物体メッシュモデルを用いた学習データ自動生成に基づく 透明物体の深度画像予測と家事支援ロボットへの応用

Depth Image Prediction for Transparent Objects and its Application to Home Assistant Robot Based on Automatic Training Data Generation Using Mesh Model

内海 佑斗	和田 健太郎	岡田 慧	稲葉 雅幸
Yuto Uchimi	Kentaro Wada	Kei Okada	Masayuki Inaba

東京大学大学院情報理工学系研究科

The University of Tokyo, Graduate School of Information Science and Technology

Home assistant robots deal with various objects in the home. However, transparent objects such as wine glasses are difficult for robots to estimate 3D position, therefore they have trouble in manipulating such objects. This paper proposes a method of semantic segmentation and depth image prediction using CNN, which enables robots to get depth of transparent objects and reconstruct them three dimensionally. It also proposes a method of collecting training data, by which a person no longer has to annotate objects in 3D space. By using our depth prediction method, robots can handle PET bottles as well as non-transparent objects.

1. はじめに

家事支援ロボットには家庭内にある多種多様な物体を扱うこ とが求められるが、その中でもペットボトルやワイングラスな ど透明物体の認識と操作は未だ解決されていない課題の一つで ある。これらは光を反射・吸収・透過するなどの性質を持つた めセンシングによる三次元的な位置推定が困難であり、家事支 援を遂行するための環境認識に際して大きな課題である.

近年,透明物体の三次元位置推定として,条件付き確 率場や多層畳み込みネットワークを利用した機械学習に よる深度画像予測の手法の有効性が明らかになってき た [Eigen 14] [Liu 15] [Wang 15].物体の特性を検知する手 法 [Eren 09] [Rantoson 10] に対し,機械学習による手法の多 くはシステムの入力としてカメラから取得した単一の色画像の みを用いることで,透明物体の認識のために特別な装置を要す ることなく深度予測を実現可能であるほか,画像中の各領域が どのような物体であるかの分類も可能である.

ところがロボットに搭載されるカメラセンサについては RGB-D の計測が可能なカメラが既に広く用いられているため,透明物体以外の深度取得にはセンサ値を採用すればよい.本研究 では,多層畳み込みネットワークを用いてラベル画像予測と深 度画像予測を同時に行い,センサから取得される深度画像のう ち,ネットワークにより透明だと判定された領域のみを予測深 度画像で置換した深度画像を出力して三次元位置推定を行うこ とを提案する.

しかしながら,教師あり学習を行うためには透明物体の 正しい深度情報を持つ画像を収集する必要がある.NYU Depth [Silberman 12] などのデータセットは、深度情報を取 得するための投射光が正常に帰ってこないという透明物体の性 質のため、これらの正しい深度情報を含んでいない。この解決 策として人間による三次元空間上でのアノテーションが挙げら れるが,物体数が増えるに従って人間の負担も増加するため, 効率的とは言えない.そこで本研究では,透明物体の三次元 メッシュモデルを利用してタスク環境を仮想的に再現し,人間 による三次元アノテーションの不要な効率的な学習データ生成 を行うことを提案する.

連絡先:内海佑斗,東京大学工学部情報システム工学研究室, 113-8656 東京都文京区本郷 7-3-1 工学部 2 号館 7 階 73B2, 03-5841-7416,uchimi@jsk.imi.i.u-tokyo.ac.jp



図 1: 学習データ生成 (上) と深度予測システム (下) の概観

本研究の提案手法の有用性を示すため、ロボットが家事支援 において透明物体を扱うタスクに提案システムを適用する実験 を行う.

2. 物体モデル配置による学習画像生成

機械学習による深度予測を行うために、透明物体の教師ラベル画像と教師深度画像を取得する必要がある.本研究では、ま ず実際の様々なタスク環境について、ロボットが各シーンにつ き複数視点から撮影することで、透明物体の置かれていない画 像データを収集する.次に、環境の点群上に透明物体の CAD モデルを仮想的に自動配置することで、教師画像を生成する. また、自動収集した教師画像に対応する色画像を生成する. また、自動収集した教師画像に対応する色画像を生成するため、本研究では光の反射・物体に生じる陰・半透明性による曇 りを再現するアルゴリズムと、GAN を組み合わせることで色 画像を自動生成することを提案する.さらに、ラベル画像予測 の精度向上のために、タスク環境上で実際の透明物体の色画像 を撮影して、二次元アノテーションを行うことで教師信号とし てラベル画像のみを持つ少数の学習データセットを用意する.

2.1 データ収集

ロボットに搭載された RGB-D カメラセンサを用いて,タ スク環境に透明物体の置かれていない色画像 *I*obj- および深度 画像 D_{obj}-を撮影する.撮影の際には色画像と深度画像の他 にも、カメラ内部パラメータ K とロボットのベースリンクか らカメラへの座標変換 R^{base} も同時に保存しておく.

2.2 深度画像およびラベル画像の生成

透明物体の写っていない色画像 *I*obj- に対し,物体メッシュ モデルを置くことのできる領域 (これを可配置領域とする)を 二次元アノテーションにより与える.続いて深度画像 *D*obj-のうち可配置領域以外を N/A として,点群を生成する領域を 限定する.これにより作られた深度画像と,収集済みのカメ ラパラメータ K および座標変換 *R*^{base}を用いて,ロボットの ベースリンク座標系での点群を生成し,その点群上に物体メッ シュモデルを仮想的に配置する.

配置の際の物体モデルの姿勢は、モデルを直方体に近似し て上を向く面を6面からランダムに決定し、その後鉛直軸周 りにランダムな角度だけ回転させることで決定される.またモ デルの位置に関しては、水平方向には可配置領域内でランダム に決定され、鉛直方向には重力に従って可配置領域内の面に接 するように決定される.

メッシュモデルを仮想配置した後にレイキャスティングを行っ てカメラ視点での物体モデルの深度画像とラベル画像 *L*_{obj-g} を取得し,そしてこの深度画像を収集済みの深度画像 *D*_{obj-} に重ねて,学習に用いる深度画像 *D*_{obj-g}を生成する.

2.3 色画像の生成

ここでは、学習時に用いるため自動生成された深度画像およ びラベル画像に対応するタスク環境の色画像を自動生成する.

2.3.1 光の反射・物体の陰・半透明による曇りの再現

色画像生成の第一段階として、光の反射・物体に生じる陰・物 体の半透明性による曇りを人工的に再現した画像 I_{obj-g} を生成 する. この段階には透明物体の写っていない色画像 I_{obj-} , 生 成された深度画像 D_{obj-g} およびラベル画像 D_{obj-g} を用いる.

再現アルゴリズムは以下のように説明される.まず最初に, 色画像 *I*_{obj} とラベル画像 *L*_{obj} を用いて,色画像のうち透 明物体の背景に相当する部分を抽出する.第二に,深度画像 *D*_{obj} に Canny アルゴリズムを適用してエッジを抽出し,各 画素のエッジからの L2 距離を計算した画像を用意する.これ と並行して,ラベル画像を透明か否かによって二値化したマス ク画像に水平右方向の Sobel フィルタを半時計回りに 45°だ け傾けたフィルタを適用し,さらに平滑化して物体左下部の エッジ付近が抽出された画像を用意する.

第三の工程では物体に陰をつける.第二の工程で作成した エッジからの距離画像の画素 i の画素値 x_i に対し,式(1)の第 2 式を用いて明度を決定する係数 v_i を計算する.ただし,第二 の工程で作成したエッジ付近の抽出された画像の画素値 y_i が 閾値 thre 以上ならば式(1)の第1式を用いて v_i を計算する.

ſ	$\frac{1}{(1 + \exp(-\beta(x_i + \theta)))^2}$	$(if y_i > thre)$	(1)
$v_i = \left\{ \right.$	$\frac{1}{1 + \exp(-\beta(x_i + \theta))}$	(otherwise)	(1)

この係数 v_i を最初の工程で抽出した色画像の R, G, B 各チャンネルに乗じることで物体に陰をつけ、その後陰をつけた画像 で入力色画像 I_{obj-} のマスク領域を置き換える.

第四の工程では透明物体の曇り,すなわち半透明性を表現 する.第三の工程で生成された画像のうちマスク領域につい て,各画素*i*のR,G,B各チャンネルの画素値*x_i*に不透明 度 opacity を考慮した式 (2) を適用して, 画素値 y_i を得る.

 $y_i = (1 - opacity) \cdot x_i + opacity \cdot 255 \tag{2}$

第五の工程では光源が画像上方向に存在すると仮定した時 の白色光源からの光の反射を表現する.垂直上方向の Sobel フィルタを深度画像 $D_{obj,g}$ に適用し,さらに平滑化を行って 光が反射して見える領域を抽出する.物体マスク領域に関し て,これにより得られる画像の各画素 i の画素値 x_i にパラメー タ reflect を乗じた上で,第四の工程で生成した色画像の R, G, B 各チャンネルの画素値 y_i を加えることで,出力色画像 $I_{obj,g}$ の R, G, B 各チャンネルの画素値 z_i を得る.

2.3.2 敵対的生成ネットワークを用いた人工色画像の変換

色画像生成の第二段階として、CycleGAN [Zhu 17] を用い て色画像 *I*_{obj-g} を変換し、変換後の色画像 *I*_{obj-g2} を訓練デー タセットの色画像として採用する.

本研究では、人工画像ドメインとして透明物体再現アルゴ リズムにより生成された色画像 *I*_{obj-g} から透明物体に注目し て切り出した ROI 画像を用いる.一方,実在画像ドメインに は実在する透明物体に注目した色画像群 *I*_{obj+} を用意する.こ れらを用いて学習したネットワークによって色画像 *I*_{obj-g} を 変換し、物体マスク領域のみを変換された画像で置き換えるこ とで、訓練データセットに用いる色画像 *I*_{obj-g} を得る.

以上により得られる色画像 *I*_{obj-g}, ラベル画像 *L*_{obj-g}, 深 度画像 *D*_{obj-g} の例を図 2 に示す.



図 2: 自動生成されたデータセットの例

3. ラベル画像予測と深度画像予測

生成された学習データセットを用いて,透明物体のラベル画 像予測および深度画像予測を行う多層畳みこみネットワーク (CNN)に学習させる.本研究の提案ネットワークは,図3に 示すようにセンサから取得できる色画像と深度画像の2つを 入力とし,透明物体のラベル画像と深度画像を出力する.

3.1 ラベル画像予測

本項および次項では、本研究で提案するネットワーク構造に ついて説明する.

ネットワークモデルのうちラベル画像予測を行う部分につい ては、生成された色画像 *I*obj-g2 を入力として学習させ、FCN-8s at-once [Shelhamer 17] を用いてフォワーディングを行う. なお各層の重みやバイアスについては、事前学習された VGG16 モデル [Simonyan 14] を用いて初期化することとする.

最終層の出力 y を用いて計算を行うセグメンテーションに 対する損失関数 $Loss_{seg}$ は,式 (3) に示すようにラベル画像 L_{obj-g} を教師信号とする Softmax Cross Entropy を用いる. $I(L_{obj-g}^{i, j})$ は指示関数で,教師ラベル画像の i 番目の画素にラ ベル $j \in \{0, 1, \dots, cls - 1\}$ が付されていれば 1,それ以外な ら 0 を返す.また $y^{i, j}$ は最終層の出力 y の i 番目の画素の第



図 3: 提案ネットワークモデル. 左上:入力深度画像, 左下:入力色画像, 右上:出力深度画像, 右下:出力ラベル画像

j チャンネルの値を表す.

$$Loss_{seg} = -\frac{1}{n_{pix}} \sum_{i=0}^{n_{pix}-1} \sum_{j=0}^{cls-1} I(L_{obj_g}^{i, j}) \log(Softmax(y^{i, j}))$$
(3)

3.2 深度画像予測

深度画像予測を行う部分については、生成された色画像 *Iobj_g2* と透明物体の置かれていない深度画像 *Dobj* を入力 として学習させる.入力深度画像 *Dobj* は JET カラーマップ を用いて3チャンネルの色画像に変換しておく.この画像に畳み 込みを行って得られる512チャンネルの特徴量画像 *pool5depth* を得た後に、教師ラベル画像 *Lobj_g*(テスト時は予測ラベル画 像 *Lpred*)をこれと同じサイズになるよう 1/32 にリサイズす る.そして予測ラベルが 0(背景)の画素のみを残してそれ以 外の画素値を0とするよう、*pool5depth*の各チャンネルごとに マスキングを行う.その後、入力色画像に畳み込みおよび最大 プーリングを順次施して画像サイズ 1/32 になった段階で最大 プーリングを行った後の 512 チャンネル特徴量画像とマスキ ングされた画像を連結させて、1024 チャンネル特徴量画像を 生成し、この特徴量画像から深度画像予測を行う.

深度画像予測の損失関数は,式(4)に示す画像全領域の損 失関数 $Loss_{depth}^{all}$ と,式(5)に示す教師マスク領域の損失関数 $Loss_{depth}^{mask}$ からなる. D_{pred}^{i} および D_{obj-g}^{i} はそれぞれ予測深 度画像,教師深度画像のi番目の画素を表す.また M^{i} は教師 ラベルが1以上(透明物体)ならば画素値1,それ以外(背景) ならば0となる教師マスク画像のi番目の画素を表す.

$$Loss_{depth}^{all} = \frac{1}{n_{pix}} \sum_{i=0}^{n_{pix}-1} (D_{obj_g}^{i} - D_{pred}^{i})^{2}$$
(4)

$$Loss_{depth}^{mask} = \frac{1}{sum(M^{i})} \sum_{i=0}^{n_{pix}-1} (D_{obj-g}^{i} - D_{pred}^{i})^{2} M^{i}$$
(5)

3.3 ネットワークの重み更新

ネットワーク全体の損失関数として,式(6)に示すように損 失関数 *Loss_{total}* を設計する.

$$Loss_{total} = Loss_{seg} + \lambda_1 (Loss_{depth}^{all} + \lambda_2 Loss_{depth}^{mask}) \quad (6)$$

本研究ではパラメータ $\lambda_1 = \lambda_2 = 10$ とする. この損失関数を 用いて誤差逆伝播を行い,ネットワークの重みを更新する. た だし,本研究ではラベル予測精度の向上のため,自動生成した 主データセットと実環境を二次元アノテーションして用意した 少数の副データセットを交互に学習する. 副データセットは教 師深度画像を持たないので,これの学習時は $\lambda_1 = \lambda_2 = 0$ とし て深度予測ネットワークの誤差逆伝播を行わないことにする.

深度画像予測ネットワークのマスキング処理に対しては誤差 逆伝播を行わないこととするが,特徴量の連結処理については 誤差逆伝播を行う.従って,ラベル予測ネットワークでは教師 ラベル画像と教師深度画像を用いた特徴量抽出が可能となる.

4. 検証実験

検証実験では、透明物体としてペットボトルを用いる.生成 したデータセットおよびネットワークモデルの有用性を検証す る実験を行い、さらに透明物体のピッキング実験を行う.

4.1 生成色画像データの検証実験

自動生成した色画像での学習が実環境画像を用いたテスト 時にも有効であることを検証する.本検証実験では,主データ セット (Main)1610 組,副データセット (Sub)120 組を用意す る.これらのデータを既存の FCN-8s at-once [Shelhamer 17] ネットワークモデルに学習させ,実環境画像を三次元アノテー ションして用意した 130 組のテストデータに対してラベル画 像予測を行った結果を表 1 に示す.なお,ここでは有色と無 色を分けて結果を記載している.130 組のテストデータの内訳 は, (A) 有色透明物体 22 組, (B) 無色透明物体 108 組となっ ている.

表 1: FCN-8s at-once による主データセットと副データセット のラベル画像予測結果. (A) 有色透明物体, (B) 無色透明物体.

Dataset	Mean IU			
	(A)	(B)	(A) + (B)	
Main	0.582	0.511	0.523	
Sub	0.796	0.713	0.727	
Main + Sub	0.793	0.673	0.693	

表1から,有色透明・無色透明のいずれの場合においても, 生成した主データセットと少数の副データセットを交互に学習 することで,実環境画像である副データセットのみの学習時に 近い Mean IU を得られたと言える.

4.2 ラベルと深度を予測するネットワークの検証実験

本研究の提案ネットワークとの比較対象として, 色画像か らラベル画像を予測する CNN および深度画像とラベル画像か ら深度画像を予測する CNN の二つを用意する. これら二つの ネットワークが予測したラベル画像および深度画像の精度と, 本研究の提案ネットワークが予測したラベルおよび深度画像 の精度を比較することで評価を行う. 比較対象のうち深度予測 ネットワークの学習時には提案手法同様にマスキング処理を行 うが,特徴量画像の連結処理については行わないこととする.

表2に,比較対象の二段階ネットワークと提案ネットワー クによるラベル予測の Mean IU と深度画像予測の精度を評価 した実験の結果を示す.深度画像予測精度は,ネットワークの 出力深度画像のうち教師マスク領域について,教師深度画像と の誤差の絶対値がx未満である画素数の割合を表している.

表 2: 二段階予測ネットワークと提案ネットワークのラベル画 像および深度画像予測結果

Network	Mean IU	Depth accuracy $(error < \mathbf{x})$		
		x = 1 cm	x = 2cm	x = 5cm
Staged	0.693	0.264	0.487	0.850
Ours	0.670	0.242	0.454	0.850

表2から二段階予測ネットワークを用いても本研究の提案 手法を用いても精度に大きな差はないと言える.しかし,二段 階予測のうち深度予測ネットワークの学習はラベル予測の学習 が終わってからでないと進行できないため,提案ネットワーク は二段階予測に比べ学習時間を低減できるという長所がある.

4.3 透明物体の把持への適用実験

最後に、ロボットによる物体把持行動および把持収納行動に 対し本研究で提案した深度画像予測システムを適用する実験を 行う.透明物体の代表としてペットボトルを家庭用冷蔵庫から 取り出すタスクを実現し、5回中4回成功した.また、ペット ボトルを含む5つの家庭内物体を連続して籠から荷箱へ移し 替えるタスクを実現した.図4および図5にそれぞれペット ボトルを取り出している様子、移し替えている様子を示す.



図 4: 冷蔵庫からペットボトルを取り出す様子



図 5: 籠からペットボトルを取り出し、荷箱へ収納する様子

5. 結論

本研究では透明物体のピッキング動作を行うための物体三 次元位置推定として、CNNを用いてラベル画像と深度画像を 同時に予測し、センサの深度画像のうち透明領域のみを予測 された深度画像で置換して出力することを提案した.また学 習データを効率的に収集するため、透明物体の三次元メッシュ モデルを用いてタスク環境を再現し、さらに光の反射・陰・物 体の曇りの再現アルゴリズムと GAN を組み合わせることで、 人間による三次元アノテーションの不要なデータセット生成手 法を提案した.検証として、透明物体であるペットボトルを家 庭用冷蔵庫から取り出す実験を行い、5回中4回成功した.ま た、ペットボトルを含む5つの家庭内物体を連続して籠から 荷箱へ移し替える実験を行い、成功した.

透明なショーケースの中にある物体のマニピュレーションに は、ショーケースと内部の物体の両方の位置を認識することが 必要となる.今後は本研究を応用して、これらの三次元位置推 定を実現することを目指す.

参考文献

- [Eigen 14] Eigen, D., et al.: Depth Map Prediction from a Single Image using a Multi-Scale Deep Network, in Advances in Neural Information Processing Systems 27, pp. 2366–2374, Curran Associates, Inc. (2014)
- [Eren 09] Eren, G., et al.: Scanning from heating: 3D shape estimation of transparent objects from local surface heating, *Opt. Express*, Vol. 17, No. 14, pp. 11457– 11468 (2009)
- [Liu 15] Liu, F., et al.: Deep Convolutional Neural Fields for Depth Estimation From a Single Image, in *The IEEE Conference on Computer Vision and Pattern Recognition* (CVPR) (2015)
- [Rantoson 10] Rantoson, R., et al.: 3D reconstruction of transparent objects exploiting surface fluorescence caused by UV irradiation, in 2010 IEEE International Conference on Image Processing, pp. 2965–2968 (2010)
- [Shelhamer 17] Shelhamer, E., et al.: Fully Convolutional Networks for Semantic Segmentation, *IEEE Transac*tions on Pattern Analysis and Machine Intelligence, Vol. 39, No. 4, pp. 640–651 (2017)
- [Silberman 12] Silberman, N., et al.: Indoor Segmentation and Support Inference from RGBD Images, in ECCV (2012)
- [Simonyan 14] Simonyan, K., et al.: Very Deep Convolutional Networks for Large-Scale Image Recognition, *CoRR*, Vol. abs/1409.1556, (2014)
- [Wang 15] Wang, P., et al.: Towards Unified Depth and Semantic Prediction From a Single Image, in *The IEEE Conference on Computer Vision and Pattern Recognition* (CVPR) (2015)
- [Zhu 17] Zhu, J., et al.: Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks, *CoRR*, Vol. abs/1703.10593, (2017)