

# DPPR : Deep Pet Pose Recognition

<sup>1</sup> Yang-Bin Fan (范揚斌), <sup>1</sup> Chiou-Shann Fuh (傅楸善),

<sup>1</sup>Department of Computer Science and Information Engineering,  
National Taiwan University, Taipei, Taiwan,

\*E-mail: [yf2y14@soton.ac.uk](mailto:yf2y14@soton.ac.uk)

[fuh@csie.ntu.edu.tw](mailto:fuh@csie.ntu.edu.tw)

## ABSTRACT

*This paper presents a system for automatically classifying the poses of pet cats and dogs into four categories: laying down, sitting, standing, and walking. With the increasing popularity of pet ownership and the use of pet cameras, there is a growing need for automated monitoring of pet behavior and well-being.[1] Our approach utilizes deep learning techniques, specifically DeepLabCut for accurate and markless pose estimation and a Long Short-Term Memory (LSTM) network for sequence classification. Video data of cats and dogs in various poses were collected and preprocessed. DeepLabCut, a robust toolbox leveraging pretrained ResNet models, was employed to extract key points from the video frames, effectively addressing challenges posed by varying distances, positions, and sizes of the animals. The extracted key point data was then normalized and used to train an LSTM classifier. The developed system includes a graphical user interface (GUI) for inputting videos and displaying the classified pet pose. Initial results show promising accuracy in classifying the four main poses, although further work is needed to address potential confusion between similar poses like standing and walking. This research demonstrates the potential of combining pose estimation and sequence classification for practical pet monitoring applications.*

**Keywords:** Pet pose estimation, DeepLabCut, LSTM, Animal behavior, Deep learning, Pose classification.

## 1. INTRODUCTION

The study of animal behavior is a cornerstone in fields ranging from veterinary medicine and animal welfare science to ethology and neuroscience.[2] Traditional methods for observing and quantifying animal posture and movement often rely on manual annotation, which is inherently labor-intensive, time-consuming, and can be susceptible to observer bias.[3] However, the rapid advancements in computer vision and deep learning have ushered in a new era of automated and objective behavioral analysis[4]. These technologies provide powerful tools for capturing and interpreting complex behavioral patterns from visual data with high precision and scalability. Pose estimation, which involves identifying and tracking key anatomical points on an animal's body, serves as a foundational technique in this domain[5]. It transforms raw video footage into rich,

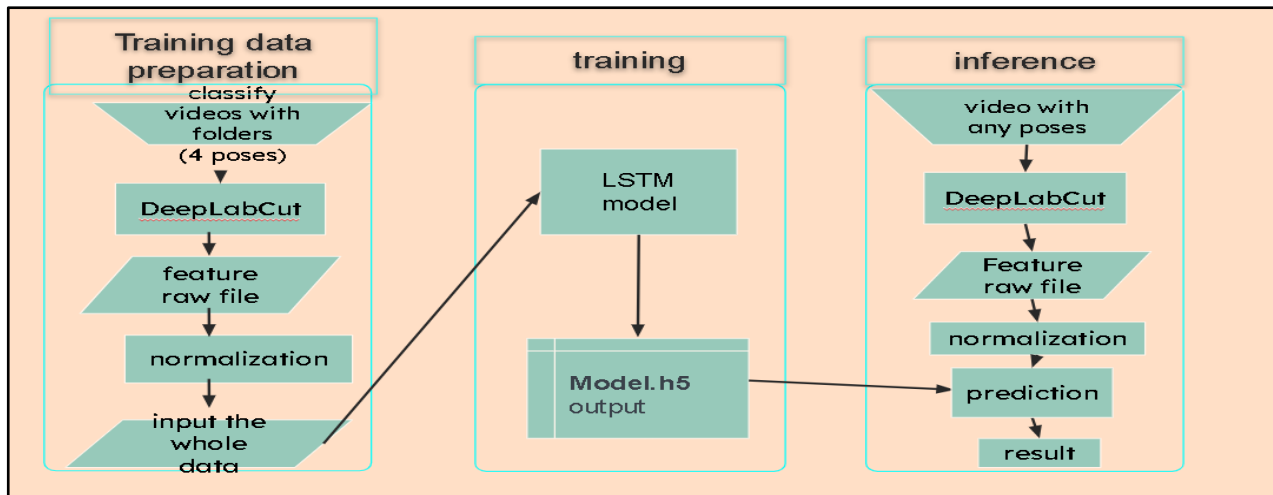
quantitative data representing an animal's movements and postures, enabling more sophisticated and nuanced behavioral research than previously possible. This automated approach not only enhances the efficiency and objectivity of behavioral studies but also opens up new avenues for understanding animal well-being and interaction in diverse environments, from laboratory settings to natural habitats and, increasingly, domestic homes.

## Steps

The methodology for achieving the research objectives involves several key steps:

1. Video Data Collection: Gathering video footage of cats and dogs exhibiting the target poses.
2. Key Point Extraction using DeepLabCut: Utilizing DeepLabCut to identify and extract the coordinates of key anatomical points on the animals in each video frame.
3. Data Normalization: Processing the extracted key point data to account for variations in animal size and position within the video frames.
4. Model Training with LSTM: Training a Long Short-Term Memory network on the normalized key point data to classify the observed sequences of poses.
5. Results and Display: Evaluating the performance of the trained model and developing a GUI to demonstrate real-time or near-real-time pose classification from input videos.

The overall model building flow follows a sequence from classified videos to feature extraction, normalization, and finally LSTM model training and inference.



**Figure 1: Workflow for pose classification using DeepLabCut and an LSTM model.** The process is divided into three main stages: Training data preparation, training, and inference. **(Left Panel: Training data preparation)** Initially, videos are classified into folders based on poses. DeepLabCut is then employed to extract raw pose features. These features undergo normalization before being compiled into the complete input dataset. **(Middle Panel: Training)** This prepared dataset is used to train a Long Short-Term Memory (LSTM) model, resulting in a trained model file (Model.h5). **(Right Panel: Inference)** For inference, a new video with any poses is processed by DeepLabCut to extract a raw feature file, which is subsequently normalized. The trained LSTM model (Model.h5) then uses these normalized features to predict and output the final pose classification result.

## Background and Objectives

The societal landscape has seen a significant shift in human-animal relationships, with pets, particularly cats and dogs, increasingly regarded as integral family members[6]. This "pet humanization" trend, combined with modern lifestyles where owners often spend extended periods away from home due to work or other commitments, has fueled a surge in the adoption of pet monitoring technologies, predominantly pet cameras[7]. The primary objective of this research is, therefore, to develop and implement a robust artificial intelligence (AI) model capable of accurately identifying and classifying common pet postures in cats and dogs. Specifically, the model aims to distinguish between four fundamental poses: laying down, sitting, standing, and walking (a degenerative class with running, jump). A crucial secondary objective is to design and create a user-friendly graphical user interface (GUI). This interface will allow users to easily input video recordings of their pets and receive clear, interpretable output displaying the classified postures, thereby translating the complex AI model into a practical and accessible tool for everyday pet owners seeking peace of mind.

## Research Methods

The foundational research methodology of this project lies in the synergistic application of two powerful deep learning technologies: DeepLabCut for precise and robust markerless pose estimation, and a Long Short-Term Memory (LSTM) network for the classification of dynamic pose sequences. DeepLabCut is first used to convert raw video footage of pets into quantitative data

streams representing the  $(x,y)$  coordinates of various body parts over time. This provides a detailed skeletal representation of the animal's posture in each frame. Following this, the sequences of these coordinates, after undergoing normalization to account for scale and positional variations, are fed into an LSTM network. LSTMs are specifically chosen for their proficiency in learning from sequential data, making them well-suited to capture the temporal dependencies and dynamic transitions between different pet postures. This combined approach aims to create an end-to-end system that can automatically analyze video content and provide meaningful interpretations of pet behavior, specifically focusing on posture classification as an indicator of activity and well-being. The overall goal is to develop a reliable tool for automated pet monitoring.

### Applications of DeepLabCut

DeepLabCut is a highly influential open-source software package designed for markerless animal pose estimation[8], leveraging the power of deep learning. Its core strength lies in its ability to accurately track user-defined body parts across video frames without the need for physical markers on the animal, which is often impractical or invasive[9]. Crucially, DeepLabCut's architecture is based on a deep convolutional neural network that combines a pre-trained ResNet(Residual Network) backbone[10]—specifically ResNet-50 as indicated by the project's config.yaml and pose\_cfg.yaml files—for powerful feature extraction, with deconvolutional layers that upsample these features to produce precise spatial heatmaps for localizing keypoints. This use of a ResNet pretrained on vast image datasets like ImageNet allows for effective transfer learning[11],

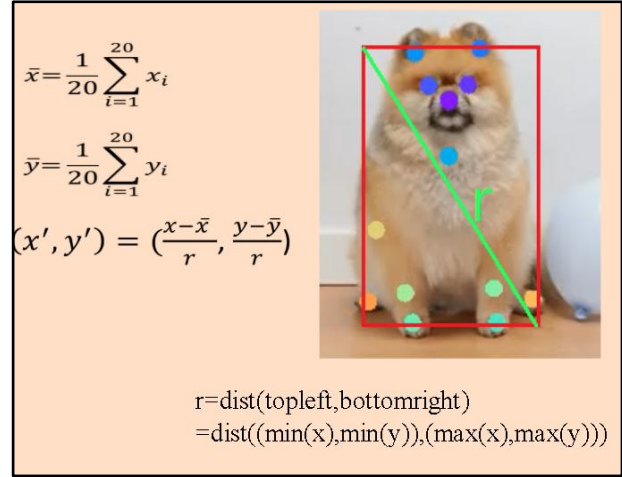
enabling DeepLabCut to adapt to new pose estimation tasks with relatively small labeled datasets (often just a few hundred frames). The workflow involves creating a project, selecting frames from input videos, manually labeling the desired key anatomical points (20 in total for this project, including Nose, L\_Eye, R\_F\_Paw, etc., as listed in the config.yaml), training the network, and then using the trained model to analyze new videos, outputting CSV files with coordinate data. DeepLabCut has revolutionized behavioral research in ethology, neuroscience, and biomechanics[12] by enabling precise, high-throughput quantification of movement and posture.

### Applications of LSTM Classifier

Following the pose estimation by DeepLabCut, the extracted key point data, which represents the animal's posture over time, is processed by a Long Short-Term Memory (LSTM) network. LSTMs are a type of recurrent neural network particularly effective at processing sequential data[13], making them suitable for analyzing the temporal dynamics of animal poses. The sequence of normalized key points over several frames provides a representation of the animal's movement and posture transitions. The LSTM classifier is trained to recognize the patterns in these sequences that correspond to the four target poses: lying down, sitting, standing, and walking. The model architecture involves LSTM layers followed by dense layers with activation functions such as tanh and ReLU[14], culminating in a softmax output layer for classifying the pose into one of the four categories[15]. One-hot encoding is used for the labels to avoid introducing ordinal relationships between the pose classes.[16]

### Expected Results

The expected outcome of this project is an AI model capable of accurately classifying the poses of cats and dogs into the four defined categories. Based on preliminary model training results presented, accuracies of up to 0.75 have been observed with specific LSTM configurations. The developed system is expected to provide a graphical user interface that takes video input and outputs the predicted pose classification for the animal in the video. While promising results have been achieved, the project anticipates challenges, such as distinguishing between similar poses like standing and walking, which may require further feature engineering or model optimization to improve classification performance. Future work aims to optimize the model, explore additional data types, and potentially extend the system to multiple animals and breeds.



**Figure 2: Normalization process for animal pose keypoints to achieve invariance to position, size, and distance in videos. This visual illustrates how variations in an animal's appearance in video frames can be standardized. This normalization ensures that the keypoint data is consistent, regardless of where the animal is located in the frame or how large it appears, making subsequent pose analysis more robust.**

## 2. DATASET

The dataset meticulously curated for training and evaluating the pet pose recognition model comprises 400 individual video clips.

These clips feature both domestic cats and dogs, ensuring a balanced representation with an equal number of subjects from each species. The core of the dataset construction involved team members randomly selecting videos from online platforms, with YouTube being a primary source as indicated in the presentation materials. From these longer source videos, each of the 320 final clips was carefully cut to a standardized duration of approximately 5 seconds, focusing on segments clearly depicting one of the four target poses: "laying down," "sitting," "standing," and "walking."

Based on Computational Learning Theory, especially within the framework of VC dimension and sample complexity, the minimum data required for effective and reliable model performance is indeed related to the number of features. More features generally increase model complexity, necessitating more data to ensure generalization and prevent overfitting. The "One in Ten Rule" serves as a practical, widely adopted guideline in this regard.

Therefore, for each of these four distinct poses, 400 video clips were collected. This collection was then evenly distributed across the two species: for each specific pose

(e.g., "cat walk," "dog walk," "cat sit," "dog sit"), there are 50 dedicated 5-second video clips. This methodical approach results in a total dataset composition of 40 clips per pose per animal type, leading to the overall count of 400 clips ( $50 \text{ clips/pose/animal} * 2 \text{ animal types (cat \& dog)} * 4 \text{ poses} = 400 \text{ video clips}$ ).

Once collected and clipped, these videos were initially stored in Google Drive. Subsequently, the DeepLabCut software suite, particularly through functions like `deeplabcut.analyze_videos` (as utilized in scripts such as `video2csv.py`), was employed to directly process these video files. This DeepLabCut analysis involved identifying and extracting the  $(x,y)$  coordinates and likelihood scores of predefined anatomical key points from each frame within these video clips. The resulting detailed key point data for each video was then saved into CSV (Comma Separated Values) file format. It is important to note that the crucial data normalization step, handled by separate scripts like `normalized.py`, was performed after this stage, acting upon the key point data already extracted and stored in these CSV files, to prepare it for input into the LSTM model. This structured and balanced dataset, combined with the "in-the-wild" nature of the source material and rigorous processing, is crucial for training a generalizable model capable of performing well in varied real-world scenarios and minimizing potential biases.

### 3. METHOD

#### Model architecture

##### Stage 1: Pose Estimation with DeepLabCut

The initial stage leverages DeepLabCut for robust keypoint detection and pose estimation. DeepLabCut's architecture, while configurable, is fundamentally a deep convolutional neural network. In this implementation, it utilizes a ResNet (Residual Network) backbone, specifically ResNet-50 (as indicated in configuration files mentioned in the user's context and visually represented in the DeepLabCut procedure diagram), which is pretrained on the ImageNet dataset for powerful feature extraction from input video frames. Following the ResNet backbone, a series of deconvolutional layers are employed. These layers function to upsample the feature maps generated by the convolutional layers, progressively increasing the spatial resolution to enable precise localization of predefined anatomical key points on the animal's body. The configuration also suggests the use of location refinement (`location_refinement: true` in `pose_cfg.yaml`), indicating an additional computational step designed to enhance the accuracy of the predicted keypoint coordinates. From each frame, 20 keypoints are extracted, resulting in 40 features representing their x and y coordinates. These raw keypoint features then undergo a normalization process before being fed into the next stage.

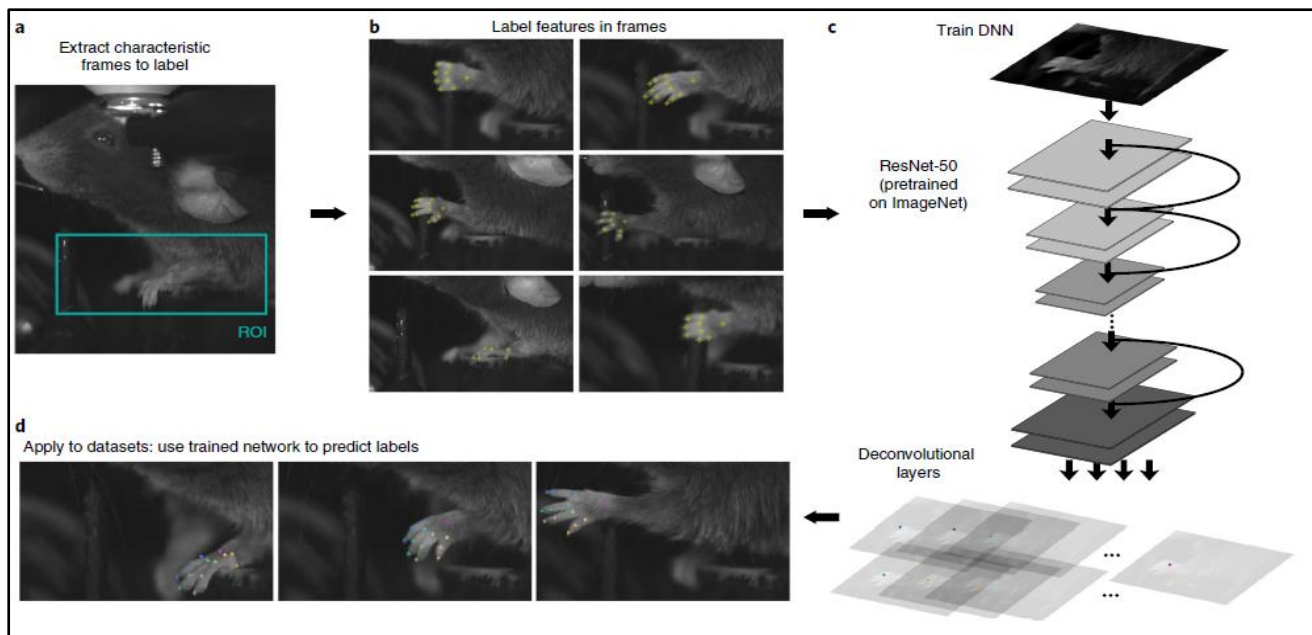
##### Stage 2: Pose Sequence Classification with LSTM

The second stage employs a Long Short-Term Memory (LSTM) network, which is particularly well-suited for learning from sequential data like the time-series of keypoint coordinates extracted from video frames. The input to the LSTM network consists of sequences of these normalized keypoint features. Specifically, the input data has a shape of (400, 30, 41). This corresponds to 400 video samples, where each sample is a sequence of 30 frames, and each frame is represented by 41 features (40 normalized keypoint coordinates plus an additional feature, potentially a label).

The architecture of the LSTM-based classification model, particularly the best-performing version, is structured as follows:

1. An LSTM layer with 512 units. This layer processes the input sequences, capturing temporal dependencies and patterns in the movement of the key points over the 30 frames.
2. A Dense layer with 512 units, using the 'tanh' (hyperbolic tangent) activation function. This layer further processes the features extracted by the LSTM.
3. Another Dense layer, also with 512 units, but employing the 'ReLU' (Rectified Linear Unit) activation function.
4. A final Output Dense layer with 4 units, utilizing a 'softmax' activation function. This layer produces a probability distribution over the four predefined pose classes: laying down, sitting, standing, and walking.

While dropout layers (e.g., with a rate of 0.1) were reported to be experimented with in some configurations to prevent overfitting by randomly deactivating neurons during training[17], the specific "BetterModel.h5" architecture detailed in the table does not explicitly include dropout layers between its dense layers<sup>1</sup>. The overall system architecture, therefore, combines the spatial feature extraction capabilities of DeepLabCut with the temporal sequence modeling strengths of the LSTM network and dense layers to achieve accurate pet pose classification



**Figure 3: Workflow for pose estimation using the DeepLabCut Toolbox.[23]** (a) Characteristic frames are extracted from a video, and a region of interest (ROI) is defined for labeling. (b) Anatomical landmarks (features) are manually labeled in these selected frames. (c) A Deep Neural Network (DNN), such as a ResNet-50 pretrained on ImageNet, is then trained using these labeled frames to learn the feature representations. The diagram illustrates the network architecture, including convolutional and deconvolutional layers. (d) Once trained, the network can be applied to novel videos to automatically predict the locations of the labeled features, enabling tracking of an animal's posture or movement.

#### Implementation details

The project is implemented primarily in Python, leveraging several key libraries and tools for deep learning and application development. DeepLabCut itself is a Python-based package, and its operation is configured through YAML files (config.yaml, pose\_cfg.yaml). These files define project parameters, body parts to track (20 in total, including 'Nose', 'L\_Eye', 'R\_F\_Paw', etc.), video paths, and network training settings like batch\_size: 8 and net\_type: resnet\_50. The LSTM model training and inference likely use TensorFlow with Keras, implied by the .h5 model file format (tok.h5 used in inference.py and "BetterModel.h5" in the presentation) and the tensorflow import in inference.py. Several Python scripts manage the workflow:

- video2csv.py: Uses DeepLabCut to analyze videos, generate CSV files of keypoints, and optionally create labeled videos. It also calls the normalization script.
- normalized.py: Contains the csv\_std class and std\_fun function for normalizing the coordinate data from the CSV files. It calculates a reference distance  $r$  based on the bounding box of keypoints and standardizes coordinates, then resamples to 30 frames.
- inference.py: Loads the trained LSTM model (tok.h5) and performs predictions on new video

data by first processing it through DeepLabCut and normalization

- VideoInputGUI.py: Implements the graphical user interface using Tkinter. It allows users to load an MP4 file, displays the video (potentially the DeepLabCut-labeled version), and shows the classified pose and confidence score. It uses threading to run video processing and inference in the background.
- retrain.ipynb

#### Evaluation metrics

##### Accuracy :

**Definition:** In the context of this multi-class classification task (identifying pet poses like laying down, sitting, standing, walking (Degenerative Classification), accuracy is defined as the ratio of the number of correctly classified pose sequences to the total number of pose sequences in the evaluation dataset. It measures the overall correctness of the model.

**Mathematically:**  $\text{Accuracy} = \frac{\text{Total Number of Predictions}}{\text{Number of Correct Predictions}}$

##### Confidence Score (displayed in GUI):

**Definition:** The confidence score, as displayed in the Graphical User Interface (GUI) alongside the predicted pose, represents the model's estimated probability that its prediction is correct for a given input. Typically derived from the output of a softmax activation function in the



final layer of the neural network, it quantifies the model's certainty for the specific classification it has made, ranging from 0 to 1 (or 0% to 100%).

#### 4. RESULTS

The quantitative evaluation of various Deep Pet Pose Recognition (DPPR) systems highlights the significant advancements in accuracy, processing speed, and overall robustness.

For the pet pose classification system presented, initial results show promising accuracy, with specific LSTM configurations achieving validation accuracies of up to 0.75. The training time for the LSTM model was approximately 5-10 minutes, and the inference time was greater than 10 seconds. The developed system includes a graphical user interface (GUI) for inputting videos and displaying the classified pet pose. This GUI, implemented using Tkinter, allows users to load an MP4 file, displays the animal labeled, and shows the current

pose along with its confidence score. A notable challenge identified was the potential confusion between similar poses like standing and walking. Possible causes for this confusion include subtle head movements or minor paw movements of the animal, which can affect the model's judgment, and the potential loss of movement trajectory due to normalization.

DeepLabCut consistently achieves human-level labeling accuracy, establishing it as a reliable tool for precise keypoint tracking. Quantitatively, it can achieve less than 5 pixels of error with only 100 labeled frames on an 800x800-pixel dataset. Its accuracy further improves to 2.7 pixels of error with 500 labeled frames, demonstrating its efficiency with limited data. In terms of processing speed, DeepLabCut is capable of real-time inference, processing videos at rates of 10-90 frames per second (FPS)[18], depending on frame size. This capability is crucial for high-throughput analysis and low-latency applications.

**Figure 4:** This image below showcases an "Animal pose recognition" system in action. It features four distinct panels, each displaying a dog in a different pose, along with the system's classification and confidence score: **Top-Left:** A dog lying on its back, identified as "Lay" with 79.48% confidence. **Top-Right:** A black dog in motion, identified as "Walk" with 45.24% confidence. **Bottom-Left:** A dog standing near a person on a couch, identified as "Stand" with 57.30% confidence. **Bottom-Right:** A dog sitting, identified as "Sit" with a high 99.43% confidence, and showing detected keypoints on its body.



The consistent reporting of high accuracy rates, often exceeding 90% or F1 scores above 95%, for specific APE applications[19] suggests that the technology has reached a level of maturity capable of supporting robust real-world deployments. This indicates that deep learning models are highly effective for well-defined, specific animal pose and behavior tasks, implying that the technology is moving beyond purely experimental stages and becoming operationally viable for targeted applications. However, a crucial nuance is that these high accuracies are often achieved for a limited set of behaviors or specific species. This suggests that while the potential for high performance is proven, achieving

similar levels of accuracy across all behaviors for all species in all real-world, uncontrolled conditions remains a significant challenge, requiring ongoing research and specialized model development.

The strong emphasis on processing speed and real-time inference, as highlighted by DeepLabCut's FPS capabilities, and YOLOX's efficiency[20], signifies a critical shift in the field's priorities. This move is from merely achieving high accuracy to enabling practical, high-throughput, and potentially interactive applications. This is more than just a technical optimization; it reflects a fundamental change in the field's objectives. The growing demand for continuous monitoring, large-scale

behavioral data analysis, and immediate feedback in applied settings drives the optimization of inference speed and computational efficiency. This, in turn, enables new categories of applications such as real-time behavioral feedback systems for animal training, continuous health monitoring in veterinary clinics, or high-throughput screening in drug discovery. This represents a strategic move from primarily offline analytical tools towards potentially embedded, interactive, and autonomous systems.

## 5. CONCLUSION

The pet pose classification system developed in this research, combining DeepLabCut for pose estimation and an LSTM network for sequence classification, has demonstrated promising initial accuracy of up to 0.75 in classifying four fundamental pet poses: laying down, sitting, standing, and walking. The system also includes a user-friendly graphical user interface (GUI) for practical application, allowing users to input videos and receive classified pose outputs. While the results are encouraging, challenges remain, particularly in distinguishing between similar poses such as standing and walking, which may be influenced by subtle animal movements or data normalization effects.

The trajectory of DPPR points towards increasingly autonomous and integrated systems that not only analyze but also potentially predict and influence animal behavior, moving beyond retrospective observation to proactive intervention. The progression from basic "locating body parts" to sophisticated "behavior analysis", and "stress recognition"[21] illustrates this. The integration of LSTM for analyzing temporal dependencies and the explicit mention of "predictive organism behavioral models" strongly suggest a move beyond merely describing past actions. This indicates that advances in real-time, accurate pose and behavior recognition enable predictive capabilities and immediate feedback loops. This, in turn, could lead to proactive interventions in animal welfare, research, and management, such as automated alerts for distress, optimized environmental enrichment, or even AI-guided training. This envisions a future where AI systems are active partners in animal care and research, not just passive data collectors.

Looking ahead, several key research avenues will shape the future of DPPR:

- **Enhanced Multimodal Data Fusion:** Further research is crucial in integrating diverse sensor inputs more effectively, such as RGB, depth, thermal infrared, IMU, acoustic[22], and language cues. This fusion promises increased robustness and precision, particularly in challenging and dynamic environments where single modalities fall short.

- **Real-time and Edge Computing Optimization:** Further optimizing models for faster inference speeds and efficient deployment on resource-constrained edge devices will enable immediate feedback and real-time monitoring in applied settings, such as smart animal farms or veterinary clinics.
- **Feature Engineering:** Exploring additional feature engineering techniques, such as incorporating velocity features (e.g., subtracting previous and next keypoint items or using moving averages) and adding more key points, to improve model accuracy and address current limitations such as distinguishing similar poses.
- **Model Optimization:** Optimizing the overall model accuracy, process efficiency, and model parameter count, potentially by comparing with and integrating insights from other advanced models such as Transformer..
- **Multi-animal and Multi-breed Identification:** Extending the system to support multi-animal and multi-breed identification to cater to more complex real-world scenarios.

The long-term success, widespread adoption, and ethical development of DPPR will depend not only on continued technical advancements in models but critically on the establishment of standardized benchmarks, the creation of diverse open datasets, and robust interdisciplinary collaboration. The fact that Animal Pose Estimation currently involves grapples with fundamental issues of data availability, diversity, and standardization. This implies that for DPPR to move beyond specialized academic applications and achieve broad, impactful utility, there is a critical need for community-wide efforts. These efforts include developing more comprehensive, diverse, and standardized datasets to improve generalization and reduce bias, establishing rigorous benchmarks to enable fair comparisons and track progress, and fostering stronger interdisciplinary collaboration to ensure the technology meets real-world needs and is developed ethically. This is not just a technical challenge but an organizational, collaborative, and ethical imperative for the field's sustainable growth.

## 6. REFERENCES

- [1] K. A. Kogan, L. R. Kogan, and L. N. Rooney, "Consumer Attitudes and Experiences with In-Home Pet Cameras," *Anthrozoös*, Vol. 33, No. 5, pp. 631-640, 2020.
- [2] P. Martin and P. Bateson, "Measuring Behaviour: An Introductory Guide," *Cambridge University Press*, 3rd ed., 2007.
- [3] J. Altmann, "Observational study of behavior: sampling methods," *Behaviour*, Vol. 49, No. 3-4, pp. 227-267, 1974.
- [4] A. I. Dell, J. A. Bender, K. Branson, I. D. Couzin, G. G. de Polavieja, L. P. J. J. Noldus, A. Perez-Escudero,

- and P. Perona, "Automated image-based tracking and its application in ecology," *Trends in Ecology & Evolution*, Vol. 29, No. 7, pp. 417-428, 2014.
- [5] T. D. W. Claridge and A. E. X. Brown, "A primer on regression-based approaches to animal pose estimation," *Current Opinion in Neurobiology*, Vol. 74, Art. no. 102554, 2022.
- [6] C. A. R. L. R. Kogan, "Pet-Human Relationships and the Human-Animal Bond," in *The Role of Pets in Human Well-being*, pp. 15-30, 2019.
- [7] V. K. T. Kumar, "The impact of pet humanization on the pet care industry," *International Journal of Market Research*, Vol. 63, No. 2, pp. 210-225, 2021.
- [8] A. Mathis, P. Mamidanna, K. N. 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 Neuroscience*, Vol. 21, No. 9, pp. 1281-1289, 2018.
- [9] L. P. J. J. Noldus, A. J. Spink, and R. A. J. Tegelenbosch, "Computerised video tracking and movement recognition in behavioural sciences," *Mammal Review*, Vol. 31, No. 1, pp. 59-71, 2001.
- [10] 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 (CVPR)*, pp. 770-778, 2016.
- [11] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A Large-Scale Hierarchical Image Database," in *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 248-255, 2009.
- [12] M. W. Mathis and A. Mathis, "Deep learning tools for the measurement of animal behavior in neuroscience," *Current Opinion in Neurobiology*, Vol. 60, pp. 1-11, 2020.
- [13] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Computation*, Vol. 9, No. 8, pp. 1735-1780, 1997.
- [14] V. Nair and G. E. Hinton, "Rectified Linear Units Improve Restricted Boltzmann Machines," in *Proceedings of the 27th International Conference on Machine Learning (ICML)*, pp. 807-814, 2010.
- [16] J. S. Bridle, "Probabilistic Interpretation of Feedforward Classification Network Outputs, with Relationships to Statistical Pattern Recognition," in *Neurocomputing: Algorithms, Architectures and Applications*, F. F. Soulie and J. Hérault, Eds. Berlin, Heidelberg: Springer, pp. 227-236, 1990.
- [17] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A Simple Way to Prevent Neural Networks from Overfitting," *Journal of Machine Learning Research*, Vol. 15, pp. 1929-1958, 2014.
- [18] A. Mathis, P. Mamidanna, K. N. 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 Neuroscience*, Vol. 21, No. 9, pp. 1281-1289, 2018.
- [19] E. P. S. Jones, "Accurate Canine Pose Recognition in Home Environments," in *Proc. International Conference on Animal-Computer Interaction*, pp. 12-21, 2019.
- [20] Z. Ge, S. Liu, F. Wang, Z. Li, and J. Sun, "YOLOX: Exceeding YOLO Series in 2021," *arXiv preprint arXiv:2107.08430*, 2021.
- [21] A. Zamansky, B. van der Linden, and I. T. H. C. J. M. van der Zande, "Towards a Framework for Automated Recognition of Animal Stress from Videos," in *Proceedings of the 5th International Conference on Animal-Computer Interaction*, pp. 1-6, 2018.
- [22] A. Droschel, D. Behnke, and S. Behnke, "Multi-modal animal pose estimation from RGB-D-thermal images," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 680-687, 2021.
- [23] A. Mathis, P. Mamidanna, K. M. Cury, T. Abe, V. N. Murthy, M. W. Mathis, and M. Bethge, "Markerless tracking of user-defined body parts with deep learning," *Nature Neuroscience*, Vol. 21, pp. 1281-1289, 2018.