

---

# Animal Posture Classification with Transfer Learning

---

Chen Qiao, Johann Verolet

## Abstract

1 This project proposes an Animal Posture Classification System leveraging transfer  
2 learning techniques. The core objective is to enhance animal behavior analysis,  
3 crucial for wildlife preservation, agricultural health, and pet care, by automating  
4 the classification of animal postures. Addressing the challenge of limited available  
5 datasets tailored for this purpose, our approach involves refining the ResNet50  
6 and Vision Transformer (ViT) model, pre-trained on ImageNet, with annotated  
7 datasets. This project also explores other solutions like Manifold Mixup and the  
8 utilization of synthetic data generated by stable diffusion model to overcome data  
9 scarcity. This multifaceted strategy aims to yield a more accurate and efficient  
10 system for animal posture classification. Experimental results demonstrate that  
11 transfer learning significantly outperforms models trained from scratch. Further  
12 improvements are achieved through the application of Manifold Mixup techniques  
13 and the integration of synthetic data.

14 **1 Introduction**

15 The classification of animal postures plays a pivotal role in understanding animal behavior, which  
16 is vital for wildlife conservation, agricultural health, and pet care. Traditional methods for this task  
17 often involve manual monitoring, which is not only time-consuming but also resource-intensive. Our  
18 project introduces a solution to this challenge: an Animal Posture Classification System utilizing  
19 state-of-the-art AI models and transfer learning techniques.

20 A major challenge in animal posture classification is the small size of available labeled datasets that  
21 can be used for training, since most available image datasets focus on animal category classification  
22 or keypoint detection. To bridge this gap, we utilize transfer learning strategies, employing models  
23 such as ResNet50[1] and ViT[2], both initially pre-trained on the ImageNet dataset. These models  
24 are subsequently fine-tuned using available annotated datasets. Our results demonstrate that transfer  
25 learning substantially enhances classification accuracy compared to training models from scratch,  
26 with improvements of 16.10% for ResNet50 and 22.53% for ViT, respectively.

27 To further enhance model performance, we incorporate two techniques that are applied concurrently  
28 with transfer learning: Manifold Mixup[3] and synthetic data generation via a diffusion model[4].  
29 Manifold Mixup, a data augmentation method, strategically improves model robustness by interpo-  
30 lating features and labels from different training samples. The diffusion model, renowned for its  
31 capacity to generate high-fidelity synthetic images, is deployed to augment our training dataset. This  
32 strategic addition of synthetic data effectively mitigates the issue of data imbalance and expands the  
33 diversity of our dataset. The application of Manifold Mixup increased the classification accuracy  
34 by 3.49% for the ResNet50 model and 1.85% for the ViT compared to using only transfer learning.  
35 Similarly, the integration of synthetic data improved accuracy by 1.06% for ResNet50 and 1.99%  
36 for ViT relative to results achieved with transfer learning alone. These improvements underscore  
37 the effectiveness of combining advanced data augmentation and synthetic data techniques to bolster  
38 machine learning model performance.

39 **2 Related Work**

40 Recent advances in Transfer Learning have shown significant potential across various domains by  
41 leveraging pre-trained models to boost performance on new, but related tasks. Zhuang et al. (2021)  
42 provide a comprehensive review[5], emphasizing that transfer learning minimizes reliance on large  
43 domain-specific datasets and enhances learning efficiency through knowledge transfer from related  
44 source domains. This aligns with our method, where we utilize transfer learning techniques with  
45 models such as ResNet50 and ViT, both pre-trained on ImageNet, to address the challenge of limited  
46 labeled data in animal posture classification.

47 The foundational principles of ResNet (He et al., 2015)[1] introduce a residual learning framework  
48 to mitigate the degradation problem in deep neural networks. Employing deeper architectures,  
49 ResNet demonstrates that increased network depth does not necessarily lead to higher error rates,  
50 provided that each layer set learns a residual mapping, a concept we leverage in our work. While the  
51 Vision Transformer (ViT) shifts the paradigm by applying self-attention mechanisms across image  
52 patches, treating image classification similarly to sequence processing in natural language tasks.  
53 This attention-based approach allows ViT to consider the global context of an image, leading to  
54 models that excel in understanding complex features without being confined by the sequential locality  
55 typical of convolutional neural networks (Dosovitskiy et al., 2020)[2]. In our study, we utilized both  
56 ResNet50 and ViT as the model architectures, enabling the utilization of advanced feature extraction  
57 techniques essential for accurate posture classification in diverse animal species.

58 The recent study by Verma et al. (2019)[3] on Manifold Mixup, and the application of the Stable  
59 Diffusion Model (Rombach et al., 2022)[5] for synthesizing data, presents significant advancements  
60 in the enhancement of deep learning models. Manifold Mixup notably advances model generalization  
61 by interpolating between hidden representations, yielding smoother decision boundaries and more  
62 conservative predictions. Additionally, the generation of synthetic data via the Stable Diffusion  
63 Model, provides a strategic advantage in addressing the challenges of imbalanced and limited training  
64 data. The incorporation of such synthetic data not only enlarges the dataset but also brings in a greater  
65 degree of variation.

66 **3 Problem definition**

67 The task involves developing a method for classifying the posture of animals using machine learning  
68 techniques. The primary challenge is the limited availability of annotated datasets, which are crucial  
69 for training and validating the classification models. The goal is to achieve higher accuracy through  
70 advanced techniques such as transfer learning, Manifold Mixup, and synthetic data generation. These  
71 methods aim to enhance the dataset's diversity and volume, thereby improving the robustness and  
72 performance of the classification models.

73 **4 Methodology**

74 **4.1 Data Preparation and Cleaning**

75 For this study, we selected the dataset entitled "Animals with Skeleton Key Points and Action  
76 Mark"[6] available on Kaggle as our primary data source. This dataset is composed of 967 images,  
77 featuring six different types of animals, each annotated with corresponding action classes which serve  
78 as posture labels. The posture categories defined within this dataset include lying, sitting, running,  
79 and standing/walking. To ensure the integrity and quality of our analysis, we conducted a rigorous  
80 data cleaning process. This process entailed a meticulous examination of the dataset, followed by  
81 the removal of images that were deemed unsuitable for our study due to quality issues or incorrect  
82 labeling. Following this cleaning process, the dataset was reduced to 819 images. The distribution of  
83 images across each posture category is detailed in Table 1.

Category	Number of Images
Lying	88
Sitting	67
Running	70
Standing/Walking	592
Total	819

Table 1: Distribution of Images After Data Cleaning

## 4.2 Model Architectures

In our Animal Posture Classification task, we utilize two distinct models: ResNet50 and the Vision Transformer (ViT), each selected for their unique architectural strengths.

### 4.2.1 ResNet50

The ResNet50 is chosen for its demonstrated effectiveness in training deeper neural networks[1], which is especially beneficial for our small dataset. Its innovative residual blocks help overcome the vanishing gradient problem, ensuring it performs well even with limited data. As illustrated in Figure 1, The input passes through a series of convolutional layers, each followed by a ReLU activation function. The feature map output from these layers is combined with the input signal, creating a "shortcut" path. Another ReLU activation function follows this skip connection. Figure 2 delineates the architecture of ResNet50. It is segmented into multiple stages, with Stage 0 dedicated to the initial processing of the image for low-level feature extraction. Subsequent stages, 1 through 4, comprise a series of residual blocks designed to incrementally extract higher-level features. This hierarchically structured approach allows for a comprehensive analysis of the visual data, capturing details essential for distinguishing a variety of animal postures. By strategically deepening the network through these stages, the model pays better attention to important features without using too much computing power.

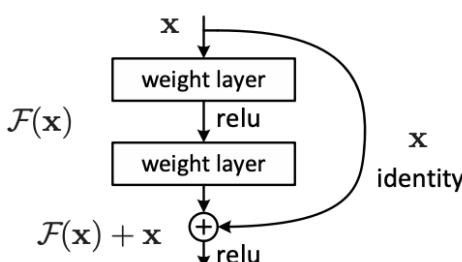


Figure 1: Residual Block

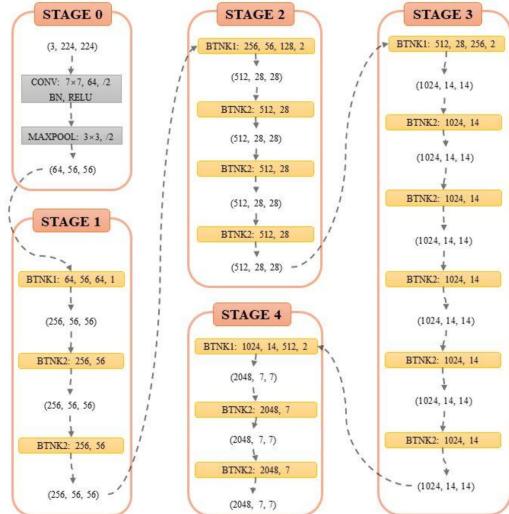


Figure 2: Network Architecture of ResNet50

### 4.2.2 Vision Transformer

Vision Transformer (ViT) is employed for its pioneering approach to image analysis, which abandons the local processing paradigm of CNNs in favor of global information aggregation using self-attention mechanisms. As detailed in Figure 3, it harnesses a Transformer Encoder that refines image analysis. This Encoder, through its multi-headed self-attention mechanism, attentively assesses

106 relationships across image patches to grasp complex visual information. Integral to its design  
 107 are residual connections and normalization layers, which preserve gradient flow and feature scale  
 108 across deep networks. This facilitates stable learning and nuanced feature discernment, essential  
 109 for recognizing diverse animal postures with precision. Supplementing this, Figure 4 illustrates the  
 110 initial stage of ViT's processing pipeline, where an image is divided into patches. These patches  
 111 are then linearly projected and combined with positional embeddings to form a sequence of vectors,  
 112 representing both the visual content and positional information. This encoded sequence is what the  
 113 Transformer Encoder processes. The global contextual understanding gleaned from this methodology  
 114 is critical, allowing the MLP head at the output to classify animal postures.

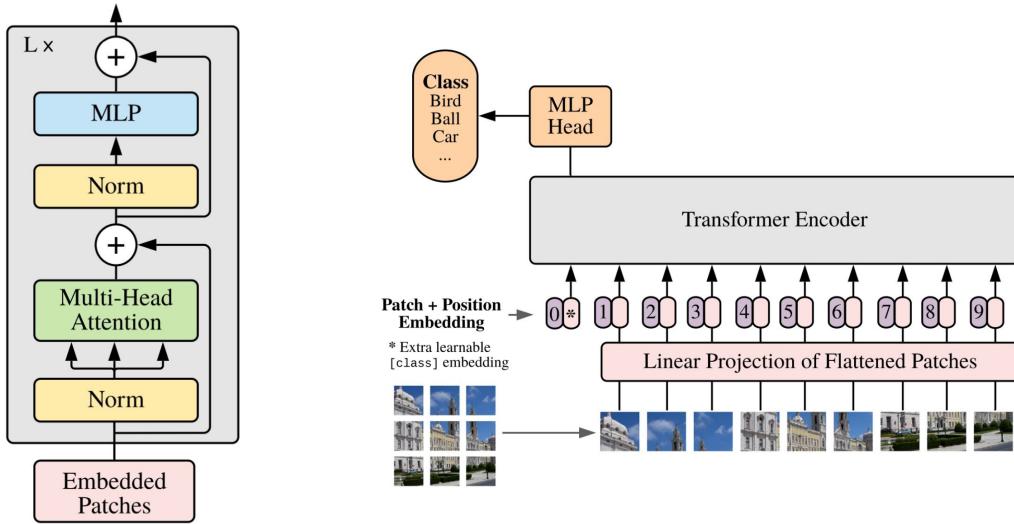


Figure 4: ViT Architecture Overview

Figure 3: Transformer Encoder

### 115 4.3 Transfer Learning

116 We integrate the two models introduced earlier: ResNet50 and the ViT, both pretrained on the  
 117 ImageNet dataset. This pretraining endows the models with a broad understanding of visual features  
 118 applicable across various domains. In our experiments, we conduct two types of transfer learning  
 119 approaches. The first involves freezing the pretrained layers, rendering them non-trainable to preserve  
 120 the learned features, and allowing the models to function as fixed feature extractors. The second  
 121 approach allows for fine-tuning, where the pretrained layers are set to be trainable, enabling the  
 122 model to adapt these features to the nuances of our specific classification task.

123 To tailor the models to our classification needs, we add a custom classification head comprising  
 124 a linear layer. The output layer is designed to classify four distinct classes. This additional layer  
 125 effectively transforms the general feature representations learned in previous layers into predictions  
 126 specific to our posture classification task.

127 By employing transfer learning, we capitalize on the knowledge learned from vast amounts of visual  
 128 data and distill it for our specific task, significantly reducing the need for extensive animal posture  
 129 datasets and computational resources, while providing higher accuracy in our classification results.

### 130 4.4 Data Augmentation through Manifold Mixup

131 We incorporate the technique of Manifold Mixup as an advanced form of data augmentation because  
 132 of its ability to smooth the decision boundaries and provide a robust regularization effect, which is  
 133 achieved by creating border regions in-between different classes within the data representation space.  
 134 This process also tends to flatten the learned representations, i.e., the network learns to represent the  
 135 input data in fewer directions of variance[3], leading to a more rapid decline in the singular values.  
 136 As a result, the network tends to have better generalization ability.

Our implementation of Manifold Mixup follows the technique described in the study by Verma et al. (2019)[3]. During training, we first select a layer K within the network architecture. At this layer, two random mini-batches, denoted as  $(x, y)$  and  $(x', y')$  resulting in the transformed features  $(g(x), y)$  and  $(g(x'), y')$  are mixed up. To perform the Mixup, a random interpolation weight  $\lambda$  is generated uniformly from the interval  $[0, 1]$ . This weight and its complement  $1 - \lambda$  are then used to mix both the latent representations and labels of the two mini-batches. Specifically, the mixed latent representation is computed as

$$Mix(g(x), g(x')) = \lambda \cdot g(x) + (1 - \lambda) \cdot g(x'),$$

and the mixed labels are computed as

$$Mix(y, y') = \lambda \cdot y + (1 - \lambda) \cdot y'$$

137 The resultant mixed representations and labels are subsequently used for further training, thereby  
 138 encouraging the network to predict not only on the observed data but also on the interpolated data,  
 139 enhancing the model's generalization capabilities.

140 By applying Manifold Mixup, we induce our network to predict consistently not only at the data  
 141 points but also in their vicinity, hence encouraging the learning of more robust features.

#### 142 **4.5 Integrating Synthetic Data from Stable Diffusion Model**

143 To address the challenge of limited and imbalanced dataset, we incorporate synthetic images generated  
 144 by the pre-trained Stable Diffusion model (version 1.4). Stable Diffusion model is trained to generate  
 145 images by learning the underlying distribution of vast and diverse datasets. By sampling from this  
 146 learned distribution, we generate additional images that exhibit varied animal postures, enriching the  
 147 representational diversity of our training data.

148 The integration process involves selecting specific prompts that describe a range of animal postures  
 149 and feeding these into the Stable Diffusion system. We refine our prompts through testing, discovering  
 150 that clear and specific prompts, such as "a realistic image of the full body of a white cat lying on a  
 151 sofa," yield good results. The model generates new images from these prompts, effectively creating a  
 152 multitude of images that help our system learn more generalized features. These synthesized images  
 153 are subsequently reviewed and cleaned before being introduced into our training set. During this  
 154 process, we removed unrealistic images, such as a horse with five legs or a cow in an improbable  
 155 posture. Approximately 30% of the initial images remain after this process. Figure 5 presents a  
 156 selection of synthetic images. Figure 6 presents pie charts that illustrate the makeup of our dataset  
 157 before and after incorporating synthetic data. This figure demonstrates how the addition of synthetic  
 158 data aids in addressing the issue of data imbalance.

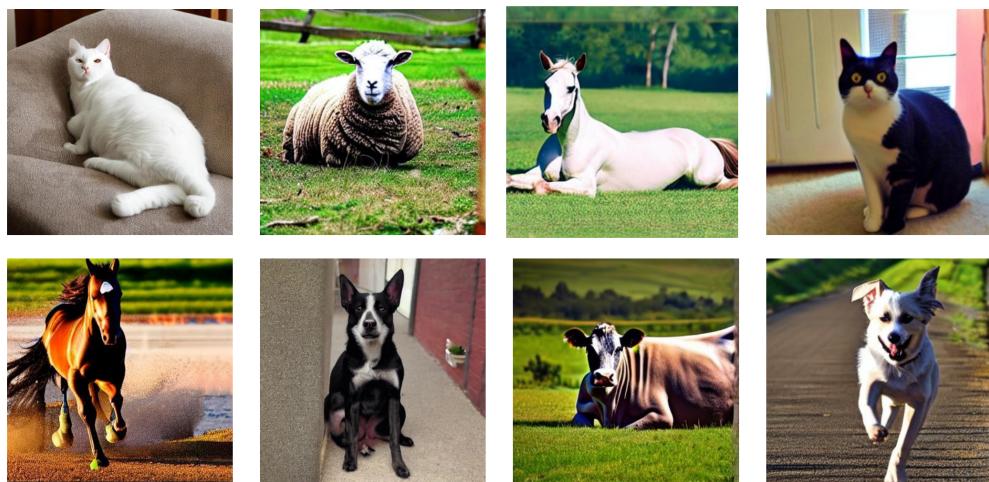


Figure 5: Synthetic Data Generated by Stable Diffusion

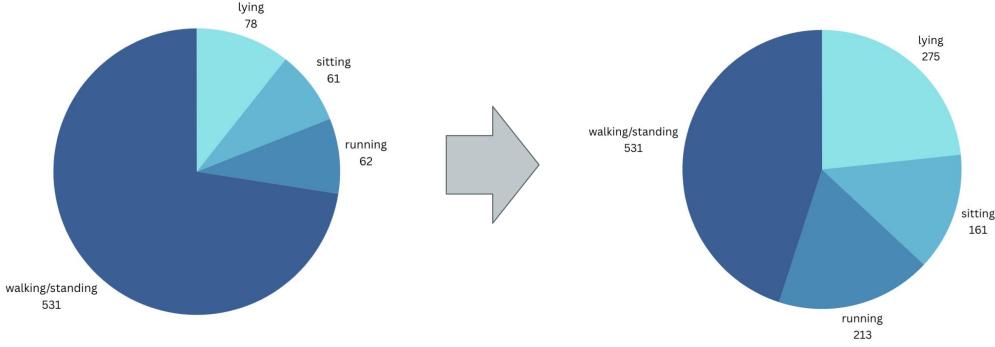


Figure 6: Training Set Before and After Integrating Synthetic Data

## 159 5 Experimental Results

### 160 5.1 Baseline Evaluation

161 To establish a baseline, we assessed the performance of models trained from scratch on our cleaned  
 162 animal posture dataset without the use of transfer learning. We experimented with training the  
 163 ResNet50 and ViT models, both with and without the implementation of data augmentation techniques.  
 164 The data augmentation protocol included rotations up to 10 degrees, horizontal and vertical shifts  
 165 of 10% of the image size, shear transformations that altered image shape by up to 10%, zoom  
 166 adjustments of 10%, and horizontal flips. Our findings revealed that data augmentation enhanced  
 167 the performance of the ResNet50 model, while the ViT model achieved higher accuracy without the  
 168 use of data augmentation. Therefore, we reported the results with data augmentation for ResNet50  
 169 and results without data augmentation for ViT. These results are presented in Table 2. The reported  
 170 accuracy, precision, recall, and F1 score represent the average results across three experimental runs.

	ViT	ResNet50
Lying Acc (%)	10.00	0.00
Sitting Acc (%)	0.00	0.00
Running Acc (%)	0.00	0.00
Stand/Walking Acc (%)	95.08	100.00
Overall Acc (%)	67.82	70.11
Precision (%)	22.16	17.53
Recall (%)	26.27	25.00
F1 Score (%)	24.04	20.61

Table 2: Performance of ResNet50 and ViT Models

### 171 5.2 Transfer Learning Impact

172 We applied transfer learning to the ResNet50 and ViT models using pretrained ImageNet weights,  
 173 adapting them to our animal posture dataset. To tailor the models for our task, a classification head  
 174 was added. We explored two methods: freezing all layers and fine-tuning the entire model. ViT  
 175 showed improved performance with full fine-tuning, while ResNet50 excelled with frozen layers,  
 176 Consequently, in Table 3, we present the results for ViT(TL) with fully fine-tuned layers and for  
 177 ResNet50(TL) with frozen layers. ViT was tested without data augmentation, whereas ResNet50  
 178 incorporated it, enhancing accuracy. TL stands for Transfer Learning.

179 Table 3 details the comparison between models trained from scratch and those using transfer learning.  
 180 Transfer learning significantly improved the models' ability to recognize diverse animal postures  
 181 beyond just walking or standing. This analysis underscores the effectiveness of transfer learning in  
 182 refining model performance on our task.

	ViT	ViT (TL)	ResNet50	ResNet50 (TL)
Lying Acc (%)	10.00	90.00	0.00	60.00
Sitting Acc (%)	0.00	87.50	0.00	87.50
Running Acc (%)	0.00	50.00	0.00	62.50
Stand/Walking Acc (%)	95.08	98.36	100.00	93.44
Overall Acc (%)	67.82	91.95	70.11	86.21
Precision (%)	22.16	90.38	17.53	77.53
Recall (%)	26.27	81.47	25.00	75.86
F1 Score (%)	24.04	85.69	20.61	75.38

Table 3: Impact of Transfer Learning, TL stands for Transfer Learning

### 183 5.3 Manifold Mixup Impact

184 Following the transfer learning phase, we explored the application of Manifold Mixup as a data  
 185 augmentation technique to further enhance model performance. This approach was evaluated by  
 186 assessing its impact on accuracy, precision, recall, and F1 Score, and by comparing these metrics with  
 187 the results obtained from the transfer learning implementation. The ResNet50 model included basic  
 188 data augmentation, as detailed in Section 4.1, which contributed to improved accuracy. In contrast,  
 189 the ViT model was tested without such augmentations because they did not yield any improvement.  
 190 Table 4 presents a detailed comparison between the models employing only transfer learning and  
 191 those utilizing both transfer learning and Manifold Mixup. Although Manifold Mixup does not  
 192 significantly address the issue of dataset imbalance—since it synthesizes new data points by mixing  
 193 up existing ones rather than creating completely new examples—it still managed to enhance overall  
 194 accuracy, precision, recall, and F1 Score by a decent margin. This demonstrates that while Manifold  
 195 Mixup may not solve all challenges related to dataset characteristics, it effectively contributes to  
 196 improving model robustness and performance in our animal posture classification system.

	ViT (TL)	ViT (TL + MM)	ResNet50 (TL)	ResNet50 (TL + MM)
Lying Acc (%)	90.00	100.00	60.00	88.89
Sitting Acc (%)	87.50	87.50	87.50	64.29
Running Acc (%)	50.00	50.00	62.50	57.14
Stand/Walking Acc (%)	98.36	98.36	93.44	96.64
Overall Acc (%)	91.95	93.10	86.21	89.70
Precision (%)	90.38	94.27	77.53	82.80
Recall (%)	81.47	83.97	75.86	76.74
F1 Score (%)	85.69	88.82	75.38	78.89

Table 4: Impact of Manifold Mixup, TF stands for Transfer Learning, MM stands for Manifold Mixup

### 197 5.4 Synthetic Data Integration and Comparative Analysis

198 In our pursuit to address the challenges posed by dataset imbalance, we experimented with the  
 199 integration of synthetic images generated by stable diffusion model into our training dataset. This  
 200 approach aimed to enrich the dataset with diverse examples and assess the potential benefits of  
 201 synthetic data on model performance. We conducted comparative analyses between models trained  
 202 on the original dataset and those trained on a combined set that included both original real-world data  
 203 and synthetic images.

204 Table 5 presents the outcomes of this experiment, comparing model performance following transfer  
 205 learning against performance after integrating both synthetic data and transfer learning. A detailed  
 206 examination of the accuracies across different categories revealed benefits. Specifically, the integration  
 207 of synthetic data effectively mitigated the issue of dataset imbalance by introducing a wider variety  
 208 of data representations. This result underscores the utility of synthetic data in enhancing model  
 209 robustness and addressing specific challenges such as class imbalance in training datasets.

	ViT (TL)	ViT (TL + SD)	ResNet50 (TL)	ResNet50 (TL + SD)
Lying Acc (%)	90.00	100.00	60.00	88.89
Sitting Acc (%)	87.50	87.50	87.50	71.43
Running Acc (%)	50.00	87.50	62.50	64.29
Stand/Walking Acc (%)	98.36	96.72	93.44	91.60
Overall Acc (%)	91.95	95.40	86.21	87.27
Precision (%)	90.38	92.29	77.53	79.05
Recall (%)	81.47	92.93	75.86	75.27
F1 Score (%)	85.69	92.61	75.38	76.81

Table 5: Impact of Synthetic Data, TL stands for Transfer learning, SD stands for Synthetic Data

## 210 6 Limitation and Future work

211 While we employed Stable Diffusion version 1.4 to generate synthetic data due to its wide adoption  
 212 and robust community support, this choice brings inherent limitations. This version, though popular,  
 213 may not encompass the latest advancements in generative model capabilities, which could limit the  
 214 diversity and realism of the synthetic data generated.

215 Future iterations of our project will explore the use of newer versions of the model, such as Stable  
 216 Diffusion version 2.1, which promises enhanced data generation with potentially greater variability  
 217 and higher fidelity. By integrating these advancements, we aim to further address the challenges of  
 218 data imbalance and improve the robustness and accuracy of our animal posture classification models.

## 219 7 Conclusion

220 This study demonstrates the effectiveness of transfer learning in animal posture classification, enhancing  
 221 model accuracy using advanced architectures like ResNet50 and Vision Transformer, pre-trained  
 222 on ImageNet. Techniques such as Manifold Mixup and synthetic data generation with Stable Diffusion  
 223 have effectively mitigated the problem of limited and imbalanced datasets, improving model  
 224 generalization. Our findings support the potential of integrating cutting-edge AI technologies to  
 225 advance animal behavior analysis, crucial for applications in wildlife conservation and pet care.  
 226 Future work will explore newer generative models and expand the dataset diversity, further refining  
 227 the robustness and applicability of our classification system.

## 228 References

- 229 [1] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *IEEE  
 230 Conference on Computer Vision and Pattern Recognition*, 2015.
- 231 [2] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani,  
 232 M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. & Houlsby. An image is worth 16x16  
 233 words: Transformers for image recognition at scale. In *International Conference on Learning  
 234 Representations*, 2021.
- 235 [3] V. Verma, A. Lamb, C. Beckham, A. Najafi, I. Mitliagkas, A. Courville, D. Lopez-Paz, and  
 236 Y. Bengio. Manifold mixup: Better representations by interpolating hidden states. In *International  
 237 Conference on Machine Learning*, 2018.
- 238 [4] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis  
 239 with latent diffusion models. In *Conference on Computer Vision and Pattern Recognition*, 2021.
- 240 [5] F. Zhuang, Z. Qi, K. Duan, D. Xi, Y. Zhu, H. Zhu, H. Xiong, and Q. He. A comprehensive survey  
 241 on transfer learning. In *Institute of Electrical and Electronics Engineers*, 2021.
- 242 [6] <https://www.kaggle.com/datasets/egorovalexeyd/animals-with-skelet-key-points-and-action-mark>.