

Contributors: M. Tzelepi, E. Kakaletsis, V. Noussi, A. Tefas, N. Nikolaidis (Aristotle University of Thessaloniki, Greece)

> Presenter: Nikos Nikolaidis Aristotle University of Thessaloniki nikolaid@aiia.csd.auth.gr <u>www.multidrone.eu</u> Presentation version 1.0

This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 731667 (MULTIDRONE)





- Safety is an extremely important aspect in drone operations:
  - Protecting the **pilot** and **other people** on the ground
  - Protecting **property** (buildings, cars)
  - Protecting the **drone** or other drones





- Drone safety has many aspects:
  - Obstacle detection and avoidance
  - Crowd or individual persons detection so as not to fly near or over them (illegal)
  - Drone to drone collision avoidance
  - Drone hardware and communications reliability, fault tolerance
  - Early warning for potentially dangerous faults and situations





- Assisting the pilot during flight or in emergency situations so as to reduce safety risks due to e.g. human errors:
  - Autonomous or semi-autonomous (assisted) flight modes
  - Emergency landing site detection
  - Automated emergency landing procedures
- Imaging, computer vision, artificial intelligence can have a crucial role in drone safety.





- Human crowd detection for safe autonomous drones
- Visual drone detection for collision avoidance
- Emergency landing site detection.





- Aim: detect and localize crowds in drone videos.
- Provide binary crowd maps or heatmaps depicting the probability of crowd presence or the crowd density in each image location.
- **Challenging task**: occlusions, viewing angle variations, small size of persons, complex background





- Why crowd detection:
  - To **prevent accidents** caused by e.g. a drone falling while flying over a crowd
  - To enrich the **flight map** with **no-fly** zones towards crowd avoidance
  - To detect/confirm **emergency landing sites** 
    - Such sites shall not be occupied by crowds
  - To comply with **legislation** 
    - No flights are allowed over or near crowd gathering locations





- There are limited efforts on crowd detection per se, especially in aerial videos.
- Research works involving crowds, e.g., crowd dynamics understanding, crowd counting, human detection and tracking in crowds, consider scenes where the crowd has already been detected.

Zhang, Yingying, et al. "*Single-image crowd counting via multi-column convolutional neural network.*", Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016.

Sam, Deepak Babu, Shiv Surya, and R. Venkatesh Babu. "*Switching convolutional neural network for crowd counting.*", Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Vol.1, No.3, 2017.



- In [1] a real-time method that detects moving crowds based on spatio-temporal analysis of a video sequence is proposed
- The method looks at the crowd motion patterns in the spatio-temporal domain.
- Each person moving in a crowd forms a line in the spatiotemporal domain
- A moving crowd generates multiple lines that intersect with each other, since people move in opposite directions or with different speed
- These intersecting lines are used as a cue for detecting a moving crowd

[1] Reisman, Pini, et al. "Crowd detection in video sequences." Intelligent Vehicles Symposium, 2004 IEEE. IEEE, 2004.



- The proposed system was developed as part of a driver assistance system and aims to detect pedestrian crowds crossing the road from a camera mounted on a vehicle
- It includes additional modules to distinguish between crowds and moving vehicles or single pedestrians
  - Higher level classifiers are used to determine if the object is a vehicle or pedestrian, and optic-flow analysis is used to determine if the moving region corresponds to a crowd
- Its applicability to drone-captured videos is questionable



- In [2] a method that segments an image into crowded and non-crowded regions is proposed.
- **Core idea:** to capture two key attributes of a crowd:
  - on a narrow scale, its basic element should weakly look like a human
  - on a larger scale, a crowd contains repetitive appearances of elements (humans)
  - These two properties are captured by a feature vector evaluated over a window

[2] Arandjelovic, Ognjen. "Crowd detection from still images." *BMVC 2008: Proceedings of the British machine vision association conference 2008.* BMVA Press, 2008.



- The method builds for each image location a pyramid of windows of various dimensions
- Local features which correspond to a sparse set of interest points are evaluated on each window
- Due to the nature of crowd regions a large number of interest points are generated in such regions





- SIFT descriptors are used to describe each point's neighborhood
- These descriptors are then quantized into SIFT words, estimated by K-means clustering
- The method employs a statistical Poisson model of occurrences of quantized SIFT words to quantify how "crowd-like" the content around the examined location is
- A properly trained **SVM** is used to classify locations into crowd and non-crowd ones







- In [3] the authors address the problem of crowd detection in aerial images by using texture-classification methods.
- Two texture-classification methodologies are applied:
  - a Bag-of-Words (BoW) model with Improved Fisher Vectors
  - features based on a Gabor filter bank
  - Since the core methodology has been originally designed for the task of texture classification it is implied that the crowd appearance is texture–like

[3] Meynberg, Oliver, Shiyong Cui, and Peter Reinartz. "Detection of high-density crowds in aerial images using texture classification.", Remote Sensing, 2016



- 1. Bag-of-Words (BoW) model with Improved Fisher Vectors
- The image is split into **patches**
- Two types of local texture features are extracted: "Local Binary Patterns" (LBP) and "Sorted Random Projections" (SRP)
- Codewords are calculated using a Gaussian Mixture Model (GMM)
- Features are encoded using Improved Fisher Vectors (IFV)
- Final result: a histogram-like feature vector for each patch



- 2. Crowd Features Using a Gabor Filter Bank
  - The image patches are convolved with a filter bank and each patch is eventually represented by one feature vector
  - A Gabor filter encodes the orientation and scale of edges of the input image
  - A high filter response in a specific orientation and of a specific scale is received if the input image contains edges of buildings or other regular structures
  - An image patch containing a crowd exhibits heterogeneous texture and generates a high filter response in every direction
  - Gabor filter responses can thus be used to separate crowd from non-crowd regions

input patch

filter response







medium-dense crowd sp

sparse crowd

- An SVM classifier is used to classify feature vectors of the image patches
- Two problems have been examined :
  - Classify patches into **crowd** and **non-crowd** ones (2 classes)
  - Classify patches into 4 classes: dense crowd, medium dense crowd, sparse crowd, no-crowd
    Image: Image:





This project has received and innovation pro;

class 3—sparse crowd

class 2—medium dense crowd



class 4- no crow

- The experimental results show that a classifier using either BoW or Gabor features can detect crowded image regions in the 2-class problem with 97% classification accuracy
- Distinguishing between the 4 crowd density classes proved to be much more difficult.
- The BoW approach has 5%–12% better accuracy than Gabor in this case



Figure 11. Multi-class classification with BoW or Gabor features. The four classes are visualized as an overlay over a typical aerial image. Color code: dense crowd (red), medium dense crowd (yellow), sparse crowd (green), no crowd (not colored). (a) original image; (b) manually labeled image; (c) BoW features; (d) Gabor features.

- In [4] a Fully Convolutional Model for crowd detection in drone videos is proposed
- The output of the method is a heatmap denoting the estimated probability of crowd existence in each frame location



[4] M. Tzelepi, A. Tefas. "Human Crowd Detection for Drone Flight Safety Using Convolutional Neural Networks." in European Signal Processing Conference (EUSIPCO), Kos, Greece, 2017.





- The authors utilized the BVLC Reference CaffeNet model trained on the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) 2012 to classify 1M+ images to 1000 ImageNet classes.
- A new model was built by discarding the fully-connected portion of the network, and by attaching an extra convolutional layer
  - Having no fully connected layers means:
    - Fewer parameters thus **low computational complexity**: the algorithm can be deployed even on drone
    - Allows for handling input **images with arbitrary dimension**



•Inspired by Linear Discriminant Analysis (LDA) that aims at best separating training samples of different classes, the proposed model combines:

•A softmax loss layer (classification layer) that preserves the **between class separability**,

•An Euclidean loss layer that aims to bring the training samples of the same class closer to the class centroid





#### **Visualization by t-SNE**



- Samples representations in 2D at 1, 10, and 20 epochs of Softmax training
- As the training proceeds samples from the two classes (crowd / non-crowd) become separated
- t-SNE: t-distributed Stochastic Neighborhood Embedding (dimensionality reduction)





- Samples representations in 2D at 1, 10, and 20 epochs of **Discriminant Analysis training**
- As the training proceeds samples from the two classes (crowd / non-crowd) are brought closer to the class centroid



#### **Experimental Results**

#### MultiDrone

• The authors compare the Discriminant Analysis (DA) regularizer with the standard L1 and L2 regularization schemes

• The **DA** regularizer improves the classification performance, while being also **superior** over the L1 and L2 regularizers

Training Approach	Test Accuracy
Softmax	$0.9435 \pm 0.0079$
Softmax & L1	0.9435 ± 0.009
Softmax & L2	0.9422 ± 0.005
Softmax & DA	0.9541 ± 0.0072
Softmax & MEB	0.9546 ± 0.0061

#### **Crowd-Drone Dataset**

This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 731667 (MULTIDRONE)



#### Speed



- The proposed crowd detector was tested on a GeForce GTX 1080 GPU for various input sizes
- •The performance was also tested on a **NVIDIA Jetson TX2 module**, which is a state of the art GPU used for **on-board drone perception**.
- •The method was compared in terms of processed frames per second (fps) with a common baseline model (i.e. VGG-16)



#### Speed





Model	Input	Jetson TX2	GeForce GTX 1080
VGG	224x224	9.36 fps	89.52 fps
Proposed	224x224	49.7 fps	416.66 fps
Proposed	512x512	13.1	99.4
Proposed	1024x1024	2.1	23.45



#### **Deep Detectors on Drones**





#### **Deep Detectors on Drones**







#### **Crowd Projection onto a 3D map**

- The 2D heatmap with crowd existence probabilities can be projected on a 3D map of the flight environment to localize the crowds on the map
- The heatmap is thresholded to obtain a binary image
- Contour following is performed to find the crowd region boundaries



RGB frame from drone camera





crowd detection heatmap crowd thresholded regions crowd regions boundaries

#### **Crowd Projection onto a 3D map**

- Assuming that drone position, camera intrinsic parameters (e.g. focal length) and camera orientation are known, the crowd contours (boundaries) can be projected from the image onto the 3D map (octomap)
- A ray is cast from each of the image contour points towards the octomap.
- A series of voxels defining the points of the crowd contour on the octomap is found
- Crowd contours generated as the drone flies are **merged** using OR operator
- Tests were performed on synthetic and real data



Octree in memory: 130 MB Octree file: 50 MB (2 MB .bt) 3D Grid: 649 MB





# **Tests on simulation videos**

- Synthetic crowd videos captured by a drone were generated in AirSim Drone Simulator
- AirSim provides drone and camera parameters, thus projection of the detected crowds on to the 3D map (also generated by AirSim) is possible
- The crowd regions are visualized in Google Maps



### **Tests on real data**

- Drone flew over a crowd gathering location in the University
- GPS data were used for drone localization
- A 3D map of the location was constructed from drone images using photogrammetry software
- Crowd detected on the drone video was projected on the 3D map and visualized in Google Maps





- Human crowd detection for safe autonomous drones
- Visual drone detection
- Emergency landing site detection.





 The widespread use of drones raises concerns regarding droneto-drone collision or collision with other aerial vehicles







- We need a way to detect drones to address such safety concerns as well as other issues:
  - **privacy concerns**: detect drones spying on people and collecting personal data, flying low near private properties
  - people/ buildings safety concerns: detect drones flying near people or buildings (especially sensitive ones such as government offices), risking deliberate or unintentional collision and injuries or damages
  - **Drone cinematography concerns:** when filming with multiple drones a drone shall not be visible in the video feed of another drone





#### **Radar-based Drone Detection**

- MultiDrone
- Radars are typically designed to detect big and speedy vehicles
- Radars designed for UAV detection are usually bulky, costly and targeted towards ground usage by governments, commercial venues, airports, etc.









- MultiDrone
- Lightweight radars for on-drone usage have recently emerged but are not yet broadly used or available
- Thermal /IR cameras mounted on drones can also be used to detect other drones in certain circumstances
  - + can operate during night or in low light conditions.
  - - affected by noise and other thermal sources



Fortem TrueView radar, 650g, 32W, 2018



DJI Zenmuse XT thermal camera

This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 731667 (MULTIDRONE)





- Visible light cameras can be also used for drone detection
  - +All drones are equipped with such cameras
- However, visual drone detection is a difficult task
  - At long distances drones appear **very small** in the frame
  - Easy to **confuse** them with other aerial objects, such as distant **planes** or even **birds**
  - They may be hard to distinguish from the

**background**, especially in urban areas







- Still, drone detection on RGB videos is not impossible especially with the advent of Machine and Deep Learning techniques for small object detection
- There are already **commercial products** deploying (among others) ML models upon image data for drone detection





- Ctrl+Sky is one such product, which uses multiple types of sensors which communicate with one another
- Drone detection is triggered by a signal captured by a sensor
  - The detection signal can be wifi signal, electromagnetic signal (radar), noise (microphone), RF signal (spectrum tool), or image (camera)

Stationary

Mobile







- Then regions of interest are detected on the video with a detector network
- These areas are tracked using various trackers (including CNNbased ones)
- Areas are classified into drones or distractors





- State-of-the-art object detectors such as YOLO, SSD, Faster-RCNN and their variants, can be trained to detect drones
- The more lightweight ones, based on YOLO or SSD, can even run on-board the drone, e.g., for collision avoidance or to avoid a drone entering the field of view of another one,
- Motion cues can also be used to aid the detection, especially in situations with much background noise







- On-drone cameras may be affected by the drone's motion
- In [1], the camera's motion as well as the target's motion are compensated for by first using neural networks to align objects of interest in successive images



• Then, classifiers trained to distinguish between drones and planes are used to detect objects of interest

[1] A. Rozantsev, V. Lepetit and P. Fua. Detecting Flying Objects using a Single Moving Camera, in IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 39, p. 879 - 892, 2017.







- In a sliding window fashion, patches from consecutive video frames are fed into a Convolutional Neural network (CNN) trained for coarse central alignment
- The coarsely aligned patches are then fed into a second CNN trained to correct for small motions
- The result is consecutive patches in which the object of interest is centered







- The resulting centered patches constitute the input of yet another CNN trained for classification
- The network is trained to differentiate between drones, planes and background







 In 2017, the SafeShore EU Project (<u>http://safeshore.eu/</u>) announced the Dronevs-Bird detection challenge [2], along with a new dataset



[2] Coluccia, et al, Drone-vs-Bird detection challenge at IEEE AVSS2017. In 14th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS), 2017







 The best performing methods in the contest used Convolutional Neural Networks and some variation of state-of-the-art object detectors

#### • [1]: two stage detection

- In the first stage, regions of the image which possibly contain drones or birds are detected, either by median background subtraction or by a Neural Network
- In the second stage, these regions are fed into a CNN classifier, to differentiate between drones and distractors (e.g., birds)

[1] Arne Schumann, Lars Sommer, Johannes Klatte, Tobias Schuchert, and Jurgen Beyerer. Deep cross-domain flying object classification for robust UAV detection. In IEEE International Workshop on Small-Drone Surveillance, Detection and Counteraction Techniques, Lecce, Italy, Aug. 2017.







- Papers [2] and [3] both use state-of-the-art object detectors with modifications to facilitate small object detection
  - [2] uses YOLOv2 with heavy dataset augmentation, even using crawled images
  - [3] uses Faster R-CNN with VGG16 backbone, pretrained on ImageNet for classification

"[2]" Cemal Aker and Sinan Kalkan. Using deep networks for drone detection. In IEEE International Workshop on Small-Drone Surveillance, Detection and Counteraction Techniques, Lecce, Italy, Aug. 2017.

"[3]" Muhammad Saqib, Nabin Sharma, Sultan Daud Khan Makkah, and Michael Blumenstein. A study on detecting drones using deep convolutional neural networks. In IEEE International Workshop on Small-Drone Surveillance, Detection and Counteraction Techniques, Lecce, Italy, Aug. 2017.







- Human crowd detection for safe autonomous drones
- Visual drone detection
- Emergency landing site detection.





- Automatic detection of landing sites is important for drone safety
- Detected landing sites can be used for normal or emergency landing
- A landing site detection algorithm shall preferably run onboard the drone to ensure that it can be used in emergency situations.
  - Drone might have lost contact with the ground
- Data used for landing site detection:
  - Videos, images
  - Lidar (Light Detection And Ranging) data / point clouds
  - 3D terrain data, e.g. Digital Elevation Models (DEM)





Maker

# UAV Landing Site Detection

- In [1] feature vectors based on Histogram of Oriented Gradient (HOG) features are evaluated on image patches
- 'Very dangerous' and 'Not recommended' patch detectors (linear SVMs) are trained using 1,200 labeled Google maps images
- Detection parameters are tuned on a further 400 labeled Google maps images
- The output of these detectors is combined to generate a 'heat map' that describes the level of danger for a given area.
- This 'heat map' can then be used to identify a 'safe' landing site.

[1] A Robust UAV Landing Site Detection System Using Mid-level Discriminative Patches,X. Guo, S. Denman, C. Fookes, S. Sridharan, ICPR 2016





- A method for landing site detection on building rooftops is proposed in [2]
- The method can run on-drone and uses a single camera
- It creates a 3D reconstruction of the environment based on dense motion stereo
- **Frame pairs** with appropriate temporal distances are selected from the videos.
- Features (STAR features and MSURF descriptor) are matched between selected frames, in order to perform **rectification**.
  - Rectification: a transformation process used to project images onto a common image plane

[2] Vision-based Landing Site Evaluation and Trajectory Generation Toward Rooftop Landing, V. Desaraju, N. Humenberger, R Brockers, S. Weiss, L. Matthies, Robotics: Science and Systems 2014







- Stereo disparity-based 3D reconstruction is performed
- The 3D model serves as input to the landing site detection algorithm.
- Suitable landing sites shall be:
  - On the rooftop: locations with low height are rejected
  - Approximately planar and level: the variance of the disparity map is used
  - Sufficiently large to permit UAV approach and landing





- The algorithm labels pixels as:
  - below roof top,
  - on roof top but unsafe,
  - insufficient space
  - safe landing area.
- A confidence map is also constructed
- Both the "landing" and the confidence map are updated over time
- An algorithm for defining a **landing trajectory** in the most suitable landing site is also proposed













\* \* \* \* \* \* \* \*



- [3] proposes a UAV forced landing site detection system operating on image data
- The input image is sampled into small overlapping patches
- Feature extraction is performed: Color (HSV) and texture (HOG, LBP) features are calculated
- A **modified footprint operator** is also employed to better describe the geometric characteristics of the local area surrounding a pixel
- Classification of each patch is performed using either a Gaussian Mixture Model (GMM) or a Support Vector Machine (SVM) classifier, trained on manually labeled data
- The result is a binary (safe / unsafe) image, filtered to remove small areas

[3] Automatic UAV Forced Landing Site Detection using Machine Learning, X. Guo, S. Denman, C. Fookes, L. Mejias, S. Sridharan, 2014 Int. Conference on Digital Image Computing: Techniques and Applications (DICTA)











- A method for landing zone detection in vegetated areas using point clouds from Lidar data is proposed in [4]
- The method can detect landing areas even in low vegetation areas that are detected as unsafe by methods that use surface planarity

and roughness



[4] 3D Convolutional Neural Networks for Landing Zone Detection from LiDAR D. Maturana and S. Scherer, 2015 IEEE International Conference on Robotics and Automation (ICRA)





 The input of the system is a stream of registered LiDAR point clouds and candidate landing sites to be evaluated.

- A volumetric density map that describes the density /spatial occupancy of each voxel in a grid is constructed
- Space containing vegetation (safe for landing) is "seen" as relatively "porous" (less dense) by the LiDAR sensor, compared to space containing solid objects (obstacles, area not safe for landing)





- A 3D Convolutional Neural Network (CNN) is used to predict the safety of subvolumes within this map
  - A volumetric CNN similar to those in [5], extended from 2D data to 3D data is used.
- The CNN model consists of an input layer, one or two convolutional layers, a single fully connected layer, and a classification output layer



 $\operatorname{Prob}(\mathit{safe}|\mathbf{x}) = 0.2$ 

[5] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," JMLR, vol. 15, pp. 1929–1958, 2014.





- MultiDrone
- The model is trained using semi-synthetic point clouds
  - real point cloud data for vegetation and ground
  - simulated point clouds for solid obstacles.



• The output is probabilistic safety prediction for each landing site.







DSM

Digital Terrain Relief Mode

- 3D terrain information in the form of Digital Elevation Models (DEM) can be used for landing site detection :
  - Raster DEMs: images where the pixel value denotes the elevation
  - Digital Surface Models (DSM) include information for both the ground and the manmade structures or vegetation
  - Digital Terrain Models (DTM) include bare ground information, without any man-Z. elevation made structures or vegetation. Digital Surface Mode DSM

he European Union's Horizon 2020 researc grant agreement No 731667 (MULTIDRON

Z, elevation

DTM



- In [1] the authors use the average height and height variance inside quadtree based DEM partitions.
- Partitions whose height variance is below a limit (flat) are selected as landing sites and merged with neighboring partitions if they have similar average heights.



[1] M. Garg, A. Kumar, and P. B. Sujit. "*Terrain-based landing site selection and path planning for fixed-wing UAVs*." International Conference on Unmanned Aircraft Systems (ICUAS). IEEE, 2015.







- In [2] surface fitting on coarse elevation models using Least Squares Error is performed.
- The slope of the fitted surface is used to specify landing areas.
  - Low slope areas are selected



[2] Aydin, Musa, and Emin Kugu, "*Finding smoothness area on the topographic maps for the unmanned aerial vehicle's landing site estimation.*", Sixth International Conference on Digital Information and Communication Technology and its Applications (DICTAP), IEEE 329946 on programme under grant agreement No 731667 (MULTIDRONE)



- In [3] an algorithm that detects potential landing sites is proposed
- The method receives as input two digital elevation models in raster format
  - the digital surface model (DSM) and
  - the digital terrain model (DTM) of a region.



[3] E Kakaletsis, N. Nikolaidis, Potential UAV Landing Sites Detection Through Digital Elevation Models Analysis, Technical



- Step 1. Detection of man-made structures and vegetation:
  - By subtracting DTM from DSM and applying a threshold to the outcome, a binary image is derived, marking pixels depicting man-made structures and vegetation (in black)







- Step 2. Terrain slope determination:
  - Local slope of the terrain is evaluated by applying the well known Sobel operator (edge detector) on the DSM.
  - The operator evaluates the gradient of the DSM image in each pixel, essentially the terrain slope.
  - Dark areas correspond to small slope







- Step 3. Slope image thresholding:
  - The elevation slope image is thresholded
  - DSM pixels are classified into flat or non-flat based on the local slope.
  - Near flat (small slope/gradient) areas (in white) are retained as potential landing areas.







- Step 4: Binary image connected components evaluation:
  - Connected components analysis is applied on the binary image resulting from the previous step.
    - Finding groups of connected pixels
  - Sufficiently large connected components i.e. of sufficient area for landing are retained (dark blue).









- Step 5. Creation of the final landing map:
  - Areas overlapping with buildings and vegetation, found in step 1, are removed from the large, flat areas found in the previous step.
- Blue pixels: landing zones, i.e., small slope and enough pixels
- Light blue pixels: no landing zones, i.e., large slope or very few pixels
- Yellow pixels: no landing zones due to buildings and vegetation













#### Thank you very much for your attention!

#### Contact: Nikos Nikolaidis nikolaid@aiia.csd.auth.gr www.multidrone.eu

