

A cow structural model for video analytics of cow health

He Liu, Amy R. Reibman, Jacquelyn P. Boerman
Purdue University, 501 Northwestern Ave, West Lafayette, IN 47907, USA

Abstract—In livestock farming, animal health directly influences productivity. For dairy cows, many health conditions can be evaluated by trained observers based on visual appearance and movement. However, to manually evaluate every cow in a commercial farm is expensive and impractical. This paper introduces a video-analytic system which automatically detects the cow structure from captured video sequences. A side-view cow structural model is designed to describe the spatial positions of the joints (keypoints) of the cow, and we develop a system using deep learning to automatically extract the structural model from videos. The proposed detection system can detect multiple cows in the same frame and provide robust performance under practical challenges like obstacles (fences) and poor illumination. Compared to other object detection methods, this system provides better detection results and successfully isolates the keypoints of each cow even when they are close to each other.

Index Terms—video analytics, pose estimation, cows, video processing

I. INTRODUCTION

Monitoring animal health is a critical component of livestock farming, because healthy animals are more productive. Such monitoring is often performed visually, because animal appearance and behavior are key indicators of health changes. For example, trained farm personnel can analyze a dairy cow’s health condition based on visual appearance [1], and can detect potential illnesses such as lameness [2]. However, time and labor limitations preclude a human routinely watching for these changes, especially in commercial farms which house a large number of cows. Thus there is increasing interest in substituting automated video analytics for human observations. Indeed, video analytic techniques have been applied to different industries including animal agriculture. With the help of the latest computer vision and image processing algorithms, visual animal biometrics has become an emerging research topic [3]. By applying video analytics methods, it is possible to develop a camera system that automatically detects the cow’s health condition with low cost.

The first step to analyze cows using visual data is to detect and segment the cows within the video sequences. This is a straightforward task when each cow walks individually on a well-lit pathway with a very clear background and no obstructions. However, if a camera system is installed on a commercial farm, with the goal of not interrupting daily operations, unrelated objects such as humans and fences are often captured. So identifying the spatial locations of cows is a fundamental first step for further analysis. However, finding the location of the cow is not enough. For the purpose of

assessing an aspect of the animal (i.e. body size or gait), simply having a binary mask that labels the cow’s location is inadequate. Further information of the cow’s structure is required, such as the locations of all body parts or joints¹. This information can then be converted to human interpretable knowledge, or further processed by autonomous health monitoring systems. In summary, we need a video system that not only isolates the cow’s spatial location, but also detects its body structure and tracks its movements.

Designing a video processing system that satisfies these two requirements is not trivial. Many previous segmentation methods focus on object detection, which generates a binary mask of the objects and their corresponding labels. But this is not enough for further cow health analysis. Recently, additional methods have been proposed to detect keypoint-based object structures, like human skeleton models. These models are formed with a series of keypoints or joints connected in a particular order. But these methods are designed with the knowledge of human structure, which is difficult to adapt to other animals like cows. While new methods such as DeepLabCut [4] focus on animal-related keypoint detection methods, this method requires clear video sequences with a single object and a clear background. It cannot be directly applied for practical cow applications on the farm. Finally, there are some visual applications [5]–[7] for cows in the literature, but they are designed for videos captured in a specifically-designed environment, which requires extra efforts and costs for the farm to collect video data.

Processing cow videos collected from a commercial farm also poses specific challenges. First, the environment in which video is captured cannot be fully controlled without interrupting the daily operation of the farm. The cameras need to be installed with specific positions and viewing angles, so that the cows can be clearly observed with few obstructions. Issues such as poor illumination [6] and heavy obstacles [8] largely influence the performance of existing detection algorithms. In addition, the environment also limits the choice of capturing devices. Surveillance cameras are the most suitable devices to install and deploy on typical farms, but the quality of the videos is limited. Distortions such as low color saturation, low frame rate, and heavy compression deteriorate the performance of detection algorithms.

In this paper, we combine deep learning with domain

¹In this paper, we use the term body parts or keypoints to refer to the points that represent joints or specific regions on the cow.

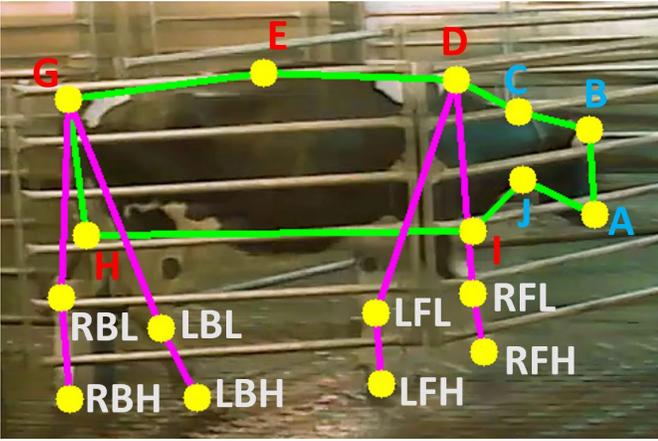


Fig. 1: The proposed cow structural model. 4 blue head region points: A:nose, B:head, C:top of neck, J:bottom of neck. 5 red body region points, D:shoulder, E:spine, G:tailhead, H:mid-thigh, I:bottom of shoulder. 8 white leg and hoof points, with name format: Right/Left + Front/Back + Leg/Hoof.

knowledge about cows to develop a cow structure detection system that operates on videos captured from a practical dairy farm. Our system estimates the number of cow objects in a frame and detects the body parts of every individual cow. For each cow object, the detected body parts compose a side-view cow structural model as shown in Figure 1. This model describes both the spatial location of the cow and additional structural information such as the body contour, the positions of major joints, and the trajectories of their movement. This detailed information provides interpretable knowledge for further health analysis.

This method also overcomes some practical issues. First, all the videos are captured without interfering with the daily work on the dairy farm, requiring only surveillance cameras and no specialized hardware. Second, by incorporating domain knowledge about cows, the video processing algorithm overcomes practical video challenges, such as poor video quality, bad lighting conditions, and heavy occlusion from fences. Later experiments show that our method provides robust results under these conditions.

There are three main contributions in this work. First, we design a cow model with keypoints that presents the structural information about a cow that enables further interpretation for subsequent cow health analysis. Second, a system is developed to extract the cow structural models from videos that are captured from practical dairy farms. Third, for the cow structural model, we also develop corresponding evaluation metrics that operate with limited ground-truth labels. In later experiments, we use multiple video datasets captured from different cameras to test the robustness of the proposed detection system. This system is also compared with other popular object detection algorithms and demonstrates clear advantages when presented with practical challenges. We note that dairy cows are just one type of four-legged livestock, and we anticipate our method

can be easily extended to similar animals.

This paper is organized as follows. Section II reviews previous object detection methods and visual-related applications for cows. Section III introduces the proposed cow structural model including the keypoints and their spatial constraints. Section IV presents the cow structure detection system, followed by detail explanations of the detection module and the post-processing module. The related cow structure evaluation metrics are described in Section V. Next, Section VI describes three experiments of our detection system: evaluation of the individual components of our system, performance on different datasets, and a comparison between different methods. These three experiments demonstrate the effectiveness of each component of our system, the robustness of our system, and the advantages of our method relative to existing methods, respectively. Finally, Section VII summarizes this work.

II. PREVIOUS WORK

Object detection is one of the most popular topics in image and video processing. Traditional Video Object Segmentation (VOS) methods such as [9]–[11] detect objects using motion information from video sequences. With the development of Convolutional Neural Network (CNN), new learning-based methods achieve much better results. Methods such as Mask R-CNN [12], DeepLab [13], and You Only Look Once (YOLO) [14], are widely applied to solve the problem of image semantic segmentation, which requires detection and classification of the objects in an image. CNNs are also applied for video object segmentation. One Shot Video Object Segmentation (OSVOS) [15] does an online fine-tuning process which is trained on one frame of the video sequence and applied to all the other frames in the sequence. This method is further extended with image-based detection methods for semantic guidance [16]. Other methods such as [17]–[19] apply different models using either temporal information or memory for object detection. There are also some popular public VOS datasets available for benchmarks, such as the YouTube VOS [20] and DAVIS dataset [21]. However, all these methods generate bounding boxes or pixel-level masks to represent detected objects, but the structural information of the object is not identified. Additional processing would be required to extract further detailed information from these masks.

Apart from spatial segmentation, there are also research focusing on object structural information such as keypoint detection and pose estimation. Benefiting from the public human pose datasets such as MPII [22] and COCO human skeleton [23], advanced methods are developed for human skeleton detection. DeepPose [24] first applies CNN for human body parts detection based on images, and the stacked hourglass network [25] extends it to detect humans at multiple spatial scales. To solve multi-human detection, the relationship between human joints are considered. ArtTrack [26] generates a simplified human body-part model; OpenPose [27], [28] and Deepcut [29] use part affinity fields to model the joint relationships. However, all these methods are designed by

incorporating different levels of knowledge about the human body, and they are not easily altered or fine-tuned for other objects like cows.

Recently, new methods such as the DeepLabCut [4] toolbox, LEAP [30] and DeepFly3D [31] extend keypoint detection to animals. One advantage of these methods is that they provide a means for users to define body parts; this allows the algorithm to adapt to different animal structures. The DeepLabCut toolbox also provides simple access to fine-tune the networks, and it can achieve promising results with a small amount of training data. However, there are two major limitations of these methods. First, they only support the labeling of one object per frame, and they do not work with multiple objects. This limits their usefulness in many situations. Second, they are designed for video sequences that have been captured under laboratory conditions, with clean background and clear illuminations. In later sections, our experiment shows that the DeepLabCut method does not work well on our cow videos.

Apart from general detection algorithms, there are some previous work on cow-related visual applications; most focus on lameness detection. However, their processing techniques are developed for a specially-designed environment where the captured images are clear enough to process. Normally their detection targets are limited to a specific region instead of the entire cow body structure. For example, methods like [7], [32] only detect the cow’s back curvature while other applications such as [6], [33] only track the trajectories of the legs and hooves. As a result, these methods are not general enough to provide a complete body structure. Apart from lameness, more researchers focus on cow identification with visual data [8], [34]. Their target is to extract cow features such as traditional image features [35] or CNN-based features [36], to distinguish different cow identities. But all these methods need a fundamental step that detects and locates the cows within the images or videos.

III. STRUCTURAL COW MODEL

In this section, we first introduce the keypoints in our cow structural model in detail, and then describe the spatial constraints between the keypoints. These constraints are further used in the detection system for separating multiple cows and detecting missing parts.

A. Cow body keypoints

This proposed structural model is designed to represent a detected cow object in the frame more effectively than using a binary cow mask. It is designed to provide both the spatial location and cow structural details, such as the body shape and positions of the body parts. For consecutive video frames, this model should also provide information so that we can track the movement or motions of these body parts. Inspired by recent approaches to model the human skeleton [24], we combine some anatomical cow joints with other spatial keypoints to represent the cow pose, and the cow structural model is built by connecting the keypoints.

Figure 1 shows our proposed side-view cow structural model. There are 17 points in total to describe the important locations of a cow object from this angle. The upper body region has 9 points, including the head region (blue) and the main body region (red). Connecting these points forms the contour of the upper body region (green lines). Another 8 points are in leg-hoof regions which represent the four limbs, and each limb has a pair of leg and hoof joints. Comparing to the anatomical 51-point cow skeleton model [37], we only select visually-observable joints. Some joints such as the elbow and stifle joints are neglected because their positions are not readily visible and thus difficult to isolate visually. In addition, we also add some points such as the two bottom corners H and I show in Figure 1. Even though they are not physical joints, connecting them with other joints forms a closed contour which spatially locates the body region. The point E on the spine is also an added point, because connecting three spine points provides information about the back curvature which is useful for lameness detection.

There are two general observations about the keypoints in this cow structural model. First, the points in the main body region (red in Figure 1) are always visible from the side-view, and their relative spatial locations do not change dramatically when the cows are walking. Second, the leg and hoof points are more difficult to detect compared to the upper body region points because of the practical issues such as bad illumination, shadows, and fast leg movement. Distinguishing between the points from the left or the right leg is also difficult when there are obstacles in front, for example the horizontal fences shown in Figure 1.

B. Keypoint constraints

Practical constraints limit the potential relationships among the keypoints in both space and time. When the cows are walking between the fences, the cameras located at a fixed position on the side wall always capture the side-view of the cow. In this case, all the cows shown in the video have the similar pose, and the keypoints of their upper body region are always located at relatively fixed spatial positions. For example, the cow’s head always appears on the right side of the body, and the body does not change size. As a result, we can compute general relationships that constrain the keypoints in the cow structural model.

To model the constraints, we first define the center of the cow’s body. This center point is computed as the spatial center of all the keypoints from the cow’s upper body region. Note that the points in the leg-hoof region are not used to compute the center point because their positions are not relatively fixed when the cow moves. Then we can estimate the relative spatial relationship between the center and all the keypoints.

Figure 2 visualizes the keypoint constraints. The middle X shows the cow center c , and the relative spatial locations of the upper body parts appear surrounding the center. Notice each body part mapping function F_j is a 2D Gaussian probability distribution, which is shown as the ellipse in the figure.

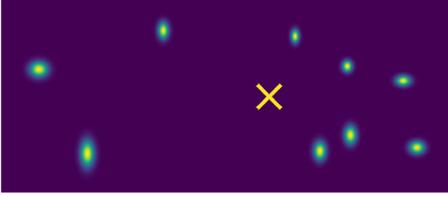


Fig. 2: The constraints between the upper body keypoints of the cow structural model. The yellow X is the cow center, and the surrounding points shows the relative positions of the keypoints in the upper body region.

Formally, for a fixed cow center point $c^* = (x, y)$, we define a set of mapping functions $F_j(\cdot)$ that describe the relative spatial locations of every upper body-part point p_j^* to the center,

$$p_j^* = F_j(c^*) \quad (1)$$

where j is the index of the body part. Each mapping function F_j is characterized by a 2D Gaussian model, and the parameters are trained using all ground-truth labels. During the training process, the approximate cow center c is computed first by averaging all labelled body parts, and the parameters in each F_j are estimated individually based on their relative spatial locations to c .

In the next section, we show how these constraints can be used to separate cows which are spatially close together in the frame. They also provide a reference when assigning body-part candidates to each individual cow object in the post-processing module.

IV. SKELETON DETECTION SYSTEM

This section introduces our proposed system to detect the structure of cows. We first review one popular work for keypoint extraction and then describe the components of our proposed system. Then we explicitly introduce two main processing components: the body part extraction module and the post-processing module.

A. The DeepLabCut toolbox

The DeepLabCut toolbox [4] is a recent popular method to extract keypoints from video sequences. The inputs are color images from videos, and it applies a CNN to generate confidence maps that represent the potential keypoint locations. One advantage of the DeepLabCut toolbox is that it provides simple access for users to manually define the output body parts, and the toolbox automatically alters the last layer of the CNN based on the number of body parts. For example, there are 17 confidence maps generated in our case because we have 17 keypoints in our cow structural model. In our system, we apply the network created by the toolbox to extract the keypoints of our cow structural model.

However, other modules from the toolbox are less suitable for our application because of two major limitations. First, this platform is designed and evaluated with videos captured from a laboratory environment with clear objects and background.

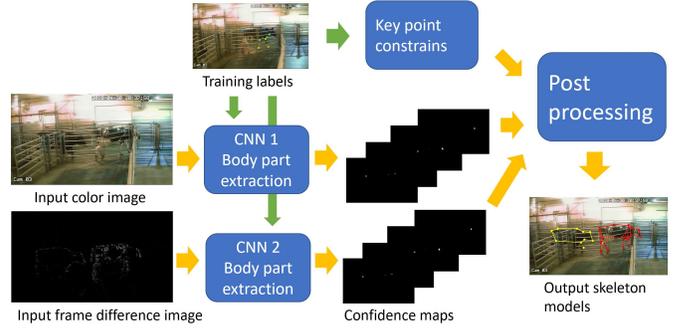


Fig. 3: A diagram of the proposed system. The green arrows show the training process and yellow arrows present the process during operation.

But our cow videos, generated from a commercial farm, have low video quality and the view of the cows are often blocked by obstructions. Later experiments show that the original DeepLabCut does not provide robust detection results on our videos. Second, this method assumes there is only one object in a frame, so it only chooses one body part from each confidence map. If there are multiple body-part candidates detected, only the position with the highest confidence score will be selected. But in videos generated from commercial farms, there could be multiple cows and obstructions like fences that easily cause false detection. We address these two limitations and build a general keypoints detection system which extracts robust keypoints on our cow videos.

B. Proposed system

This detection system is targeted to extract the structural model for every cow object from video sequences. Figure 3 presents the overall system; its primary components are two CNNs for the extraction and a post-processing module. The body part extraction module uses trained networks to convert each single image into a group of confidence maps. Each map shows the potential locations of a particular body part, and the values of the map represent their detection confidence. The post-processing module generates the final structural model based on two groups of confidence maps and the trained keypoint constraints. Both modules are discussed in detail in the next two sections.

In this figure, both the training process and the testing process are labelled using colored arrows. During the training process (indicated by the green arrow in Figure 3), the ground-truth labels are used to fine-tune both CNNs and the keypoint constraints. During operation (indicated by the yellow arrow), the system takes the input of both the color image and the frame difference image on the left, and generates the cow structural model for a single image. After all the frames from a video sequence are processed, the post-processing module refines all the detected cow structures based on temporal information.

C. The body part detection module

The goal of this module is to find the spatial locations of all potential keypoints from raw images. In our system, we apply the original DeepLabCut network [4], labeled CNN1, to extract keypoints from color images. This network structure follows DeeperCut [38], and is implemented using ResNet [39] for the convolution stages, followed by one de-convolution layer before the output layer to recover the target spatial locations of the keypoints. The last two convolution layers apply atrous convolution, which increases effective fields-of-view of the applied convolution and preserves spatial resolution [13]. By default, the DeepLabCut network is pre-trained on ImageNet [40] for image classification tasks, and we use our own cow labels to fine-tune the last de-convolution layer for keypoint detection.

However, as mentioned above, low video quality and heavy obstacles influence the performance. To overcome this issue, we add an extra network, CNN2, into the system. The architecture of this network is same as the first, but it processes frame difference images. There are three major advantages of using frame difference images for our cow videos. First, because we have fixed cameras, the frame difference image better captures the moving objects and eliminates the stationary obstacles such as fences. Second, many of our target keypoints are on the contour of the cow body, and the frame difference highlights these edges of a walking cow.

Third, frame difference also reduces the influence of color variation. This is useful, because the color responses of different cameras are not the same especially under poor illumination. In addition, the majority of the cows have color variations introduced by the patterns on the cows, but some cows only have a single coloring, such as pure white, black or brown.. If these patterns are not included in the training frames, then the color-based CNN methods would likely fail to detect cows with unseen colors. As a result, using frame differences provides robustness to these factors.

However, using the frame difference images alone is not enough because they eliminate too much spatial information, especially for legs and hooves. This is because most of the legs are stationary even when the cows are moving. As a result, our system merges both networks together to improve the body part detection accuracy.

D. The post-processing module

The post-processing module collects and merges the confidence maps from the two CNNs, and assigns the cow body-part candidates to each cow object instead of just to one cow per frame. This step enables the system to detect multiple cows together and track their temporal movements. There are three major steps in this post-processing module: body part extraction, spatial clustering, and temporal filtering.

1) *Body part extraction*: This step extracts the spatial locations of all body-part candidates from the confidence maps generated by the CNNs. Notice that at this stage, the number of cow objects in the image is unknown and we want to extract all possible candidates. For each body part, the confidence

map from the two networks are merged together, and we use non-maximum suppression to select all the points whose confidence scores are higher than their neighbors.

The output of this step are lists of body-part candidates. Formally, for a given frame at time t , all these body-part candidates can be represented as $p_j^{i,t} = (x, y)$, where j is the index of that body part, and $i \in \{1, 2, \dots\}$ indicates the count of all possible keypoints extracted for this body part. The total number of i is not determined because the number of cow objects is unknown at this stage, and there could be some incorrectly-detected candidates. All these candidates are further selected and clustered in the next step.

The confidence maps from the two networks are treated differently during this process. For keypoints from the upper body region, the two confidence maps are directly merged to find body parts. But since color information is more useful to detect the legs and hooves, only the confidence map from the color image network is used unless this map contains no candidates.

2) *Spatial clustering*: The second step in the post-processing module is spatial clustering. This step selects the correct body parts and clusters them into different cow objects. The first task before clustering is to determine the number of cows in the frame by counting cow centers. Given a set of extracted keypoint candidates $p_j^{i,t}$ from the upper body parts, the corresponding cow center positions can be estimated based on the constraints of the keypoints, shown in Equation 2.

$$c_j^{i,t} = F_j^{-1}(p_j^{i,t}) \quad (2)$$

Then a mean-shift clustering method is applied to the 2D spatial positions of all the cow centers $c_j^{i,t}$. Based on the clustering results of the center points, the corresponding body parts are labelled into separate cow objects. We ignore a cow object if the system cannot detect more than half of its keypoints.

The cow centers are also used to predict the location of missing body parts that the network fails to detect. After all keypoints are clustered into distinct cow objects, then for each cow object, we compute the averaged cow centers based on the detected points, and the miss-detected keypoints can be estimated using the keypoint constraints F_j . The predicted body parts based on the constraints may not be always accurate, but they provide a rough estimate of the cow's spatial location, which is useful for searching for keypoints in leg-hoof regions.

The final process in this step is to match the leg-hoof points. Similar to [6], we indicate the region of all possible leg-hoof points using a rectangle that is one-third wider than the rectangle of the upper body. Candidates outside this region will not be considered. The search process relies on the structural model. We follow the order of *shoulder/tailhead*, *leg*, *hoof* along each limb, and search the joints from among the candidates that lie in the search range. We also reject inappropriate points by applying the rule that each limb should have a certain rotation range; the angle between *shoulder* to *leg* and *leg* to *hoof* must be greater than 90 degrees for

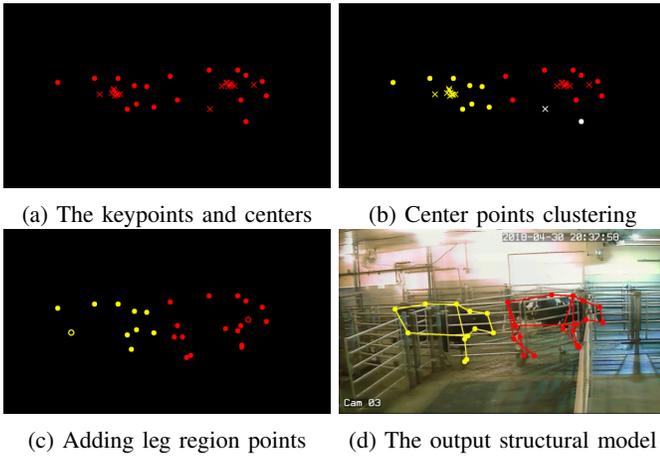


Fig. 4: The procedure of spatial clustering during post-processing. Circles represent the body parts p and crosses are the estimated cow centers c . Empty circles are the predicted body parts. Each color indicates a different cow object.

valid keypoints. Finally, all the selected leg-hoof joints are connected to the body contour to complete the final cow structural model.

Figure 4 illustrates the procedures of the spatial clustering step. The top left image shows the original extracted body parts from the previous step. The red circles are the extracted candidates and each is converted to a corresponding cow center, shown as crosses. Then in the top right image, all center crosses are clustered using mean shift to produce three clusters shown in distinct colors. Here the incorrect cluster (white) is eliminated because there are not enough candidates. Next in the bottom left image, points in the leg and hoof region are assigned to each cow object. Notice the empty circles are predicted points; the yellow one is blocked by the fences. Finally, by connecting all keypoints together, we form two cow structural models as shown in the bottom right image.

3) *Temporal filtering*: The final step in post-processing module refines the detection results using temporal information and matches cow objects across different frames. The two previous steps each operate on a single image, but the relationship between neighboring video frames is helpful to refine keypoint positions. It is reasonable to assume that the cows walk on an identical path between the fences and that they move steadily and slowly. This means that for a specific keypoint in the upper body, its trajectory over time should be smooth and any points far from the trajectory line can be considered outliers.

Based on this idea, we refine the positions of every upper body-part point across time to improve the temporal smoothness of the output. Before this step starts, all the frames in a video have been processed, so we know the number of cow objects in each frame. Then for every body part in the upper body region, we temporally filter each trajectory to remove and correct the outliers. In our experiment, we use a median filter, which is simple and provides robust prediction. Other filters

such as the Kalman filter do not work well especially when there are too many missing points from the previous steps. Notice that the leg-hoof region points are not involved in this process, since their trajectories are much more complicated.

Based on the trajectories of each cow object, the cow objects can be matched between neighboring frames. After this process, the system detects the total number of cows shown in a complete video sequence, and parameters about how every cow moves can be inferred, including the speed and rhythm [41].

V. EVALUATION METRICS

This section introduces our evaluation metrics. Although our method uses few ground-truth labels for training, ground truth is typically also required for performance evaluation. Therefore, in this paper we propose to use both supervised measures, which compare the detected results with ground-truth labels, and unsupervised measures, which directly evaluate the results without labels. Adding unsupervised measures to the evaluation process improves its thoroughness in the presence of insufficient labels. We first discuss the supervised measures for the cow structural model, and then introduce two unsupervised metrics.

A. Supervised measures

Quantifying the performance of the cow structural model requires more than the typical measures used to quantify object detection. As mentioned in Section I, the cow structural model is designed to provide two types of information: the spatial location of the body region, and the detailed positions of body parts. Both information is represented in terms of the keypoints of the cow body parts, and our ground-truth labels are also in terms of keypoints. As a result, we separately evaluate the area of the cow body region and the points in the leg-hoof region. Two metrics are developed and described below in detail: the *Body F1 score* and the *Leg-hoof F1 score*. In each case, the F1 score is harmonic mean of precision and recall when comparing the detection results to the ground truth. Notice that both metrics compare accuracy at the keypoint level. In a later experiment in Section VI-D, we also propose a method to convert the cow structural model to a binary mask with both body region and extended limbs, for the sole purpose of comparing our detected keypoint model with other mask-based segmentation methods.

1) *Body F1*: This metric measures the spatial area formed by the body region points. We connect the keypoints in the upper body region and generate one polygon mask for both the detected structural models and the ground-truth keypoints. Then we compare the two masks using the typical Intersection Over Union (IOU) metric and report the F1 score.

2) *Leg-hoof F1*: For the legs and hooves, a single pixel position represents each keypoint. However, physical joints typically extend for a larger spatial region. Therefore, the evaluation metric must accommodate this discrepancy, which may introduce systematic errors to both the labelling and detection process. For this reason, when measuring the distance between

ground truth and the detected leg and hoof keypoints, we set a threshold distance of 30 pixels. If the distance between the points is less than this threshold, we consider the joint to be detected, and points further away are considered to be missed. After thresholding, we determine how many leg-hoof points are successfully detected, and summarize this using the F1 score computed from the precision and recall. Since we do not create ground-truth labels for keypoints that are completely blocked by obstacles, these blocked joints do not affect the evaluation result.

B. Unsupervised measures

Unsupervised measures allow performance evaluation without ground-truth labels. This is particularly critical for video, where exhaustive application-specific labeling becomes even more onerous. Without labels, previously proposed metrics such as mean of region similarity, contour accuracy [42], and temporal stability metric [21] cannot be computed. Here, we apply prior knowledge to evaluate the performance when the ground-truth labels are not provided.

We consider two rules for the cow structural model. First, the spatial locations of the keypoints in a model should always form a cow-shaped object. Second, the shape of the cow body should be stable during the walk and the keypoints should have similar smooth trajectories. Based on these two constraints, we introduce two unsupervised metrics: the valid cow percentage and temporal consistency.

1) *The Valid Cow Percentage (VCP)*: This metric counts the fraction of detected cow models that are valid. Here valid means that the positions of the keypoints in the structural model can form a cow-shaped object. Like the supervised measure, we validate the upper body region and leg-hoof region separately.

For the upper body region, we use the trained keypoint constraints (Figure 2) as a reference, and compute the similarity between the detected contour and the reference using the Fréchet distance [43]. We choose this distance because it better captures the similarity between two curves, which are the body contour in our case. The computed distance is thresholded to form a binary decision whether the upper body region is valid or not. For points in the leg-hoof region, we define two interpretable rules to validate their spatial positions: all leg-hoof points should be lower than the body region points, and all hoof points should be lower than their corresponding leg points. If all leg-hoof points satisfy these two rules and the upper body region contour is also validated, the cow structure is considered valid.

This validation scheme is applied to all the detected cow objects in a video sequence, and the Valid Cow Percentage (VCP) is computed as the number of valid cow objects divided by the number of detected cows. The absolute VCP score is directly related to the actual number of cows in the testing video sequence, so the score is only meaningful when compared with other methods on the same testing dataset.

2) *Temporal Consistency (TC)*: The second unsupervised metric evaluates the Temporal Consistency (TC), which re-

flects the smoothness of the motion of moving objects in a video sequence. It is reasonable to assume that at a certain camera angle, the points from the body region always share the same translational motion because the shape of the cow body is stable. So ideally, the motion vector between every keypoint generated from one frame to the next frame should be the same. The Temporal Consistency (TC) metric evaluates this co-movement and computes the difference between the motion vectors generated by the body parts.

Formally, for each body part p_j^t in a cow object, we compute its motion vector from time t to $t + 1$ and summarize the variations d between all the motion vectors as

$$d^t = std(p_j^{t+1} - p_j^t), \forall j \in \{j_1, j_2, \dots\} \quad (3)$$

where std is the standard deviation, and j_i represents the index of the body part from the upper body region. Then the temporal consistency is computed as the average motion vector differences for all the frames in a video sequence.

$$TC = mean(d^t), \forall t \quad (4)$$

Notice this measure is applied to every individual cow object in a video sequence, and smaller TC values imply smoother object movements.

VI. EXPERIMENTS

This section presents the validation experiments. We first give a high-level summary of how we collect and prepare the video data from a commercial farm. Then we present three different experiments. The system output experiment compares the results of every stage in the proposed method, to demonstrate the importance of each component. Next, the dataset robustness experiment is performed on three different sets of video, to demonstrate the robustness of the method. Finally, we compare our method with other popular object segmentation methods, to demonstrate the advantages of our proposed method for cow detection.

A. Data collection

All cow videos in these experiments are collected from the Purdue Animal Sciences Research and Education Center located in West Lafayette, IN, USA from 2018 to 2019. All procedures were approved by the Institutional Animal Care and Use Committee (PACUC #1803001704). The cameras are mounted at fixed positions include a side-view of the path where cows walk every day. This path has fences on both sides and only allows one cow to walk through at a time. This limits the amount of cow-overlap; however the dense fences partly block the view of the cows, and some body parts are not visible behind the fences, as shown in Figure 1. This walking path is a typical component of many dairy farms.

During the course of data collection, we used three different capture devices: a commercial surveillance camera with Digital Video Recorder (DVR), a GoPro camera, and a high-quality IP camera. Table I shows the detailed information of the three video sets captured from the three cameras. The DVR videos have the worst quality with low frame rate and low

TABLE I: Summary of three sets of video data used in the experiments. The # Pixel per cow is in units of millions.

	Set 1	Set 2	Set 3
Capture Device	DVR	GoPro	IP camera
Video info	1280*720 @12fps	1232*384 @30fps	1920*1088 @30fps
# Pixels per cow	0.88m	0.29m	1.35m
Image Quality	low	low	high
Field of view	narrow	wide	wide
# cow per clips	single	multiple	multiple
# video clips (# for training)	87 (5)	18 (2)	114 (5)
# training frames	100	40	100
# testing frames	585	59	611

resolution. The GoPro videos provide higher frame rate, but they are spatially cropped with less spatial details. The IP camera captures high quality videos with both high frame rate and rich spatial information.

Table I compares several factors among the cameras that will influence detection performance. As noted, the video resolution and frame rate are different between the three sets, and Set 3 has the best quality. The number of pixels per cow refers to the average number of pixels that each cow occupies in an image, which is an indication of the spatial detail in each set. Notice that Set 2 only has 0.29 million pixels per cow, which is less than a third of the other two sets. The field-of-view each camera are also different. Set 1 videos only capture the center of the walking path where there are fewer fences, while the other two sets capture a wider view which includes two sides that have denser fences. In addition, the typical number of cows in one video are different across the sets. Narrow field-of-view videos normally captures a single cow in the frame, but the wider-angle videos could contain multiple cows, which challenges the detection method. In general, Set 3 has more video clips than the others with the greatest variety, so we will further divide this set into subsets in a later experiment described in Section VI-D.

To prepare the videos, we temporally segment the hours-long sequences into 10-second clips, on average, where all cows walk from left to right. In each set, we separate the clips into training and testing groups, where the number of training clips per set are shown in parentheses in Table I after the number of video clips. All multiple-cow clips are testing clips, so the training clips all contain only a single cow object. Non-consecutive frames are chosen randomly for labeling from both training and testing clips.

B. System component evaluation

This experiment compares all the internal outputs from our proposed system shown in Figure 3, to demonstrate the importance of each individual module. We choose the output of CNN1 as the baseline method, which is the original method in the DeepLabCut (DLC) toolbox [4]. However, this method can only detect one object per frame, so for a fair comparison, we only use the videos in Set 1 since these only contain one cow object. We compare the baseline method with four other

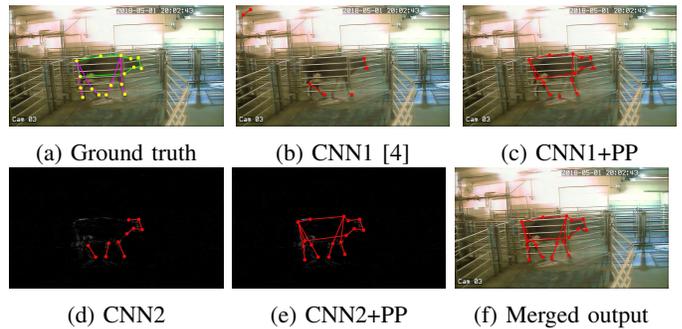


Fig. 5: The outputs of different stages in the proposed system.

internal outputs from the system: the CNN2 output from the difference videos, the CNN1 output plus the Post-Processing (PP) stages, the CNN2 output plus the PP stages, and the final merged result.

The implementation details are explained below. The frame difference images are generated by the sum of differences between the current frame and both the previous and next frame. The training labels from Set 1 are used to fine-tune both CNNs in the system. Recall that CNN1 processes the color images and CNN2 processes the frame-difference images. Both networks are pre-trained on ImageNet [40] and their final upsampling layers are fine-tuned with our cow images. For the two CNN methods without PP stage, we follow the extraction method from the DeepLabCut toolbox by setting a hard threshold and finding the location in the confidence maps with the maximum probability.

Both supervised and unsupervised evaluation metrics are used, but their testing data are different. For unsupervised measures, we compare the Valid Cow Percentage (VCP) and Temporal Consistency (TC) for all the frames in the testing videos because no labels are required. But for supervised measures, only the 585 labelled testing frames are used for evaluation. Among these labelled images, we report the body F1 score and leg-hoof F1 score, and the VCP score is also computed to compare the cow detection capability of each module in the system. Both qualitative and quantitative results are presented below.

Figure 5 shows an example of all five outputs of one testing image in Set 1. The direct outputs from the two CNNs without post-processing (top middle and bottom left) miss-detect some body parts, because they apply the strategy from the original DLC method that only selects one maximum point. Our proposed post-processing module uses non-maximum suppression to select all local maximum values from the confidence map, and all body-part candidates are detected (see bottom right of Figure 5). Considering the leg-hoof points, some joints of the swing leg are missed by CNN1 based on color image, because of motion blurriness and heavy compression. But these points are detected by CNN2 using the frame difference image, and the merged result generates a complete cow structural model.

The numerical comparison results are presented in Table II. In general, our complete system (last row) improves the

TABLE II: Comparison of the outputs of the system components on Set 1 videos (single-cow). Notice smaller TC value means smoother object movement in the video. Bold numbers show the best performance method in each column.

	Unsupervised		Supervised		
	VCP	TC	VCP	Body F1	Leg-hoof F1
CNN1 (DLC [4])	0.447	102.8	0.714	0.260	0.391
CNN2	0.408	155.0	0.673	0.366	0.252
CNN1+PP	0.632	8.92	0.846	0.772	0.373
CNN2+PP	0.667	10.19	0.929	0.841	0.333
Merged output	0.705	9.0	0.960	0.879	0.434

performance compared to the method in the DLC toolbox (first row). It can be observed that adding a Post-Processing (PP) module largely improves the system performance. The temporal and spatial prediction in the PP module improves the cow-detection ability demonstrated by the increasing VCP scores. Notice the two VCP scores from supervised measure and unsupervised measures are not comparable because their test sets are different. In addition, the temporal filtering process in the PP module largely improves the Temporal Consistency (TC), because the original CNN method purely operates on an image without considering temporal information. Comparing two F1 scores in the supervised measures, the PP step improves the detection accuracy for the cow structural model because more body-part candidates are selected from the intermediate CNN output.

Comparing the first two rows from the table, we can see CNN2 has better performance than CNN1 for the cow body region but works poorly on the leg and hoof regions. As explained in Section IV-C, CNN2 operates on gray-scale edges generated by the frame difference and better captures smoothly moving objects like the body region. But it cannot work in isolation because it eliminates too much information contained in the original images, such as the stationary legs. As a result, merging the two networks together obtains better detection for the leg-hoof region points.

C. Dataset robustness evaluation

This experiment evaluates the system robustness with different datasets. Training-based detection methods normally perform worse when they are applied to testing data that is substantially different from the training set. In this experiment, we evaluate the performance of our system when testing on frames collected from the three different cameras, that capture the same region of the farm but with different capture angles. This experiment also explores the influence of image quality on our system, since the video qualities from the 3 sets are also different.

For the training images in each video set, we fine-tune three detection systems, S_1 , S_2 , and S_3 , based on each individual corresponding datasets, respectively. An extra system S_{all} is trained on all the training frames together. In the testing phase, each trained system is applied to the images from the three sets separately. We also test each system on all testing images together for an overall comparison. All training and

TABLE III: System performance comparison on different video sets. The bold numbers show the best performance of each column.

Trained system	Body F1 score on				Leg-hoof F1 score on			
	Set1	Set2	Set3	All	Set1	Set2	Set3	All
S_1	0.80	0.42	0.51	0.64	0.61	0.18	0.35	0.46
S_2	0.72	0.65	0.58	0.65	0.16	0.59	0.33	0.26
S_3	0.82	0.56	0.59	0.69	0.61	0.52	0.56	0.58
S_{all}	0.82	0.64	0.61	0.71	0.62	0.65	0.56	0.59

testing data are separated regardless of their dataset, and no images used for both training and testing. In total, there are 4 trained models testing on 4 groups of test sets, which forms 16 training/testing pairs. For each pair, we measure the final system output using supervised metrics: body F1 score and leg-hoof F1 score. Table III shows the comparison results.

In Table III, each row represents a system trained from one dataset, and each column shows the system performance on one corresponding test set. Comparing the four systems, it can be observed that S_{all} achieves similar and slightly better performance than the others, and this merged system even works better than when each individual system is both trained and tested on its own videos (diagonal values). This demonstrates that adding training data from other similar video sets helps to improve the detection performance.

The results in Table III also allow us to examine the performance of the method when the input videos have different qualities. While both Set 1 and Set 2 have low quality, the images in Set 2 (see Figure 6) have a small spatial resolution while the images in Set 1 (see Figure 5) are blurry with poor illumination. Therefore, the results of system S_1 on Set 2 and of system S_2 on Set 1 images are poor, especially for the leg and hoof regions. However, system S_3 , which is trained on high quality images, provides better results on both these two datasets. This demonstrates that using higher quality images or increasing the variation of training data can improve system performance.

A final observation from the table is that the body region F1 scores are more stable across different systems than the leg-hoof F1 scores, due to the fact that the post-processing module that only operates on the body region. The spatial and temporal prediction in the post-processing model improve the estimation of missing and incorrectly detected points, which compensates for poor CNN performance. Since the legs and hooves are estimated directly from the CNN outputs, the performance variation is primarily due to the variation of training data.

In addition to numerical comparison, in Figure 6 we also present some visual results from all 4 trained systems applied to a test image from Set 2 that contains two cows. Comparing the outputs, system S_1 fails to detect two cow objects and S_3 is confused with some body parts between the two cow objects. However, system S_2 and S_{all} both detect two cow objects and present an accurate cow shape, because these two systems are both trained with data from Set 2. But the merged result from S_{all} is more accurate on some body parts, for example the points on each cow's back, because

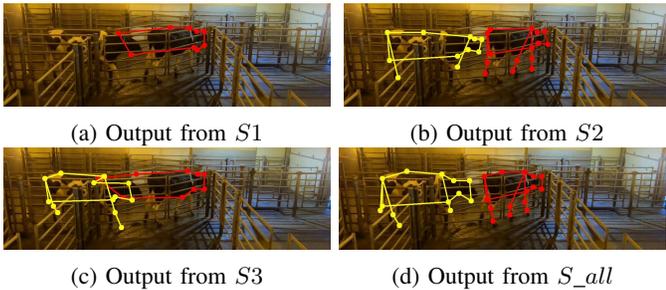


Fig. 6: The detection comparison between systems trained on different video sets. This example image is from Set 2.

of the additional training data involved. However, for the leg and hoof region, none of the systems detect all the points, due to the difficulty of observing them and the lack of post-processing process.

D. Segmentation methods comparison

This experiment compares the detection performance between our system and other popular object detection methods. Recall that the motivation for our system is not only to segment the spatial location of the cow, but also to detect critical keypoints about its body parts. Therefore, ideally comparison methods should also target these two goals. However, as mentioned in Section II, most previous keypoint detection methods focus on human objects and incorporate knowledge about human body parts, and it is difficult to adapt them to cow bodies for a comparison. On the other hand, there are many popular object detection methods which can be fine-tuned to segment cows, and these make for an effective comparison. In this experiment, we compare the cow object detection performance between our system and other three popular pixel-wise object detection methods: One Shot Video Object Segmentation (OSVOS) [15], DeepLab [13], and Mask R-CNN [12].

To create a performance comparison that does not disadvantage the object detection methods, we convert the output of our structural model into a binary cow mask, with two steps. First, all keypoints from the upper body region are connected to form a closed area representing the cow body. Second, every leg-hoof limb is expanded from a line into a polygon with a horizontal width of 20 pixels, as shown in the second column of Figure 7. This expansion process is applied to both the ground-truth labels and the detection results. The newly expanded ground-truth masks are then used to fine-tune the object detection methods, as well as to compute performance metrics. Still the point-to-mask conversion is not perfect. Notice the approximated masks cannot exactly cover the cow object from the original image; see for example the inaccurate edges of the cow body and the straight legs.

We use all the training and testing data from the three video sets in this experiment. In total, there are 240 single-cow frames for training and 1255 images for testing. Each of the three comparison methods are fine-tuned with the approximate

TABLE IV: Comparison of methods on different test sets.

	Set 1	Set 2	Set 3	All
OSVOS [21]	0.571	0.580	0.570	0.571
DeepLab [13]	0.655	0.513	0.577	0.610
Mask R-CNN [12]	0.735	0.692	0.630	0.682
ours	0.750	0.668	0.662	0.703

cow masks, with different implementation details. For OSVOS [15], we use the parent network pre-trained on the DAVIS 2016 [21] dataset and fine-tune it with our data. The output results are binarized using Otsu [44] threshold. For DeepLab [13], we use the pre-trained network from the COCO dataset [23], and we modify the last layer to produce two classes: cows and background. The fine-tuning process is applied only on the last atrous spatial pyramid pooling layers with binary entropy loss. For Mask R-CNN [12], we use the network pre-trained on the COCO dataset and fine-tune its region proposal network and feature pyramid network. The classifier outputs are also adjusted to the two classes of cows and background.

Figure 7 shows some visual examples of the detection results. From left to right are the original image, ground truth, and the results from OSVOS, DeepLab, Mask R-CNN, and our system, respectively. Each row shows an example which is selected from a different test set. Example (a) includes a human wearing black clothes who is walking right behind the cow. This confuses OSVOS which considers it to be part of the moving foreground object. Example (c) shows a special case which contains a pure white cow, and this color is not present in the training data. The DeepLab method completely misses the cow, because it directly extracts information from the color image and this rare color has not been seen before. The OSVOS method detects part of this cow using motion information, but Mask R-CNN works well because its region proposal network determines there is an object candidate and segments the cow object correctly.

Examples (b) and (d) contain multiple cows, and each method does detect multiple cow objects. However, the three masked-based methods merge all detected cow objects together because the objects are close to each other, and we need further effort to count the number of cows or to extract other detailed information. But our result provides a clear delineation between the cow objects, due to the use of the structural model. Another observation about these two examples is that the cow positions in these two images are different. Some cows are in the middle with fewer fences and others are on either the left or right side with denser fences blocking the view. Every method can detect the middle cow, but the cows on the sides are more challenging to detect due to the obstacles. We further analyze the influence of fences in later paragraphs.

Numerical comparison results among the methods are also reported using the F1 scores of the IOU between the detection results and the ground-truth masks. The measures are reported based on every test set separately in Table IV, and on distinct subsets of Set 3 in Table V.

From Table IV, it can be observed that our method achieves the highest accuracy for most sets, although its performance

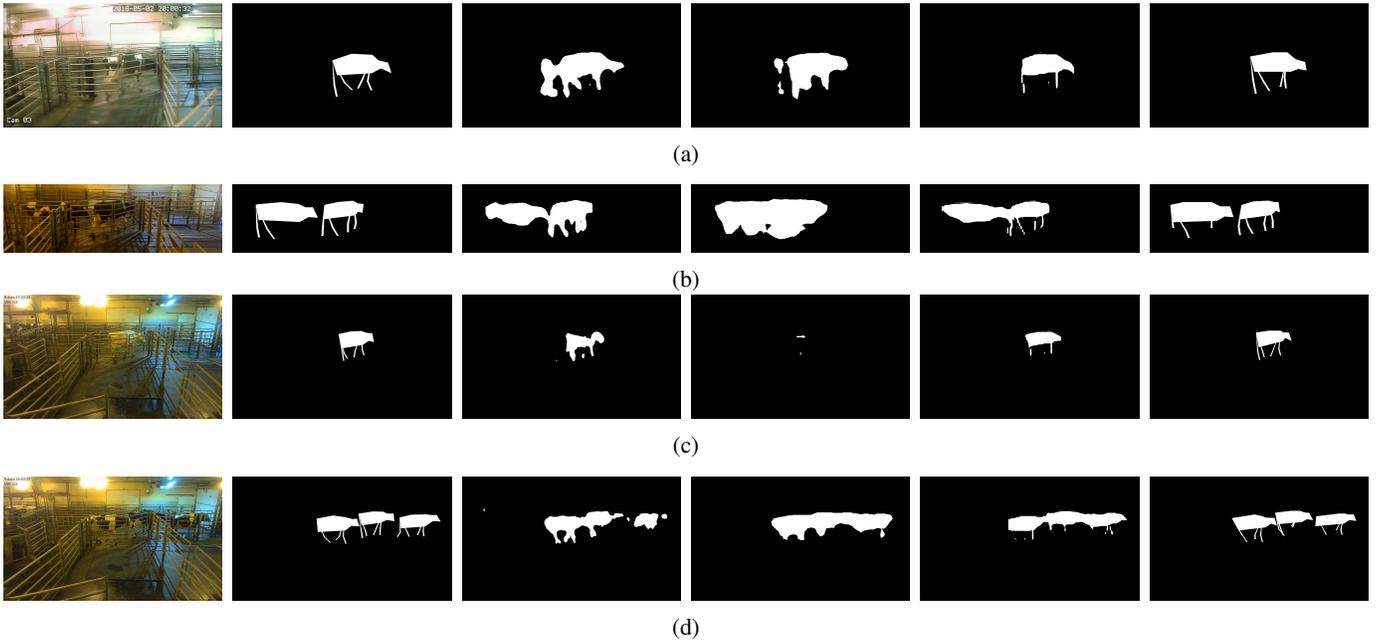


Fig. 7: Results using different detection methods. From left to right: original image, ground truth, OSVOS [15], DeepLab [13], Mask R-CNN [12], and ours. Example (a) and (b) are from Set 1 and Set 2, respectively; example (c) and (d) are both from Set 3.

TABLE V: Comparison of methods on subsets of Set 3. *Middle* means the cow is in the image center which has fewer obstacles, while *Side* means the cows are on the two sides with denser fences.

	Middle	Side	Single-cow	Multiple-cows
OSVOS [21]	0.672	0.589	0.650	0.547
DeepLab [13]	0.644	0.537	0.616	0.518
Mask R-CNN [12]	0.749	0.574	0.703	0.520
ours	0.734	0.645	0.711	0.587

relative to the fine-tuned Mask R-CNN is similar. There are three factors which may influence these scores. First, when comparing the masks using IOU, we use a merged mask containing both the cow body and leg regions. Since the body region occupies a larger area of the ground-truth mask, the IOU score can still be high even if the legs are miss-detected. Second, because the masks for our method and the ground truth are both converted from keypoints, it is highly sensitive to the positions of the keypoints, especially for the narrow leg regions. Small position shifts can lead to a large change to the converted mask, which will influence the IOU score. Third, when our system does not detect a leg or hoof point, the mask will be empty in this region. This will also decrease the IOU of our system. Nonetheless, our system performs well in comparison.

As mentioned above, a main consideration of our system is to obtain acceptable performance even when there are multiple cows, and when there are obstacles like fences. We use Set 3 videos to further explore the influence of these issues, to eliminate any performance variations due to video quality. As

Figure 7 shows, Set 3 images have a wider view of the walking path, and cows in the center have fewer fences while cows on the left or right sides are blocked with denser fences. So we separate the testing frames from Set 3 into four subsets: cows in the middle, cows on either side, single-cow frames, and multiple-cow frames. Among the four subsets, images with cows in the middle and with a single cow set will be easier than images from the other two subsets. The qualitative comparison F1 scores of these subsets are shown in Table V. From the table, Mask R-CNN has better performance on the easier test case when the cows are blocked by fewer fences. But for difficult test sets like denser obstacles, our proposed system works better. The OSVOS method also performs well when there are more obstacles because this method only considers the foreground and background, which allows it to separate the stationary fences from the moving cows.

In general, compared to the other three mask-based object detection methods, our proposed system has three advantages. First, based on the keypoints detection, our method can correctly detect the cow structure even when the cows are behind the fences or there are humans nearby. Second, when there are multiple cow objects, this system can explicitly isolate each cow even when they are close to each other. Third, it can detect cows with color patterns that do not exist in the training data through the use of frame difference images. However, our system also has two limitations. First, the cow structural model completely depends on the accuracy of the body parts, and one inaccurate detection can cause large errors for the body contour and influence the overall spatial location. Second, the prediction system in our method is based on the

keypoint constraints from the cow structural model, which is fixed after the training process. If there are not enough cow body parts detected, the prediction system still forces the results to conform to a particular shape, which could cause incorrect results.

VII. CONCLUSION

In this work, we design a practical system to detect the structural information for cows recorded in video. We use keypoints to form a cow structural model, which represents both the cow's overall spatial location and the positions of its specific body parts, such as the joints from the leg and hoof. The proposed detection system applies two CNNs to extract the keypoints from raw images, and a post-processing model is developed to select individual points and convert them into cow structural models. This system can detect and track multiple cow objects at the same time, and it also works with different quality videos which are captured on commercial farms during normal operation.

REFERENCES

- [1] L. Fleishman and J. Endler, "Some comments on visual perception and the use of video playback in animal behavior studies," *Acta ethologica*, vol. 3, no. 1, pp. 15–27, 2000.
- [2] N. B. Cook and K. V. Nordlund, "The influence of the environment on dairy cow behavior, claw health and herd lameness dynamics," *The Veterinary Journal*, vol. 179, no. 3, pp. 360–369, 2009.
- [3] S. Kumar and S. K. Singh, "Visual animal biometrics: survey," *IET Biometrics*, vol. 6, no. 3, pp. 139–156, 2016.
- [4] A. Mathis, P. Mamidanna, K. M. Cury, T. Abe, V. N. Murthy, M. W. Mathis, and M. Bethge, "DeepLabcut: markerless pose estimation of user-defined body parts with deep learning," Nature Publishing Group, Tech. Rep., 2018.
- [5] A. Pluk, C. Bahr, T. Leroy, A. Poursaberi, X. Song, E. Vranken, W. Maertens, A. Van Nuffel, and D. Berckmans, "Evaluation of step overlap as an automatic measure in dairy cow locomotion," *Transactions of the ASABE*, vol. 53, no. 4, pp. 1305–1312, 2010.
- [6] K. Zhao, J. Bewley, D. He, and X. Jin, "Automatic lameness detection in dairy cattle based on leg swing analysis with an image processing technique," *Computers and Electronics in Agriculture*, vol. 148, pp. 226–236, 2018.
- [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.
- [8] A. Ter-Sarkisov, R. Ross, and J. Kelleher, "Bootstrapping labelled dataset construction for cow tracking and behavior analysis," in *2017 14th Conference on Computer and Robot Vision (CRV)*. IEEE, 2017, pp. 277–284.
- [9] W. Wang, J. Shen, and F. Porikli, "Saliency-aware geodesic video object segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3395–3402.
- [10] Y.-H. Tsai, M.-H. Yang, and M. J. Black, "Video segmentation via object flow," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 3899–3908.
- [11] Y. J. Lee, J. Kim, and K. Grauman, "Key-segments for video object segmentation," in *IEEE International Conference on Computer Vision*, 2011, pp. 1995–2002.
- [12] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask R-CNN," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2961–2969.
- [13] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "DeepLab: Semantic image segmentation with deep convolutional nets, Atrous convolution, and fully connected CRFs," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 4, pp. 834–848, 2017.
- [14] J. Redmon and A. Farhadi, "YOLO v3: An incremental improvement," *arXiv preprint arXiv:1804.02767*, 2018.
- [15] S. Caelles, K.-K. Maninis, J. Pont-Tuset, L. Leal-Taixé, D. Cremers, and L. Van Gool, "One-shot video object segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 221–230.
- [16] K.-K. Maninis, S. Caelles, Y. Chen, J. Pont-Tuset, L. Leal-Taixé, D. Cremers, and L. Van Gool, "Video object segmentation without temporal information," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 6, pp. 1515–1530, 2018.
- [17] P. Tokmakov, K. Alahari, and C. Schmid, "Learning video object segmentation with visual memory," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 4481–4490.
- [18] J. Cheng, Y.-H. Tsai, S. Wang, and M.-H. Yang, "SegFlow: Joint learning for video object segmentation and optical flow," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 686–695.
- [19] P. Voigtlaender and B. Leibe, "Online adaptation of convolutional neural networks for video object segmentation," *arXiv preprint arXiv:1706.09364*, 2017.
- [20] N. Xu, L. Yang, Y. Fan, J. Yang, D. Yue, Y. Liang, B. Price, S. Cohen, and T. Huang, "Youtube-VOS: Sequence-to-sequence video object segmentation," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 585–601.
- [21] F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van Gool, M. Gross, and A. Sorkine-Hornung, "A benchmark dataset and evaluation methodology for video object segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 724–732.
- [22] M. Andriluka, L. Pishchulin, P. Gehler, and B. Schiele, "2D human pose estimation: New benchmark and state of the art analysis," in *IEEE Conference on Computer Vision and Pattern Recognition*, June 2014.
- [23] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: Common objects in context," in *European Conference on Computer Vision*. Springer, 2014, pp. 740–755.
- [24] A. Toshev and C. Szegedy, "DeepPose: Human pose estimation via deep neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1653–1660.
- [25] A. Newell, K. Yang, and J. Deng, "Stacked hourglass networks for human pose estimation," in *European Conference on Computer Vision*. Springer, 2016, pp. 483–499.
- [26] E. Insafutdinov, M. Andriluka, L. Pishchulin, S. Tang, E. Levinkov, B. Andres, and B. Schiele, "ArtTrack: Articulated multi-person tracking in the wild," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6457–6465.
- [27] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, "Realtime multi-person 2D pose estimation using part affinity fields," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 7291–7299.
- [28] Z. Cao, G. Hidalgo, T. Simon, S.-E. Wei, and Y. Sheikh, "OpenPose: realtime multi-person 2D pose estimation using part affinity fields," *arXiv preprint arXiv:1812.08008*, 2018.
- [29] L. Pishchulin, E. Insafutdinov, S. Tang, B. Andres, M. Andriluka, P. V. Gehler, and B. Schiele, "DeepCut: Joint subset partition and labeling for multi person pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 4929–4937.
- [30] T. D. Pereira, D. E. Aldarondo, L. Willmore, M. Kislín, S. S.-H. Wang, M. Murthy, and J. W. Shaveitz, "Fast animal pose estimation using deep neural networks," *Nature Methods*, vol. 16, no. 1, p. 117, 2019.
- [31] S. Günel, H. Rhodin, D. Morales, J. Campagnolo, P. Ramdya, and P. Fua, "Deepfly3D: A deep learning-based approach for 3D limb and appendage tracking in tethered, adult drosophila," *bioRxiv*, p. 640375, 2019.
- [32] S. Viazzi, C. Bahr, A. Schlageter-Tello, T. Van Hertem, C. Romanini, A. Pluk, I. Halachmi, C. Lokhorst, and D. Berckmans, "Analysis of individual classification of lameness using automatic measurement of back posture in dairy cattle," *Journal of Dairy Science*, vol. 96, no. 1, pp. 257–266, 2013.
- [33] X. Song, T. Leroy, E. Vranken, W. Maertens, B. Sonck, and D. Berckmans, "Automatic detection of lameness in dairy cattle: vision-based trackway analysis in cow's locomotion," *Computers and Electronics in Agriculture*, vol. 64, no. 1, pp. 39–44, 2008.
- [34] W. Andrew, C. Greatwood, and T. Burghardt, "Visual localisation and individual identification of Holstein Friesian cattle via deep learning," in

Proceedings of the IEEE International Conference on Computer Vision, 2017, pp. 2850–2859.

- [35] K. Zhao, X. Jin, J. Ji, J. Wang, H. Ma, and X. Zhu, “Individual identification of Holstein dairy cows based on detecting and matching feature points in body images,” *Biosystems Engineering*, vol. 181, pp. 128–139, 2019.
- [36] W. Shao, R. Kawakami, R. Yoshihashi, S. You, H. Kawase, and T. Naemura, “Cattle detection and counting in UAV images based on convolutional neural networks,” *International Journal of Remote Sensing*, pp. 1–22, 2019.
- [37] G. Aujay, F. Hétroy, F. Lazarus, and C. Depraz, “Harmonic skeleton for realistic character animation,” in *Proceedings of the 2007 ACM SIGGRAPH/Eurographics Symposium on Computer Animation*. Eurographics Association, 2007, pp. 151–160.
- [38] E. Insafutdinov, L. Pishchulin, B. Andres, M. Andriluka, and B. Schiele, “DeeperCut: A deeper, stronger, and faster multi-person pose estimation model,” in *European Conference on Computer Vision*. Springer, 2016, pp. 34–50.
- [39] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [40] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in Neural Information Processing Systems*, 2012, pp. 1097–1105.
- [41] H. Whay, “Locomotion scoring and lameness detection in dairy cattle,” *In Practice*, vol. 24, no. 8, p. 444, 2002.
- [42] X. Li, Y. Qi, Z. Wang, K. Chen, Z. Liu, J. Shi, P. Luo, X. Tang, and C. C. Loy, “Video object segmentation with re-identification,” *arXiv preprint arXiv:1708.00197*, 2017.
- [43] H. Alt and M. Godau, “Computing the Fréchet distance between two polygonal curves,” *International Journal of Computational Geometry & Applications*, vol. 5, no. 01n02, pp. 75–91, 1995.
- [44] P.-S. Liao, T.-S. Chen, P.-C. Chung *et al.*, “A fast algorithm for multilevel thresholding,” *JOURNAL OF INFORMATION SCIENCE AND ENGINEERING*, vol. 17, no. 5, pp. 713–727, 2001.