

# Transferability of labels between multilens cameras

Ignacio de Loyola Páez-Ubieta<sup>1</sup><sup>a</sup>, Daniel Frau-Alfaro<sup>1</sup><sup>b</sup> and Santiago T. Puente<sup>1</sup><sup>c</sup>

<sup>1</sup>*AUtomatics, Robotics, and Artificial Vision (AUROVA) Lab, University Institute for Computer Research (IUII), University of Alicante, Crta. San Vicente s/n, San Vicente del Raspeig, E-03690, Alicante, Spain*  
{*ignacio.paez, daniel.frau, santiago.puente*}@ua.es

Keywords: Multispectral Imagery, Labeling, Phase Correlation, Label Transfer, Pills

Abstract: In this work, a new method for automatically extending Bounding Box (BB) and mask labels across different channels on multilens cameras is presented. For that purpose, the proposed method combines the well known phase correlation method with a refinement process. During the first step, images are aligned by localizing the peak of intensity obtained in the spatial domain after performing the cross correlation process in the frequency domain. The second step consists of obtaining the best possible transformation by using an iterative process maximising the IoU (Intersection over Union) metric. Results show that, by using this method, labels could be transferred across different lens on a camera with an accuracy over 90% in most cases and just by using 65 ms in the whole process. Once the transformations are obtained, artificial RGB images are generated, for labeling them so as to transfer this information into each of the other lens. This work will allow users to use this type of cameras in more fields rather than satellite or medical imagery, giving the chance of labeling even invisible objects in the visible spectrum.

## 1 INTRODUCTION

When training a detection (Wang et al., 2023) or segmentation (Wang et al., 2020) Neural Network (NN), a large amount of data is required to adapt an already trained model to a specific task, such as household waste (Páez-Ubieta et al., 2023) (Páez-Ubieta et al., 2023).

However, times have changed and automatic labeling models for objects in RGB images appeared recently, with Segment Anything Model (SAM) model (Kirillov et al., 2023) becoming a reference in a very short period of time.


MultiSpectral Imaging (MSI) consists of using sensors that provide images in different frequency ranges compared to the traditional RGB cameras. In several fields such as in agriculture (Hoffer et al., 1966) (Mia et al., 2023), medical (Andersson et al., 1987) (Ma et al., 2023) or remote sensing (Maxwell, 1976) (Yuan et al., 2021) these type of cameras have given promising results during the last two centuries.


However, labeling images outside the RGB domain could be a difficult task. However, little by little, more and more articles dealing with the labeling pro-


cess on Light Detection And Ranging (LiDAR) (Ošep et al., 2024) or multilens (Gallagher et al., 2024) images are being published.

For instance, (Gómez and Meoni, 2021) introduced a semisupervised learning approach to automatically classify scenes from land datasets such as EuroSAT (Helber et al., 2019) or the aerial UC Merced land use (UCM) (Yang and Newsam, 2010). For that purpose, they label between 5 and 300 images per class, which are then feed into a Graphics Processing Unit (GPU) for training a model. In our case, we do not need any kind of training phase since we directly obtain the transformation between the camera lenses in order to have a more detailed object recognition rather than just a scene classification. Also, 15 images were used during the transformation phase, but less images could also be used for the proposed method.

Another example is (Ulku et al., 2022), in which authors aim to segmentate semantically trees using satellite and aerial images from the DSTL Satellite Imagery Feature Detection Image (Benjamin et al., 2016) and RIT-18 (The Hamlin State Beach Park) Aerial Image (Kemker et al., 2018) datasets. To this end, they use several segmentation NNs to perform the task of labeling trees on the images. In our case, we do not require any kind of semantic segmentation

<sup>a</sup> <https://orcid.org/0000-0001-9901-7264>

<sup>b</sup> <https://orcid.org/0009-0000-4098-3783>

<sup>c</sup> <https://orcid.org/0000-0002-6175-600X>

NN to label our images. Also, their trees occupy big areas on the images, making it easier for the NN to find and label them, meanwhile the objects in our case are far smaller - tougher to label as they require far more precision and attention to fine detail.

Other works such as (Park et al., 2021), use a multispectral and RGB cameras for detecting sick pine trees through aerial photographs. For performing the image alignment for later labeling process they use the Scale-Invariant Feature Transform (SIFT) method (Lowe, 1999). However, this aforementioned method only works well when keypoints and descriptors could be obtained from the images, being useless with uniform objects. Also, RGB and 6 channel multispectral images are analysed by the NN, making the process not very efficient since some of the 9 channels could contain no info at all.

In this work, the transformation between the images captured from a multispectral camera (which is a type of multilens camera) will be obtained with the purpose of extending the Bounding Box (BB) or mask labels from one image into the others, as well as allowing the users to label in fake RGB images, saving a significant amount of time. The process consists of a two step process (displacement calculation and refinement) that uses traditional computer vision techniques and Central Processing Unit (CPU) resources, leaving aside time and resource consumption on GPUs.

The main contributions of this work are:

- A new method for obtaining the transformation between the lens of a multispectral camera that is proven to be highly accurate in no time.
- The possibility of generating fake RGB images from combining its components by applying the previous transform.
- Transforming labels in both BB and mask formats across images so as to label objects that disappear in certain frequencies.

This work is organised as follows: Section 2 introduces the proposed method, which is divided in two steps, Section 3 presents the setup that was used during the experiments, as well as the transformations between each lenses and the fake RGB labeling process and Section 4 summarises the article and introduces future works that will be done using as core project this method.

## 2 METHODOLOGY

In this Section, the method for obtaining the transformation between different lens on the camera is

presented. It is composed of two steps: displacement calculation, using the phase correlation method, and refinement, using a sliding window across several scales.

### 2.1 Displacement calculation

The different lens on a camera provide images that are not aligned. As the lenses are at the same height, a 2D transformation (rotation, translation, scale and/or skew) is the most probable conversion to relate them. However, the assumption that there is only a displacement between the captured images by the different lenses was done. If results proved otherwise, other kind of transformation would be applied.

The obtained images are in the space domain, in which each pixel represents the intensity. However, we switch into the frequency domain, in which images are reorganised by frequencies, distributing them according to its periodicity (high periodicity will be placed in the center of the image, meanwhile low periodicity will be placed far from it).

By getting advance of how images are distributed in the frequency domain, the displacement between two images is a linear phase shift, which is the core idea of the phase correlation algorithm.

It receives as input two images  $i_1$  and  $i_2$ . The first step is removing sharp discontinuities at the image borders, since they produce a high frequency component, reducing the accuracy of the method. This problem is called spectral leakage, but applying a Hanning window (Eq. 1) will make it disappear, smoothing the image and removing artifacts and edges.

$$w(x, y) = \left( 0.5 \left( 1 - \cos \left( \frac{2\pi x}{M-1} \right) \right) \right) \cdot \left( 0.5 \left( 1 - \cos \left( \frac{2\pi y}{N-1} \right) \right) \right) \quad (1)$$

where  $M$  and  $N$  represent the dimensions of the image and  $x$  and  $y$  represent the pixel coordinates. If applied to the aforementioned images,  $i_{1h}(x, y)$  and  $i_{2h}(x, y)$  are obtained. The second step consists of in transforming these spectral leakage free images into the frequency domain by using the Discrete Fourier Transform (DFT) (Eq. 2), obtaining  $I_{1h}(u, v)$  and  $I_{2h}(u, v)$ , respectively.

$$I_{1h}(u, v) = \sum_{x=0}^{M-1} \sum_{y=0}^{N-1} i_{1h}(x, y) \cdot e^{-2\pi i \left( \frac{ux}{M} + \frac{vy}{N} \right)} \quad (2)$$

Once the images are in the frequency domain, the phase shift between them encode the translational shift we are looking for in the space domain. For that

purpose, the third step is to isolate this phase information by using the cross-power spectrum (Eq. 3), normalising the magnitude and retaining this phase information.

$$CP(u, v) = \frac{I_{1h}(u, v) \cdot I_{2h}^*(u, v)}{|I_{1h}(u, v) \cdot I_{2h}^*(u, v)|} \quad (3)$$

where  $I_{2h}^*(u, v)$  is the complex conjugate of  $I_{2h}(u, v)$ . The fourth step is coming back to the spatial domain by using the Inverse Discrete Fourier Transform (IDFT) applied to the calculated cross-power spectrum  $CP(u, v)$ , obtaining the correlation matrix  $c(x, y)$  (Eq. 4).

$$c(x, y) = \sum_{u=0}^{M-1} \sum_{v=0}^{N-1} CP(u, v) \cdot e^{2\pi i \left( \frac{ux}{M} + \frac{vy}{N} \right)} \quad (4)$$

Finally, peak location  $(\Delta x, \Delta y)$  in the correlation matrix  $c(x, y)$  is obtained (Eq. 5) by performing a  $5 \times 5$  weighted centroid around the peak, so as to obtain subpixel accuracy, normalizing the result between 0 and 1.

$$(\Delta x, \Delta y) = \text{weightedCentroid} \left\{ \arg \max_{(x, y)} \{c(x, y)\} \right\} \quad (5)$$

## 2.2 Refinement

Once the relative displacement  $(\Delta x, \Delta y)$  between the two input images  $i_1, i_2$  is obtained using the phase correlation method, a refinement process is required in order to get more precision.

For that purpose, centered in the relative displacement numbers, a cascade of possible better displacements are obtained. First,  $(\Delta x, \Delta y)$  are rounded and are added or subtracted a value  $RV : i \in 1 \dots n$  with  $n$  being the number of refinement steps in both axis  $x$  and  $y$  in different scales  $s$ . This  $s$  value will represent different orders of magnitude, varying it inside the discrete values on  $[1, 0.1, 0.01]$  to check pixel and subpixel precision (Eq. 6), having as a result several possible combinations.

$$\begin{aligned} \Delta x_p &= [\Delta x - RV \cdot s, \dots, \Delta x, \dots, \Delta x + RV \cdot s] \\ \Delta y_p &= [\Delta y - RV \cdot s, \dots, \Delta y, \dots, \Delta y + RV \cdot s] \end{aligned} \quad (6)$$

These obtained possible numbers  $(\Delta x_p, \Delta y_p)$  are inserted into a homogeneous transformation (Eq. 7) and applied to each of the different labels on the image to check if a better solution is obtained. For that purpose, some labeled images serve as basis to compare with these newly obtained labels.

$$\begin{bmatrix} l_{N:nx} \\ l_{N:ny} \end{bmatrix} = \begin{bmatrix} 1 & 0 & \Delta x_p \\ 0 & 1 & \Delta y_p \end{bmatrix} \begin{bmatrix} l_{nx} \\ l_{ny} \end{bmatrix} \quad (7)$$

These labels  $l = [(l_{1x}, l_{1y}), \dots, (l_{nx}, l_{ny}), \dots, (l_{Nx}, l_{Ny})]$  are a collection of  $N$  points that represent the borders of the labeled objects. In case of having a mask,  $N$  could be any positive number, meanwhile a BB is defined by  $N = 2$ , representing the top left and bottom right coordinates of the box.

The final transform will convert the original mask or BB coordinates  $(l_{nx}, l_{ny})$  into the new reference frame, obtaining  $l_M = [(l_{M:1x}, l_{M:1y}), \dots, (l_{M:nx}, l_{M:ny}), \dots, (l_{M:Nx}, l_{M:Ny})]$ . They will be compared against the aforementioned labeled images  $l_{GT} = [(l_{GT:1x}, l_{GT:1y}), \dots, (l_{GT:nx}, l_{GT:ny}), \dots, (l_{GT:Nx}, l_{GT:Ny})]$ , trying to achieve the highest Intersection Over Union (IOU).

IOU, also called Jaccard index, is a metric that returns how much two labels coincide, being the result a value between  $[0, 1]$ . It is represented by Eq. 8.

$$IoU = \frac{l_M \cap l_{GT}}{l_M \cup l_{GT}} \quad (8)$$

## 3 EXPERIMENTATION

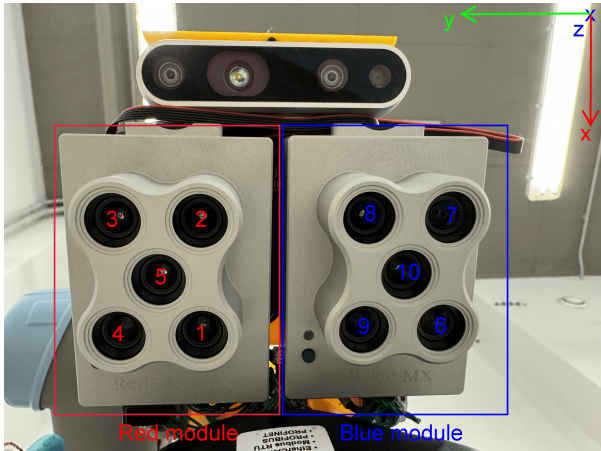
In this Section, hardware and software setup is presented, as well as the experiments that prove the effectiveness of the proposed method for obtaining labels across multiple multispectral images.

### 3.1 Setup

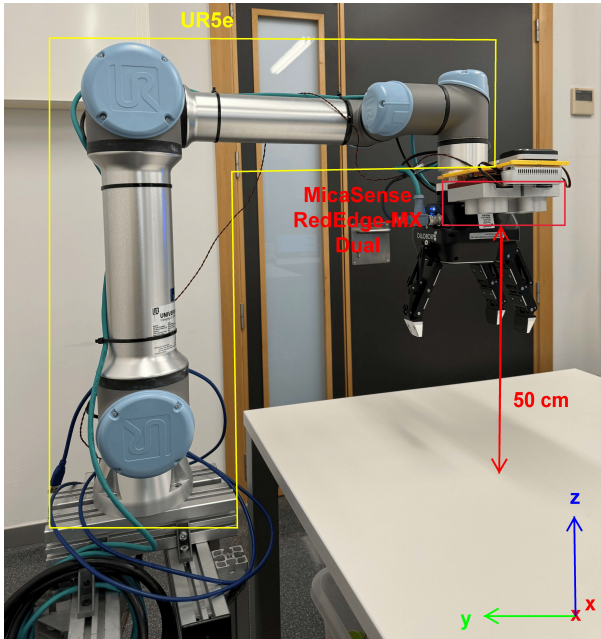
In order to perform image alignment, several instances are required.

In terms of hardware (Fig. 1), a MicaSense RedEdge-MX Dual multispectral camera is used for obtaining the aforementioned images. This camera counts with 10 cameras divided in two modules: a red and a blue one (see Fig. 1a). A brief description of the frequencies for each band available in both of them is presented in Table 1. All 10 bands produce 12 bit images at a resolution of 1280 x 960. This camera is mounted on the wrist of an Ur5e with 6 Degrees of Freedom (DoF) robotic arm, which allow us to precisely place the camera in a still position on the space. Concretely, the camera is positioned parallel to a table 500 mm over it (see Fig. 1b). Regarding the objects used for performing the experiments, 16 small pills in the size range of 8 to 22 mm are used. They provide

several different shapes and colors to allow all the experiments that will be performed.



(a) Detail of the multispectral camera lenses (modules and order).



(b) Robotic arm in which the MicaSense RedEdge-MX Dual is attached to, allowing to capture objects 50 cm away.

Figure 1: Hardware used during the experiments.

In terms of software, images were labeled using LabelMe tool (Torralba et al., 2010). There are two different labeling approximations used across all the experiments: BB and mask. The computer used to get the results works in Ubuntu 20.04.4 with Python 3.8.10 and OpenCV 4.7.0, running an 11th Generation Intel® Core™ i9-11900H with 8 physical and 16 logical cores respectively. They work at a frequency

Table 1: Band numbers, frequencies and color name from each channel available on the MicaSense RedEdge-MX Dual multispectral camera (see Fig. 1a).

Module	Band	$f \pm A$ (nm)	Color name
Red	1	$475 \pm 16$	Blue
	2	$560 \pm 13.5$	Green
	3	$668 \pm 7$	Red
	4	$717 \pm 6$	Red Edge
	5	$842 \pm 28.5$	Near IR
Blue	6	$444 \pm 14$	Coastal Blue
	7	$531 \pm 7$	Green
	8	$650 \pm 8$	Red
	9	$705 \pm 5$	Red Edge I
	10	$740 \pm 9$	Red Edge II

of 2.50 GHz, which allow to carry out all necessary operations in a short period of time.

### 3.2 Transformations

Although the camera has 10 lenses, only the ones in the red half part of the camera will be used.

Several images from the 5 lenses that compose the red part of the multispectral camera are obtained. Concretely, 15 images were captured with each camera, from which 12 were used for obtaining the transform and 3 for proving the accuracy of the obtained transform. From them, the training images from band 5 (lenses in the middle) and the test images from all 5 lenses were labeled. Also, the refinement steps  $n$  was set to 5, meaning 121 possible matrices in 3 different levels of pixel accuracy are used.

After applying the phase correlation method and the refinement for both BB and mask labels, the following results shown in Table 2 and Table 3 were obtained.

Table 2: Transformations in pixel level, IOU and time for BB labels to transfer band 5 labels into the other lenses.

Band	Transform (px)	IoU (%)	Time (ms)
1	${}^1_5T_{BB} = \begin{pmatrix} -52.0 \\ 47.0 \end{pmatrix}$	98.58	66.31
2	${}^2_5T_{BB} = \begin{pmatrix} 53.9 \\ 46.1 \end{pmatrix}$	100	53.52
3	${}^3_5T_{BB} = \begin{pmatrix} 52.9 \\ -23.4 \end{pmatrix}$	95.95	55.66
4	${}^4_5T_{BB} = \begin{pmatrix} -52.1 \\ -18.9 \end{pmatrix}$	93.53	56.55

In order to show the active refinement process and how it improves the IOU step by step, band 1 from mask labeling is shown in Table 4. Step 0 consists of applying phase correlation, step 1 consists of refining in pixel level with  $s$  equal to 1, step 2 consists of re-

Table 3: Transformations in pixel level, IOU and time for mask labels to transfer band 5 labels into the other lenses.

Band	Transform (px)	IoU (%)	Time (ms)
1	$\frac{1}{5}T_{MK} = \begin{pmatrix} -52.05 \\ 47.2 \end{pmatrix}$	97.49	82.81
2	$\frac{2}{5}T_{MK} = \begin{pmatrix} 54.78 \\ 46.52 \end{pmatrix}$	94.36	72.58
3	$\frac{3}{5}T_{MK} = \begin{pmatrix} 53.5 \\ -23.8 \end{pmatrix}$	93.77	73.05
4	$\frac{4}{5}T_{MK} = \begin{pmatrix} -53.24 \\ -19.01 \end{pmatrix}$	89.91	58.83

fining in subpixel level with  $s$  equal to 0.1 and step 3 consists of refining in 2 levels of subpixel with  $s$  equal to 0.01.

Table 4: Phase correlation and refinement steps applied to band 1 of mask labeled images.

Step	Transform (px)	IoU (%)	Time (ms)
0	$\frac{1}{5}T_{MK-0} = \begin{pmatrix} -51.85 \\ 47.02 \end{pmatrix}$	96.73	76.75
1	$\frac{1}{5}T_{MK-1} = \begin{pmatrix} -52.0 \\ 47.0 \end{pmatrix}$	96.99	1.96
2	$\frac{1}{5}T_{MK-2} = \begin{pmatrix} -52.0 \\ 47.2 \end{pmatrix}$	97.40	2.05
3	$\frac{1}{5}T_{MK} = \begin{pmatrix} -52.05 \\ 47.2 \end{pmatrix}$	97.49	2.05

Several images of BB labeled images after applying  $\frac{1}{5}T_{BB}$ ,  $\frac{2}{5}T_{BB}$ ,  $\frac{3}{5}T_{BB}$  and  $\frac{4}{5}T_{BB}$  from Table 2 can be seen in Fig. 2. As input Fig. 2a is provided, which represents band 5. Once the labels of this image are transformed, images in Figs. 2b, 2c, 2d and 2e are obtained. As it can be seen, the transformed labels adjust almost perfectly into the objects of the other bands, saving the user the need of manually annotate them. For a quick comparison of the quality, ground truth human labeled images are provided in Figs. 2f, 2g, 2h and 2i. The worst result is obtained in band 4, since the pill starts mimicking with the background, making it difficult for both the proposed method and the user to discern it.

Moving onto a more difficult problem, mask labeled images are tested. These images are transformed using matrices  $\frac{1}{5}T_{MK}$ ,  $\frac{2}{5}T_{MK}$ ,  $\frac{3}{5}T_{MK}$  and  $\frac{4}{5}T_{MK}$  from Table 3. Results of the process are shown in Fig. 3. The layout follows the same pattern as in Fig. 2. Taking that into account, the worst result is again band 4, being the reason the same as in the previous case. The pill starts to disappear (compared to the other 3 bands). However, the labeling process is still successful.

### 3.3 RGB label transferability

Once the labeling process is shown to be successful, a new experiment will be performed - labeling in RGB images created from the multispectral images and then move these labels into the other bands.

For that purpose, bands representing the red, green and blue frequencies need to be mixed together. According to Table 1, bands 1-3 from the red module are the ones required. So the images are converted from bands 1-3 to band 5 by using Eq. 9, obtaining an artificial RGB image  $im_{RGB} = [R, G, B]$ .

$$\begin{aligned} R &= \frac{3}{5} T_{BB|mask}^{-1} \cdot im_{band3} \\ G &= \frac{2}{5} T_{BB|mask}^{-1} \cdot im_{band2} \\ B &= \frac{1}{5} T_{BB|mask}^{-1} \cdot im_{band1} \end{aligned} \quad (9)$$

Then, this fake RGB image is labeled in BB or mask format  $l_{RGB} = [(l_{RGB:1x}, l_{RGB:1y}), \dots, (l_{RGB:nx}, l_{RGB:ny}), \dots, (l_{RGB:Nx}, l_{RGB:Ny})]$  by the user. Once the label is ready, it is transferred back into the other bands in the camera, obtaining labels in all frequencies (Eq. 10).

$$\begin{aligned} l_{band1} &= \frac{1}{5} T_{BB|mask} \cdot l_{RGB} \\ l_{band2} &= \frac{2}{5} T_{BB|mask} \cdot l_{RGB} \\ l_{band3} &= \frac{3}{5} T_{BB|mask} \cdot l_{RGB} \\ l_{band4} &= \frac{4}{5} T_{BB|mask} \cdot l_{RGB} \\ l_{band5} &= I_{2 \times 3} \cdot l_{RGB} \end{aligned} \quad (10)$$

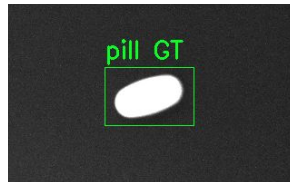
The first process for generating the  $im_{RGB} = [R, G, B]$  can be seen in Fig. 4. For it, a combination of blue (Fig. 4a), green (Fig. 4b) and red (Fig. 4c) images is performed, generating a fake RGB image (Fig. 4d).

Once the fake RGB image is generated, objects have been labeled in BB (Fig. 5a) and mask (Fig. 6a) formats. After applying the transformations obtained from Section 3.2, the images with BB (Figs. 5b-5f) and masks (Figs. 6b-6f) were obtained. As it can be seen, obtained results convincingly demonstrate the accuracy of the proposed approach. For example, objects in band 4 with BB (Fig. 5e) as well as with mask (Fig. 6e) labeling disappear partially. However, due to the designed method, objects in those positions are labeled even though they are not visible.

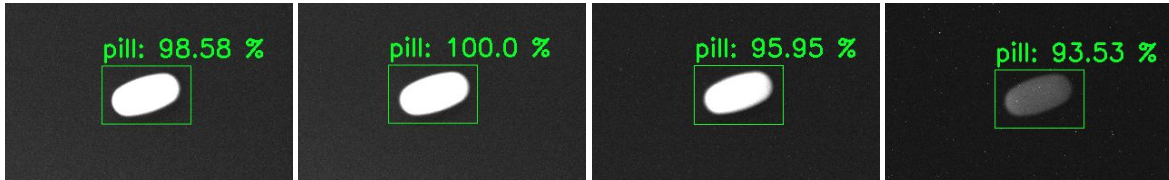
## 4 CONCLUSION

In this work, a method for labeling automatically multispectral images starting with a single band image BB or mask labeled is presented.





(a) Band 5 input.

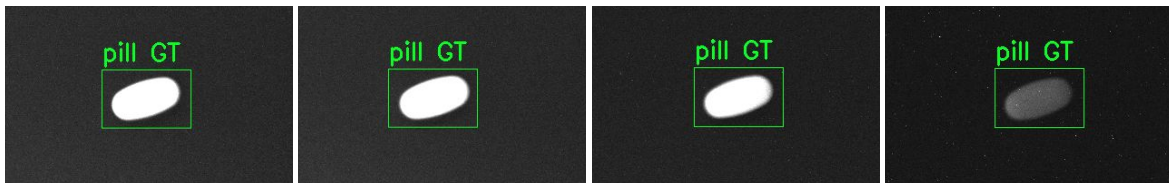


(b) Band 1 output.

(c) Band 2 output.

(d) Band 3 output.

(e) Band 4 output.



(f) Band 1 reference.

(g) Band 2 reference.

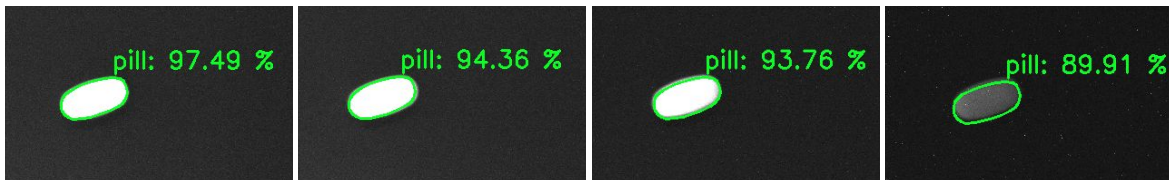
(h) Band 3 reference.

(i) Band 4 reference.

Figure 2: BB labeled experiment: (a) reference image (band 5) to start with, (b, c, d, e) transformed labels (bands 1-4) and (f, g, h, i) ground truth labels for comparison purposes (bands 1-4).



(a) Band 5 input.

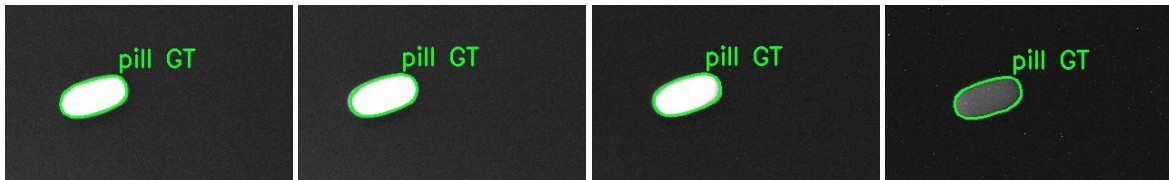


(b) Band 1 output.

(c) Band 2 output.

(d) Band 3 output.

(e) Band 4 output.



(f) Band 1 reference.

(g) Band 2 reference.

(h) Band 3 reference.

(i) Band 4 reference.

Figure 3: Mask labeled experiment: (a) reference image (band 5) to start with, (b, c, d, e) transformed labels (bands 1-4) and (f, g, h, i) ground truth labels for comparison purposes (bands 1-4).

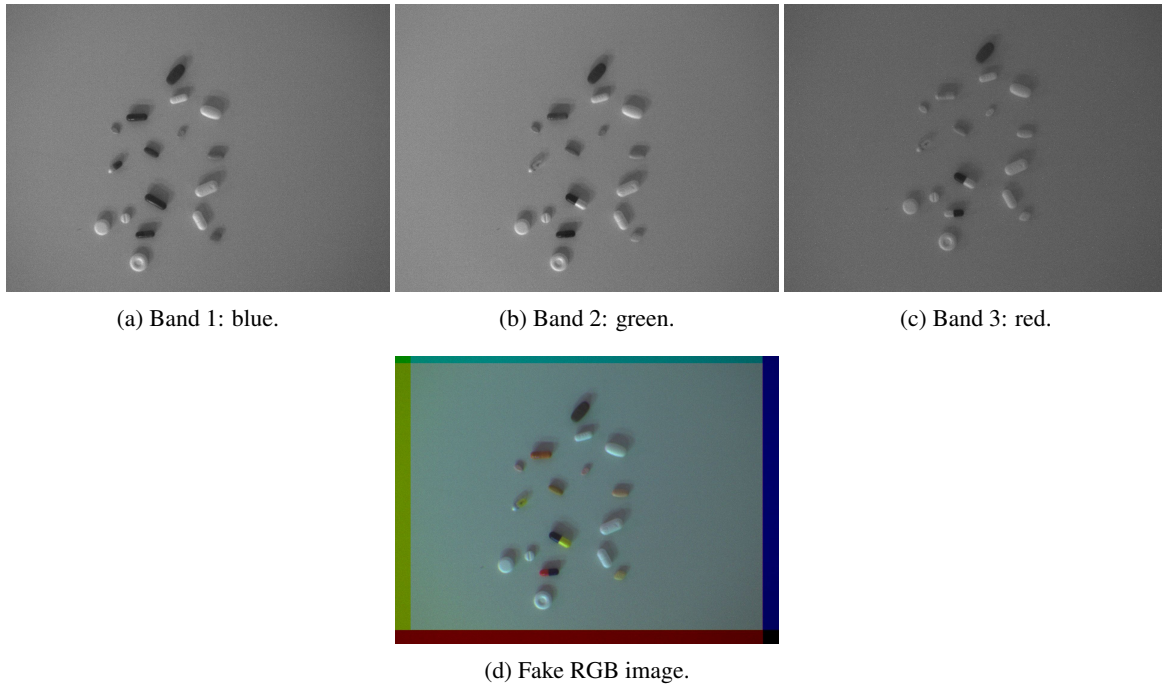


Figure 4: Combination of 3 channels to create fake RGB image: (a,b,c) bands 1, 2, 3 respectively, (d) Generated fake RGB image.

In order to achieve it, a process with two steps is used: phase correlation and refinement. During the first one, the transformation between two images is obtained by smoothing the image with a Hanning window, transforming the spatial domain images into the frequency domain with the Fourier discrete transform, applying the cross-power spectrum formula to retain just the phase information of the images, converting the cross-power spectrum back to the spatial domain, so as to finally look for the peak location, which is no other than the translation between the two analysed images. The second process consists of refining the transformation obtained from the previous step by searching in a proximity window for a better one throughout an iterative process through pixel and two levels of subpixel accuracy, saving the best transformation as the one that provided the highest percentage in the IOU index.

In order to test the method, the transformation between 5 multispectral lenses from a MicaSense RedEdge-MX Dual camera were obtained. Just by labeling 12 images from band 5 with a high contrast allowed to obtain the transformation of BB and mask label types with an accuracy of 97% and 94% and just 58 ms and 72 ms, respectively. Right after that, by using the inverse of the obtained transformations, an artificial RGB image is generated, allowing the labeling process to be performed in colored images. After it,

the labels are transformed back into each lenses so as to have the labels in all 5 channels of the multispectral camera.

Future works will consist of testing the proposed method into more multispectral cameras with different morphology, as well as testing it with all the 10 lenses that the camera used in the present paper has. Also, a RGB camera could be added in order to avoid generating fake RGB images from the multispectral lenses and accumulating a small error during the process. In a different path, a dataset of domestic waste could be created aiming to train different NNs and test if the extra information provided by 10 lenses and 12 bit images could help discerning better between categories when compared to the same NNs using 8 bit RGB images.

## ACKNOWLEDGEMENTS

Research work was funded by grant PID2021-122685OB-I00 funded by MICIU/AEI/10.13039/501100011033 and ERDF/EU.

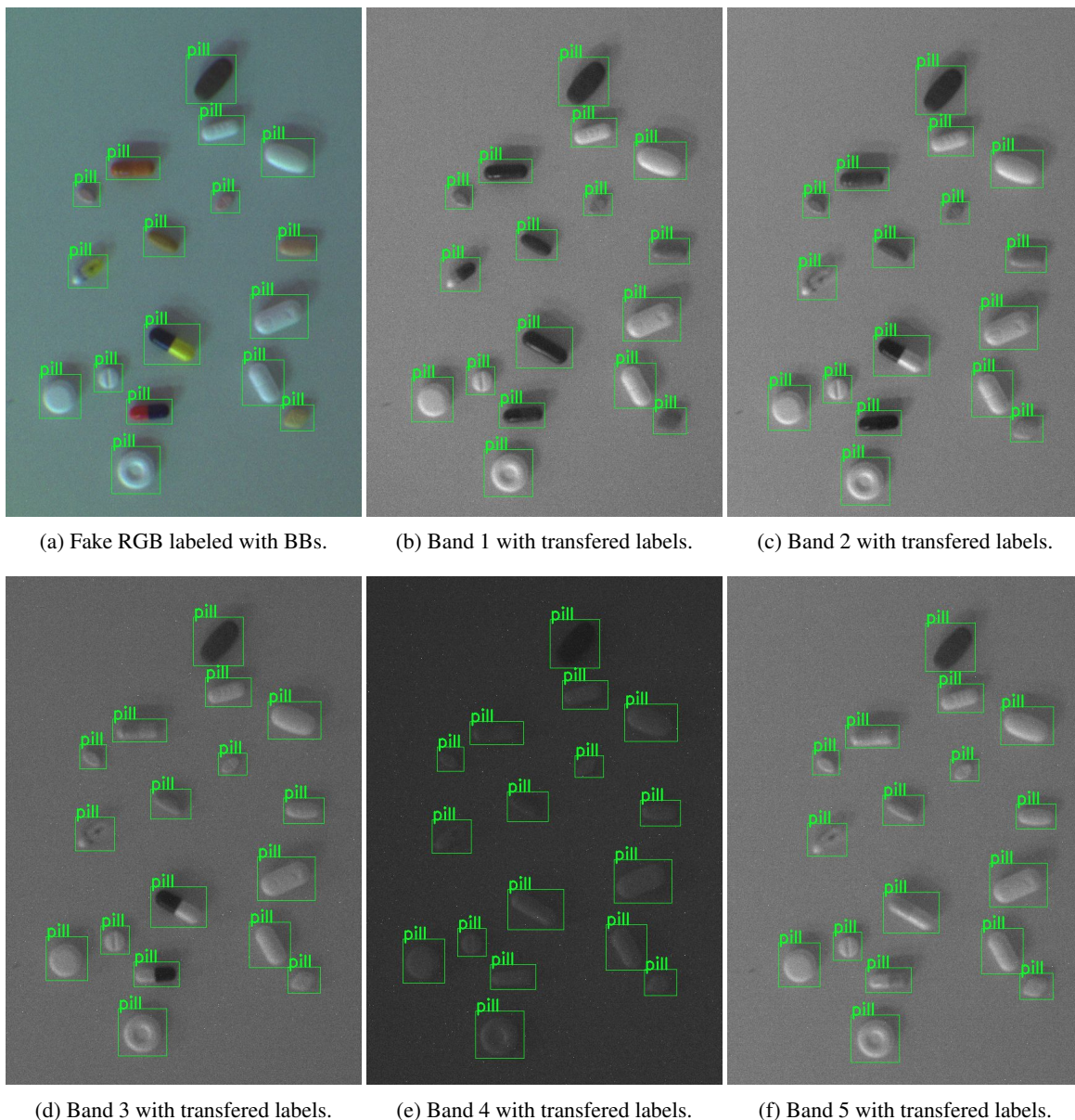
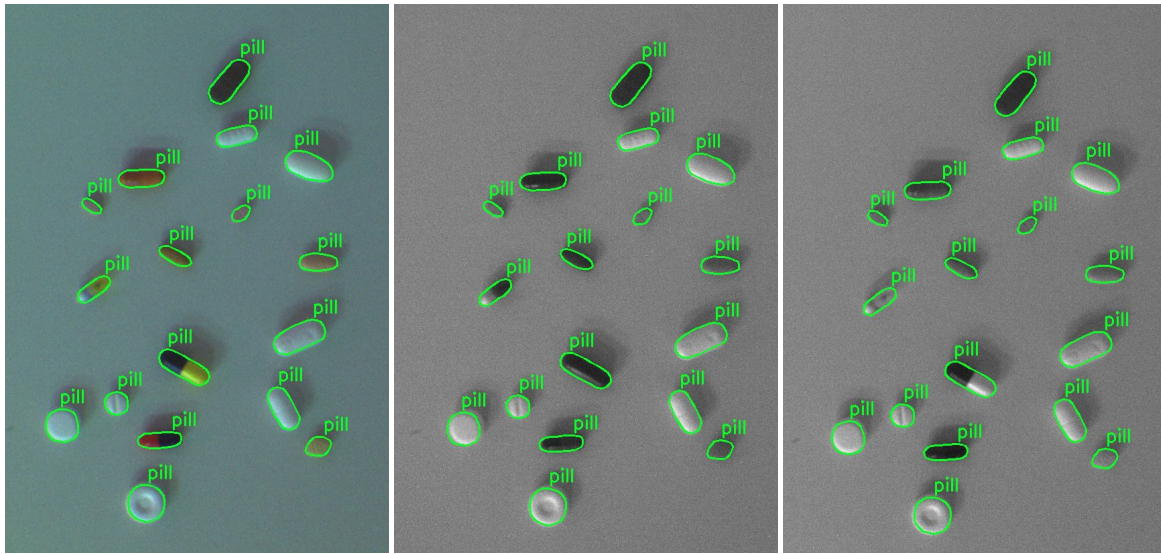


Figure 5: BB labeled fake RGB image and transferred labels: (a) Fake RGB image with BB labels, (b) labels transferred to band 1, (c) labels transferred to band 2, (d) labels transferred to band 3, (e) labels transferred to band 4 and (f) labels transferred to band 5.

## REFERENCES

- Andersson, P., Montan, S., and Svanberg, S. (1987). Multispectral system for medical fluorescence imaging. *IEEE Journal of Quantum Electronics*, 23(10):1798–1805.
- Benjamin, MatvL, midaha, PGibson, RMcKinlay, and Kan, W. (2016). Dstl satellite imagery feature detection.
- Gallagher, J. E., Gogia, A., and Oughton, E. J. (2024). A multispectral automated transfer technique (matt) for machine-driven image labeling utilizing the segment anything model (sam). *arXiv preprint*.
- Gómez, P. and Meoni, G. (2021). Msmatch: Semisupervised multispectral scene classification with few labels. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 14:11643–11654.
- Helber, P., Bischke, B., Dengel, A., and Borth, D. (2019). Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observa-*

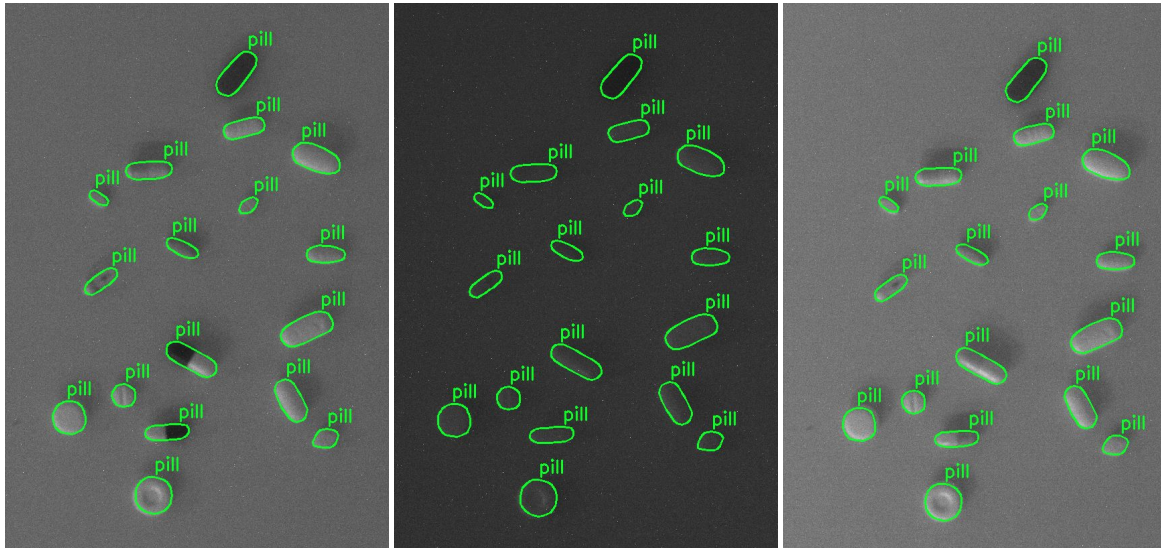




(a) Fake RGB labeled with masks.

(b) Band 1 with transfered labels.

(c) Band 2 with transfered labels.



(d) Band 3 with transfered labels.

(e) Band 4 with transfered labels.

(f) Band 5 with transfered labels.

Figure 6: Mask labeled fake RGB image and transfered labels: (a) Fake RGB image with mask labels, (b) labels transfered to band 1, (c) labels transfered to band 2, (d) labels transfered to band 3, (e) labels transfered to band 4 and (f) labels transfered to band 5.

*tions and Remote Sensing*, 12(7):2217–2226.

Hoffer, R., Johannsen, C., and Baumgardner, M. (1966). Agricultural applications of remote multispectral sensing. In *Proceedings of the Indiana Academy of Science*, volume 76, pages 386–396.

Kemker, R., Salvaggio, C., and Kanan, C. (2018). Algorithms for semantic segmentation of multispectral remote sensing imagery using deep learning. *ISPRS Journal of Photogrammetry and Remote Sensing*.

Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A. C.,

Lo, W.-Y., et al. (2023). Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026.

Lowe, D. G. (1999). Object recognition from local scale-invariant features. In *Proceedings of the seventh IEEE international conference on computer vision*, volume 2, pages 1150–1157. IEEE.

Ma, F., Yuan, M., and Kozak, I. (2023). Multispectral imaging: Review of current applications. *Survey of Ophthalmology*, 68(5):889–904.

Maxwell, E. L. (1976). Multivariate system analysis of mul-

- tispectral imagery. *Photogrammetric Engineering and Remote Sensing*, 42(9):1173–1186.
- Mia, M. S., Tanabe, R., Habibi, L. N., Hashimoto, N., Homma, K., Maki, M., Matsui, T., and Tanaka, T. S. T. (2023). Multimodal deep learning for rice yield prediction using uav-based multispectral imagery and weather data. *Remote Sensing*, 15(10).
- Ošep, A., Meinhardt, T., Ferroni, F., Peri, N., Ramanan, D., and Leal-Taixé, L. (2024). Better call sal: Towards learning to segment anything in lidar. *arXiv preprint*.
- Páez-Ubieta, I. d. L., Castaño-Amorós, J., Puente, S. T., and Gil, P. (2023). Vision and tactile robotic system to grasp litter in outdoor environments. *Journal of Intelligent & Robotic Systems*, 109(2):36.
- Park, H. G., Yun, J. P., Kim, M. Y., and Jeong, S. H. (2021). Multichannel object detection for detecting suspected trees with pine wilt disease using multispectral drone imagery. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 14:8350–8358.
- Páez-Ubieta, I. d. L., Velasco-Sánchez, E., Puente, S. T., and Candelas, F. (2023). Detection and depth estimation for domestic waste in outdoor environments by sensors fusion. *IFAC-PapersOnLine*, 56(2):9276–9281. 22nd IFAC World Congress.
- Torralba, A., Russell, B. C., and Yuen, J. (2010). Labelme: Online image annotation and applications. *Proceedings of the IEEE*, 98(8):1467–1484.
- Ulku, I., Akagündüz, E., and Ghamisi, P. (2022). Deep semantic segmentation of trees using multispectral images. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 15:7589–7604.
- Wang, C.-Y., Bochkovskiy, A., and Liao, H.-Y. M. (2023). Yolov7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7464–7475. IEEE.
- Wang, X., Zhang, R., Kong, T., Li, L., and Shen, C. (2020). Solov2: Dynamic and fast instance segmentation. *Advances in Neural information processing systems*, 33:17721–17732.
- Yang, Y. and Newsam, S. (2010). Bag-of-visual-words and spatial extensions for land-use classification. In *Proceedings of the 18th SIGSPATIAL international conference on advances in geographic information systems*, pages 270–279.
- Yuan, K., Zhuang, X., Schaefer, G., Feng, J., Guan, L., and Fang, H. (2021). Deep-learning-based multispectral satellite image segmentation for water body detection. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 14:7422–7434.