# Effective Data Fusion with Generalized Vegetation Index: Evidence from Land Cover Segmentation in Agriculture

Hao Sheng[1], Xiao Chen[1], Jingyi Su[2], Ram Rajagopal[1], and Andrew Ng[1]

[1]Stanford University
[2]Chegg, Inc

[1]{haosheng,markcx,ramr}@stanford.edu   ang@cs.stanford.edu   [2]jingyi.su.js@gmail.com

## Abstract

*How can we effectively leverage the domain knowledge from remote sensing to better segment agriculture land cover from satellite images? In this paper, we propose a novel, model-agnostic, data-fusion approach for vegetation-related computer vision tasks. Motivated by the various Vegetation Indices (VIs), which are introduced by domain experts, we systematically reviewed the VIs that are widely used in remote sensing and their feasibility to be incorporated in deep neural networks. To fully leverage the Near-Infrared channel, the traditional Red-Green-Blue channels, and Vegetation Index or its variants, we propose a Generalized Vegetation Index (GVI), a lightweight module that can be easily plugged into many neural network architectures to serve as an additional information input. To smoothly train models with our GVI, we developed an Additive Group Normalization (AGN) module that does not require extra parameters of the prescribed neural networks. Our approach has improved the IoUs of vegetation-related classes by $0.9 - 1.3$ percent and consistently improves the overall mIoU by $2$ percent on our baseline.*

## 1. Introduction

Deep learning has been widely adopted in computer vision across various applications such as diagnosing medical images[1], classifying objects in photos[2], annotating video frames[3], etc. However, recognizing the visual patterns in the context of agriculture, especially segmenting the multi-labeled masks, has not been explored extensively in detail. One primary reason that hinders the progress is the difficulty of handling complex multi-modal information inside the images[4] because the sensing imagery in agriculture contains Near Infrared band and other thermal bands

that are distinguished from traditional images spanning over red, green, and blue (RGB) visual bands. Such multi-band information is crucial for understanding the land cover context and field conditions, e.g., the vegetation of the land.
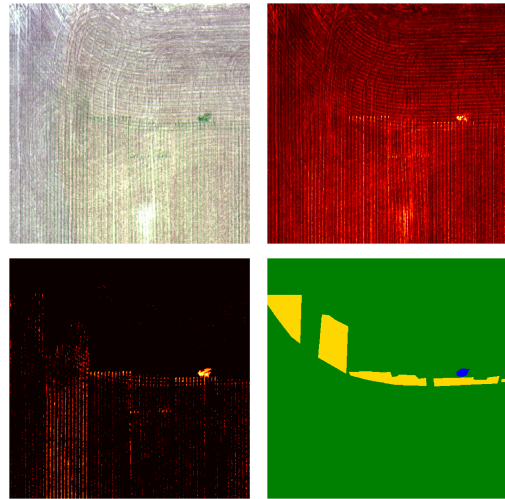


Figure 1. An example of an NRGB image and its Vegetation Index (VI) and ground-truth labels. **Top-left:** Input RGB channels; **Top-right:** Input near-infra red (NIR) channel; **Bottom-left:** Vegetation Condition Index (VCI)[5] calculated based on RGB and NIR channel; **Bottom-right:** Ground-truth labels, where yellow denotes *Double Plant* and blue denotes the *Weed Cluster*. VCI is able to pick up both *Weed Cluster* (a cluster of very high VI values) and *Double Plant* (lanes of different VI values compared to the background crops).

To leverage the information of multiple distinct bands in the images, researchers in the last several decades have focused on developing different algorithms and metrics to perform the land segmentation[6, 7]. As discussed in the literature review in Section 2, the design of Vegetation Index (VI) has been essential for studying land cover segmentation[8–

1

11]. The key idea of VI is to assess the vegetation of a region based on the reflectances from multiple bands, including the Near-Infrared band and other thermal bands, and hence ultimately approximate the region's land cover segments. Nevertheless, in the context of deep learning, we have yet to investigate how to leverage the domain knowledge of VI while making use of models learned or transferred from non-agriculture data to segment the land accurately.

To tackle this question, we describe a general form of VI that serves as an additional input channel for image segmentation. Such a general form of VI covers many specific VI variants in existing studies[10–13], which motivate us to develop a generalized learnable VI block that fuses the VIs and images in a convolution fashion. Based on the fused input, we also propose a new additive group normalization, a natural generalization of the instance normalization and layer normalization, because the VI channel and RGB channels can be considered as different groups.

Our work contributes to the research of agriculture land cover segmentation in three ways. Firstly, we systematically compare the vegetation indices that primarily depend on the Near-Infrared, red, green, and blue channels. We highlight the key idea of calculating VIs and disclose the connections among them. Secondly, we propose a model-agnostic module named General Vegetation Index (GVI) that captures many existing VI features. This module particularly fits convolutional neural networks, even for the pretrained models very well, because it doesn't need to change model structures too much. Thirdly, we introduce the additive group normalization (AGN) that helps to fine tune models smoothly when GVI is introduced to a pretrained model. With these components in place, we modified a model based on DeepLabV3[14] and ran experiments on land segmentation in agriculture. With careful evaluations, we achieved an mIoU of 46.89% which exceeds the performance of the baseline model by about 2 percent.

## 2. Related Work

**Vegetation Index.** Vegetation Indices (VIs) are simple and effective metrics that have been widely used to provide quantitative evaluations of vegetation growth[9]. Since the light spectrum changes with plant type, water content within tissues and so on[15, 16], the electromagnetic waves reflected from canopies can be captured by passive sensors. Such characteristics of the spectrum can provide extremely useful insights for applications in environmental and agricultural monitoring, biodiversity conservation, yield estimation, and other related fields[17]. Because the land vegetation highly correlates with the land cover reflectance, researchers have built more than 60 VIs in the last four decades with mainly the following light spectra: (i) the ultraviolet region (UV, 10-380 nm); the visible spec-

tra, which consists blue (B, 450-495 nm), green (G, 495-570 nm) and red (R, 620-750 nm); (iii) the near and mid-infrared band (NIR, 850-1700 nm)[18, 19]. Such VIs are validated through direct or indirect correlations with the vegetation characteristics of interest measured *in situ*, such as vegetation cover, biomass, growth, and vigor assessment[9, 20].

To our best knowledge, the first VI, i.e., the Ratio Vegetation Index (RVI), was proposed by Jordan[8] in 1969. RVI was developed with the principle that leaves absorb relatively more red than infrared light. Widely used at high-density vegetation coverage regions, RVI is sensitive to atmospheric effects and noisy when vegetation cover is sparse (less than 50%)[7]. The Perpendicular Vegetation Index (PVI)[21] and the Normalized Difference Vegetation Index (NDVI)[10] followed the same principle but to normalize the output, having a sensitive response even for a low vegetation coverage. To eliminate the effects of atmospheric aerosols and ozone, Kaufman and Tanre [22] proposed the Atmospherically Resistant Vegetation Index (ARVI) in 1992, and Zhang *et al.* [11] improved the ARVI by eliminating its dependency to a 5S atmospheric transport model[23]. Another direction was to improve VI's robustness against different soil backgrounds[21]. The Soil-Adjusted Vegetation Index (SAVI)[24] and modified SAVI (MSAVI)[25, 26] turned out to be much less sensitive than the RVI to changes in the background. Based on ARVI and SAVI, Liu and Huete introduced a feedback mechanism by using a parameter to simultaneously correct soil and atmospheric effects, which they called the Enhanced Vegetation Index (EVI)[27]. With the recent progress in remote sensing (increasing number of bands and narrower bandwidth)[28], more VIs are being built to capture not only the biomass distribution and classification, but also chlorophyll content (Chlorophyll Absorption Ratio Index(CARI))[29], plant water stress (Crop Water Stress Index (CWSI))[30], and light use efficiency (Photochemical Reflectance Index (PRI))[20, 31]. With these aforementioned studies, a summary of VIs that derives from NIR-Red-Green-Blue (NRGB) images can be found in Table 1. Although we refer to a full literature review of VIs in [32] and [9], many VIs share similar form that motivates us to find a generalized formula to capture the essence of VIs.

**Remote Sensing with Transfer Learning and Data Fusion.** As an emerging interdisciplinary field, remote sensing (on both aerial photographs and satellite images) with deep learning has experienced quite a few benchmark datasets that have been released in recent years, such as FMOW[33], SAT-4/6[34], EuroSat[35], DeepGlobe 2018[36], Agriculture-Vision [4] and so on. Most of those datasets come with more than the visible band (i.e., RGB), including near and mid-infrared band (NIR) and sometimes shortwave red (SW). The different input structure, together with the context switch from a human-eye dataset (such as

ImageNet[37]) to a bird's-eye dataset, makes Transferring Learning less straightforward. Penatti *et al.* [38] systematically compared ImageNet pretrained CNN with other descriptors (feature extractors, e.g. BIC) and found it achieve comparable but not the best performance in detecting coffee scenes. Xie *et al.* [39] has shown simply adopting the ImageNet pretrained model while discarding the extra information does not achieve the best result in predicting the Poverty Level. Zhou *et al.* [40] has also observed the similar phenomena in their Road Extraction task. In addition, a two-stage fine-tuning process is proposed in [41], where an ImageNet pretrained network is further fine-tuned on a large satellite image dataset with the first several layers frozen. An alternative direction in exploring the large-scaled but not well-labeled data is to construct satellite-image-specified geo-embedding through weakly supervised learning[42], or unsupervised learning with Triplet Loss[43]. These aforementioned steps motivate us to use a pretrained model based on ImageNet, which has been demonstrated to have a good performance empirically in transfer learning.

In [44], Sidek and Quadri defined data fusion as "dealing with the synergistic combination of information made available by different measurement sensors, information sources, and decision-makers." Studies in the deep learning community have also proposed data fusion approaches that are specific to satellite images at a different level in practice. For example, [45] concatenates LiDAR and RGB to predict roof shape better. In DeepSat [34], Basu *et al.* achieves the state of the art performance on SAT-4 and SAT-6 land cover classification problems by incorporating NDVI[10], EVI[12] and ARVI[11] as additional input channels. A recent study[46] proposed a novel approach to select and combine the most similar channels using images from different timestamps. Apart from the multi-channel data fusion, fusions at multi-source [47, 48] and multi-temporal[49] levels have also shown their empirical value. Such an idea of fusing the multiple input channels also inspired our design of the fusion module of General Vegetation Index.

**Multi-spectral Image Data Fusion.** Multi-spectral image data fusion is also widely used in robotics, medical diagnoses, 3D inspection, etc[50]. *Color related techniques* represent color in different spaces. The Intensity-Hue-Saturation (IHS fusion)[51] transforms three channels of the data into the IHS color space, which separates the color aspects in its average brightness (intensity). The values in IHS space correspond to the surface roughness, its dominant wavelength contribution (hue), and its purity (saturation) [52, 53]. Then, one of the components is replaced by a fourth channel that needs to be integrated. *Statistical/numerical methods* introduce a mathematical combination of image channels. The Brovey algorithm[54] calculates the ratio of each image band by summing up the cho-

sen bands, followed by multiplying with the high-resolution image.

In addition to concatenating multi-spectral channels, several deep learning architectures were proposed for multi-spectral images. [55] pretrains a Siamese Convolution Network to generate a weighted map for infrared and RGB channels in the inference time. Li *et al.* [56] first decomposes the source images into base background and detail content and then applies a weighted average on the background while using a deep learning network to extract multi-layer features for detail content.

These studies gave an initial attempt to tackle the image classification problem using multiple spectral inputs in deep learning models. But we have yet to investigate how the multi-spectral image can be translated into the VI-related input in the context of agriculture segmentation.

## 3. Proposed Method

### 3.1. Overview

In general, our approach hinges on fusing Vegetation Index with raw images. We first introduce using well-known VIs as another input channel, and then we generalize the idea of VI to a fully learnable data fusion module. Last but not least, we propose an Additive Group Normalization (AGN) to handle the warm-start with a pretrained model. We describe the technical details in the following subsections.

### 3.2. Vegetation Index for Neural Nets

According to [9], during the practice of remote sensing, more than 60 VIs have been developed in the last four decades. However, not all VIs are derived from NIR and RGB channels, few of which generalize across datasets without tuning their sensitive parameters manually. For example, the Perpendicular Vegetation Index (PVI) [21] is defined as follows:

$$\text{PVI} = \sqrt{(\rho_{\text{soil}} - \rho_{\text{veg}})^2_{\text{R}} - (\rho_{\text{soil}} - \rho_{\text{veg}})^2_{\text{NIR}}} \quad , \quad (1)$$

where $\rho_{\text{soil}}$ is the soil reflectance and $\rho_{\text{veg}}$ is the vegetation reflectivity. However, PVI is sensitive to soil brightness and reflectivity, especially in the case of low vegetation coverage, and needs to be re-calibrated for this effect[22]. Such sensitivity introduces semantic difficulty as we try to feed the VI into the neural network as another input channel. There are also VIs designed for a specific dataset in the first place. On top of the Landsat Multispectral Scanner (MSS), Landsat Thematic Mapper (TM) and Landsat 7 Enhanced Thematic Mapper (ETM) data, Cruden *et al.* [19] applied a Tasseled Cap Transformation and came up with the empiri-

cal coefficients for Green Vegetation Index (GVI) as:

$$\text{GVI} = -0.290\text{MSS}_4 - 0.562\text{MSS}_5 + 0.600\text{MSS}_6$$
$$+ 0.49\text{MSS}_7, \tag{2}$$

where $\text{MSS}_i$ denotes the $i$th band of Landsat MSS. Landsat TM and Landsat 7 ETM are not usually available for satellite and aerial imagery outside this product family. This Green VI is composed by a linear combination of the multi-channel input, which shares a similar concept among many other VIs. To better understand the popular format of different indices, we summarized some representative VIs, shown in Table 1, that are derived from NIR-Red-Green-Blue (NRGB) images, together with their definitions and value ranges. Based on the definitions, we calculate the pixel-wise correlation matrix for all 12 VIs (Figure 2). The correlation coefficients are calculated at the pixel level using all data released for training. For SAVI, we choose $L$ to be 0.5. Except for certain pairs (such as NDVI v.s. SAVI),

Figure 2. Pair-wise correlation coefficients of all 12 available vegetation indices.

the correlation between different VIs are within the range of $(-0.2, 0.9)$. We include all 12 VIs as extra input channels in our experiments when leveraging the information from existing vegetation indices.

### 3.3. Learnable Vegetation Index

Some high correlations between VIs stem from not only the fundamental vegetation status, but also the empirical function that researchers have introduced. We notice that 9 out of the 13 VIs from Table 1 share the following general form:

$$\text{VI} = \frac{\alpha_0 + \alpha_\text{R}\text{R} + \alpha_\text{G}\text{G} + \alpha_\text{B}\text{B} + \alpha_\text{NIR}\text{NIR}}{\beta_0 + \beta_\text{R}\text{R} + \beta_\text{G}\text{G} + \beta_\text{B}\text{B} + \beta_\text{NIR}\text{NIR}}, \tag{3}$$

where $\alpha_c, \beta_c, \forall c \in \{\text{R, G, B, NIR}\}$ are parameters to be determined, and could be learnable in deep learning models. As suggested in [5], we can normalize the response (output) by nearby regions to supress outliers. By extending the pixel-wise operation to each neighborhood of image channels, we introduce a learnable layer of Generalized Vegetation Index (GVI):

$$GVI(\mathbf{x}, \alpha, \beta) = \frac{\mathbf{x} \circledast \alpha}{\mathbf{x} \circledast \beta}, \tag{4}$$

where $\circledast$ denotes the convolution operation, $\mathbf{x}$ is our NRGB inputs, and $\alpha, \beta$ are the learnable weights. In practice, we clip both the numerator and denominator to avoid numerical issues. Depending on the output channels, this layer has the capacity to express a variant number of VIs when learned. An illustrative example can be found in Figure 3.
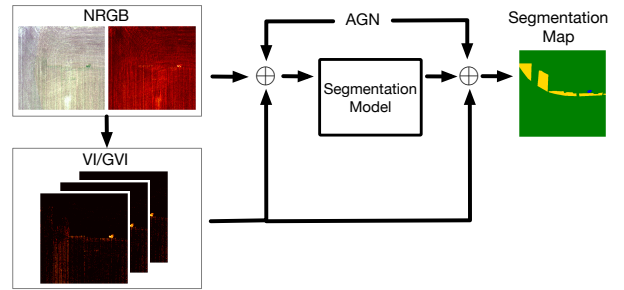
Figure 3. Our data fusion module is model-agnostic. The VI or GVI input channel is compatible to any segmentation model trained on NRGB image. Additive Group Normalization (AGN) is applied in the Near-Infrared channel with a linear combination of the batch normalization.

### 3.4. Additive Group Normalization Index

In contrast to explicitly normalizing the constructed indices using a ratio (e.g., in Equation (4)), we could normalize the value using nearby regions and channels, as we saw in VCI[5]. Fortunately, the deep learning community has already developed the counterpart approaches, such as the Batch Normalization (BN)[61], Layer Normalization (LN)[62], Instance Normalization (IN)[63] and Group Normalization (GN)[64]. However, we found that the neural network, even equipped with the most widely used BN, has an internal difficulty in fitting existing VIs, as shown in Figure 4.

The relative high errors in prediction indicate that BN is not able to captured channel-normalized features, while VIs are usually normalizing the inputs across the spectrum. Motivated by such observations, we introduce the Additive Group Normalization, which combine the BN and GM together in an additive fashion. Unlike BN, which normalizes each channel of features using the mean and variance computed over the mini-batch, GN splits channels into groups

| Index | Definition | Meaningful Range |
|-------|------------|------------------|
| NDVI[*] [10] | $\dfrac{\text{NIR} - \text{R}}{\text{NIR} + \text{R}}$ | $[0, 1]$ |
| IAVI[*†] [11] | $\dfrac{\text{NIR} - (\text{R} - \gamma(\text{B} - \text{R}))}{\text{NIR} + (\text{R} - \gamma(\text{B} - \text{R}))}, \gamma \in (0.65, 1.12)$ | $[-1, 1]$ |
| MSAVI2[*] [21] | $0.5 \left( (2\text{NIR} + 1) - \sqrt{(2\text{NIR}+1)^2 - 8(NIR - R)} \right)$ | $[0, 1]$ |
| EVI[*] [12] | $2.5 * \dfrac{\text{NIR} - \text{R}}{\text{NIR} + 6\text{R} - 7.5\text{B} + 1}$ | $(-\infty, \infty)$ |
| VDVI[*] [13] | $2 * \dfrac{2\text{G} - \text{R} - \text{B}}{2\text{G} + \text{R} + \text{B}}$ | $[-1, 1]$ |
| WDRVI[*] [6] | $\dfrac{0.2\text{NIR} - \text{R}}{0.2\text{NIR} + \text{R}}$ | $[-1, 1]$ |
| MCARI [57] | $\dfrac{1.5 * (2.5 * (\text{NIR} - \text{R}) - 1.3 * (\text{NIR} - \text{G}))}{\sqrt{(2\text{NIR} + 1)^2 - (6\text{NIR} - 5\text{R}) - 0.5}}$ | $(-1.6, 4.88)$ |
| GDVI [58] | $\text{NIR} - \text{G}$ | $[-1, 1]$ |
| SAVI[*†] [24] | $(1 + L) * \dfrac{\text{NIR} - \text{R}}{\text{NIR} + \text{R} + L}, L \in \{0, 0.5, 1\}$ | $[0, 1]$ |
| RVI[*] [59] | $\dfrac{\text{R}}{\text{NIR}}$ | $[0, \infty)$ |
| VCI [5] | $\dfrac{\text{NDVI} - \text{NDVI}_{\min}}{\text{NDVI}_{\max} + \text{NDVI}_{\min}}$ | $[0, 1]$ |
| GRVI[*] [58] | $\dfrac{\text{NIR}}{\text{G}}$ | $[0, \infty)$ |
| NDGI[*] [60] | $\dfrac{\text{G} - \text{R}}{\text{G} + \text{R}}$ | $[-1, 1]$ |

Table 1. Summary of vegetation indices that are derived from NIR-Red-Green-Blue (NRGB) images
[*]: These vegetation indices share the general format as equation 3
[†]: Parameters need to be calibrated and this VI cannot be fed into the neural network directly

and uses the within-group mean and variance to normalize the particular group:

$$\hat{x}_{nchw}^{GN} = \frac{x_{nchw} - \mu_{nc}^{(GN)}}{\sqrt{\sigma_{nc}^{2(GN)} + \epsilon}}$$

$$\mu_{nc}^{(GN)} = \frac{1}{HWG} \sum_{c \in G_c} \sum_H \sum_W x_{nchw}$$

$$\sigma_{nc}^{2(GN)} = \frac{1}{HWG} \sum_{c \in G_c} \sum_H \sum_W \left( x_{nchw} - \mu_{nc}^{(GN)} \right), \quad (5)$$

where $G$ is the number of groups, $G_c$ is the group assignment of channel $c$, and $\hat{\mathbf{x}}^{(GN)} = \{\hat{x}_{nchw}^{(GN)}\}$ is the GN response. Depending on the number of groups, such a normalization can be reduced to either Instance Normalization ($G = C$) or Layer Normalization ($G = 1$).

Inspired by the adaptive Instance-Batch Normalization [65], we designed our Additive Group Normalization

(AGN) as follows:

$$\hat{\mathbf{x}}^{(AGN)} = \sigma(\rho) \cdot \hat{\mathbf{x}}^{(GN)} + \hat{\mathbf{x}}^{(BN)}, \quad (6)$$

where $\rho$ is a learnable parameter controlling the contribution of Group Normalization in each layer and $\hat{\mathbf{x}}^{(AGN)} \in \mathcal{R}^{N \times C \times H \times W}$ is the response of AGN. This normalization does not introduce extra parameters (except for the running mean and standard deviation) but leverages the existing capacity of the underlying network.

When $\rho$ is a large negative number, the term $\hat{\mathbf{x}}^{(G)}$ gets a negligible weight and $\hat{\mathbf{x}}^{(AGN)} \approx \hat{\mathbf{x}}^{(B)}$. This property makes fine-tuning of experiments much smoother on a pretrained model with an architecture of Batch Normalization: To control the "ramping up" of Group Normalization, we initialize $\rho$ with a negative number, e.g., $-10$, and the model weights are updated gradually. We show experimental results for both training from scratch and fine-tuning in Section 4.
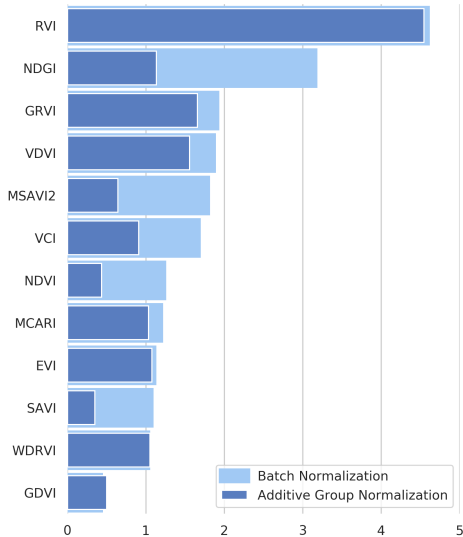
Figure 4. Mean L1 error over standard deviation (%). At each pixel, we trained a two-layer, fully-connected neural network with the NRGB channels to fit the Vegetation Indices using a batch size of 16. We plot the relative error for each vegetation index, i.e., the L1 error over the mean standard deviation in percentage. The additive group normalization fits almost all VIs better compared to batch normalization.

## 4. Experiments

### 4.1. Architecture Setup

We use EfficientNet-B0 / EfficientNet-B2[66] as our base encoder in the DeepLabV3[14] framework. They are parameter-efficient networks that achieve the same performance of ResNet-50 / ResNet-101 respectively with a much lower number of parameters.[67]

### 4.2. Training Details

We used backbone models pretrained on ImageNet in all our experiments. During initialization, we copied the pre-trained weights for the red channel filter to the one for the NIR channel in the first layer. We trained each model for 80 epochs with a batch size of 64 on eight GeForce GTX TITAN X GPUs. Unless specified, we used a combination of Focal Loss[68] and Dice Loss[69] with weights 0.75 and 0.25 respectively. We did not weigh classes differently, albeit the dataset is unbalanced. We also masked all the pixels that are either not valid or not within the region of the farmland. We use the Adam optimizer [70] with a base learning rate of 0.01 and a weight decay of $5 \times 10^{-4}$. During the training, we monitored the validation loss and stopped the experiments if the loss didn't decrease within ten epochs. Once a model was trained, we fine-tuned it with VI, GVI, or AGN modules. We adopted the cosine annealing strategy [71], with the learning rate ranges from 0.0001 to 0.01 and a cycle length of 10 epochs. For a fair comparison, we also fine-tuned our baseline model in this stage.

### 4.3. Dataset and Evaluation Metric

We evaluated our approach on Agriculture-Vision[4] with mean Intersection-over-Union(IOU) across classes. Since our annotations may overlap while we modeled the segmentation as a multi-class classification problem pixel-wisely, we also described the mean IOU calculation as follows.

**Agriculture-Vision.** Agriculture-Vision is an aerial image dataset that contains 21,061 farmland images captured throughout 2019 across the US. Each image is of size $512 \times 512$ and with four color channels, namely, RGB and Near Infrared (NIR). By the time the experiments are done, the labels for the test set have not been released yet, so we used the verication set to test the trained model.

**IOU with overlapped annotations.** We followed the protocol from the data challenge organizer to accommodate the evaluation for overlapped annotations. For pixels with multiple labels, a prediction of either label was counted as a correct pixel classification for that label, and a prediction that did not contain any ground truth labels was counted as an incorrect classification for all ground truth labels.

### 4.4. Results

Table 2 presents the validation results of the baseline model, together with several proposed methods to leverage the information from the NIR band. The average mIoU of the baseline model yeilds 44.92% accuracy. When plugging the GVI module into our model, we achieve 46.05% accuracy on mIoU, and highest accuracy in some categories such as *background*, *double plant*, and *standing water*. Moreover, when we use the Additive Group Normalization (AGN), the model performs the best in terms of mIoU at the accuracy of 46.87%. Our model consistently outperforms a) only using the NIR bands without extra information from vegetation indices; b) adding vegetation indices directly as inputs. And we saw gains in vegetation-related classes (e.g., Weed Cluster) as well as non-vegetation classes (e.g., Waterway). We include some examples in Figure 5.

## 5. Conclusion

In this work, we introduced the *General Vegetation Index* that enhanced the power of neural networks in agriculture and highlighted the connection between this GVI and other existing VIs. When starting from a pretrained model with minimal modifications, our proposed GVI and additive group normalization can achieve, and in some cases, exceed state-of-the-art performances. Our best result of mIOU is about 2% better than the baseline model. In addition, our method doesn't require sophisticated network architecture with the increase of model parameters. Such a result is a promising step forward when incorporating VI related in-
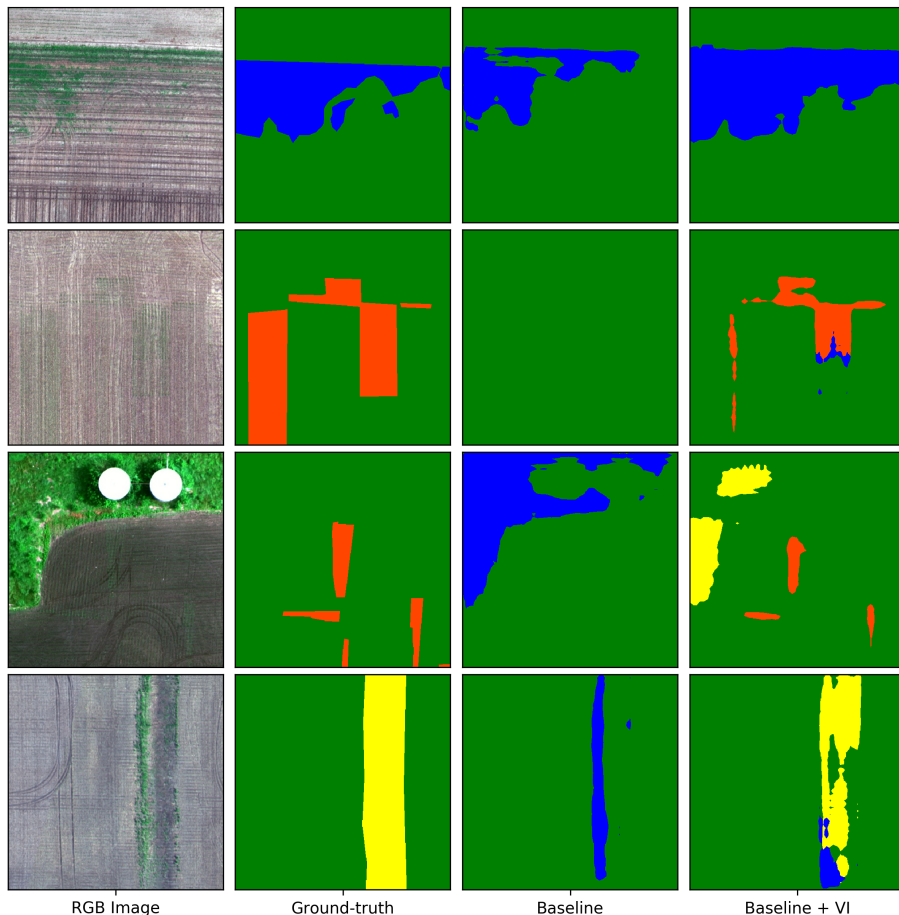
Figure 5. Examples of segmentation results. We include four examples (rows) of their RGB input, ground-truth labels, predictions of the baseline model, and predictions of the baseline model with VIs inputs (columns). Segmentation labels: Green for *background*, Blue for *Weed Cluster*, Red for *Double Plant* and Yellow for *Waterway*. Including VIs helps the model perform better in vegetation-related classes (e.g. *Weed Cluster*) as well as non vegetation classes (e.g. *Waterway*).

| Architecture | Method | mIoU (%) | Background | Cloud Shadow | Double Plant | Planter Skip | Standing Water | Waterway | Weed Cluster |
|---|---|---|---|---|---|---|---|---|---|
| DeepLabV3 | Baseline | 44.92 | 78.84 | 40.59 | 33.14 | 0.74 | 51.03 | 60.64 | 49.48 |
| | Baseline + VI | 46.04 | 78.75 | 41.17 | 33.66 | 0.46 | 56.67 | 62.06 | 49.50 |
| | Baseline + GVI | 46.05 | **79.81** | 34.58 | **35.24** | 0.83 | **58.08** | 63.47 | 50.32 |
| | AGN | **46.87** | 79.28 | **41.22** | 34.56 | **1.05** | 57.14 | **63.53** | **51.28** |

Table 2. mIoUs and class IoUs of baseline models, baseline models with Vegetation Index as additional models and our proposed generalized vegetation index model.

formation with multi-band images for segmentation tasks in agriculture.

While our approach sheds a promising light on segmenting lands in agriculture, we believe several potential directions could be valuable for future work. Firstly, how the model architecture can affect the result is still open for exploring. It is not clear if the segmentation results are sensitive to different models with VI inputs. Secondly, we would like to incorporate some additional training techniques, e.g.,

virtual adversarial training, which is orthogonal to our data fusion approach to improve the model performance further. Lastly, the ability to generalize our method on a larger scale dataset remains open to investigate.

## References

[1] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silviana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, et al. Chexpert:

A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 590–597, 2019.

[2] Barak Oshri, Annie Hu, Peter Adelson, Xiao Chen, Pascaline Dupas, Jeremy Weinstein, Marshall Burke, David Lobell, and Stefano Ermon. Infrastructure quality assessment in africa using satellite imagery and deep learning. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 616–625, 2018.

[3] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. Large-scale video classification with convolutional neural networks. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 1725–1732, 2014.

[4] Mang Tik Chiu, Xingqian Xu, Yunchao Wei, Zilong Huang, Alexander Schwing, Robert Brunner, Hrant Khachatrian, Hovnatan Karapetyan, Ivan Dozier, Greg Rose, et al. Agriculture-vision: A large aerial image database for agricultural pattern analysis. *arXiv preprint arXiv:2001.01306*, 2020.

[5] Felix N Kogan. Application of vegetation index and brightness temperature for drought detection. *Advances in space research*, 15(11):91–100, 1995.

[6] Anatoly A Gitelson. Wide dynamic range vegetation index for remote quantification of biophysical characteristics of vegetation. *Journal of plant physiology*, 161(2):165–173, 2004.

[7] J Grace, C Nichol, M Disney, P Lewis, Tristan Quaife, and P Bowyer. Can we measure terrestrial photosynthesis from space directly, using spectral reflectance and fluorescence? *Global Change Biology*, 13(7):1484–1497, 2007.

[8] Carl F Jordan. Derivation of leaf-area index from quality of light on the forest floor. *Ecology*, 50(4):663–666, 1969.

[9] Jinru Xue and Baofeng Su. Significant remote sensing vegetation indices: A review of developments and applications. *Journal of Sensors*, 2017, 2017.

[10] JW Rouse, RH Haas, JA Schell, and DW Deering. Monitoring vegetation systems in the great plains with erts. *NASA special publication*, 351:309, 1974.

[11] Renhua Zhang, XN Rao, and NK Liao. Approach for a vegetation index resistant to atmospheric effect. *Acta Botanica Sinica*, 38(1):53–62, 1996.

[12] Alfredo Huete, Kamel Didan, Tomoaki Miura, E Patricia Rodriguez, Xiang Gao, and Laerte G Ferreira. Overview of the radiometric and biophysical performance of the modis vegetation indices. *Remote sensing of environment*, 83(1-2):195–213, 2002.

[13] Wang Xiaoqin, Wang Miaomiao, Wang Shaoqiang, and Wu Yundong. Extraction of vegetation information from visible unmanned aerial vehicle images. *Transactions of the Chinese Society of Agricultural Engineering*, 31(5), 2015.

[14] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017.

[15] L Chang, S Peng-Sen, and Liu Shi-Rong. A review of plant spectral reflectance response to water physiological changes. *Chinese Journal of Plant Ecology*, 40(1):80–91, 2016.

[16] Chunhua Zhang and John M Kovacs. The application of small unmanned aerial systems for precision agriculture: a review. *Precision agriculture*, 13(6):693–712, 2012.

[17] David J Mulla. Twenty five years of remote sensing in precision agriculture: Key advances and remaining knowledge gaps. *Biosystems engineering*, 114(4):358–371, 2013.

[18] Hazli Rafis Bin Abdul Rahim, Muhammad Quisar Bin Lokman, Sulaiman Wadi Harun, Gabor Louis Hornyak, Karel Sterckx, Waleed Soliman Mohammed, and Joydeep Dutta. Applied light-side coupling with optimized spiral-patterned zinc oxide nanorod coatings for multiple optical channel alcohol vapor sensing. *Journal of Nanophotonics*, 10(3): 036009, 2016.

[19] Brett A Cruden, Dinesh Prabhu, and Ramon Martinez. Absolute radiation measurement in venus and mars entry conditions. *Journal of Spacecraft and Rockets*, 49(6):1069–1079, 2012.

[20] Alex Haxeltine and IC Prentice. A general model for the light-use efficiency of primary production. *Functional Ecology*, pages 551–561, 1996.

[21] Arthur J Richardson and CL Wiegand. Distinguishing vegetation from soil background information. *Photogrammetric engineering and remote sensing*, 43(12):1541–1552, 1977.

[22] Yoram J Kaufman and Didier Tanre. Atmospherically resistant vegetation index (arvi) for eos-modis. *IEEE transactions on Geoscience and Remote Sensing*, 30(2):261–270, 1992.

[23] D Tanré, C Deroo, P Duhaut, M Herman, JJ Morcrette, J Perbos, and PY Deschamps. Technical note description of a computer code to simulate the satellite signal in the solar spectrum: the 5s code. *International Journal of Remote Sensing*, 11(4):659–668, 1990.

[24] Alfredo Huete. Huete, ar a soil-adjusted vegetation index (savi). remote sensing of environment. *Remote sensing of environment*, 25:295–309, 1988.

[25] Jiaguo Qi, Abdelghani Chehbouni, Alfredo R Huete, Yann H Kerr, and Soroosh Sorooshian. A modified soil adjusted vegetation index. 1994.

[26] Jing M Chen. Evaluation of vegetation indices and a modified simple ratio for boreal applications. *Canadian Journal of Remote Sensing*, 22(3):229–242, 1996.

[27] Hui Qing Liu and Alfredo Huete. A feedback based modification of the ndvi to minimize canopy background and atmospheric noise. *IEEE transactions on Geoscience and Remote Sensing*, 33(2):457–465, 1995.

[28] Eija Honkavaara, Heikki Saari, Jere Kaivosoja, Ilkka Pölönen, Teemu Hakala, Paula Litkey, Jussi Mäkynen, and Liisa Pesonen. Processing and assessment of spectrometric, stereoscopic imagery collected using a lightweight uav spectral camera for precision agriculture. *Remote Sensing*, 5(10): 5006–5039, 2013.

[29] Moon S Kim, CST Daughtry, EW Chappelle, JE McMurtrey, and CL Walthall. The use of high spectral resolution bands for estimating absorbed photosynthetically active radiation (a par). 1994.

[30] SB Idso, RD Jackson, PJ Pinter Jr, RJ Reginato, and JL Hatfield. Normalizing the stress-degree-day parameter for environmental variability. *Agricultural meteorology*, 24:45–55, 1981.

[31] A Ruimy, L Kergoat, Alberte Bondeau, and ThE Participants OF ThE Potsdam NpP Model Intercomparison. Comparing global models of terrestrial net primary productivity (npp): Analysis of differences in light absorption and light-use efficiency. *Global Change Biology*, 5(S1):56–64, 1999.

[32] A Bannari, D Morin, F Bonn, and AR Huete. A review of vegetation indices. *Remote sensing reviews*, 13(1-2):95–120, 1995.

[33] Gordon Christie, Neil Fendley, James Wilson, and Ryan Mukherjee. Functional map of the world. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6172–6180, 2018.

[34] Saikat Basu, Sangram Ganguly, Supratik Mukhopadhyay, Robert DiBiano, Manohar Karki, and Ramakrishna Nemani. Deepsat: a learning framework for satellite imagery. In *Proceedings of the 23rd SIGSPATIAL international conference on advances in geographic information systems*, pages 1–10, 2015.

[35] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(7):2217–2226, 2019.

[36] Ilke Demir, Krzysztof Koperski, David Lindenbaum, Guan Pang, Jing Huang, Saikat Basu, Forest Hughes, Devis Tuia, and Ramesh Raska. Deepglobe 2018: A challenge to parse the earth through satellite images. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 172–17209. IEEE, 2018.

[37] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.

[38] Otávio AB Penatti, Keiller Nogueira, and Jefersson A Dos Santos. Do deep features generalize from everyday objects to remote sensing and aerial scenes domains? In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 44–51, 2015.

[39] Michael Xie, Neal Jean, Marshall Burke, David Lobell, and Stefano Ermon. Transfer learning from deep features for remote sensing and poverty mapping. In *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.

[40] Lichen Zhou, Chuang Zhang, and Ming Wu. D-linknet: Linknet with pretrained encoder and dilated convolution for high resolution satellite imagery road extraction. In *CVPR Workshops*, pages 182–186, 2018.

[41] Zhao Zhou, Yingbin Zheng, Hao Ye, Jian Pu, and Gufei Sun. Satellite image scene classification via convnet with context aggregation. In *Pacific Rim Conference on Multimedia*, pages 329–339. Springer, 2018.

[42] Burak Uzkent, Evan Sheehan, Chenlin Meng, Zhongyi Tang, Marshall Burke, David Lobell, and Stefano Ermon. Learning to interpret satellite images in global scale using wikipedia. *arXiv preprint arXiv:1905.02506*, 2019.

[43] Neal Jean, Sherrie Wang, Anshul Samar, George Azzari, David Lobell, and Stefano Ermon. Tile2vec: Unsupervised representation learning for spatially distributed data. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3967–3974, 2019.

[44] Othman Sidek and SA Quadri. A review of data fusion models and systems. *International Journal of Image and Data Fusion*, 3(1):3–21, 2012.

[45] Jeremy Castagno and Ella Atkins. Roof shape classification from lidar and satellite image data fusion using supervised learning. *Sensors*, 18(11):3960, 2018.

[46] Yady Tatiana Solano Correa, Francesca Bovolo, and Lorenzo Bruzzone. Vhr time-series generation by prediction and fusion of multi-sensor images. In *2015 IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*, pages 3298–3301. IEEE, 2015.

[47] Michael Schmitt and Xiao Xiang Zhu. Data fusion and remote sensing: An ever-growing relationship. *IEEE Geoscience and Remote Sensing Magazine*, 4(4):6–23, 2016.

[48] Feng Gao, Jeff Masek, Matt Schwaller, and Forrest Hall. On the blending of the landsat and modis surface reflectance: Predicting daily landsat surface reflectance. *IEEE Transactions on Geoscience and Remote sensing*, 44(8):2207–2218, 2006.

[49] Paola Benedetti, Gaetano Raffaele, Osé Kenji, Ruggero Gaetano Pensa, Dupuy Stephane, Dino Ienco, et al. $M^3$fusion: Un modèle dapprentissage profond pour la fusion de données satellitaires multi-{Echelles/Modalités/Temporelles}. In *Conférence Française de Photogrammétrie et de Télédétection CFPT 2018*, pages 1–8, 2018.

[50] Martin Liggins II, David Hall, and James Llinas. *Handbook of multisensor data fusion: theory and practice*. CRC press, 2017.

[51] Barbara Anne Harrison and David Laurence Barry Jupp. *Introduction to remotely sensed data: Part one of the microBrian resource manual*. East Melbourne, Vic: CSIRO Publications, 1989.

[52] Alan R Gillespie, Anne B Kahle, and Richard E Walker. Color enhancement of highly correlated images. i. decorrelation and hsi contrast stretches. *Remote Sensing of Environment*, 20(3):209–235, 1986.

[53] WJOSEPH CARPER, THOMASM LILLESAND, and RALPHW KIEFER. The use of intensity-hue-saturation transformations for merging spot panchromatic and multispectral image data. *Photogrammetric Engineering and remote sensing*, 56(4):459–467, 1990.

[54] Thierry Ranchin and Lucien Wald. Fusion of high spatial and spectral resolution images: The arsis concept and its implementation. 2000.

[55] Jingchun Piao, Yunfan Chen, and Hyunchul Shin. A new deep learning based multi-spectral image fusion method. *Entropy*, 21(6):570, 2019.

[56] Hui Li, Xiao-Jun Wu, and Josef Kittler. Infrared and visible image fusion using a deep learning framework. In *2018 24th International Conference on Pattern Recognition (ICPR)*, pages 2705–2710. IEEE, 2018.

[57] CST Daughtry, CL Walthall, MS Kim, E Brown De Colstoun, and JE McMurtrey Iii. Estimating corn leaf chlorophyll concentration from leaf and canopy reflectance. *Remote sensing of Environment*, 74(2):229–239, 2000.

[58] Ravi P Sripada, Ronnie W Heiniger, Jeffrey G White, and

Randy Weisz. Aerial color infrared photography for determining late-season nitrogen requirements in corn. *Agronomy Journal*, 97(5):1443–1451, 2005.

[59] Robert Lawrence Pearson. Remote mapping of standing crop biomass for estimation of the productivity of the shortgrass prairie. In *Eighth International Symposium on Remote Sensing of Enviroment*, pages 1357–1381. University of Michigan, 1972.

[60] Fred Baret and Gerard Guyot. Potentials and limits of vegetation indices for lai and apar assessment. *Remote sensing of environment*, 35(2-3):161–173, 1991.

[61] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.

[62] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.

[63] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016.

[64] Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 3–19, 2018.

[65] Hyeonseob Nam and Hyo-Eun Kim. Batch-instance normalization for adaptively style-invariant neural networks. In *Advances in Neural Information Processing Systems*, pages 2558–2567, 2018.

[66] Mingxing Tan and Quoc V Le. Efficientnet: Rethinking model scaling for convolutional neural networks. *arXiv preprint arXiv:1905.11946*, 2019.

[67] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[68] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.

[69] Carole H Sudre, Wenqi Li, Tom Vercauteren, Sebastien Ourselin, and M Jorge Cardoso. Generalised dice overlap as a deep learning loss function for highly unbalanced segmentations. In *Deep learning in medical image analysis and multimodal learning for clinical decision support*, pages 240–248. Springer, 2017.

[70] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[71] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.