走査型 LiDAR のための時刻情報を用いた超解像

林佑介1 田口賢佑1 森田渉吾1 今枝航1 藤吉弘宣2

概要: LiDAR は広い範囲にわたって高精度に距離や位置,形状を検出し,3次元での把握が可能である.その反面, カメラよりも解像度およびフレームレートが低いという課題がある.また,レーザー光をスキャンする走査型 LiDAR では,各 Depth 測定値の測定時刻が異なり,この測定時刻の違いに起因して移動物体の Depth がずれるという課題が ある.本研究では,高解像度で移動体に対してずれのない Depth を得るため,走査型 LiDAR で得られた Depth 測定値 を時間方向に拡張した3次元の Voxel 表現とし,これを入力データとしてニューラルネットワークを用いて推定する. 空間的には遠方や物体の境界付近,時間的には目的の時刻に遠い時刻や移動物体のように高解像度化が困難な Depth 測定値に対して,背景と移動体の切り分けや不確実性を考慮することで精度向上を図る.

キーワード:走査式 LiDAR,超解像,不確実性,CG データ

1. はじめに

突風による道路への落下物や集中豪雨による土砂の流 入の迅速な検知など,屋外への常時設置が必要となるイン フラ監視において、より正確に状況を把握するために LiDAR は有用なセンシングデバイスの一つとなる. レー ザー光を反射させ, 戻ってくるまでの時間を計測する LiDAR は、広い範囲にわたって高精度に距離や位置、形状 を検出し、3次元での把握が可能になる.その反面、カメ ラよりも解像度およびフレームレートが低いという課題が ある.また、レーザー光の照射方法は大きく分けて走査型 とフラッシュ方式の2つがあるが、フラッシュ方式は面で 光源を照射するため検出距離が走査型よりも短くなり、外 光で測定結果がばらついてしまうために屋外での使用にも 向いていない.一方, 走査型 LiDAR は受光面積当たりのレ ーザーの光量を大きくできるために検出可能距離を長くす ることができるのがメリットであるが、モーターを用いた 回転機構や MEMS 等によりレーザー光をスキャンする必 要があるため、各 Depth の測定時刻が異なる.この測定時 刻の違いに起因して、車やバイクなどの移動物体の Depth が実際の値とずれてしまうという課題がある.

本研究は、走査型 LiDAR においてレーザー光のスキャンによる時刻ずれを解消し、移動体に対してずれのない正確な Depth を得るための超解像技術として以下の3つの項目を提案する.

①高解像度で移動体に対してずれのない Depth を得るため、走査型 LiDAR で得られた Depth 測定値を時間方向に 3 次元に拡張した Voxel 表現で表し、これを入力データとして全層畳み込みネットワークで推定

②畳み込みが有効に特徴を把握できるようにするため, 画像から背景を抽出するマスクを作成することで時空間上 の Depth 測定値を高密度化 ③時空間方向の不確実性を考慮することで移動体の Depth 推定の精度向上

ただし、本研究で用いる LiDAR はカメラとの光軸が一 致したフュージョンセンサで、Depth 測定値と画像の視差 がない.また、本手法はこのフュージョンセンサを固定し て設置した場合に限定し、静的なシーンの Depth マップの 超解像を扱う.

2. 関連研究

2.1 超解像

超解像には1枚のフレームを用いるものと複数フレーム を用いるものの2種類がある.1枚のフレームの超解像と しては、ニューラルネットワークを超解像に適用し end-toend で学習したモデル[1]がある. また, [2]では ResNet[3]の 構造を用いて入力層と差分値のみを学習し、スキップ接続 により入力情報を保持して復元するという提案がされてい る. 複数のフレームを用いた超解像としては、現フレーム と前フレームを入力し、出力として各ピクセルがどの方向 に動いたかを推定したフローマップを使って前フレームの 特徴量マップを現フレーム用に変換する提案[4]がある.ま た、1 枚の画像を時関方向に拡張した Voxel 表現とし、コ ード化した露出と組み合わせることにより時空間上での超 解像を行う研究[5]がある. Voxel 表現を用いた別の研究と して、[6]ではイベントカメラを用いて疎なイベント情報を 時空間上に表し、リカレントニューラルネットワークを用 いて高速高解像なビデオを生成する.本研究では、これら の Voxel 表現を用いた研究と同様に, Depth 測定値を時空 間上に表すとともに背景をマスクにより抽出することで高 密度化することで超解像を行う.

2.2 不確実性

不確実性の推定は、ベイジアンニューラルネットワーク [7]をはじめとして、コンピュータビジョンを対象としたニ

¹ 京セラ(㈱) 先進技術研究所

Advanced Technology Research Institute, Kyocera Corp. 2 中部大学

Chubu University

ューラルネットワークの研究においても様々な報告がなさ れている.ベイジアンニューラルネットワークにおいては, 重みの分布から異なるネットワークを複数回サンプリング することにより,ターゲットとなる分布の平均と分散を算 出する.[8]では,サンプリングが変分推論に置き換えられ ている.別の方法としては,アンサンブル[9]やモンテカル ロドロップアウト[10]が用いられる.異なるアプローチと しては,推論の際に不確実性を算出する研究[11]がある.不 確実性を応用する研究として,焦点ずれによるレンズのぼ けがニューラルネットワークに与える影響を解析するため にモデルの知識不足か否かを示す認識論的不確実性

(Epistemic Uncertainty) を分析した研究[12]がある.本研究 では[11]の手法を用いて学習の際に Depth 推定するととも にデータの難しさを表す偶然的不確実性 (Aleatoric Uncertainty) を算出することで Depth の超解像の性能が向 上することを検証する.

2.3 CG データ

運転シーンをシミュレーションした CG データとして SYNTHIA[13]や VirtualKiTTI[14]がある. SYNTHIA では画 像にセマンティックセグメンテーションのラベルが付けら れている一方, VirtualKITTI では KITTI[15]の動画がシミュ レーションされ,画像に完全に一致した Depth が画素ごと に対応付けされている. CARLA[16]は自動運転のアルゴリ ズム開発,検証を支援するためのシミュレーターである. カメラや LiDAR のセンサ出力に対応しており,シミュレ ーター空間上の任意の位置にセンサを設置することが可能 である.しかしながら,これらのデータセットではフレー ムレートはいずれも遅く,1 フレーム分の時間内の各時刻 に対応する教師データとなる Depth 情報を得ることができ ない.そこで本論文では CARLA でも使用されている 3D シ ミュレーターの Unreal Engine4[17]と 3D CG データベース として AUTOCity[18]を用いて CG データを作成した.

2.4 ローリングシャッターおよび LiDAR ノイズの解消

従来の LiDAR に関する研究では, Depth 測定値はノイズ がないものとされている.しかしながら,前述した KITTI のような公開データセットにおいても,たとえエゴモーシ ョンを見込んでローリングシャッターの影響を補正しても, 透明物体や反射強度の高い物体においてノイズが生じるこ とが報告されている[21].3D モデルを用いて人間が目視で これらのノイズを除去することは可能だが,アノテーショ ンコストが非常に高い.カメラのローリングシャッターを 解消する研究に注目すると,[19]では平面運動を仮定して 画像の歪みを補正する方法が提案されている.また,列ご との読み出しタイミングと露出の長さを制御して3次元の 時空間のシーン情報をサンプリングすることで,高速撮影 や動きブレのない画像生成の提案[20]がなされている. LiDAR の Depth 測定値のノイズ除去に関連する研究とし て,LiDAR のローリングシャッターの効果とキャリブレー ションミスによる Depth 測定値のずれに対して, [21]では ステレオカメラの輝度情報の整合性(Photometric Consistency) およびステレオカメラと LiDAR の Depth 情報 の整合性(Depth Consistency)を用いてノイズ除去を行う提 案がされている. [22]では, Depth Completion のタスクに対 して, 実空間上のノイズを含んだ疎な Depth を CG 空間上 で再現するドメイン適応の手法が提案されている.本研究 では LiDAR の Depth 測定値に含まれるノイズがないよう, 教師データを CG 上で作成することによってノイズ除去さ れた高解像度な Depth 推定が可能であることを示す.

3. 提案手法

3.1 概要

我々の目的は、1フレーム $F = \{f_t\}$ for $t \in [0, T-1]$ の Depth 測定値からある時刻の Depth マップを推定するこ とである. 今回の検証では、推定する Depth マップの時刻 は1フレームのうちの中央値の時刻 $T/_2$ とする. また、推定 する Depth マップの解像度は元の1フレームで得られる解 像度よりも高い解像度とする. この超解像のタスクを達成 するために、1フレームの Depth 測定値を時空間の Voxel と して表現し、この Voxel を入力として全層畳み込みネット ワークで推定を行う. ネットワークの学習は教師あり学習 の手法を使い、学習データとして大量の CG で作成した距 離データを用いる.

3.2 入力データの表現方法及びマスクによる前処理

F(x,y,t)は $M \times N$ ピクセルと1フレーム分の時間 T に対応する時空間の体積を表現した Voxel である. 従来の LiDAR は、この Voxel を時間軸上に投影し、その結果 $M \times N$ の解像度を得る. 我々は、時空間分解能でL倍のゲインを達成したい、つまり、 $M \times N \times L$ の時空間体積を実現する超解像を行いたい.

全層畳み込みネットワークが Voxel を効率よく処理できる ようにするため,各時刻 t に対応した背景を抽出するマス クM(x,y,t)と1フレームの累積との要素積を算出し, Voxel の背景部分はこの要素積となるよう重ね合せる.これによ り背景の Depth を高密度化する.各時刻のマスクの作成は, 例えば LiDAR の1フレーム分の時間で得られる画像2枚 を用いてオプティカルフローを算出し,等速直線運動を仮 定して各時刻の背景を抽出する.本研究の実験では,CGデ ータで作成したマスクを用いた.

3.3 学習データ

我々のネットワークは1フレーム分の時間内の各時刻に 対応する Depth 測定値が教師データとして必要である.し かしながら,そのような Depth 測定値および画像を含む公 開データセットはない.そのため,CGで作成した画像およ び Depth を用いて学習することを提案する.

我々は、3D シミュレーターとして Unreal Engine4 を、3D CG データベースとして AUTOCity を用いた. AUTOCity は

物理ベースのレンダリングが可能で,カメラで撮影される 画像とLiDAR で得られるDepth 情報を作成でき,仮想的に 光軸を一致させることできる.そのため,我々が想定して いるカメラとLiDAR の光軸が一致したフュージョンセン サから得られる画像とDepth 測定値の視差がない情報を得 ることが可能である.画像およびDepthの解像度は240×48, フレームレートは300fpsと設定した.カメラとLiDAR は 固定した状態で,車やバイク,自転車,人が動く状況をシ ミュレーションし,200シーンを作成した.各シーンは5秒 で,画像とDepth はそれぞれ 300,000 枚から構成される.

学習データのフレームレートおよび空間解像度は,カメ ラは6fpsで240×48, LiDARは3fpsで120×24と設定する. CGデータの時間分解は3msecであるため,LiDARの1フ レームは100枚で1セットとなり,Voxelとしては12×2×100 で表される.学習データの総数は3,000セットとなる.

3.4 ネットワーク構成および損失関数

ネットワーク構成はUNet[23]の構造を模した全層畳み込 みネットワークを用いる.ネットワークの概要を図1に示 す.



ネットワークはヘッド層 (Head layer), エンコーダー層 (Enc. layer), 残差ブロック (Res. layer), デコーダー層 (Dec.

layer), 出力層 (Pred. layer) で構成し, 各エンコーダーと デコーダー層をスキップ接続[24]する. 実験ではエンコー ダー層は4層, 残差ブロックは2ブロック, ヘッド層の出 カチャネル数は32と設定する. [11]にならって, ネットワ ーク $f^{\hat{W}}$ は式 (1) の通り Depth 推定値 \hat{D} と不確実性 $\hat{\sigma}^2$ を出力する.

$$\begin{bmatrix} \hat{D}, \hat{\sigma}^2 \end{bmatrix} = f^{\widehat{W}}(F)$$
(1)
損失関数は式 (2) を用いてネットワークを学習する.
$$L(\theta) = \frac{1}{N} \sum_i \frac{1}{2} \exp(\log \hat{\sigma}_i^2) \|y_i - \hat{y}_i\| + \frac{1}{2} \log \hat{\sigma}_i^2$$
(2)

式(2) でモデリングした不確実性は,推定が難しいデー タの程度を示す Aleatoric Uncertainty[12]を表しており,こ れを用いて学習することにより,空間的には遠方や物体の 境界付近,時間的には目的の時刻に遠い時刻や移動物体の ように Depth 推定が困難な Depth 測定値に対して頑健な出 力となることを狙う.

3.5 学習条件

学習データの総数は 3,000 セットのうち 2,700 セットを 学習用,300 セットをテスト用に用いる.ネットワークは PyTorch[25]を用いて実装した.最適化アルゴリズムは ADAM[26]を用い,学習率は 0.00005,バッチサイズは 1, エポック数は 16 と設定した.

4. 評価

定量的,定性的に我々の超解像の提案手法(以下,Voxel +マスク+不確実性)を評価する.また,1フレーム分の 累積(以下,累積),時刻情報を加えたVoxel(以下,Voxel), Voxelにマスクを加えた場合(以下,Voxel+マスク),累積 に不確実性を考慮した場合(累積+不確実性),Voxelに不 確実性を考慮した場合(Voxel+不確実性)について比較を 行う.評価指標は平均平方2乗誤差(RMSE)を用いる.

4.1 結果

定量的評価を図2,3,4に示す.



図 2 画面全体の RMSE

図2は画面全体のRMSEである.累積とVoxelを比較す るとほぼ同様の精度であることが分かる.これらに対して Voxel+マスク,Voxel+マスク+不確実性はいずれも精度 が改善しており,提案手法がもっともいい結果であること が分かる.また,累積, Voxel, Voxel+マスクの3つの手法 に対して,それぞれ不確実性を加えるといずれも性能が改 善している.







図 4 背景の RMSE

図 3,4 はそれぞれ移動体,背景の RMSE である.累積 と Voxel と比較すると、Voxel は累積よりも移動体がよくな っているが,背景は悪化している.これは Voxel では時刻 情報を用いることにより移動体の軌跡を類推することがで きる一方,背景は累積よりも Depth 測定値が疎であるため であると考えられる.そこで Voxel にマスクを加えて背景 を高密度化した Voxel+マスクに注目すると、移動物体と 背景のいずれも累積と Voxel よりも性能改善していること が分かる.さらに不確実性を加えた Voxel+マスク+不確 実性は、Voxel+マスクよりも移動体の RMSE を改善する ことができている.損失関数に不確実性を加えたことで移 動体に注目して学習していると考えられる.

図5,6に定性評価を示す.

図5はGTのDepthおよび累積, Voxel, Voxel+マスク, Voxel+マスク+不確実性を入力とした場合のそれぞれの Depth 推定結果である.赤は近方,緑に近づくにつれ遠方 を表している.画面中央の遠方のポールと画面右側の自転 車に注目する.GTに対して,累積やVoxelは自転車のエッ ジが不明瞭であることが分かる.また,Voxelのみでは累積 よりもポールのエッジが不明瞭で,図4に示す背景の RMSEがVoxelのみでは精度が悪い結果と同じ傾向である ことが分かる. Voxel+マスクおよび Voxel+マスク+不確 実性に注目すると, Voxel よりもポールのような遠方のエ ッジが若干改善している. また, 自転車のような複雑な形 状の Depth が再現されていることが分かる.



図 5 Depth 推定結果



図 6 不確実性

図6は不確実性を可視化したもので,移動物体が分かる よう時刻 1/300sec と 50/300sec に注目する.参考に不確実 性と併せて推定結果を示す.推定結果については図5と同 様に赤は近方,緑に近づくにつれ遠方を表している.不確 実性については赤が小さく,黄色は大きいことを表してい る.遠方を除く背景および移動体の中央付近では不確実性 が小さいのに対して,遠方の背景および移動体のエッジ付 近で不確実性が大きくなっている.この結果から,本研究 での狙い通り Depth 推定がそもそも難しい遠方や移動体の エッジ付近に対しては不確実性が大きすぎる場合には損失 関数が小さくなり無視できるようにすることにより,Depth 推定の超解像に関する精度について最適化が進む効果が出 ていると考えられる.また,図3および図6の結果から, 学習時に不確実性を考慮することにより移動体,特にその 中央付近の精度が向上したことが分かる.

5. まとめ

走査型 LiDAR においてレーザー光のスキャンによる時 刻ずれを解消し,移動体に対してずれのない正確な Depth を得るための超解像技術について提案を行った.そのため の手法として,走査型 LiDAR で得られた Depth 測定値を時 間方向に3次元に拡張した Voxel 表現を用いて入力データ として全層畳み込みネットワークで推定する.また,畳み 込みが有効に特徴を把握できるようにするため,画像から 背景を抽出するマスクを作成することで時空間上の Depth 測定値を高密度化する.さらに,ネットワークの出力とし て不確実性をモデリングすることで移動体の Depth 推定の 精度を向上する.自作の CG データセットを用いた実験に より,累積や Voxel を入力とした場合よりも移動体,背景 に対して提案手法が有効であることを示した.

参考文献

- C. Dong, C. C. Loy, K. He, and X. Tang, "Image super-resolution using deep convolutional networks," in IEEE Transactions on Pattern Analysis and Machine Intelligence, 38(2):295–307, 2016.
- [2] J. Kim, J. K. Lee, and K. M. Lee, "Accurate image superresolution using very deep convolutional networks," in IEEE International Conference on Computer Vision and Pattern Recognition, 1646–1654, 2016.
- [3] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in IEEE International Conference on Computer Vision and Pattern Recognition, 2016.
- [4] M. S. Sajjadi, R. Vemulapalli, and M. Brown, "Frame-recurrent video super-resolution," in IEEE International Conference on Computer Vision and Pattern Recognition, 6626–6634, 2018.
- [5] Y. Hitomi, J. Gu. Gupta, T., Mitsunaga, and S. Nayar, "Video from a single coded exposure photograph using a learned over-complete dictionary," in IEEE International Conference on Computer Vision, 287–294, 2011.
- [6] H. Rebecq, R. Ranftl, V. Koltun, and D. Scaramuzza, "High speed and high dynamic range video with an event camera," IEEE Transactions on Pattern Analysis and Machine Intelligence, 2019.
- [7] D. J. C. MacKay, "A practical Bayesian framework for backprop networks," Neural Computation, 4:448–472, 1992.
- [8] A. Graves, "Practical variational inference for neural networks," in Advances in neural information processing systems, 2348–2356, 2011.
- [9] B. Lakshminarayanan, A. Pritzel, and C. Blundell, "Simple and scalable predictive uncertainty estimation using deep ensembles," in Advances in Neural Information Processing Systems, 6402– 6413, 2017.
- [10] Y. Gal and Z. Ghahramani, "Dropout as a Bayesian approximation: Representing model uncertainty in deep learning," in International Conference on Machine Learning, pages 1050–1059, 2016.
- [11] D. A. Nix and A. S. Weigend. "Estimating the mean and variance of the target probability distribution," in Proceedings of 1994 IEEE International Conference on Neural Networks, 1:55–60, 1994.
- [12] M. Carvalho, B. L. Saux, P. T. Peloux, A. Almansa, and F. Champagnat, "Deep depth from defocus: how can defocus blur

improve 3d estimation using dense neural networks?," in European Conference on Computer Vision, 307–323, 2018.

- [13] G. Ros, L. Sellart, J. Materzynska, D. V'azquez, and A. L'opez. "The SYNTHIA dataset: A large collection of synthetic images for semantic segmentation of urban scenes," in IEEE International Conference on Computer Vision and Pattern Recognition, 2016.
- [14] A. Gaidon, Q. Wang, Y. Cabon, and E. Vig, "Virtual worlds as proxy for multi-object tracking analysis," in IEEE International Conference on Computer Vision and Pattern Recognition, 2016.
- [15] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the KITTI vision benchmark suite," in IEEE International Conference on Computer Vision and Pattern Recognition, 3354– 3361, 2012.
- [16] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, "CARLA: An open urban driving simulator," in Conference on Robot Learning, 2017.
- [17] Epic Games. Unreal Engine 4. https://www.unrealengine.com.
- [18] VERTechs. AUTO City. https://www.vertechs.jp/products.
- [19] D. Bradley, B. Atcheson, I. Ihrke, and W. Heidrich, "Synchronization and rolling shutter compensation for consumer video camera arrays," in IEEE International Workshop on Projector-Camera Systems, 2009.
- [20] J. Gu, Y. Hitomi, T. Mitsunaga, and S. Nayar, "Coded rolling shutter photography: Flexible space-time sampling," in IEEE International Conference on Computational Photography, 1–8, 2010.
- [21] X. Cheng, Y. Zhong, Y. Dai, P. Ji, and H. Li, "Noise-aware unsupervised deep lidar-stereo fusion," in IEEE International Conference on Computer Vision and Pattern Recognition, 2019.
- [22] A. Lopez-Rodriguez, B. Busam, and K. Mikolajczyk, "Project to adapt: Domain adaptation for depth completion from noisy and sparse sensor data," in Proceedings of the Asian Conference on Computer Vision, 2020.
- [23] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in International Conference on Medical Image Computing and Computer-Assisted Intervention, 2015.
- [24] A. Z. Zhu, L. Yuan, K. Chaney, and K. Daniilidis, "EV-FlowNet: Selfsupervised optical flow estimation for event-based cameras," in Robotics: Science and Systems, 2018.
- [25] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in PyTorch," in Conference on Neural Information Processing Systems Workshops, 2017.
- [26] D. P. Kingma and J. L. Ba, "Adam: A method for stochastic optimization," in International Conference on Learning Representations, 2015.