

## Problem

Human share the world with and learn the behavior from with billions of animals. **However, low cooperativity and high species diversity prevent us from building comprehensive large-scale datasets to know our neighbors.** Further, the small data training is like cultivating on a tundra.

- Small data training makes the model **lack of robustness** and thus difficult to cope with more free-ranging movements, occlusions and environments.
- The synthetic animal data used as a supplement **lacks realism of poses** and it is difficult to be **blended into the real background.**

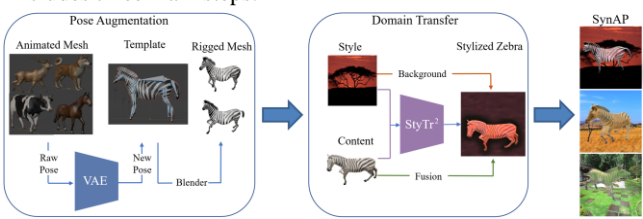
## Approach

We present a cost-effective and generic prior-aware synthetic data generation pipeline, called **PASyn**, for animals pose estimate tasks that suffer from severe data scarcity.

- A novel variational autoencoder (VAE)-based synthetic animal data generation pipeline PASyn to generate probabilistically-valid pose data
- A style transfer strategy to militate the inconsistency between synthetic animal and real background
- A **synthetic animal pose (SynAP)** dataset, containing 3,000 zebra images and 3,600 images of six common quadrupeds

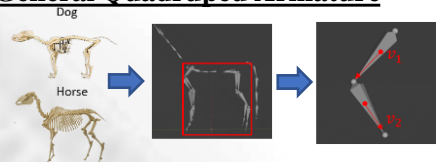
## Prior-Aware Synthetic Data Generation

Our Prior-Aware Synthetic Data Generation (PASyn) Pipeline includes three main steps:



## Pose Augmentation

### General Quadruped Armature



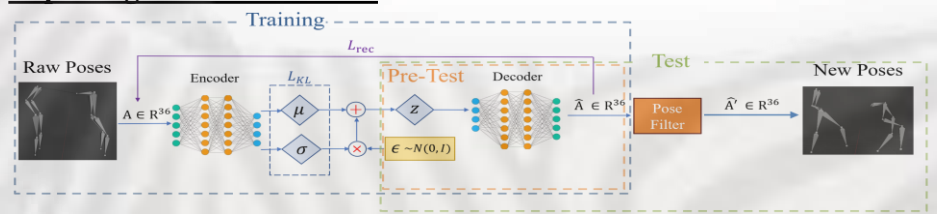
### Training

$$\mathcal{L}_{total} = w_1 \mathcal{L}_{KL} + w_2 \mathcal{L}_{rec},$$

$$\mathcal{L}_{KL} = KL(q(z|A) || \mathcal{N}(0, I)),$$

$$\mathcal{L}_{rec} = \|A - \hat{A}\|_2^2$$

### Capturing Animal Pose Prior

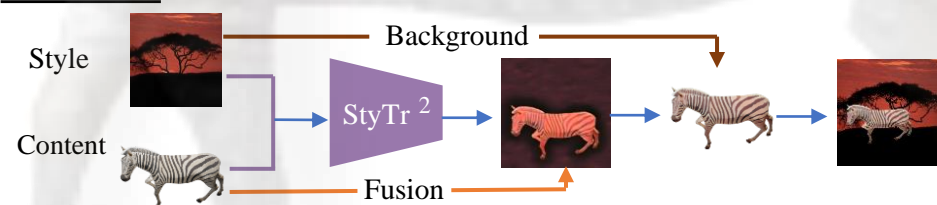


## Stylization

### Why style transfer

- Enrich the texture diversity in a reasonable way
- Blend the synthetic animal image into the real backgrounds

### Architecture



## Synthetic Animal Pose (SynAP) Dataset

### Size

SynAP contains **3,000** synthetic zebra images and SynAP+ extends the SynAP with **3,600** images of horses, cows, sheep, dogs, giraffes and deer.

### Background

**300** grass, savanna, and forest real scenes are collected from Internet to stylize the synthetic animal.

## Results

### Evaluation Over SynAP

Method	Backbone	Training Set	PCK@0.05 Pose Estimation Accuracy on Zebra-300 Set										
			Eye	Nose	Neck	Shoulders	Elbows	F-Paws	Hips	Knees	B-Paws	RoT	Average
MMPose [8]	HRNet-w32	R(99)	97.3	95.8	83.2	78.8	77.1	62.6	86.0	74.9	59.8	82.4	78.7
		R(99)+S(3K)	<b>97.8</b>	<b>98.3</b>	81.1	<b>94.0</b>	<b>93.5</b>	<b>92.0</b>	93.7	93.5	89.0	<b>87.6</b>	<b>92.4</b>
		R(99)+S(5K)	97.5	96.9	81.8	89.6	91.3	90.7	<b>94.1</b>	<b>94.1</b>	<b>90.4</b>	<b>86.0</b>	91.6
DeepLabCut [20]	EfficientNet-B6	R(99)	93.7	96.2	<b>82.5</b>	<b>91.4</b>	80.8	67.4	88.1	84.5	71.8	83.2	83.6
		R(99)+S(3K)	<b>95.1</b>	<b>97.9</b>	81.5	90.1	83.3	75.5	<b>93.2</b>	89.3	83.9	86.8	87.6
		R(99)+S(5K)	94.1	92.6	80.8	90.8	<b>87.0</b>	<b>85.7</b>	90.5	<b>93.3</b>	<b>88.3</b>	<b>86.8</b>	<b>89.2</b>
MMPose [8]	ResNet-50	R(99)	96.2	96.9	<b>80.8</b>	59.0	71.3	71.2	88.5	78.2	59.3	85.2	76.9
		R(99)+S(3K)	95.6	95.8	69.9	<b>87.3</b>	<b>84.6</b>	84.3	90.8	91.2	84.4	77.2	86.7
		R(99)+S(5K)	<b>97.0</b>	<b>97.9</b>	74.5	85.3	84.2	<b>84.5</b>	<b>94.6</b>	<b>91.4</b>	<b>88.0</b>	<b>85.6</b>	<b>88.4</b>

### The effect of SynAP with large real data

Training Set	Real Zebra	PCK@0.05 Pose Estimation Accuracy on Zebra-300 Set										
		Eye	Nose	Neck	Shoulders	Elbows	F-Paws	Hips	Knees	B-Paws	RoT	Average
R(8K) (SOTA)		97.5	97.2	79.4	87.8	90.3	93.8	95.3	94.1	89.5	86.4	91.4
R(8K) + S(3K)	✓	97.3	<b>98.3</b>	79.0	93.1	94.9	96.0	95.3	<b>96.7</b>	<b>93.3</b>	<b>89.6</b>	<b>93.8</b>
R(8K) + S(5K)		97.3	97.6	<b>81.1</b>	<b>93.7</b>	<b>95.7</b>	<b>96.0</b>	<b>96.6</b>	96.0	<b>94.3</b>	87.6	<b>94.2</b>
R(8K)		79.7	87.7	37.4	77.6	80.0	87.6	82.0	86.4	81.3	67.2	78.3
R(8K) + S(3K)	✗	94.8	96.2	<b>67.1</b>	90.8	87.9	90.2	87.6	91.6	89.7	77.6	88.2
R(8K) + S(5K)		<b>97.5</b>	<b>96.2</b>	66.1	<b>91.6</b>	<b>89.5</b>	<b>93.8</b>	<b>93.9</b>	<b>93.5</b>	<b>91.1</b>	<b>77.6</b>	<b>90.2</b>

### Ablation Study

Index	Training Set	VAE	Style Transfer	$\sigma^2$	Zoo Zebra Set	Zebra-300 Set
a	R(99)	✗	✗	✗	76.0	78.7
b	S(3K)	✗	✗	2I	38.7	30.0
c	S(3K)	✓	✗	2I	44.2	36.7
d	S(3K)	✓	✓	2I	42.9	46.6
e	R(99)+S(3K)	✗	✗	2I	89.8	88.0
f	R(99)+S(3K)	✓	✗	2I	90.4	89.8
g	R(99)+S(3K)	✓	✓	I	90.5	91.1
h	R(99)+S(3K)	✓	✓	2I	<b>91.5</b>	<b>92.4</b>

### Visualized Results



**In conclusion**, our synthetic animal pose dataset SynAP and its extended version SynAP+, and th of real data is verified on different backbones and achieves state-of-the positive effect of them on pose estimation task of animals with a small amount-art performance.