

# 深層學習俯瞰



AVILEN

- 1 深層学習とは
- 2 画像処理
- 3 自然言語処理
- 4 強化学習

# 1 深層学習とは

## 2 画像処理

## 3 自然言語処理

## 4 強化学習

## はじめに

Q

深層学習ってなに？

## はじめに

Q

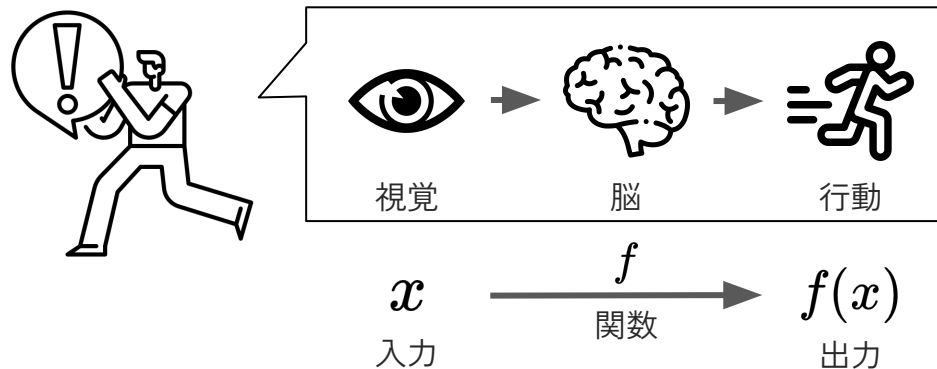
深層学習ってなに？

A

脳内の情報処理を模したニューラルネットワークを使う学習のこと。

## 脳内の情報処理をコンピュータで再現

人間の脳は1種の関数ともみなすことができる

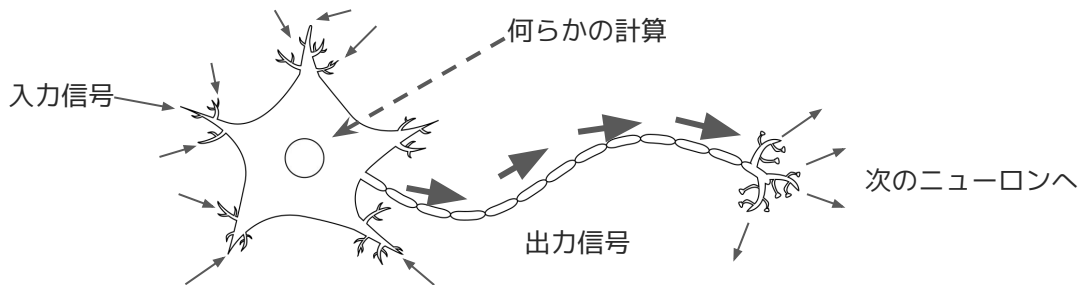


脳内では**電気信号で情報処理**している

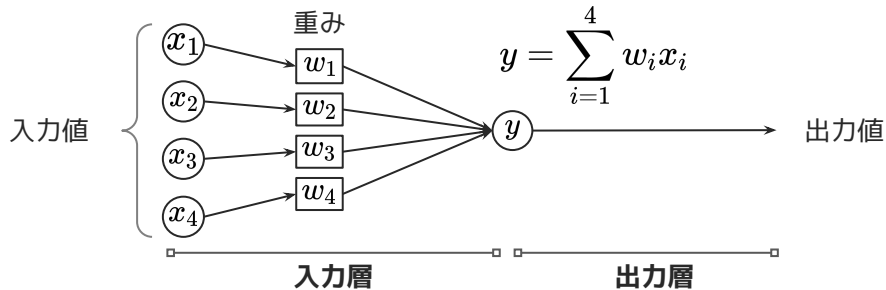
→ これを真似して、コンピュータで再現できそう

ニューラルネットワーク (NN) では、  
脳内のニューロンを模倣した計算を行っている

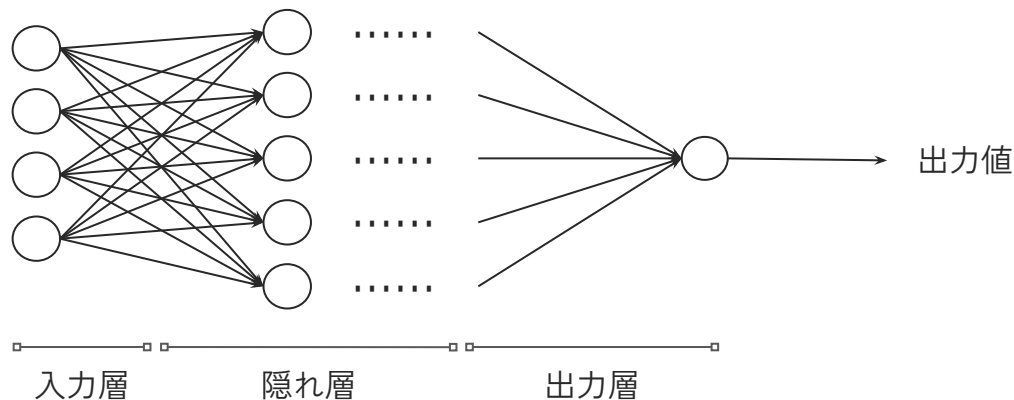
### ニューロン



### 単純パーセプトロン



## DNNと深層学習



### ディープニューラルネットワーク(DNN)

単純パーセプトロンを組み合わせで入力層、隠れ層、出力層を形成したもの

### 深層学習(ディープラーニング)

DNNが尤もらしい出力をするための最適な重みの組み合わせを勾配法により探索していく学習法



## まとめ

01

ニューラルネットワークとは  
**脳内のニューロン**を模倣し計算を行う数理モデル

02

深層学習では  
**ディープニューラルネットワーク**を使って学習を行う

1 深層学習とは

2 画像処理

3 自然言語処理

4 強化学習

## はじめに

Q

画像処理ってなに？

## はじめに

Q

画像処理ってなに？

A

画像を入力とするタスクを指し、  
4種類に分けられる。

# 画像処理タスクは大きく分けて4つある

## 画像認識

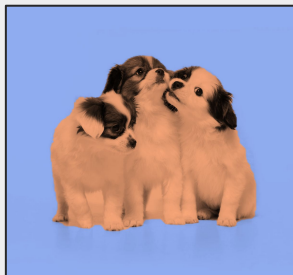
画像に何が写っているかを認識する



犬

## セグメンテーション

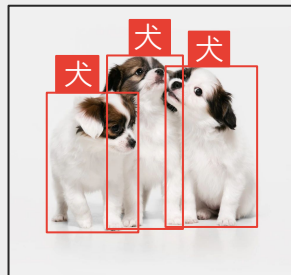
画像のどこに何が写っているかを  
**画素レベル**で分類する



■ 背景 ■ 犬

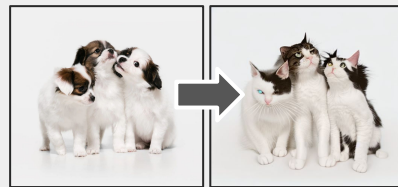
## 物体検出

画像のどこに何が写っているかを  
**長方形**で分類する



## 画像生成

画像を変換するなど  
新しい画像を作る



元画像

新しい画像

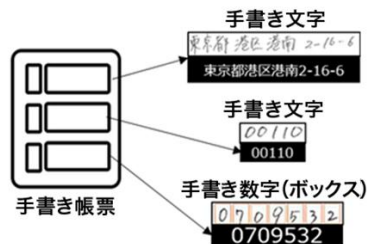
# 画像処理の応用例

## ～ 画像認識・セグメンテーション ～

### 画像認識の応用例

画像内の文字を読み取り、テキストデータとして抽出する「AI-OCR」

深層学習を用いることで文字認識率の精度が向上した。

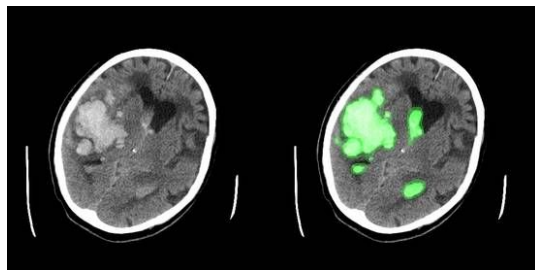


[1] 手書き文字/数字を認識するAI-OCR

### セグメンテーションの応用例

エルピクセルの医用画像解析ソフトウェア「EIRL Brain Segmentation」

頭部CT画像から脳出血部位を自動抽出する。



[2] エルピクセルの「EIRL Brain Segmentation」

[1] Canon 「AI OCRとは？成功のポイントは認識精度を理解し業務全体を見直すこと」 (2023-01-20) <https://canon.jp/business/trend/ai-ocr>

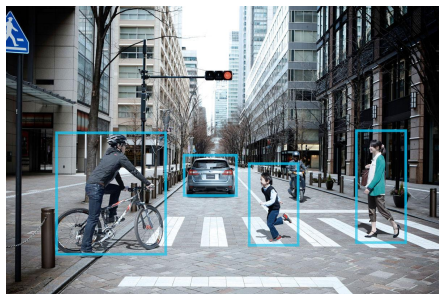
[2] 日経メディカル 「頭部CT画像の脳出血部位をAIが自動抽出」 (2021-06-18) <https://medical.nikkeibp.co.jp/leaf/mem/pub/report/t285/202106/570785.html>

# 画像処理の応用例 ～ 物体検出・画像生成～

## 物体検出の応用例

自動車製造会社SUBARUの運転支援システム  
「アイサイト」

自動車や歩行者、2輪車の場所を検知する。



[1] SUBARUの「アイサイト」

## 画像生成の応用例

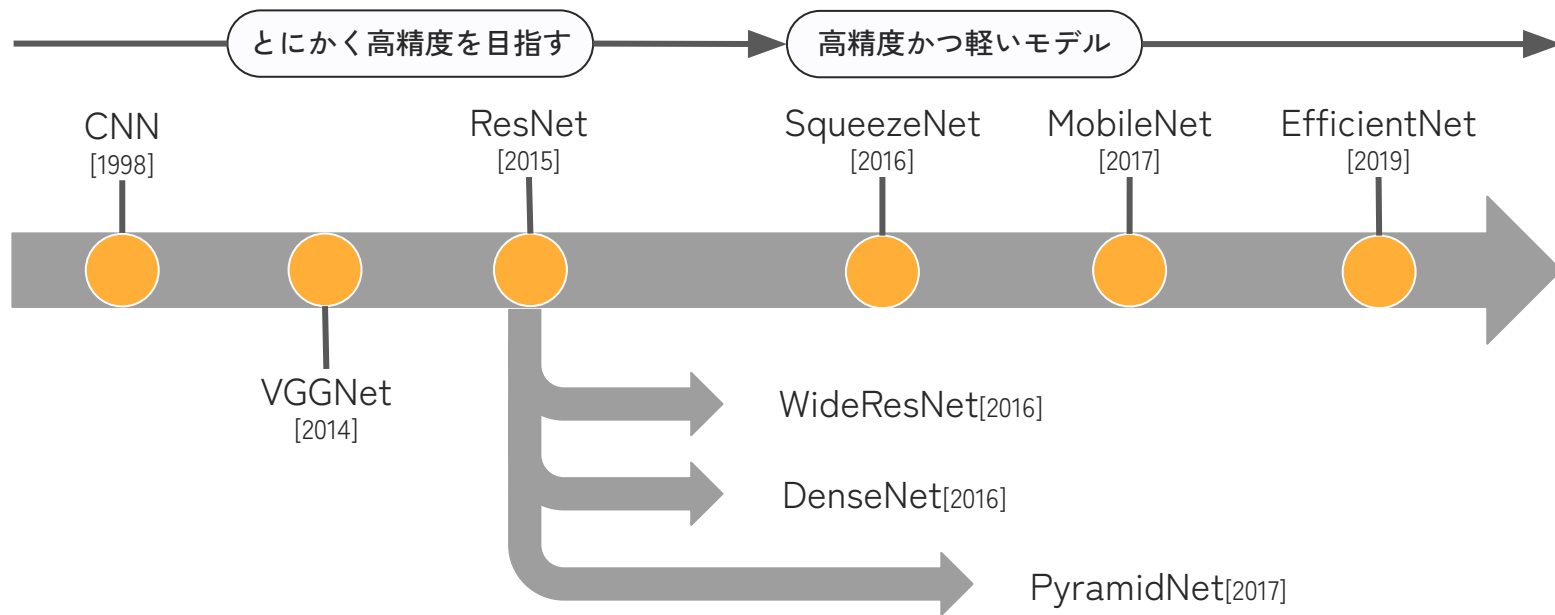
テキストから画像を生成するサービス  
「Stable Diffusion」

右図は、  
「white dog with glasses」  
というテキストから  
Stable Diffusionが生成  
したイラスト。



[1] オートブレーブ「スバル「アイサイト」搭載車が500万台達成 ～その歴史と未来～」(2022-09-07) <https://autoprove.net/japanese-car/subaru/211253/>

## 画像認識モデルの歴史





# 画像処理の代表的手法

## ～ CNN(Convolutional Neural Network) ～

### 特徴

---

隠れ層として**畳み込み層**や**プーリング層**を用いた  
ニューラルネットワーク

#### 畳み込み層の役割

画像の特徴を抽出する

#### プーリング層の役割

位置が多少ズレても出力が変わらないようにする

### メリット

---

- 隣接するピクセルどうしの位置関係を崩さずに処理できる

### デメリット

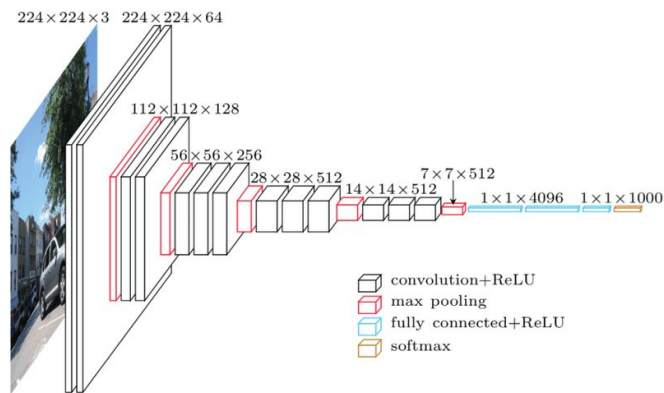
---

- 視点や背景が変わると精度が低下しやすい

# CNNを使った代表的なモデル ～ VGGNet ～

## 特徴

Visual Geometry Groupが2014年に提案したCNNモデル



## メリット

- シンプルなので汎用性がある
- モデルサイズが小さい

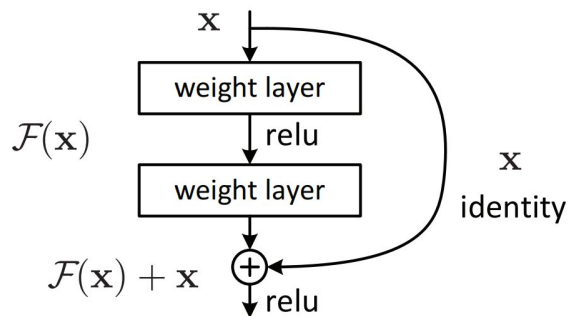
## デメリット

- ResNet(後述)よりは精度が低い

# “深いモデル”が実現可能に ～ ResNet ～

## 特徴

**残差接続**(residual connection)という機構を導入して  
勾配消失問題を解決し、大幅に深いモデルが可能に



## メリット

- 精度が高く、汎用性も高い

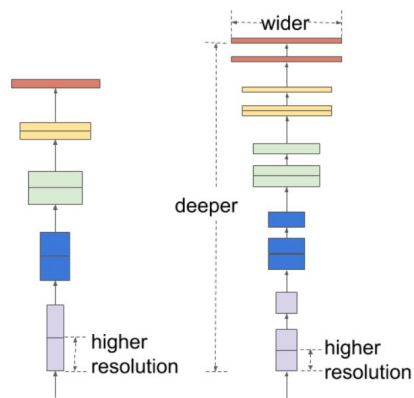
## デメリット

- 層を深くすると学習時間が長くなる

# ResNetの改良版 ～ WideResNet ～

## 特徴

幅を広くすることで、層を深くせずとも  
高精度が出る



## メリット

- 少ない層数、少ないパラメータでも高精度を出すことができる
- 短時間でResNetと同精度を達成

# ピラミッド型の構造 ～ PyramidNet ～

## 特徴

層の幅を”一気に”ではなく  
”徐々に”増やしていくようにしたモデル

ResNet	Wide ResNet	PyramidNet
3x3, 16	3x3, 16	3x3, 16
3x3, 32	3x3, 64	3x3, 42
3x3, 32	3x3, 64	3x3, 69
3x3, 32	3x3, 64	3x3, 96
3x3, 64	3x3, 128	3x3, 122
3x3, 64	3x3, 128	3x3, 149
3x3, 64	3x3, 128	3x3, 176
3x3, 128	3x3, 256	3x3, 202
3x3, 128	3x3, 256	3x3, 299
3x3, 128	3x3, 256	3x3, 256

## メリット

- 精度が高く、汎用性も高い

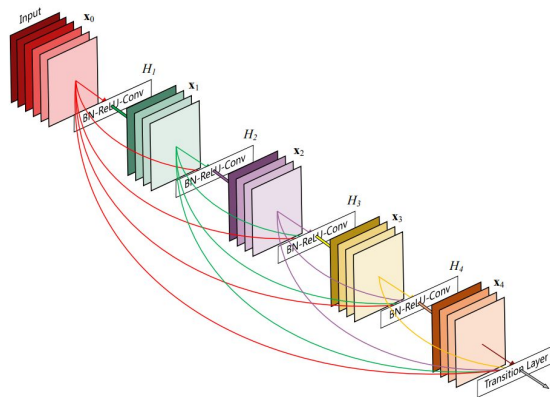
## デメリット

- 層を深くすると学習時間が長くなる

# 残差接続を密に ～ DenseNet ～

## 特徴

ResNetにおける**残差接続**を、  
全ての畳み込み層の間に施したネットワーク



## メリット

- ❑ ResNetよりも層を深くできる
- ❑ ResNetよりもモデルが小さくしかも精度も高い

## デメリット

- ❑ 計算コストが大きい

# スマホで動くほど軽量 ～ MobileNet ～

## 特徴

精度をなるべく保ちつつも、  
モデルを軽量化することを目指して設計された

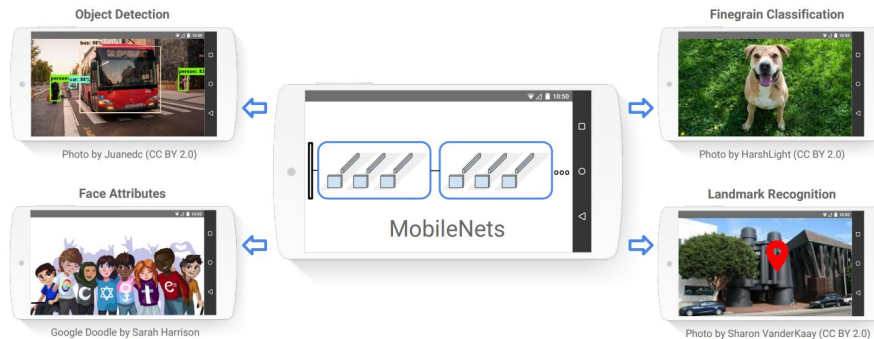


Figure 1. MobileNet models can be applied to various recognition tasks for efficient on device intelligence.

Howard, A. G. et, al. (2017). Mobilenets: Efficient convolutional neural networks for mobile vision applications. arXiv preprint arXiv:1704.04861.

## メリット

- ❑ 学習時間が圧倒的に短い
- ❑ モデルサイズも小さい

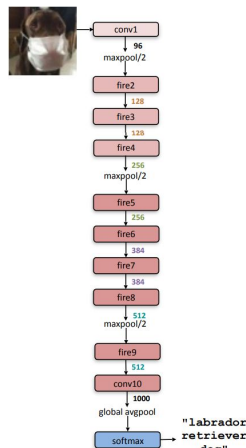
## デメリット

- ❑ 精度はResNetに劣る

# さらなる軽量化を実現 ～ SqueezeNet ～

## 特徴

Fire Moduleを導入して計算コストを削減



## メリット

- ❑ 学習時間が圧倒的に短い
- ❑ モデルサイズも小さい

## デメリット

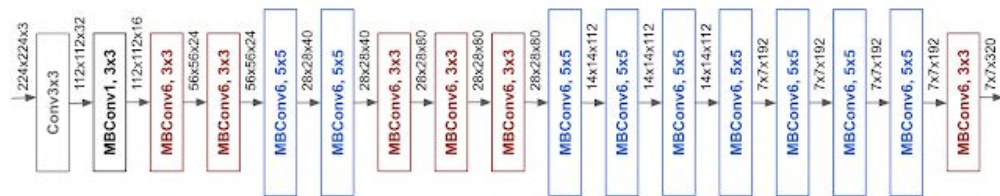
- ❑ 精度はResNetに劣る



# 軽量かつ高精度 ～ EfficientNet ～

## 特徴

少ないパラメータ数で、当時(2019年)の  
最先端モデルを超える精度を達成



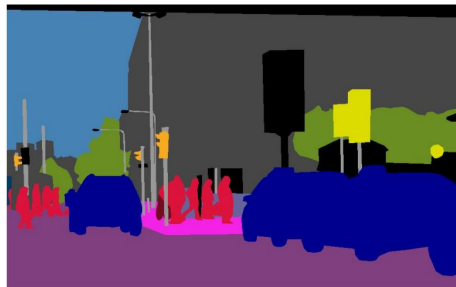
## メリット

- 小さいモデルサイズで高精度が出せる
- 転移学習によく利用される

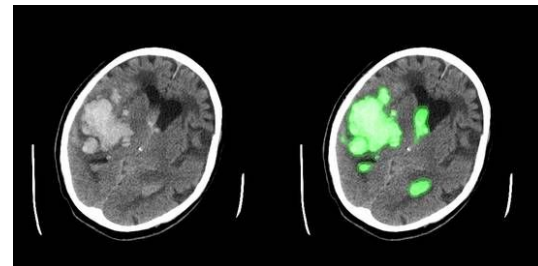
## セグメンテーションの概要

### 応用例

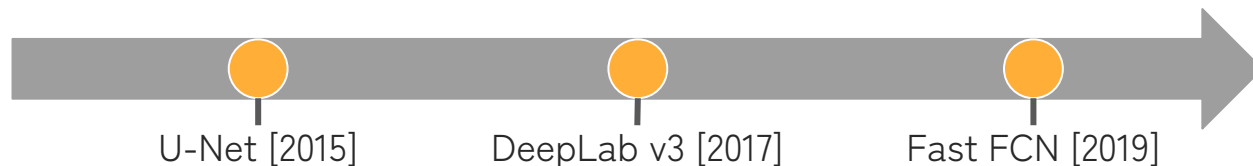
- ❑ 自動運転
- ❑ 医療画像分析
- ❑ 異常検知



[1] 自動運転



[2] 脳出血部位の検出



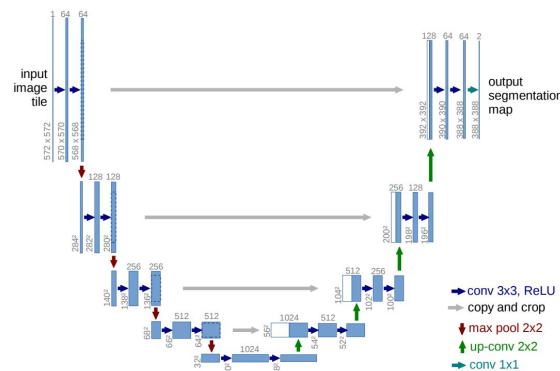
[1] A. Kirillov, K. He, R. Girshick, C. Rother, and P. Dollar: Panoptic segmentation, In "Computer Vision and Pattern Recognition, (2019) の図1から引用

[2] 日経メディカル「頭部CT画像の脳出血部位をAIが自動抽出」(2021-06-18) <https://medical.nikkeibp.co.jp/leaf/mem/pub/report/t285/202106/570785.html>

# セグメンテーションの代表的なモデル ～ U-Net ～

## 特徴

**スキップ接続**(skip connection)によって  
入力側の層で得た特徴を、出力側の層へ伝える



## メリット

- 局所的な特徴を捉えることができる

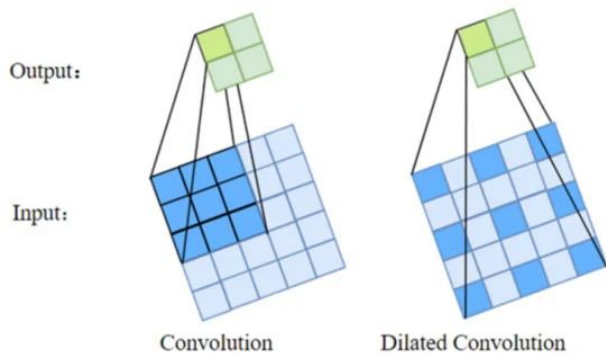
## デメリット

- 計算コストが高い

# 物体の境界を正確に分離 ～ DeepLab v3 ～

## 特徴

**膨張畳み込み**(Dilated Convolution)によって  
画像の広範囲から情報を集めることができる



## メリット

- 物体の境界線付近の予測精度が改善

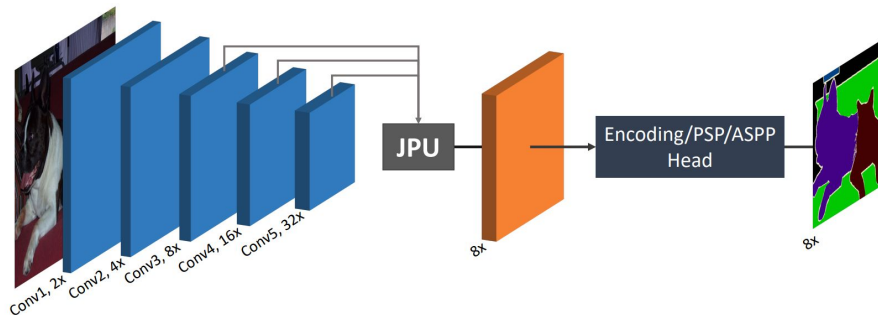
## デメリット

- 計算コストが高い
- メモリ使用量が多い

# 計算コストを削減 ～ FastFCN ～

## 特徴

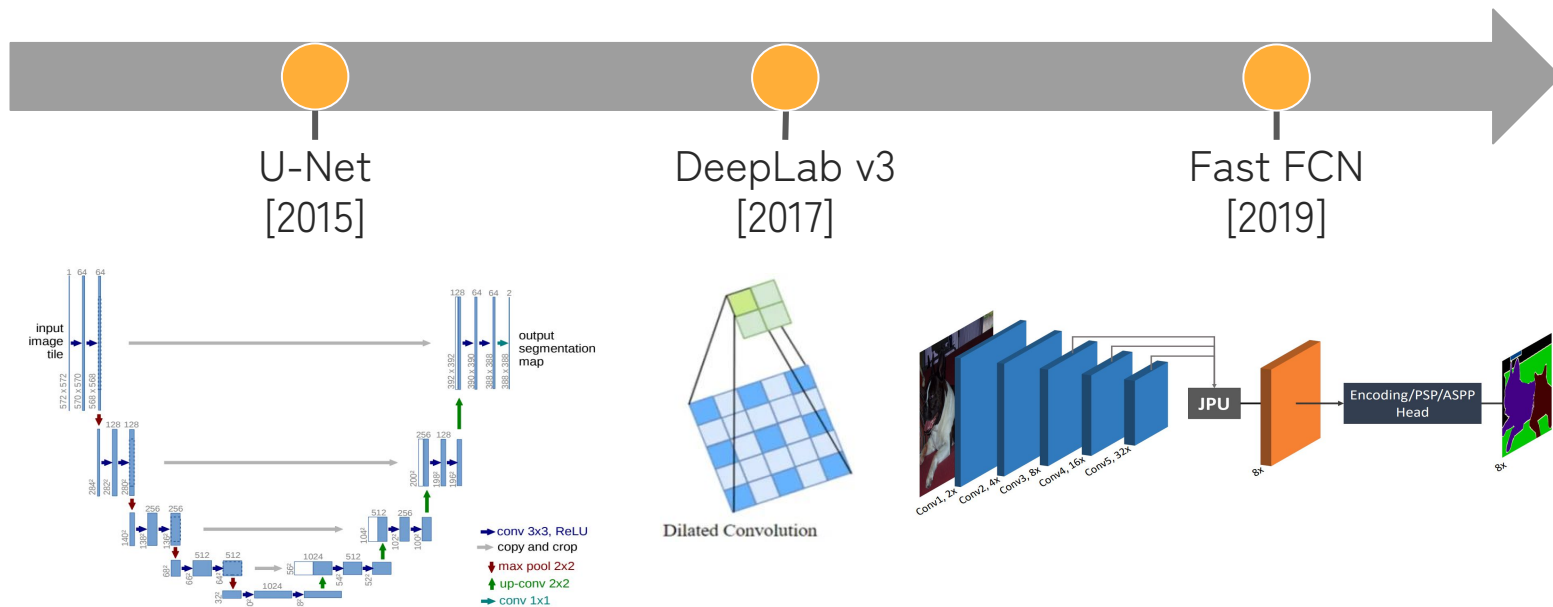
**Joint Pyramid Upsampling(JPU)**を  
膨張畳み込みの代わりに用いて、計算コストを削減



## メリット

- ❑ 精度はDeepLabと同等以上
- ❑ 計算コストはDeepLabの1/3以下

## セグメンテーションのまとめ



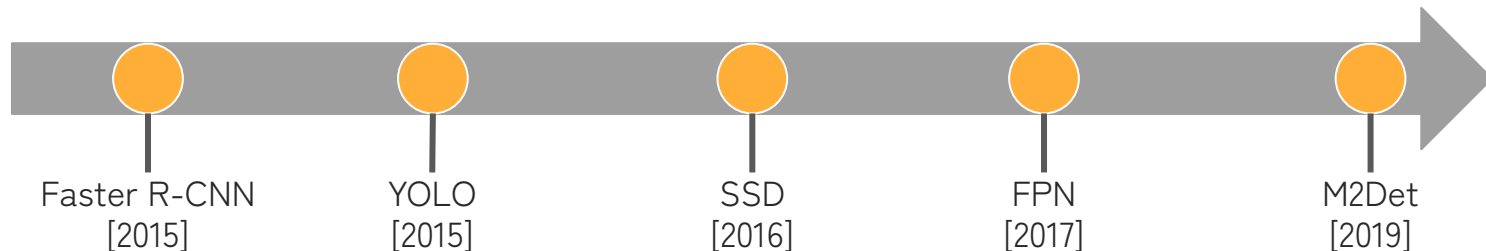
## 物体検出手法の概要

リアルタイムでの検出が求められることが多い。(例:自動運転、ロボット制御など)

→ **高速な処理**が必要

しかし、検出速度と分類精度は**トレードオフ**の関係にあることが多い。

→ 目的に応じたモデルの選択



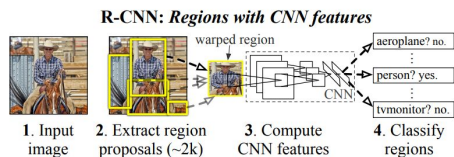
# End-to-Endな学習に初めて成功 ～ Faster R-CNN ～

## 特徴

物体検出における**End-to-End**学習に初めて成功

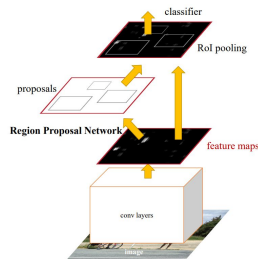
従来(R-CNNなど)<sup>[1]</sup>

各工程を別々のモデルで処理



Faster R-CNN<sup>[2]</sup>

1つのモデルで処理



## メリット

- ❑ 従来のモデル(R-CNN)に比べて検出速度が速く、精度も高い

## デメリット

- ❑ まだ処理速度に問題が残る

[1] Girshick, R. et al. (2014). Rich feature hierarchies for accurate object detection and semantic segmentation. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 580-587).

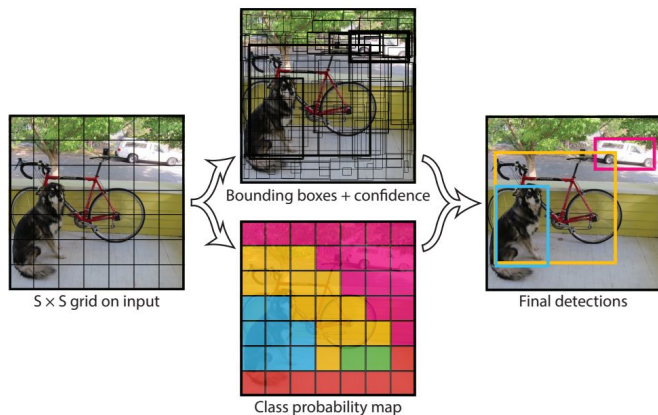
[2] Ren, S., et al. (2015). Faster r-cnn: Towards real-time object detection with region proposal networks. Advances in neural information processing systems, 28.



# リアルタイムな検出 ～ YOLO(You Only Look Once) ～

## 特徴

検出と識別を同時に行い、処理を高速化



## メリット

- ❑ 検出速度が非常に速い
- ❑ 背景の誤検出が少ない

## デメリット

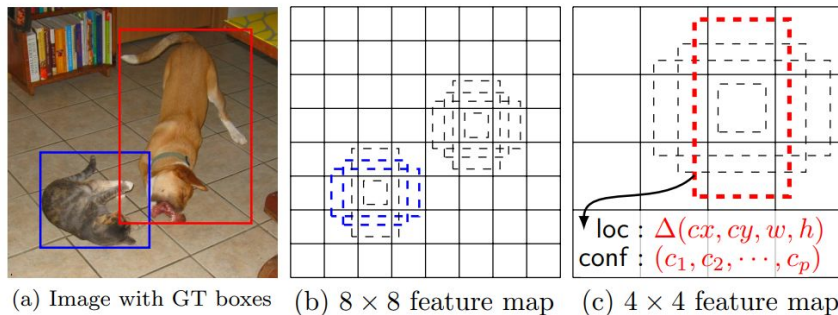
- ❑ 物体が密集した画像に弱い
- ❑ 色々なサイズの物体に弱い

# YOLOの改良モデル

## ～ SSD(Single Shot Multibox Detector) ～

### 特徴

様々なサイズに対応できるようにYOLOを改良



### メリット

- Faster R-CNNと同等の精度
- 低解像度な画像でも検出可能

### デメリット

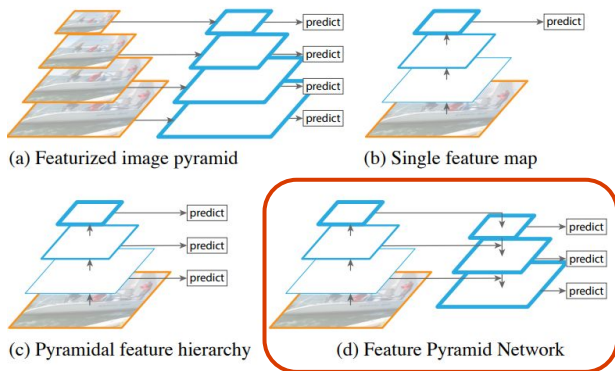
- 処理速度はYOLOにやや劣る

# 小サイズの物体検出

## ～ FPN(Feature Pyramid Network) ～

### 特徴

スキップ接続を採用し、  
より多様な物体を検出可能に



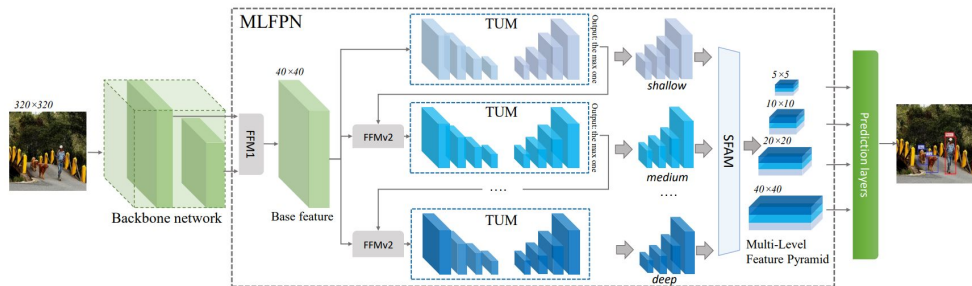
### メリット

- 小さな物体の検出精度が向上
- 計算コストが比較的低い

# 高精度でしかも高速 ～ M2Det ～

## 特徴

Muti-Level Feature Pyramid Network(MLFPN)を  
組み込んだモデル



## メリット

- 精度が非常に高く、検出速度も速い

# 画像生成手法の概要

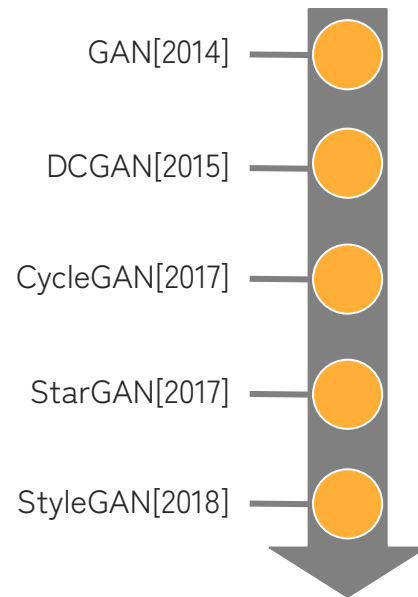
## 応用例

### データ拡張

生成した画像を画像認識AIの学習データとして追加する

### 画像編集

画像の加工や変換を行う



# 画像生成の代表的手法

## ～ GAN(Generative Adversarial Network) ～

### 特徴

**Generator(生成器)**と**Discriminator(識別器)**とを競わせることで、本物らしい画像を生成する

#### Generator(生成器)

Discriminatorが偽物と見抜けないよう画像生成する

#### Discriminator(識別器)

本物の画像かGeneratorが生成した画像かを識別する

### メリット

- ❑ 従来の生成手法(VAEなど)よりも高画質の画像が生成できる

### デメリット

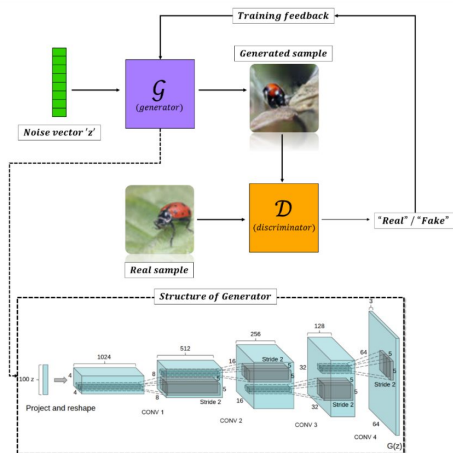
- ❑ Generator/Discriminatorの一方が優位になりやすく、学習が不安定

# CNNをGANに導入

## ～ DCGAN(Deep Convolutional GAN) ～

### 特徴

GeneratorとDiscriminatorにCNNを導入



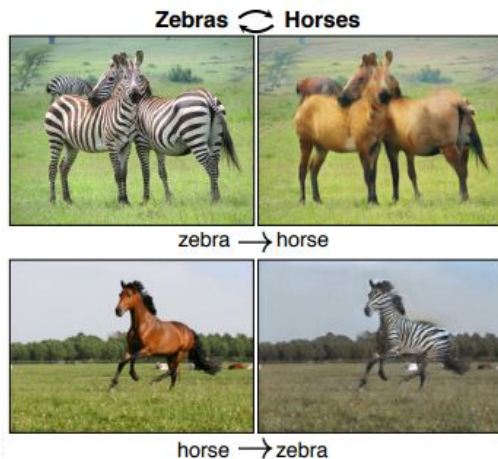
### メリット

- ❑ 従来のGANよりも解像度が高い画像が生成できる
- ❑ CNNの特徴マップを利用したモデルの解釈が可能

# 画像のスタイルを相互変換 ～ CycleGAN ～

## 特徴

2つのスタイル間で画像を相互変換



## メリット

- ❑ 双方向の変換が可能
- ❑ 学習画像としてペア画像を用意する必要がない

## デメリット

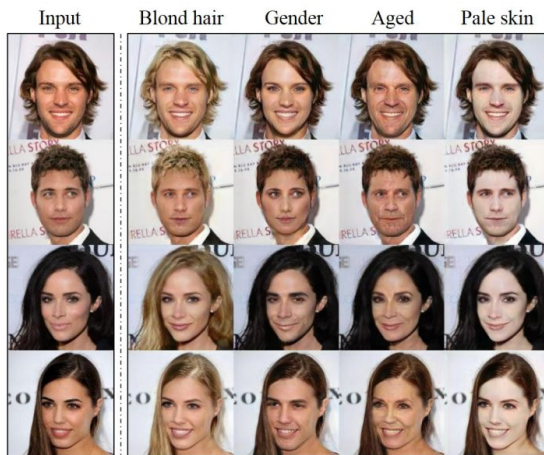
- ❑ 1つのモデルで1つのペア変換しかできない



# 複数のスタイル変換 ～ StarGAN ～

## 特徴

1つのモデルで複数のスタイル変換を可能に



## メリット

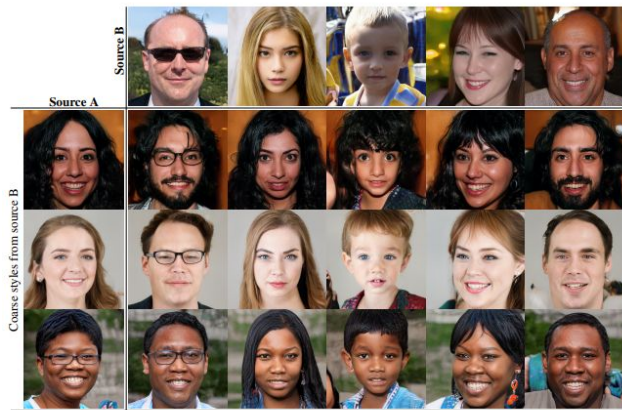
- 複数のモデルを作る必要がない

Choi, Y., Choi, M., Kim, M., Ha, J. W., Kim, S., & Choo, J. (2018). Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 8789-8797).

# 高解像度な画像を生成 ～ StyleGAN ～

## 特徴

本物の写真と見分けられないほど  
高解像度で自然な画像を生成



## メリット

- 2つの画像のスタイルを混ぜた画像を生成
- 高解像度で自然な画像を生成

Karras, T., Laine, S., & Aila, T. (2019). A style-based generator architecture for generative adversarial networks. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (pp. 4401-4410).

## まとめ

### 01

画像処理タスクは

- ❑ 画像認識
- ❑ セグメンテーション
- ❑ 物体検出
- ❑ 画像生成

の4つに分けられる

1 深層学習とは

2 画像処理

3 自然言語処理

4 強化学習

## はじめに

Q

自然言語処理ってなに？

## はじめに

Q

自然言語処理ってなに？

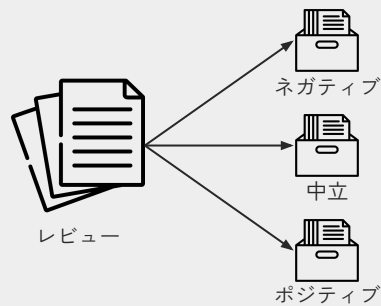
A

人間が普段使っている言語（自然言語）をコンピュータが処理する技術のこと。

## 自然言語処理でできること

### 文書分類

文書をいくつかの  
カテゴリーに分類



### 機械翻訳

文章を別の言語へ翻訳

私は英語を学んでいます。



I study English.

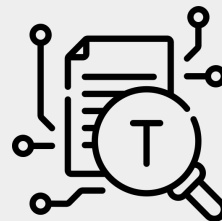
### 質問応答

- ・チャットボット
- ・音声アシスタント



### その他

- ・文章要約
- ・感情分析
- ・文書検索



## 自然言語処理モデルの歴史

### 再帰型NNの提案

**RNN**

[1986]

**LSTM**

[1997]

### ベクトル表現

**word2vec**

[2013]

**seq2seq**

[2014]

### Attentionの時代

**Attention**

[2015]

**Transformer**

[2017]

**BERT**

[2018]

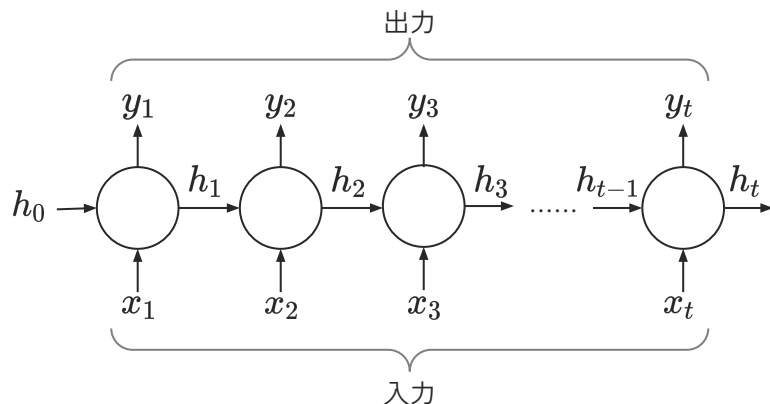


# 時系列データで活躍

## ～ RNN(Recurrent Neural Network) ～

### 特徴

1つ前までの情報処理の結果と新たな入力とを  
合わせて処理し、次のセルへ送る



### メリット

- 直前の文脈を考慮できる

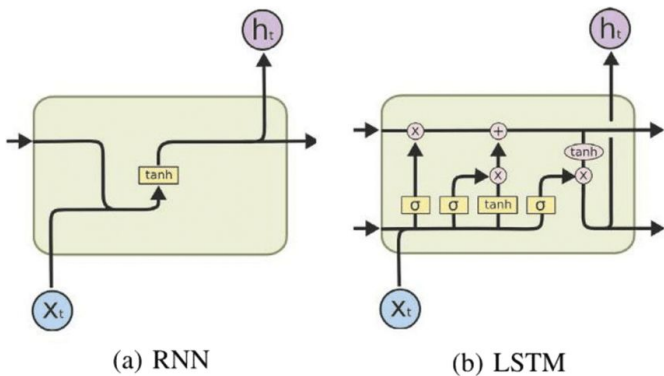
### デメリット

- 記憶を長期間保持できない
- 逐次的に計算するため、  
並列計算ができない

# より前の文脈まで考慮できる ～ LSTM(Long Short Term Memory) ～

## 特徴

何を記憶し何を忘却するかを決定する  
ゲート機構を導入



## メリット

- ❑ 長期間記憶を保持できる

## デメリット

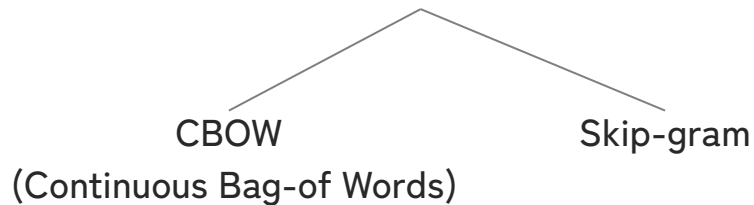
- ❑ RNNより計算量が多い

# 単語をベクトルに変換 ～ word2vec ～

## 特徴

単語をベクトルに変換するための  
分散表現を学習するモデルの総称

主なモデルは2種類



## メリット

- ベクトル表現を使った単語同士の演算ができる

(例)

「王」－「男」＋「女」＝「女王」

「パリ」－「フランス」＋「日本」＝「東京」

# 文字列から文字列を生成 ～ seq2seq ～

## 特徴

**文字列から文字列**を生成するためのモデルの総称。  
Encoder-DecoderモデルやAttention(後述)がある。

「seq → seq」を  
「seq → vec」「vec → seq」に分けて考える。

### (例)機械翻訳

This is a book. → [0.2, 0.5, -1.5, 2.6, ..., 0.1] → これは本です。  
↑ (意味を表すベクトル、固定長) ↑  
**Encoder**で変換 **Decoder**で変換

## デメリット

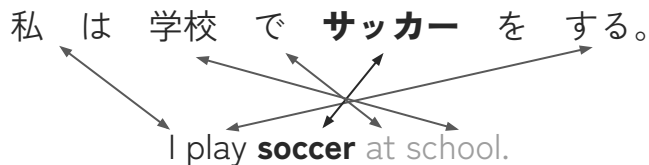
- ベクトルの次元が固定されているため、長文の場合ベクトルに入りきらない。

## 重要な部分に注意を向ける ～ Attention ～

### 特徴

入力側のどの部分に**注目**するかを判断し、注目したデータを基に出力を計算する機構。

#### (例) 機械翻訳



次に「soccer」を出力するときに、入力側の「**サッカー**」に注目している

### メリット

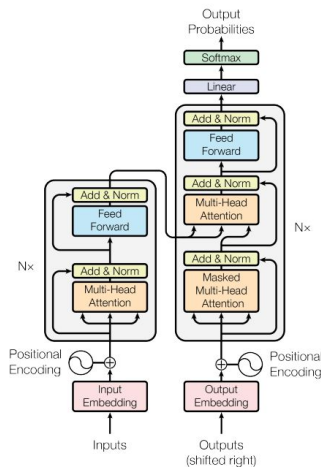
- 入力側(Encoder)のベクトルの次元に関係なく、Decoderに情報伝達できる。

# RNN/CNNを使わずにAttentionだけ ～ Transformer ～

## 特徴

RNNやCNNを使わずに、Attentionだけを使った機構。

自然言語分野だけではなく、  
画像認識、強化学習の分野にも  
応用されている。



## メリット

- 逐次的に計算するRNNを使わないため、並列計算が可能



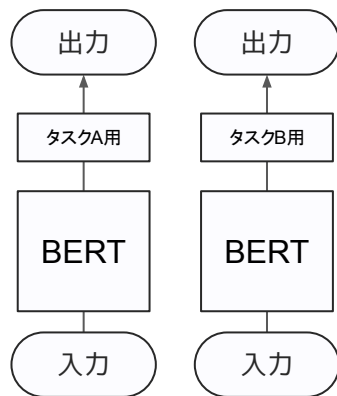
計算速度が向上

# 汎用性が高いモデル

～ BERT (Bidirectional Encoder Representations from Transformers) ～

## 特徴

Transformerをベースにしたモデル。幅広く応用されており、2019年には**Google検索エンジン**にも導入

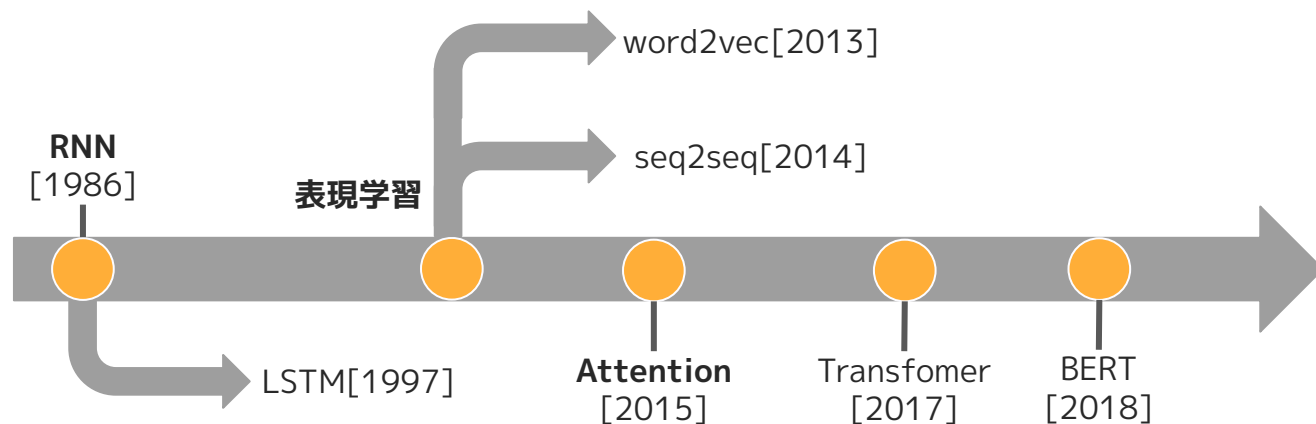


最後の1層だけタスクごとに用意  
→ ファインチューニングして完成！

## メリット

- 様々な自然言語処理タスクに対応可能
- ファインチューニングが簡単
- 圧倒的な精度

## 自然言語処理のまとめ





## まとめ

01

①再帰型NNや②分散表現を学習するモデル、  
③Attentionを用いたモデルが提案されている

02

Attentionを用いたモデルは  
Google翻訳や検索エンジンにも導入されている

1 深層学習とは

2 画像処理

3 自然言語処理

4 強化学習

## はじめに

Q

強化学習ってなに？

## はじめに

Q

強化学習ってなに？

A

明確な正解がわからないときに、  
行動を起こした結果から学習する仕組み

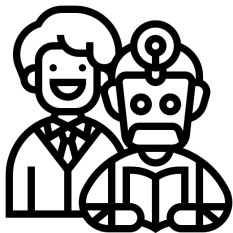
## 教師あり学習との違い

### 教師あり学習

学習データに**正解**を与える。

→ パターンを抽出して学習

これは「1」って読むんだよ



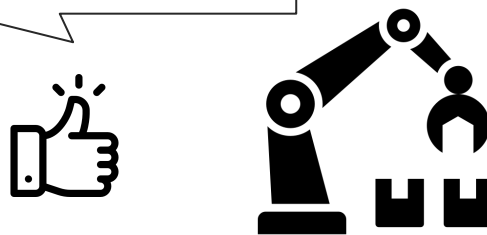
明確な正解が存在する

### 強化学習

正解を与える代わりに、**報酬**を与える。

→ 報酬を最大化するように学習

よし、ちゃんと運べたね！



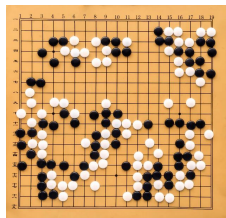
運び方にもいろいろある  
状況によっても変わる

# 強化学習でできること

## ゲーム



ゲーム攻略

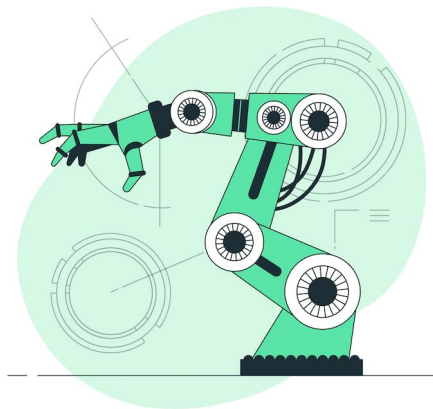


囲碁・将棋

## 自動運転



## ロボット制御



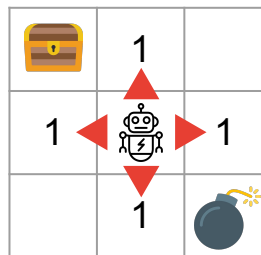
# Q学習 (Q-learning)

**Q値**とよばれる値を、プレイヤーが行動するごとに更新していく手法。

## Q値

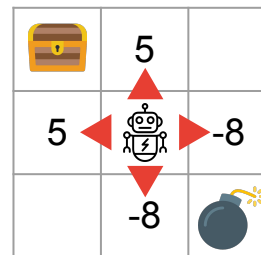
現在の状態sから行動aをしたときに、**報酬をどれくらいもらえそうか**を表す値。

Q値は最初はわからない → 実際に行動していく中でQ値を推定する。



学習初期

→ たくさん行動 →



十分な学習後

# DQN (Deep Q-Network)

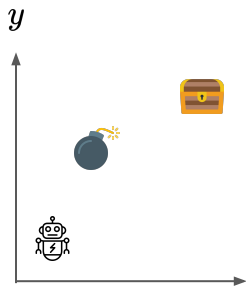
ニューラルネットワークを用いて、Q値を推定する手法。

## Q学習の問題点

状態や行動のパターンが多いと、推定すべきQ値の個数が爆発的に増え、連続の場合は無限個になる。

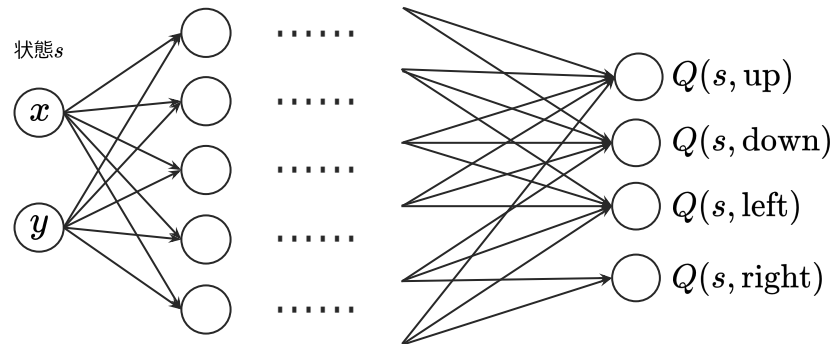


状態の個数：9個



状態の個数： $\infty$ 個

Q値を**Q関数**としてあつかう  
→ 深層学習で関数を予測する

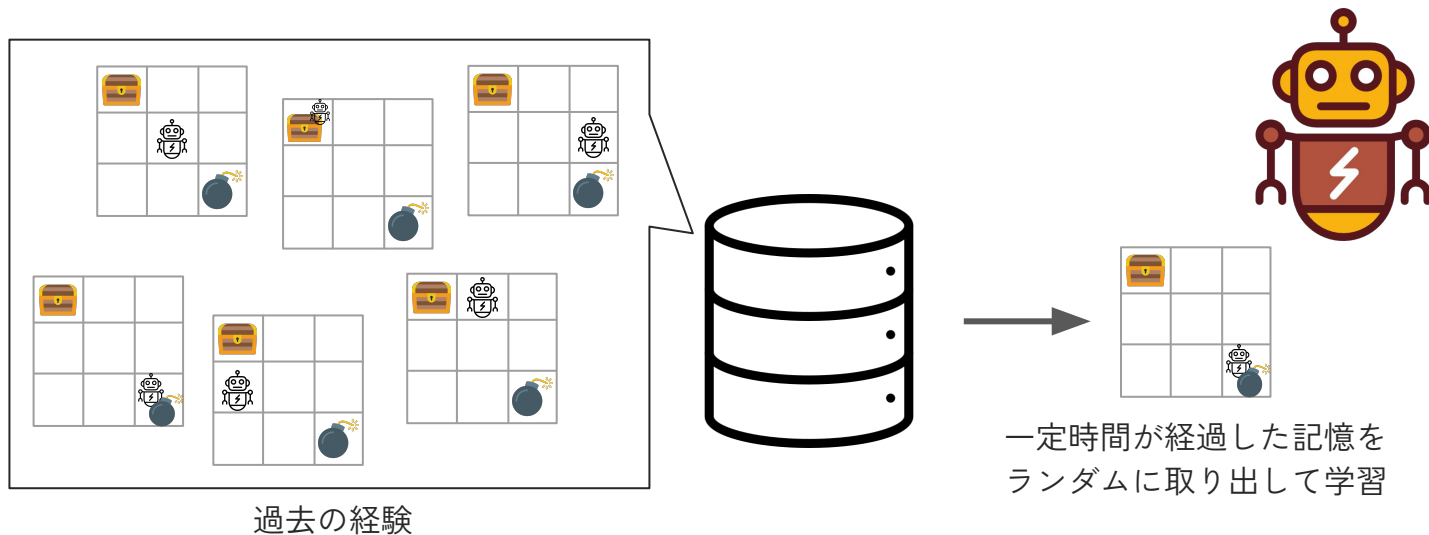




# 経験再生と分散学習

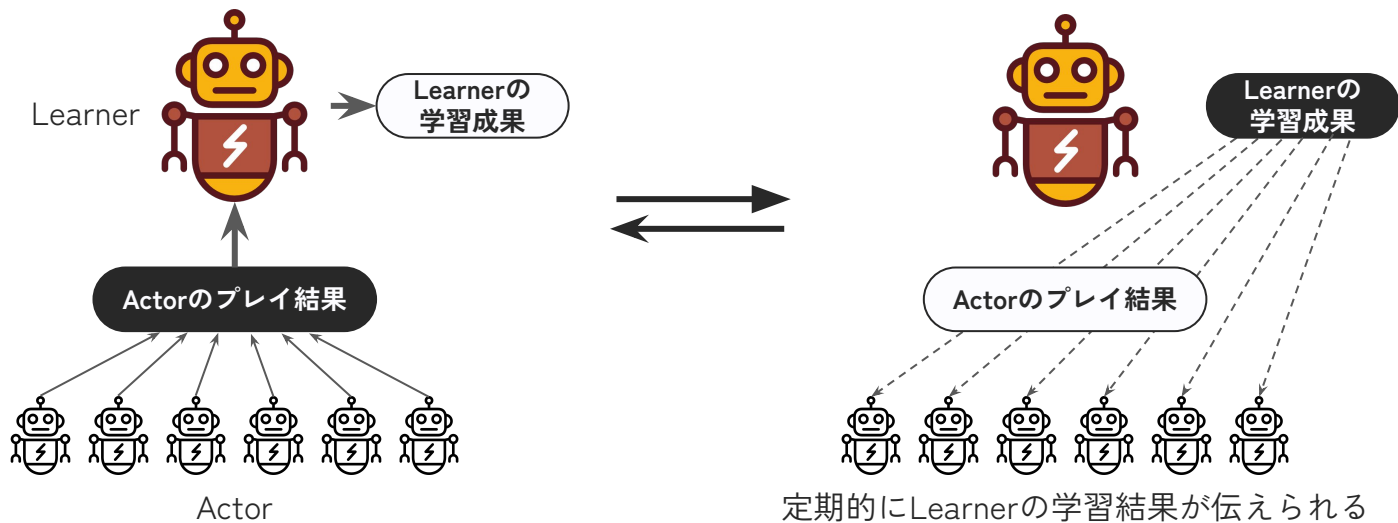
## ～ 経験再生（Experience Replay）～

過去にとった行動を記憶しておき、その行動を定期的にランダムに取り出して学習する手法。



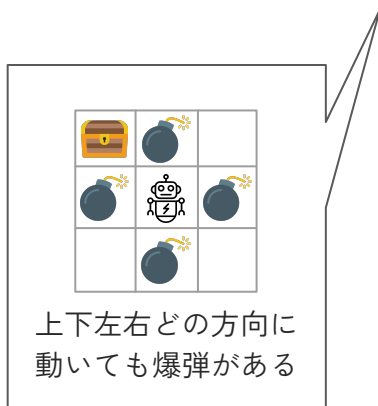
# 経験再生と分散学習 ～ 分散学習 ～

複数のActor(プレイヤー)を並列に行動させることで、学習効率を向上させる手法。



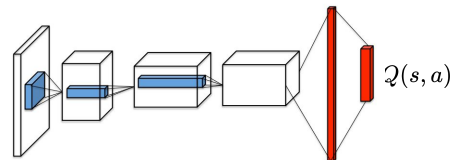
# Dueling Network

行動によってQ値が**変わらない**状態と**変わる**状態を分けるように学習するネットワーク



## 通常のDQN

$Q(s, a)$  を直接学習する

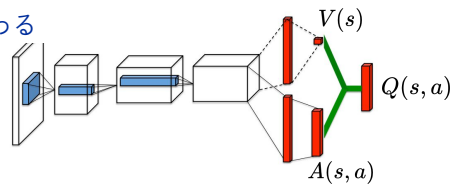


## Dueling Network

行動しだけで変わる

$$Q(s, a) = V(s) + A(s, a)$$

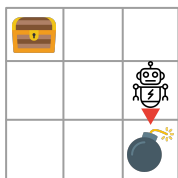
どう行動しても同じ  
と分割して学習する



## Ape-X

優先順位つき経験再生・分散学習・Dueling Networkを活用したモデル。

過去の行動をサンプリングする際、**学習が進みやすくなるものを優先**

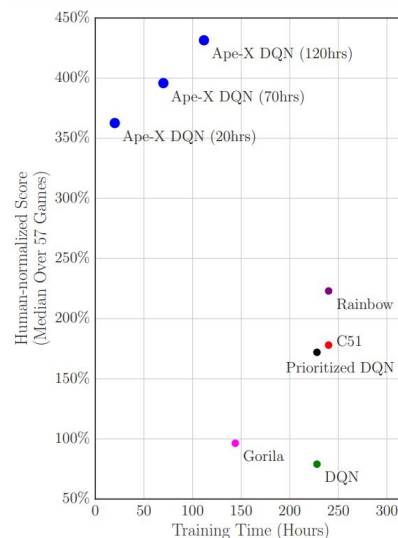


あれ？報酬もらえる  
と予測したのに...

← 予想した報酬と  
実際の報酬の差  
(=TD誤差)が大きい

### 特徴

- ❑ 学習時間が短い
- ❑ 従来のモデルを圧倒する精度

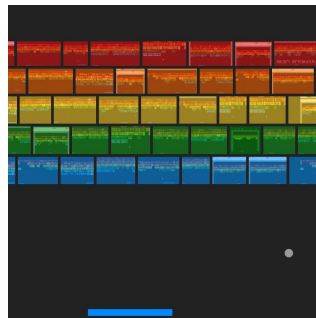


## R2D2(Recurrent Replay Distributed DQN)

Ape-XにLSTMを導入し、時系列の影響を学習するようにしたモデル。

例:ブロック崩しゲーム

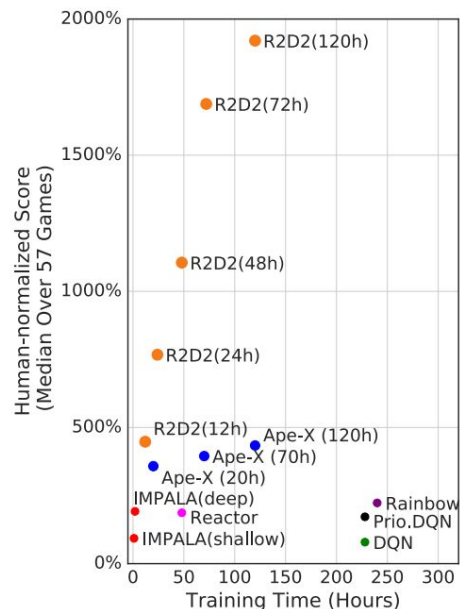
- ボールの進行方向・速さが重要となる。
- 1つの状態を知るだけでは適切な行動がとれない。
- RNNやLSTMを使って、前の状態も考慮。



ブロック崩しゲーム

LSTMを使うと学習が不安定になる問題

→ Stored state や Burn-in という手法で解消

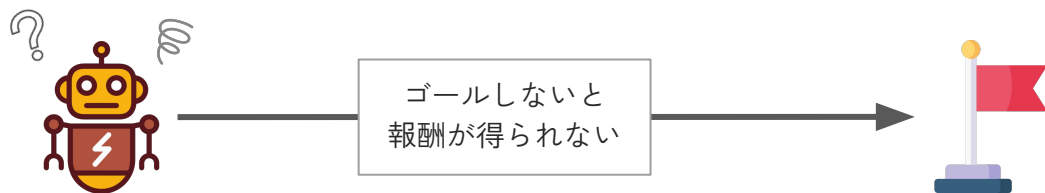


## R2D3(Recurrent Replay Distributed DQN from Demonstrations)

人間のデモプレイ(demonstrations)を参考にして学習させたモデル。

報酬がなかなか得られないような設計

→ 正しい行動がとれるまで暗中模索する必要があり、学習が進まない。



人間のプレイ結果を混ぜることで正しい行動ができるようにガイド

→ 学習がうまく進むようになる。

## まとめ

01

強化学習では  
**報酬**を与え、**報酬を最大化**するように学習する

02

DQN（Deep Q-Network）は  
**ニューラルネットワーク**を用いて**Q値**を推定する手法