

物体検出に用いられるニューラルネットワークモデル

最新モデルのサーベイと目的に応じたモデルの選択

Neural Network Models for Object Detection

A Survey of the Latest Models and Optimal Model Selections for Specific Tasks

金子 純也

Junya Kaneko

Morning Project Samurai 株式会社

Morning Project Samurai Inc.

junya@mpsamurai.com, <http://www.mpsamurai.com>

山田 貢己

Miki Yamada

(同上)

m.yamada@mpsamurai.com

keywords: survey, neural network, object detection, instance segmentation, deep learning

Summary

「ショートノート」は 200 ワード、それ以外は 200 ~ 500 ワード以内の英文で summary を記す (ここは、論文執筆後に書く。)

ここにチートシートを出力する。

1. ま え が き

この論文の狙い: このサーベイは、理工系大学 2 年生程度の数学の知識を前提に、物体検出をこれから始めるにはどうすればよいかという道すじを伝えることを目的として執筆したものである。物体検出でできることは何? から始まり、物体検出をするために必要なもの (ハード、ソフト、データ、知識、明確な目的) を簡潔に纏めてある。また、読者にとっての理想的なサーベイ (即時性、分かりやすさ、一言で説明、参考文献は充実) というものの一つの解として、随時更新される GitHub 上の (日本語で書かれた) 物体検出まとめサイトとライブラリを紹介する。

世の中の状況: 近年、「AI (人工知能)」という言葉が国内外に蔓延しており、技術者のみならず一般の人の日常生活にもすっかり浸透した。常に手の届くところに AI があり、AI に囲まれて生活していると言っても過言ではない。テレビやインターネットの画像は、本物と見間違えるほどの人工画像で溢れ、スマホや机上のスピーカーに話しかけるとあらゆる情報を教えてくれるばかりでなく、電化製品を操作することもできるようになった。高速道路を自動運転する車も増えている。

AI 分野における物体検出: ここ数年で飛躍的に性能を向上させた AI 関連技術は、画像認識、物体検出、ロボット制御、音声認識、機械翻訳、ビッグデータ分析などであり、これらの多くの領域で深層ニューラルネットワー

ク (Deep Neural Network (DNN)) が使われている。とりわけ画像認識分野でこの DNN が注目されるようになったのは、2012 年に開催された最先端の一般物体認識の性能を競うコンテスト ILSVRC において DNN を使った手法が他の手法に大差をつけて優勝したことが発端である。「画像中の物は何か?」に答える物体認識をさらに進めて、「画像中のどこに何があるか?」に答えようとするものが本論文のテーマ「物体検出」であり、現在の AI ブームを巻き起こした源流がここにあると言ってよい。

物体検出でできること: 物体検出とは、カメラで撮影された画像データを電子的に処理し、予め登録しておいた物体 (例えば、人、犬、猫、自動車、飛行機、...) を見つけ出し、その正確な画像上の位置と物体の種類を予測するものである。「予測 (predict)」は「推定 (estimate)」, 「推論 (inference)」などとも呼ばれ全て同じ意味で使われる。

現在の標準的方法においては、予め、検出したい対象の学習データ (物体が写っている画像、物体の種類、物体の位置を示す矩形の座標) を大量に用意し、画像を入力すれば種類と位置を出力するように、ニューラルネットワーク等の予測モデルを学習させる。通常、これに数時間から数日要すると言われている。

学習が完了した予測モデルの能力は、タスクの種類によっては人間の能力 (予測結果のスコアの平均値) を超えたと言われているものもある (物体認識など)。ただし、物体が写し出された画像の品質 (解像度、ノイズ、露出不

足/過多) や撮影アングル (遮蔽物, 変形, 大き (小) 過ぎる) に問題がある場合は性能が低下することは避けられない。

物体検出と似た技術として, 次の 3 つがある:

- セマンティックセグメンテーション (全画素の物体の種類を認識するが, 同種の物体同士は区別しない)
- インスタンスセグメンテーション (同種の異なる個体を区別して物体検出を行い, 且つ, 画素単位で個体の識別をする)
- パノプティックセグメンテーション (全画素の物体の種類を認識し, 同種の異なる個体も区別する)

本論文では上記のセグメンテーション技術も含めた広い意味での物体検出について述べる。

物体検出をするために必要なもの:

ソフトウェア

- PyTorch/TensorFlow 等の深層学習ライブラリとその稼働環境 (Linux/Windows/Mac 上の Python, jupyter notebook 環境など)。
- 物体検出を行うソフトウェア (予測モデルの作者, または, 物体検出を行おうとする担当者が作ったもの)。

ハードウェア

- 前記ソフトウェアが実行できる環境 (PC (GPU があると良い), 或いは, Google Corabulatory などのサーバ上の実行環境)。

データ

- 学習データ (事前学習用, 並びに, fine-tuning 用の入力と出力のペア)
- 本来処理したいデータ (入力)。

これらは, 使用するソフトウェアで読み取ることでできる状態にしておく (データの前処理)。

知識 物体検出のソフトウェアを使うには, 入出力データの意味を理解する必要がある。特に, 予測モデルの出力データは通常は誤差を含むものとなるため, 出力が表す数値が確率値を表すのか, 何らかの物理量を表すのか, 分類のカテゴリを表すのか, 正確に把握する必要がある。また, 学習時の損失関数の値から, 予測モデルの推定誤差を見積もることができるのだが, それには出力結果を正しく解釈できる統計学の知識が必要となる。

明確な目的 何がしたいのかということを明確化しておくことが, 物体検出を行おうとするときに重要となる。物体検出は新しい技術であり, 標準的な統計解析の手法よりも手間と計算コストが大きくなりがちである。他の方法では解決できないのか? と問いかけて, 本当にこれが必要であることを確認しておくべきである。

理想的なサーベイとは: 最新技術のサーベイ論文は, 有用であり様々な分野で昔から (論文雑誌が生まれた頃から) 活用されていると思われる。しかしながら, 進歩の速

い分野においてはサーベイが出た頃には既に内容が古くなってしまっているという問題が往々にして起こる。また, とても良く書かれたサーベイほど内容が濃く多くなり, 執筆に時間と労力を要するのはもちろん, それを読み解くのにも時間を要するということがよくある。

我々は, github 上に随時更新される形式でサーベイを公開することを試みた。出版されたときには既に古くなっているという懸念を取り払える可能性を期待している。また, この分野に新規参入しようとしている人になるべく短時間で必要な情報にたどり着き, 取り組んでいる問題を解決する最適な方法を見つけたり, 或いは, 新たな研究に取り組めることを目指し, 内容の拡充性や緻密性よりも, なるべく視覚的に解りやすいコンパクトな内容になるよう心掛けた。

2. 物体検出 (Object detection)

2.1 物体検出器 (object detector) の働き

基本動作: 物体検出器は, 画像 (1 枚の静止画をファイルにしたもの) を処理し, 処理結果 (検出個数, 検出物体の画像座標, 検出物体の種類) を出力する。セグメンテーションの場合は, 出力解像度 (画像のサイズ) に応じた各画素のクラス分類結果も出力される。画像を読み込ませる際には, 幅と高さを含む情報も与える必要がある。検出器によっては画像サイズや画像ファイル形式が指定されているものもあるのでその場合は, 予め画像ファイルを変換する前処理が必要となる。

出力結果の見かた: 検出座標は物体を囲む矩形 (bounding box) の座標を 4 つの数値で表すことが多い。また, 検出の信頼度が $[0, 1]$ の実数で出力される場合はそれが推定正答確率を表すように設定されている。

学習のしかた: 学習データは, 推論実行に必要なデータに正解データを加えたものであり, いわゆる教師あり学習を行わせる。ただ, 推論時には無い学習に関するパラメタの設定をしなければならない。検出器の構成 (中間層の層数, 特徴量次元のサイズ, 検出器独自のパラメタなど) もこの段階で設定する。通常は, 確率的降下法でモデルのパラメタを学習させることが多く, 検出器の重みパラメタの初期化方法 (平均, 分散, 値を指定など), 最適化方法 (SGD, Adam, 他), 学習率のスケジューリング, 学習打ち切り基準, ミニバッチサイズなどを指定する。学習時には交差検証を行わせて未学習データに対する誤差も計算させることができるので, その誤差を控えておくことにより推論時の予測精度を見積もることができる。

事前学習 (pre-training) と事後学習 (post-training, fine-tuning): 世界最高性能を出すほどの検出器の学習は, たいてい事前学習と事後学習の 2 段階で行われる。事前学習は主に公開データベースなどの大量データを用いて検出器の前半部分を学習させて適切な内部表現を獲得する

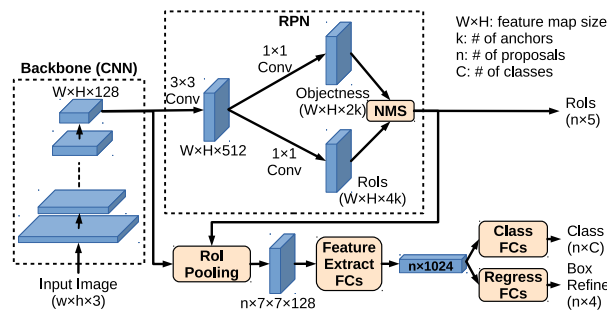


図1 Faster R-CNN の構造 .

ために行われる．検出器によっては事前学習済みの重み係数パラメタが公開されているものもある．事後学習は，最終目的に合致したデータを追加して，場合によっては検出器の最終層を追加して，所望の検出処理を実行できるように最終調整の意味合いで実施するものである．事前/事後学習のやり方は各検出器によって異なるため，説明書や論文等に記された方法を参考にして実行する必要がある．

物体検出器は次の2種類に分類できる [Jiao 19]：

Two-stage 検出器 高い位置決め精度と物体認識精度をもつ (Faster R-CNN など)．最初の CNN ベースの物体検出器 R-CNN とその改良手法はこの方式である．次の2つの stage で処理される：

- 第1 stage: 物体の bbox の候補を選出する．Faster R-CNN では Region Proposal Network (RPN) と呼ばれている．
- 第2 stage: RoI Pooling が各候補 box から特徴量を切り出し，後続の分類器と bbox 回帰器に送る．

One-stage 検出器 推論速度が高速である (YOLO, SSD など)．Region proposal をせずに直接，入力画像から bbox を検出して出力する．近年では，改良により検出精度も向上している．

2.2 Two-stage 検出器

§1 Faster R-CNN

代表的な two-stage 検出器であり，物体検出に対して当時の最高性能を達成しながら，処理速度の点でも大きく改善したモデルである [Ren 15]．PASCAL VOC 2007 に対して mAP=69.9% の精度を 5fps で達成した．

Faster R-CNN の構造: 図1に示すように，backbone，Region Proposal Network (RPN)，RoI Pooling，特徴抽出ネットワーク，分類器，回帰器で構成される．

Backbone は CNN で構成され， $W \times H \times 128$ のサイズの特徴量を出力する．提案当時は VGG-16 [Simonyan 15] などが用いられた．

RPN は backbone が出力した特徴量をさらに CNN で処理してカテゴリに依存せずに物体を検出してその候補 (proposal と呼ぶ) の RoI (bbox 座標) と物体らしさ

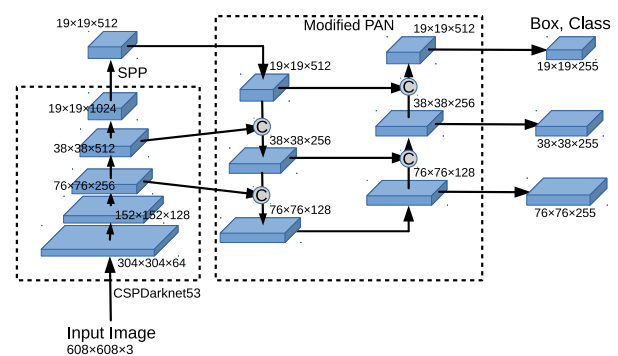


図2 YOLOv4 の構造 .

(objectness) を n 個ずつ出力する．RPN 内部では特徴量を CNN で変換し，特徴量の各空間座標点毎に k 個の bbox(4 個) と物体らしさ (2 個) の数値を同時に出力している (計 $WH(4+2)k$ 個)．これらの bbox は non-maximum suppression (NMS) によって，物体らしさが閾値を超えているものだけを残し，また，1 個の物体を複数回検出したものも 1 個だけに集約することにより，RPN が出力する proposal を n 個以下に抑える．ここで，各座標点に割り当てられた k 個の初期 bbox を anchor box と呼んでいる．

RoI Pooling は，backbone からの特徴量と RPN から n 個の proposals を受け取り，bbox 内部の特徴量を切り出して物体の大きさにかかわらずに空間サイズを 7×7 の固定サイズに変換する．

固定サイズの特徴量は，共通の特徴抽出を行う FC 層により $n \times 1024$ のデータに変換され，分類器と回帰器に送られて別々の FC 層で処理されて，クラス分類確率と bbox 調整量がそれぞれ出力される．(bbox 調整量が全クラス分出力されるのかどうか確認要！)

Faster R-CNN の学習:

- (1) CNN による特徴量生成は，画像認識などの学習で構築された重みを初期値として fine tuning を行う．
- (2) RPN だけを End-to-End で学習．anchor と教師データに基づき，物体か背景か (物体らしさ， $2k$ 個) と，anchor からのずれ ($4k$ 個) を学習．
- (3) RPN 固定で全体を学習する．RPN の RoI 出力は全体の出力及び RoI pooling へ送る．特徴抽出器で分類と回帰の共通部分の FC 計算を行う．分類器はクロスエントロピー，softmax などで計算し，回帰器が行う線形回帰は L1 ノルムで学習し，bbox の位置の補正を行う．

§2 TFANet

§3 Few-Shot Object Detection

2.3 One-stage 検出器

§1 YOLOv4

代表的な one-stage 検出器であり，高速，高性能で軽量であることが特長である [Bochkovskiy 20]．YOLOv3 までの特徴を継承しつつ，採用可能な backbone，neck，

head から最適な組み合わせを選択し、また、適用可能な各種技術を取り込んで構成された (図 2)。

YOLOv4 の構造: 入力画像は、backbone(画像の特徴を抽出する処理 (縦横サイズを小さくしながら特徴量を含むチャンネルを増やしていくこと)), neck(FPN や PAN のように、特徴量ピラミッドを昇り降り横断する処理), head(「backbone + neck」が出力する特徴量から、クラス信頼度と bounding box オフセットを推論するネットワーク) の順に処理される。

backbone である CSPDarknet53 は、YOLOv3 で使われた Darknet53 に Cross Stage Partial Network (CSP) を導入したものである。Darknet53 は、residual 結合を持ち Convolution 2 層で構成されるブロックが多数積み重なった構造であり、名前の 53 は Convolution が 53 個あることから来ているとしている (ただし、53 番目の層は全結合層である)。CSP network は予測性能を向上させる工夫であり、注目するブロックへの特徴量入力を 2 つに分割し、一方はそのブロックで処理し、もう一方は処理をスキップしてそのまま送り、これら 2 つを連結 (concatenation) して出力することを行う。ただし、YOLOv4 においては、この「分割」処理の代わりに、分岐させた直後に stride=2 の Convolution で処理して両方のチャンネルサイズを 1/2 にしている。

CSPDarknet53 の Top の出力は SPP モジュール (kernel size=1, 5, 9, 13, stride=1 として、4 つ並列に max pooling を行い、これらを concat するもの) に送られる。比較的大きな $k \times k$ max-pooling が効果的に backbone 特徴量の受容野を増加させてから、neck の Top に渡される。

neck である Modified PAN は、Path Aggregation Network (PAN)(3 つのピラミッドを行き来して高解像度情報を特徴マップに効果的に伝えるモデル) における bottom-up path の加算計算を concatenation に変更したモデルが用いられている。

ここではさらに、Modified Spatial Attention Module (SAM)(Convolution の出力を 2 つに分岐し、一方に Convolution + sigmoid 処理を行い、元信号に掛け算する) の演算を行ったものが PAN の 3 個の出力として head に渡される。

3 つの head は、それぞれ独立に convolution 計算を 2 回行い、最終的なクラス信頼度と bounding box オフセットを出力する。

YOLOv4 の学習: 一般的な確率的降下法を用いた学習に加えて、いくつか効果的な学習の工夫を導入している:
CutMix 学習画像の一部を切り取ったパッチを別の学習画像の一部に貼り付けて、正解ラベルもパッチの面積に比例させてミックスしたものを生成してそれで学習する。

Mosaic データ拡張 CutMix を、4 つの画像を用いるように拡張したもの。

DropBlock 特徴マップに対して、無作為に選出した矩

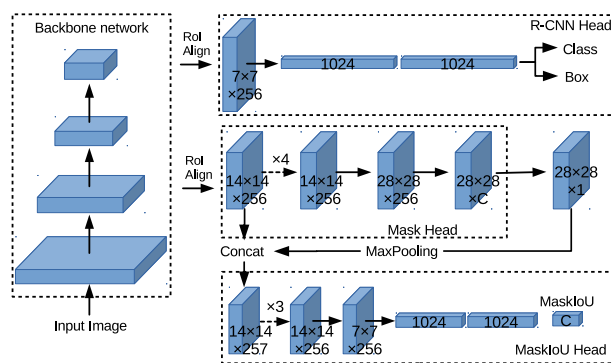


図 3 Mask Scoring R-CNN の構造。

形範囲 (block) にマスクを掛けたもので学習する。

Class label smoothing クロスエントロピー損失で学習する際に、logit が発散しないように正解ラベルを $q'(k|x) = (1 - \epsilon)\delta_{k,y} + \epsilon/K$ とした損失を用いて学習する。

Complete IoU (CIoU) 損失 二つの bounding box が離れていても、近づきすぎても適切な損失関数になるように、IoU 損失を改良した損失関数。

CmBN batch normalization を (ミニバッチ複数回につき 1 回の重み更新を行う場合に) ミニバッチをまたいで「平均、分散」を蓄積して、重みの更新のタイミングで、bias と scale を更新するようにしたものである。

Self-adversarial-training (SAT) まず、物体を写した元画像を改変し、物体が存在しない画像であると騙す画像を生成する。次に、この改変された画像中の物体を検出するように通常の学習を行う。

Cosine annealing scheduler コサイン関数の半周期の形状を利用して、学習率を少しずつ減少させる。これを周期を増やしながら複数回繰り返す。

§ 2 EfficientDet

3. インスタンスセグメンテーション (Instance segmentation)

3.1 Mask Scoring R-CNN (MS R-CNN)

マスク品質 (インスタンスマスクと正解マスクとの IoU として定量化されるもの) を分類スコアと明示的に関連付けたモデルである [Huang 19]。MS R-CNN は、予測マスクの品質を学習するためのブロック (MaskIoU Head) を、Mask R-CNN[He 17] に導入したモデルになっている (図 3)。MaskIoU Head はインスタンスの特徴量と対応する予測マスクを一緒に取り込み、それを元に Mask IoU を回帰推定する。そして、推論時に予測 MaskIoU を分類スコアに掛け算して補正する。

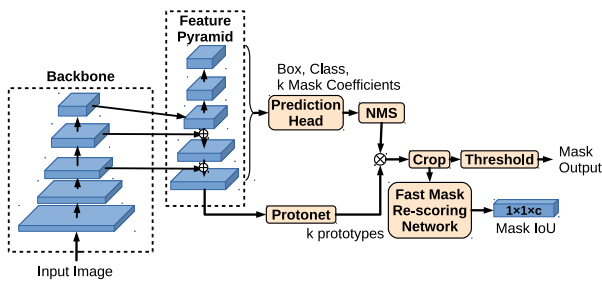


図4 YOLACT++の構造。

§1 MS R-CNN の学習

学習サンプルとして RPN proposals を使う。proposal box と正解 box との IoU が 0.5 以上の学習サンプルが必要となる。これは Mask R-CNN の Mask head の学習サンプルの場合と同じである。各学習サンプルに対する回帰目標を生成するために、まず目標クラスの予測マスクを取得し、予測マスクを閾値=0.5 で 2 値化する。そして、2 値化マスクと正解との MaskIoU を使う。MaskIoU を回帰するのは L2 損失を使い、損失重みは 1 にする。ネットワーク全体は end-to-end で学習する。

§2 MS R-CNN の推論処理

MaskIoU Head は分類スコア (R-CNN head の出力) の調整に使う。推論の手順は次のようになる：

- (1) R-CNN head が N 個の bounding box を出力する。
- (2) N 個の bounding box のうち、SoftNMS[Bodla 17] で上位 k 個のボックスを選択する。
- (3) 上位 k 個のボックスを Mask Head に入力し、 k 個のマルチクラスマスクを生成する (ここまでは標準的 Mask R-CNN の手順)。
- (4) これら k 個のマスクを目標として MaskIoU Head に入力し、予測 MaskIoU を出力する。
- (5) 予測 MaskIoU を、分類スコアに掛け算し、上位 k 個の修正された分類スコアを得る。

3.2 YOLACT++

実時間 (>30fps) で動作するインスタンスセグメンテーションのモデルであり、MS COCO に対して当時の最高性能に匹敵する性能 (34.1 mAP at 33.5fps) を達成した [Bolya 20]。fully-convolution モデルであり、deformable convolution を YOLACT[Bolya 19] の backbone に導入する等の改良をしている。物体検出モデルの RetinaNet[Lin 17] をもとに、インスタンスセグメンテーション向けに改良したものである (図 4)。

§1 YOLACT++ の学習

easy negative が多くて学習が困難になる問題は、OHem 法^{*1}を用いて negative:positive=3:1 にして学習することに対応する。Class 信頼度は分類損失 (クロスエントロ

^{*1} 入力画像に対する全 RoI をミニバッチと考えて、損失の値でソートして識別が難しい negative を選択して学習させる方法。

ピー)、bbox は L1 損失、mask (「mask 係数 × Prototype」で得られるもの) は pixel-wise binary cross entropy でそれぞれ学習する。Re-Scoring Net は Mask IoU(係数) を回帰する学習を行う。

Semantic Segmentation Loss は、学習時のみ接続されるネットワークの学習であり、この学習の実施により mAP が 0.4 ポイント向上する。P3 特徴量出力に 1×1 convolution 1 層で処理して c チャネルの出力をさせて最後にシグモイド関数をかける。これが正解 mask になるように学習する。

§2 YOLACT++ の推論処理

Prediction Head が各 anchor の Class 信頼度、bbox、 k 個の mask 係数を出力し、Protonet が k 個の Prototype (mask) を出力する。そして、Prediction Head 出力を NMS 処理して選ばれた結果に対して、mask 係数と Prototype を積和した結果 (全画面の mask) を、予測 bbox の外側を 0 で埋めたものを 2 値化して最終的な予測 mask 出力を得る。

並行して、2 値化する前の mask を Re-Scoring Net に入力して、mask IoU 出力を得る。分類スコアは、この mask IoU を掛けて補正される。

4. パノプティックセグメンテーション (Panoptic(?) segmentation)

5. む す び

謝 辞
謝辞について
a

◇ 参 考 文 献 ◇

- [Bochkovskiy 20] Bochkovskiy, A., Wang, C.-Y., and Liao, H.-Y. M.: YOLOv4: Optimal Speed and Accuracy of Object Detection, preprint arXiv:2004.10934 (2020)
- [Bodla 17] Bodla, N., Singh, B., Chellappa, R., and Davis, L. S.: Soft-NMS Improving Object Detection with One Line of Code, in *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 5562–5570 (2017)
- [Bolya 19] Bolya, D., Zhou, C., Xiao, F., and Lee, Y.: YOLACT: Real-time Instance Segmentation, in *International Conference on Computer Vision (ICCV 2019)*, Seoul, pp. 9156–9165 (2019)
- [Bolya 20] Bolya, D., Zhou, C., Xiao, F., and Lee, Y.: YOLACT++: Better Real-time Instance Segmentation, *IEEE Transactions on Pattern Analysis and Machine Intelligence (Early Access)*, No. DOI: 10.1109/TPAMI.2020.3014297 (2020)
- [He 17] He, K., Gkioxari, G., Dollr, P., and Girshick, R.: Mask R-CNN, in *2017 IEEE International Conference on Computer Vision (ICCV)*, Venice, pp. 2980–2988 (2017)
- [Huang 19] Huang, Z., Huang, L., Gong, Y., Huang, C., and Wang, X.: Mask Scoring R-CNN, in *Proc. of IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6409–6418 (2019)
- [Jiao 19] Jiao, L., Zhang, F., Liu, F., Yang, S., Li, L., Feng, Z., and Qu, R.: A Survey of Deep Learning-Based Object Detection, *IEEE Access*, Vol. 7, pp. 128837–128868 (2019)
- [Lin 17] Lin, T.-Y., Goyal, P., Girshick, R., He, K., and Dollr, P.: Focal Loss for Dense Object Detection, in *2017 IEEE International*

Conference on Computer Vision (ICCV), Venice (2017)

[Ren 15] Ren, S., He, K., Girshick, R., and Sun, J.: Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks, in *Advances in Neural Information Processing Systems (NIPS)*, Vol. 28 (2015)

[Simonyan 15] Simonyan, K. and Zisserman, A.: Very Deep Convolutional Networks for Large-Scale Image Recognition, in *International Conference on Learning Representations (ICLR)* (2015)

〔担当委員：×× 〕

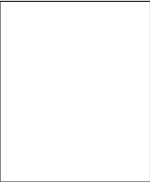
19YY 年 MM 月 DD 日 受理

◇ 付 録 ◇

A. 付録のタイトル 1

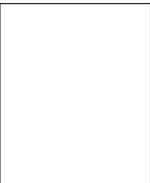
付録の本文 1

—— 著 者 紹 介 ——



金子 純也(正会員)

著者 1 の略歴



山田 貢己(正会員)

1989 年東京大学大学院物理学専攻修了・理学博士。同年株式会社東芝入社。ニューラルネットワークの研究開発，セキュリティ技術，画像認識技術，テレビの高画質化技術，車載画像認識プロセッサ等の開発業務に従事。2020 年ジャパニクス株式会社に入社。現在，Morning Project Samurai 株式会社において AI 開発業務に従事。