

仮想空間の農業ハウス内の画像を教師データとした 画像セグメンテーションの性能評価

Evaluation of Image Segmentation Model using Training Data from Virtual Agricultural Houses

伊藤亮¹⁾, 小田川晴奎¹⁾, 安達武範²⁾, 中野智三²⁾, 山本聡史³⁾

Ryo Ito, Haruki Odagawa, Takenori Adachi, Tomomi Nakano, and Satoshi Yamamoto

1) 秋田県立大学システム科学技術学部 (〒015-0055 秋田県由利本荘市土谷字海老ノ口84-4,
E-mail: ryo.ito@akita-pu.ac.jp)

2) 株式会社プラスプラス (〒020-0857 岩手県盛岡市北飯岡1-10-85 B1)

3) 秋田県立大学生物資源学部 (〒010-0044 秋田県南秋田郡大潟村南2-2)

We investigated the performance of object recognition when learning solely from images in virtual space without using actual photographs. Additionally, to enhance object recognition performance, we employed image generation AI to process images and examined its effects. Results showed that significant overfitting occurred when using images from virtual space without any processing. However, applying image processing with image generation AI effectively suppressed overfitting and improved the mean Intersection over Union (mIoU) score.

Key Words : VR, Image segmentation, Agricultural house, Neural network, Machine learning

1. はじめに

近年の機械学習技術の飛躍的な発展に伴い、特に画像から物体認識を行うタスクでは機械学習の応用が必須になっていると言える。このようなタスクで新規に学習を行う際、通常は多数の写真などを用意し、それぞれの写真のどこに何が映っているかを示すデータセットを作成する作業（アノテーション作業）が必要であるが、これには多大な労力が必要になる。これに対し、可視光外の波長の光に反応する塗料を用いてアノテーションを自動化する手法[1]が提案されているが、課題として特殊なデバイスや暗室などの環境が必要になることや、対象の物体によっては光が拡散または透過しまう点などが挙げられる。一方でVR技術を応用する手法[2]が提案されている。この手法では、仮想空間さえできてしまえば学習用のデータの取得やアノテーションは自動的に行えるが、仮想空間の画像と実物の映像（写真）との差異が学習に悪影響を及ぼすことが考えられる。

本研究では画像からの物体検出技術の中でも画像セグメンテーションを取り上げ、実際の写真を用いずに仮想空間の農業ハウスをもとに得られた画像から学習を行っ

た場合の、物体認識の性能について調査した。また物体認識の性能を向上させることを目的として画像生成AIを用いた画像の加工を実施し、その効果を調査した。

2. モデルの作成と自動アノテーション

本研究では仮想空間の農業ハウスとしてトマトのハウス栽培を取り上げる。秋田県立大学アグリイノベーション教育研究センターの施設を対象とし3Dモデルを作成した[3]。オープンソース3Dソフトウェア Blender を使い、その標準機能である Geometry Nodes と Python スクリプトによってランダムな果樹の樹形を生成する。さらにゲーム開発プラットフォームとして広く用いられている Unity 上にそれらを配置し、その後それぞれに果実モデルを配置した（図1）。画像セグメンテーションのモデルを学習させるためのデータセットを自動で生成するため、Unity上で葉や茎、実など異なるマテリアル毎にID値を数値設定可能なシェーダーを作成した。カメラ位置をランダムに変えながら同じ画角でカラー画像とID値をもとにしたグレースケール画像（タグ付け画像）を出力することで、自動でデータセットを作成することができる（図2）。



図 1. 仮想空間における農業ハウス（トマト）の作成



図 2. 学習データの自動生成

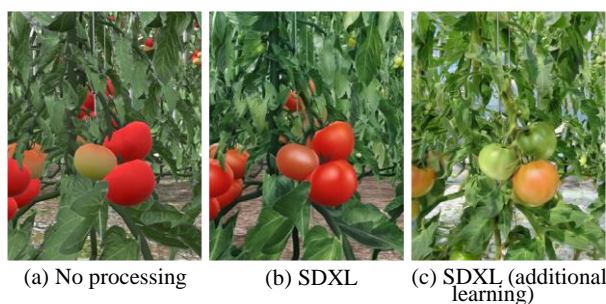


図 3. Stable Diffusion XL による画像の加工

3. 画像生成AIによるオーギュメンテーション

自動アノテーションによって得られた画像データは現実の農業ハウス内の写真と比較し多様性がなく色調が単調であるため、このデータを用いて学習したニューラルネットワークで現実のトマトハウスで推論を実施した場合、画像認識の性能（ここではmean Intersection over Union: mIoUで評価）に悪影響が出ることが想定される。これに対し近年急速に進歩した画像生成AIの技術を応用することで、タグ付け画像はそのまま（検出対象のトマトの実や萼などの位置や形状を変えない）の状態で学習用画像をより現実の写真に近いものに加工することを試みた。画像生成AIとしては Stable Diffusion XL (SDXL) を使い、その標準機能であるControl Net (canny)を用いることで、できるだけ元の画像の構造が変化しないようにしている。図3に、未加工の画像、SDXLで加工した画像、および追加学習を施したSDXLで加工した画像の例を示す。

4. 学習結果および考察

学習には仮想空間の農業ハウスから自動アノテーションによって得られた3000枚の画像 (256×256pixel) およびそれらに対応するタグ付け画像を用い、評価には現実のトマトハウスの写真を手動でアノテーションし作成したデータセットを用いた。なおここでは実、萼および花柄の3つのみを認識対象とした。ニューラルネットワークとしてはU-Net (エンコーダ部分は MobileNetV3 Large) を用い、ImageNet により学習済みの状態からファインチューニングを行った。また学習率を80Epoch周期で変動させる。

図4に学習データに対してのmIoU、図5に評価データ（写真）に対してのmIoUを示す。両者を比較すると、加工していない画像では学習データのmIoUが向上していくのに対し評価データのmIoUが悪化しており、過学習が発生している。これに対しSDXLで加工したものは過学習が抑制されていることが分かる。また、追加学習を施したSDXLの結果から今回の追加学習の効果はほとんどなかったと考えられるが、SDXLで加工した画像1500枚と追加学習を施したSDXLで加工した画像1500枚を混合して学習データとしたもの (Mixed) はmIoUの値が最も良い結果となった。これは混合によってデータに多様性が生じたためと考えている。しかしながら学習回数を増やしてもなか

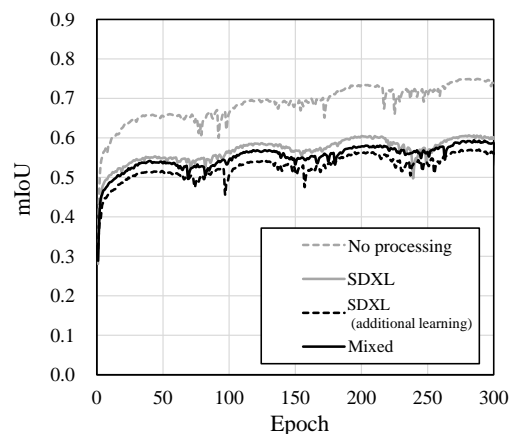


図 4. 学習データに関する mIoU

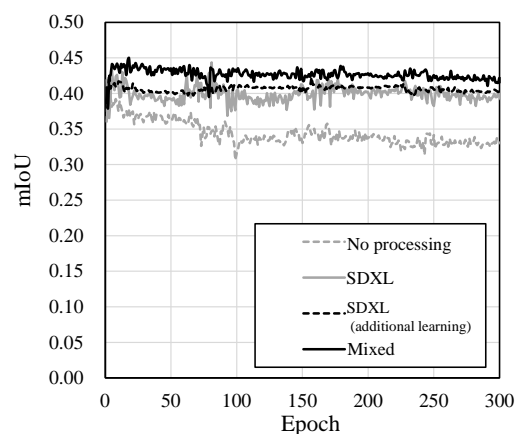


図 5. 評価データ（写真）に関する mIoU

なか学習が進んでいかない状況であるため、実用に向けては今後も認識精度の向上に取り組んでいく必要がある。

5. まとめ

本研究では画像セグメンテーションにおいて、実際の写真を用いずに仮想空間の画像のみから学習を行った場合の、物体認識の性能について調査した。また物体認識の性能を向上させることを目的として画像生成AIを用いた画像の加工を実施しその効果を調査した。仮想空間の画像を加工せずに使用した場合は顕著な過学習が発生してしまったのに対し画像生成AIを用いた画像の加工を実施した場合は過学習が抑制され、よりmIoUが良くなることが示された。しかしながら実用に十分な性能は得られておらず、引続き最適な画像データの加工方法やモデル構造などの検討が必要と考えられる。

参考文献

- [1] K. Takahashi and K. Yonekura: Invisible Marker: Automatic Annotation of Segmentation Masks for Object Manipulation, 2020 IEEE/RSJ International Conference on Intelligent Robots and Systems, pp.8431-8438, 2020.
- [2] S. Borkman, et al.: Unity Perception: Generate Synthetic Data for Computer Vision, arXiv: 2107.04259v2, 2021.
- [3] 山本聡史ほか: 農業におけるデジタルツイン: アグリデジタルツインの可能性, 計算工学講演会論文集, Vol.27, F-11-01, 2022.