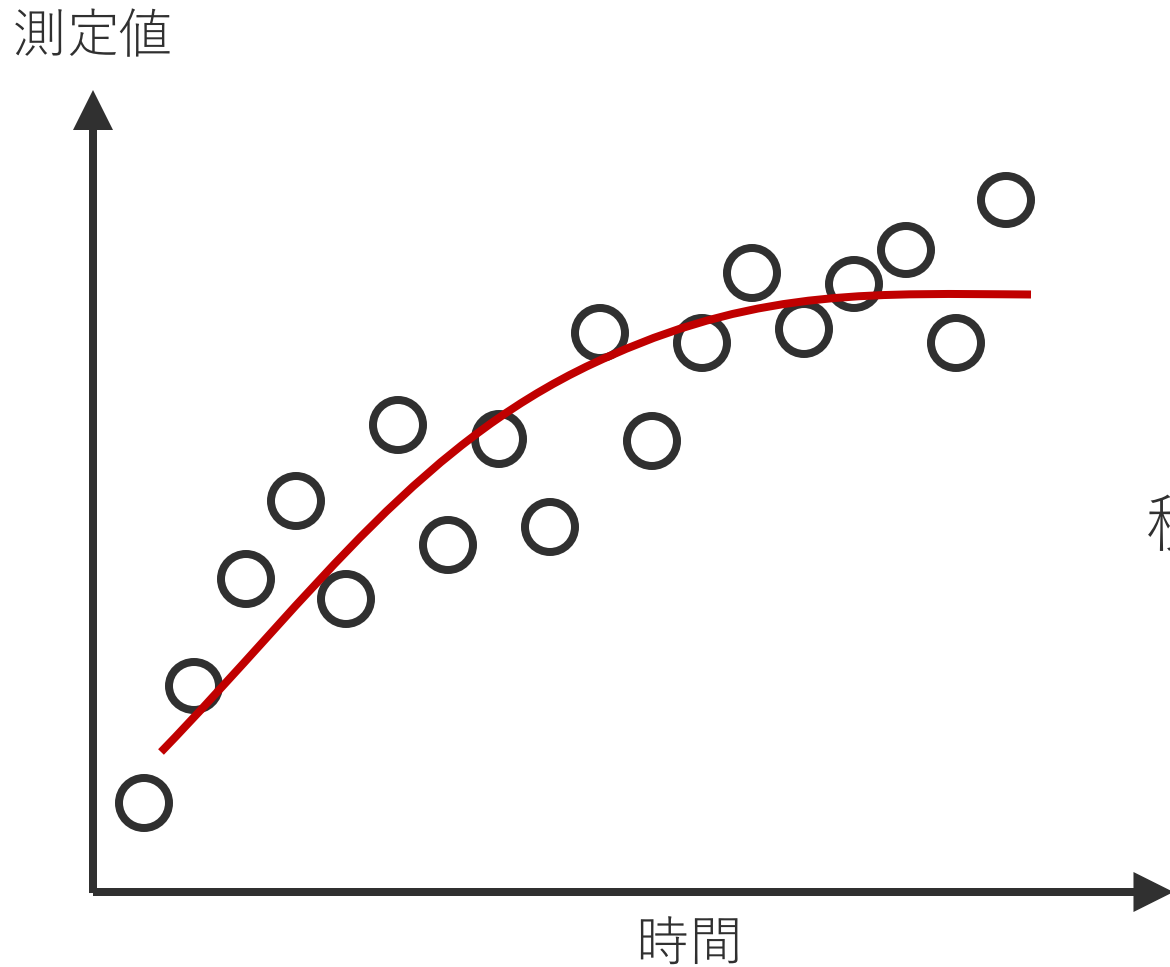
# 移動平均



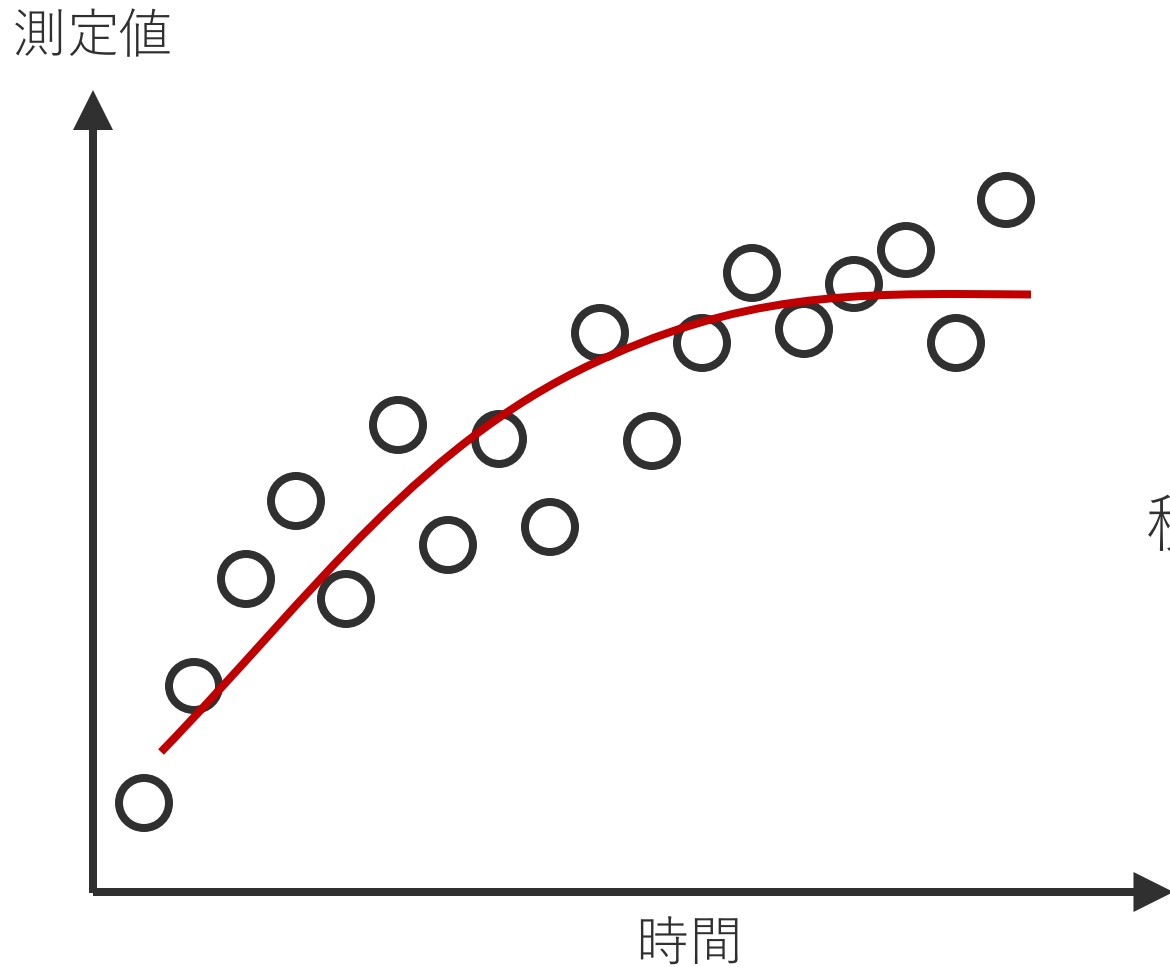
# 移動平均



# 移動平均

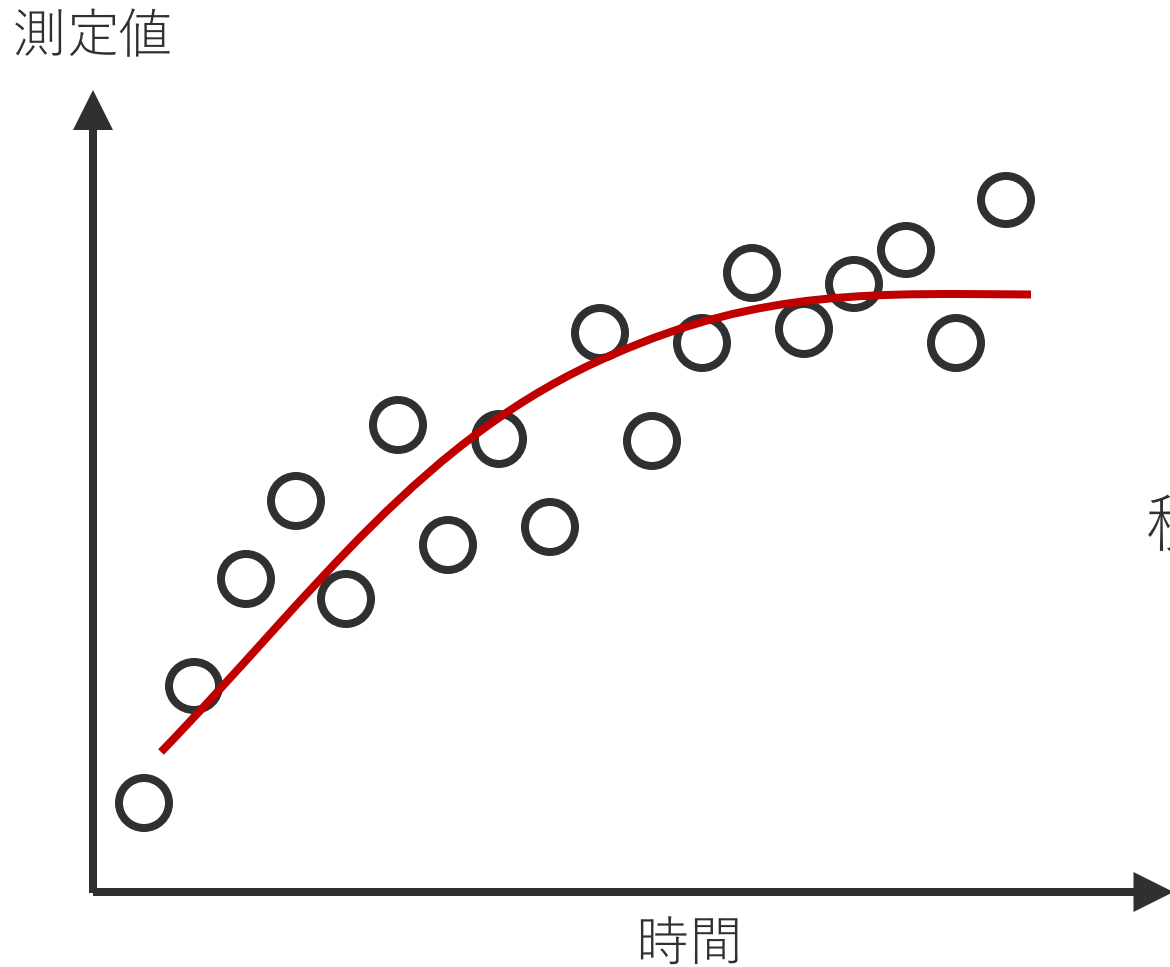


# 移動平均



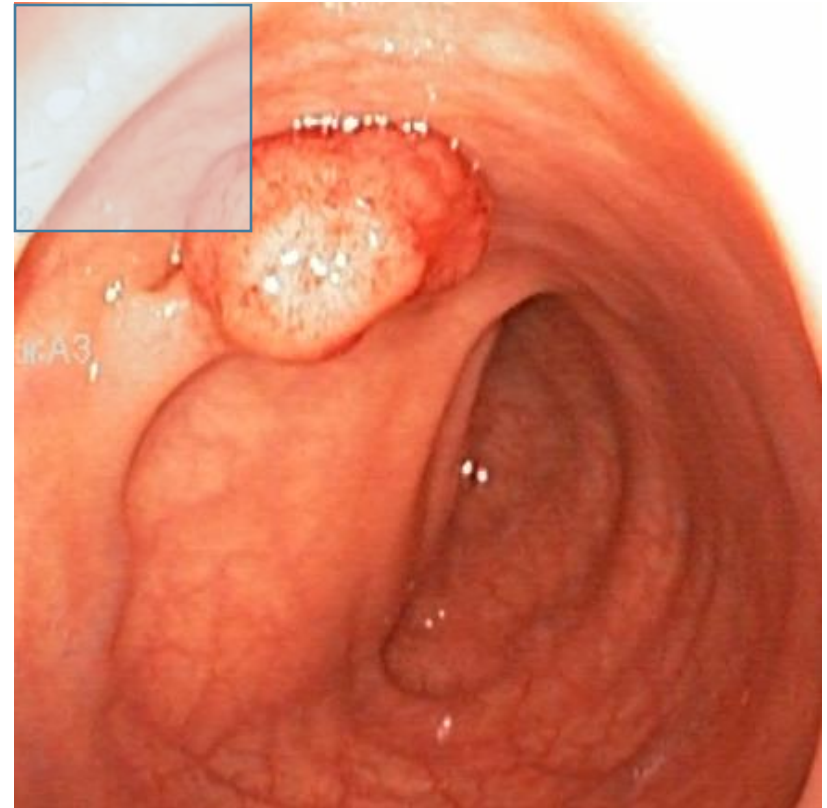
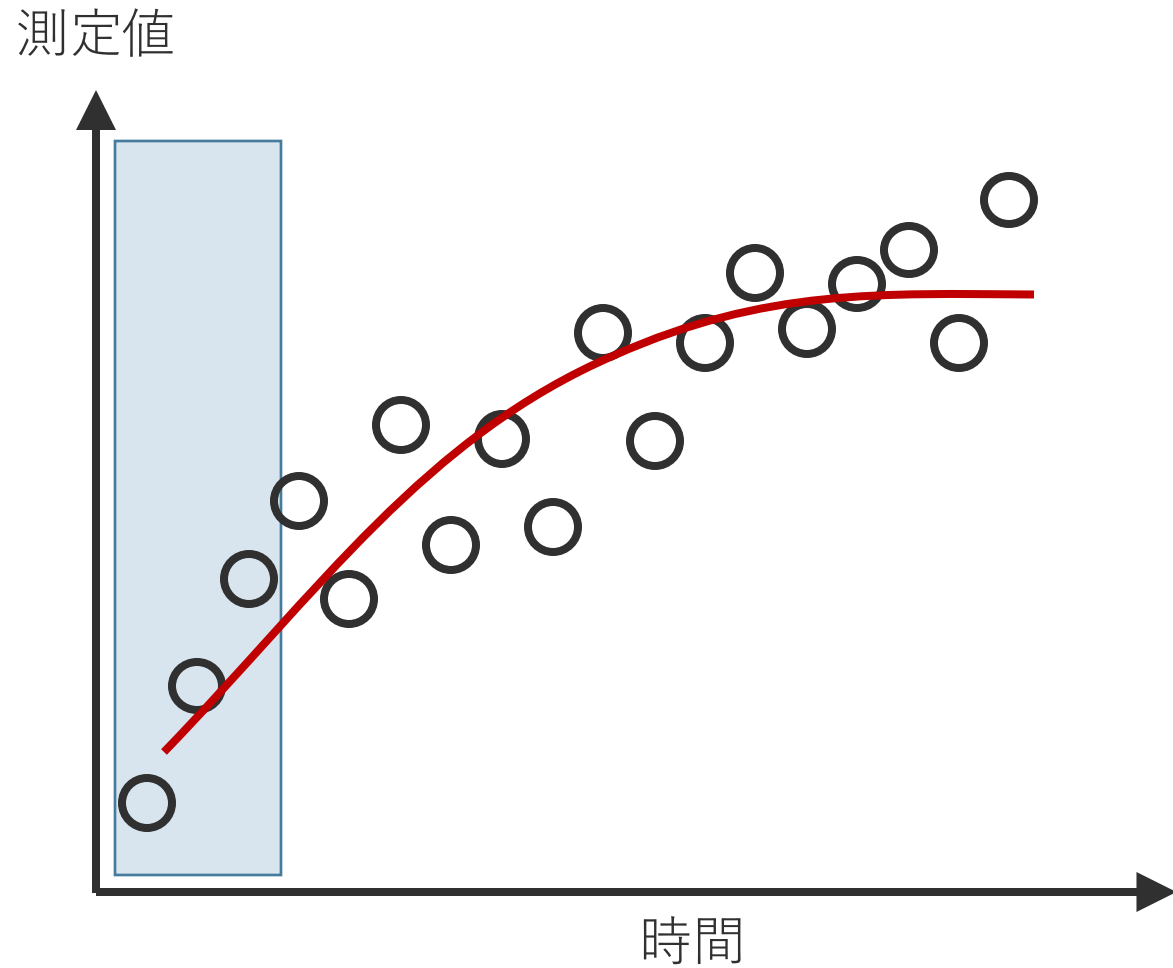
入力値	2	3	4	2	6	7	5	9	4
							1	1	1
移動平均		3	3	4	5	6	7	6	

# 移動平均



入力値	2	3	4	2	6	7	5	9	4
							1	0	1
移動平均	3.0	2.5	5.0	4.5	5.5	8.0	4.5		

# 畳み込み (移動平均)



# 畳み込み演算

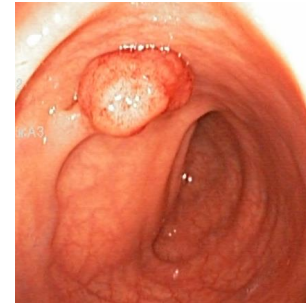
## Grayscale Photo



```
162 149 144 131 124
163 202 123 101 143
132 121 146 150 142
178 178 183 129 126
144 125 135 112 171
241 201 191 122 120
132 152 152 137 121
145 150 171 104 143
152 151 154 114 141
151 149 123 101 153
```

1 channel

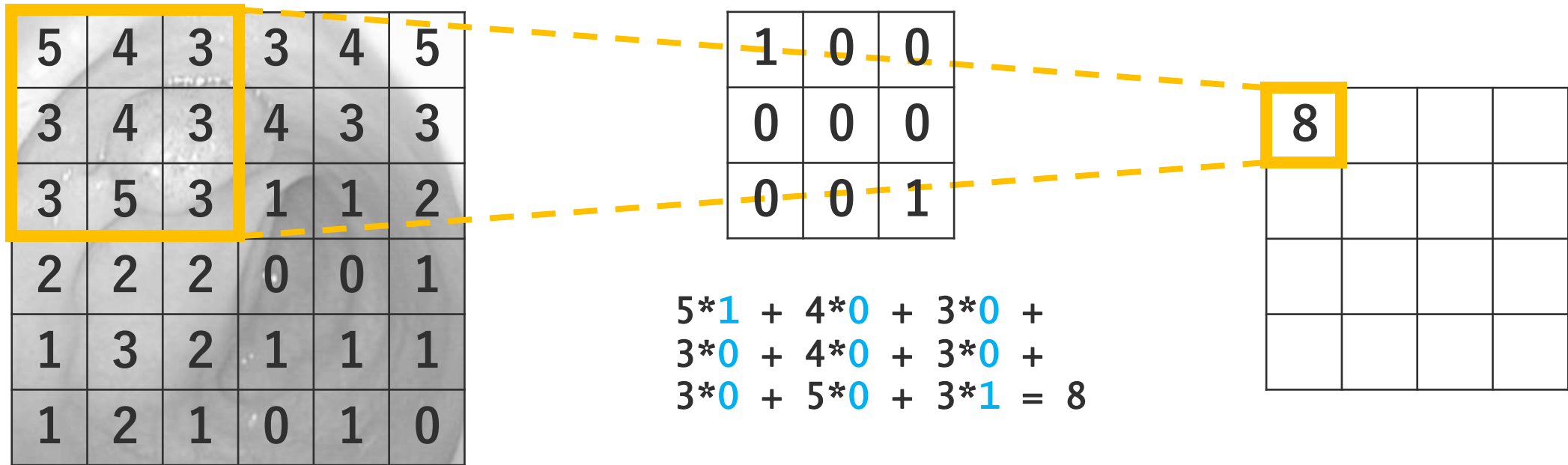
## Color Photo



```
021 031 022 063 042
041 065 034 044 033
162 149 144 131 124
082 024 056 072 038
163 202 123 101 143
046 046 033 049 029 0224
025 061 010 091 030 052 1243
178 178 183 129 126
031 051 053 081 002 10217
144 137 133 112 1171
093 049 058 092 026 0236
241 201 191 122 120
054 104 061 002 102 11216
132 152 152 137 121
109 023 020 053 022 0220
143 130 117 104 1143
111 042 024 102 102 1231
152 151 154 114 141
151 149 123 101 153
251 254 255 245 253
254 253 250 249 251
```

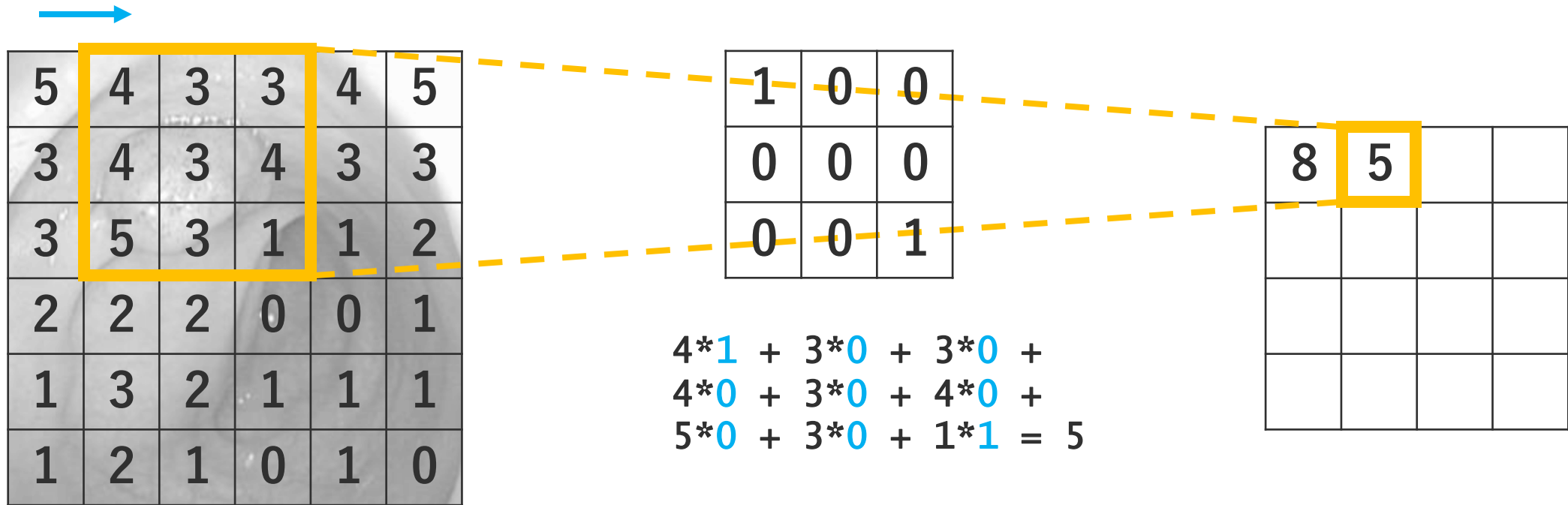
3 channels

# 畳み込み演算

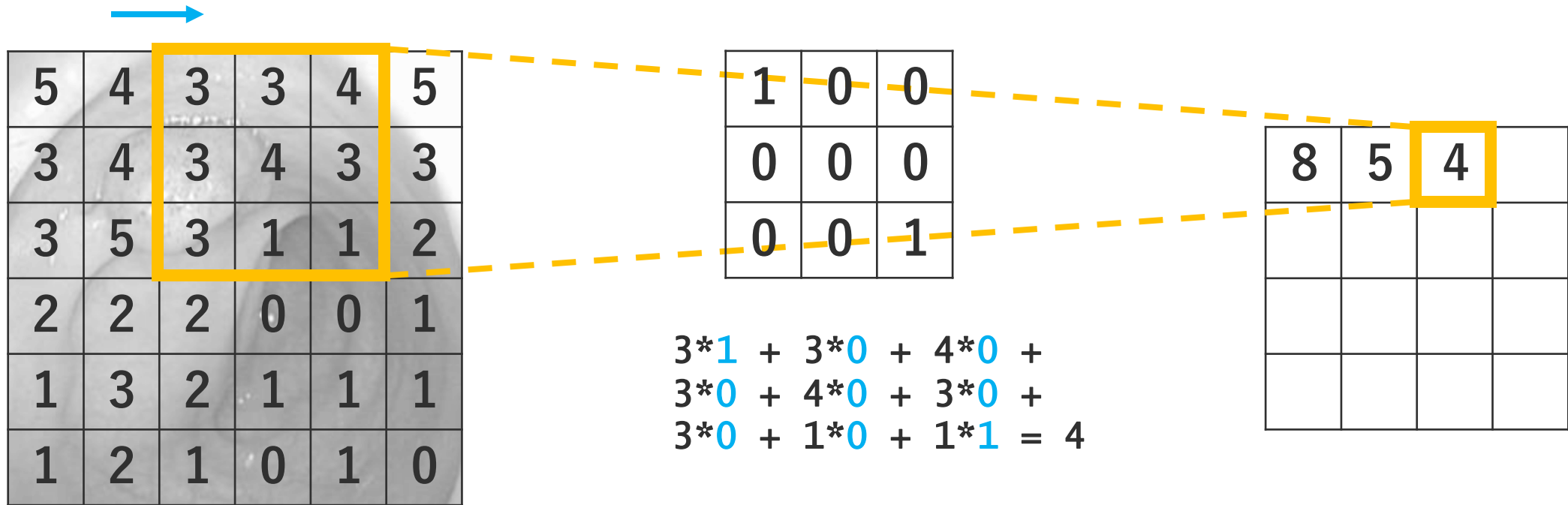




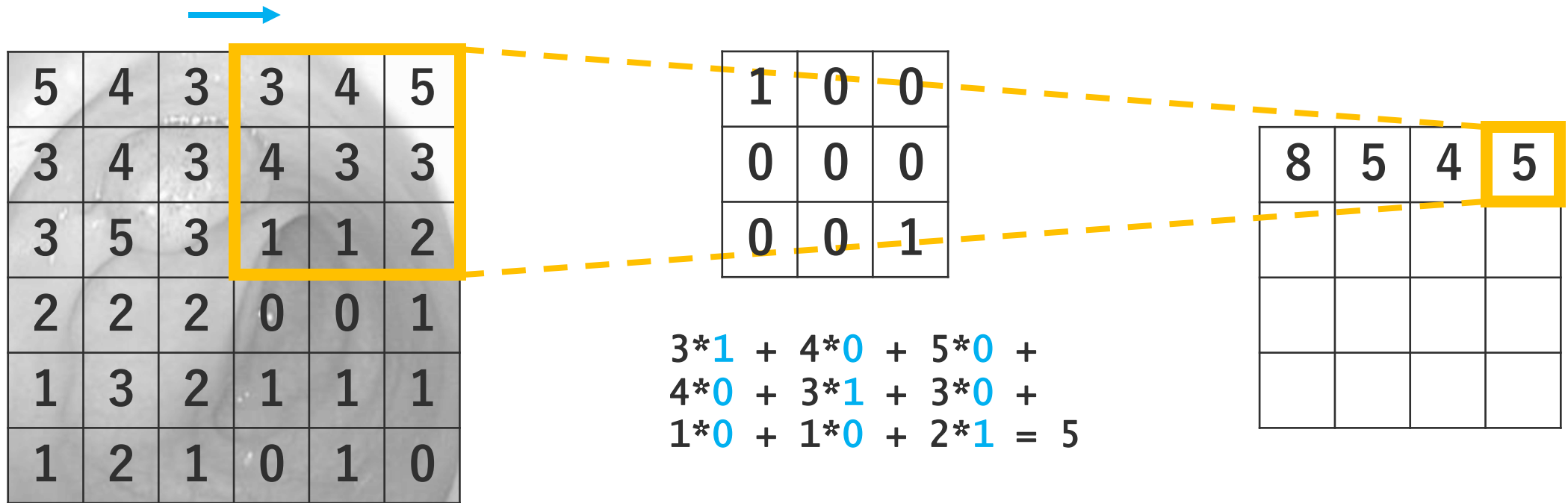
# 畳み込み演算



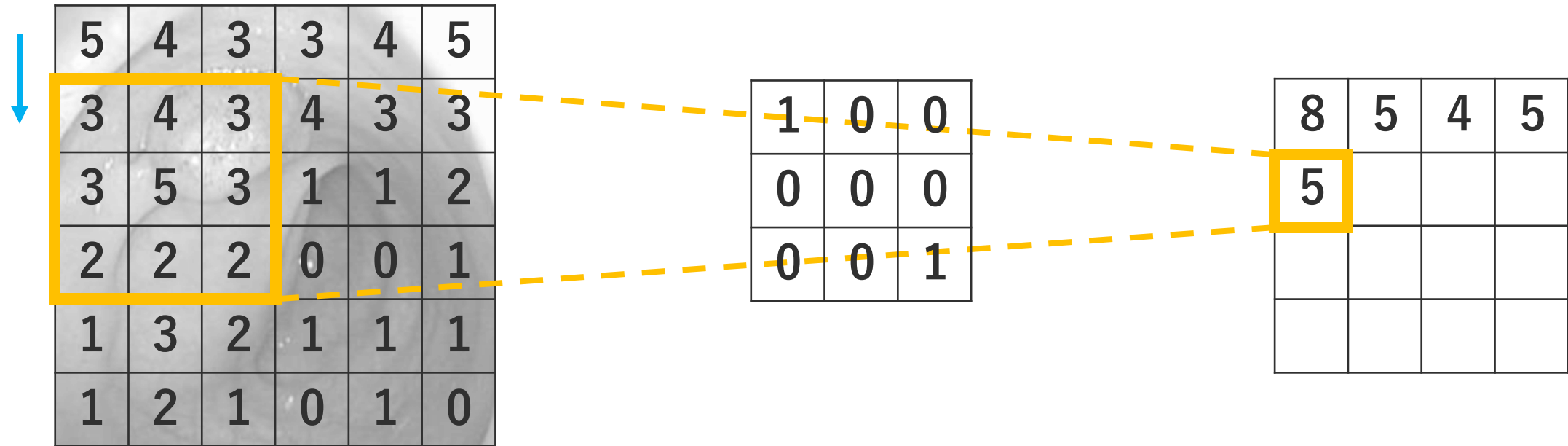
# 畳み込み演算



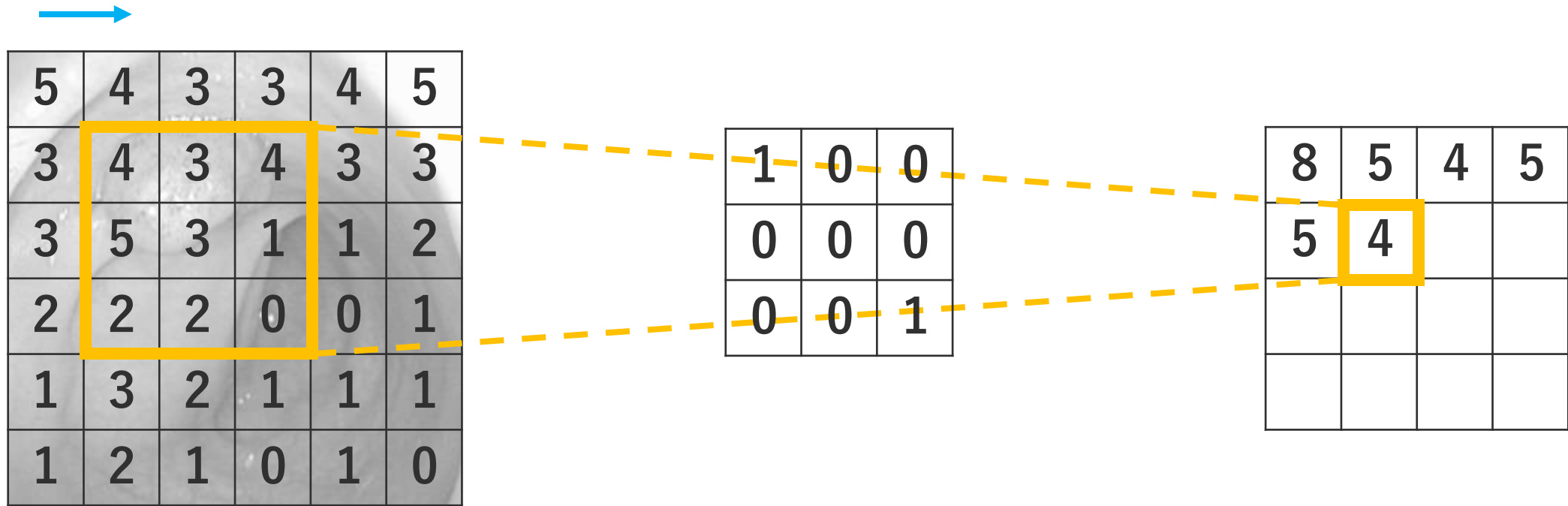
# 畳み込み演算



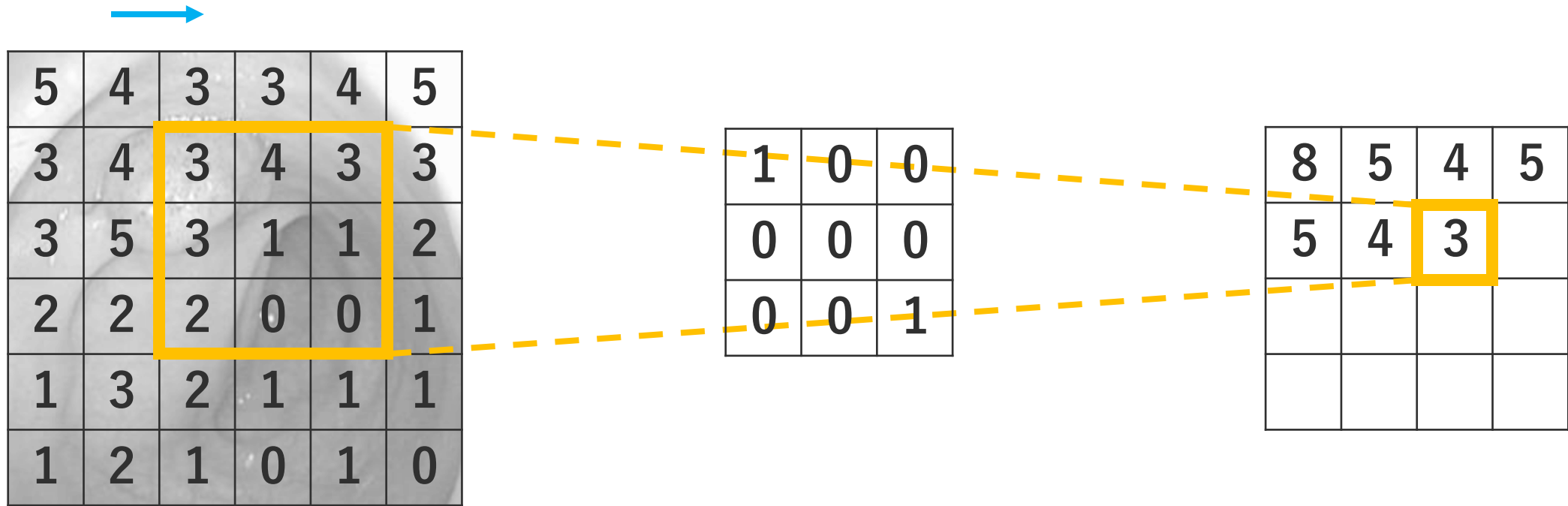
# 畳み込み演算



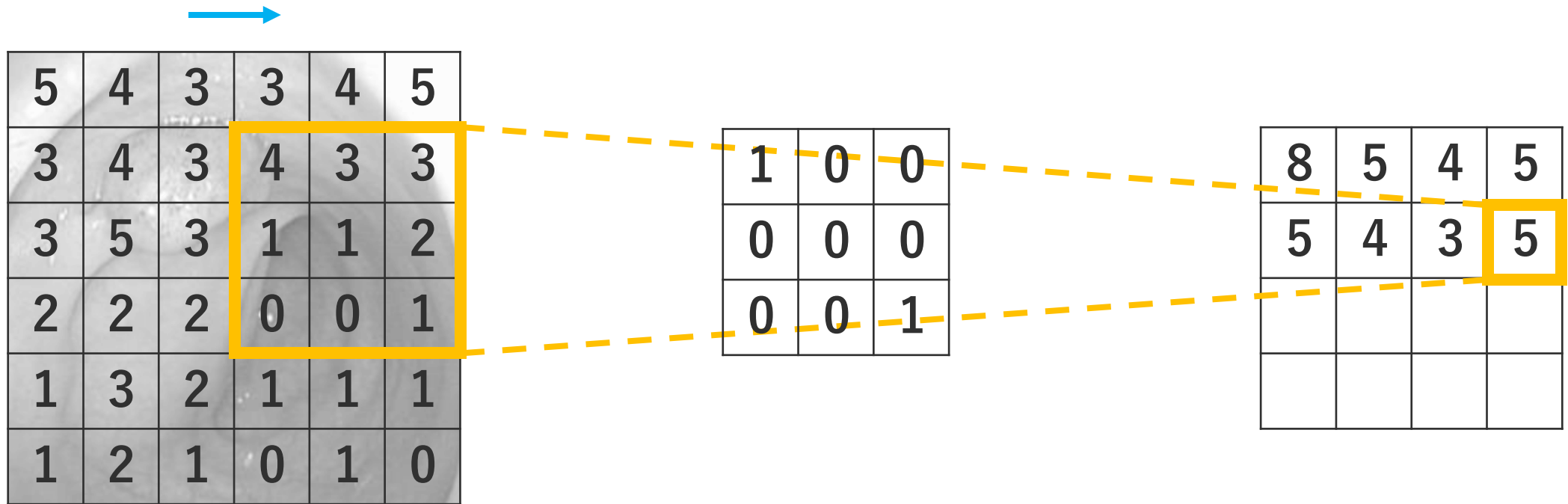
# 畳み込み演算



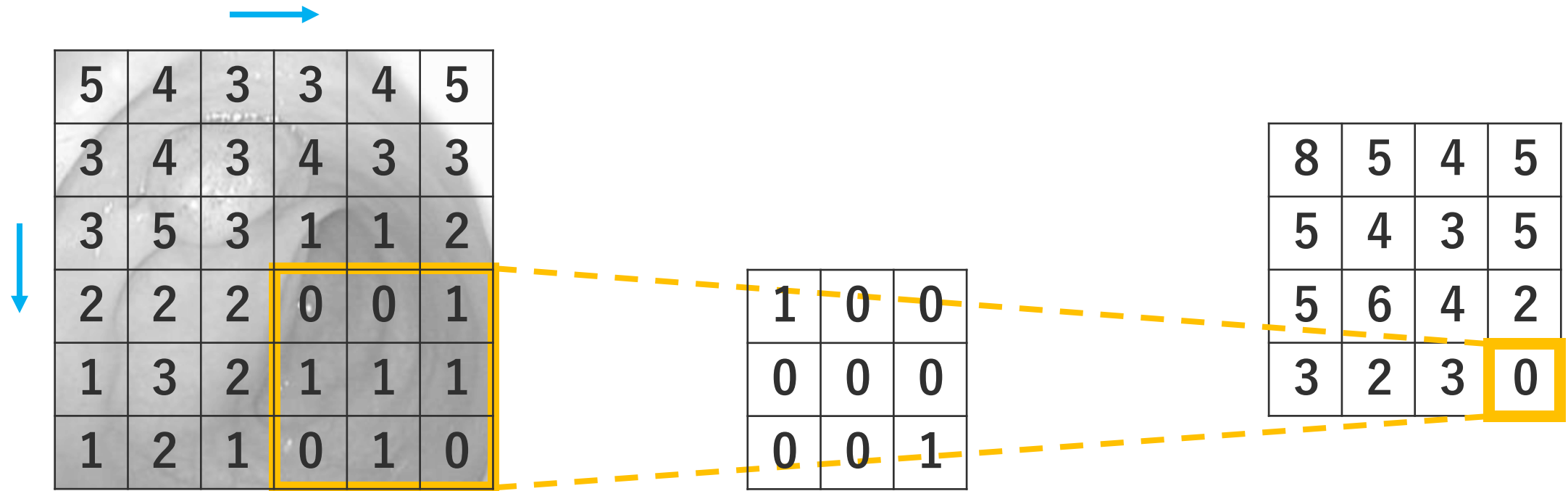
# 畳み込み演算



# 畳み込み演算

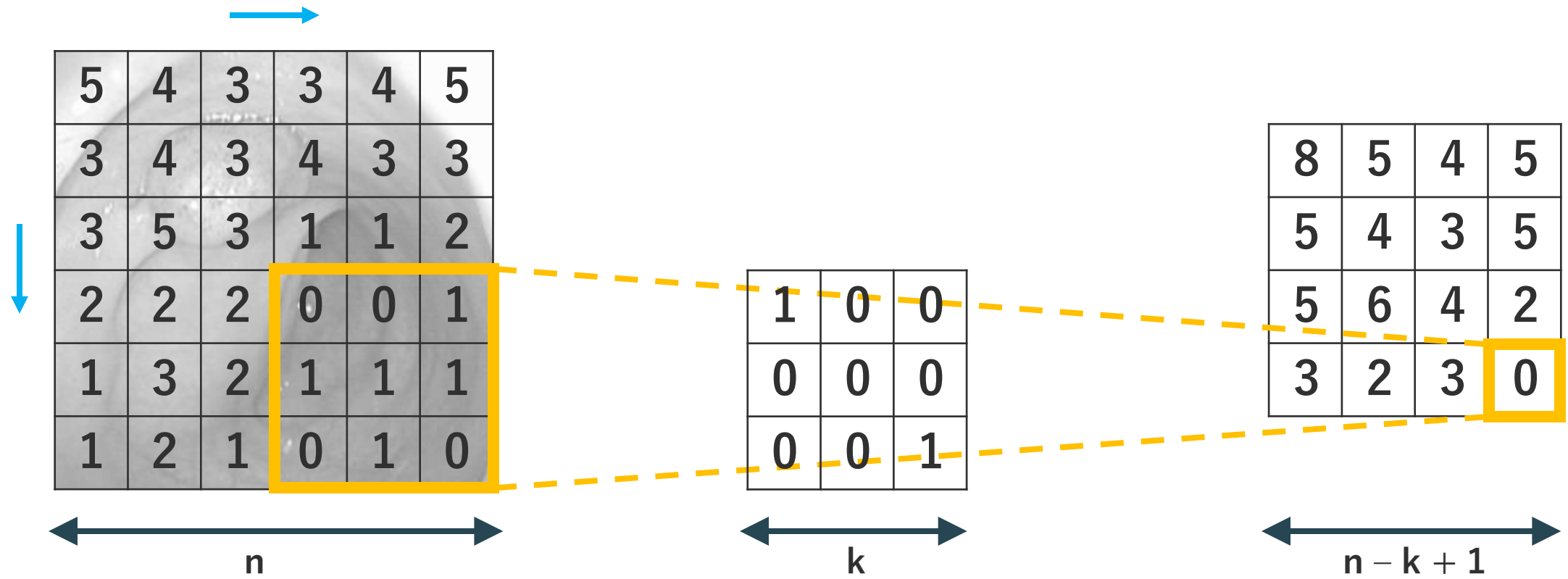


# 畳み込み演算

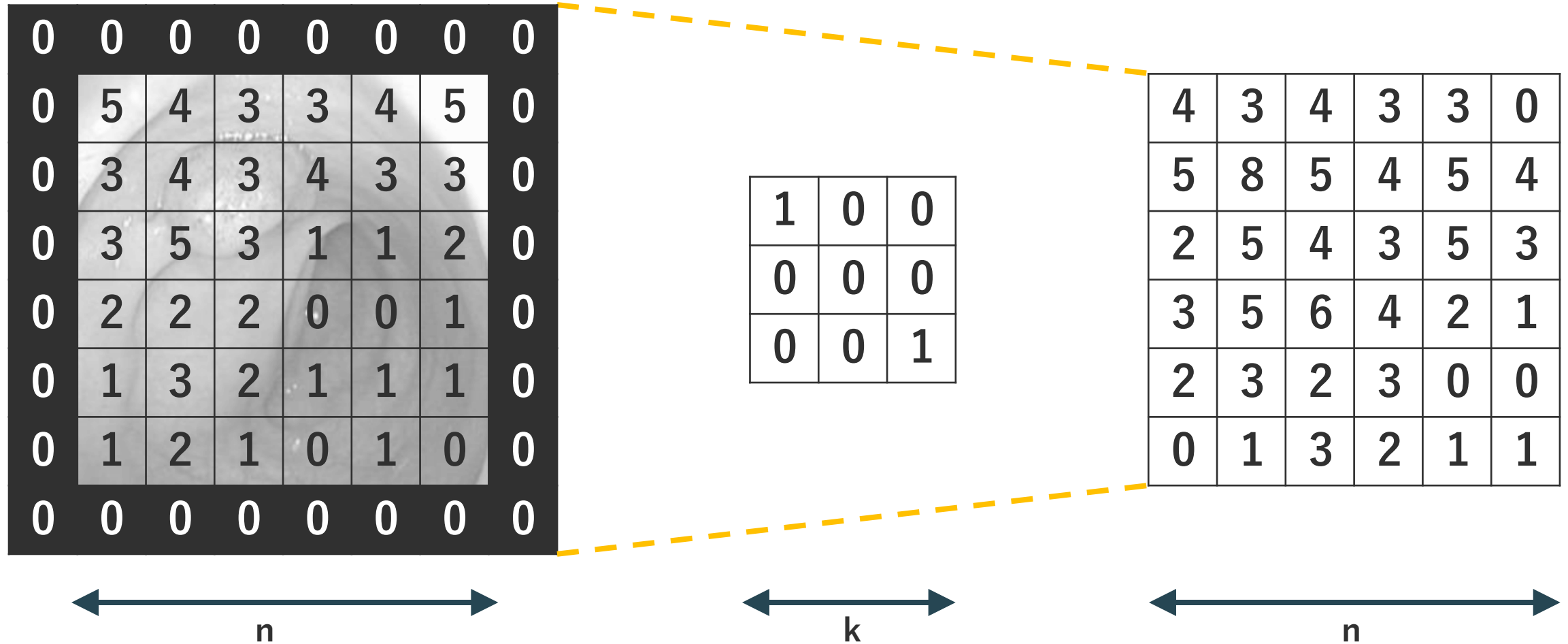




# 畳み込み演算



# 畳み込み演算 (zero-padding)



# 畳み込み演算

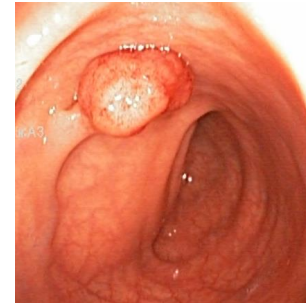
## Grayscale Photo



```
162 149 144 131 124
163 202 123 101 143
132 121 146 150 142
178 178 183 129 126
144 125 135 112 171
241 201 191 122 120
132 152 152 137 121
145 150 171 104 143
152 151 154 114 141
151 149 123 101 153
```

1 channel

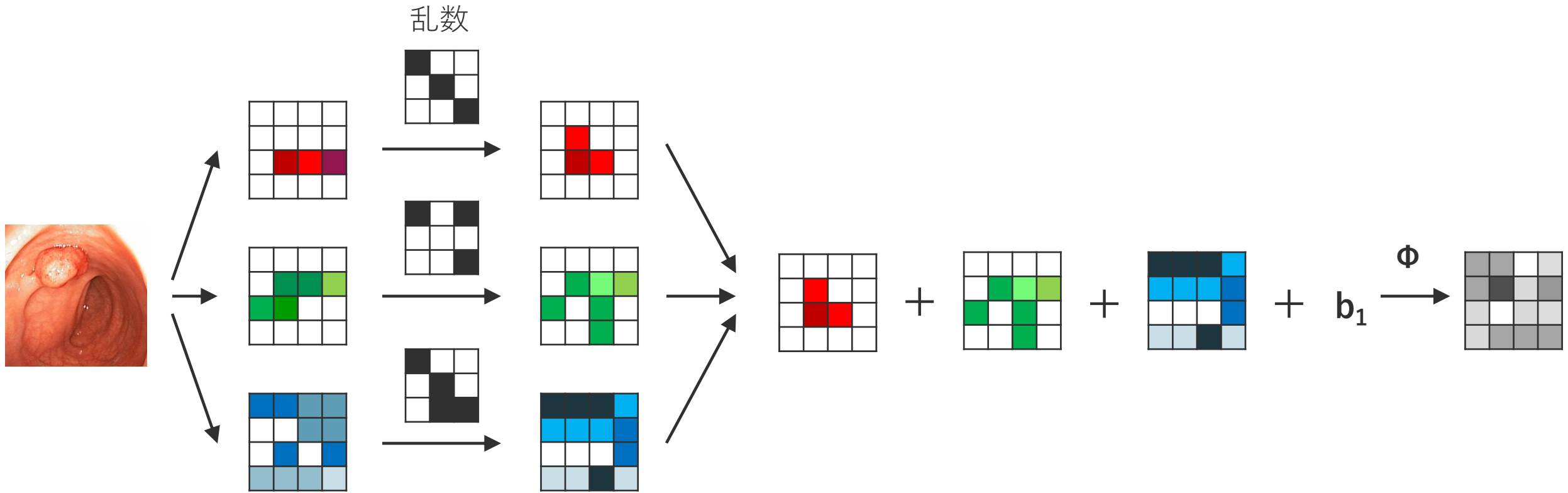
## Color Photo



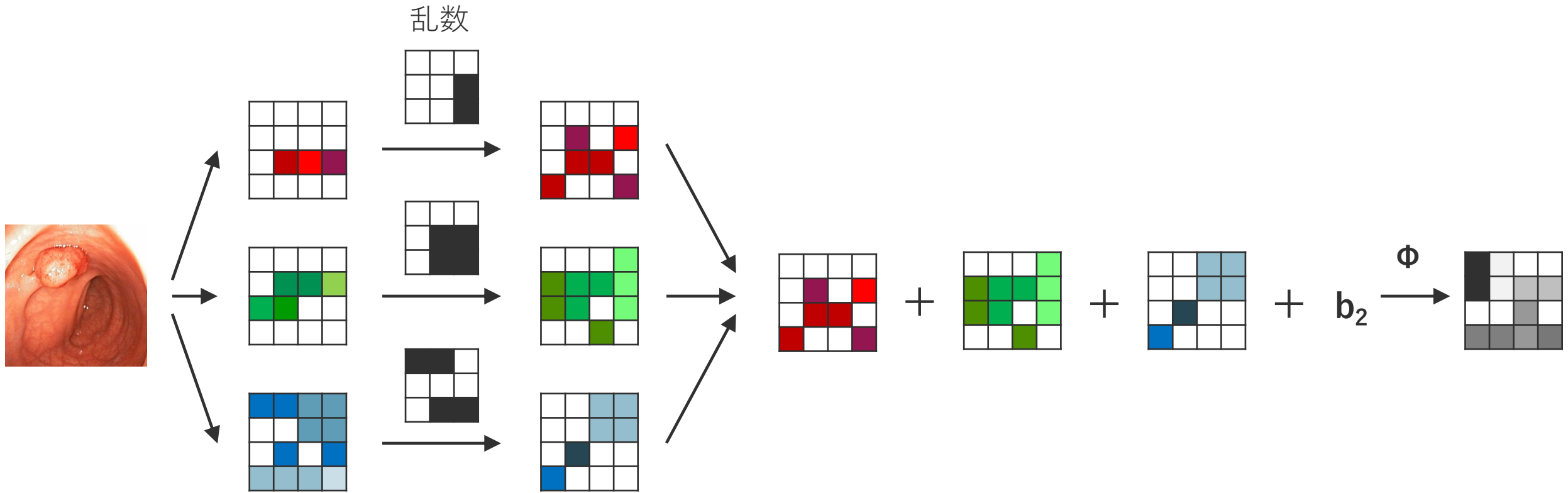
```
021 031 022 063 042
041 065 034 044 033
162 149 144 131 124
082 024 056 072 038
163 202 123 101 143
046 046 033 049 029
025 061 010 030 052
178 178 183 129 126
031 051 053 081 041
144 125 135 112 171
093 049 058 094 026
241 201 191 122 120
054 104 061 024 041
132 152 152 137 121
109 023 020 053 022
111 042 024 021 024
152 151 154 114 141
151 149 123 101 153
251 254 255 245 253
254 253 250 249 251
```

3 channels

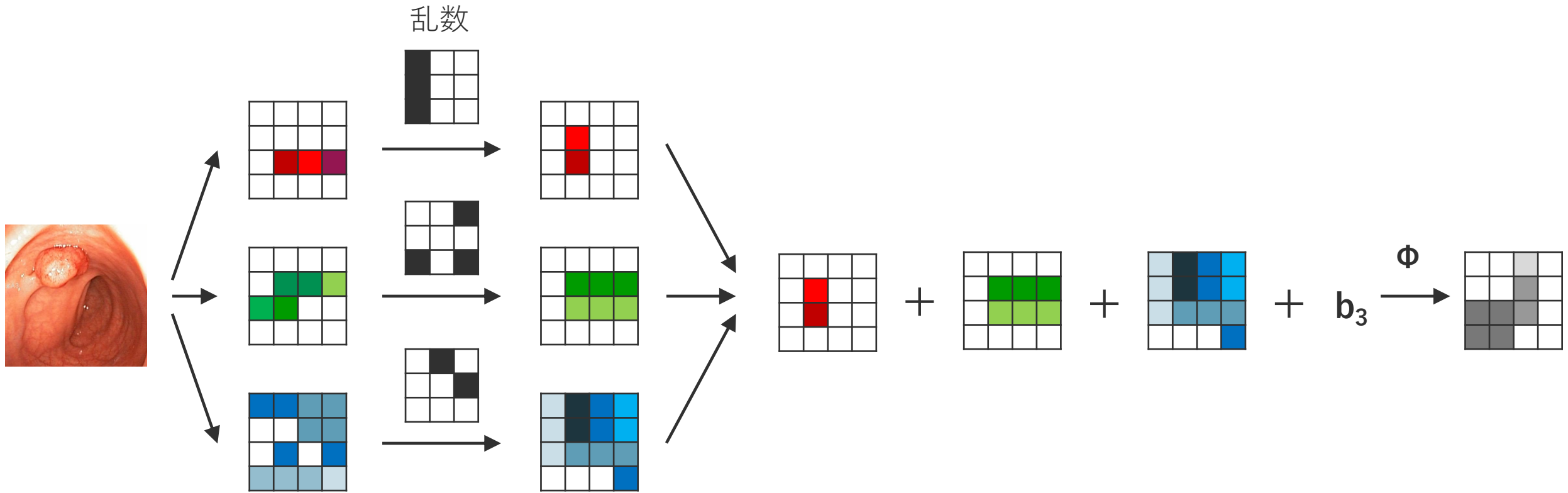
# 畳み込み演算



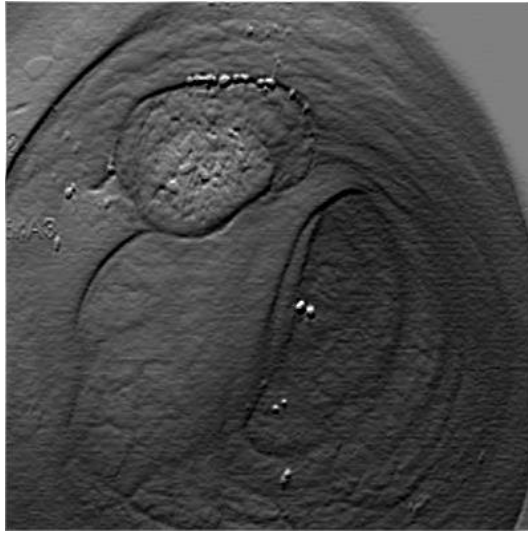
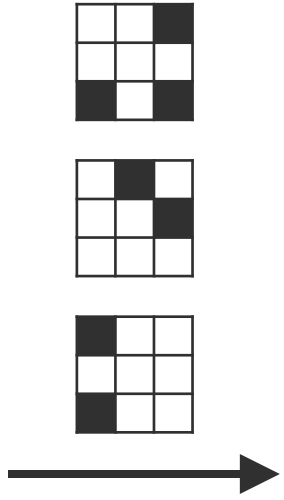
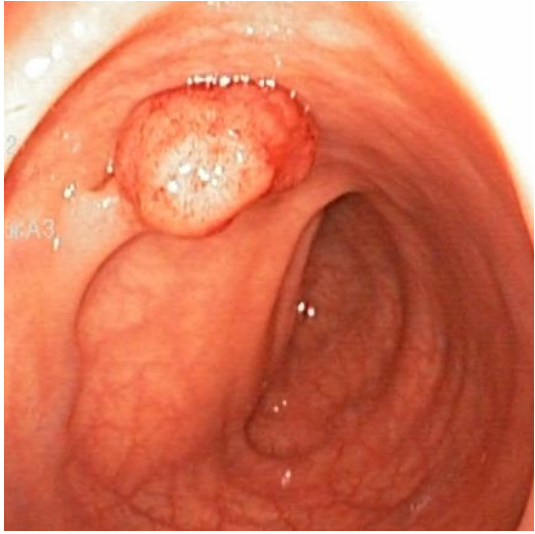
# 畳み込み演算



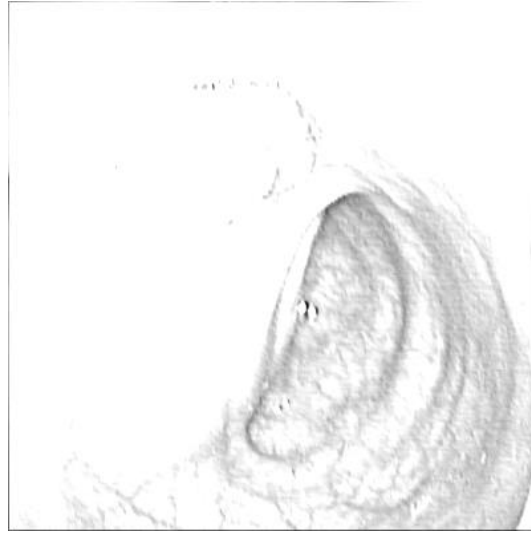
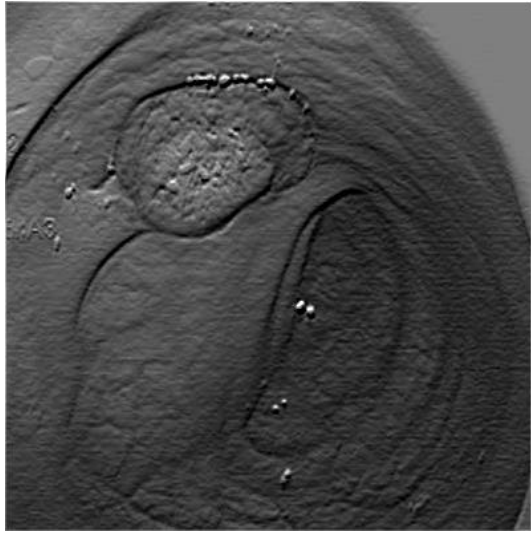
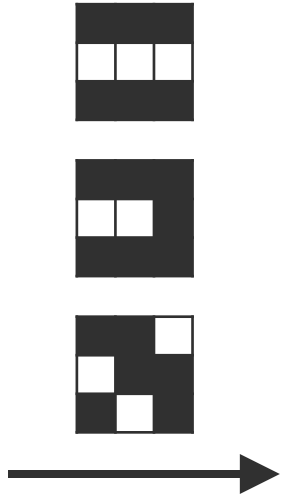
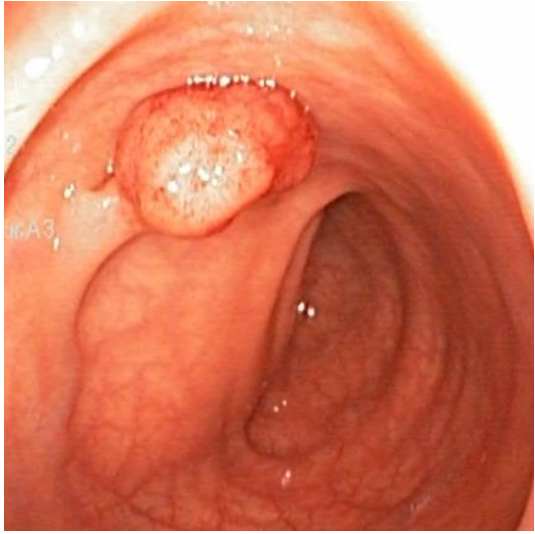
# 畳み込み演算



# 畳み込み演算

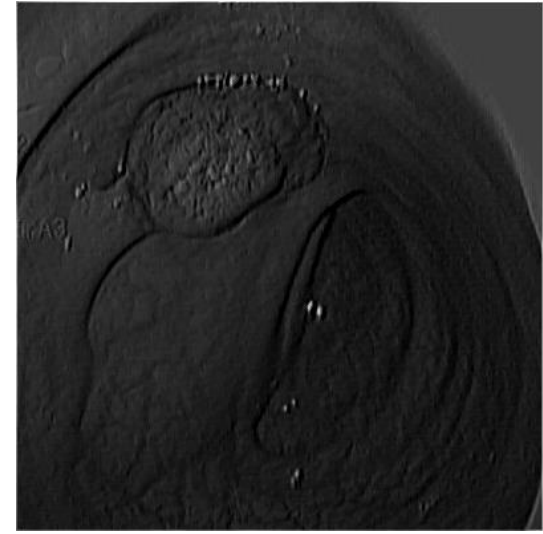
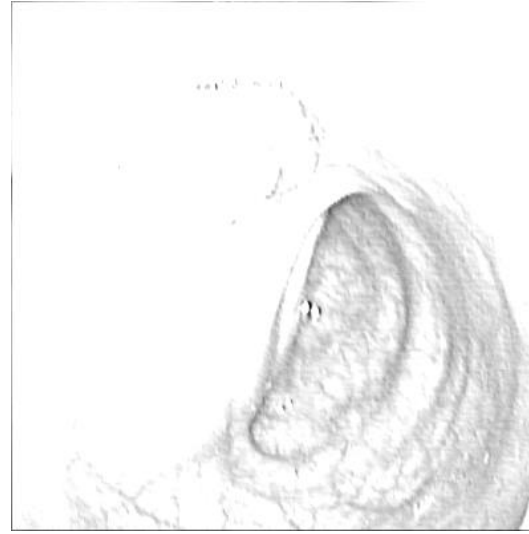
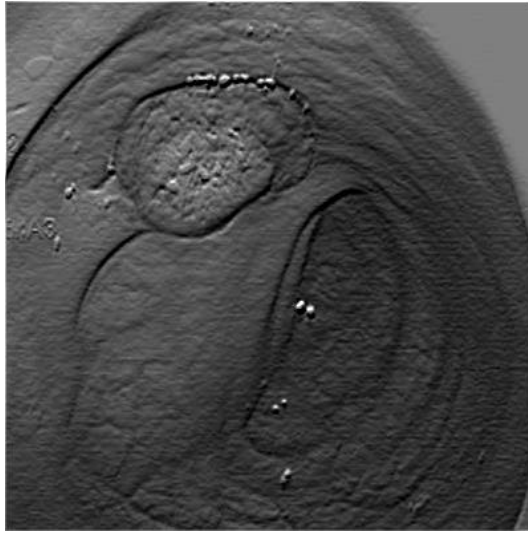
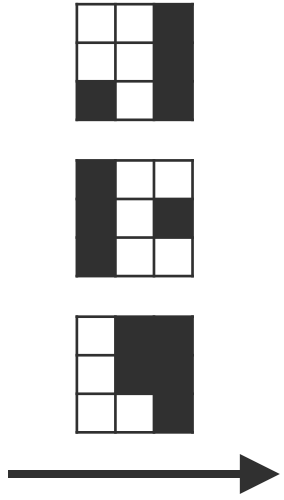
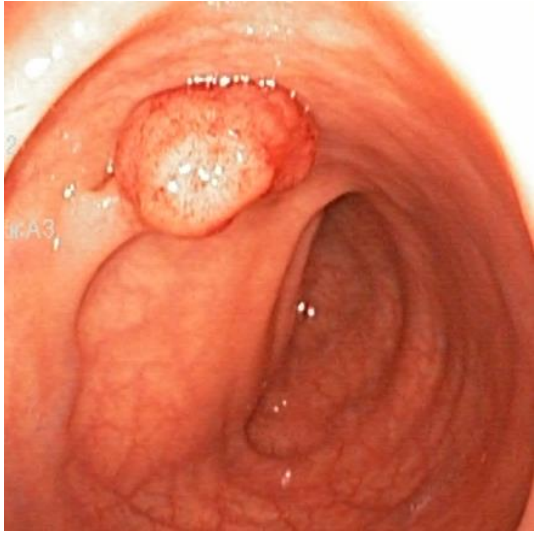


# 畳み込み演算

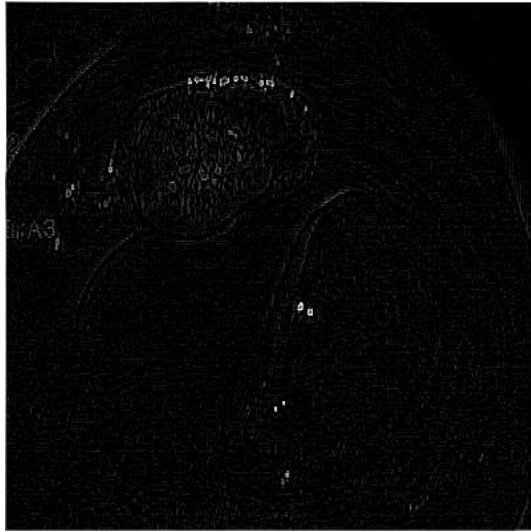
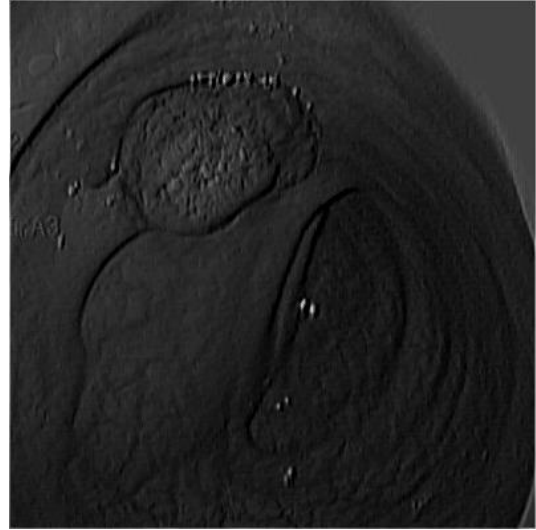
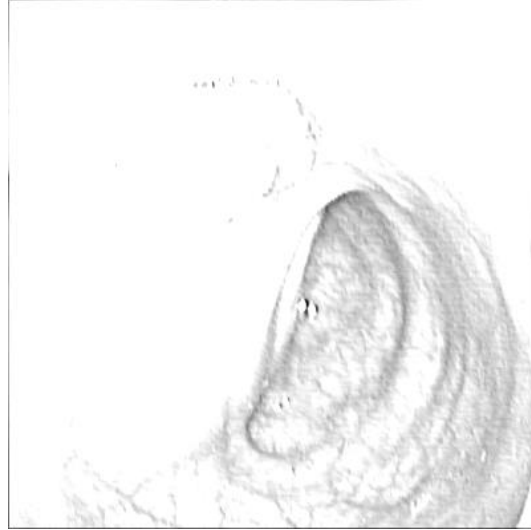
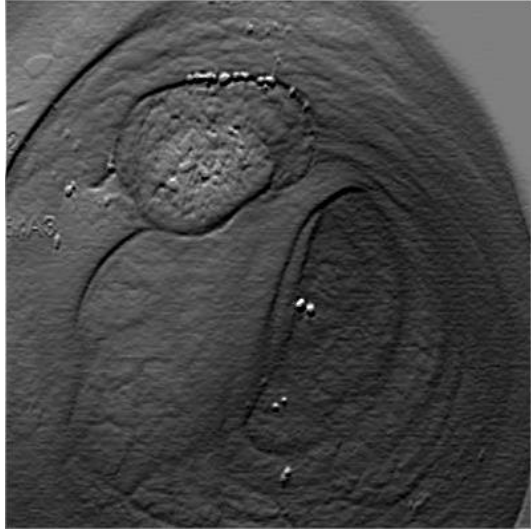
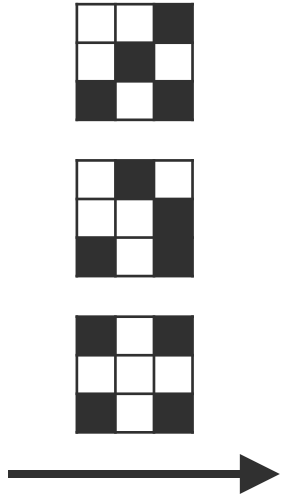
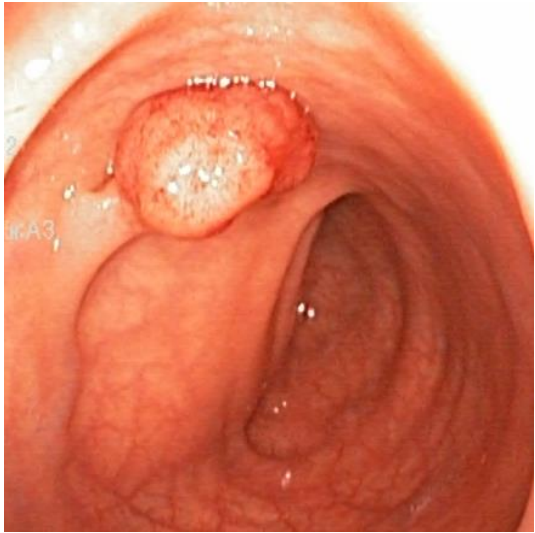




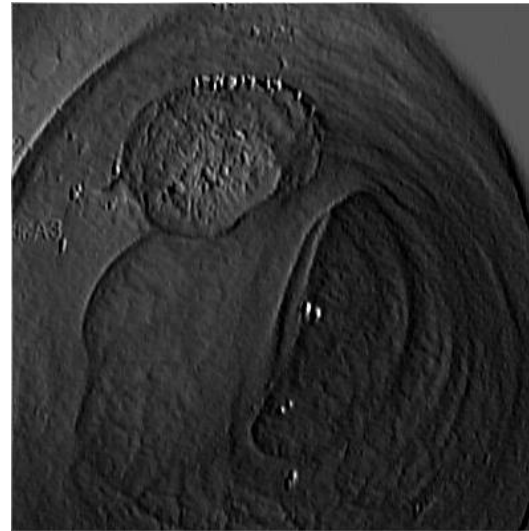
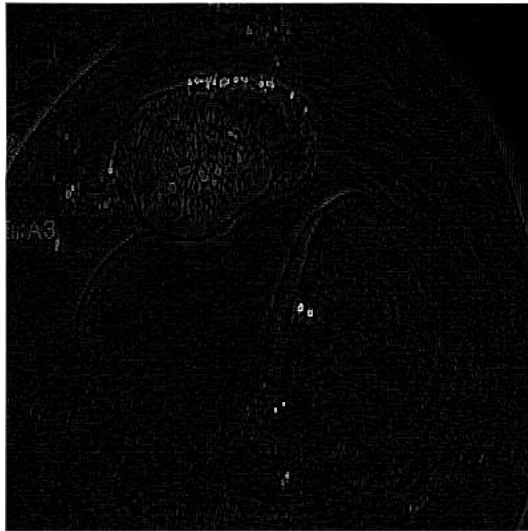
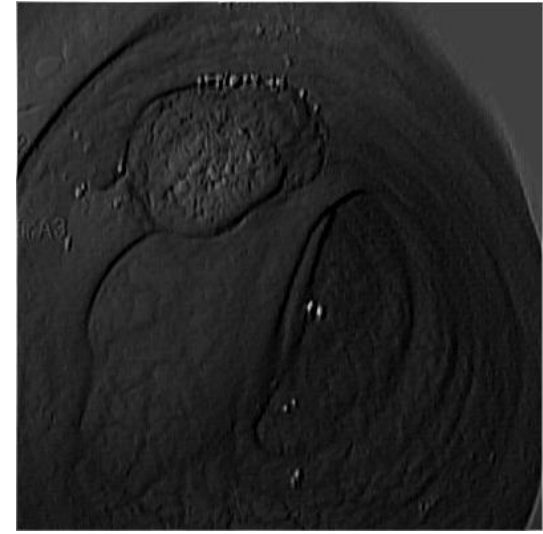
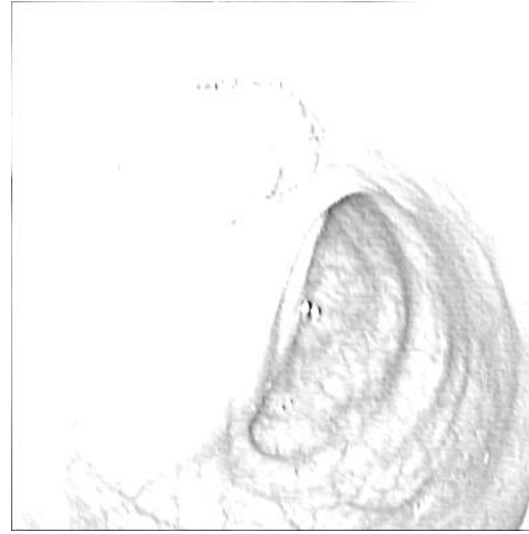
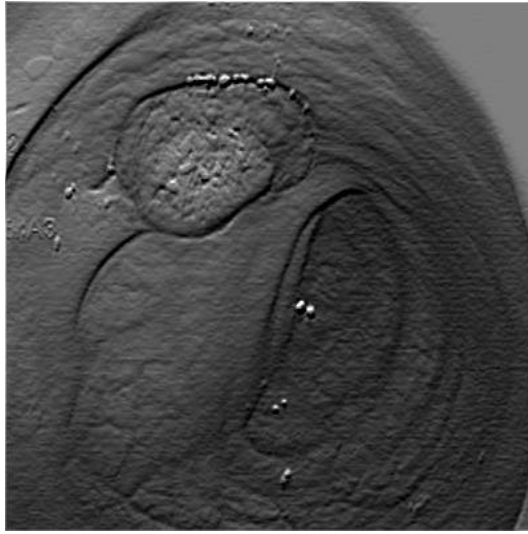
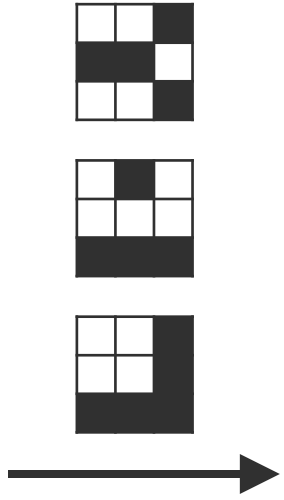
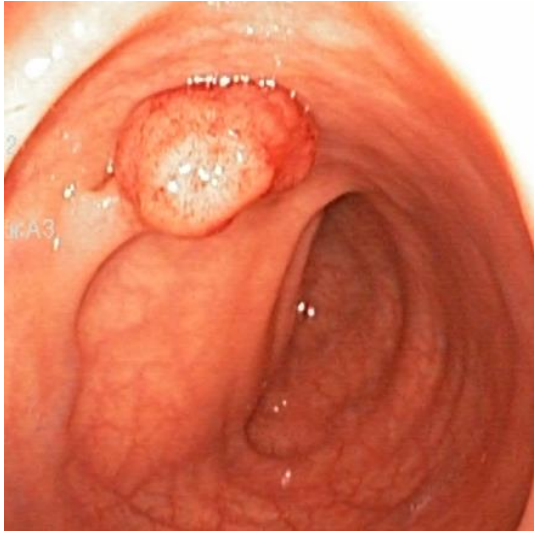
# 畳み込み演算



# 畳み込み演算

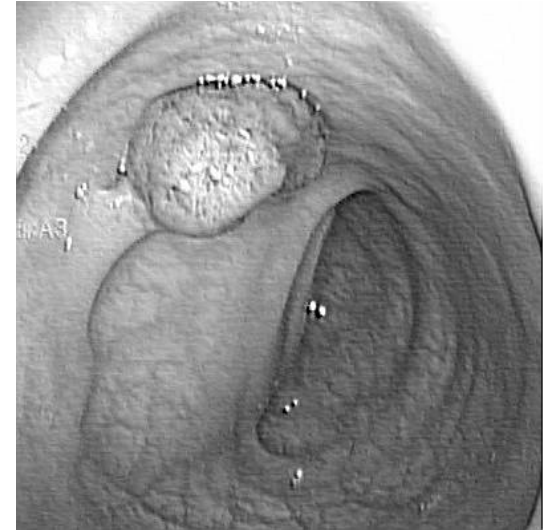
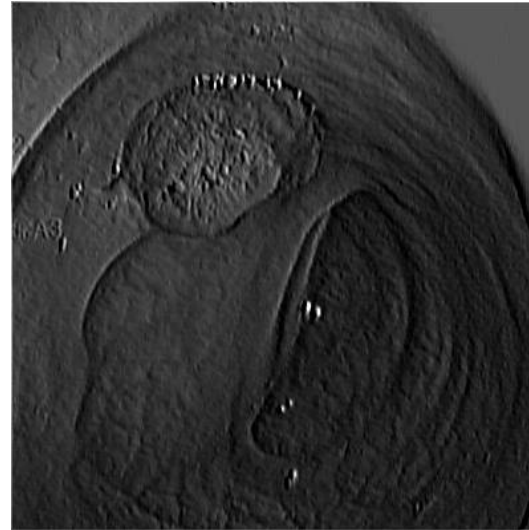
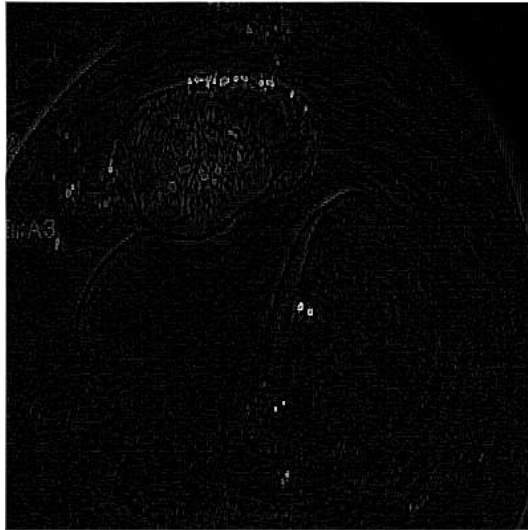
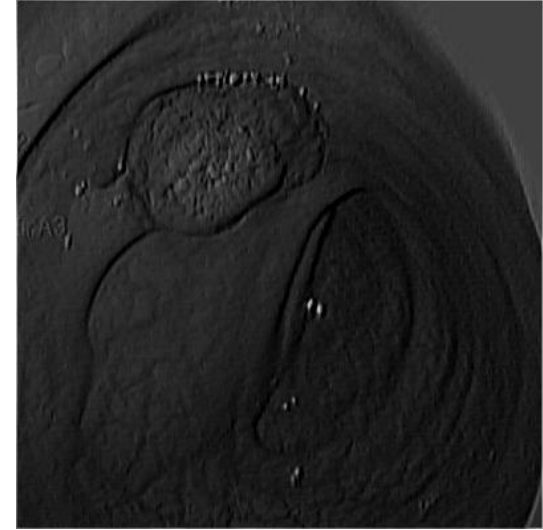
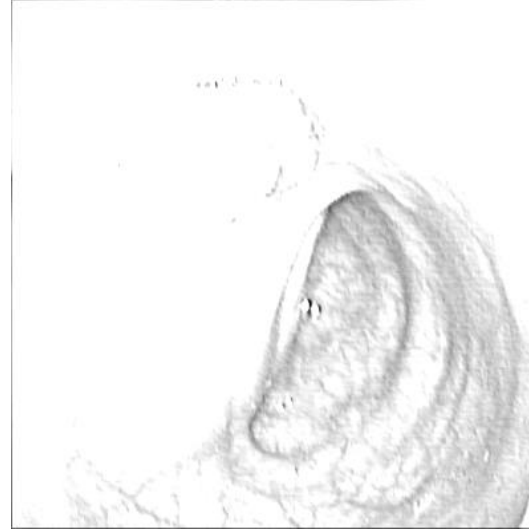
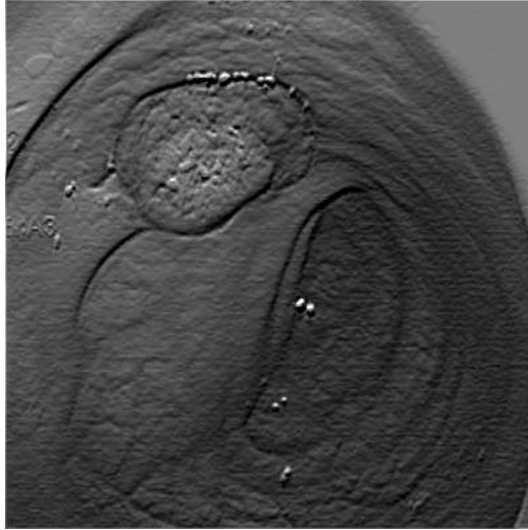
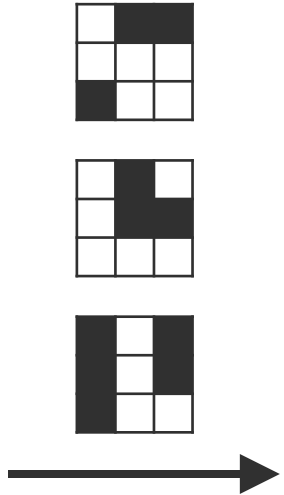
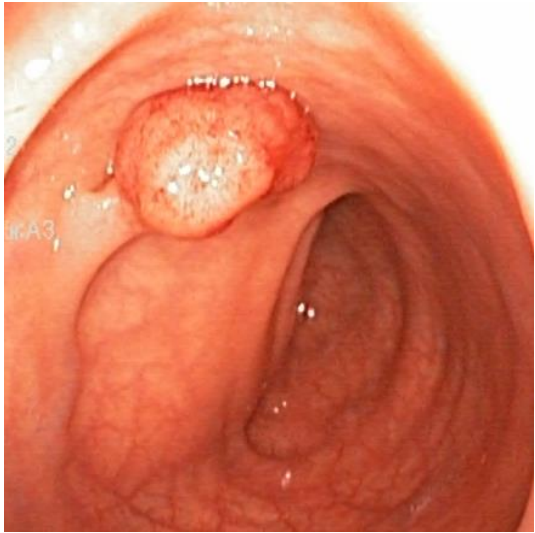


# 畳み込み演算

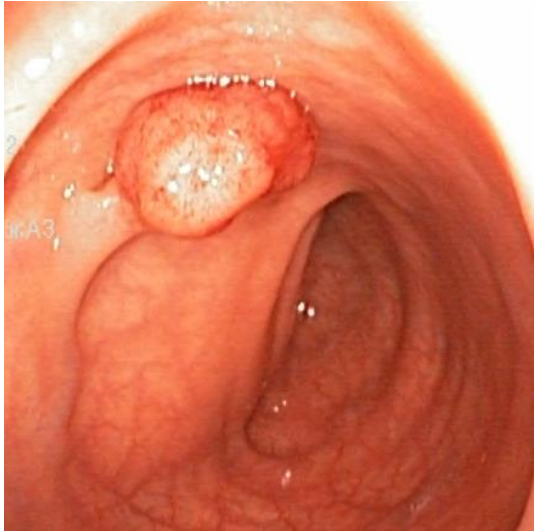




# 畳み込み演算



# 畳み込み演算

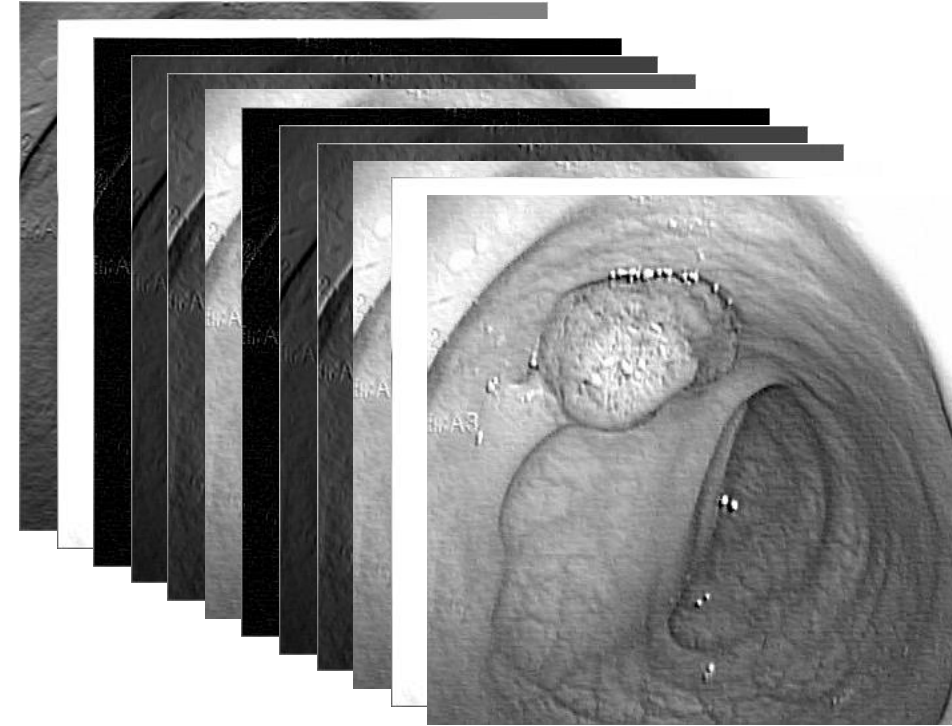


入力画像

畳み込み演算

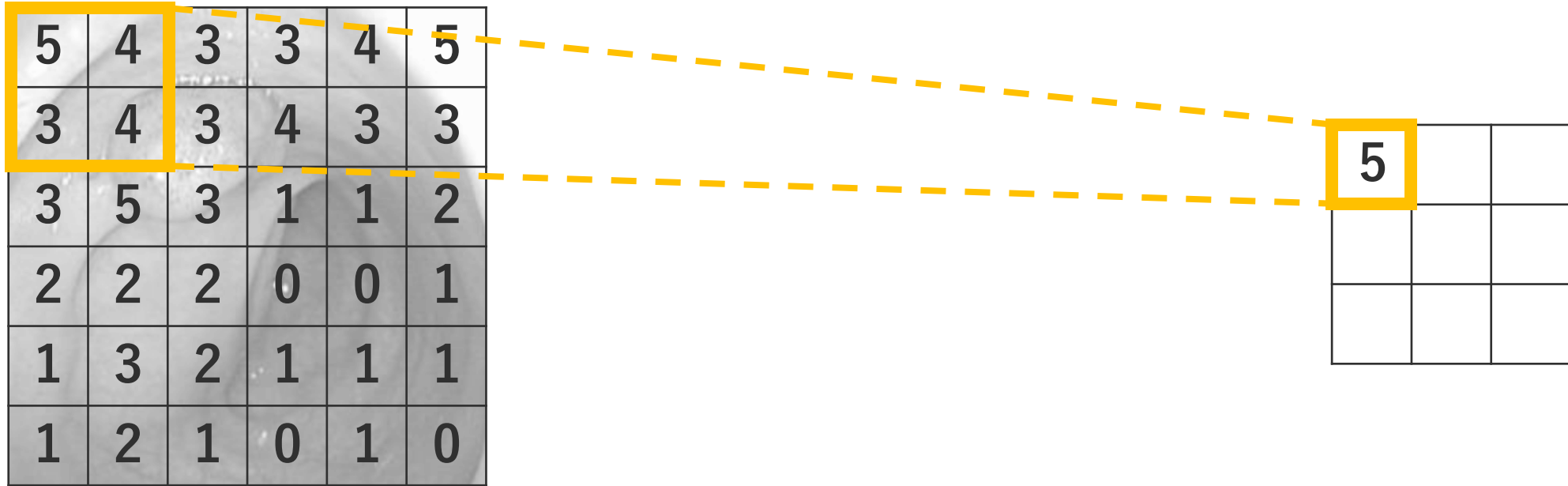


- 特徴の抽出
  - 輪郭・模様など
  - 抽象的な概念など
- 空間情報の保存
- パラメータの削減

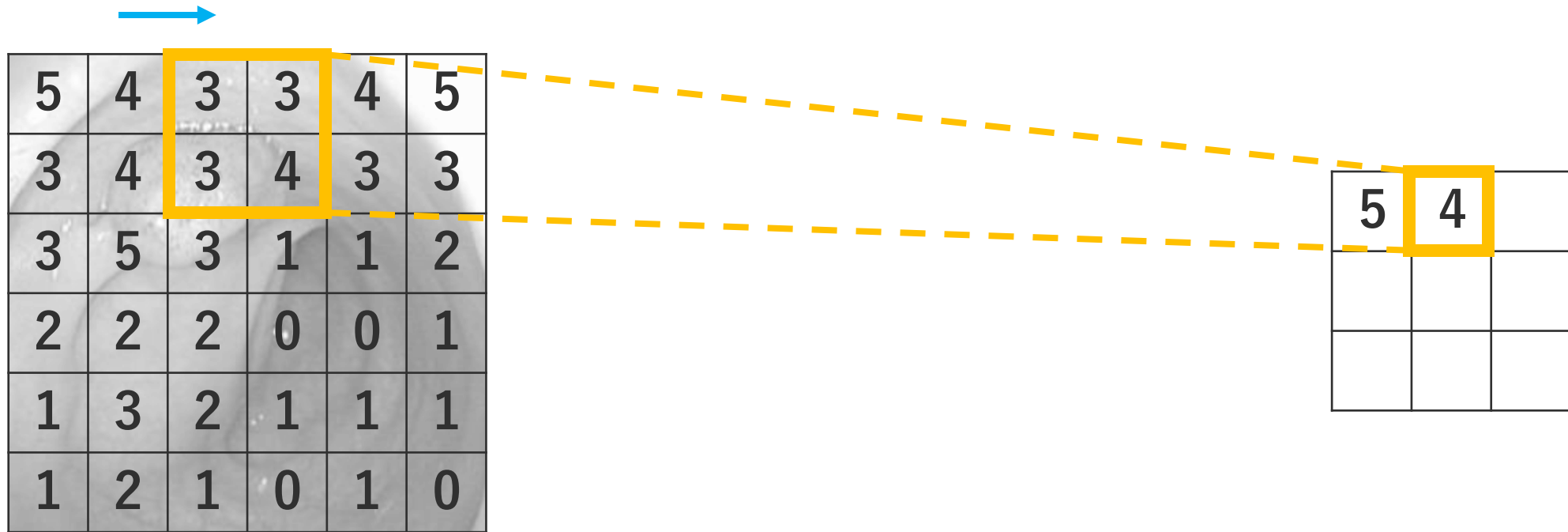


出力画像  
(特徴マップ)

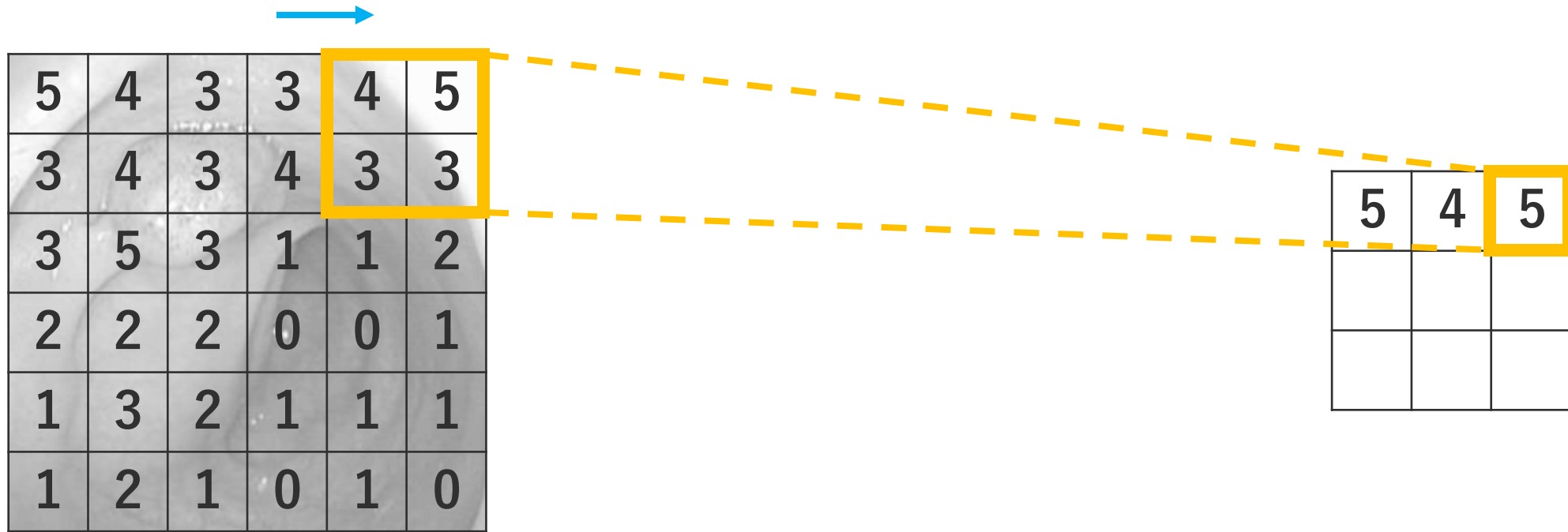
# 最大プーリング演算



# 最大プーリング演算

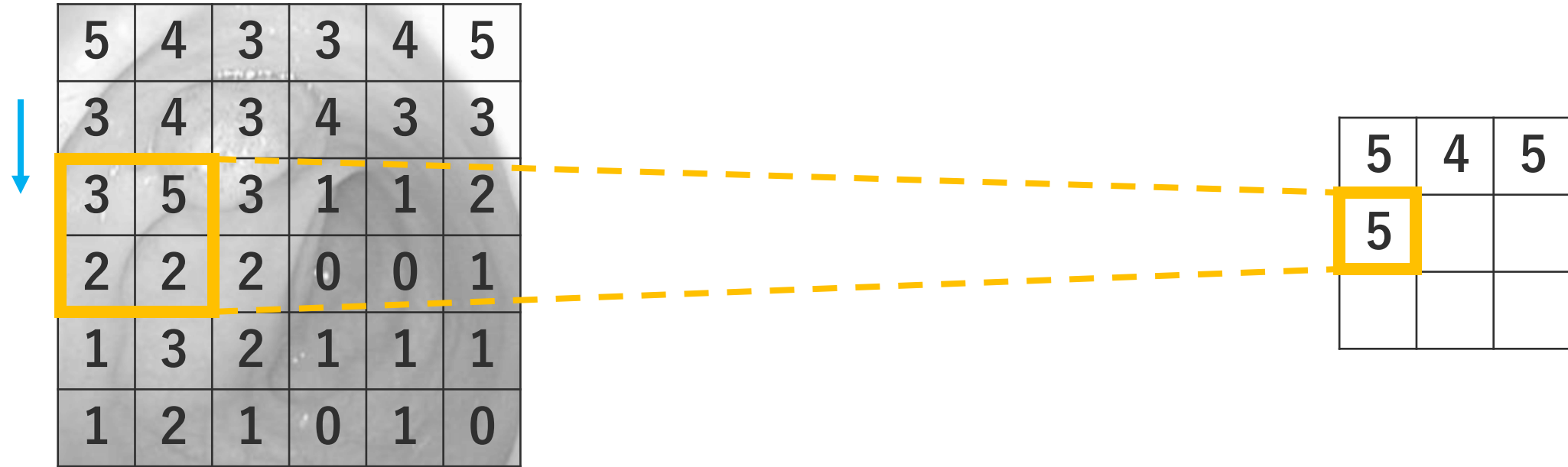


# 最大プーリング演算

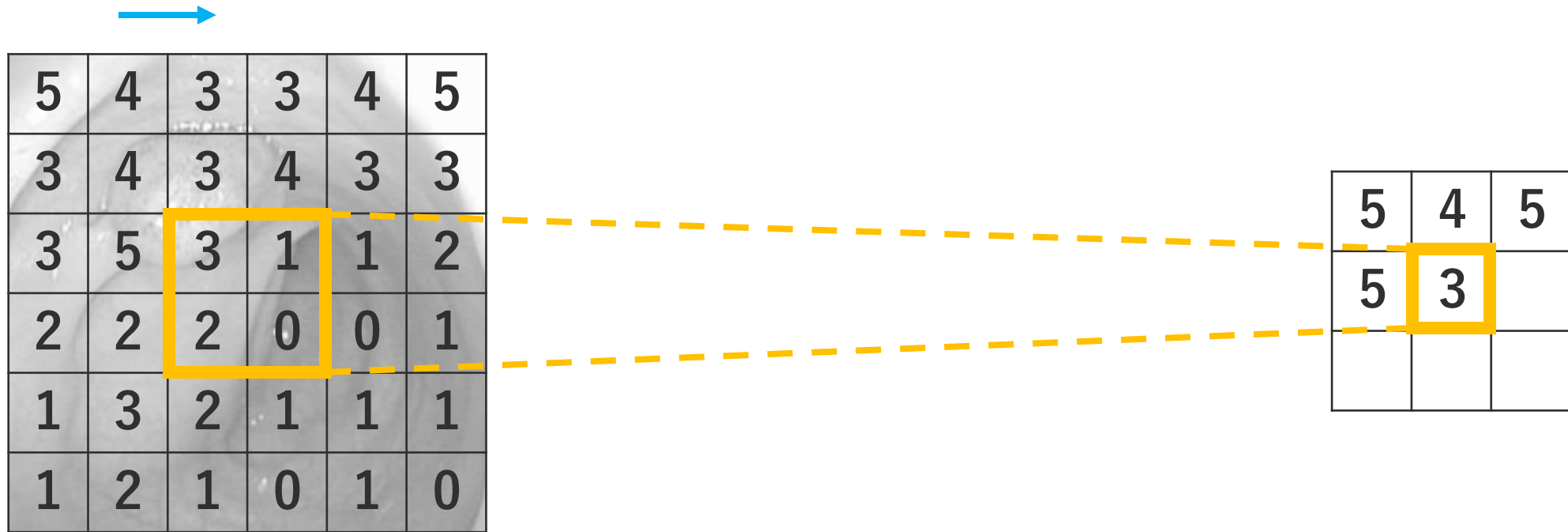




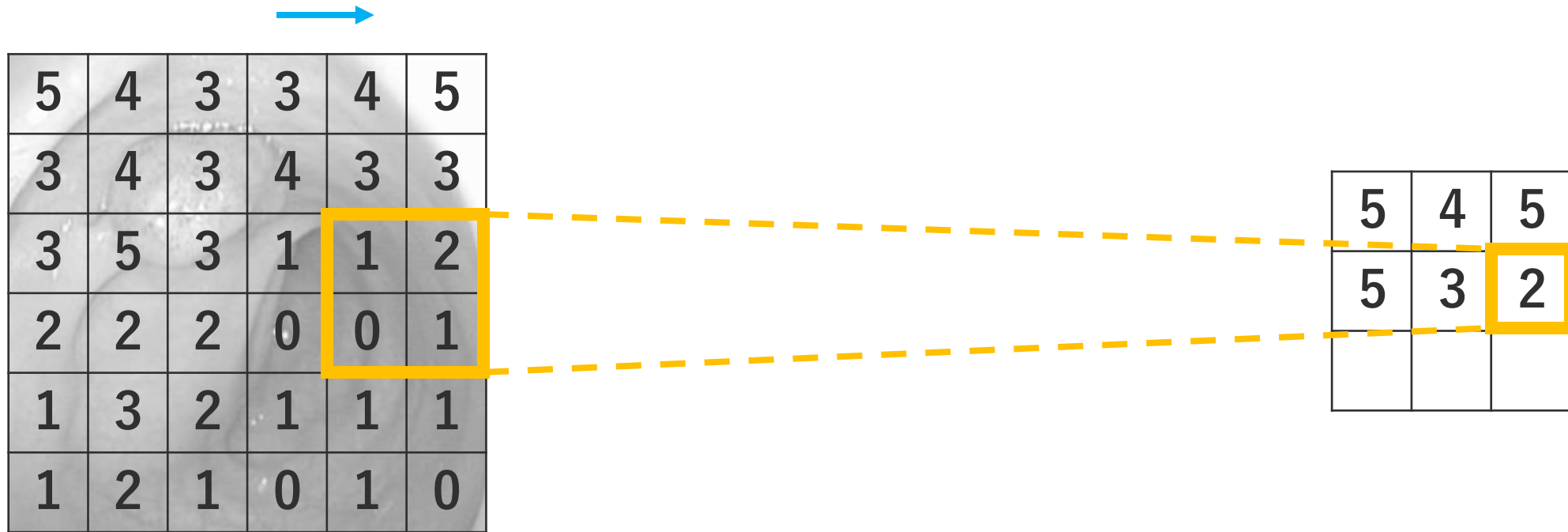
# 最大プーリング演算



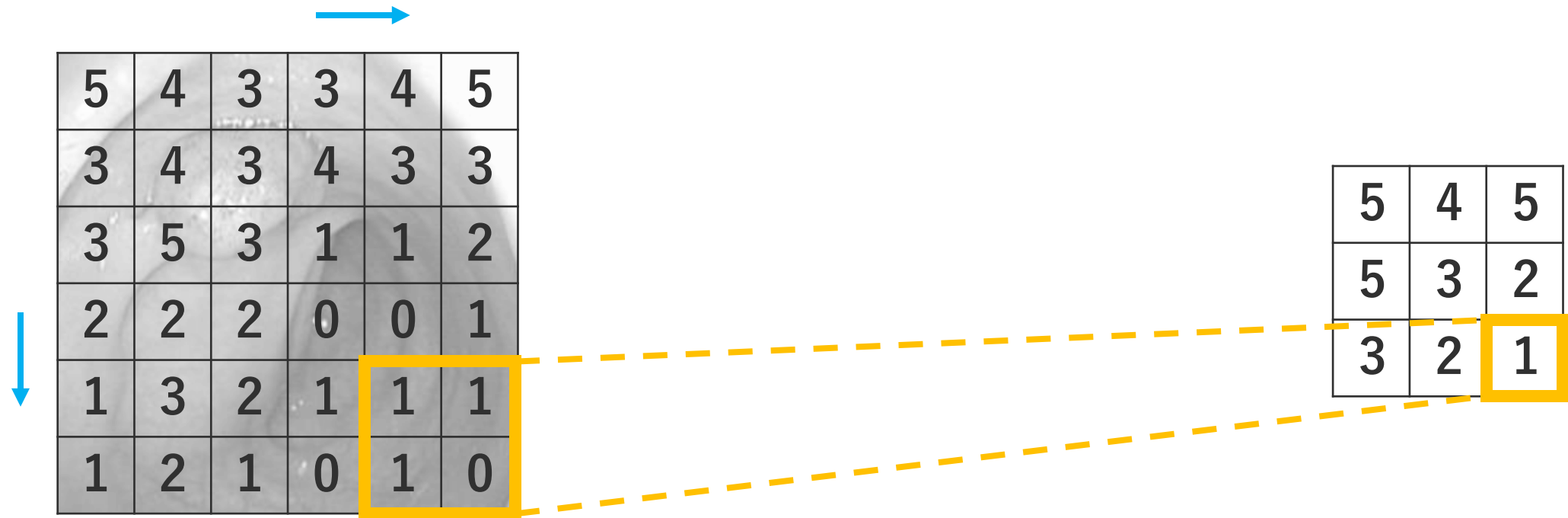
# 最大プーリング演算



# 最大プーリング演算



# 最大プーリング演算



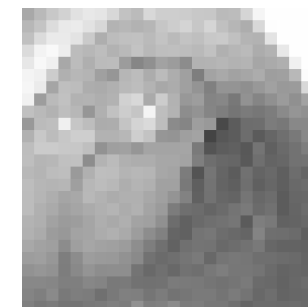
# 最大プーリング演算

---



入力画像

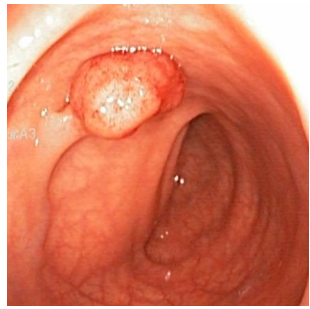
最大プーリング演算



出力画像

- 特徴マップの圧縮
- パラメータの削減
- 平均プーリング演算、グローバルプーリング演算などもある

# 物体分類



128 213 154 214 211  
 251 023 049 044 033 024  
 224 052 002 028 021 030 032  
 186 012 023 008 030 031 075  
 215 043 073 028 049 020 032  
 241 134 028 031 017 030 078  
 215 149 058 242 226 020 223  
 223 142 052 051 150 012 198  
 228 223 220 152 122 043 193  
 227 242 253 249 244 254 132  
 211 249 228 251 243 186  
 021 133 111 034 213

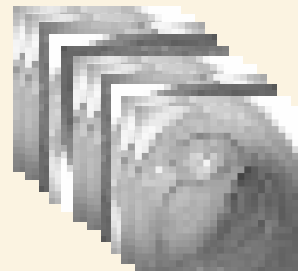
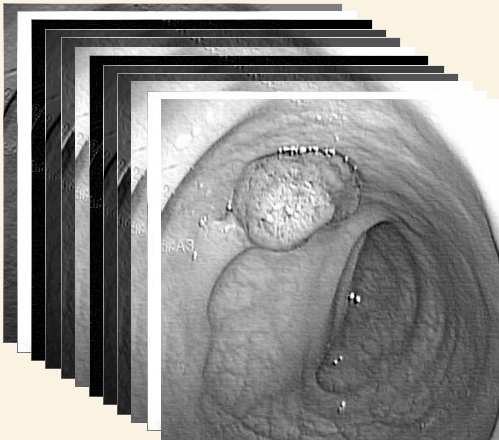
特徴抽出



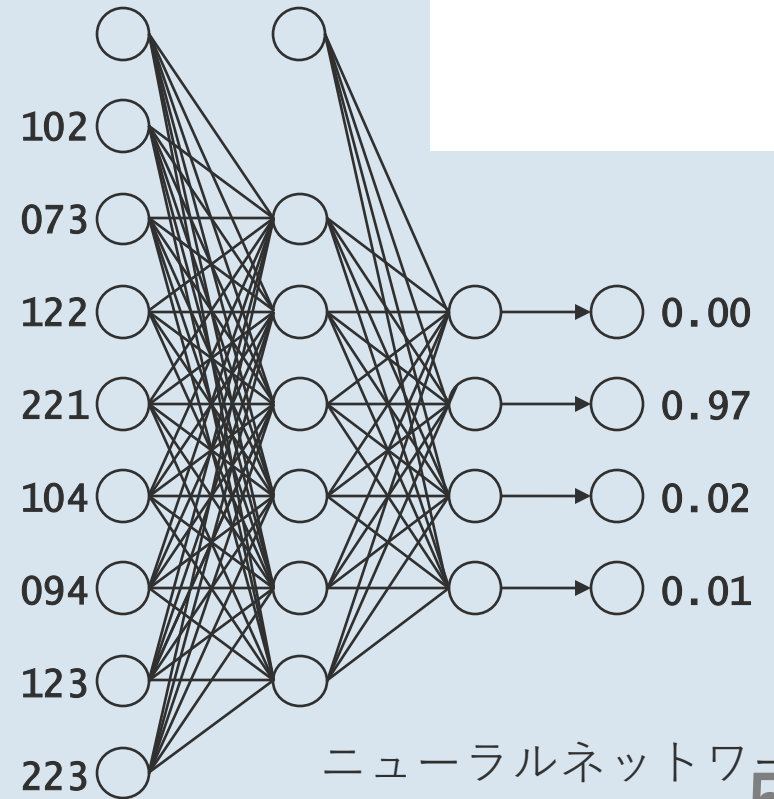
102  
073  
122  
221  
104  
094  
123  
223

分類

normal 0.00  
 polyps 0.97  
 esophagitis 0.02  
 colitis 0.01



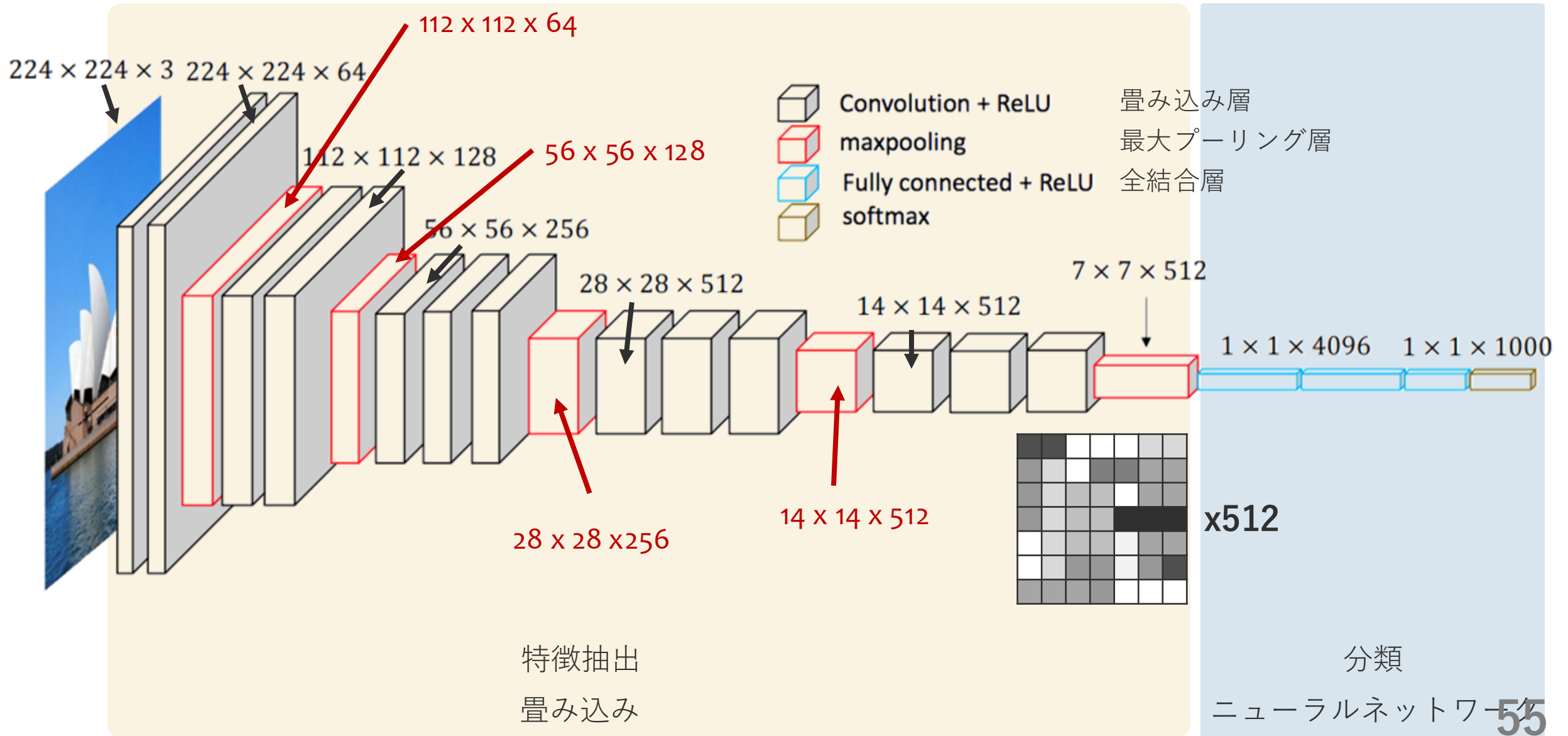
畳み込み



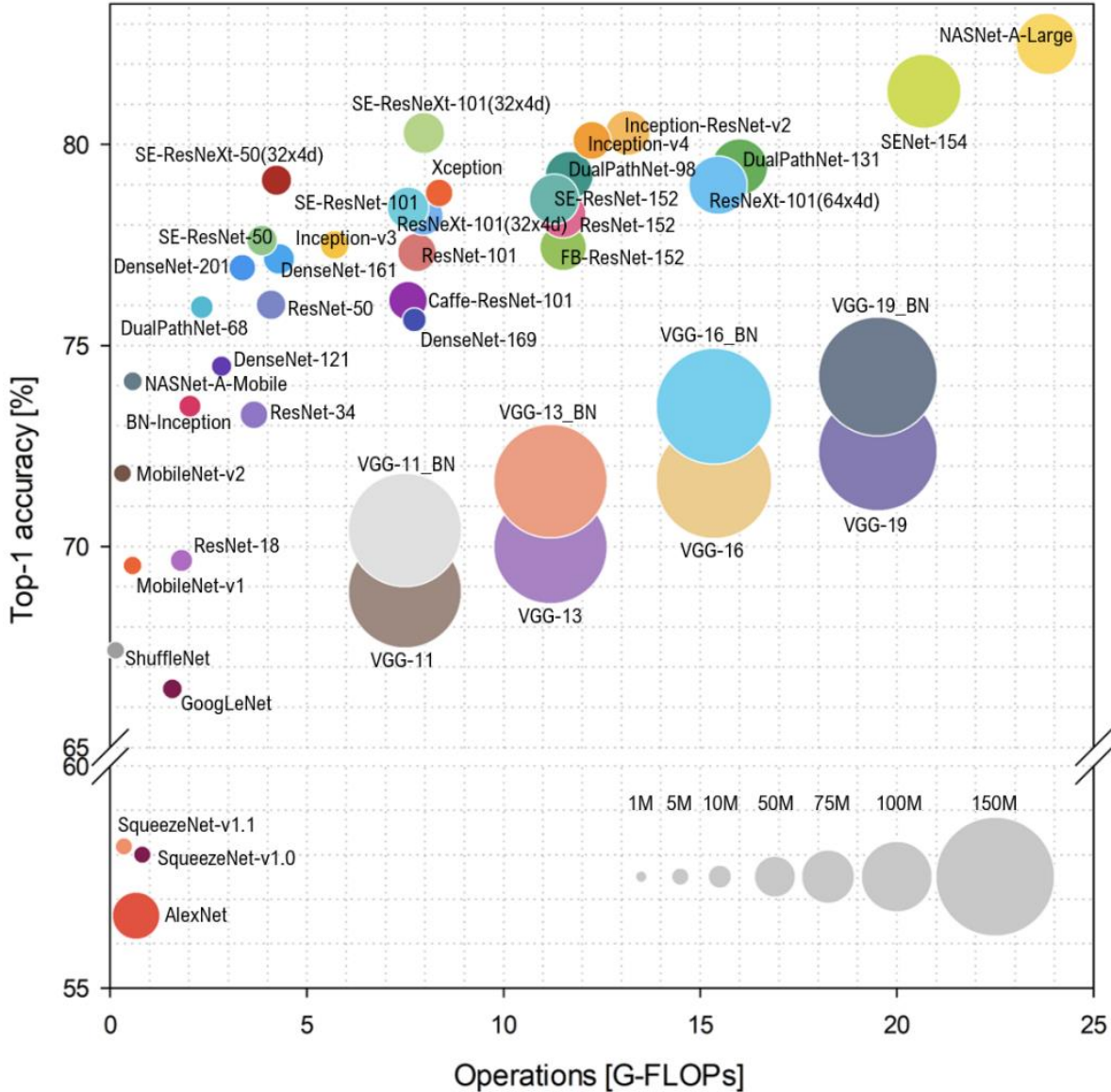
ニューラルネットワーク

# VGG16

<https://doi.org/10.1145/3038912.3052638>



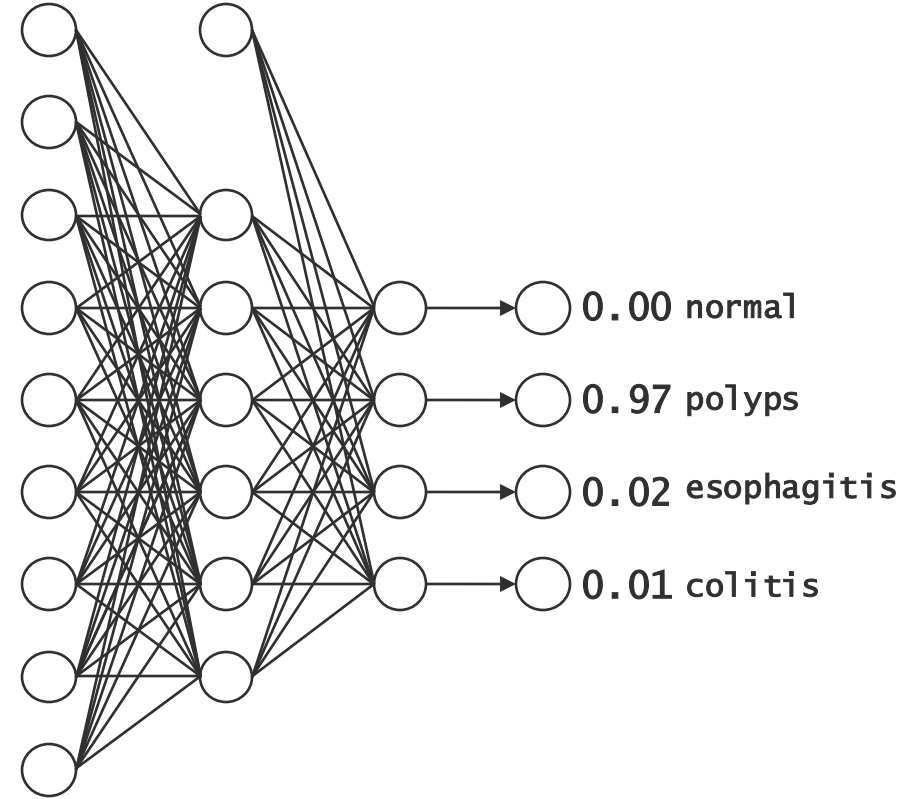
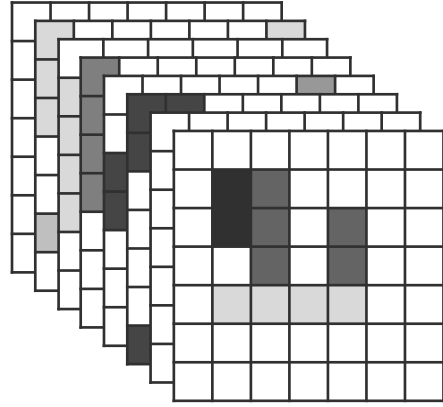
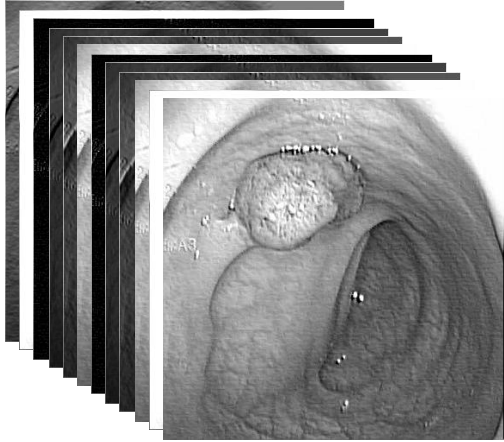
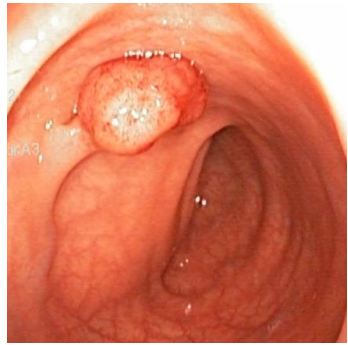
# CNNアーキテクチャ



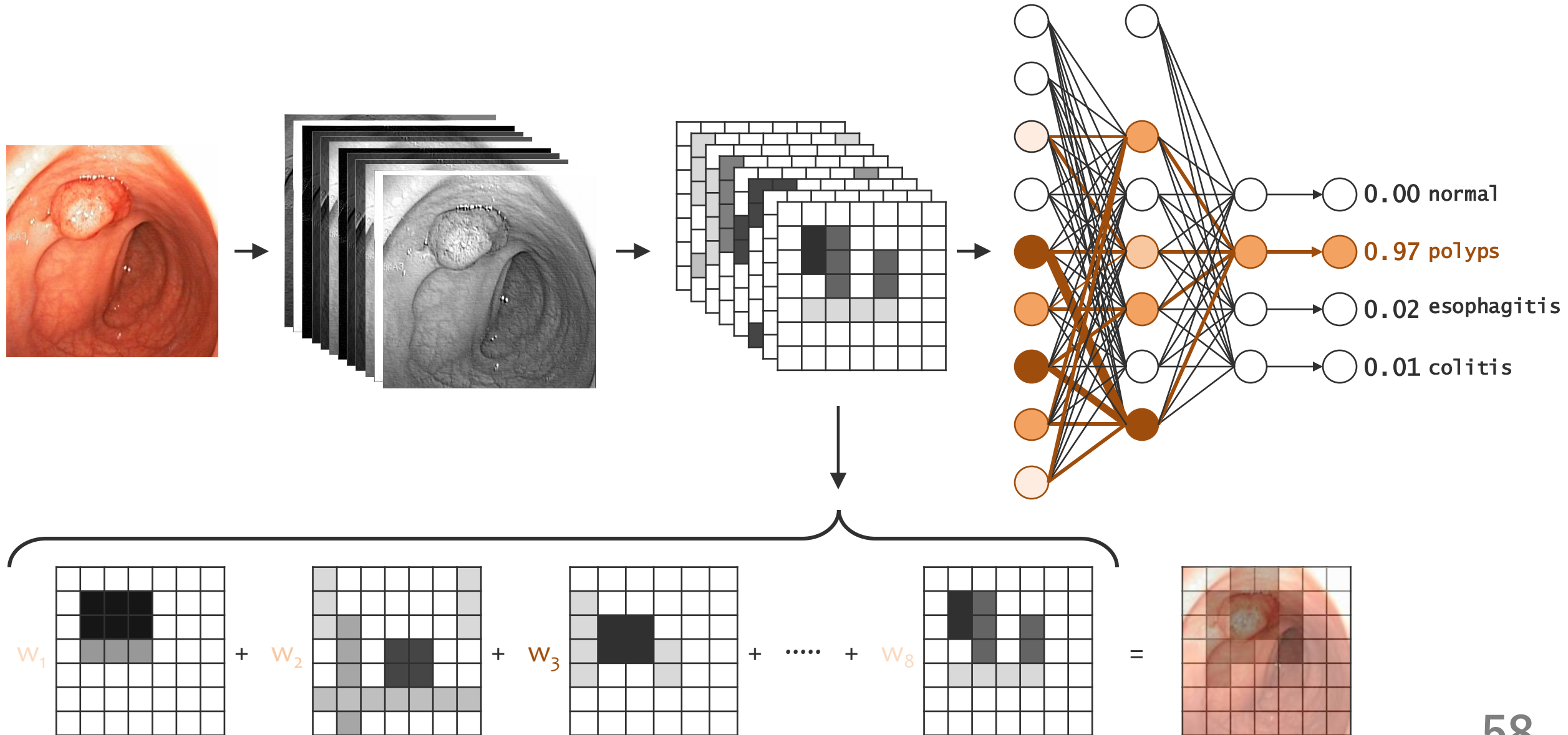
Bianco et al., 2018 <https://doi.org/10.48550/arXiv.1810.00736>



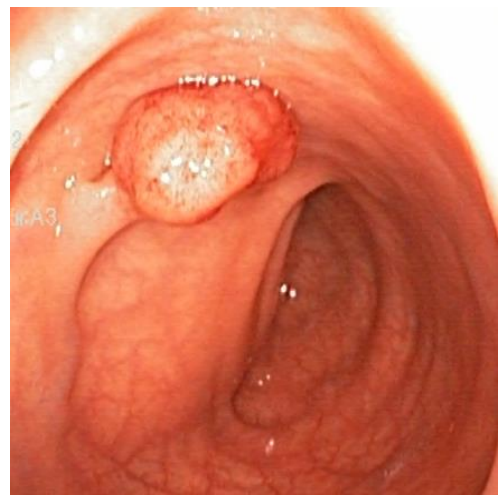
# 判定根拠の可視化



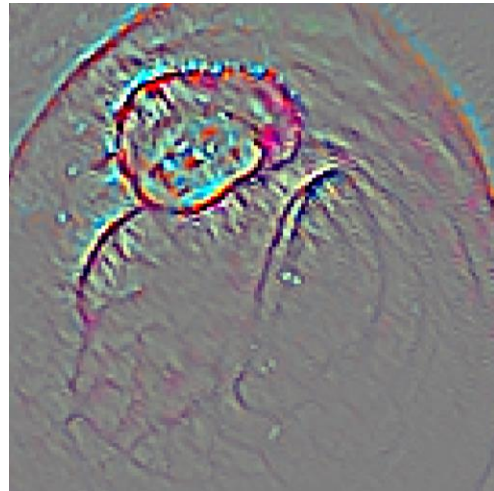
# 判定根拠の可視化



# 判定根拠の可視化

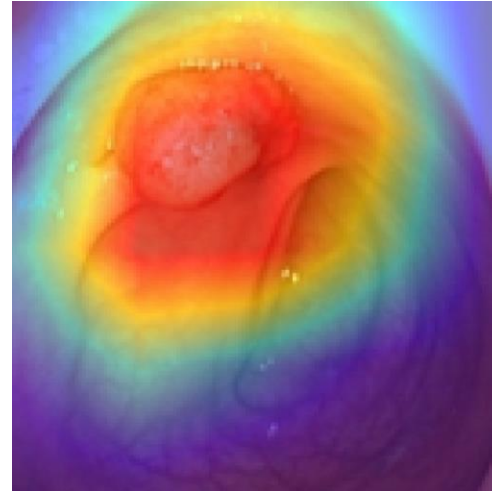


guided backpropagation



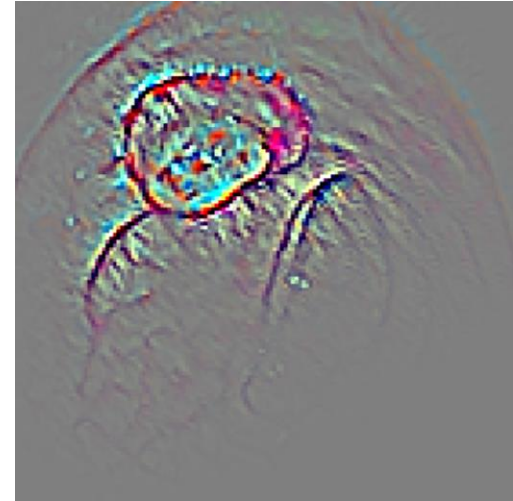
+

Grad-CAM



=

Guided Grad-CAM



# 深層学習

---

ニューラルネットワーク

**畳み込みニューラルネットワーク**

物体分類

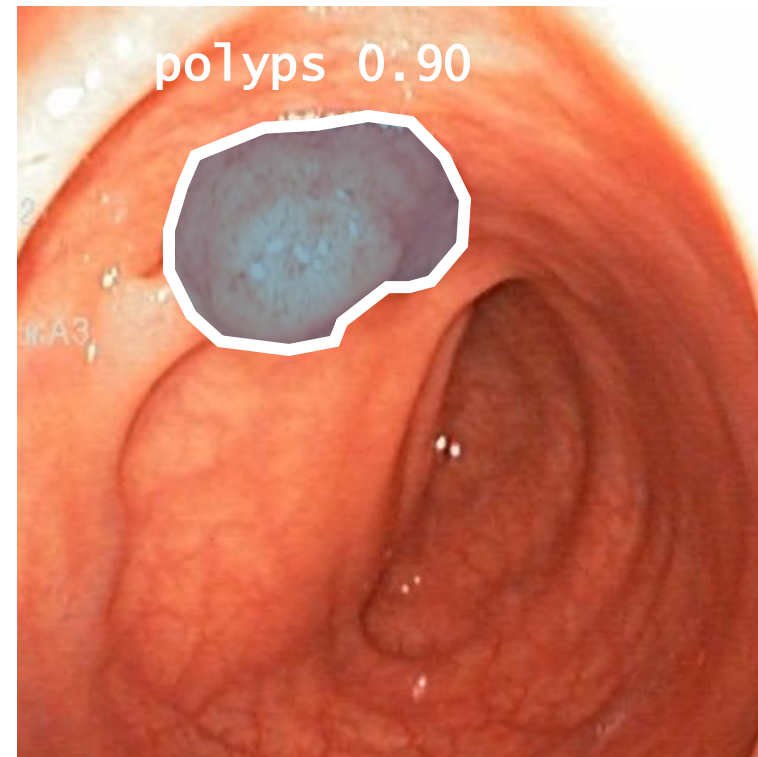
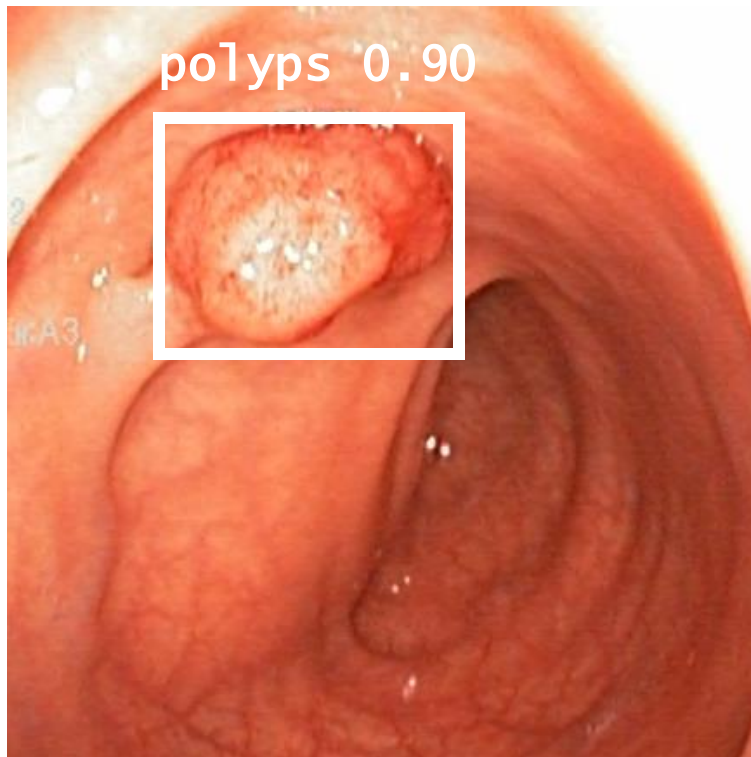
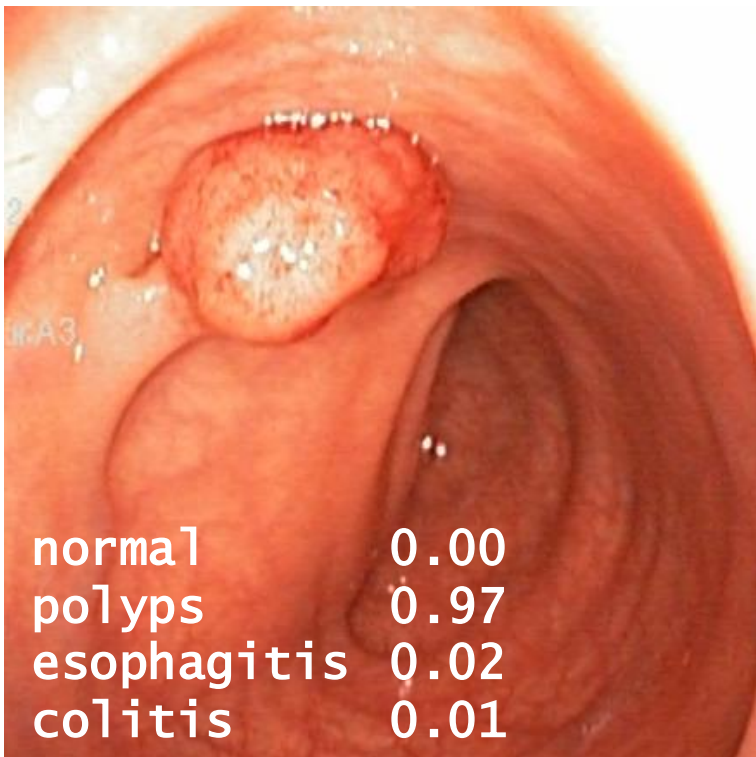
**物体検出・セグメンテーション**

敵対的生成ネットワーク

## 物体分類

## 物体検出

## セグメンテーション



# 物体検出アルゴリズム

Two-stage methods

RCNN

Fast RCNN

SPPNet

**Faster RCNN**

FPN

Pyramid Networks

Mask R-CNN

DETR

Mask2Former

2014

2016

2018

2020

2025

One-stage methods

SSD

RetinaNet

**YOLO**

SquuzeDet

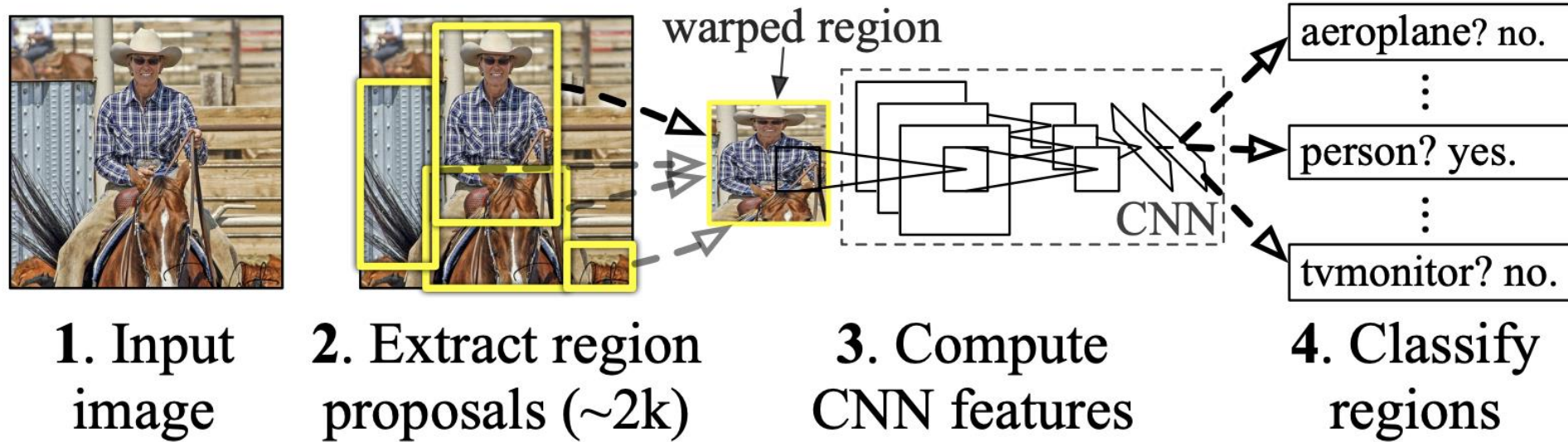
ConerNet

U-Net

YOLACT

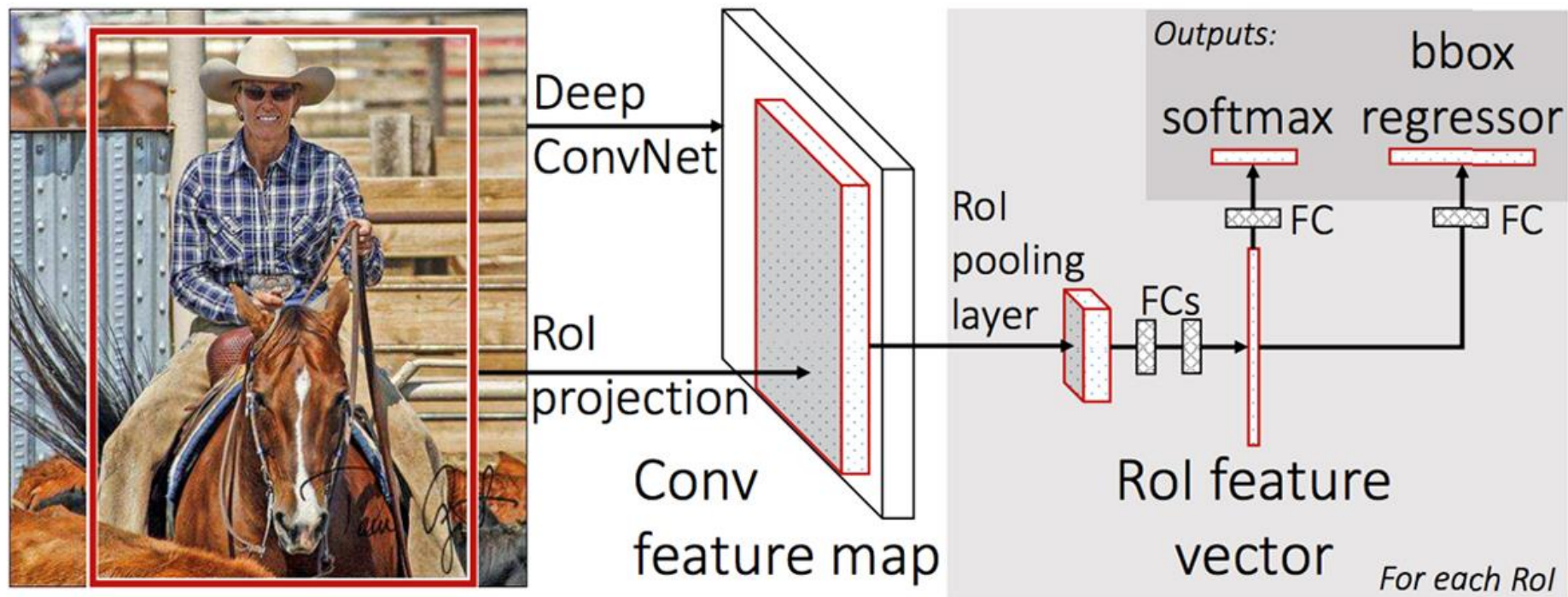


# R-CNN



- Select search アルゴリズムを利用し、入力画像にあるオブジェクトが存在していそうな領域候補 (Region of Interest; RoI) を決定する。
- 各候補領域を CNN に入力し特徴抽出を行う。
- CNN から出力される特徴を SVM に入力し分類を行う。

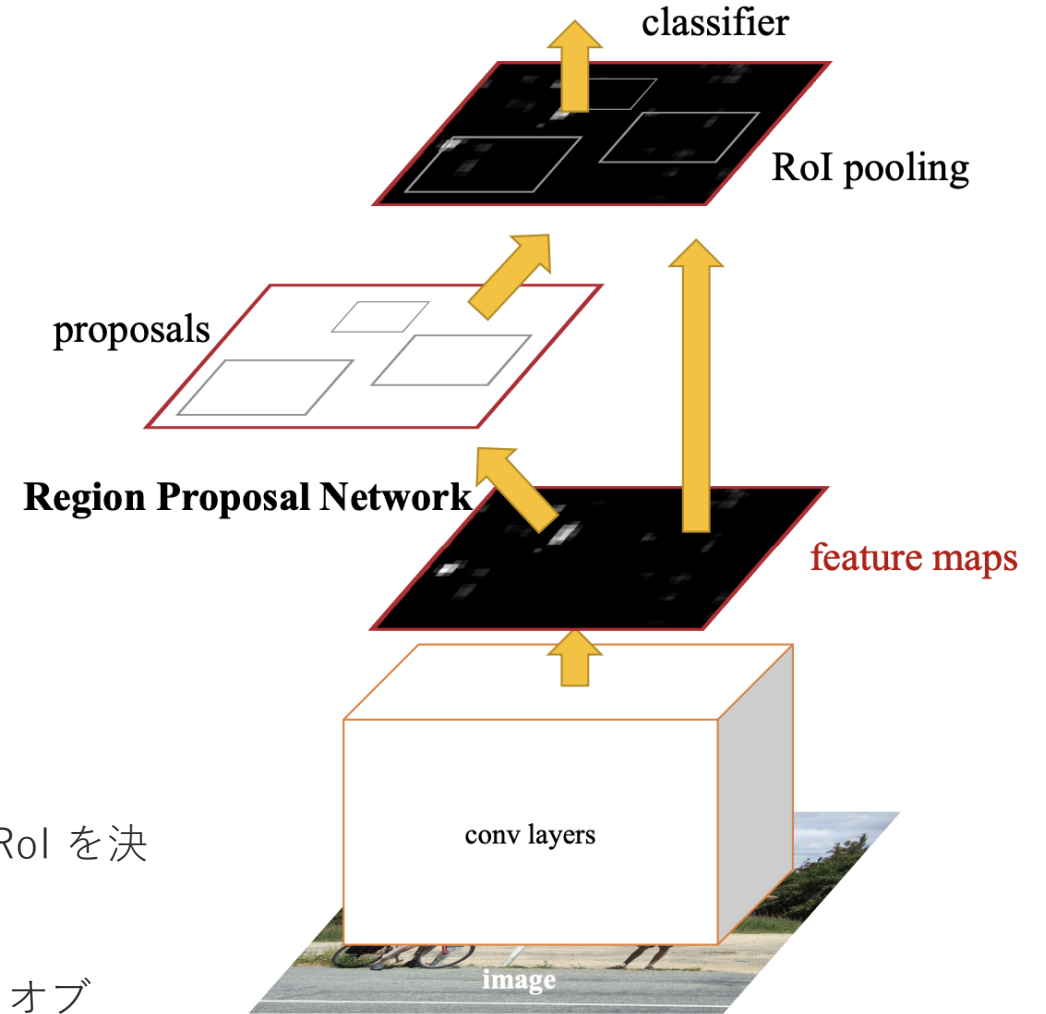
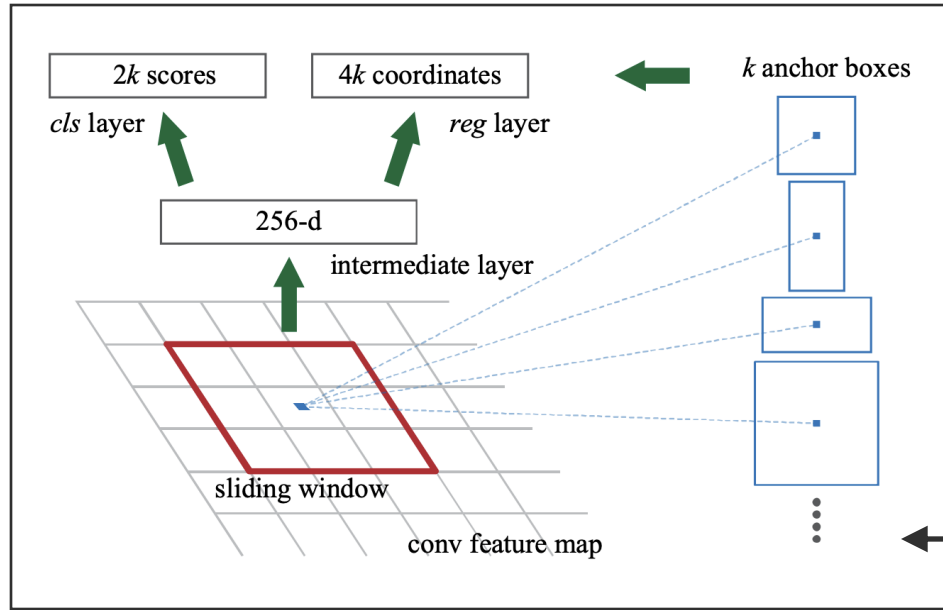
# Fast R-CNN



- 入力画像を CNN に入力し特徴マップを計算する。
- 入力画像をもとに決定した RoI を特徴マップに射影する。
- 特徴マップ上の RoI を処理し、物体分類とバウンディングボックスの座標を予測する。

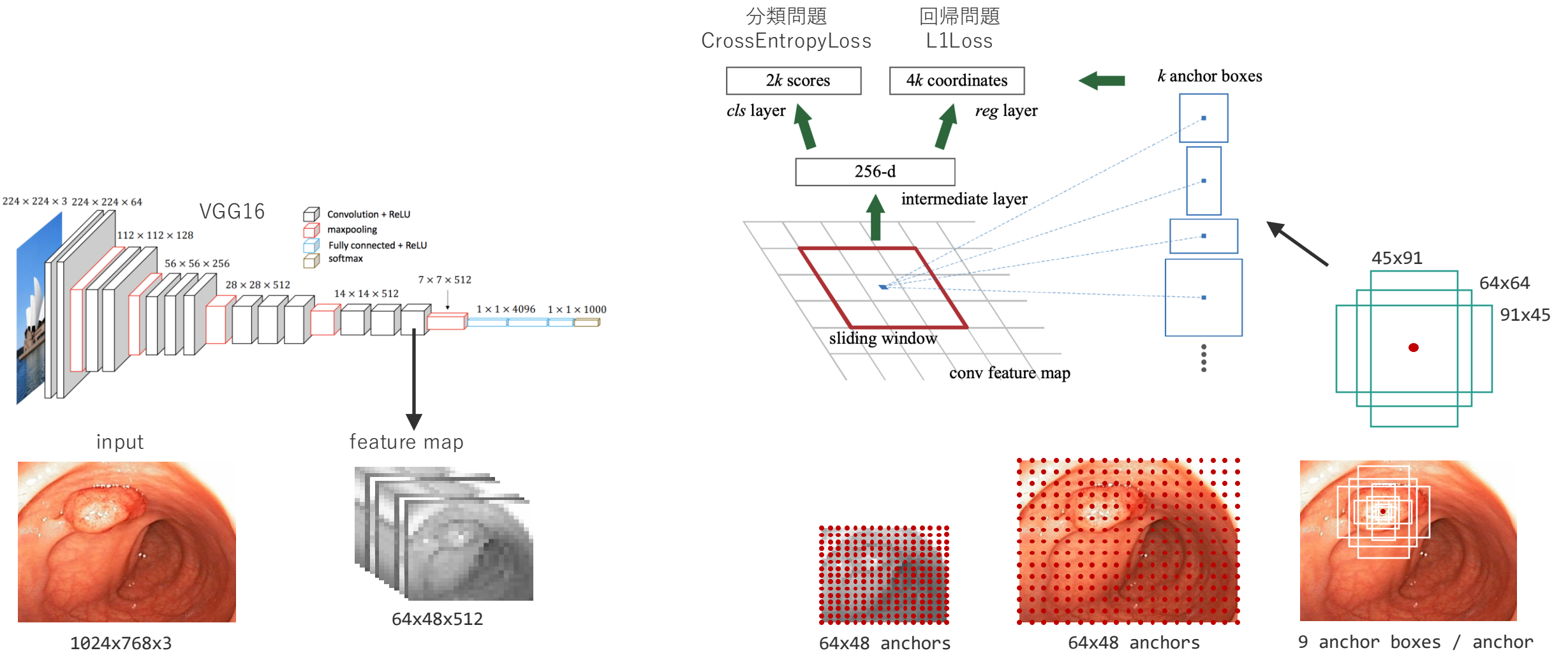


# Faster R-CNN



- 入力画像を CNN に入力し特徴マップを計算する。
- 特徴マップを Region Proposal Network に入力し、RoI を決定する。
- 特徴マップと RPN から出力される RoI を利用して、オブジェクトの分類を行う。

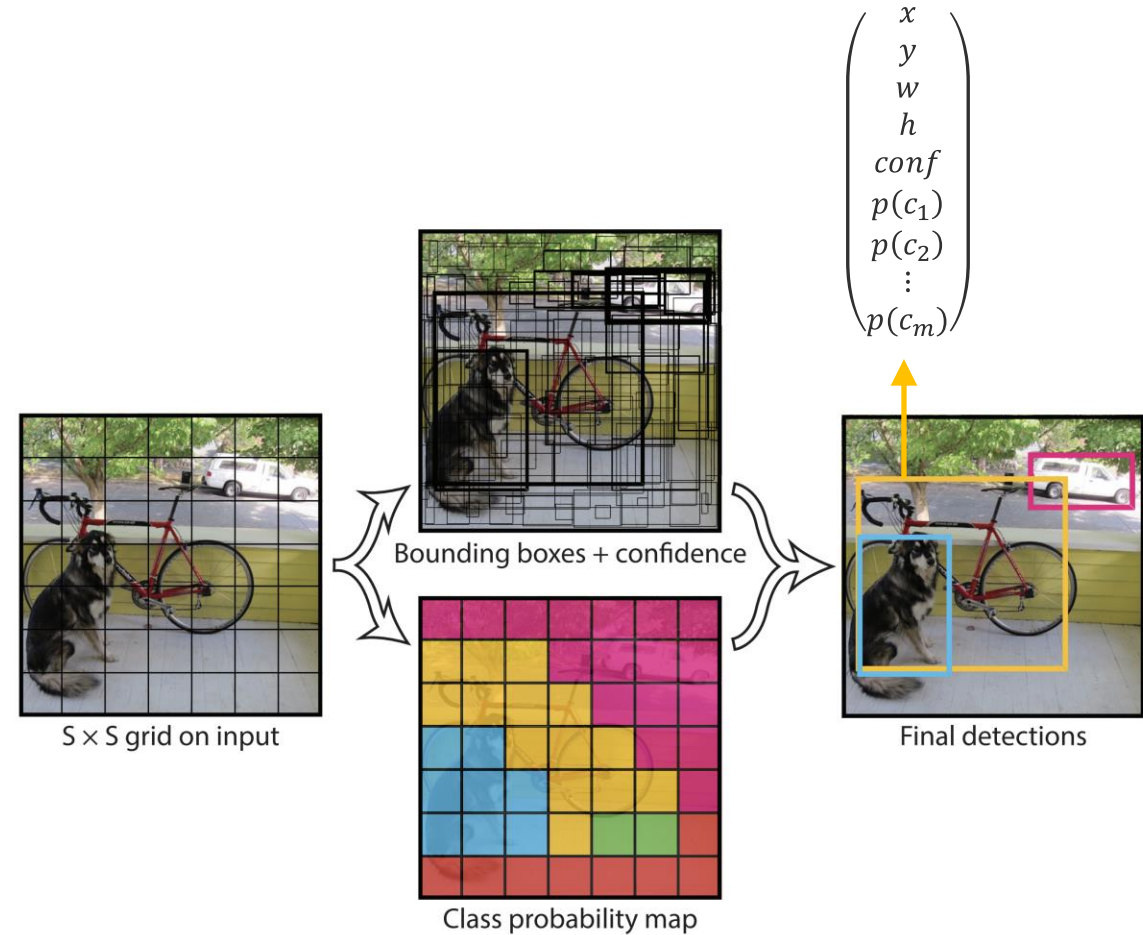
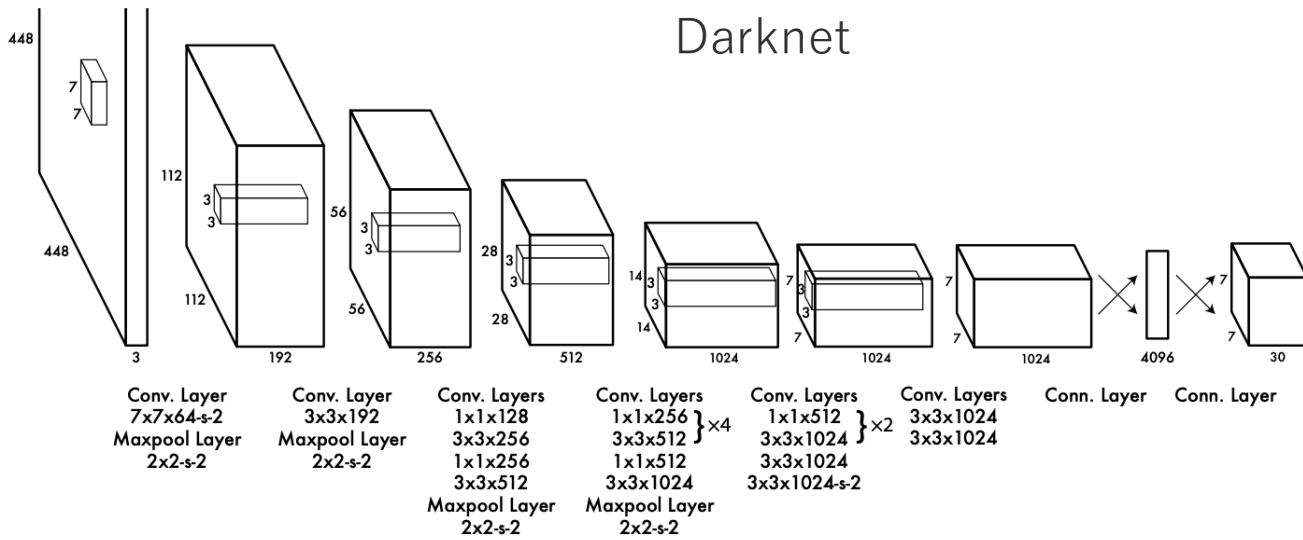
# Resion Proposal Network



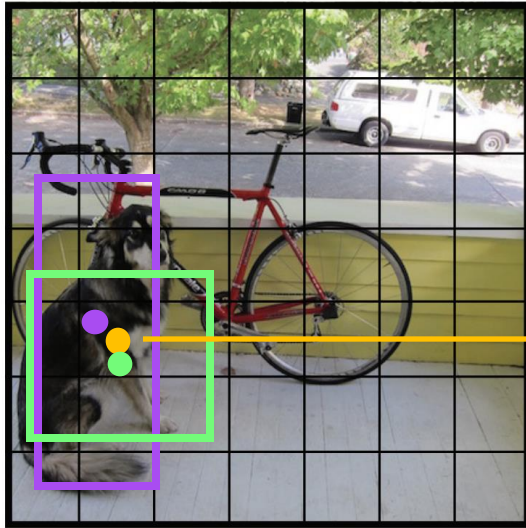
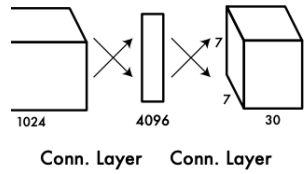
<https://arxiv.org/abs/1506.01497>  
<https://doi.org/10.1145/3038912.3052638>

parameters  
 base width: 64, 128, 256  
 ratio: 1:1, 1:2, 2:1

# You Only Look Once (YOLO)

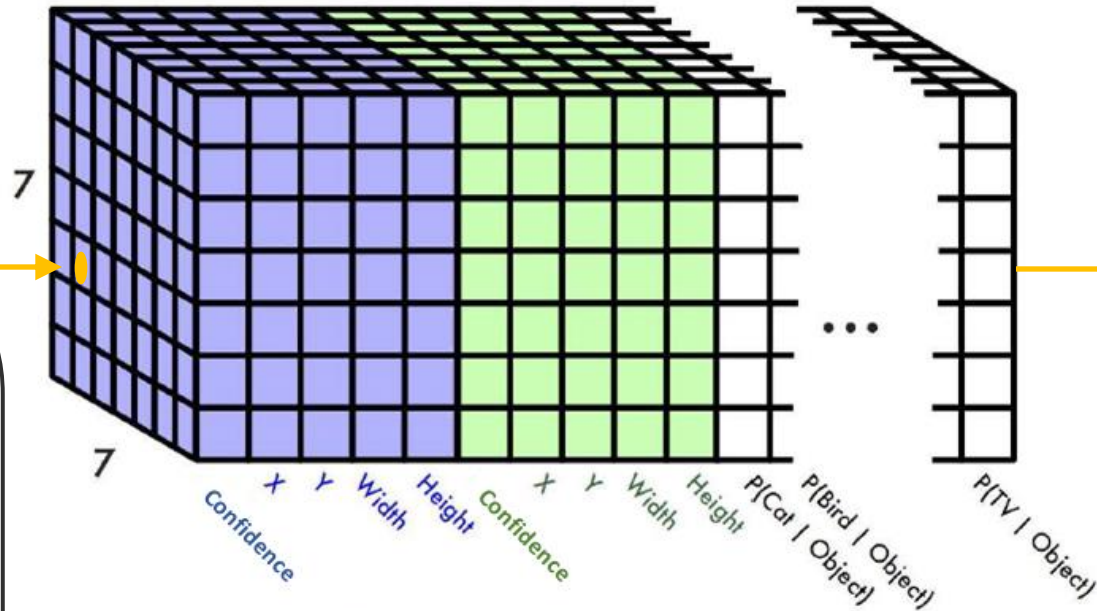


# YOLO – モデル訓練



$S \times S$  grid on input

- $x$
- $y$
- $w$
- $h$
- $conf$
- $x$
- $y$
- $w$
- $h$
- $conf$
- $p(c_1)$
- $p(c_2)$
- $\vdots$
- $p(c_m)$



1st - 5th  
Box #1

6th - 10th  
Box #2

11th - 30th  
Class Probabilities

$$conf = \operatorname{argmax}_{\hat{b} \in \{b_1, b_2\}} \{p(object) \cdot \text{IoU}_{\hat{b}}^{truth}\}$$

- 0.286
  - 0.714
  - 0.332
  - 0.584
  - 0.845
  - 0
  - 1
  - $\vdots$
  - 0
- $p(\text{dog})$
- 教師ラベル

# YOLO – 損失関数

訓練過程で、オブジェクトのあるセルではバウンディングボックスの中心位置および幅と高さが学習され、信頼度も増加する。背景のセルでは学習されずに出鱈目なバウンディングボックスを低信頼度で生成される。

1: オブジェクト  
0: 背景

bbox 中心座標の予測誤差

$$\lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right]$$

bbox の幅と高さの予測誤差

$$+ \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ \left( \sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left( \sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right]$$

confidence の予測誤差

$$+ \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left( C_i - \hat{C}_i \right)^2$$

confidence の予測誤差 (背景セルならば  $C_i = 0$ )

$$+ \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} \left( C_i - \hat{C}_i \right)^2$$

クラスの予測誤差

$$+ \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} \left( p_i(c) - \hat{p}_i(c) \right)^2$$

予測誤差が非常に小さいので、定数倍  $\lambda_{\text{coord}} = 5$  をかけて大きくする。

自信を持って生成した bbox が偶然に教師ラベルと重なった場合、損失がゼロになる。

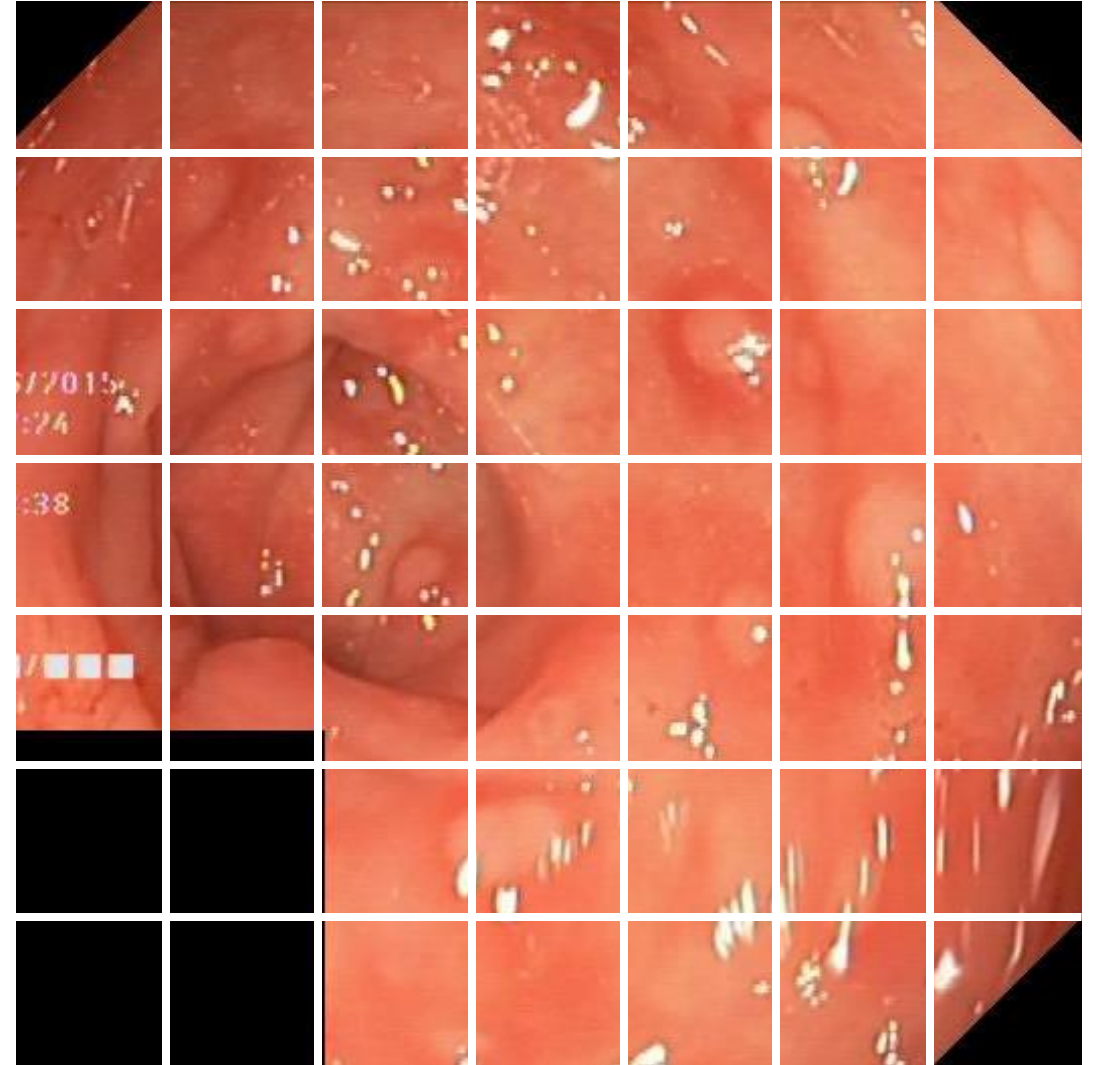
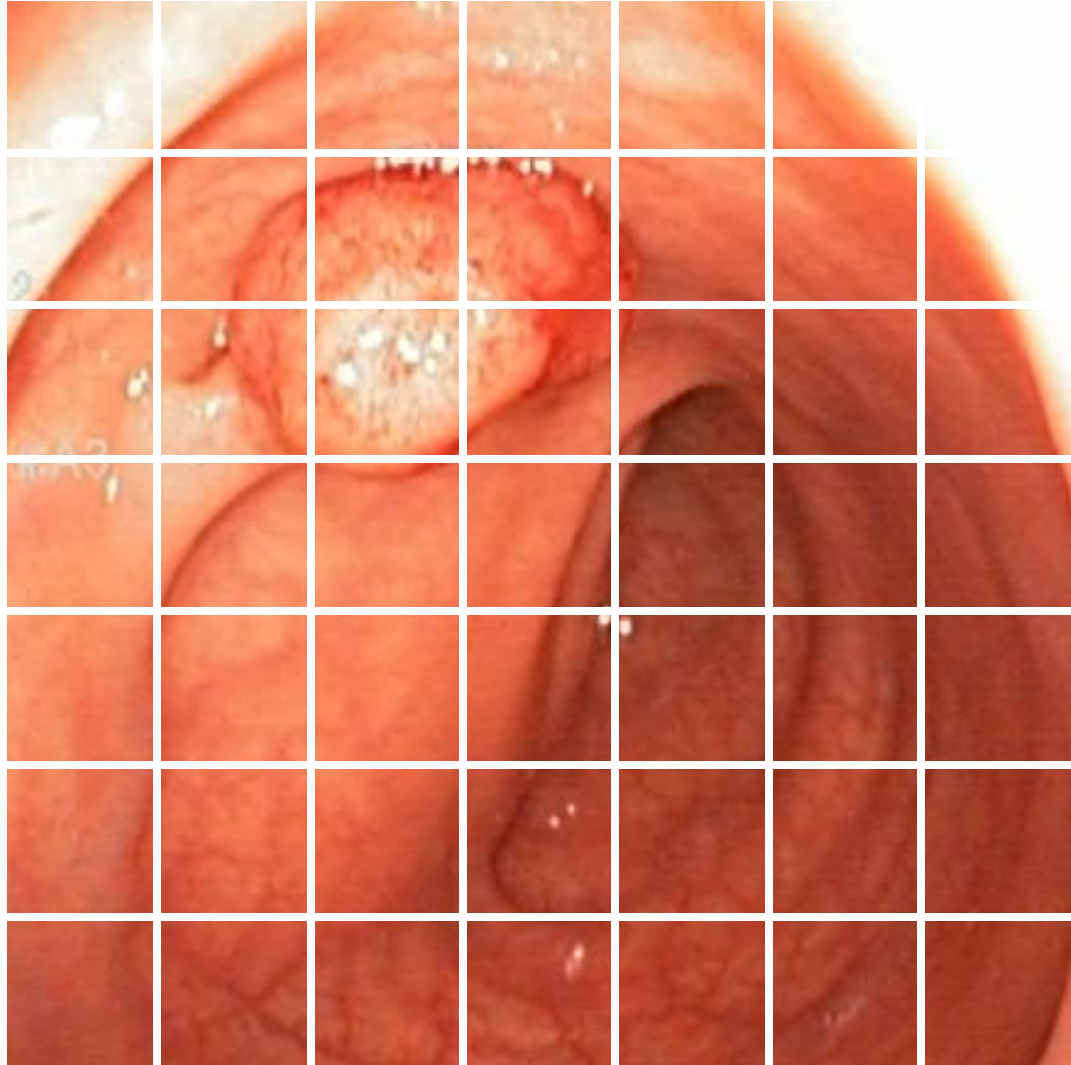
自信を持って生成した bbox にオブジェクトが含まれていなければ、その分だけ損失を計上する。

画像の大部分が背景であるため、この項が大きな値を取りやすく、他の項の値が相対的に小さくなり、学習が進まない。この効果を弱めるために、定数倍  $\lambda_{\text{noobj}} = 0.5$  をかける。

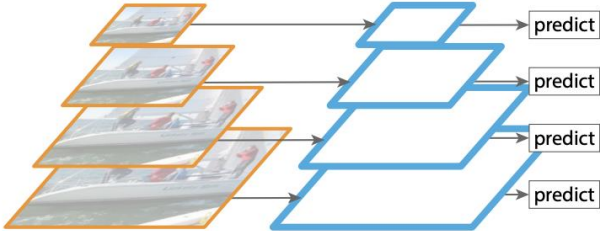


# YOLO

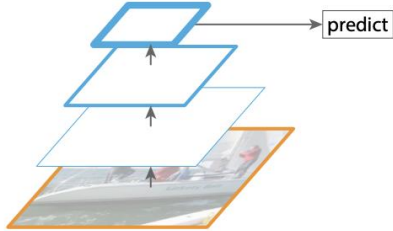
---



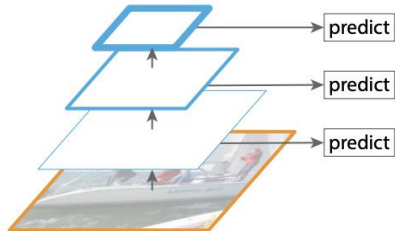
# Feature Pyramid Network (FPN)



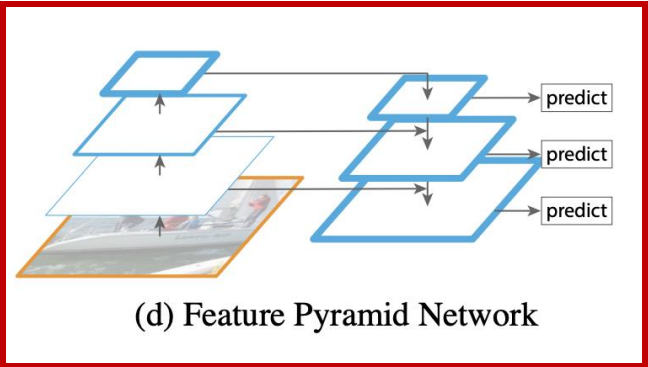
(a) Featurized image pyramid



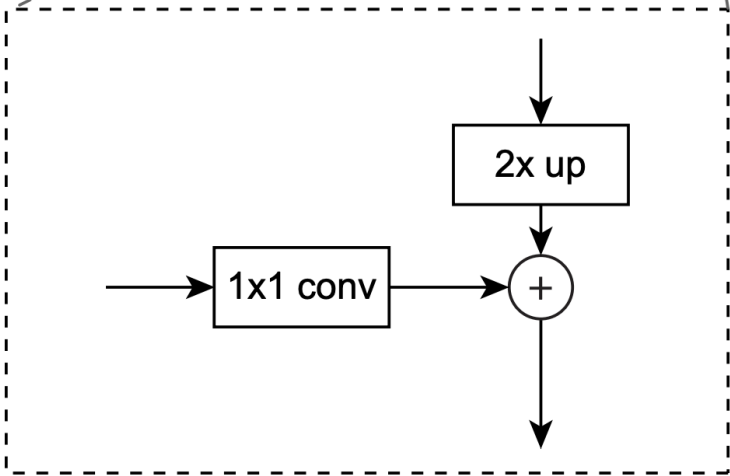
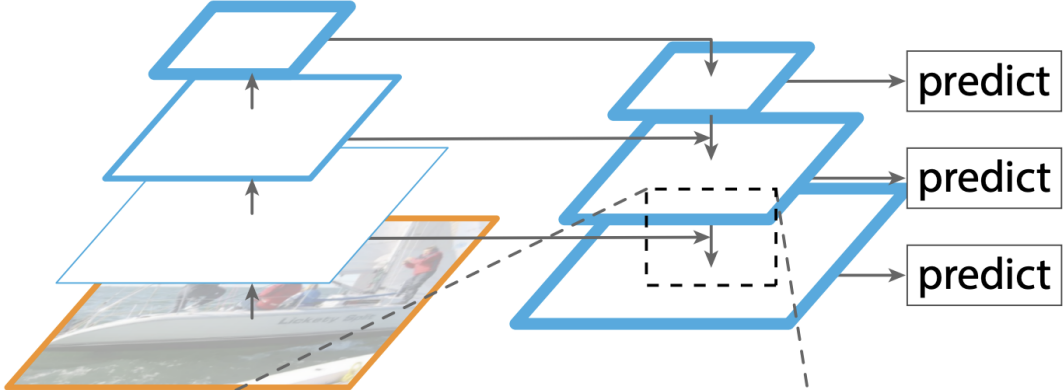
(b) Single feature map



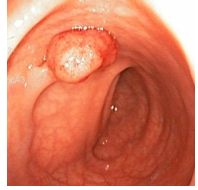
(c) Pyramidal feature hierarchy



(d) Feature Pyramid Network

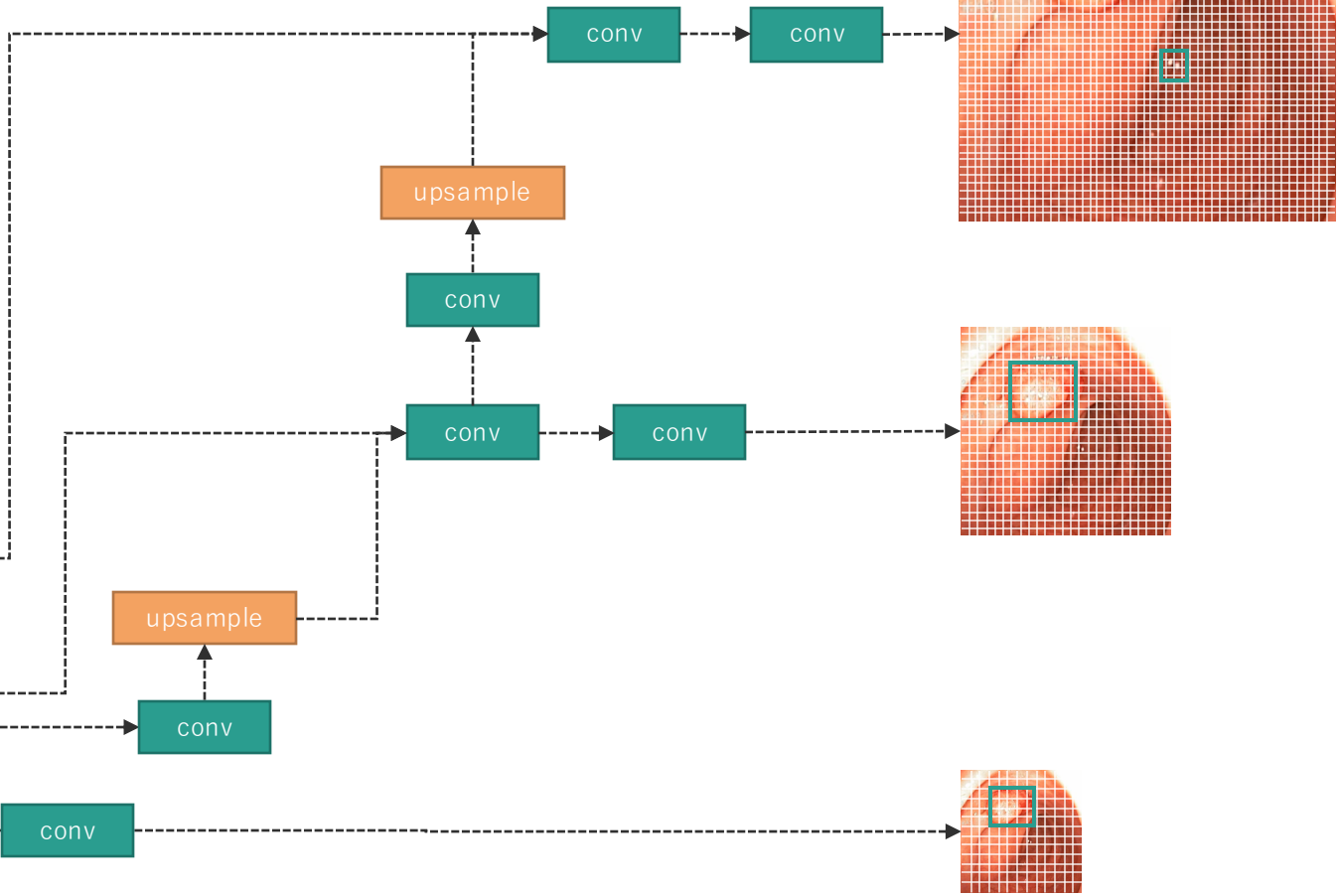


# YOLO3



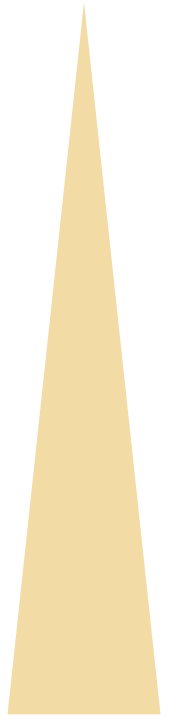
Darknet-53

Type	Filters	Size	Output
Convolutional	32	3 × 3	256 × 256
Convolutional	64	3 × 3 / 2	128 × 128
Convolutional	32	1 × 1	128 × 128
Convolutional	64	3 × 3	
Residual			128 × 128
Convolutional	128	3 × 3 / 2	64 × 64
Convolutional	64	1 × 1	64 × 64
Convolutional	128	3 × 3	
Residual			64 × 64
Convolutional	256	3 × 3 / 2	32 × 32
Convolutional	128	1 × 1	32 × 32
Convolutional	256	3 × 3	
Residual			32 × 32
Convolutional	512	3 × 3 / 2	16 × 16
Convolutional	256	1 × 1	16 × 16
Convolutional	512	3 × 3	
Residual			16 × 16
Convolutional	1024	3 × 3 / 2	8 × 8
Convolutional	512	1 × 1	8 × 8
Convolutional	1024	3 × 3	
Residual			8 × 8
Avgpool		Global	
Connected		1000	
Softmax			



small objects

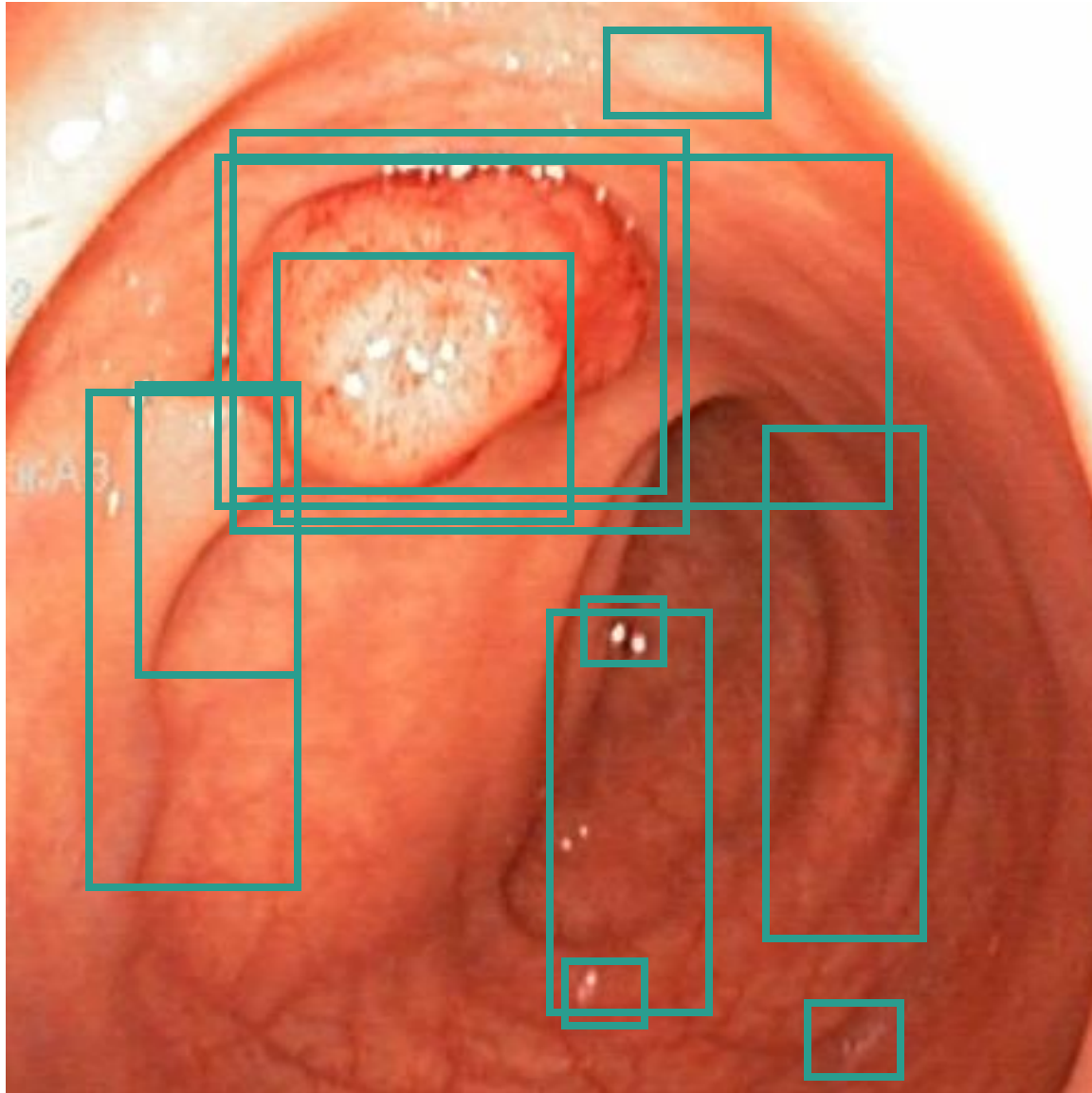
large objects





# Non-maximum Suppression (NMS)

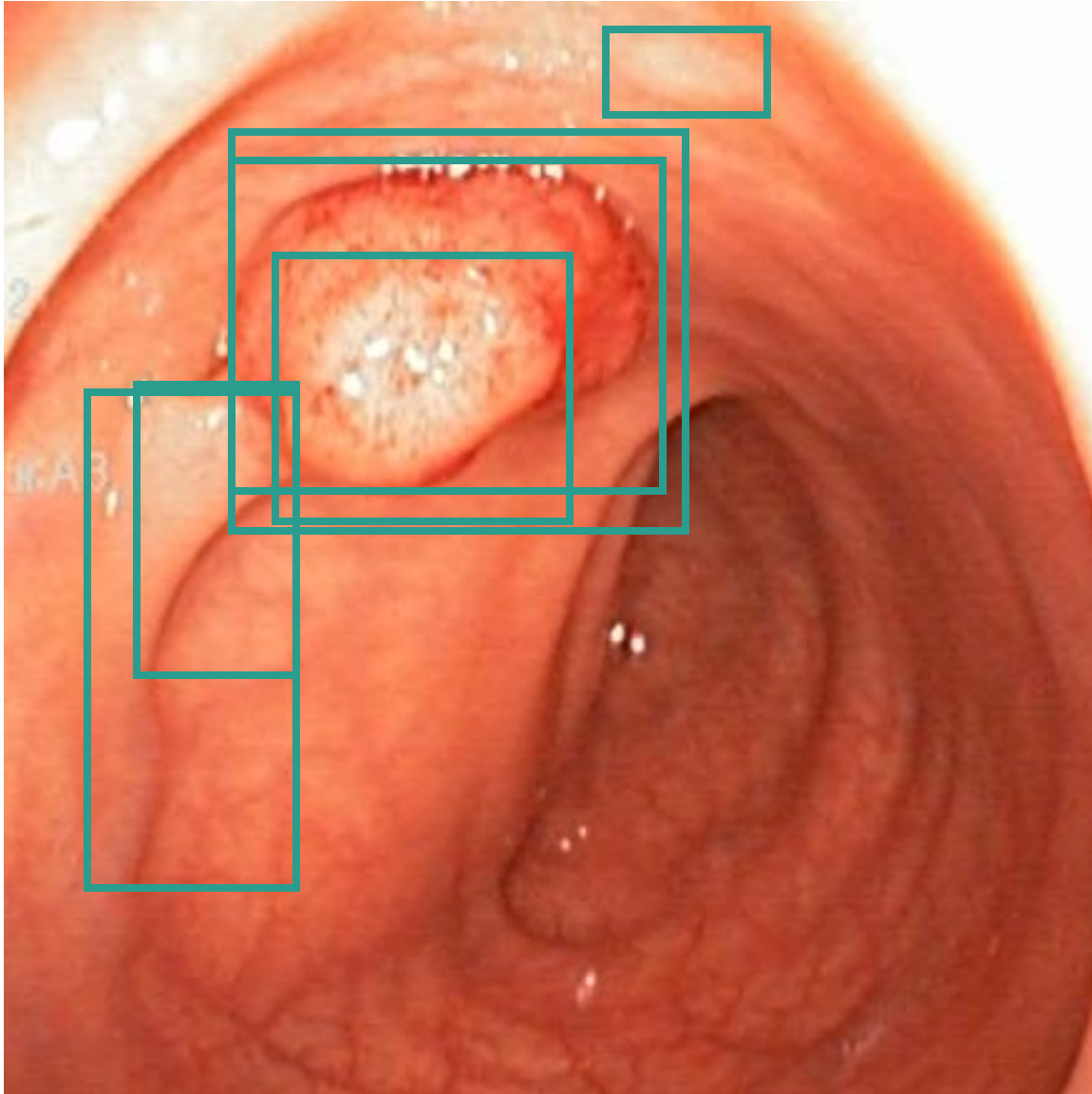
---



1. 予測において数多くの bbox が出力される。

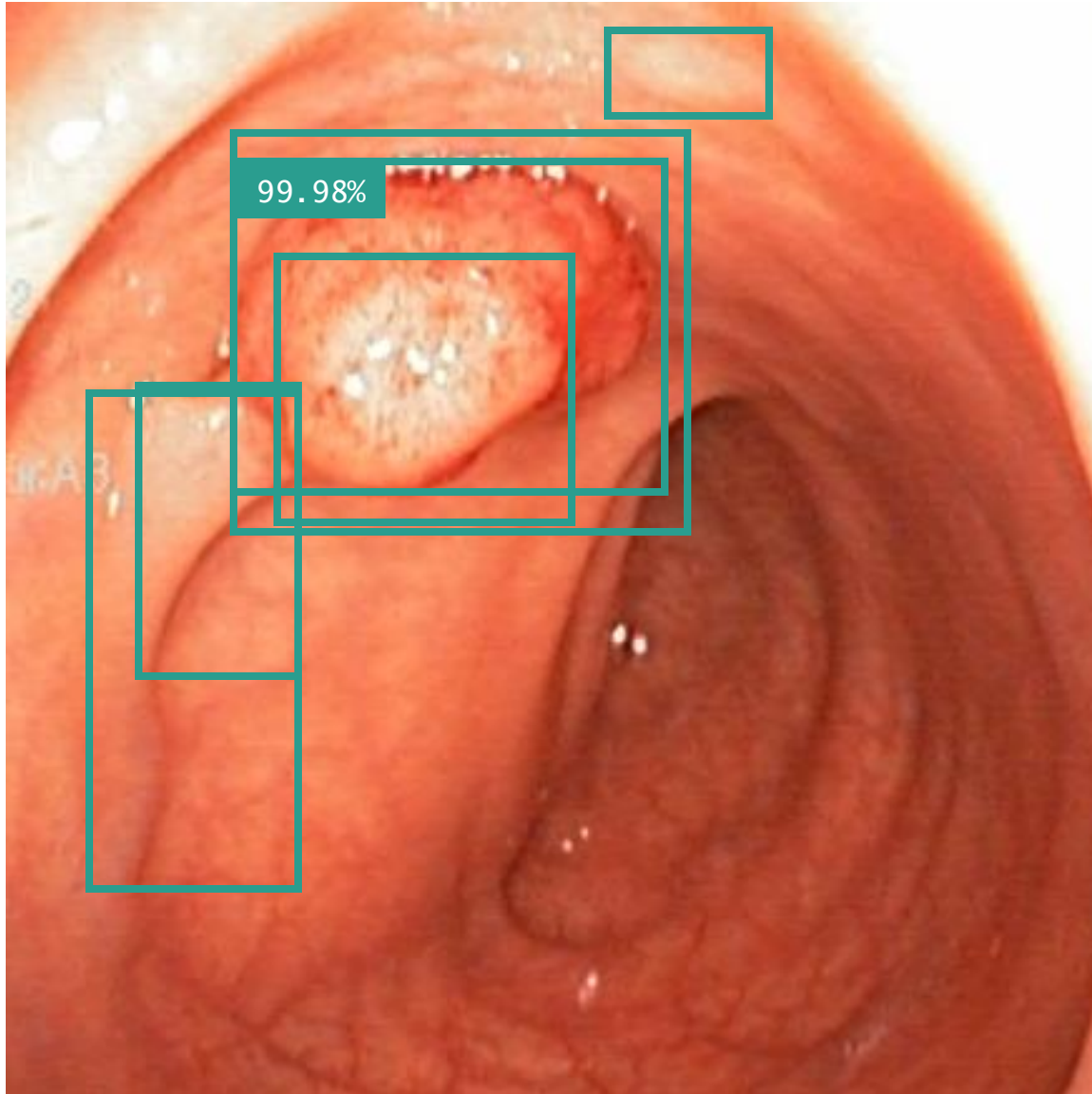
# Non-maximum Suppression (NMS)

---



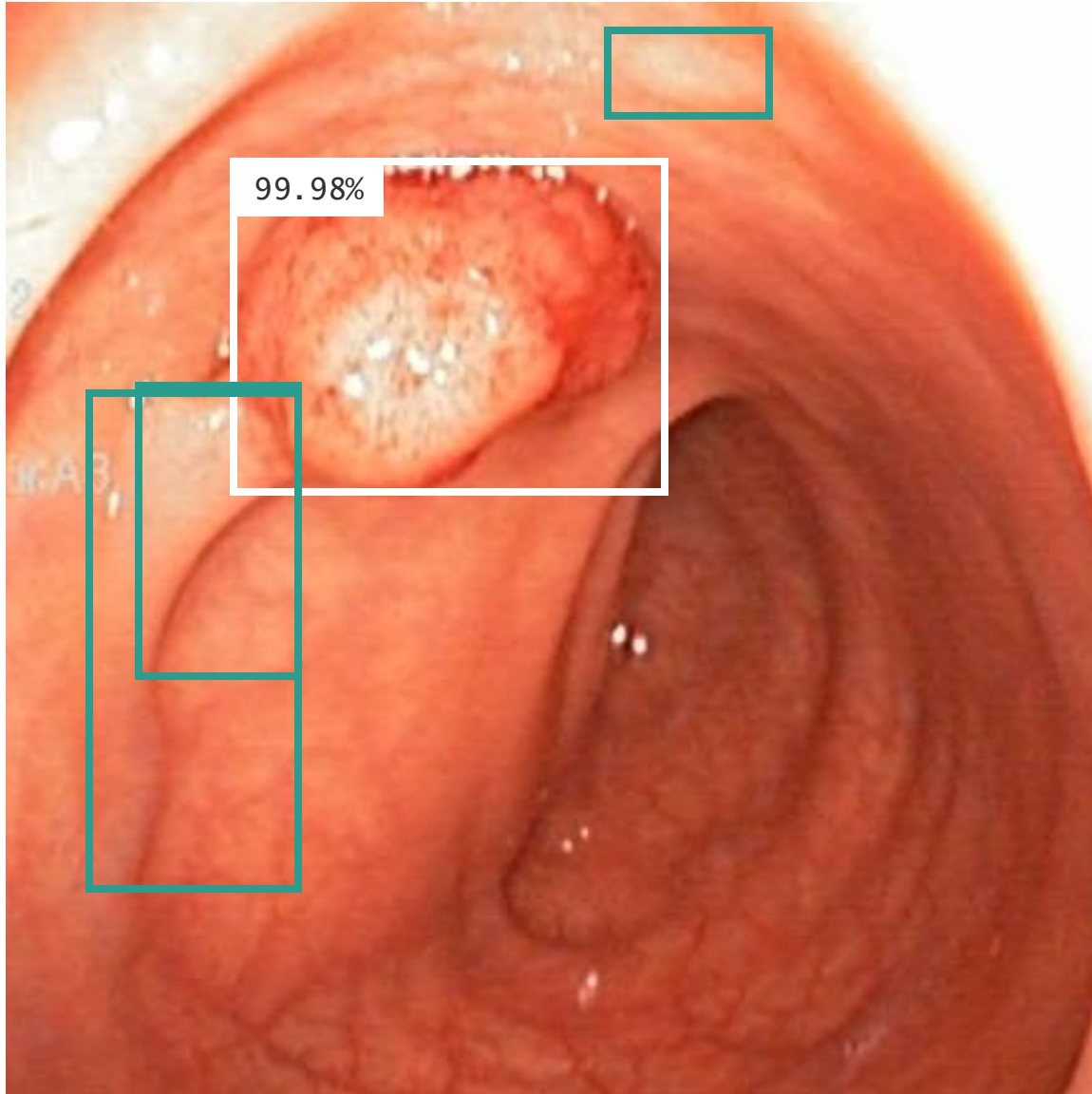
1. 予測において数多くの bbox が出力される。
2. 信頼度が最も高い k 個の bbox を残す。

# Non-maximum Suppression (NMS)



1. 予測において数多くの bbox が出力される。
2. 信頼度が最も高い  $k$  個の bbox を残す。
3. 信頼度が最も高い bbox を基準として、他の bbox との IoU を計算する。

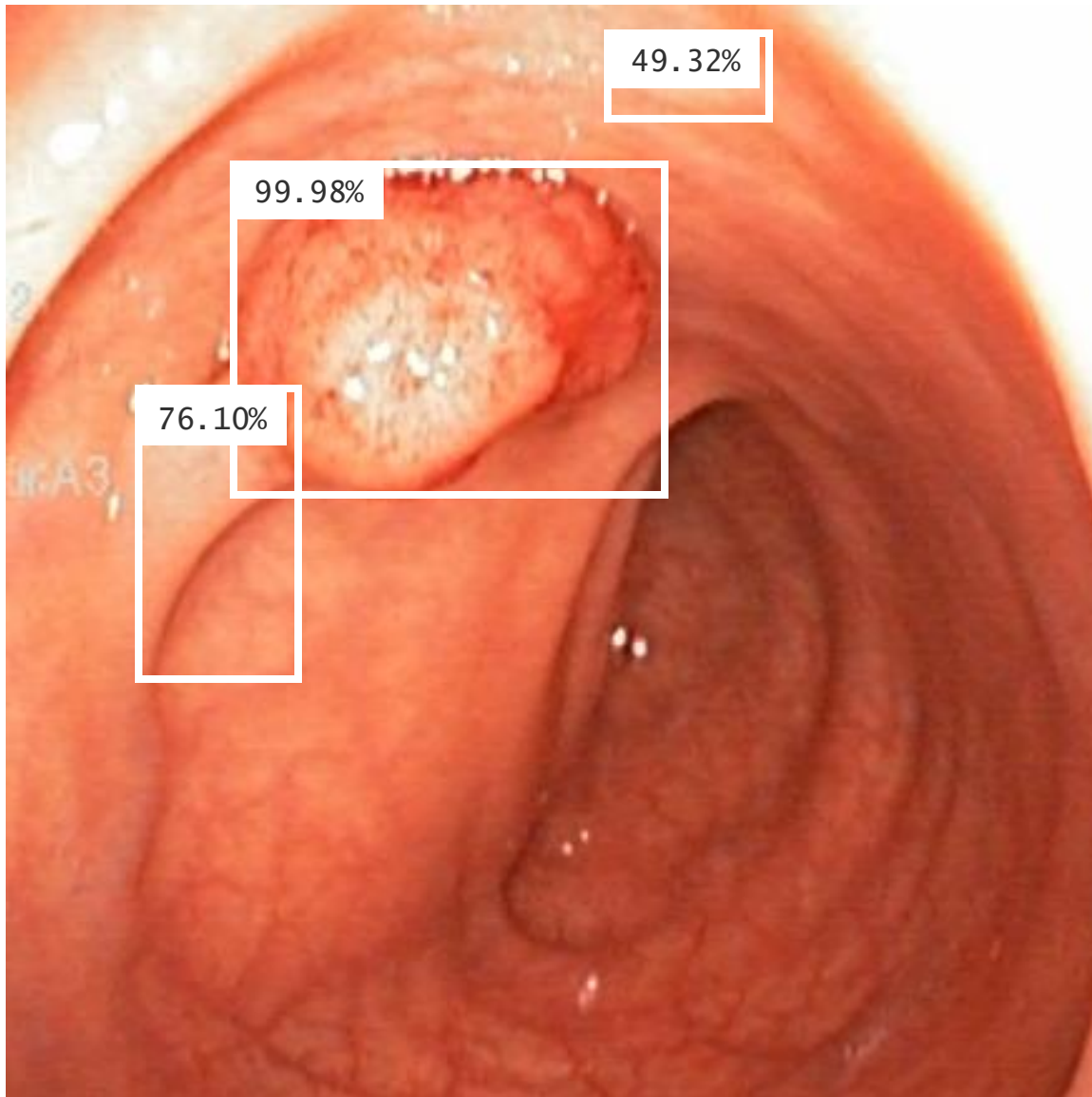
# Non-maximum Suppression (NMS)



1. 予測において数多くの bbox が出力される。
2. 信頼度が最も高い  $k$  個の bbox を残す。
3. 信頼度が最も高い bbox を基準として、他の bbox との IoU を計算する。
4.  $\text{IoU} > 0.5$  ならばその bbox を除去し、基準となった bbox を予測結果として出力する。

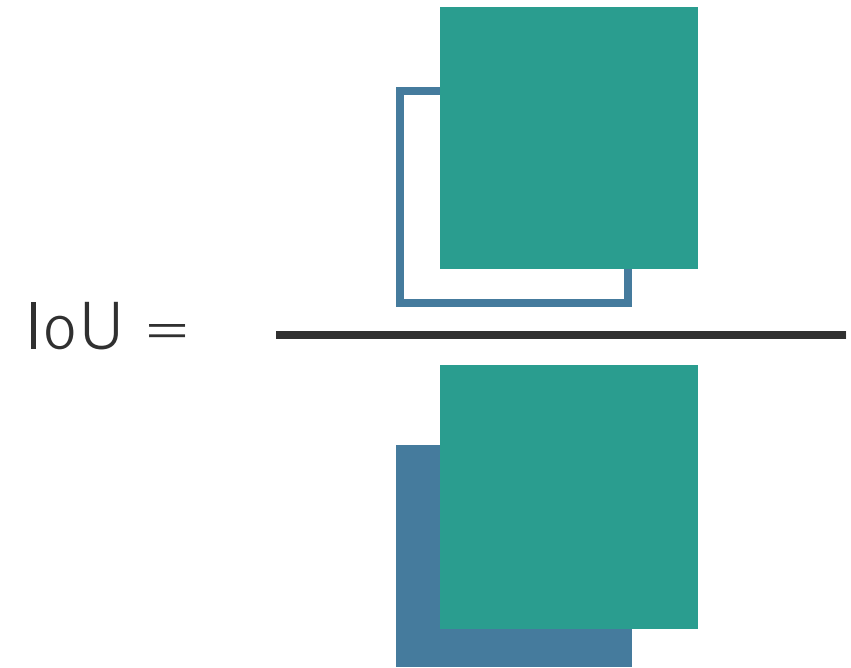
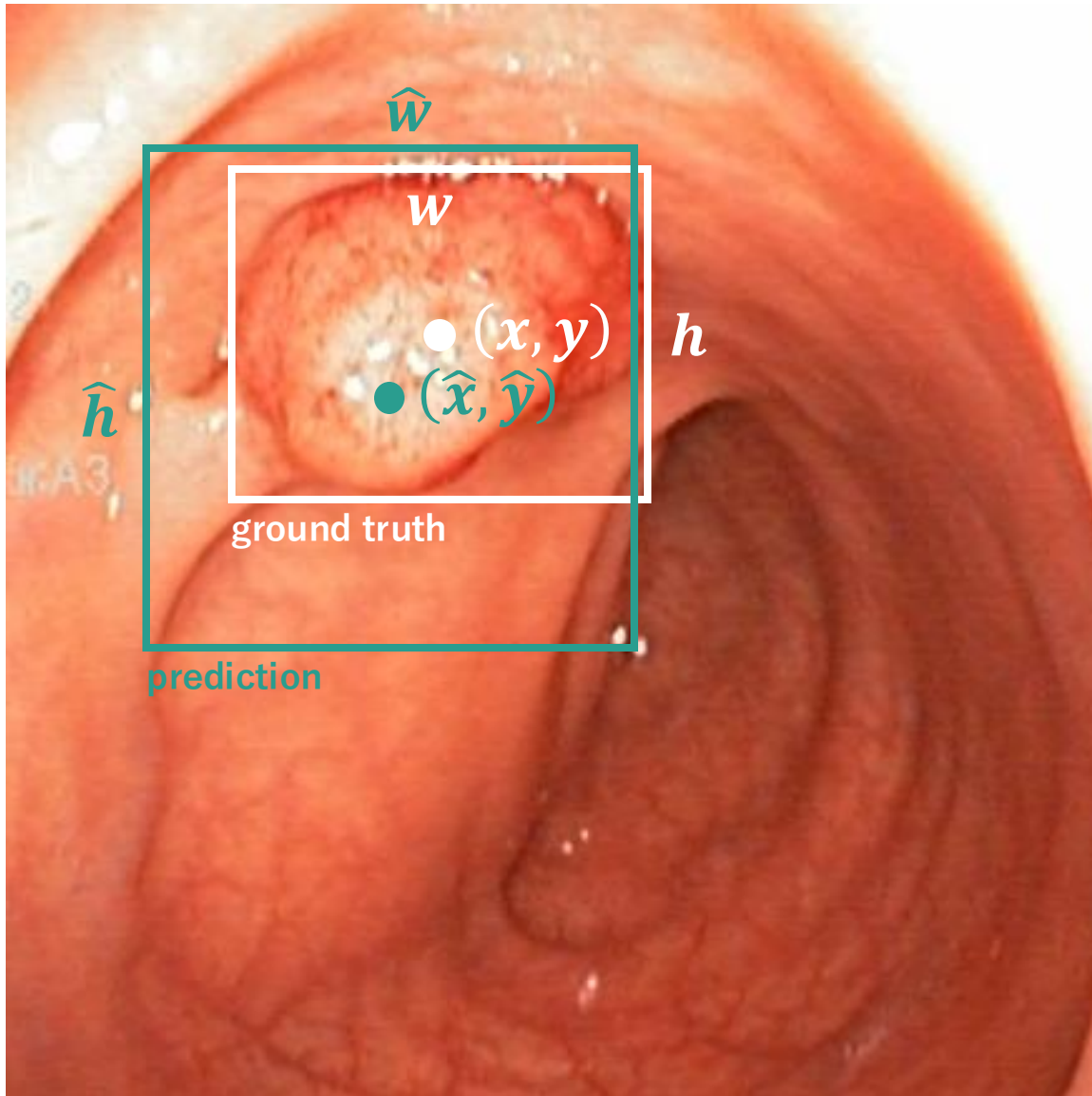


# Non-maximum Suppression (NMS)



1. 予測において数多くの bbox が出力される。
2. 信頼度が最も高い  $k$  個の bbox を残す。
3. 信頼度が最も高い bbox を基準として、他の bbox との IoU を計算する。
4.  $\text{IoU} > 0.5$  ならばその bbox を除去し、基準となった bbox を予測結果として出力する。
5. 残った bbox に対して、手順 3-4 を繰り返す。

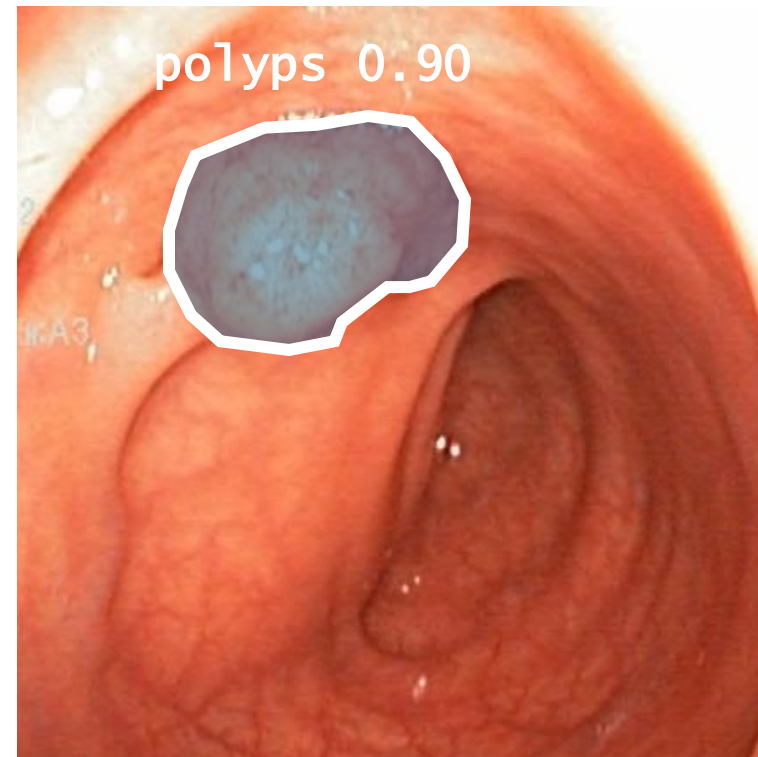
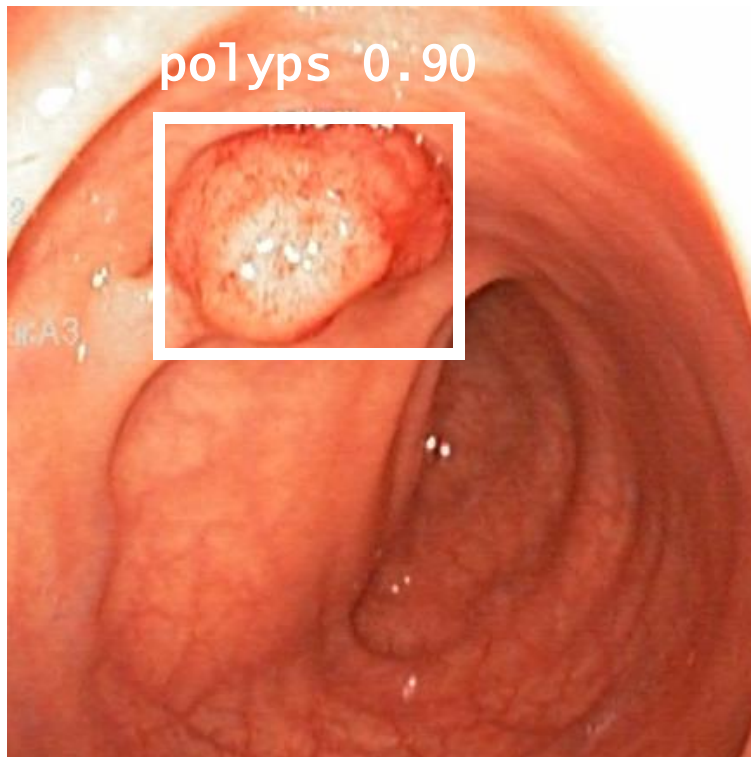
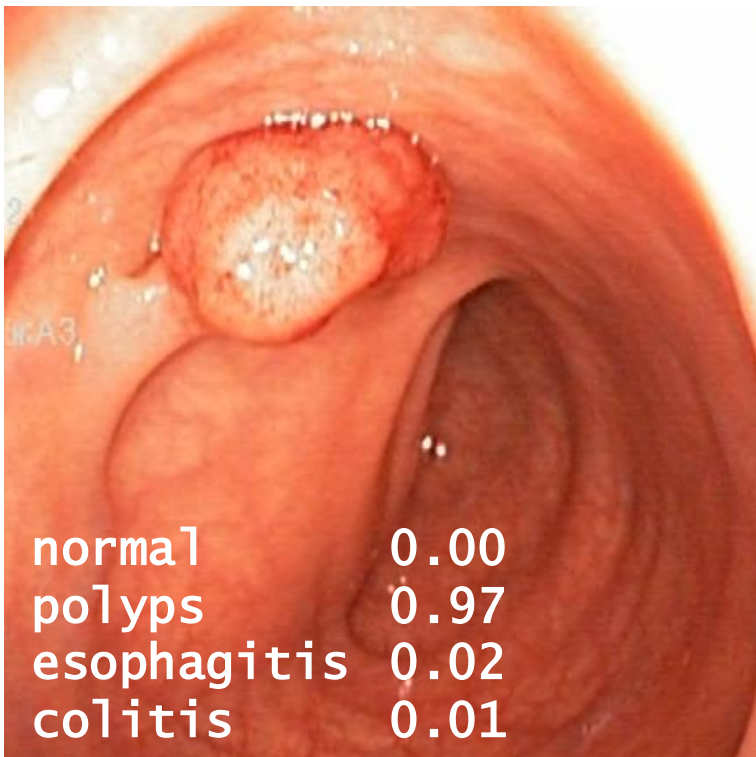
# Intersection over Union (IoU)



## 物体分類

## 物体検出

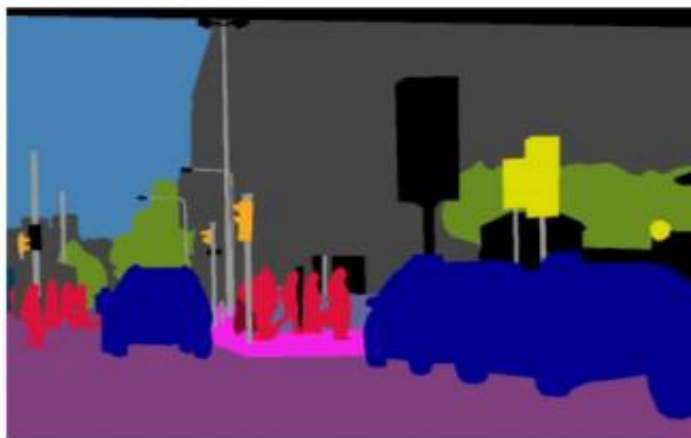
## セグメンテーション



# セグメンテーション



(a) image



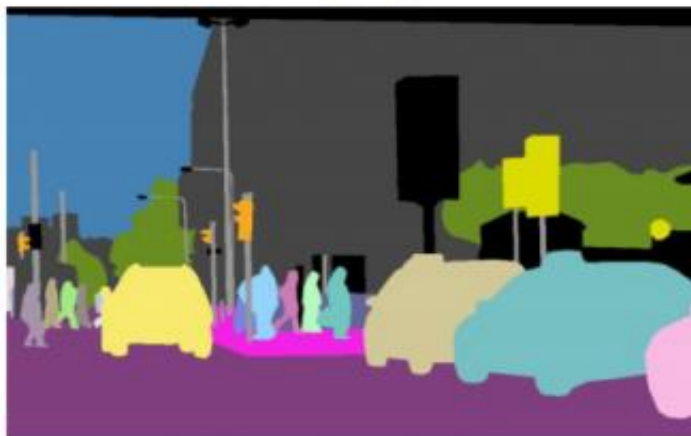
(b) semantic segmentation

● すべてのピクセルをクラスレベルで検出。



(c) instance segmentation

● 個々のオブジェクトをピクセルレベルで検出し、オブジェクトごとに ID を付与。

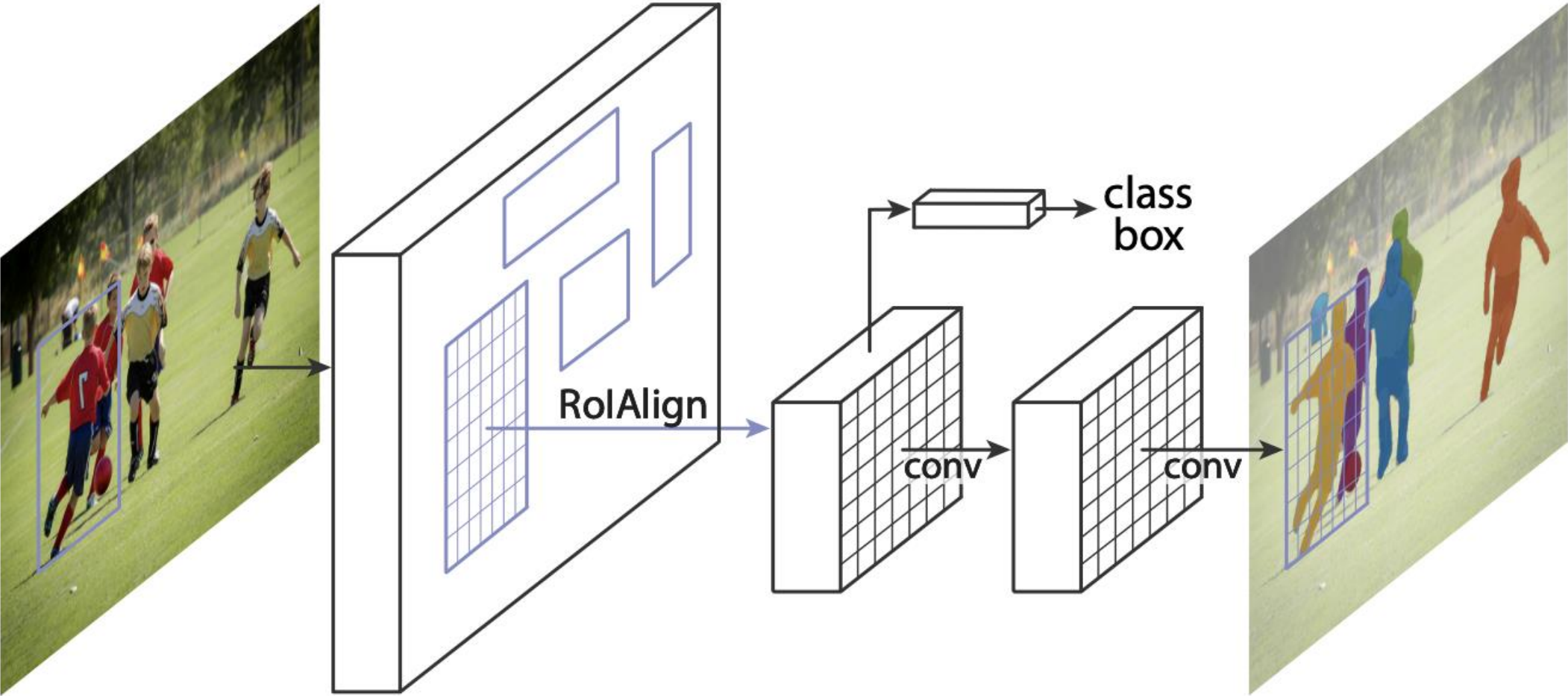


(d) panoptic segmentation

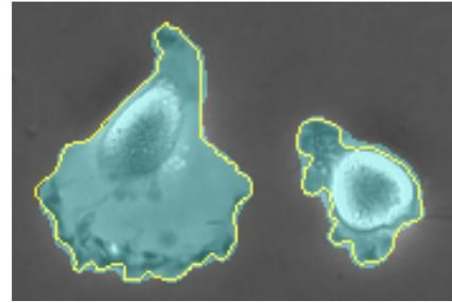
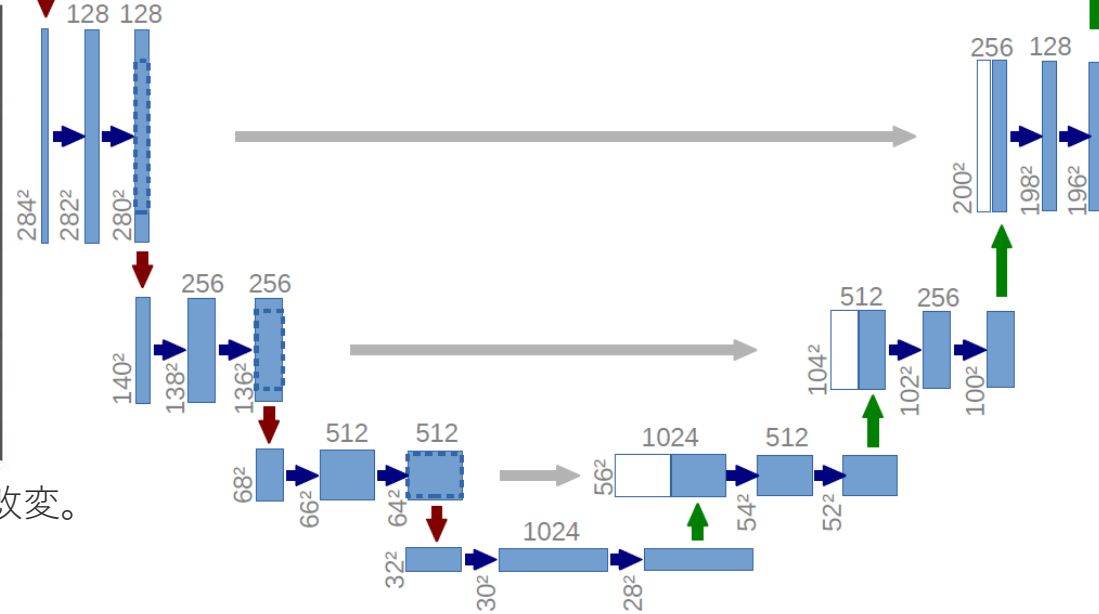
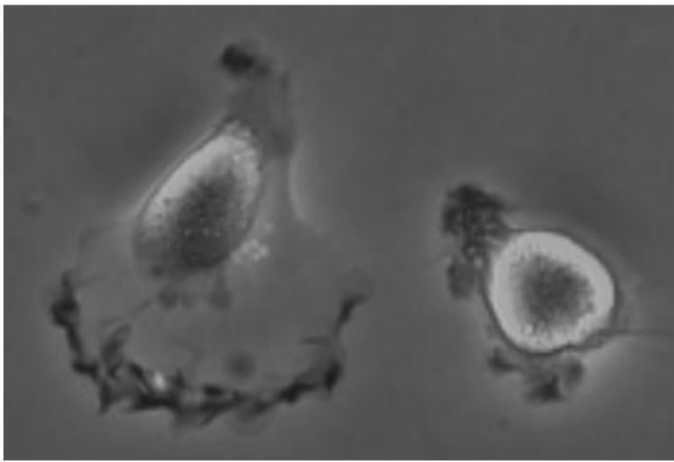
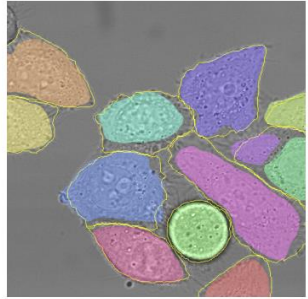
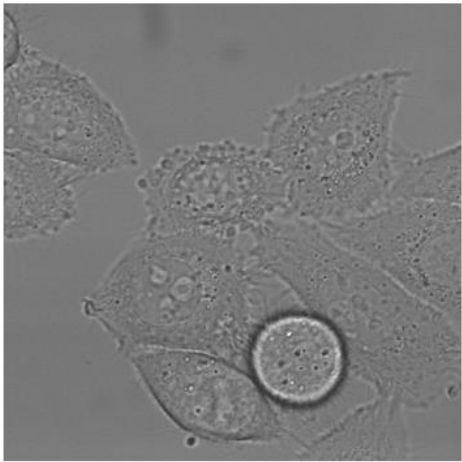
● すべてのピクセルをクラスレベルで検出し、オブジェクトごとに ID を付与。



# Mask R-CNN

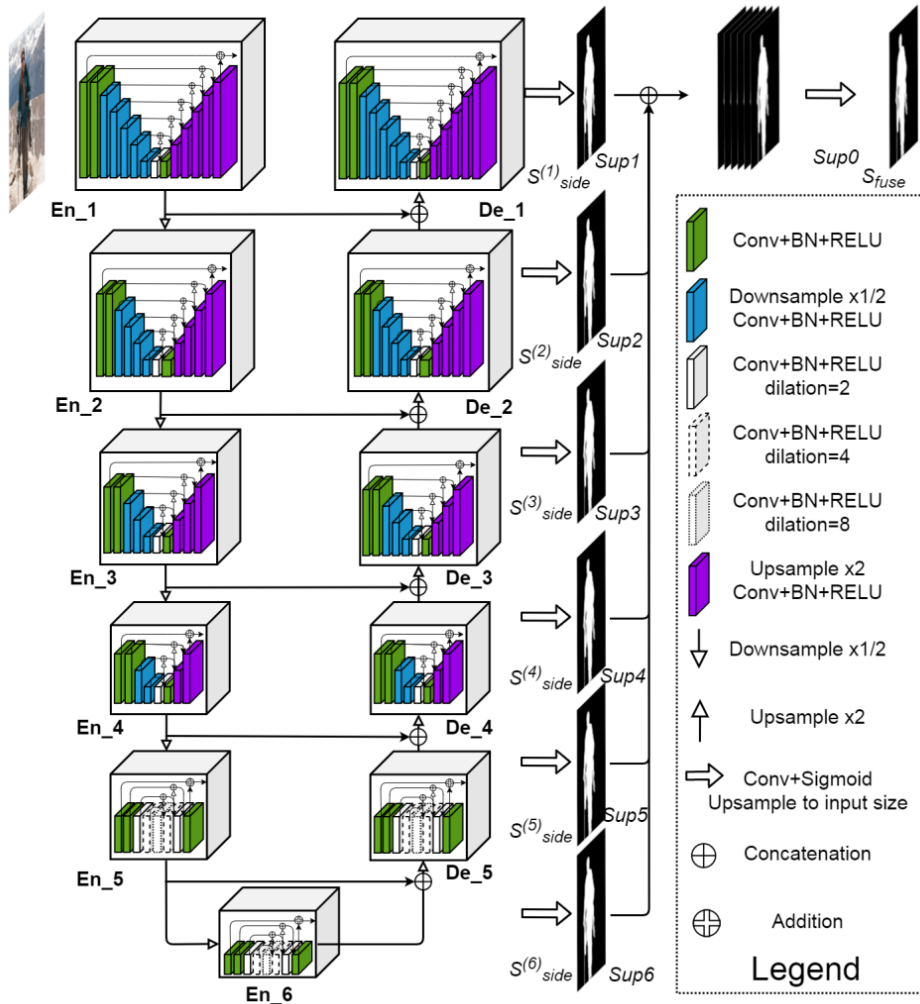


# U-Net

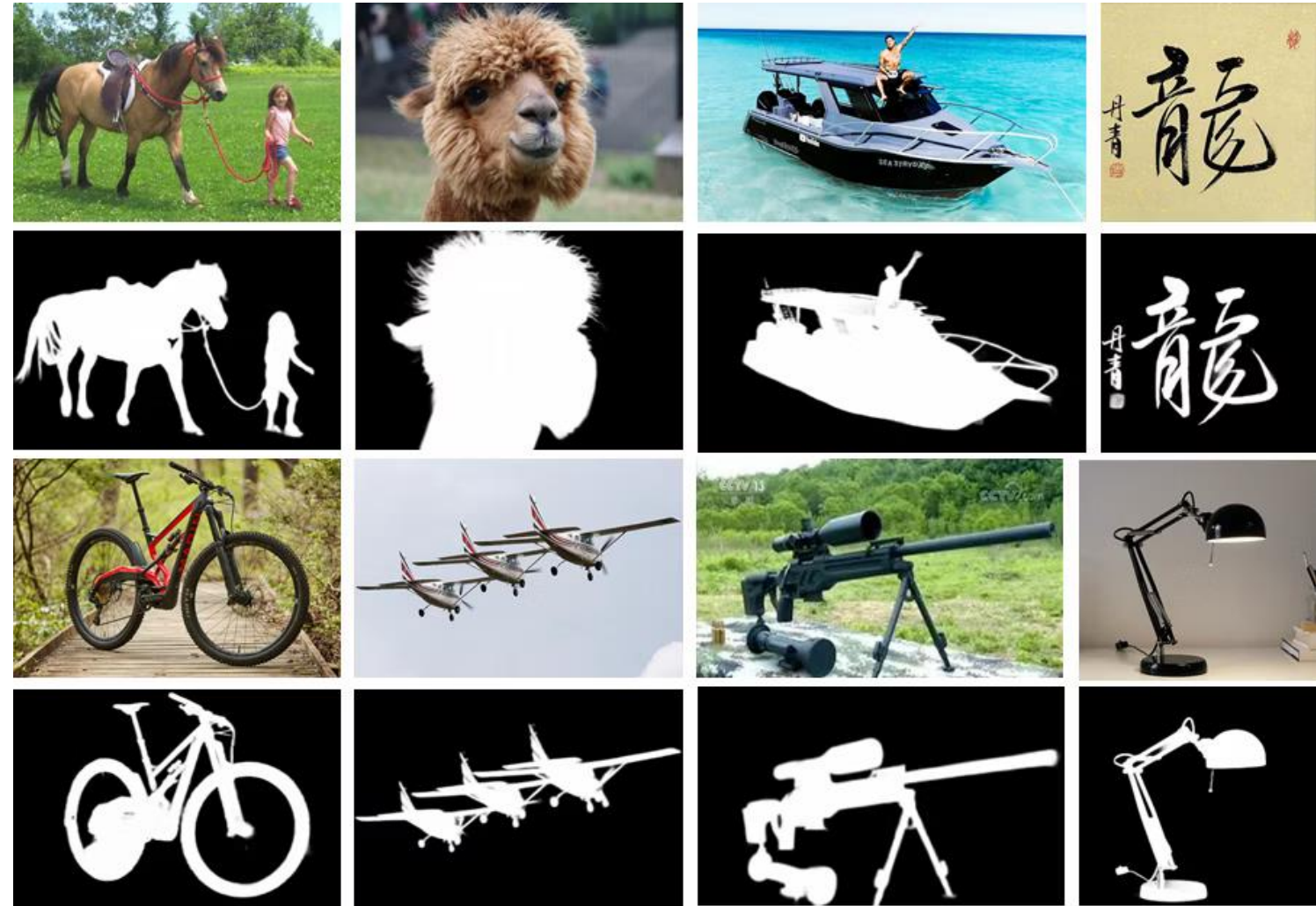


Ronneberger et al., 2015. Fig 1, 4 より改変。

# Salient Object Detection



Qin et al., 2020. Fig 5



<https://github.com/xuebinqin/U-2-Net>





# 深層学習

---

ニューラルネットワーク

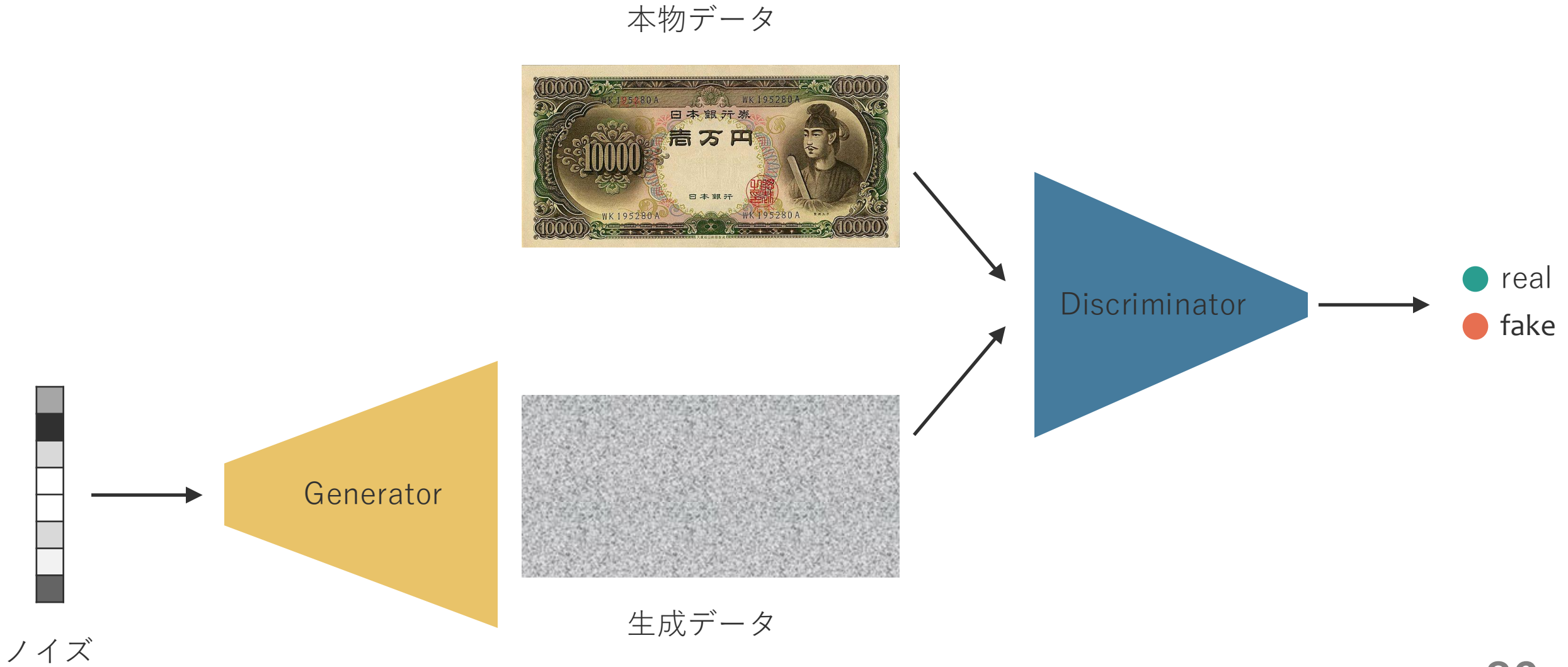
畳み込みニューラルネットワーク

物体分類

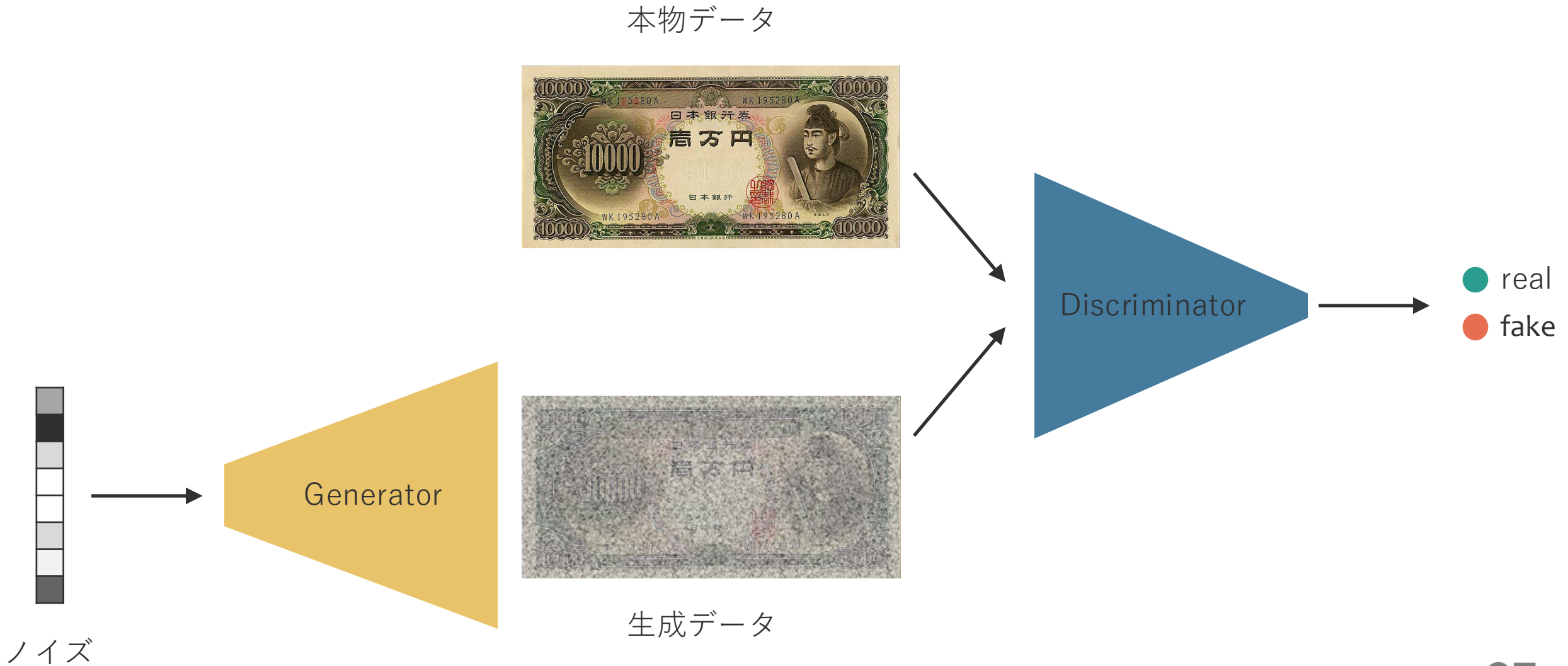
物体検出・セグメンテーション

敵対的生成ネットワーク

# 敵対的生成ネットワーク

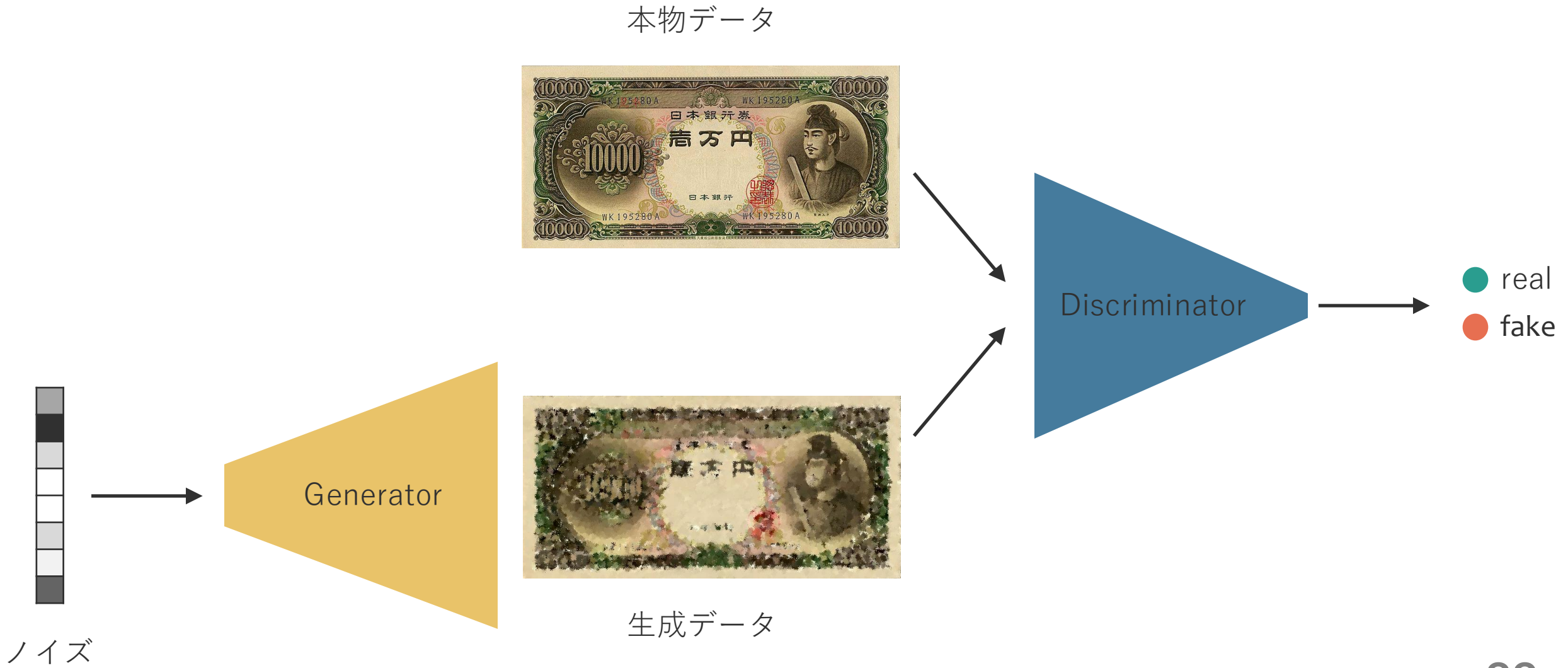


# 敵対的生成ネットワーク



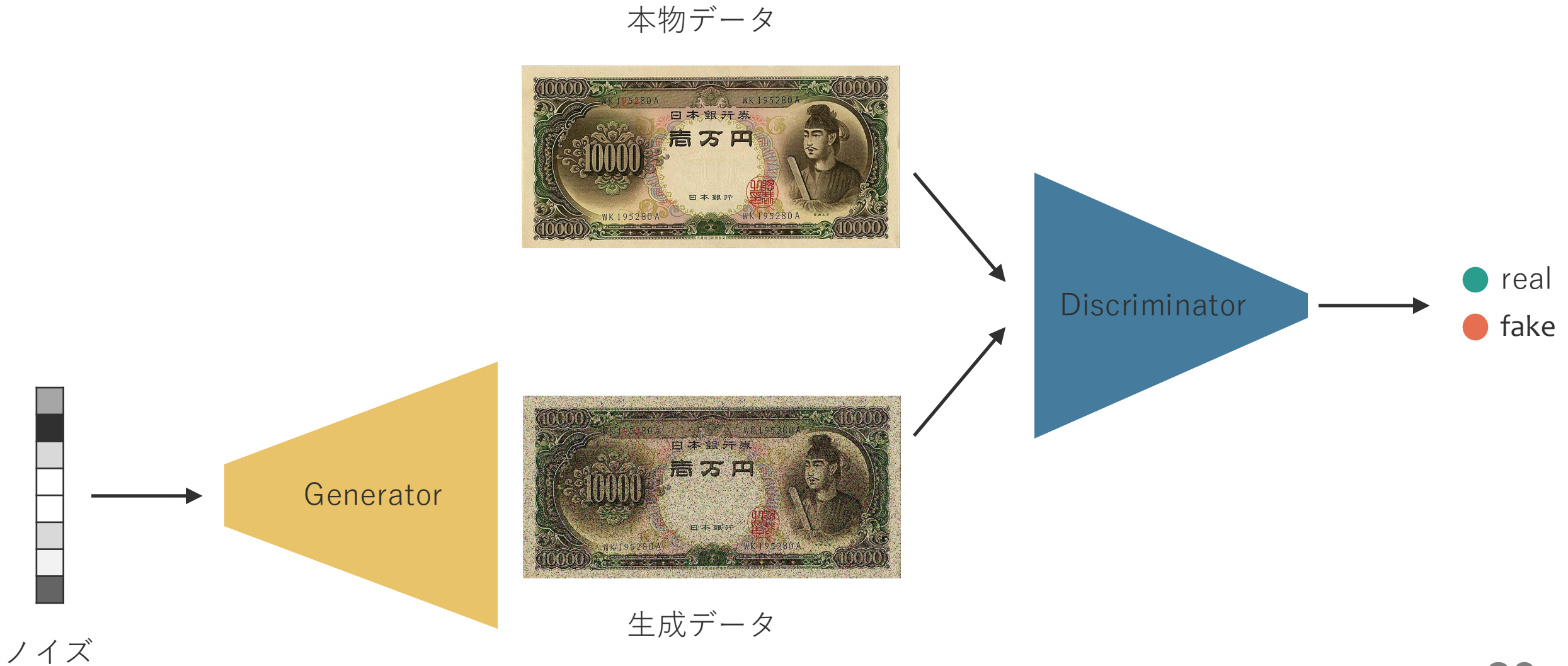


# 敵対的生成ネットワーク

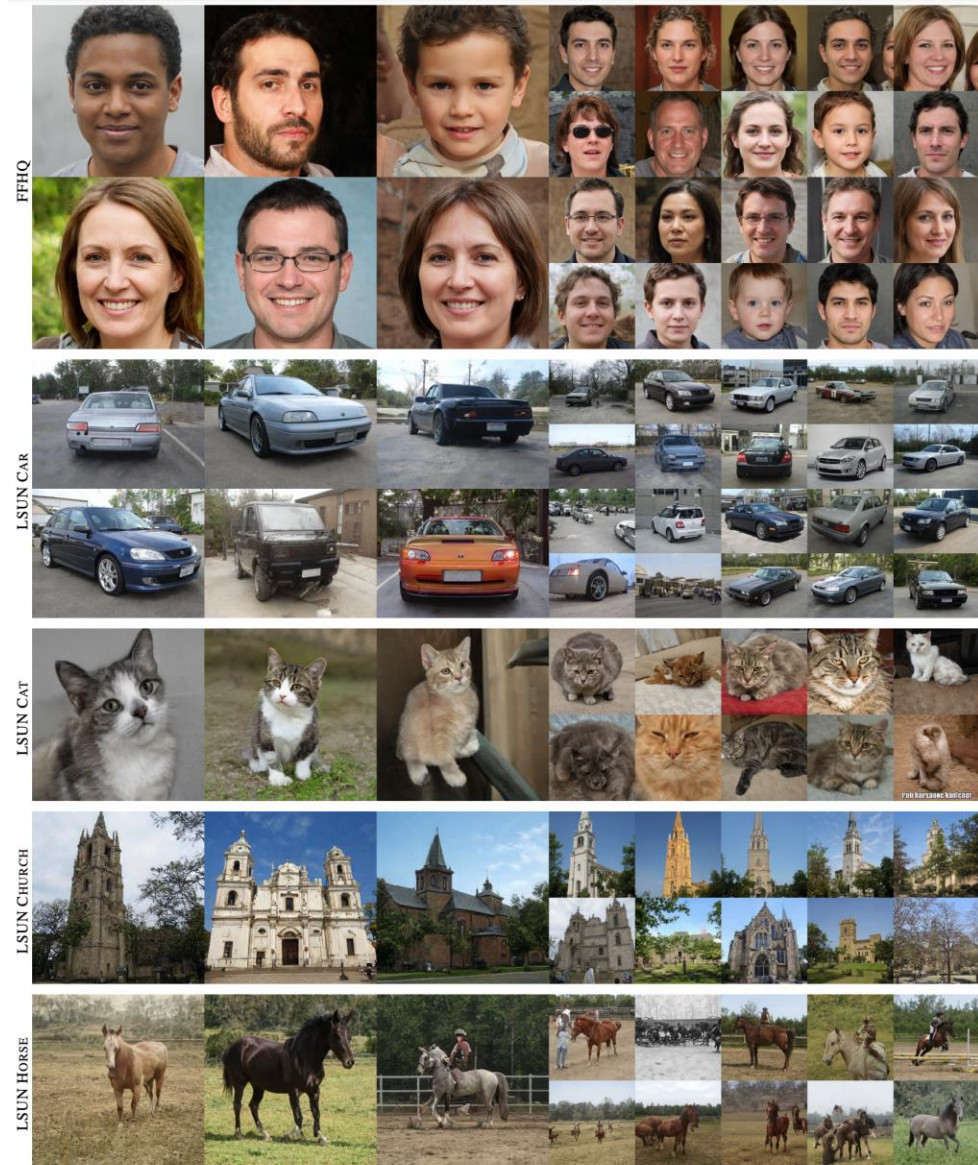




# 敵対的生成ネットワーク



# StyleGAN2



Karras et al., 2020. Fig 12



# CycleGAN

Monet  $\leftrightarrow$  Photos



Monet  $\rightarrow$  photo

Zebras  $\leftrightarrow$  Horses



zebra  $\rightarrow$  horse

Summer  $\leftrightarrow$  Winter



summer  $\rightarrow$  winter



photo  $\rightarrow$  Monet



horse  $\rightarrow$  zebra



winter  $\rightarrow$  summer



Photograph



Monet



Van Gogh



Cezanne



Ukiyo-e



# HDGAN

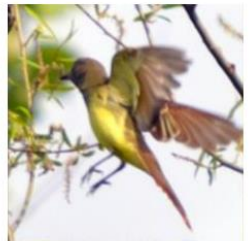
Input

This little bird has a white breast and belly, with a gray crown and black secondaries



GroundTruth

A small, dull brown backed and yellow breasted bird, with a brown head



GroundTruth

7 samples

