# グラフ畳み込みを用いたカラー画像と距離画像による 高速な3次元物体検出

○高橋 正裕 (中央大学), Alessandro Moro (RITECS), 梅田 和昇 (中央大学)

## Fast 3D Object Detection with Color and Range Images Using Graph Convolutional Network OMasahiro TAKAHASHI (Chuo University), Alessandro Moro (RITECS) and Kazunori UMEDA (Chuo University)

Abstract: In this paper, a lightweight 3D object detection model using color and depth images is proposed. In recent years, several studies have been done on the application of deep learning to 3D object detection. However, many of them are computationally expensive and difficult to run in real time because they deal with dense point clouds. In the proposed model, after feature extraction from the color image, a sparse point cloud is created from the range image to achieve fast object detection. Graph convolution for point clouds and feature extraction with depth information are also used. As a result, the proposed model achieved 56.4 fps when using ResNet34.

## 1. 研究背景

物体を検出・認識することは、防犯やマーケティン グ等、様々な分野に応用が可能である.そのため、こ れらを高速かつ正確に行うことができれば、カメラを 用いたシステムの能力全体を向上させることができる ので、重要なタスクであると考えられる.また、ステ レオカメラを始めとした距離画像も計測可能なセンサ が安価に手に入れることができるようになってきてい る.

近年、この分野に対しては、深層学習を3次元物体 検出に応用する研究が多く行われている。その中でも PointNet[1]や PointNet++[2]は点群からの特徴抽出を可 能とし、これらを物体検出や Semantic Segmentation タ スクに応用した VoteNet[3]や YOLO3D[4]といったモデ ルが提案されている。しかし、これらの研究の多くは 密な点群を扱うために計算コストが高く、リアルタイ ムでの運用が困難であり、求められる GPUの能力も高 い.また、LiDAR 等によって得られた点群を用いるこ とを前提にしているなど、大規模な環境を想定したも のが多く、手軽に導入できるとは言えない。

そこで本研究では、RGB-Dカメラから得られたカラ ー画像と距離画像を用いた軽量な 3 次元物体検出モデ ルを提案する.具体的には、カラー画像から特徴抽出 を行ったのち、距離画像から疎な点群を作成すること で、高速な物体検出を実現する.また、作成された点 群に対してさらに GCN(Graph Convolutional Network)[5] を用い、奥行き情報を考慮した特徴抽出を行う.

本稿では、3次元物体検出用データセットである SUN RGB-Dデータセット[6]を用いて提案モデルの物体検出

精度を評価する.また,処理速度についても Backbone Network 別に比較し,評価を行う.結果として,提案モ デルの性能は VoteNet 等の物体検出手法に及ばないも のの,非常に高速かつ軽量なモデルにより 3 次元物体 検出を実現することができた.

#### 2. 提案手法

#### 2.1 ネットワーク構造

本研究で提案するネットワーク構造を Fig. 1 に示す. ネットワークは大きく分けて, 色特徴抽出部分 (Backbone Network), グラフ作成部分, 3 次元特徴抽出 部分の 3 つに分けられる.

色特徴抽出部分:既存の特徴抽出モデルである VGG16[7]やResNet[8]を用い,カラー画像から特徴抽出 を行う.これにより得られた高次元の特徴をそのまま 以降のモデルで用いることで,クラス分類や物体検出 に必要なカラー画像からの特徴を、3次元物体検出に活 用する.また、この部分は一般的なカラー画像から得 られた特徴を用いることが目的であるため、物体検出 器[9-12]と同様の学習済みモデルを用いることが可能 である.

グラフ作成部分:まず,カラー画像から得られた高 次元の特徴を属性値として付与した点群を, KNN(K-Nearest Neighbor)によって16近傍に対してエッ ジを持つグラフ構造に変換する.この時,特徴マップ に合わせてスパースな点群を生成し,そこからランダ ムサンプリングを行う.これにより,以降のモデルに 入力する点の数を固定するとともに,ランダムサンプ リングを行うことで Data Augumentation と同様の効果



**Fig.1** 提案モデルの構造

を得ることができると考えられる.

3 次元特徴抽出部分: グラフ構造となった点群を, 6 層の GCN(Graph Convolutional Network)により、近傍点 の位置を考慮した、特徴抽出やバウンディングボック スの推定を行う. 今回のモデルは点群の座標を扱うた め、入力に負の値が存在する. そのため、活性化関数 には、LeakyReLU[14]と同様に負の値を扱うことができ る TanhExp[15]を用いる. 出力形式は YOLO(You Only Look Once)や SSD(Single Shot Detection)に共通する手法 である Unified Detection を元に作成した. Unified Detection とは、バウンディングボックスの座標やクラ ス分類結果等をチャンネル毎に格納することで、クラ ス分類と領域特定を同時に出力することが可能となる 手法である、Fig.2のように、提案モデルでは、出力の 各3次元座標が、注目点からバウンディングボックス の中心点までのベクトル(dx, dy, dz), バウンディングボ ックスの大きさ(width, height, depth), バウンディングボ ックスの角度(phi, theta),出力ベクトルの信頼度,クラ ス分類結果をチャンネル別で保持するように出力する. よって、出力サイズは、バッチサイズ×3次元点の数× (9+クラス数)となる.これにより、モデルに大きな分岐 がなく、よりコンパクトなモデルとなるので、高速な 物体検出が可能となる.

以上のようなネットワーク構造により,提案モデル はカラー画像と距離画像から3次元のバウンディング ボックスを出力する.

## 2.2 学習と推論

提案モデルでは、Backbone Network に用いる VGGや ResNet は、ImageNet[16]によって事前学習を行ったもの を用いる.これにより、画像中で事前に特徴のある領 域がある程度提示された状態で3次元物体の推定を行 うことができる.

提案モデルの出力形式はUnified Detection であるため, 誤差関数も YOLO に類似したものを学習に用いる.この誤差関数は、バウンディングボックスの中心座標へ 向かうベクトルの MSE(Mean Squared Error), バウンデ ィングボックスの大きさの MSE, バウンディングボッ クス の回転角度 の MSE, ベクトルの 信頼度 の BCE(Binary Cross Entropy) 誤差の合計で表される.よっ て, 誤差関数は式(1)のようになる.

## Loss =

 $\lambda_{bb} \{MSE((dx, dy, dz)_{out}, (dx, dy, dz)_{tar}) + MSE((phi, theta)_{out}, (phi, theta)_{tar})$ (1)

 $+ BCE(cos\theta_{out}, cos\theta_{tar}) + BCE(cls_{out}, cls_{tar})$ 

+  $\lambda_{nobb}$  { BCE( $cos\theta_{out}, cos\theta_{tar}$ ) }

ここで、下付き文字 out はモデルの出力、tar は教師デ ータを表している.  $\cos\theta$  は、バウンディングボックス の中心点へ向かうベクトルと正解のベクトルの角度  $\theta$ から得られ、式(2)によって与えられる.  $\cos\theta$  =

 $dx_{out}dx_{tar} + dy_{out}dy_{tar} + dz_{out}dz_{tar}$ 

 $+ dz_{out} dz_{tar}$  (2)

 $\sqrt{dx_{out}^2 + dy_{out}^2 + dz_{out}^2} \sqrt{dx_{tar}^2 + dy_{tar}^2 + dz_{tar}^2}$ また、 $\lambda_{bb}$ と $\lambda_{nobb}$ は、それぞれ物体が存在する、もしく はしないときに計算される誤差に対する係数で、今回 は $\lambda_{bb}=1$ 、 $\lambda_{nobb}=10$ を用いる.ここで $\lambda_{nobb}$ を大きめに設 定する理由としては、ここを $\lambda_{bb}$ と同じ値を設定すると 誤検出が多く発生してしまうためである.cls はクラス 分類結果であり、one-hot ベクトルで出力されるため、 BCE によって誤差計算を行う.

得られた候補から最適なバウンディングボックスを 選択する代表的な手法としては、R-CNNに用いられて いる NMS(Non Maximum Suppression) がある. これは IoU (Intersection over Union)と呼ばれる領域の重なり度 合いを表すスコアをもとに、同じ物体に対して推定さ れたバウンディングボックスを消去する方法である. YOLO3D 等の 3 次元情報を扱う物体検出器は、センサ 正面方向と垂直方向の 2 方向から 2 次元に対する IoU を計算し、NMSによるバウンディングボックスの選択 を行っている. しかし、この方法では同じ 3 次元のバ ウンディングボックスに対して 2 回 IoU を計算してい ることとなり、GPU による並列計算を行う上で非効率



### Fig.2 提案モデルの出力形式



Fig.4 成功例の Ground Truth と提案モデルの出力結果

的である. そこで,提案モデルでは,Fig.3のような体 積の重なり度合いを表す 3DloUを定義し,これをもと に NMS を実行する.これにより,IoU計算は GPU上 で1回しか行わず,処理速度の改善が期待できる.NMS における IoU のしきい値は,YOLO 等においては 0.5 が用いられていたが,提案モデルでは最も良いスコア であった0.6を用いた.

#### 3. 物体検出実験

提案モデルの有効性を確認するため、SUNRGB-Dデ ータセットを用いて精度を検証した.また、リアルタ イム性についても検証するため、Backbone Network を 変えたときの処理時間の比較を行った.検証に用いた マシンのスペックは、CPUが Intel Core i7 8770K、GPU が RTX2080 であった.

#### 3.1 物体検出精度検証

物体検出精度の検証として, SUN RGB-Dデータセッ トのうち, bed, table, sofa, chair, toilet, desk, dresser, night stand, bookshelf, bathtubの計10クラスを100エポック分 学習させた. バッチサイズは3,入力画像サイズは224 ×224[pixel],学習係数は0.0001に設定し,最適化手法 には Adam(Adaptive moment estimation)を用いた.

物体検出の Ground Truth と提案モデルによる物体検 出結果の成功例を Fig. 4 に示す.このシーンでは, bed と night stand が正解として与えられているが,どちらも かなり正確に検出できているといえる.特に bed に関し ては角度や物体のサイズも正確に検出できており,大 きい物体に関しても検出が可能であることがわかる. night stand に関しては物体の中心推定に誤差が生じて いるが,サイズを正しく推定できていることがわかる.



Fig. 3 3D IoU



Fig.5 失敗例の Ground Truth と提案モデルの出力結果

続いて,失敗例を Fig.5 に示す. このシーンでは,それ ぞれの物体は検出できているものの,各物体の中間に 誤検出が発生していることがわかる.これら2つの例 を比較すると,点群の質に差があることがわかる.成 功例の点群は各物体に対して得られている点群にあま りノイズがなく,点群がまとまっており,その結果近 傍点を考慮した GCNによる特徴抽出がうまくいったと 言える.しかし失敗例においては物体毎に得られてい る点群にまとまりがなく,物体が存在しない部分に関 しても点群のばらつきが多いため,近傍点探索による グラフの作成がうまくいかなかったと考えられる.

Backbone Network を変えたときの 3D loUによる結果 の平均値を Table 1 に示す.この結果によると、ResNet34 を Backbone として用いたときの結果が最も良いことが わかる. ResNet34 より深い層を持つモデルでうまくい かなかった理由としては、今回用いた入力画像のサイ ズが小さく,出力サイズが足りなかったことが考えら れる. 3D loUのスコアとしては, 05 を超えているため, 物体をある程度正確に検出ができているといえる. し かし同じ3次元物体検出用モデルである VoteNet のス コアの0.83と比較すると、精度面では及ばなかったこ とがわかる.この理由としては、VoteNetでは密な点群 を用いて検出しているのに対し、提案モデルではスパ ースな点群を用いた検出であるため、その分3次元推 定が困難になっていると考えられる.しかし, Fig.5の ように誤検出が多いものの物体位置の特定はできてい るため、学習エポック数の増加や学習係数の調整次第 では解決可能であると考えられる.

### 3.2 処理速度検証

続いて,提案モデルの処理速度の検証を行った.検

Table 1 モデルの物体検出精度の結果

Backbonename	Mean 3D IoU
VGG16	0.586
ResNet18	0.603
ResNet34	0.627
ResNet50	0.606
ResNet101	0.610
ResNet152	0.581

証方法としては、合計 200 フレームを推論させ、その 結果を用いて統計量を算出した.

Table 2 は, Backbone Network を変えたときの処理速 度の結果である.この結果によると,最も軽量な VGG16 を用いたモデルは,処理速度の中央値が 85[fps]を超え ており,物体検出精度検証において一番良い結果であ った ResNet34 を用いたモデルも、564[fps]を出せてい るため,十分にリアルタイム性を保持できていると言 える.ここで,中央値をベンチマークとした理由とし ては,中央値や標準偏差からわかるように,最小値が 外れ値を取っているためである.

3 次元物体検出を行うことが可能な他のモデルと比較しても、処理速度に関しては、提案モデルがシンプルなネットワークで構成されていることもあり、高速な検出を実現することができた.最も高速な YOLO3Dも、TITAN Xを用いて 40[fps]であるため、処理速度に関してはこれらを大きく上回ったといえる.処理速度を維持したまま精度向上を行う方法としては、シンプルな数層の GCNを追加することでバウンディングボックスの候補数を絞ることが考えられる.また、これによりバウンディングボックス中心点の推定精度の向上も見込める.

#### 4. 結言

本研究では、RGB-Dから3次元のバウンディングボッ クスを出力可能で軽量なモデルを作成した.提案モデ ルでは、カラー画像から得た特徴を点群とともにGCN へ入力することで、バウンディングボックスの奥行き 方向の位置と長さを取得した.

今後の展望としては,更なる深層化やパラメータ調 整を行うことで,精度向上を図る.

#### 参考文献

[1] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3d classification and segmentation," in Proc. of

Table 2 処理速度検証結果

Backbone	Speed [fps]				Standard
name	max	min	mean	median	deviation
VGG16	91.6	4.6	79.4	86.7	13.6
ResNet18	69.5	4.6	60.1	64.7	9.7
ResNet34	60.7	4.6	53.0	56.4	7.6
ResNet50	21.9	4.0	20.6	21.1	1.5
ResNet101	19.0	2.7	17.9	18.1	1.3
ResNet152	16.8	4.0	16.1	16.3	1.0

the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 652-660, 2017.

- [2] C. R. Qi, H. Su, and L. J. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," Conference on Neural Information Processing Systems (NeurIPS), pp. 5099-5108, 2017.
- [3] C. R. Qi, O. Litany, K.He, and L. J. Guibas, "Deep Hough Voting for 3D Object Detection in Point Clouds," in arXiv preprint arXiv:1904.09664v2, 2019.
- [4] W. Ali, S. Abdelkarim, M. Zahran, M. Zidan, and A. E. Sallab, "YOLO3D: End to end real time 3D Oriented Object Bounding Box Detection from LiDAR Point Cloud," arXiv preprint arXiv: 1808.02350, 2018.
- [5] T. N. Kipf and M. Welling, "Semi-Supervised Classification with Graph Convolutional Networks," in 5<sup>th</sup> International conference on Learning Representations (ICLR), 2016.
- [6] S. Song, S. Lichtenberg, and J. Xiao, "Sun rgb-d: A rgb-d scene understanding benchmark suite," in Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 567-576, 2015.
- [7] K. Simonyan and A. Zisserman, "Very Deep Convolutional Networks for Large Scale Image Recognition," in arXiv preprint arXiv:1409.1556, 2014.
- [8] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016.
- [9] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 580-587, 2014.
- [10] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R CNN: Towards real time object detection with region proposal networks," in Proc. of Conference on Neural Information Processing Systems (NeurIPS), 2015.
- [11] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You Only Look Once: Unified, Real

Time Object Detection," in Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 779 788, 2016.

- [12] J. Redmon and A. Farhadi, "YOLOv3: An Incremental Improvement," in arXiv preprint arXiv:1804.02767, 2018.
- [13] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C. Y. Fu, and A. C. Berg, "SSD: Single Shot MultiBox Detector," in arXiv preprint arXiv:1512.02325, 2015.
- [14] A. L. Maas, A. Y. Hannun, and A. Y. Ng, "Rectifier nonlinearities improve neural network acoustic models," in International Conference on Machine Learning (ICML), p. 3, 2013.
- [15]X. Liu and X. Di, "TanhExp: A Smooth Activation Function with High Convergence Speed for Lightweight Neural Networks," in arXiv preprint arXiv:2003.09855v2, 2020.
- [16] J. Deng, W. Dong, R. Socher, L. J. Li, K. Li, and L. F. Fei, "ImageNet: A large-scale hierarchical image database," in Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 248-255, 2009.