



## City Research Online

### City, University of London Institutional Repository

---

**Citation:** Ter-Sarkisov, A., Ross, R. J. and Kelleher, J. D. (2017). Bootstrapping Labelled Dataset Construction for Cow Tracking and Behavior Analysis. In: 2017 14th Conference on Computer and Robot Vision (CRV). (pp. 277-284). IEEE. ISBN 9781538628188

This is the accepted version of the paper.

This version of the publication may differ from the final published version.

---

**Permanent repository link:** <https://openaccess.city.ac.uk/id/eprint/24660/>

**Link to published version:** <http://dx.doi.org/10.1109/CRV.2017.25>

**Copyright and reuse:** City Research Online aims to make research outputs of City, University of London available to a wider audience. Copyright and Moral Rights remain with the author(s) and/or copyright holders. URLs from City Research Online may be freely distributed and linked to.

---

City Research Online:

<http://openaccess.city.ac.uk/>

[publications@city.ac.uk](mailto:publications@city.ac.uk)

---

# Bootstrapping Labelled Dataset Construction for Cow Tracking and Behavior Analysis

Aram Ter-Sarkisov, Robert Ross, John Kelleher  
School of Computing,  
Dublin Institute of Technology  
Dublin, Republic of Ireland  
453615@dit.ie

**Abstract**—This paper introduces a new approach to the long-term tracking of an object in a challenging environment. The object is a cow and the environment is an enclosure in a cowshed. Some of the key challenges in this domain are a cluttered background, low contrast and high similarity between moving objects – which greatly reduces the efficiency of most existing approaches, including those based on background subtraction. Our approach is split into object localization, instance segmentation, learning and tracking stages. Our solution is benchmarked against a range of semi-supervised object tracking algorithms and we show that the performance is strong and well suited to subsequent analysis. We present our solution as a first step towards broader tracking and behavior monitoring for cows in precision agriculture with the ultimate objective of early detection of lameness.

**Index Terms**—machine learning; animal behavior; machine vision

## I. INTRODUCTION

In the domain of modern animal husbandry the early detection and treatment of lameness is a serious and challenging problem. Lack of adequate treatment can lead to substantial losses for farms and reduced well-being for animals [1], [2]. Detection and treatment of animal lameness has traditionally involved the hiring of expensive specialists after the disease has already become highly pronounced. Due to the negative consequences of late detection, there has recently been an increased interest in applying statistical and machine learning methods to lameness detection. These methods range from regressions and longitudinal studies ([3], [4]) through to neural networks and support vector machines [5], [6], [7]. Analysis such as these mainly rely on 4 main sensor types: accelerometers, weight platforms, remote sensors and video cameras. In our work we are particularly focused on video data since it is convenient for both humans and cows; there is no need for a lengthy installation of equipment, the equipment is cheap when viewed over the long term, and importantly for the animal the method is non-invasive.

The main drawback of video based analysis in this domain is the complexity in information retrieval: one needs to extract the animal’s shape and behavior over a period of time and in relatively complex environments. As such, to the best of our knowledge, previous video recordings for lameness detection were performed in open space with good contrast between the

background and the object, as in [7], [8] over a short period of time. In Section II we discuss these and other results in greater detail.

It has been shown elsewhere that there are multiple features that correlate with lameness, i.e., gait, head tilt, weight distribution, and behavior. While each of these can in principle be analyzed through video, we suggest that analysis of behavior is particularly interesting, because difference in behavior between lame and so-called sound cows has been observed extensively ([9], [10], [3], [11]). For example it has been shown that the frequency and duration of actions like lying and walking correlates with lameness onset.

Our long term research goal is to monitor animal behavior directly from video to predict the early onset of lameness, which, given the above, is harder than with the use of accelerators or weight platforms from both scientific and technical points of view. Among other things, it cannot be approached in a straightforward manner (raw data → features → classification), because extraction of raw features and generation of a labelled corpus from video data full of challenges (poor lighting conditions, frequent occlusion, bad contrast) are problems that do not have a straightforward solution. In the current article we set out a method that combines deep learning algorithms, heuristic methods and an ensemble learning algorithm to track the movement of a cow and construct a labelled dataset.

Our method is constructed and benchmarked against a new video corpus of animal behavior. We detail the content and construction of this video corpus in Section IV. We then present the details of our novel approach to object tracking in Section V. In Section VI experimental results are presented together with a comparison to other relevant algorithms. Finally in Section VII we outline our conclusions and plans for future work.

## II. PREVIOUS WORK

The use of Computer Vision in monitoring animal behavior has a varied history. Rather than providing a comprehensive review of all such work here, we instead focus on the challenge at hand, i.e., the specific case of behavior tracking

for bovines. To the best of our knowledge, there has to date been no direct overlap in the literature between object tracking and behavior-based cow lameness detection from video. However, there have been attempts to use machine learning algorithms like support vector machines (SVM) to directly predict cow lameness stage and lameness identification from video. In [7] cow's backposture was extracted from two sets ( $n_1 = 28, n_2 = 66$ ) of cows shot on video while walking and scored by observers based on the presence of lameness features between 1 (no lameness) and 3 (severe lameness). The authors report a 96% accuracy (percentage of correctly classified observations) of estimation using this method. Background subtraction in the publication was done in a relatively straightforward manner, because the video was shot outside at a convenient angle with little to no clutter in the background and very mild partial occlusion. Most importantly, contrast between the cow and the background is quite sharp, and there were no other cows in the video.

Earlier work [8] showed that cow's movement (locomotion) was examined in a manner similar to [7]. That work established high correlation (94.2%) between hoof positions estimated by the camera and humans, and the difference in the position of the hooves between sound and lame cows. In addition to these two studies, there have been other publications on lameness detection based on video, but all of them were shot outside, with good contrast between background and object, making the task much easier than the problem we address in this paper.

Beyond the specific application to cow tracking, there has of course been significant research into object tracking. Of this work the kernel-based, semi-supervised and ensemble tracking algorithms have shown much promise. Tracking algorithms in these categories are some of the most popular, because they fulfil two important requirements: they do not require fully-labelled datasets for supervised learning and they generalize well to different problems. Semi-supervised algorithms include Tracking-Learning-Detection (TLD), developed by Kalal et al, [12], a related Median Flow algorithm, also by Kalal et al, [13], Multiple Instance Learning (MIL) ([14], [15]) with an extension built upon random forests, [16]. Kernel-based algorithms use a Kernel filter for similarity measure: [17]. One of the most recent approaches is Kernelized Correlation filter (KCF) by Henriques et al, [18], [19], that employs circulant matrices, several types of kernels (Gaussian, Linear, Ridge) and Fast Fourier Transform to learn a set of dense samples (all subwindows) from the tracked object. A big advantage of all these algorithms is that the user only needs to define the starting coordinates and the size of the bounding box (hence semi-supervised), therefore they are easy to test. Also recently, deep learning algorithms, such as Convolutional Neural Nets (ConvNN) were adapted for tracking tasks by Ma et al in [20]. These ConvNNs use a popular VGG-Net-19 architecture

([21]). A major disadvantage to these algorithms and the reason that we cannot apply them directly at this point is that they require a significant amount of labelled data.

### III. APPROACH OVERVIEW

In our work we have been interested in the automated bootstrapping of labelled dataset construction to cut down on the cost and technical challenges associated with building a large labelled dataset suitable to use with for example Deep Learning methods. Our approach to building this dataset is to build a predictive model that can track a particular cow through the video. We use a Random Forest classifier to do this tracking. Here we will give a high-level overview of the stages in the approach so as to help reader understand how the different components in our approach interact.

Our approach to building a tracking model for a given cow involves a number of steps. The first step is that we choose the cow in the video that is to be tracked. We do this by drawing a bounding box around that cow in the first frame of the video and label this bounding box as containing the target cow. We then also draw bounding boxes around all the other cows in the video and label these bounding boxes as distractor cows. We then sequentially process each frame in the video.

When processing each frame we first apply a model called CRFasRNN [22] to localize blobs of pixels that the model predicts as belonging to a cow or multiple cows. Because of this ambiguity, we then apply an edge detector called HED [23] followed by a thresholding method (ISODATA [24]) and this process isolates instances of cows within each blob. In other words, this process may segment a blob further into multiple cow instances. Note, that for now we do not address the problem of merging blobs but we will discuss this in future work.

Once we have extracted a set of cow instances from the frame, we then label each of these cow instances in the frame as belonging to either the target cow or one of the distractor cows. We do this by labelling each instance in the current frame with the label of the nearest instance in the previous frame. Using this process of frame analysis followed by the nearest neighbor instance labelling, we can track the target cow successfully through a short well-behaved video sequence. However, this approach doesn't scale to longer noisier videos. To do this we use the labelled short video to construct a training set for a random forest model that can track the cow through the longer more difficult video sequences. We build this training set for the random forest model by extracting 9 features from each instance in each frame. We then construct the training set by having one-row-per-instance-per-frame with each row labelled as being either the target instance or distractor instance. We then train and validate the random forest model on this constructed dataset and then test this model on the portions of the video

not used to construct the dataset. Having obtained the results, we manually identify true and false positives on a frame by frame basis.

#### IV. DATASETS

Our raw datasets consist of video data collected over a period of 14 days in a cowshed environment. Cameras observed enclosures which contained 10 individual animals. Cameras were mounted at a fixed angle to the animals, and in total 3TB of video data was collected. The data was collected and provided by the Irish National Agriculture and Food Development Authority (Teagasc), No labelling of the raw data was provided.

From this wealth of data we extracted a number of snippets, listed in Table I. As a video sequence is a time series, i.e. ordered data, each snippet is split into two consecutive subsets: training/validation and testing. Test datasets are usually a few times longer than training and validation datasets. Since cows move slowly, we only took every 10th frame from the video sequence. Therefore, in Table I we report both the number of frames and the length of the video. Names of datasets consist of the channel/camera id, date (day, month, year format) and the time in hour:minute:second format, with the actual time corresponding to the first second of the snippet. For example, **ch0106092015115543** means camera number 1 shooting on June, 9, 2015 starting at 11:55:43.

Training and testing datasets were selected with the following objectives in mind: in the training set the cow that we want to track has to be fairly well visible for the whole duration (therefore they are quite short), so that the learning algorithm has enough correctly labelled features to train with. In the test set, on the other hand, most or at least some challenges should be present. While we recognize that these are simplifying assumptions in comparison to making a purely random selection of sequences from the data, we believe that this method is useful at this stage.

Given the nature of the problem and the recording environment, the dataset includes a number of challenges:

- 1) Background: background in the video is generally dark and suffers from the low contrast, it is easy to confuse the background with animals, especially dark-skinned (black and brown), certain cows (especially black ones) are often not discernible even by a human eye,
- 2) Lighting: the lighting is generally low and uneven due to the presence of narrow and long gaps in the walls and ceilings. As a result, many cows have bright rectangular patches on their skin, often splitting the object in two or more parts.
- 3) Objects: In every enclosure there are 10 cows of approximately the same size and different skin color, usually black, brown, white and striped (white and black), hence

they are easy to confuse with each other. As cows are malleable objects, throughout the video their appearance changes substantially, from small while facing away from the camera and blending with the background to large and contrasting when standing perpendicular to the camera's direction.

- 4) Occlusion: there are two types of occlusion in the dataset. First, the components of the cowshed and enclosure, like metal bars and concrete troughs that serve as boundaries of the enclosure. Secondly, due to the size of the enclosure, cows block each other from view much of the time, thus if a cow changes its action (walking→standing, standing→lying) while blocked away from the view, it is very challenging to identify this change automatically.

These features cause considerable trouble for existing tracking algorithms. This is exhibited in Fig. 1, where we used five tracking algorithms implemented in OpenCV 3.1.0 library for Python 2 and mentioned in Section II: MIL, TLD, KCF, Boosting and MeanFlow. Their drawbacks become evident after about first 70-100 frames (35-50 s) as the cow starts to move from its starting position. Trackers fail to account for the changing shape of the cow as it turns around and instead learn from other objects in the bounding box: a similar cow and the background. The second challenge are the metal bars (enclosure boundary) serving as a partial occlusion as the cow moves behind it.

#### V. METHOD

An essential feature of many tracking algorithms is the dependence on the contrast between the tracked object and background, however cluttered it is. This poses a particular challenge for our project, because pixel intensity (hence the color) of large areas of the background, such as the floor, drinking troughs, metal bars and pathways between enclosures are very similar to that of many dark-skinned cows. Similarly, light-skinned cows are easily confused with patches of light passing through gaps in the wall and ceiling. Another problem is that cows of the same color are essentially similar, hence it is enormously challenging for a tracking algorithm to tell between two brown cows, especially if one of them blocks the other from camera view. The main idea of our approach for this reason is to extract contours of the tracked object instead of the background. As explained in Section III, the algorithm consists of two steps: the first one does instance segmentation, feature extraction and training and validation, it is presented in Fig. 2 and the second does feature extraction and testing of a learning algorithm, and is presented in Fig. 3. The instance segmentation phase is presented in Fig. 4.

##### A. Framewise cow instance segmentation

This is the first important phase in both steps of our approach. We start by localizing potentially interesting areas in the frame, and for this purpose we use a pre-trained deep learning algorithm: Conditional Random Fields as Recurrent Neural Network (CRFasRNN), recently introduced in a paper

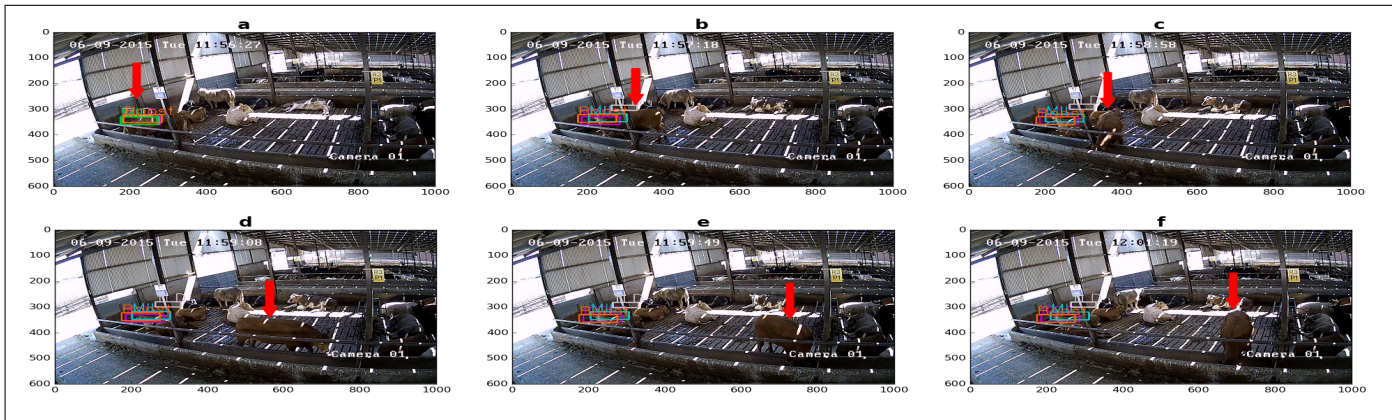


Fig. 1: Performance of KCF (purple box), MIL (light blue box), TLD (white box), MedianFlow (green box) and Boosting (orange box) tracking algorithms on the test dataset **ch0106092015115543**. The red arrow points to the cow that these algorithms must track for the whole duration of the snippet. Each algorithm is initialized in the frame (a) with the features sampled from the user-defined bounding box that includes the greater part of the tracked cow. All six of them break down as the cow starts moving away from the starting position. The cow first ruminates close to the starting position in (b), then moves to the right of the vertical bar and stops to drink in (c), then continues its movement to the right in (d) in order to finally stop and start eating/drinking from the trough, in (e) and (f). As detailed in Table I, the total length of this snippet is 300 seconds. Best viewed in color.

TABLE I: Size (in frames) and length (in seconds) of training/holdout and test sets.

Title	Size of training and holdout sets	Size of test set and holdout sets	Length of training and holdout sets	Length of test set
<b>ch0106092015115543</b>	107	600	43	300
<b>ch0406272015143027</b>	150	360	150	360
<b>ch0106292015090316</b>	92	750	91	765
<b>ch0710062015201033</b>	89	362	89	364
<b>ch0915062015120155</b>	47	600	25	605

by Zheng et al in [22]. This is a combination of a fully convolutional network (FCN, see [25]) and a conditional random field (CRF) with a layout of a recurrent neural network, RNN. This novel object segmentation algorithm returns a mask with the color of the pixels corresponding to the identified class, see Fig. 4. Yellow blobs in the mask correspond to the identified cow pixels (some cows in this video were misclassified as sheep and horses; to avoid further complications, we relabelled these pixels as cows and note that further training of the CRFasRNN is likely to increase accuracy). So far CRFasRNN has shown very high performance on our video data, compared to other algorithms, including FCN. ConvNN’s architecture used for our approach is FCN-8, which is based on VGG-16 ([21]). This identification of elements in the raw image is by far the most computationally intensive part of our model, taking  $\sim 50$ s/frame on the CPU and  $\sim 7$ s/frame on GPU (Tesla K40).

Once we have localized the blobs with potentially interesting objects, we get a bounding box around those larger than an optimal threshold empirically found to be 800 pixels and extract this patch from the original image. What we want now is to identify instances of cows, since CRFasRNN does not report the number of instances, but only that the observations in these blobs are consistent with the

class ‘cow’. As explained above, background subtraction is intractable in cowsheds due to very noisy and low-contrast background. Instead, we attempt to extract contours of objects. This comes from the observation that, even if the objects have a similar pixel intensity, the value of the derivative at the boundary has the potential for delineating between them.

For contour detection we found that the Holistically-Nested Edge Detector (HED), recently introduced by Xie and Tu in [23] performs strongly by filtering out most of the noise in the image. Therefore in the next stage we run each localized area in the image through HED to obtain edges. The output of HED is an image with the darker pixels corresponding to more important edges. Since we only need the most important of them, we want to split each contour image into two disjoint subsets: objects and background. We do this by using the threshold detection algorithm ISODATA, introduced by Ridler and Calvard in [24]. A combination of these two methods produces a set of isolated objects that we label as such and add to the training set should their area exceed 800 pixels (also found empirically). We also get the objects’ convex hulls (to track the change in the cows’ shape) and bounding boxes. This step is much less computationally expensive, taking 5-7s on a CPU or  $< 1$ s on a GPU (Tesla K40).

Our approach takes CRFasRNN’s output a step further by offering a solution for instance segmentation: blobs from CRFasRNN’s output merely tell us that there are cows in this locality. We extract the hypothesis for the number of cows and the approximation for their shapes and locations. In Fig. 4 we present the flowchart of the instance segmentation step.

### B. Feature Extraction

The process detailed in the previous subsection is applied to every frame in the video sequence. We combine that process with a simple 1-Nearest Neighbor (1NN) tracking method to automate feature extraction from instances identified in the first phase of processing. In the very first frame of the sequence, we select the cow we want to track and assign it label 1; all other objects are labelled with 0. For all other frames, once we have all the instances, we get a simple distance matrix, where the number of rows is equal to the number of instances in the current frame, and the number of columns to the number of labelled instances in the previous frame, hence each entry in the matrix is a distance from every instance in this frame to every labelled instance in the previous frame. Once we have this matrix, we assign every instance in this frame the label of the nearest instance from the previous frame. This heuristic approach is too simple for any serious tracking problem, and therefore our training databases are very small (50-100 frames) and well-behaved, i.e. the cow we want to track is well visible throughout the video of the training set. Once we have labelled all instances in the frame, we extract features from them, which are added to the training dataset in the correct order (i.e. concatenated with the previous data). In total, we use 9 features from three types:

- 1) Pixel intensity: We use the instance’s centroid as the mean and  $\sigma^2 = 5$  to sample 100 5x5 patches in each object; after averaging over their pixel intensities we get a vector of features: overall mean, maximum and three quartiles.
- 2) Size: we use the size of the bounding box around the object. This is motivated by the fact that in the previous stage we ignored small objects (under 800 pixels), and that our instance segmentation algorithm tends to find large portions of cows.
- 3) Location: we store the  $(x, y)$  coordinates of the centroid of the bounding box as the distance feature. During the training and test phases, we find the distance between every instance’s centroid and the centroid of the previous observation of the tracked cow. This coordinate difference is the actual feature used.

Finally, once we have collected all the data from the training video sequence, we manually clean up the training dataset by removing mislabelled (false positive) observations and pass the correctly labelled dataset to the training algorithm.

### C. Training

With features extracted for each cow, we trained a classifier to automatically identify an individual animal in a video frame. We selected the Random Forest (RF) classifier, attributed to Breiman [26] to learn the features. Originally in each training set the proportion of positive (tracked cow) observations is about 16 %. To provide the classifier with more data, we sampled out about 50 % of negative (distractor cows) observations, thus increasing the positive data points to about a third of the training database. As the data is essentially a time series, we train the algorithm on the first  $K$  observations and validate on the remaining  $n - K$ . The second contribution of this paper, after instance segmentation algorithm, is optimizing a classifier based on the training, validation and testing output. We found that RF with 300 trees, cross-entropy error function, using all features during training, with bootstrap samples and out-of-bag samples for generalization do the best job on our data. Training of a single forest takes a very small amount of time,  $\sim 5$  seconds.

## VI. RESULTS

Results from testing the RF classifier, precision and recall, are summarized in Tables II and III. Test sets, which are taken from the same video, are different to the training set in a number of ways: there are many issues that an algorithm must handle, such as full and partial occlusion, bad lighting, low contrast, cluttered background and other. For comparison we use the five trackers mentioned above: TLD, MIL, KCF, MedianFlow and Boosting on each test set. In three datasets our approach achieves the highest precision and in two - the highest recall rate. Its strength is particularly well visible on **ch0106092015115543** and **ch0915062015120155**, where the tracked cow moves around. Other either immediately loose it (as TLD in **ch0106092015115543**) or confuse it with the background as soon as the cow leaves the area where the tracking started. We consider this to be a specific strength of our algorithm. On two sets where our approach underperformed (e.g. **ch0710062015201033**) the problem is related to the generalization capacity of the classifier: the cow does not move much, but its features are too easily confused with those of other objects.

## VII. DISCUSSION AND FUTURE WORK

In this article we have presented a new tracking algorithm developed for tracking malleable objects (cows) in a challenging environment (enclosures in a cowshed). The ultimate goal of this project is to identify lameness in cows at an early stage; successful cow tracking is the first stage in this project. This article has three main contributions:

- 1) Framewise instance segmentation,
- 2) Optimal Random Forest algorithm,
- 3) Construction of a large dataset for further analysis of cows’ behavior

From here there are three main directions in which we would like to take the development of this algorithm:

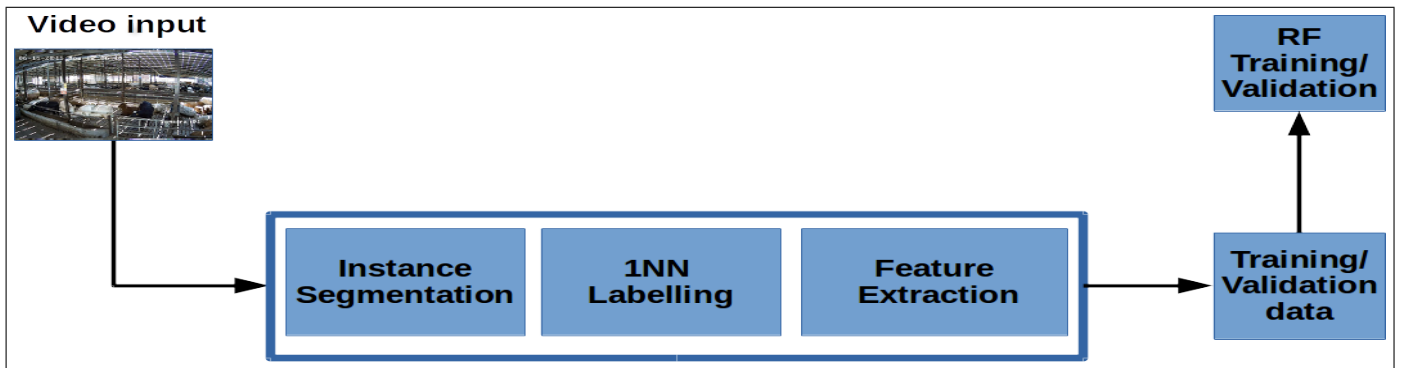


Fig. 2: Flowchart of the first step of the algorithm. Video data is processed framewise: phases within the bounding box are repeated for the full duration of the video, adding labelled data points to the training dataset. Once the training dataset is fully built and manually cleansed of mislabelled observations, it is passed on to the Random forest (RF) algorithm for training and cross-validation. The final output of this stage is an RF classifier with optimal parameters for tracking a particular cow.

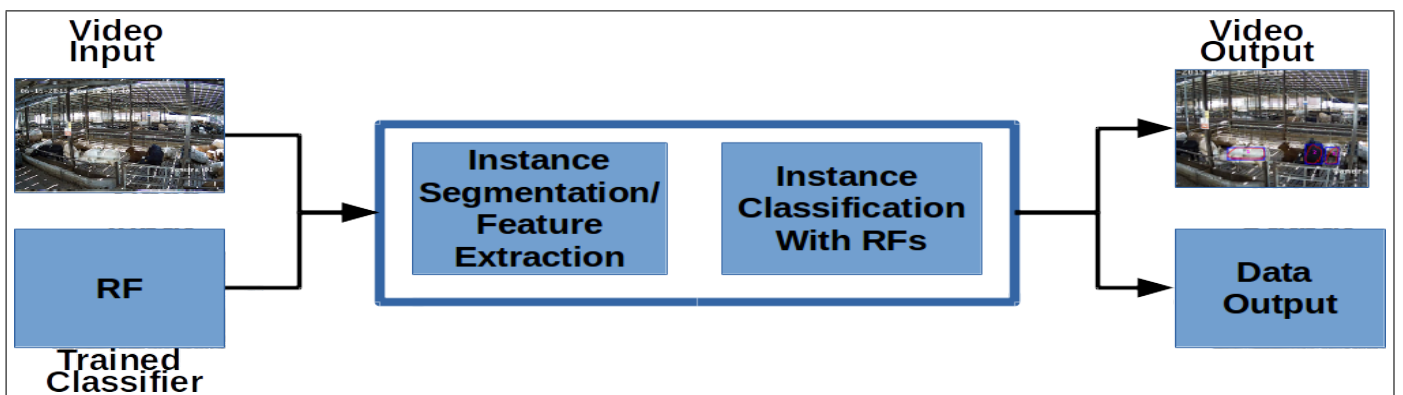


Fig. 3: Flowchart of the second step of the algorithm. Video data is processed framewise: phases within the bounding box are repeated for the full duration of the video. RF classifier optimized in the first step processes objects segmented in the instance segmentation stage in each frame and labels them based on the features extracted from instances. Labels and instances are added to the video output. Some data is stored and output for further analysis.

TABLE II: Performance of the five algorithms available in the OpenCV library on the test datasets compared to ours. Values in each column are the percentage of frames the algorithm identified the tracked cow, i.e. recall. Bold are the best algorithm for this test dataset.

Video	KCF	TLD	MIL	Boost	MF	Ours
ch0106092015115543	0.21/0.2	0/0	0.17/0.2	0.2	0	<b>0.98/0.54</b>
ch0406272015143027	<b>1/0.76</b>	0.56/0.41	<b>1/0.88</b>	0.93/0.17	0.4/0.41	<b>1/0.77</b>
ch0106292015090316	<b>1/1</b>	0.02/0	<b>1/0.56</b>	1/0.35	0.5/0.13	0.7/0.47
ch0710062015201033	<b>1/1</b>	0.44/0.16	<b>1/1</b>	0.78/1	0.4/0.16	0.83/0.08
ch0915062015120155	0.36/0.34	0.25/0.4	0.36/0.34	0.36/0.34	0.28/0.33	<b>0.56/0.85</b>

- 1) Improvement of instance segmentation step. Currently we use a combination of two deep learning, thresholding and a labeling algorithm. Although they do the job reasonably well, there is enough space for improvement.
- 2) Improvement of generalization. Although Random Forest does a good job with the tracking, it does not always confidently generalize to any angle at which the cow faces the camera. We therefore intend to retrain algorithm like CRFasRNN to get it to track objects (transfer learning) using the segmentation information from the previous step.
- 3) Construction of cow behavior dataset. In addition to tracking the cows' movement, we need to track their behavior, as it correlates with the physiological condition (overall, lame cows lie for longer periods of time). For this purpose we will be training a large deep learning algorithm like LSTM ([27]) or a similar recurrent neural network.

#### REFERENCES

- [1] L. Warnick, D. Janssen, C. Guard, and Y. Gröhn, "The effect of lameness on milk production in dairy cows," *Journal of dairy science*, vol. 84, no. 9, pp. 1988–1997, 2001. 1

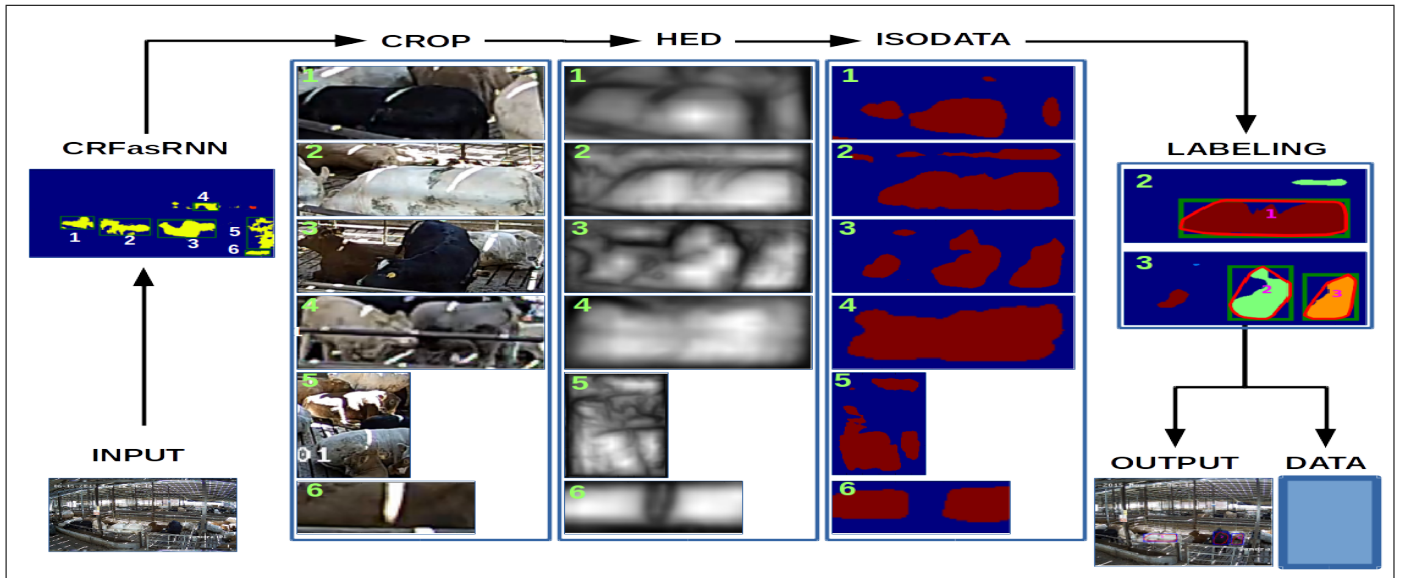


Fig. 4: One full step of cow instance segmentation. The input is a frame 1000x600 pixels. The first step is localization with CRFasRNN. In the second step we crop all promising image insets that exceed a pre-set threshold value. In each inset we extract contours with HED. To obtain isolated objects from these contours we use an ISODATA thresholding function. Each object with size exceeding another pre-set threshold value is assigned an object status: we obtain its convex shape and a bounding box around it. We add these contours and bounding boxes to the input image and sample the data from the image, which is added to the dataset. The output is an image with overlaid bounding boxes and object contours.

TABLE III: Results on the test sets obtained using Random Forest with the best parameters trained on datasets in Table I. TP: true positives, FP:false positives, TN:true negatives, FN: false negatives, P:precision, R:recall. In the last column main challenges that the algorithm faced in the video are listed.

Video	TP	FP	TN	FN	P	R	Challenges
ch0106092015115543	328	6	4928	272	98.2%	54.6%	Partial occlusion, dim and uneven lighting, low contrast
ch0406272015143027	275	0	3119	85	100%	76.8%	Cluttered background
ch0106292015090316	355	131	9529	395	70%	47%	Partial Occlusion, very low contrast, dim lighting
ch0710062015201033	30	6	2945	330	83%	8%	Partial occlusion, very uneven lighting, low contrast
ch0915062015120155	419	300	6048	71	56%	85%	Partial and full occlusion, very uneven lighting, cluttered background

- [2] J. Hernandez, J. K. Shearer, and D. W. Webb, "Effect of lameness on the calving-to-conception interval in dairy cows," *Journal of the American veterinary medical association*, vol. 218, no. 10, pp. 1611–1614, 2001. 1
- [3] C. Kamphuis, E. Frank, J. Burke, G. Verkerk, and J. Jago, "Applying additive logistic regression to data derived from sensors monitoring behavioral and physiological characteristics of dairy cows to detect lameness," *Journal of dairy science*, vol. 96, no. 11, pp. 7043–7053, 2013. 1
- [4] L. Alban, J. Agger, and L. Lawson, "Lameness in tied danish dairy cattle: the possible influence of housing systems, management, milk yield, and prior incidents of lameness," *Preventive veterinary medicine*, vol. 29, no. 2, pp. 135–149, 1996. 1
- [5] M. Pastell and M. Kujala, "A probabilistic neural network model for lameness detection," *Journal of dairy science*, vol. 90, no. 5, pp. 2283–2292, 2007. 1
- [6] P. Martiskainen, M. Järvinen, J.-P. Skön, J. Tiirikainen, M. Kolehmainen, and J. Mononen, "Cow behaviour pattern recognition using a three-dimensional accelerometer and support vector machines," *Applied Animal Behaviour Science*, vol. 119, no. 1, pp. 32–38, 2009. 1
- [7] A. Poursaberi, C. Bahr, A. Pluk, A. Van Nuffel, and D. Berckmans, "Real-time automatic lameness detection based on back posture extraction in dairy cattle: Shape analysis of cow with image processing techniques," *Computers and Electronics in Agriculture*, vol. 74, no. 1, pp. 110–119, 2010. 1, 2
- [8] X. Song, T. Leroy, E. Vranken, W. Maertens, B. Sonck, and D. Berckmans, "Automatic detection of lameness in dairy cattle television-based trackway analysis in cow's locomotion," *Computers and electronics in agriculture*, vol. 64, no. 1, pp. 39–44, 2008. 1, 2
- [9] J. Olechnowicz, J. Jaskowski *et al.*, "Behaviour of lame cows: a review," *Veterinari Medicina*, vol. 56, no. 12, pp. 581–588, 2011. 1
- [10] M. Alsaad, C. Römer, J. Kleinmanns, K. Hendriksen, S. Rose-Meierhöfer, L. Plümer, and W. Büscher, "Electronic detection of lameness in dairy cows through measuring pedometric activity and lying behavior," *Applied Animal Behaviour Science*, vol. 142, no. 3, pp. 134–141, 2012. 1
- [11] Ö. Cangar, T. Leroy, M. Guarino, E. Vranken, R. Fallon, J. Lenehan, J. Mee, and D. Berckmans, "Automatic real-time monitoring of locomotion and posture behaviour of pregnant cows prior to calving using online image analysis," *Computers and electronics in agriculture*, vol. 64, no. 1, pp. 53–60, 2008. 1
- [12] Z. Kalal, K. Mikolajczyk, and J. Matas, "Tracking-learning-detection," *IEEE transactions on pattern analysis and machine intelligence*, vol. 34, no. 7, pp. 1409–1422, 2012. 2
- [13] —, "Forward-backward error: Automatic detection of tracking failures," in *Pattern recognition (ICPR), 2010 20th international conference on*. IEEE, 2010, pp. 2756–2759. 2
- [14] B. Babenko, M.-H. Yang, and S. Belongie, "Visual tracking with online multiple instance learning," in *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*. IEEE, 2009, pp. 983–990. 2
- [15] —, "Robust object tracking with online multiple instance learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 33, no. 8, pp. 1619–1632, 2011. 2
- [16] C. Leistner, A. Saffari, and H. Bischof, "Miforests: Multiple-instance



- learning with randomized trees,” in *European Conference on Computer Vision*. Springer, 2010, pp. 29–42. 2
- [17] D. Comaniciu, V. Ramesh, and P. Meer, “Kernel-based object tracking,” *IEEE Transactions on pattern analysis and machine intelligence*, vol. 25, no. 5, pp. 564–577, 2003. 2
- [18] J. F. Henriques, R. Caseiro, P. Martins, and J. Batista, “Exploiting the circulant structure of tracking-by-detection with kernels,” in *European conference on computer vision*. Springer, 2012, pp. 702–715. 2
- [19] —, “High-speed tracking with kernelized correlation filters,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 3, pp. 583–596, 2015. 2
- [20] C. Ma, J.-B. Huang, X. Yang, and M.-H. Yang, “Hierarchical convolutional features for visual tracking,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 3074–3082. 2
- [21] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014. 2, 4
- [22] S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, and P. H. Torr, “Conditional random fields as recurrent neural networks,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 1529–1537. 2, 4
- [23] S. Xie and Z. Tu, “Holistically-nested edge detection,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 1395–1403. 2, 4
- [24] T. Ridler and S. Calvard, “Picture thresholding using an iterative selection method,” *IEEE Trans Syst Man Cybern*, vol. 8, no. 8, pp. 630–632, 1978. 2, 4
- [25] J. Long, E. Shelhamer, and T. Darrell, “Fully convolutional networks for semantic segmentation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3431–3440. 4
- [26] L. Breiman, “Random forests,” *Machine learning*, vol. 45, no. 1, pp. 5–32, 2001. 5
- [27] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997. 6