# Lameness detection in dairy cows using pose estimation and bidirectional LSTMs

Helena Russello [a,*], Rik van der Tol[a], Eldert J. van Henten[a], Gert Kootstra [a,*]

[a]*Agricultural Biosystems Engineering group, Wageningen University & Research, Wageningen, The Netherlands*

## Abstract

This study presents a lameness detection approach that combines pose estimation and Bidirectional Long-Short-Term Memory (BLSTM) neural networks. Combining pose-estimation and BLSTMs classifier offers the following advantages: markerless pose-estimation, elimination of manual feature engineering by learning temporal motion features from the keypoint trajectories, and working with short sequences and small training datasets. Motion sequences of nine keypoints (located on the cows' hooves, head and back) were extracted from videos of walking cows with the T-LEAP pose estimation model. The trajectories of the keypoints were then used as an input to a BLSTM classifier that was trained to perform binary lameness classification. Our method significantly outperformed an established method that relied on manually-designed locomotion features: our best architecture achieved a classification accuracy of 85%, against 80% accuracy for the feature-based approach. Furthermore, we showed that our BLSTM classifier could detect lameness with as little as one second of video data.

*Keywords:* lameness, cows, locomotion, pose-estimation, deep-learning, lstm

*Corresponding authors
   *Email addresses:* `firstname.lastname@wur.nl` (Helena Russello ), `firstname.lastname@wur.nl` (Gert Kootstra )

## 1. Introduction

Lameness is a prevalent condition in dairy cows, and is characterized by an abnormal gait due to lesions in their hooves or limbs. Lameness has negative welfare and economic impacts as it affects milk production, fertility, and life quality of the cows (Enting et al., 1997; Whay and Shearer, 2017). Lameness is commonly detected during visual locomotion scoring sessions where trained observers assess the lameness prevalence of the herd, but these are time-consuming and performed sporadically (O'leary et al., 2020). Automating locomotion scoring, e.g. by means of continuous camera monitoring, could allow for earlier lameness detection, and thereby timely treatment.

Video cameras are attractive sensors for automatic locomotion scoring because they are relatively inexpensive, non-intrusive, and scale well with large herds. Automatic lameness detection from videos is commonly approached by using locomotion traits from clinical gait scoring methods (Sprecher et al., 1997) to classify the degree of lameness in cows. For example, studies by Poursaberi et al. (2010); Viazzi et al. (2014); Van Hertem et al. (2018) used the posture of the back, Song et al. (2008) used the tracking distance, Wu et al. (2020); Zheng et al. (2023) used the step size, and Kang et al. (2020) used the supporting phase to classify lameness. Other studies, such as Zhao et al. (2023); Barney et al. (2023); Taghavi et al. (2024); Russello et al. (2024) combined multiple locomotion traits as input features for lameness classification. Combining the back posture, head position, tracking distance, and stride length led to an improved classification performance over using single features (Russello et al., 2024). Although models using locomotion traits have shown promising results in detecting lameness, they have some limitations. For instance, these locomotion traits assume low-noise data and may perform poorly with noisy data Taghavi et al. (2023). Additionally, manually selecting locomotion traits as features can restrict the information provided to such models and it may fail to capture complex patterns using only hand-crafted features.

In order to address these limitations, several studies have used deep neural networks to perform both feature extraction and lameness classification. Most of these methods consist of two to three steps: first isolating the body structure of the animal from video-frames, then performing feature extraction and lameness classification (the last two are often combined as one step). For instance, Karoui et al. (2021) tracked adhesive physical markers placed on the legs of cows, and trained a Convolutional Neural Network (CNN) to predict

lameness from the motion of said markers. Wu et al. (2020) extracted leg coordinates with YOLOv3 (Farhadi and Redmon, 2018) to create a step-size feature. Sequences of step-size vectors were then used in a Long-Short-Term Memory neural network (LSTM) and in a Bi-directional LSTM (BLSTM) to detect lameness. Arazo et al. (2022) applied a segmentation model alongside the SlowFast video-recognition model (Feichtenhofer et al., 2019) to classify lameness based on time-series of the body contour of the cows. Jiang et al. (2020) extracted sequences of optical-flow maps, which captured the motion between video-frames, and used it in a BLSTM network to classify lameness. Both Arazo et al. (2022) and Jiang et al. (2020) demonstrated that their models could classify lameness effectively using RGB video data alone (a one-step approach). However, they found that the classification accuracy was improved when using a two-step approach that incorporated segmentation masks (Arazo et al., 2022) or optical-flow maps (Jiang et al., 2020) to guide the model's focus on the body structure of the animal.

Even though the use of deep learning models reduces the need for extensive pre-processing and manual feature engineering, while increasing robustness, the multi-steps methods discussed in the previous paragraph still have inherent limitations. Leg-only tracking (Wu et al., 2020; Karoui et al., 2021) potentially overlooks valuable motion patterns from other body parts, whereas whole-body analysis (Jiang et al., 2020; Arazo et al., 2022) lacks specificity in tracking critical anatomical features, and its high dimensionality typically requires more training data, which makes the learning process more challenging.

To overcome these constraints, we propose a two-step approach that combines the advantages of pose estimation and BLSTM models (Graves and Schmidhuber, 2005). We use the T-LEAP pose estimation model (Russello et al., 2021), which allows tracking of multiple keypoints, including leg joints, head position, and spine movement, thereby providing the classification model with compact yet comprehensive motion data from bio-mechanically relevant keypoints. Building upon the work of Russello et al. (2024), which demonstrated the effectiveness of keypoint trajectories for computing hand-crafted locomotion features, we use BLSTM models to directly learn locomotion features from sequences of keypoint trajectories and classify lameness. A BLSTM is a type of neural network that excels at capturing long-term relationships in sequential data and modeling complex temporal patterns. BLSTMs have been successfully used in diverse application domains, includ-

ing human motion analysis (Du et al., 2019; Battistone and Petrosino, 2019), and also lameness detection (Jiang et al., 2020; Wu et al., 2020). In these application domains (Du et al., 2019; Jiang et al., 2020; Wu et al., 2020), BLSTMs have consistently outperformed regular (unidirectional) LSTMs. Therefore, we used BLSTMs for detecting lameness based on the keypoint trajectories of walking cows.

In summary, our proposed method offers the following advantages: the pose estimation provides a compact data representation of the spatio-temporal motion from multiple body parts, while the BLSTM model eliminates the need for noise filtering as well as manual feature engineering. Our contributions are as follows:

1. We present a two-step approach using BLSTM neural networks for detecting lameness from time-series data of nine keypoints. Our research explores whether detecting lameness through learned locomotion features outperforms previously established manually-designed locomotion traits. To ensure a fair assessment, we conduct a direct comparison with Russello et al. (2024) using the same dataset.
2. We evaluate the performance of our approach across multiple BLSTM architectures (different number of layers and layer sizes). Additionally, we evaluate the performance of the BLSTM classifiers across different sequence lengths (30, 60, and 90 frames, corresponding to 1, 2, and 3 seconds of video at 30 FPS), since healthy and lame cows tend to walk at different speeds (Flower et al., 2005).
3. To promote reproducibility and facilitate future research, we will make both our data and code publicly available upon publication.

## 2. Materials and Methods

Our approach, illustrated in Figure 1, uses T-LEAP to extract keypoint trajectories from videos of walking cows. The keypoint trajectories are cropped to a fixed sequence length, and passed to a BLSTM classification model. The data and approach are further described in the following subsections.

### 2.1. Dataset

The experiments used keypoint trajectories from the dataset introduced in Russello et al. (2024), comprising 272 videos of 98 individual Holstein-
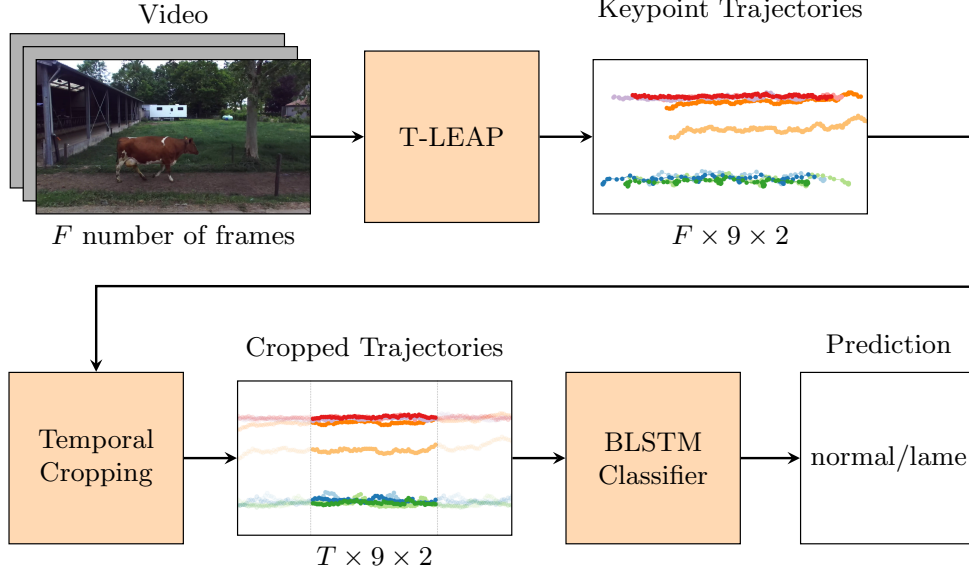
Figure 1: Outline of our lameness detection approach. The keypoint trajectories are extracted from videos using T-LEAP (Fig. 3 and 4). The keypoint trajectories are then trimmed to a fixed sequence length $T$ and passed to the BLSTM classifier.

Friesian cows. A ZED camera[1] was used to film the cows from the side as they walked freely through an outdoor walkway. The videos were recorded at 30 frames per second, with lengths ranging from 90 to 207 frames (mean length: 134 frames).

The dataset from Russello et al. (2024) consists of keypoint trajectories extracted from the 272 videos. For each video, four observers scored the gait using the Sprecher et al. (1997) scale. The gait scores were merged into binary lameness labels (*normal/lame*). Out of the 272 keypoint trajectories, 143 were labeled as *normal*, and 129 were labeled as *lame*, yielding a relatively balanced dataset.

The keypoint trajectories were automatically extracted with the T-LEAP pose estimation model introduced in Russello et al. (2021). In short, T-LEAP is a pose estimation model that was trained on videos of walking cows to track nine anatomical landmarks (keypoints) of the cow's hooves, head, and

---

[1]https://www.stereolabs.com/zed

back (Figure 2). The keypoint trajectories, illustrated in Figures 3 and 4, represent the motion in the 2D image-plane of the nine keypoints through each video frame. Figure 3 shows examples of keypoint trajectories for a normal gait, and Figure 4 for a lame gait. For a more detailed description of the dataset, we refer the reader to Russello et al. (2024).
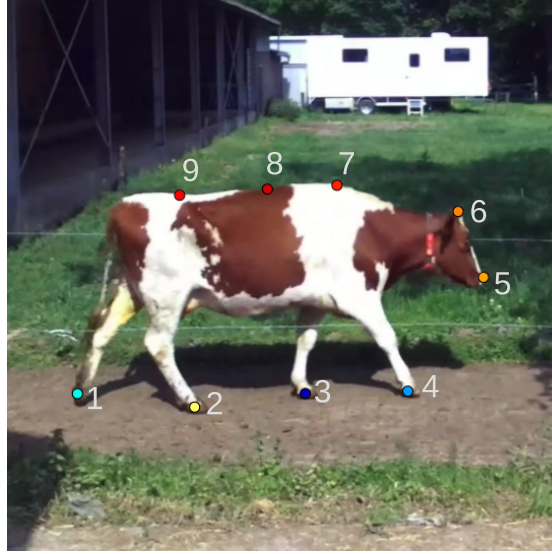


Figure 2: The 9 keypoints as described in Russello et al. (2024). The keypoints are named as follows: 1: Left-hind hoof, 2: Right-hind hoof, 3: Left-front hoof, 4: Right-front hoof 5: Nose, 6: Forehead, 7: Withers, 8: Caudal thoracic vertebrae, 9: Sacrum.
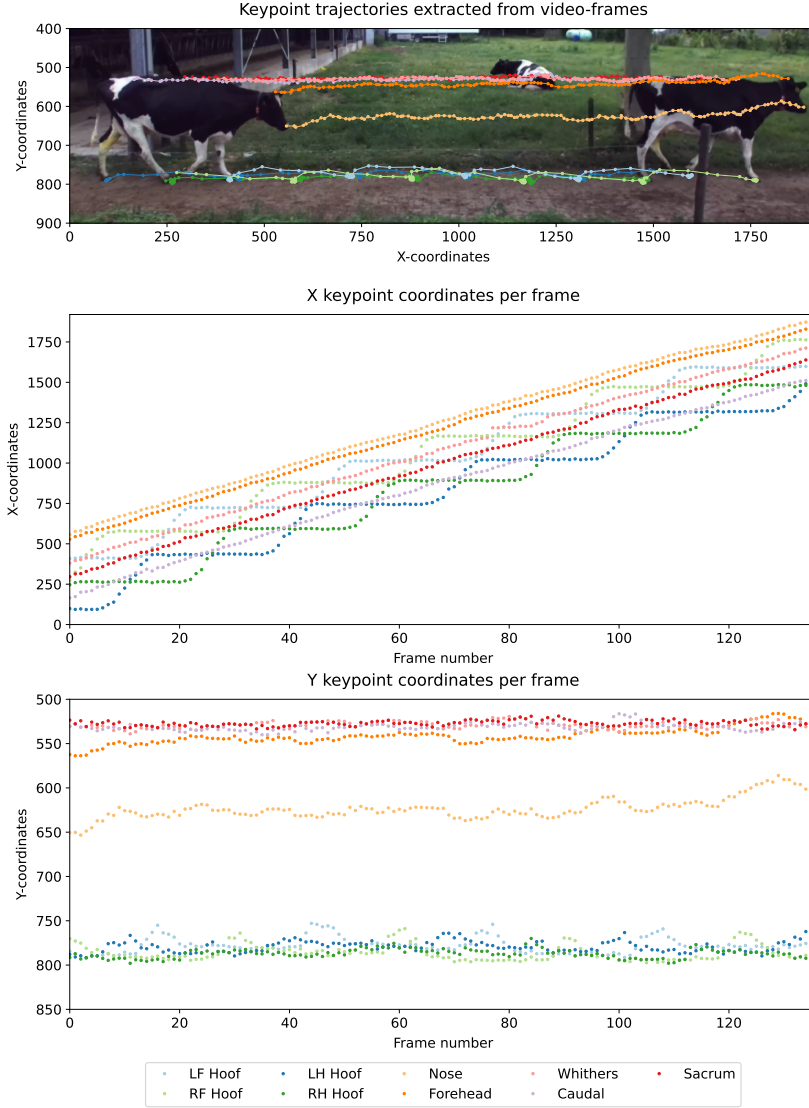
Figure 3: Example of the keypoint trajectories extracted with T-LEAP from a video of a cow presenting a *normal* gait. The top figure is augmented with the first and last frame of the video and shows the $(x, y)$ coordinates of each keypoint extracted from all the video frames. The middle and bottom figures show the keypoints' $x$ coordinates per frame, and $y$ coordinates per frame, respectively.
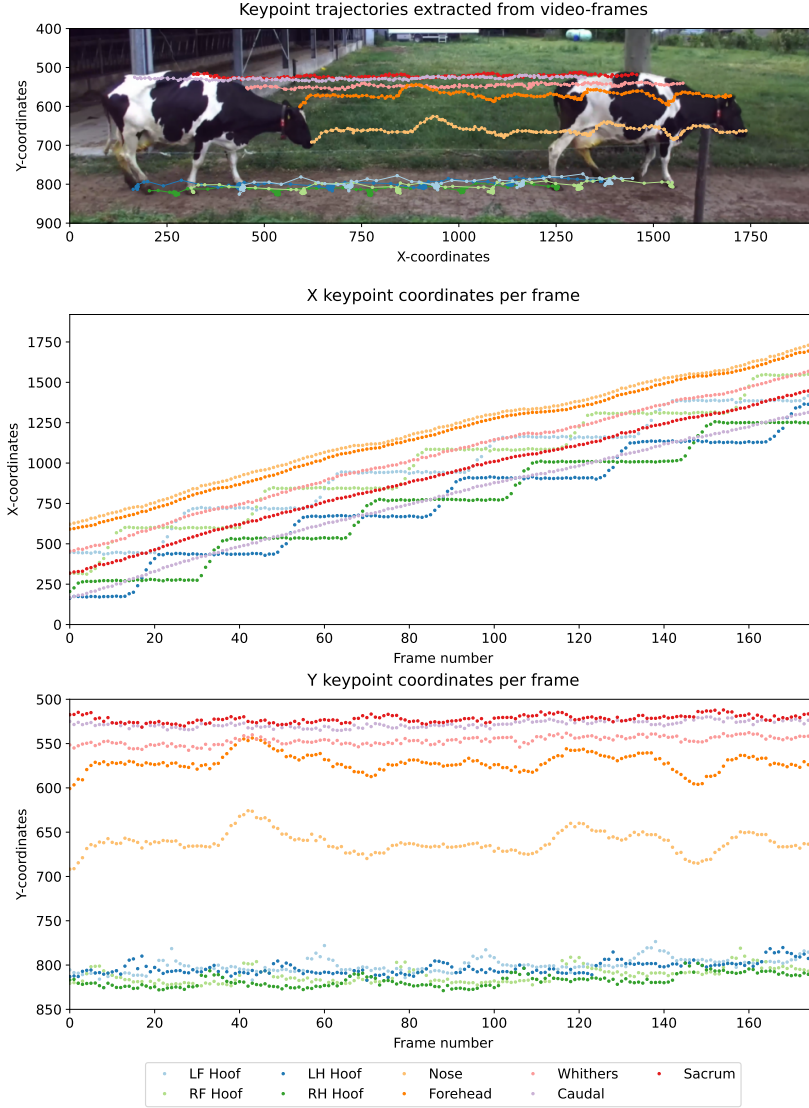
Figure 4: Example of the keypoint trajectories extracted with T-LEAP from a video of a cow presenting a severely *lame* gait. The top figure is augmented with the first and last frame of the video and shows the $(x, y)$ coordinates of each keypoint extracted from all the video frames. The middle and bottom figures show the keypoints' $x$ coordinates per frame, and $y$ coordinates per frame, respectively.

## 2.2. Data augmentation

Data augmentation was performed on the *training set* through temporal cropping and random jitter as described below. For the *test set*, no random jitter was applied and fixed sequences of 30, 60 and 90 frames long were extracted from the middle portion of the keypoint trajectories.

*Random Temporal cropping.* Given that keypoint trajectories varied between 90 and 207 frames, while the experiments at hand required fixed sequence lengths of $T \in \{30, 60, 90\}$ frames, random cropping was applied to the data, allowing fixed sequence lengths, while augmenting the data. That is, given a keypoint trajectory of $F$ frames, the starting frame was randomly selected from a discrete uniform distribution over $[0, F - T]$. This approach ensured that, for each training epoch, the model encountered different segments of the keypoint trajectories.

*Random jitter.* Similar to Karoui et al. (2021), we augmented keypoint coordinates with random Gaussian noise to increase model robustness. Specifically, we applied jitter to each keypoint coordinate by randomly sampling from a normal distribution where the mean was set to the original coordinate value and the standard deviation was defined as one percent of the head length (measured as the Euclidean distance between head and nose keypoints).

## 2.3. Lameness detection using BLSTMs

In this study, we used a Bidirectional Long-Short-Term Memory (BLSTM) classification model (Graves and Schmidhuber, 2005) to detect lameness based on the keypoint trajectories of walking cows. LSTMs (Hochreiter and Schmidhuber, 1997) are a type of Recurrent Neural Network (RNN) that excel at capturing long-term relationships in sequential data and learning complex temporal patterns. As their name suggests, RNNs use recurrent connections, that is, the output of a neuron at one time step loops back as an input to the neuron at the next time step. This allows to capture temporal dependencies and patterns within sequences. An LSTM is a type of RNN designed to overcome the limitations of RNNs and allow to capture long-term dependencies in sequential data (Hochreiter and Schmidhuber, 1997). An LSTM unit is typically composed of a cell-state controlled through three gates: a forget gate, an input gate, and an output gate. The forget gate controls what information to discard from the previous state, the input gate controls what new information to store in the cell-state , and the output gate controls what information of the cell-state to output. At each time step,

the LSTM takes the current input and previous hidden state, then uses these gates to selectively update its memory and produce an output. This gating mechanism allows LSTMs to learn long-term dependencies in sequences - they can remember important information from many time steps ago while forgetting irrelevant details (Gers et al., 2002). This makes them useful for analyzing gait patterns (Lefebvre et al., 2013; Battistone and Petrosino, 2019), since the relation between the sequential change in the pose of the animal and the presence of lameness is complex.

The main difference between unidirectional and bidirectional LSTMs lies in how they exploit the temporal information in input sequences. Unidirectional LSTMs process the input sequence in a single direction, either from the start to the end of the sequence (forward direction) or from the end to the start (backward direction). A forward LSTM only includes information from the past to predict the current timestep, whereas a backward LSTM includes information from the future to predict the current timestep. As their name suggests, bidirectional LSTMs combine a forward LSTM and a backward LSTM, both of which are connected to the same output layer. This means that for any timestep in a sequence, the BLSTM utilizes sequential information from the points before and after that timestep (Graves and Schmidhuber, 2005). This is useful if future features provide different information to the current timestep than past features, which is expected with asymmetric signals, such as gait.

### 2.4. BLSTM model architecture

The proposed BLSTM classification model, illustrated in Figure 5, receives sequences of keypoints as input and classifies the sequence as either *normal* or *lame*. The model architecture consists of two parts: a BLSTM neural network followed by a Fully Connected Network (FCN) (Figure 5). We empirically selected different BLSTM architectures, consisting of two or three layers of 128 or 256 hidden units ($h$). The BLSTM neural network presented in Figure 5 is composed of two BLSTM layers; our three-layers architecture has one additional BLSTM layer, but the FCN remains the same.

The input ($x$) consists of sequences of shape $T \times 18$, where $T$ corresponds to the length of the sequence, and 18 corresponds to the flattened vector containing the 2D image coordinates of 9 keypoints. The input is passed to the BLSTM sequentially, i.e. from $x_0$ to $x_T$. A BLSTM layer combines a forward and backward LSTM, each containing $h$ hidden units. The forward

LSTM processes the input sequence in the forward direction (from the first item in the sequence, $x_0$, to the last item, $x_T$) and the backward LSTM processes the input sequence in the backward direction (from $x_T$ to $x_0$). The outputs of the forward and backward LSTMs are concatenated and passed to the next BLSTM layer. The outputs of the last BLSTM layer are concatenated and given to the classification layer (FCN), which is composed of two fully-connected layers of input size $2h$ and $h$. The classification layer outputs a logit ($\hat{o}$), that is, an unnormalized prediction. During inference, a sigmoid function ($\hat{y} = \sigma(\hat{o})$) is used to normalize the output between 0 and 1, representing the probability of lameness. A cow is classified as lame if $\hat{y} \geq 0.5$.
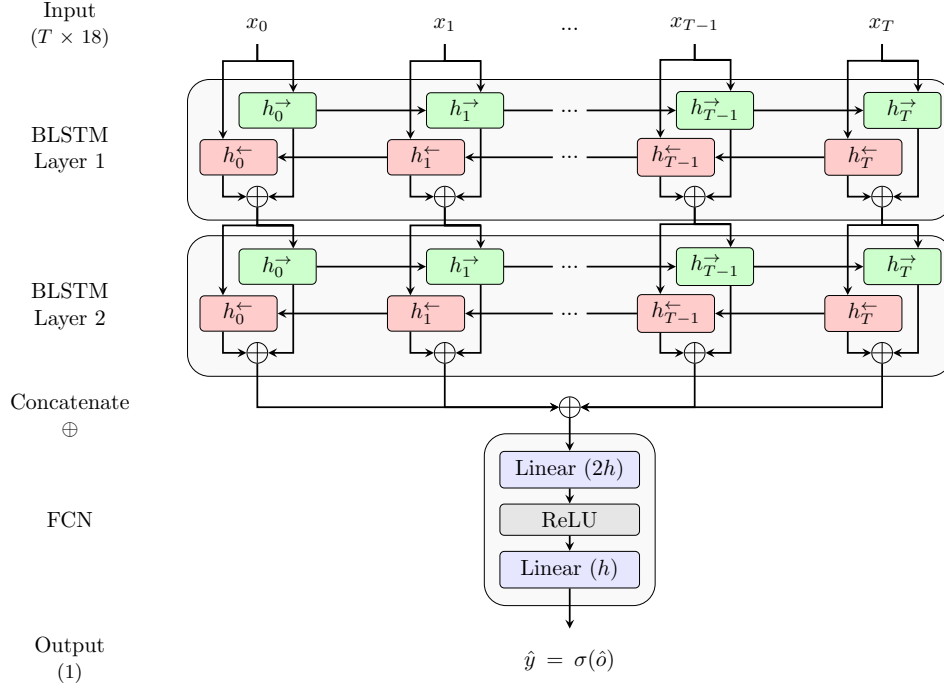


Figure 5: Architecture of the 2-layer BLSTM classifier, where $T$ corresponds to the sequence length (30, 60 or 90), $h$ corresponds to the number of hidden units (128 or 256). The green rectangles represent the forward LSTMs, and the red rectangles represent the backward LSTMs. Note that the 3-layer BLSTM architectures have an additional BLSTM layer, but their FCN remains the same.

11

## 2.5. Training procedure

The models were implemented using the PyTorch deep-learning framework (version 2.1.0) (Ansel et al., 2024). The experiments were run on a computer equipped with an AMD Ryzen 9 5900X CPU, an NVIDIA GeForce RTX 4070 Ti-SUPER GPU, and 64 GB of RAM.

Given the limited size of the dataset and to increase the reliability of the results, we followed the same procedure as Russello et al. (2024), and split the dataset into training and validation sets with a stratified 5-fold cross-validation with grouping. The folds were grouped on the ID of the cows, ensuring that individual cows appeared in the training or validation set, not both. To ensure an equal distribution of the classes in each training fold, we used PyTorch's weighted random sampler (Ansel et al., 2024) to address class imbalance through minority class oversampling. All models were trained with 5-fold cross-validation with batches of size 8 for 100 epochs. We used the AdamW-AMSGrad optimizer (Reddi et al., 2019; Loshchilov and Hutter, 2017) with a learning rate scheduler that reduced the learning rate by a factor of two every 50 epochs. Gradient clipping (Pascanu et al., 2013) with a threshold of 0.5 was applied to prevent vanishing and exploding gradients, a common problem with LSTMs (Bengio et al., 1994). The learning rate, weight decay, and dropout hyperparameters were tuned using a flat cross-validation approach. That is, the hyperparameters were first optimized per fold using a grid search over a 5-folds cross-validation, and the models were then re-trained on the same 5-folds cross-validation with the best set of hyper-parameters. We used flat cross-validation as it is computationally less expensive than nested cross-validation and generally results in the selection of an algorithm of similar quality to that selected via nested cross-validation (Wainer and Cawley, 2021).

## 2.6. Evaluation metrics

The performance of the models was evaluated using the following metrics: accuracy, macro F1-score, sensitivity, and specificity. The F1-score was macro-averaged; that is, the metric was calculated per class (*normal/lame*) and then averaged. The evaluation metrics were averaged over the five cross-validation folds. The differences between the different models were statistically tested using the McNemar's test with binomial distribution (McNemar, 1947), setting $p < 0.05$ as the level of significance.

### 2.7. Experiments

### 2.7.1. Different model architectures

We first explored different network architectures, with two and three BLSTM layers of 128 and 256 hidden units, later referred to as $2 \times 128$, $3 \times 128$, $2 \times 256$, and $3 \times 256$. Although these architectures were chosen empirically, they were inspired from Graves et al. (2013), who proposed, among others, BLSTM architectures for speech recognition consisting of 2 and 3 layers of 250 hidden units. These first experiments were conducted with sequences of $T = 90$ frames, since this was the length of the shortest video in the dataset.

### 2.7.2. Comparison with reference model

This study sought to compare the performance in detecting lameness through learned locomotion features (proposed model) to previously established manually-designed locomotion traits (reference model). We then performed a direct comparison with Russello et al. (2024), and compared our results to their best model, a Support Vector Machine (SVM) classifier trained with six locomotion features on the same dataset.

Briefly, Russello et al. (2024) used the same dataset of keypoint trajectories and applied outlier detection and filtering to correct and smoothen noisy keypoint trajectories. These keypoint trajectories were then used to compute six locomotion traits, namely back-posture, head-bobbing, tracking distance, stride length, stance duration, and swing duration. The locomotion traits were then used as input features to train multiple classifiers, such as decision trees or SVMs, to detect lameness.

Unlike in Russello et al. (2024), no filtering was applied to the BLSTM input data, since BLSTMs can typically handle some level of noise. To compare the impact of filtering, we trained the reference SVM classifier with locomotion traits computed from the filtered keypoint trajectories, as well as from the raw, unfiltered keypoint trajectories.

### 2.7.3. Different sequence lengths

The sequence duration is an important factor to take into account when detecting lameness, since healthy cows tend to walk faster (1.3 seconds per stride) than lame cows (1.5 seconds per stride) (Flower et al., 2005).

The first set of experiments were conducted with sequences of 90 frames (three seconds of video data), since this was the length of the shortest video in the dataset. We further extended our analysis by evaluating our models

on shorter sequences, namely sequences lengths of 30 and 60 frames, which corresponds to one and two seconds of video data, respectively. To this end, we re-trained and evaluated the best-performing BLSTM architecture on sequences of 30 and 60 frames.

## 3. Results

The results of the different BLSTM architectures and sequence lengths, compared to the SVM approach on the same dataset, are presented and discussed in the following subsections.

### 3.1. Different BLSTM architectures

The upper part of Table 1 presents the results of various BLSTM architectures using 90-frame sequences. The accuracy ranged from 83% to 85% and macro F1-scores between 83% and 84%. The sensitivity ranged from 78% to 83%, and the specificity ranged from 85% to 89%. The BLSTM $2 \times 256$ and BLSTM $3 \times 128$ architectures performed the best overall, and had an F1 score of 84%. The BLSTM $2 \times 256$ had a higher accuracy, while the BLSTM $3 \times 128$ had a higher sensitivity and specificity. The BLSTM $2 \times 128$ and BLSTM $3 \times 256$ architectures performed slightly worse than the other two. In general, there was minimal variation in performance between the different architectures, and only the BLSTM $2 \times 128$ was significantly different from the other architectures.

Table 1: Evaluation results of the different BLSTM architectures for sequence lengths of 90 frames, followed by the results of the SVM approach of Russello et al. (2024) on the same dataset. The results are presented in percents (%). The best results are highlighted in **bold** and the second best results are underlined.

| Model | N# frames | Accuracy | F1 (macro) | Sensitivity | Specificity |
|---|---|---|---|---|---|
| BLSTM $2 \times 128$ | 90 | 83.48 | 82.63 | 78.48 | **88.74** |
| BLSTM $3 \times 128$ | 90 | <u>84.47</u> | **83.70** | **82.74** | <u>87.78</u> |
| BLSTM $2 \times 256$ | 90 | **84.59** | **83.70** | <u>81.94</u> | 87.07 |
| BLSTM $3 \times 256$ | 90 | 83.11 | 82.37 | 81.77 | 85.07 |
| Gait features + SVM (filtering) | 124 ($\pm$ 24) | 80.07 | 78.70 | 76.78 | 81.15 |
| Gait features + SVM (no-filtering) | 124 ($\pm$ 24) | 75.17 | 73.38 | 70.70 | 76.26 |

### 3.1.1. Comparison with reference model

The lower part of Table 1, presents the results of Russello et al. (2024) for the same dataset. The evaluation results are displayed for their best model, that is, a radial-kernel SVM classifier, with and without keypoint filtering.

14

The performance of the SVM classifiers were lower across all metrics, for both approaches with and without keypoint filtering. For instance, without filtering on the keypoint trajectories, the F1-score was about nine percentage points lower than the BLSTM. When applying filtering, the performance of the SVM classifier increased dramatically, but remained significantly lower than that of the BLSTM with a F1-score lower by five percentage points. The predictions from the BLSTM models were all significantly different from the predictions of the SVM models.

### 3.2. Different sequence lengths

Table 2 presents the results of the BLSTM $3 \times 128$ model trained with sequences of 90, 60, and 30 frames. Using 90-frame sequences led to the best overall performance. Although the performance was lower accross all metrics with sequences of 60 and 30 frames, the difference was not statistically significant.

Table 2: Evaluation results of the BLSTM model trained with different sequence lengths. The results are presented in percents (%). The best results are highlighted in **bold**.

| Model | N# frames | Accuracy | F1 (macro) | Sensitivity | Specificity |
|---|---|---|---|---|---|
| | 90 | **84.47** | **83.70** | **82.74** | **87.78** |
| BLSTM $3 \times 128$ | 60 | 83.08 | 82.23 | 82.60 | 83.73 |
| | 30 | 83.07 | 82.23 | 82.07 | 83.91 |

## 4. Discussion

### 4.1. Lameness detection

The primary objective of this study was to investigate whether a learning-based end-to-end approach could improve lameness detection over conventional feature-based methods in recent literature (Russello et al., 2024). We conducted a direct comparison between our BLSTM approach and a feature-based SVM approach Russello et al. (2024), using identical experimental conditions. That is, we used the the same dataset of 272 cow videos, with the same labels, and same keypoint trajectories generated by the T-LEAP pose estimation model.

All proposed BLSTM architectures substantially outperformed the SVM model, with BLSTM $3 \times 128$ and $2 \times 256$ achieving the highest accuracy. These results suggest that learning temporal gait patterns directly through

BLSTMs offers improved performance compared to manual feature engineering approaches. A key finding emerged regarding robustness to noise: the SVM model required a prepossessing step to correct outliers and filter keypoint trajectories and its accuracy declined by 5 percentage points when using unfiltered data. In contrast, our BLSTM models achieved better results even when trained on unfiltered trajectories. This robustness to noisy keypoint detections suggests that BLSTMs would be more resilient to real-world challenges such as occlusions or corrupted video frames. Further research could explore the extent of this robustness by systematically evaluating BLSTM performance under different noise levels in the validation data. Such analysis could help establish minimum accuracy requirements for pose estimation models in lameness detection applications.

We approached lameness detection as a binary classification task, being either "*normal*" or "*lame*" walking, rather than multi-class locomotion scoring. The decision to merge the Sprecher et al. (1997) 5-level locomotion scale (normal (1), mildly lame (2), moderately lame (3), lame (4), severely lame (5)) into a binary classification (normal (1) vs. lame (2-5)) was motivated in Russello et al. (2024) as a way to address the dataset limitations: its relatively small size and skewed distribution toward lower lameness levels. As discussed in Russello et al. (2024), finer locomotion scoring would likely require a larger dataset with better representation of severe lameness cases. Furthermore, the performance of the classifiers may be limited by the quality of the ground-truth data, as some bias from the subjective nature of locomotion scoring likely persisted even after data curation.

The T-LEAP pose estimation model was originally trained on 17 keypoints (Russello et al., 2021), which included two more leg keypoints at the fetlock and carpal/tarsal joints. Here, we used nine of these keypoints (four hooves, three back keypoints, forehead, and nose) to represent the gait kinematics. We used this subset of keypoints to allow a direct comparison between our BLSTM approach and the SVM approach in Russello et al. (2024), who used these nine keypoints to compute hand-crafted locomotion features. Since BLSTMs can capture complex temporal patterns, incorporating additional keypoints could allow richer gait kinematics representations. For instance, including keypoints at the fetlock and carpal/tarsal joints would include information of the touch and release angles, as well as joint angle velocities, which may vary with lameness severity (Pluk et al., 2012). Investigating the benefits of additional keypoints for lameness detection is

16

left to future work.

The dataset used by Russello et al. (2024) contained videos that were selected when cows walked continuously, without distraction or interruption. As a result, the models were trained and evaluated under the assumption of uninterrupted locomotion. However, the impact of cows stopping on lameness detection performance remains unexplored. This potential selection bias might be a limitation, as real-world scenarios frequently involve cows pausing their locomotion. Future research should evaluate the degradation in model performance when cows exhibit natural behaviors. This would provide insights into the applicability of these detection systems in farm environments.

*4.2. Different sequence lengths*

The impact of sequence length on lameness detection was evaluated on the $3 \times 128$ BLSTM model, testing sequences of 30, 60, and 90 frames (corresponding to 1, 2, and 3 seconds of video at 30 FPS). While 90-frame sequences achieved the best performance, the results from shorter sequences were comparable, suggesting that lameness detection might be feasible with only one second of video data.

Looking at the stride duration in our dataset, healthy cows required 31 frames on average (1.03 seconds) to complete a stride (hoof strike to hoof strike), and severely lame cows requiring up to 45 frames (1.5 seconds). Although 30-frame sequences might not capture a complete gait cycle (i.e., one stride per leg), they appear to provide sufficient information for accurate detection. This can be attributed to several factors. First, even within a 30-frame window, the model can assess walking speed, which decreases with moderate to severe lameness (Blackie et al., 2013). Second, with slower, lame cows, these 30-frames sequences capture fewer strides, and while they might miss the stride of the affected leg, they might include compensatory movements in non-lame limbs that cows use to alleviate pain (Blackie et al., 2013).

Despite the promising results with shorter sequences, we recommend using longer sequences when possible. Shorter sequences risk missing critical gait anomalies, particularly given that cows, being stoic animals, often attempt to hide their pain. Longer sequences provide more opportunities to identify subtle anomalies in gait patterns and establish a more comprehensive assessment of the animal's locomotion.

Further research in understanding the minimum required sequence length could allow for an optimal camera placement in farms. For instance, in space-limited or cluttered environments, the animals might only be visible for short periods of time.

*4.3. Comparison with related work*

Performing direct comparisons between our approach and related work presents significant challenges due to variations in data, methodologies, and evaluation metrics. In subsection 4.1, we conducted a direct comparison with Russello et al. (2024) by using the same dataset and leveraging their open-source code, which enabled us to demonstrate our method's improvements objectively.

In the following paragraphs, we discuss our approach in relation to other directly relevant studies. However, these comparisons are constrained by the absence of publicly available data and code from these works. Without reproducing their experiments (which would be prohibitively time-consuming), direct one-to-one comparisons remain infeasible. Nevertheless, we compare our results and contrast our findings within these limitations. To advance reproducibility and promote open-source research practices, we commit to making our data and code publicly available upon publication and encourage similar transparency in future research.

Karoui et al. (2021) developed a lameness detection system based on keypoint trajectories obtained from physical reflective markers placed on the cows legs. They extracted four keypoint trajectories per leg, and generated additional synthetic data by adding random noise to the keypoint trajectories, resulting in a dataset of 24000 samples. The keypoint trajectories were used to train a LeNet CNN LeCun et al. (1998) for binary classification, achieving both accuracy and F1-score of 91%. Our approach differs in that rather than using physical markers, we used a markerless pose estimation model, which is less invasive and easier to deploy in practice. Furthermore, physical markers are more prone to skin-movement errors in kinematic data due to skin displacement while walking (Bergh et al., 2014). In addition to the leg keypoints, we also included keypoints on the back and head regions, allowing us to cover a broader range of gait patterns. Their experiments showed that generating new data by adding 5% noise variation yielded the best results, and while we similarly augmented our data using Gaussian noise with 1% standard deviation, we did not extensively study its impact. However, based on Karoui et al.'s findings, we believe this augmentation

18

technique contributed to improving our models' robustness and generalization capabilities.

Wu et al. (2020) developed a lameness detection system using YOLOv3 for leg localization in videos of walking cows. Their approach extracted the relative step sizes by calculating the horizontal distance between left and right legs in each frame. These sequences of step-size vectors were used as gait-features for various classification models. Their evaluation on 700 videos using 10-fold cross-validation achieved high accuracy metrics across different architectures: Decision Trees (90.5%), Support Vector Machines (94.3%), LSTM (98.6%), and BLSTM (99%). The performance gap between their SVM and BLSTM (4-5 percentage points) aligns with our findings. While their overall accuracy exceeded ours, this can be attributed to their larger training dataset: their 10-fold cross-validation used 630 videos for training, while our 5-fold cross-validation used, on average, 228 videos for training. When they reduced the training set to 350 videos using 2-fold cross-validation, performance decreased substantially (SVM: 82.9%, BLSTM: 86.7%). Furthermore, despite using significantly shorter sequences (30-90 frames versus their 500 frames), our system achieved similar performance. Our approach differs in the following aspects. While Wu et al. (2020) focused solely on the horizontal movements of the legs, we included both horizontal and vertical motion of keypoints from the legs, back, and head, providing richer spatial information. Furthermore, they first computed step-size vectors from the leg coordinates, and used this gait-feature as an input to their models, whereas we used the raw keypoint trajectories, thus eliminating the need for manual feature engineering and pre-processing.

Jiang et al. (2020) proposed a lameness classification approach combining optical-flow maps and a BLSTM classifier. Their study used a dataset of 1080 videos of walking cows (756 for training, 324 for testing) of 125 to 1000 frames long, with lameness severity scored on a 4-point scale (normal, slight, moderate, and severe). Sequences of optical-flow maps were generated from the videos frames, and a DenseNet CNN (Huang et al., 2017) was used to extract spatial features from each time-step of the optical-flow maps. A BLSTM classifier was then trained to classify 4-levels of lameness from sequences of spatial features, achieving a classification accuracy of 95%. Our approach differs primarily in how spatial data is extracted from video frames. We focused on tracking nine specific keypoint coordinates, while their method used optical-flow maps to capture motion information

from the entire body. Although their approach provided richer spatio-temporal information, it resulted in significantly higher input dimensionality to the BLSTM ($T \times 224 \times 224$ compared to our $T \times 18$). This increased complexity made their learning process more computationally demanding, requiring pre-training and a larger training dataset. Additionally, they studied the impact of sequence lengths on the classification accuracy, and had an improvement from 85% using 30-frame sequences to 95% using 60-frame sequences. Notably, they showed that longer sequences particularly benefited the detection of moderate and severe lameness cases, which can be attributed to the slower walking pace of lame cows. Our findings differ from theirs in that regard, as we observed minimal accuracy differences between the BLSTM models using 30, 60, and 90 frames. These contrastive findings may be attributed to differences in dataset composition. First, they performed classification across four, evenly balanced, lameness severity levels, whereas we performed a binary classification. Even though our dataset with binary labels was fairly balanced, the binary labels were merged from a five-point lameness scale, with skewed distribution towards non- and slightly-lame cows and with fewer severe cases. Merging the lameness levels decreased the granularity in the data, thereby potentially masking the benefits of longer sequences for detecting severe cases.

To summarize, this subsection compared our pose-estimation and BLSTM approach against three directly relevant studies: Karoui et al. (2021) used physical reflective markers on cow legs with LeNet CNN achieving 91% accuracy, but required invasive marker placement; Wu et al. (2020) employed YOLOv3 for leg localization and achieved up to 99% accuracy with BLSTM on 700 videos, though performance dropped significantly with smaller training sets; and Jiang et al. (2020) combined optical-flow maps with BLSTM for 4-level lameness classification, achieving 95% accuracy but requiring higher computational complexity and larger datasets.

## 5. Conclusion

In this paper, we developed a lameness-detection approach combining pose estimation and BLSTMs. Motion sequences of nine keypoints (located on the cows' hooves, head and back) were extracted from videos of walking cows with the T-LEAP pose estimation model. The trajectories of the keypoints were then used as an input to a BLSTM classifier that was trained to perform binary lameness classification. The experiments consisted of comparing our

proposed four BLSTM architectures to an established method relying on manually-designed locomotion traits, as well as comparing the performance across short sequence lengths (1,2 and 3 seconds of video data).

Across the proposed model architectures, the BLSTM $3 \times 128$ and $2 \times 256$ lead to the overall best performance, with an accuracy of 84.5%, against 80.1% accuracy for the feature-based approach. Furthermore, our pose-estimation and BLSTM approach achieved comparable performance to other studies that used BLSTM lameness classifiers (Wu et al., 2020; Jiang et al., 2020). BLSTM classifiers demonstrated similar performances across different sequence lengths. Our model achieved an accuracy of 84.5% when using 3-seconds sequences, and 83.1% when using both 2-seconds and 1-seconds sequences, suggesting that lameness can be detected with as little as one second of video data.

Combining markerless pose-estimation to extract the movement of multiple keypoints with BLSTMs to learn temporal locomotion features offered several advantages. Using raw keypoint trajectories eliminated the need for manual feature engineering, while including leg, back and head keypoints provided rich spatial information. Furthermore, our approach was more efficient than previous studies that used BLSTM lameness classifiers (Wu et al., 2020; Jiang et al., 2020), since we worked with with significantly shorter sequences and smaller training datasets.

While our approach offers a promising direction in automatic lameness detection from videos, further research directions include: (1) systematically evaluating BLSTM robustness under varying noise levels to establish minimum accuracy requirements for pose estimation models, (2) investigating the benefits of additional keypoints for improved detection, (3) exploring fine-grained lameness classification on a multi-point scale rather than binary classification, (4) evaluating the performance of lameness detection in "real-world" conditions, and (5) establishing the minimum duration requirements of video clips.

### Acknowledgments

**Declaration of interests**

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

**Declaration of generative AI and AI-assisted technologies in the writing process**

During the preparation of this work the authors used Claude (Sonnet 4) by Anthropic in order to refine selected sentences for clarity. After using this tool, the authors reviewed and edited the content as needed and take full responsibility for the content of the published article.

# References

Ansel, J., Yang, E., He, H., Gimelshein, N., Jain, A., Voznesensky, M., Bao, B., Bell, P., Berard, D., Burovski, E., et al. (2024). Pytorch 2: Faster machine learning through dynamic python bytecode transformation and graph compilation. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2*, pages 929–947.

Arazo, E., Aly, R., and McGuinness, K. (2022). Segmentation Enhanced Lameness Detection in Dairy Cows from RGB and Depth Video. arXiv:2206.04449 [cs].

Barney, S., Dlay, S., Crowe, A., Kyriazakis, I., and Leach, M. (2023). Deep learning pose estimation for multi-cattle lameness detection. *Scientific Reports*, 13(1):4499.

Battistone, F. and Petrosino, A. (2019). TGLSTM: A time based graph deep learning approach to gait recognition. *Pattern Recognition Letters*, 126:132–138.

Bengio, Y., Simard, P., and Frasconi, P. (1994). Learning long-term dependencies with gradient descent is difficult. *IEEE transactions on neural networks*, 5(2):157–166.

Bergh, A., Egenvall, A., Olsson, E., Uhlhorn, M., and Rhodin, M. (2014). Skin displacement in the equine neck. *Equine Veterinary Journal*, 46:36–36.

Blackie, N., Bleach, E., Amory, J., and Scaife, J. (2013). Associations between locomotion score and kinematic measures in dairy cows with varying hoof lesion types. *Journal of Dairy Science*, 96(6):3564–3572. ISBN: 0022-0302 Publisher: Elsevier.

Du, X., Vasudevan, R., and Johnson-Roberson, M. (2019). Bio-lstm: A biomechanically inspired recurrent neural network for 3-d pedestrian pose and gait prediction. *IEEE Robotics and Automation Letters*, 4(2):1501–1508.

Enting, H., Kooij, D., Dijkhuizen, A., Huirne, R., and Noordhuizen-Stassen, E. (1997). Economic losses due to clinical lameness in dairy cattle. *Livestock production science*, 49(3):259–267.

Farhadi, A. and Redmon, J. (2018). Yolov3: An incremental improvement. In *Computer vision and pattern recognition*, volume 1804, pages 1–6. Springer Berlin/Heidelberg, Germany.

Feichtenhofer, C., Fan, H., Malik, J., and He, K. (2019). Slowfast networks for video recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6202–6211.

Flower, F., Sanderson, D., and Weary, D. (2005). Hoof pathologies influence kinematic measures of dairy cow gait. *Journal of dairy science*, 88(9):3166–3173.

Gers, F. A., Schraudolph, N. N., and Schmidhuber, J. (2002). Learning precise timing with lstm recurrent networks. *Journal of machine learning research*, 3(Aug):115–143.

Graves, A., Mohamed, A.-r., and Hinton, G. (2013). Speech recognition with deep recurrent neural networks. In *2013 IEEE international conference on acoustics, speech and signal processing*, pages 6645–6649. Ieee.

Graves, A. and Schmidhuber, J. (2005). Framewise phoneme classification with bidirectional lstm and other neural network architectures. *Neural networks*, 18(5-6):602–610.

Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8):1735–1780.

Huang, G., Liu, Z., Van Der Maaten, L., and Weinberger, K. Q. (2017). Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708.

Jiang, B., Yin, X., and Song, H. (2020). Single-stream long-term optical flow convolution

network for action recognition of lameness dairy cow. *Computers and Electronics in Agriculture*, 175(April):105536. Publisher: Elsevier.

Kang, X., Zhang, X. D., and Liu, G. (2020). Accurate detection of lameness in dairy cattle with computer vision: A new and individualized detection strategy based on the analysis of the supporting phase. *Journal of Dairy Science*, 103(11):10628–10638. Publisher: Elsevier.

Karoui, Y., Jacques, A. A. B., Diallo, A. B., Shepley, E., and Vasseur, E. (2021). A Deep Learning Framework for Improving Lameness Identification in Dairy Cattle. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(18):15811–15812. Number: 18.

LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324.

Lefebvre, G., Berlemont, S., Mamalet, F., and Garcia, C. (2013). BLSTM-RNN Based 3D Gesture Classification. In *Artificial Neural Networks and Machine Learning – ICANN 2013*, pages 381–388. Springer.

Loshchilov, I. and Hutter, F. (2017). Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.

McNemar, Q. (1947). Note on the sampling error of the difference between correlated proportions or percentages. *Psychometrika*, 12(2):153–157.

O'leary, N., Byrne, D., O'Connor, A., and Shalloo, L. (2020). Invited review: Cattle lameness detection with accelerometers. *Journal of dairy science*, 103(5):3895–3911.

Pascanu, R., Mikolov, T., and Bengio, Y. (2013). On the difficulty of training recurrent neural networks. In *International conference on machine learning*, pages 1310–1318. Pmlr.

Pluk, A., Bahr, C., Poursaberi, A., Maertens, W., Van Nuffel, A., Berckmans, D., Pluk, A., Poursaberi, A., Berckmans, D., and Van Nuffel, A. (2012). Automatic measurement of touch and release angles of the fetlock joint for lameness detection in dairy cattle using vision techniques. *Journal of dairy Science*, 95(4):1738–1748. Publisher: Elsevier.

Poursaberi, A., Bahr, C., Pluk, A., Nuffel, A. V., Berckmans, D., Van Nuffel, A., and Berckmans, D. (2010). 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*, 74(1):110–119. ISBN: 0168-1699 Publisher: Elsevier B.V.

Reddi, S. J., Kale, S., and Kumar, S. (2019). On the convergence of adam and beyond. *arXiv preprint arXiv:1904.09237*.

Russello, H., van der Tol, R., Holzhauer, M., van Henten, E. J., and Kootstra, G. (2024). Video-based automatic lameness detection of dairy cows using pose estimation and multiple locomotion traits. *arXiv preprint arXiv:2401.05202*.

Russello, H., van der Tol, R., and Kootstra, G. (2021). T-LEAP: occlusion-robust pose estimation of walking cows using temporal information. *arXiv:2104.08029 [cs]*. arXiv: 2104.08029.

Song, X., Leroy, T., Vranken, E., Maertens, W., Sonck, B., and Berckmans, D. (2008). Automatic detection of lameness in dairy cattle-Vision-based trackway analysis in cow's locomotion. *Computers and Electronics in Agriculture*, 64(1):39–44. ISBN: 0168-1699 _eprint: 9809069v1.

Sprecher, D., Hostetler, D., and Kaneene, J. (1997). A LAMENESS SCORING SYSTEM

THAT USES POSTURE AND GAIT TO PREDICT DAIRY CATTLE REPRODUC-
TIVE PERFORMANCE. *Science*, (97).

Taghavi, M., Ouweltjes, W., Klandermans, B., and Kamphuis, C. (2024). Keeping an
eye on locomotion and behavior using computer vision. In *11th European Conference
on Precision Livestock Farming*, pages 1268–1275. European Association for Precision
Livestock Farming.

Taghavi, M., Russello, H., Ouweltjes, W., Kamphuis, C., and Adriaens, I. (2023). Cow
key point detection in indoor housing conditions with a deep learning model. *Journal
of Dairy Science*.

Van Hertem, T., Tello, A. S., Viazzi, S., Steensels, M., Bahr, C., Romanini, C. E. B.,
Lokhorst, K., Maltz, E., Halachmi, I., Berckmans, D., Schlageter Tello, A., Viazzi, S.,
Steensels, M., Bahr, C., Romanini, C. E. B., Lokhorst, K., Maltz, E., Halachmi, I., and
Berckmans, D. (2018). Implementation of an automatic 3D vision monitor for dairy
cow locomotion in a commercial farm. *Biosystems Engineering*, 173:166–175. ISBN:
1537-5110 Publisher: Elsevier.

Viazzi, S., Bahr, C., Van Hertem, T., Schlageter-Tello, A., Romanini, C. E. B., Halachmi,
I., Lokhorst, C., and Berckmans, D. (2014). Comparison of a three-dimensional and two-
dimensional camera system for automated measurement of back posture in dairy cows.
*Computers and Electronics in Agriculture*, 100:139–147. ISBN: 0168-1699 Publisher:
Elsevier B.V.

Wainer, J. and Cawley, G. (2021). Nested cross-validation when selecting classifiers
is overzealous for most practical applications. *Expert Systems with Applications*,
182:115222.

Whay, H. R. and Shearer, J. K. (2017). The impact of lameness on welfare of the dairy
cow. *Veterinary Clinics: Food Animal Practice*, 33(2):153–164.

Wu, D., Wu, Q., Yin, X., Jiang, B., Wang, H., He, D., and Song, H. (2020). Lameness
detection of dairy cows based on the YOLOv3 deep learning algorithm and a relative
step size characteristic vector. *Biosystems Engineering*, 189:150–163. Publisher: Elsevier
Ltd.

Zhao, K., Zhang, M., Ji, J., Zhang, R., and Bewley, J. M. (2023). Automatic lameness
scoring of dairy cows based on the analysis of head- and back-hoof linkage features
using machine learning methods. *Biosystems Engineering*, 230:424–441.

Zheng, Z., Zhang, X., Qin, L., Yue, S., and Zeng, P. (2023). Cows' legs tracking and
lameness detection in dairy cattle using video analysis and Siamese neural networks.
*Computers and Electronics in Agriculture*, 205:107618.