

Visual Mesh: Real-time Object Detection Using Constant Sample Density

Trent Houliston^[0000-0002-7744-0472] and
Stephan K. Chalup^[0000-0002-7886-3653]

School of Electrical Engineering and Computing
The University of Newcastle, Callaghan, NSW, 2308, Australia.
trent@houliston.me, stephan.chalup@newcastle.edu.au

Abstract. This paper proposes an enhancement of convolutional neural networks for object detection in resource-constrained robotics through a geometric input transformation called Visual Mesh. It uses object geometry to create a graph in vision space, reducing computational complexity by normalizing the pixel and feature density of objects. The experiments compare the Visual Mesh with several other fast convolutional neural networks. The results demonstrate execution times sixteen times quicker than the fastest competitor tested, while achieving outstanding accuracy.

Keywords: Convolutional Neural Network · Deep Learning · Ball Detection · Graph Transformation · TensorFlow · Machine Vision

1 Introduction

This paper introduces a Visual Mesh that defines an input transformation for convolutional neural networks (CNN). By normalizing object size, the Visual Mesh accounts for differences in an object's appearance when detecting and localizing it. This allows simpler network architectures to be used and reduces oversampling, improving the computational performance substantially.

CNNs require powerful hardware to perform in real-time. Despite this, some networks have been developed to run on constrained hardware with limited success. Speck et al. [13] built a CNN for detecting the coordinates of a soccer ball on an image. When implemented on their target platform it ran in 26 ms and had an accuracy of 58 % in x and 52 % in y . The accuracy dropped to less than 30 % in distances over two meters. Therefore, this approach had limited success in object localization.

Faster and more accurate systems have been developed that only perform object classification. These systems utilize color segmentation to provide proposals for a CNN to classify. As a result they were much faster than systems that localize objects, however, color segmentation is sensitive to changes in lighting conditions and must be manually calibrated. Javadi et al. [7] utilized such a system for detecting humanoid robots. The best performing network ran in 2.36 ms with 97.56 % accuracy per proposal on an Intel Core i5 2.5 GHz. Cruz et al. [5] developed a system to classify Aldebaran NAO robots. This network

executed in ≈ 1 ms per proposal. Albani et al. [2] and Bloisi et al. [3] utilized a similar technique for ball detection. This system was implemented on an Aldebaran NAO robot and processed 14-22 frames per second as the only process running. The reliance on color segmentation for proposals limits these networks to color coded environments.

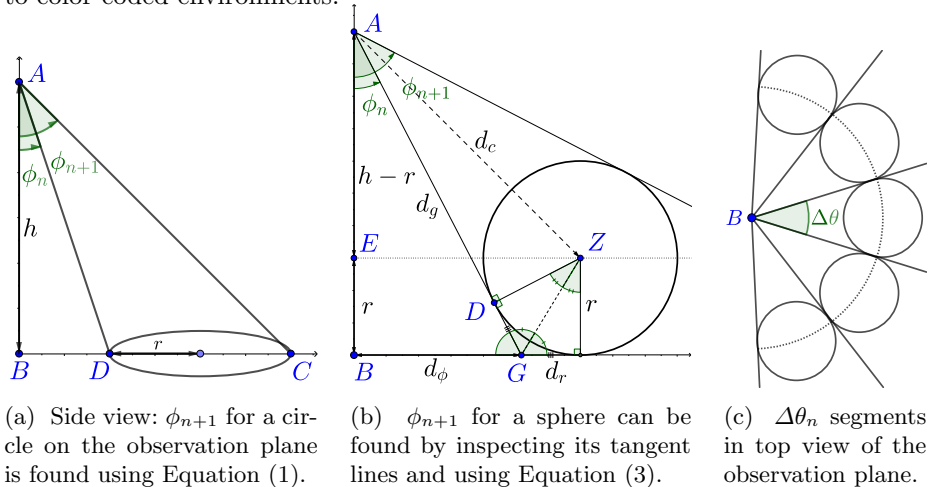


Fig. 1: Geometry for calculating ϕ_{n+1} and $\Delta\theta_n$.

2 Visual Mesh Geometry

The Visual Mesh detects objects that lie on a plane at a known distance and orientation from the camera. This plane is referred to as the observation plane. The geometry of the Visual Mesh can be described using Figure 1 where the camera is assumed to be at point A . The target object's geometry determines the pixel resolution, i.e., the placement of points in the Visual Mesh. The geometry for two target object shapes are analyzed in this paper: Circles are appropriate for detecting two-dimensional objects on the observation plane (Figure 1a). Spheres are appropriate for three-dimensional objects that have an approximately equal extension in all dimensions (Figure 1b). More complex objects such as cylinders could also be modeled.

For establishing the Visual Mesh two orthogonal angular components have to be determined. These are $\Delta\phi_n$ and $\Delta\theta_n$ and are given by the angular diameters of the target object with respect to points A and B . The height h of the camera above the observation plane and the radius r are required to calculate the mesh.

The first component is $\Delta\phi_n := \phi_{n+1} - \phi_n$ and is determined by the inclinations ϕ_n from directly below the camera. A series ϕ_n , $n = 0, \dots, N$ is given recursively by function $f : \mathbb{R} \rightarrow \mathbb{R}$, $\phi_{n+1} = f(\phi_n) = \phi_n + \Delta\phi_n$ where $\phi_0 = 0$.

The second component, $\Delta\theta_n$, is measured around point B in the observation plane and depends on ϕ_n for both, circle and sphere objects (Figure 1c).

The inclinations $(\phi_n)_{n=0,\dots,N}$ induce a series of nested concentric cones with vertex at A and center axis orthogonal to the observation plane. Each of these cones is radially segmented at its basis by $\Delta\theta_n$ and the tangent rays from B .

2.1 Circle

The geometry for circles is shown in Figure 1a. ϕ_{n+1} for a circle is calculated by adding the diameter $2r$ of the circle to its distance \overline{BD} to obtain

$$\phi_{n+1} = \tan^{-1} \left(\tan(\phi_n) + \frac{2r}{h} \right) \quad (1)$$

Figure 1c shows the geometry for $\Delta\theta_n$ within the 2D observation plane where

$$\Delta\theta_n = 2 \sin^{-1} \left(\frac{r}{h \tan(\phi_n)} \right) \quad (2)$$

This formulation of $\Delta\theta_n$ has a singularity when the center of the object is closer than its radius making it more difficult for the mesh to detect objects directly below the camera.

2.2 Sphere

For spheres $\Delta\phi_n$ is determined by the sphere's shadow from a virtual light at A and it decreases more slowly with n than for circles. Figure 1b shows how ϕ_{n+1} is calculated. Using the triangle $\triangle AEZ$ and edges \overline{AE} and \overline{EZ} gives

$$\begin{aligned} \phi_{n+1} &= 2 \tan^{-1} \left(\frac{d_\phi + d_r}{h - r} \right) - \phi_n \\ &= 2 \tan^{-1} \left(\frac{r \sec(\phi_n)}{h - r} + \tan(\phi_n) \right) - \phi_n \end{aligned} \quad (3)$$

The calculation of $\Delta\theta_n$ is the same as for circles and uses Equation (2).

2.3 Object Dependent Sample Density

The current description guarantees one point in the mesh for the target object. For use in computer vision, multiple sample points per object are required. Let's assume our object requires k pixels to be recognizable. In the Visual Mesh this k corresponds to the number of intersections of the ϕ_n rings with the object. A ϕ_n ring is obtained by rotating vector \overrightarrow{AD} about the axis \overrightarrow{AB} . If $\Delta\phi_n$ and $\Delta\theta_n$ are reduced, the spacing between the ϕ_n rings will be decreased which leads to more intersections with the target object.

An increase in the number of sample points for the circle model can be achieved by dividing $\Delta\theta_n$ in (2) by k and also the diameter of the circle by k , i.e., replacing $2r$ in (1) by $2r/k$.

In the sphere model k sample points can be achieved by creating a version of the mesh where the original sphere is replaced by smaller spheres so that the original sphere intersects with k ϕ_n rings associated with the smaller spheres (Figure 2). If k is expressed as fraction $k = \frac{p}{q}$, $p, q \in \mathbb{N} - \{0\}$, the equation relating the radii of the spheres is given by $f^q(\phi_0, r_0) = f^{pq}(\phi_0, r_1)$ where r_0 is the radius of the target and r_1 is the radius of the small spheres in the mesh. A solution for r_1 can be obtained numerically.

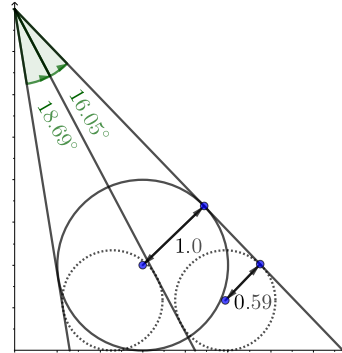


Fig. 2: Multiple sample points on a sphere can be calculated by finding the smaller sphere’s radius.

2.4 Graph Structure of the Mesh

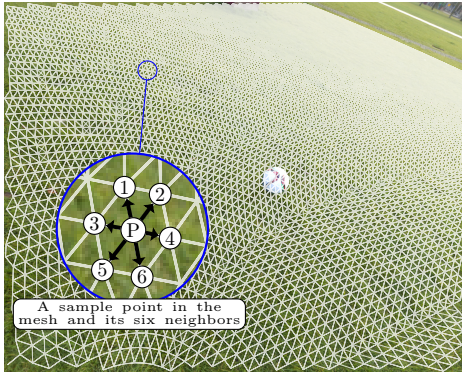


Fig. 3: The Visual Mesh projected onto an image. Note that four ϕ_n rings pass through the ball regardless of its location.

in distance within the original image, while the Visual Mesh ensures objects are always the same size.

2.5 Network

Once the image data has been transformed by the Visual Mesh, it exists as a graph, rather than a grid of pixels. The pixels no longer have nine neighbors, but six. This changes how convolutions occur when executed on the graph.

For example, a 3×3 convolution in a typical CNN accesses eight pixels around a central pixel. The equivalent operation in the graph accesses points

A mesh can be generated using the points around ϕ_n rings (see arcs in Figure 3). In each ϕ_n ring, points are separated by $\Delta\theta_n$. This ensures that the number of points within an object falls within a small range (± 1 in ϕ and θ). Each point is connected with edges to the two adjacent points on the same ϕ_n ring as well as to the two nearest points on the $\phi_{n\pm 1}$ rings. The single point below the camera is connected by six equally spaced points. Projecting these points onto an image creates a mesh structure as shown in Figure 3.

Another method to view the mesh is to project the ϕ_n rings onto concentric circles as in Figure 4. Due to perspective, the size of objects decreases

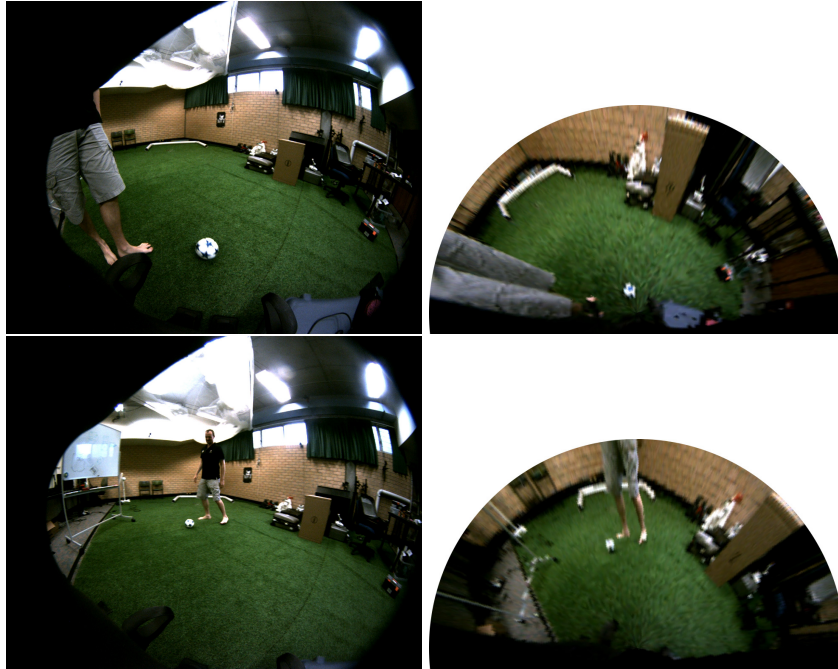


Fig. 4: The Visual Mesh projected in concentric rings. Due to perspective, the size of the ball decreases with distance in the original (left). In the Visual Mesh, the ball has always a similar size (right).

with a graph distance of one, its six neighbors. This has a positive impact on performance, as two fewer values need to be considered. Larger convolutions would be equivalent to operating on points at a larger maximum graph distance. For example, a 5×5 convolution would operate on all points that have a maximum graph distance of two to the central point.

3 Evaluation of the Visual Mesh

3.1 Dataset

A semi-synthetic dataset with masks that segment out the ball was created for training. By using 360° high dynamic range (HDR) images to provide image-based lighting, along with physics-based rendering, realistic semi-synthetic scenes were generated. From this, the mask images, as well as the camera orientation and position can be obtained.

Using a number of different HDR scenes taken from RoboCup 2017, the NUbots' laboratory and online¹, as well as over a hundred different soccer ball

¹ HDRI Haven <https://hdrihaven.com/>

designs, over 160,000 images were generated. These soccer ball designs were not limited to 50% white as per RoboCup rules and included balls of various colors. Each of these images varied the position of the soccer ball and switched between equisolid and rectilinear camera projections.

The distance of the balls from the camera varied between zero and ten meters. The intensity of the lighting varied in the scene. The rendered soccer ball was selected from a set of 140 different models. The distribution of distances was designed to provide a uniform variation in the pixel size of the ball. This allowed a consistent variation in the angular diameter of the ball in the image. It prevented a large number of visually small balls that would have occurred with a uniform distribution over distance.

3.2 Network Architecture

Each node in the Visual Mesh performed a 3×3 convolution using its six neighbors. These layers were stacked to varying depths from two to nine and with output widths varying from two to eight, resulting in a fully convolutional net.

Networks of width four performed significantly better than networks of other widths as the hardware utilized can vectorize on four elements. The results discussed in Section 3.4 only include network widths of four.

Networks were also tested with ReLU [9], ELU [4] and SELU [8]. SELU consistently outperformed ELU and ReLU in terms of training time and network accuracy. SELU is computationally more expensive than ReLU but is similar to that of ELU. Results in Section 3.4 only include those tested with SELU.

The network depths used for evaluation were three, five, and nine layers. These were chosen as their receptive fields were half, one and two ball radii, respectively. This ensured the networks had sufficient contextual information to correctly classify the ball.

In addition to these Visual Mesh networks, similar CNNs using a regular hexagonal grid were trained. These networks allowed a comparison between the Visual Mesh and a network that has equal computational cost due to selecting the same number of pixels. This network provides a comparison to an equivalent network without the constant sample density of the Visual Mesh.

3.3 Training

The training of these networks was undertaken using the TensorFlow library [1]. The pixel coordinates from the Visual Mesh and the indices of each pixel’s six neighbors were used to apply the Visual Mesh at each layer. Once this gather step was performed, the neural network steps were undertaken as normal.

When training these neural networks, the number of ball points and non-ball points were balanced. This was achieved by selecting an equal number of points from each class. The backpropagation gradients were only calculated from the selected points.

This method was chosen instead of the traditional method of weighting the gradients intentionally. The majority of non-ball points in training images are grass. As a result, the initial networks experienced over-fitting on the field.

Once the initial network was trained, the error in its classification of each point in the image was used as a probability to select that point. This resulted in fewer grass points selected in future training. This resampling was run twice, with the probabilities added together with a 5% baseline probability. This greatly improved the accuracy in subsequent training.

In addition to these networks, five convolutional network architectures were fine-tuned on this dataset. These networks were SSD MobileNet and RCNN Inception V2 trained using TensorFlow [1, 6] and YOLOv1 [10], YOLOv2 [11] and YOLOv3 [12] trained using Darknet². These networks were chosen as they were regarded as some of the fastest real-time networks.

3.4 Results

Precision In addition to the Visual and hexagonal meshes, five typical CNNs were also evaluated. Their results were measured using a 75% IoU. 75% was chosen as 50% was considered a poor match. With 50% IoU, the center of the detection can be at the edge of the object.

As shown in Figure 5a, the accuracy of the Visual Mesh consistently outperforms the hexagonal mesh of an equivalent size. Increasing the depth of the network increases its performance.

The performance of the Visual Mesh remains approximately constant as distance increases. However, as shown in Figure 5c the performance of the hexagonal mesh, as well as the other CNNs degrades with increased distance. Note that as the generated data was made uniform over pixel size rather than distance, the number of sample images falls off as distance increases. The fewer samples increase noise in the plot.

The performance of the Visual Mesh exceeds the performance of the hexagonal mesh even when the number of points in the object is the same. The number of points in both tested networks are equal at 2.5 m. Figure 5b shows the number of points in the Visual Mesh stays constant over distance, except for a peak at 0 m. This peak is when points are directly below the camera. This is a singularity point for the Visual Mesh as it is currently implemented. The hexagonal mesh has a decreasing number of points as distance increases.

Detections Figure 6a shows a typical set of detections from each of the trained networks. YOLOv1 is omitted as it performs strictly worse than YOLOv2. The Visual Mesh has a good detection while the hexagonal mesh has several false positives. The five other networks all detect the ball. YOLOv2 has a lower confidence than the other networks on the dataset. SSD MobileNet, YOLOv2 and YOLOv3's bounding boxes are less accurate across the dataset.

² Darknet <http://pjreddie.com/darknet/>

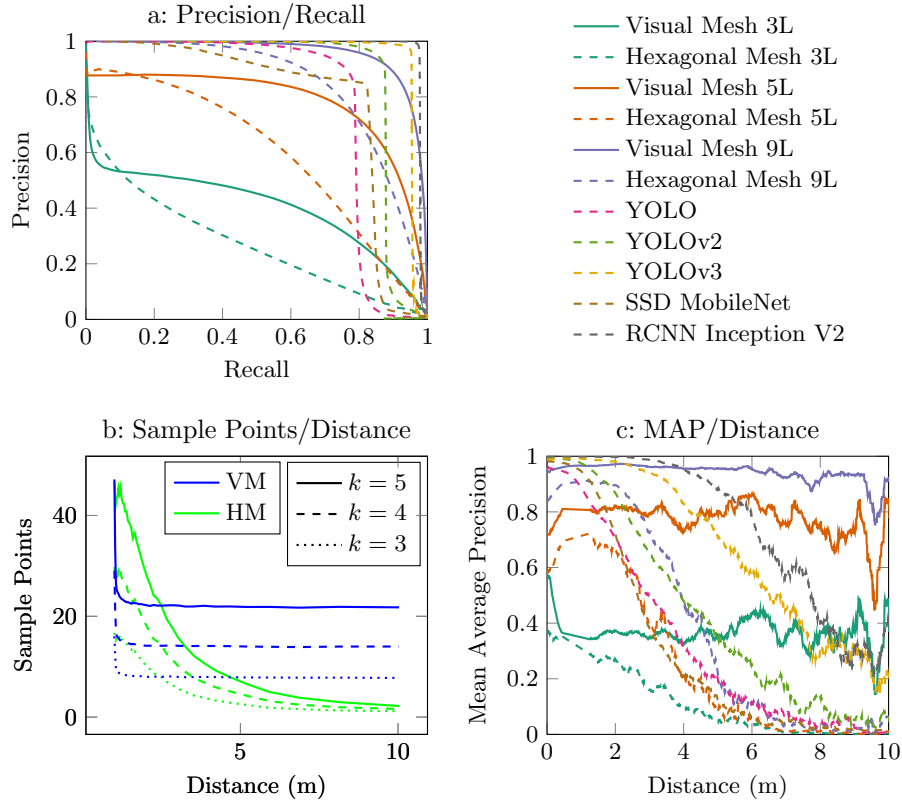


Fig. 5: (a) The Precision/Recall curve over all data. (b) The number of sampled points in an object over distance (VM is Visual Mesh HM is Hexagonal Mesh). (c) The Average Precision of the detectors over distance.

The Visual Mesh excels at distant detections as shown in Figure 6b. Except for the Visual Mesh and RCNN Inception V2, none of the other networks detect the ball. RCNN Inception V2 has a poorly fitted bounding box. This is typical of distant balls in the dataset.

Figure 6c shows how the Visual Mesh is able to use scale to identify target objects. The hexagonal mesh found many false positives on objects that had a different size than expected, but similar appearance as the target.

Execution performance Each network was tested on the CPU and GPU from the Intel NUC7i7BNH as well as on an NVIDIA 1080Ti. The input images were 1280×1024 for all networks. SSD MobileNet and RCNN Inception V2 were not

executed on the Intel GPU as TensorFlow does not support OpenCL at this time. YOLOv3 was not executed as it is not supported by the OpenCL version of Darknet. Table 1 summarizes the results with respect to execution time. The times for all networks are measured from when the image is first sent to the algorithm until the inferences are returned. Therefore, the time taken to project Visual Mesh points is included.

Table 1: Execution performance: For the Visual Mesh on the Iris Plus Graphics and the NVIDIA 1080Ti the device utilization was 70 % and 35 % respectively. For all other cases utilization was at 100 %.

	Intel Core i7 7567U	Intel Iris Plus Graphics 650	NVIDIA 1080Ti
Visual Mesh 5	1.64ms	2.10ms	2.18ms
Visual Mesh 9	2.44ms	2.48ms	2.25ms
YOLOv1	1468.24ms	721.13ms	17.55ms
YOLOv2	1221.49ms	613.73ms	16.13ms
YOLOv3	2651.33ms	N/A	19.00ms
SSD MobileNet	37.76ms	N/A	11.32ms
RCNN Inception V2	1521.32ms	N/A	47.75ms

3.5 Discussion and Conclusion

The results for the Visual Mesh show that consistent feature density improves the accuracy of the network. When the Visual Mesh and the hexagonal mesh had an equal number of points on the ball the Visual Mesh was more accurate. As distance increased, the accuracy of the hexagonal mesh degraded while the Visual Mesh remained consistent. This degradation can also be seen in other networks as accuracy declines over distance.

The nine-layer Visual Mesh is used for the following comparisons. It provided the highest accuracy and its computational performance was not significantly worse than the five-layer Visual Mesh.

RCNN Inception V2 and YOLOv3 performed the best of the other networks tested. While the other CNNs failed to detect distant objects, these networks continued to detect them. However, the bounding boxes became increasingly inaccurate. At a lower IoU threshold they have a higher detection rate. Visual Mesh exceeds their performance after 4 m.

As seen in Table 1, the execution performance of the Visual Mesh exceeds that of the other convolutional networks. Of these networks, only SSD MobileNet and the Visual Mesh could be considered for real-time use on resource-constrained systems. The performance of the Visual Mesh is fast enough that the transfer times of images is a significant factor for GPU based computation. The NUC’s CPU outperforms its GPU for the five-layer Visual Mesh because of this. The NVIDIA 1080Ti also suffers this effect, resulting in only 35 % utilization.

The Visual Mesh has a number of advantages beyond its accuracy and speed. As objects always have a similar number of points additional post-processing options are available to improve accuracy. Within detected areas, metrics such

as graph diameter can be used to filter out irrelevant areas. Additionally the best fitting subgraph can be used to remove invalid points in a detection.

Higher resolution does not increase the computational cost of the Visual Mesh. The number of points that are projected onto the image does not change for the same camera lens and orientation. However, increased resolution will allow the Visual Mesh to project points that are a greater distance from the camera. If the resolution of the camera is insufficient for the level of detail requested, the Visual Mesh will begin sampling the same pixel multiple times. The Visual Mesh is still accurate with limited amounts of this duplicated data. However, as the amount of information decreases, the accuracy of the network will decline.

As the distance to objects increases, typical networks must learn to account for the differences in scale that occur. Often these differences are not well represented in the training data or, can be biased in the training data. This can require additional training data to be generated by scaling. For the Visual Mesh, this is not necessary as the object will always appear the same size. This reduces the complexity of training as well as the complexity of the required network.

As the Visual Mesh is always oriented relative to the observation plane, the resulting network is better able to handle changes in the orientation of the camera. This form of transformational invariance only applies to rotations in the camera, not rotations on the object. This can reduce the amount of training required if the object can be assumed in a particular rotation. For example, if extended to detect the goal posts the training data would not need to be modified for different orientations of the camera as they would always be normal to the observation plane. Without this invariance, training for goal posts would have to include multiple orientations.

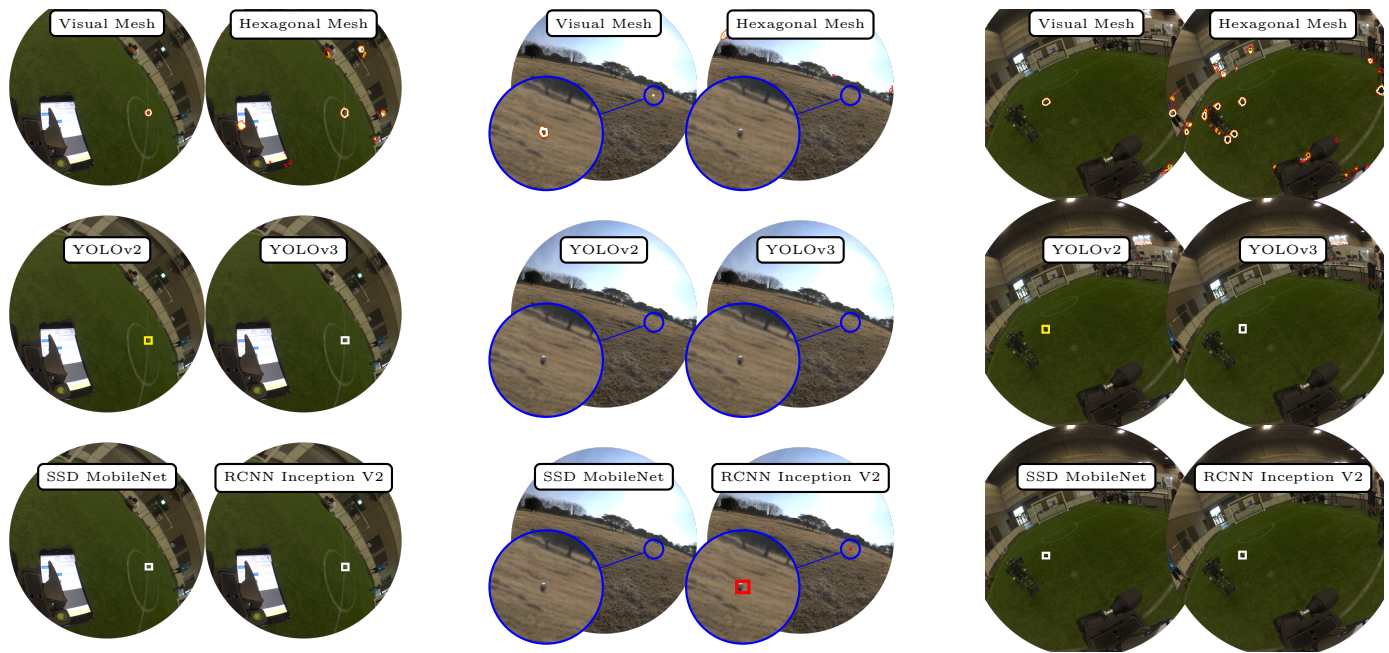
Networks based on the Visual Mesh also have an independence to the lens used. As the sample points are always in the same place in the world, changing to a different lens geometry does not change the points. This makes it easier to train using data from different lenses and apply trained networks to new lenses.

The presented formulation of the Visual Mesh has two primary limitations. One concern is that it cannot function when the height of objects are greater than or equal to the height of the camera. In these cases, the Visual Mesh correctly predicts that all objects are visible on the horizon. This results in a single line of points. In practice, this does not afford good detection performance. The second limitation is that objects that are directly below the camera fall into a singularity. When in this singularity, twice as many points intersect with the objects until they move beyond this position. This increases the complexity that the Visual Mesh must learn.

The training and execution code for the Visual Mesh is available at <https://github.com/Fastcode/VisualMesh>.

Acknowledgments

TH was supported by a Australian Government Research Training Program scholarship and a completion scholarship from 4Tel Pty. Ltd.



(a) Detection on an empty field. All networks achieve acceptable performance. However, the hexagonal mesh has some false positives.

(b) Detections at extreme distance. Only the Visual Mesh achieves a good detection. RCNN Inception V2 detects but IoU of the bounding box is $< 50\%$.

(c) Detection on a field with a robot. The hexagonal mesh has many false positives on the robot, while the other detectors perform well.

Fig. 6: Example detections where confidence intervals are represented by colors: Red $> 50\%$, Yellow $> 75\%$ and White $> 90\%$

Bibliography

- [1] Abadi, M., et al.: Tensorflow: A system for large-scale machine learning. In: Proceedings of the 12th USENIX Conference on Operating Systems Design and Implementation. pp. 265–283. OSDI’16, USENIX Association, Berkeley, CA, USA (2016)
- [2] Albani, D., Youssef, A., Suriani, V., Nardi, D., Bloisi, D.D.: A deep learning approach for object recognition with NAO soccer robots. In: Behnke, S., Sheh, R., Sarel, S., Lee, D.D. (eds.) RoboCup 2016: Robot World Cup XX. pp. 392–403. Springer International Publishing, Cham (2017)
- [3] Bloisi, D., Duchetto, F.D., Manoni, T., Suriani, V.: Machine learning for realistic ball detection in robocup SPL. CoRR abs/1707.03628 (2017)
- [4] Clevert, D.A., Unterthiner, T., Hochreiter, S.: Fast and accurate deep network learning by exponential linear units (elus). CoRR abs/1511.07289 (2015)
- [5] Cruz, N., Lobos-Tsunekawa, K., Ruiz-del Solar, J.: Using convolutional neural networks in robots with limited computational resources: Detecting NAO robots while playing soccer. CoRR abs/1706.06702 (2017)
- [6] Huang, J., Rathod, V., Sun, C., Zhu, M., Korattikara, A., Fathi, A., Fischer, I., Wojna, Z., Song, Y., Guadarrama, S., Murphy, K.: Speed/accuracy trade-offs for modern convolutional object detectors. CoRR abs/1611.10012 (2016)
- [7] Javadi, M., Azar, S.M., Azami, S., Shiry, S., Ghidary, S.S., Baltes, J.: Humanoid robot detection using deep learning: A speed-accuracy tradeoff. In: RoboCup 2017: Robot World Cup XXI. Springer International Publishing (2018), in press
- [8] Klambauer, G., Unterthiner, T., Mayr, A., Hochreiter, S.: Self-normalizing neural networks. In: Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (eds.) Advances in Neural Information Processing Systems 30, pp. 971–980. Curran Associates, Inc. (2017)
- [9] Nair, V., Hinton, G.E.: Rectified linear units improve restricted boltzmann machines. In: Proceedings of the 27th International Conference on International Conference on Machine Learning. pp. 807–814. ICML’10, Omnipress, USA (2010)
- [10] Redmon, J., Divvala, S.K., Girshick, R.B., Farhadi, A.: You only look once: Unified, real-time object detection. CoRR abs/1506.02640 (2015)
- [11] Redmon, J., Farhadi, A.: YOLO9000: Better, faster, stronger. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 6517–6525 (2017)
- [12] Redmon, J., Farhadi, A.: YOLOv3: An incremental improvement. CoRR abs/1804.02767 (2018)
- [13] Speck, D., Barros, P., Weber, C., Wermter, S.: Ball localization for robocup soccer using convolutional neural networks. In: Behnke, S., Sheh, R., Sarel, S., Lee, D.D. (eds.) RoboCup 2016: Robot World Cup XX. pp. 19–30. Springer International Publishing, Cham (2017)