



# Master Thesis

## Incorporating Temporal Smoothing Into 2D Pose Estimation

André Oskar Andersen (wpr684)  
wpr684@alumni.ku.dk

2023

**Supervisor**  
Kim Steenstrup Pedersen kimstp@di.ku.dk

## **Abstract**

## Preface

## **Acknowledgement**

# Contents

<b>1</b>	<b>Introduction</b>	<b>7</b>
1.1	Related Work . . . . .	7
1.2	Problem Definition . . . . .	7
1.3	Reading Guide . . . . .	7
<b>2</b>	<b>Deep Learning Theory</b>	<b>8</b>
2.1	Feedforward Neural Networks . . . . .	8
2.1.1	Fully-connected Layers . . . . .	8
2.1.2	Convolutional Layer . . . . .	8
2.2	Recurrent Neural Networks . . . . .	9
2.2.1	Convolutional Long Short-Term Memory . . . . .	10
2.3	Transformer . . . . .	10
2.3.1	Attention . . . . .	11
2.3.2	The Architecture . . . . .	12
2.4	General Machine Learning Terminology . . . . .	13
<b>3</b>	<b>Models</b>	<b>16</b>
3.1	Baseline . . . . .	16
3.2	Bidirectional Convolutional LSTM . . . . .	16
3.2.1	Summing the Sequence Directions . . . . .	17
3.2.2	Concatenating the Sequence Directions . . . . .	17
3.3	DeciWatch . . . . .	18
<b>4</b>	<b>Dataset</b>	<b>20</b>
4.1	The ClimbAlong Dataset . . . . .	20
4.2	The BRACE Dataset . . . . .	23
4.3	The Penn Action Dataset . . . . .	25
<b>5</b>	<b>Pretraining</b>	<b>27</b>
5.1	Data Preprocessing . . . . .	27
5.2	Training Details . . . . .	28
5.2.1	Data Configuration . . . . .	28
5.2.2	Setups . . . . .	28
5.2.3	Training Configuration . . . . .	29
5.3	Training and Validation Results . . . . .	29
5.4	Test Results . . . . .	31
<b>6</b>	<b>Finetuning</b>	<b>35</b>
6.1	Data Preprocessing . . . . .	35
6.2	Training Details . . . . .	35
6.3	Results . . . . .	36
<b>7</b>	<b>Discussion</b>	<b>37</b>
<b>8</b>	<b>Conclusion</b>	<b>38</b>
<b>9</b>	<b>References</b>	<b>39</b>

## Notation

$\boldsymbol{x}$	A vector
$\boldsymbol{X}$	A matrix
$\mathbf{X}$	A Tensor
$\boldsymbol{X}_{ij}$	Element located at row $i$ column $j$ in matrix $\boldsymbol{X}$
*	The convolution operator
o	The Hadamard product
$\sigma$	The sigmoid activation function
$\nabla_{\boldsymbol{x}} f$	Gradient of $f$ with respect to $\boldsymbol{x}$
$f(\boldsymbol{x}; \theta)$	Function $f$ with input $\boldsymbol{x}$ and parameter $\theta$
$\mathcal{U}(a, b)$	Uniform distribution with support $[a, b]$

# 1 Introduction

## 1.1 Related Work

2-dimensional pose estimation can be divided into either being image-based or video-based, where the methods in the latter case use the tempoeral information of the video to perform the pose estimation.

Image-based methods were initially based on the geometry between the joints of the taget image [17, 20, 25]. Following this, were the convolutional-based methods, that used convolutional neural networks [12] to perform the pose estimation [22, 15, 4, 8]. More recent methods use transformers [21] to deliver state-of-the-art results [23, 24].

Early video-based methods used 3-dimensional convolutions to capture the temporal information between neighboring frames [16, 5]. Other methods use LSTM's [9] to capture this temporal information [13, 2]. Like in the case of image-based methods, transformers [21] have recently been introduced to the video-based methods to capture the temporal information and deliver state-of-the-art-results [26].

## 1.2 Problem Definition

## 1.3 Reading Guide

## 2 Deep Learning Theory

The following section covers the most important background theory for the experiments in Section 6.3. This includes an introduction to various types of neural networks, as well as an introduction to the optimization of such networks.

### 2.1 Feedforward Neural Networks

**Feedforward neural networks** are the most basic type of neural networks. The aim of a feed-forward neural network is to approximate some function  $f^*$ , by defining a mapping  $\mathbf{y} = f(\mathbf{x}; \boldsymbol{\theta})$  and learning the parameters  $\boldsymbol{\theta}$ , that results in the best approximation of  $f^*$ . These models are called **feedforward** because there are no **feedback** connections in which the outputs of the model are fed back into itself. Instead, information flows through the function being evaluated from  $\mathbf{x}$ , through the intermediate computations used to define  $f$ , and finally to the output  $\mathbf{y}$ . Feedforward neural networks generally consists of multiple **layers**, arranged in a chain structure, with each layer being a function of the layer that preceded it [7].

#### 2.1.1 Fully-connected Layers

The most simple type of layer found in a feedforward neural network is the **fully-connected layer**. The fully-connected layer usually consists of some learnable parameter matrix and learnable parameter vector, as well as a **activation function**, which is a non-linear function, that makes it possible for the network to model non-linearity. Three of the most common activation functions are the **Rectified Linear Unit (ReLU)**, **Sigmoid** and **Hyperbolic Tangent (Tanh)** activation functions, defined as

$$\text{ReLU}(x) = \max\{0, x\} \quad (1)$$

$$\text{Sigmoid}(x) = \sigma(x) = \frac{1}{1 + \exp(x)} \quad (2)$$

$$\tanh(x) = \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)} \quad (3)$$

Given the learnable parameter matrix  $\mathbf{W}$ , learnable parameter vector  $\mathbf{b}$  and activation function  $g$ , the  $i$ 'th fully-connected layer is defined as

$$\mathbf{h}^{(i)} = \begin{cases} g^{(i)} (\mathbf{W}^{(i)\top} \mathbf{h}^{(i-1)} + \mathbf{b}^{(i)}) & \text{if } i > 1 \\ g^{(1)} (\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)}) & \text{if } i = 1 \end{cases}. \quad (4)$$

Thus, for a neural network with  $n$  layers, we have the mapping [7]

$$\mathbf{y} = f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{h}^{(n)}. \quad (5)$$

#### 2.1.2 Convolutional Layer

A **convolutional layer** is a specialized kind of feedforward layer, usually used in analysis of time-series or image data. If a network has at least one convolutional layer, it is called a **Convolutional neural network (CNN)** [7].

The convolutional layer consists of a set of **kernels**, each to be applied to the entire input vector, where each kernel is a learnable parameter matrix  $k \times k$  [18]. Each kernel is applied on the input to produce a **feature map**. The kernels are applied to the input by "sliding" over the input (where the step size is called **stride**). Each  $k \times k$  grid of the input is then used to

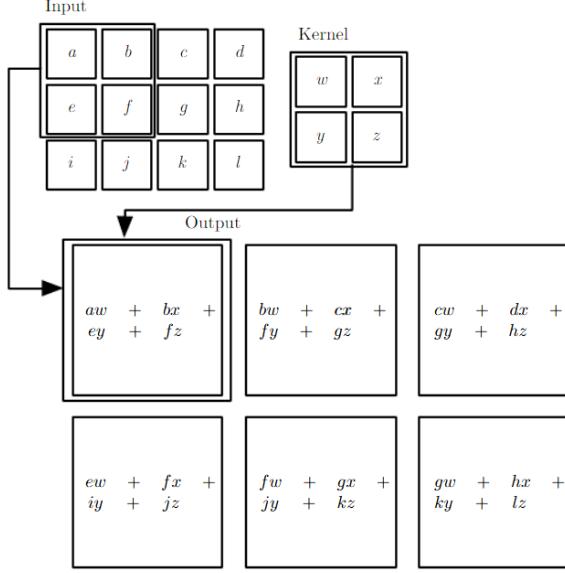


Figure 1: An example of applying a 2d kernel on an input [7].

compute the dot-product between the grid and each kernel, which is then placed in the corresponding feature map of each kernel, as visualized in Figure 1. [1]. To control the dimensions of the output, one might **pad** the sides with a constant value. Commonly, zero is used as the padding-value.

As seen in Figure 1, each kernel produces a linear combination of all pixel values in a neighbourhood defined by the size of the kernel. Thus, unlike a fully-connected layer, a convolutional layer captures the high correlation between a pixel and its neighbours. Further, by limiting the size of the kernel, the network will use much fewer parameters, than if it was replaced by a fully-connected layer [7].

## 2.2 Recurrent Neural Networks

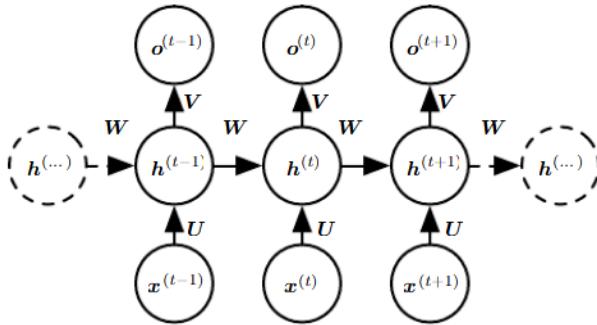


Figure 2: An illustration of an RNN [7].

**Recurrent neural networks (RNNs)** are a family of neural networks for processing sequential data. Figure 2 illustrates the general setup of such a network, which maps an input sequence of  $x$  values to a corresponding sequence of output  $o$  values. Generally, a RNN consists of three parts: (1) the input ( $x^{(i)}$ ), (2) the hidden state ( $h^{(i)}$ ), and (3) the output ( $o^{(i)}$ ). During inference, the model maps each input value to an output value in a sequential matter, where it first maps the first input value, then the second, then the third, and so forth. The network maps each input value to an output value by making use of the hidden state from the preceding step,

where the first hidden state has to be initialized [7].

### 2.2.1 Convolutional Long Short-Term Memory

One common recurrent neural network unit is the **convolutional long short-term memory (ConvLSTM)**, which is an adaptation of the standard **long short-term memory (LSTM)** for sequences of images. Both LSTMs work by potentially stacking multiple **cells** together, such that each LSTM cell is taking in outputs of the preceding LSTM cell as its input.

The idea of both LSTM cells is to create paths through time that have derivatives that neither vanish nor explode. This is done by introducing a memory cell  $C_t$ , which accumulates state information over a long duration. This cell is accessed, written and cleared by several controlling gates. The model learns during training, when to access, write and clear this memory cell. By using the memory cell and gates to control the flow of information, the gradient will be trapped in the cell and thus be prevented from vanishing too quickly [7, 19].

The ConvLSTM cell consists of three gates. The first gate is the input gate  $I_t$ , which controls whether or not to accumulate the information of a new input to the cell. This gate is defined as

$$I_t = \sigma(W_{xi} * X_t + W_{hi} * H_{t-1} + W_{ci} \circ C_{t-1} + B_i) \quad (6)$$

where  $X_t$  is the current input image,  $H_{t-1}$  is the hidden state of the previous time step,  $C_{t-1}$  is the cell output of the previous time step, and  $W_{xi}$ ,  $W_{hi}$ ,  $W_{ci}$  and  $B_i$  are learnable parameters

The second gate is the **forget gate** unit  $F_t$  (at time step  $t$ ), which controls whether or not to "forget" the status of the cell output of the previous time step. The forget gate is defined as

$$F_t = \sigma(W_{xf} * X_t + W_{hf} * H_{t-1} + W_{cf} \circ C_{t-1} + B_f) \quad (7)$$

where  $W_{xf}$ ,  $W_{hf}$ ,  $W_{cf}$  and  $B_f$  are learnable parameters.

The last gate is the **output gate** unit  $O_t$ , which controls whether or not the latest cell output will be propagated to the final state  $H_t$ . This cell is defined as

$$O_t = \sigma(W_{xo} * X_t + W_{ho} * H_{t-1} + W_{co} \circ C_t + B_o) \quad (8)$$

where  $W_{xo}$ ,  $W_{ho}$ ,  $W_{co}$  and  $B_o$  are learnable parameters.

By combining the three gates, we get the following definition of the ConvLSTM Cell. Let  $X_1, X_2, \dots, X_t$  be a sequence of input images. Then, at each time step  $t$ , we compute

$$I_t = \sigma(W_{xi} * X_t + W_{hi} * H_{t-1} + W_{ci} \circ C_{t-1} + B_i) \quad (9)$$

$$F_t = \sigma(W_{xf} * X_t + W_{hf} * H_{t-1} + W_{cf} \circ C_{t-1} + B_f) \quad (10)$$

$$C_t = F_t \circ C_{t-1} + I_t \circ \tanh(W_{xc} * X_t + W_{hc} * H_{t-1} + B_c) \quad (11)$$

$$O_t = \sigma(W_{xo} * X_t + W_{ho} * H_{t-1} + W_{co} \circ C_t + B_o) \quad (12)$$

$$H_t = O_t \circ \tanh(C_t) \quad (13)$$

For  $t = 1$ , both  $H_{t-1}$  and  $C_{t-1}$  have to be initialized [19].

### 2.3 Transformer

The sequential nature of RNNs precludes parallelization with training examples, heavily slowing down the training of these models. Another type of model for sequential data is the

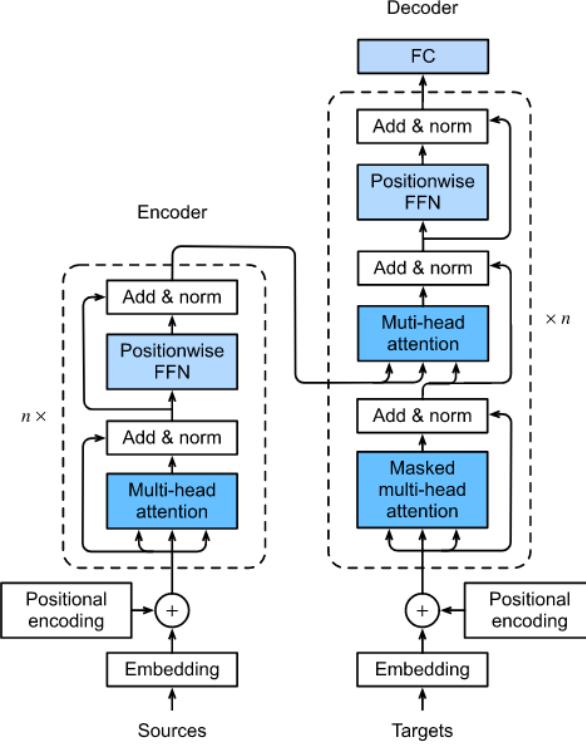


Figure 3: An illustration of the Transformer architecture [27].

**Transformer**, which eschews recurrence and instead relies on an **attention** mechanism to draw dependencies between input and output. The Transformer allows for more parallelization and can reach state of the art results [21].

### 2.3.1 Attention

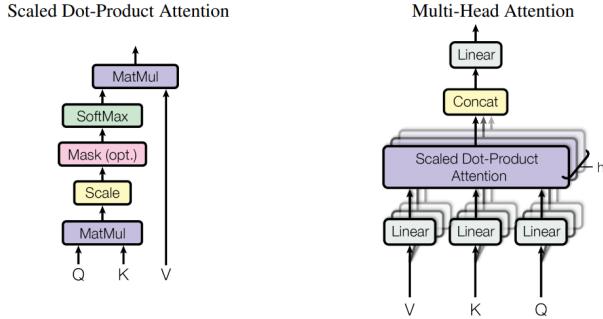


Figure 4: Illustration of (left) the Scaled Dot-Product Attention, and (right) the Multi-Head Attention, which consists of several attention layers running in parallel [21].

An attention function is a function that maps a given query and a given set of key-value pairs to an output. In this function the query, keys, values, and output are all vectors, and the output is a weighted sum of the values, where the weight of each value is computed by a compatibility function of the query with the corresponding key [21]. Figure 4 illustrates the two types of attention functions used in the Transformer.

**Scaled Dot-Product Attention:** The scaled dot-product attention is the core attention function of the transform. The function takes queries and keys of dimensions  $d_k$  as input, and values

of dimension  $d_v$ . In practice, the attention function is computed on a set of queries, keys and values simultaneously, by packing them into matrices  $\mathbf{Q}$ ,  $\mathbf{K}$ , and  $\mathbf{V}$ , respectively. Thus, the scaled dot-product attention is computed as

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V} \quad (14)$$

where the scaling of  $\mathbf{Q}\mathbf{K}^T$  is done to counteract the cases with a large  $d_k$ , which would result in the softmax-function having an extremely small gradient [21].

**Multi-Head Attention:** The **multi-head attention** is the used attention function by the Transformer and is an extention to the scaled dot-product attention. This attention function linearly projects the queries, keys and values  $h$  times with different, learned linear projections to  $d_k$ ,  $d_k$ , and  $d_v$  dimensions, respectively. For each of the projected queries, keys and values, the scaled dot-product attention is applied, yielding  $d_v$ -dimensional output values. These output values are finally concatenated and further projected, resulting in the final values. By using multi-head attention, the model is allowed to jointly attend to information from different representation subspaces at different positions [21].

### 2.3.2 The Architecture

The Transformer follow an encoder-decoder structure, where the encoder maps an input sequence  $\mathbf{x}$  to a sequence of continuous representations  $\mathbf{z}$ . The decoder then uses  $\mathbf{z}$  to generate an output sequence  $\mathbf{y}$ , one element at a time. At each step the model consumes the previously generated output element as additional input when generating the next output [21]. Figure 3 illustrates the overall architecture of the Transformer.

**Encoder:** The encoder consists of  $N$  identical layers, where each layer consists of two sub-layers. The first sub-layer is a multi-head **self-attention** layer, and the second sub-layer is a position-wise fully-connected feedforward network. Around each sub-layer is a **residual connection**, where the the input of the sub-layer is added to the output of the sub-layer. This residual connection is then followed by a round of **layer normalization**. In this self-attention layer the keys, values and queries come from the output of the previous layer in the encoder. Each position in the encoder can attend to all positions in the previous layer of the encoder [21].

**Decoder:** The decoder also consists of  $N$  identical layers. In addition to the two sub-layers in each encoder layer, the decoder also consists of a third sub-layer, which performs multi-head attention over the output of the encoder stack (called **cross-attention** [11]). Also here, is residual connections used around each of the sub-layers, followed by layer normalization. To ensure, that the predictions for position  $i$  only depends on the known outputs at positions less than  $i$ , the self-attention sub-layer in the decoder stack is modified. This self-attention sub-layer, allows each position in the decoder to attend to all positions in the decoder up to and including that position. The self-attention sub-layer is modified by masking out all values in the input of the softmax which correspond to illegal connections [21].

**Feedforward Networks:** Each of the layers in the encoder and decoder contains a fully-connected feedforward layer, consisting of two linear transformations with the ReLU activation-function applied in between, as described below [21]

$$FFN(\mathbf{x}) = \max(0, \mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2. \quad (15)$$

**Positional Encoding:** As the model does not contain any recurrences nor any convolutions, it has to carry some other type of information to know about the relative or absolute position of the elements in the input sequence. This is done by adding **positional encodings** to the input of the encoder and decoder stacks. The positional encoding is a vector of size  $d_{model}$  and is defined as the following

$$\mathbf{PE}_{pos,2i} = \sin\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right) \quad (16)$$

$$\mathbf{PE}_{pos,2i+1} = \cos\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right) \quad (17)$$

(18)

where  $pos$  is the position and  $i$  is the dimension [21].

## 2.4 General Machine Learning Terminology

The following section covers the general terminology and some algorithms, that will be used through this paper.

**Loss Function:** Training a machine learning model is generally done by minimizing a pre-defined **loss function**, which is used for measuring the error of the model. One common loss function for regression is the **Mean Squared error (MSE)**. Let  $f^*(x)$  be a vector of ground truth observations and  $f(x; \theta)$  be an estimation of  $f^*(x)$ . Then, MSE is defined as

$$MSE = \frac{1}{n} \sum_{i=1}^n \left( f^*(x^{(i)}) - f(x^{(i)}; \theta) \right)^2. \quad (19)$$

Thus, MSE measures the average squared difference between the ground truth observations and the estimated observations [10].

**Optimizer:** The minimization of the loss function is usually done by an **optimizer**, which is a type of algorithm, that is typically some type of variation of **gradient descent**. Gradient descent is a process that iteratively updates the parameters of a function by

$$\theta \leftarrow \theta - \eta \nabla_{\theta} L(f(x^{(i)}; \theta), f^*(x^{(i)})) \quad (20)$$

where  $\eta$  is called the **learning-rate**, which controls the stepsize of each update, and  $L$  is a loss function such as MSE.

One of the most used optimizers is **ADAM**, which is an optimizer that uses **momentum** and **adaptive learning-rate** to accelerate its learning. Momentum works by accumulating an exponentially decaying moving average of past gradients and continuing to move in their direction, which especially accelerates learning in the face of high curvature, small but consistent gradients, or noisy gradients. Adaptive learning rate is a mechanism for adapting individual learning rates for the model parameters. The idea is as follows; if the partial derivative of the loss, with respect to a given model parameter, remains the same sign, then the learning rate should increase. If that partial derivative changes sign, then the learning rate should decrease. The algorithm uses a hyperparameter  $\rho$  for controlling the length scale of the moving average. The pseudocode of the algorithm has been illustrated in Algorithm 1. [7].

**Online, mini-batch and batch gradient methods:** Unlike gradient descent, ADAM bases its

computations on a **mini-batch** of samples instead of basing it on a single sample. Thus, there are three ways of sampling data for the optimizer. The first way is done by sampling a single sample at each iteration. This class of methods is called **online gradient methods**. The advantage of these methods is, that each iteration is very quick, as we only have to compute the gradient of a single sample, however, the main disadvantage is, that the algorithm uses a lot of iterations. On the other hand we have **batch gradient methods** which uses the average gradient of all of the  $n$  samples of the dataset to compute  $\nabla L(f(\mathbf{x}, \mathbf{y}))$ . By doing so the algorithm only requires one iteration, however, this iteration is much slower than each iteration of the online gradient methods. The last class of methods is **mini-batch gradient methods** which uses the average gradient of  $|B|$  samples of the dataset for computing  $\nabla L(f(\mathbf{x}, \mathbf{y}))$ , where  $B$  is a predefined hyperparameter called the **mini-batch size**. This method lies in between the online gradient methods and batch gradient methods, as it uses fewer iterations than the online gradient methods, however, it uses more iterations than the batch gradient methods. Likewise, each iteration is faster to compute than it is for the batch gradient methods, but slower than it is for the online gradient methods. Further, online gradient methods can lead to a very noisy learning process due to the high variance of the computed gradients of the samples, which can make it more difficult for the optimizer to reach the minimum of the loss function. mini-batch gradient methods and batch gradient methods removes this noise by averaging the computed gradient at each iteration, essentially leading to a quicker learning process [7].

---

**Algorithm 1** Adam [7]

---

**Require:** Learning rate  $\eta$

**Require:** Exponential decay rates for moment estimates,  $\rho_1$  and  $\rho_2$  in  $[0, 1)$

**Require:** Small constant  $\delta$  used for numerical stabilization

**Require:** Initial parameters  $\theta$

Initialize 1st and 2nd moment variables  $s = \mathbf{0}$ ,  $r = \mathbf{0}$

**while** stopping criterion not met **do**

    Sample a minibatch of  $m$  random observations from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$

    Compute gradient:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

$t \leftarrow t + 1$

    Update biased first moment estimate:  $s \leftarrow \rho_1 s + (1 - \rho_1) \mathbf{g}$

    Update biased second moment estimate:  $r \leftarrow \rho_2 r + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}$

    Correct bias in first moment:  $\hat{s} \leftarrow \frac{s}{1 - \rho_1^t}$

    Correct bias in second moment:  $\hat{r} \leftarrow \frac{r}{1 - \rho_2^t}$

    Compute update:  $\Delta\theta = -\eta \frac{\hat{s}}{\sqrt{\hat{r}} + \delta}$

    Apply update:  $\theta = \theta + \Delta\theta$

---

**Layer Normalization:** One problem when optimizing a neural network is, that the gradients with respect to the weights of one layer depends on the outputs of the neurons of the previous previous layer. If these outputs change in a highly correlated way, the model could experience some undesired "covariate shift", which would make the training take longer. One way to reduce these "covariate shifts" is to make use of **layer normalization**, where the layer normalization statistics over all the hidden units in the same layer  $l$  is computed as the following

$$\mu^{(l)} = \frac{1}{H} \sum_{i=1}^H \mathbf{w}_i^{(l)\top} \mathbf{h}^{(l-1)}, \quad \sigma^{(l)} = \sqrt{\frac{1}{H} \sum_{i=1}^H \left( \mathbf{w}_i^{(l)\top} \mathbf{h}^{(l-1)} - \mu^{(l)} \right)^2} \quad (21)$$

where  $H$  is the amount of hidden units of the  $l$ th layer,  $\mathbf{w}_i^{(l)}$  is the incoming learnable weight vector to the  $i$ th hidden unit of the  $l$ th layer, and  $\mathbf{h}^{(l-1)}$  is the output of the preceding layer. Thus, the summed inputs  $\mathbf{w}_i^{(l)\top} \mathbf{h}^{(l-1)}$  is normalized and rescaled by

$$\bar{a}_i = \frac{\mathbf{w}_i^{(l)\top} \mathbf{h}^{(l-1)} - \mu^{(l)}}{\sigma^{(l)}}. \quad (22)$$

Further, the input to the current layer  $l$  also learns an adaptive bias  $b$  and gain  $g$  for each neuron after the normalization. Thus, the  $i$ th input to layer  $l$  is defined as the following

$$h_i^{(l)} = f\left(\frac{g_i(a_i - \mu^{(l)}) + b_i}{\sigma^{(l)}}\right) \quad (23)$$

where  $f$  is a non-linear activation function [3].

**Evaluation:** We often want to know how well a machine learning model performs on an unseen sample. This is commonly done by splitting the dataset into three non-overlapping subsets to avoid a biased evaluation. The first dataset is the **training** dataset which is used for training the machine learning model. The second dataset is the **validation** dataset which is used during training for evaluating various settings of the the used (hyper)parameters. The last dataset is the **testing** dataset, which is then used for evaluating the final machine learning model.

The evaluation of a machine learning model is done by making use of an **evaluation metric**, which depends on the task of the machine learning model. For pose estimation, the **Percentage of Correct Keypoints (PCK)** is commonly used. This metric measures the ratio of predicted keypoints that are closer to their corresponding ground-truth keypoints than some threshold. Often **PCK@0.2**, **PCK@0.1**, **PCK@0.05** are used, which uses 20%, 10% and 5% of the torso diameter as its threshold, respectively [2].

**Dropout:** A machine learning model is said to be **overfitting** if the loss of the training dataset is decreasing, while the loss of the validation dataset is increasing. One can make use of a group of techniques called **regularization** to avoid the overfitting of the network. One common regularization technique is **dropout**, where each parameter of the machine learning model is randomly disabled with probability  $p$  during training. By doing so the parameters are being trained on the dataset fewer times, leading to less of a probability of the model overfitting.

**Epoch:** Often, the optimizer needs multiple "rounds" of using the whole training data for reaching its minimum. Each of these rounds is called an **epoch**. One can pick the number of epochs to iterate through prior to training a model and just stop the training once the optimizer has used all of its epochs. Another method is to use **early-stopping**, where one keeps track of the loss of the model on the validation-set for each epoch. Once the validation-loss has not decreased for  $n$  consecutive epochs, the training is terminated.

### 3 Models

The following section covers the architectures of the various models that will be used in Section 6.3 when we will be performing our experiments.

#### 3.1 Baseline

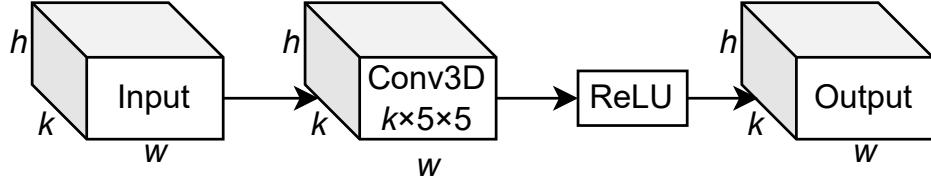


Figure 5: Illustration of the implemented baseline model.

The first model we will be using in Section 6.3 is a very simple baseline model based on a 3-dimensional convolution. Figure 5 illustrates the architecture of the model.

The model takes a sequence of  $T \in \mathbb{R}$  estimated poses  $\hat{\mathcal{P}} = \{\hat{\mathcal{P}}^t\}_{t=1}^T$  as input, where each estimated pose  $\hat{\mathcal{P}}^t \in \mathbb{R}^{K \times h \times w}$  is represented using heatmaps, such that  $K \in \mathbb{R}$  is the amount of keypoints in each estimated pose,  $h \in \mathbb{R}$  is the height of each heatmap and  $w \in \mathbb{R}$  is the width of each heatmap.

Once the data has been passed to the model, the processing of the data is very simple. As illustrated in Figure 5, the model starts by applying a 3-dimensional convolutional layer to the input data. The convolutional layer consists of  $K$  filters, each with a kernel-size of  $T \times 5 \times 5$ . To ensure the input and output of the convolutional layer has the same shape, we pad the input with zeros.

Once the convolutional layer has processed the data, the ReLU activation-function is applied element-wise to the data, resulting in the final prediction of the model.

#### 3.2 Bidirectional Convolutional LSTM

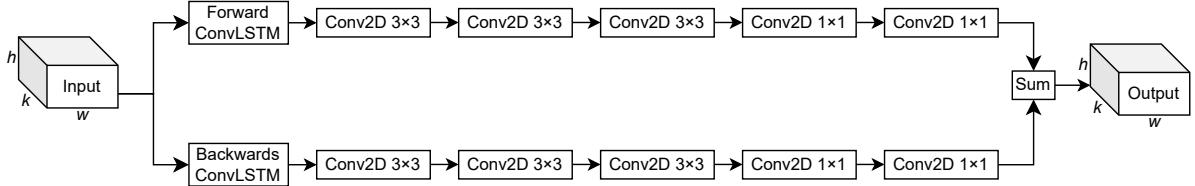


Figure 6: Illustration of the implemented bidirectional convolutional LSTM, where the two sequence orders are summed together.

Our second and third model are based on the LSTM-extension of Unipose-LSTM by Artacho and Savakis [2]. Artacho and Savakis' model was based on a unidirectional convolutional LSTM. However, as we do not require our model to work in real-time, and as we believe a bidirectional convolutional LSTM would be beneficial, our models will be based on a bidirectional convolutional LSTM. The major difference between our second and third model is how they combine the information of the two sequence directions.

Both models take a sequence of  $T \in \mathbb{R}$  estimated poses  $\hat{\mathcal{P}} = \{\hat{\mathbf{P}}^t\}_{t=1}^T$  as input, where each estimated pose  $\hat{\mathbf{P}}^t \in \mathbb{R}^{K \times h \times w}$  is represented as a set of heatmaps, where  $K \in \mathbb{R}$  is the amount of keypoints,  $h \in \mathbb{R}$  is the height of each heatmap and  $w \in \mathbb{R}$  is the width of each heatmap.

### 3.2.1 Summing the Sequence Directions

Figure 6 illustrates the architecture of the second model. The model starts by branching into two separate branches, that processes the estimated poses in opposite sequence order. Each branch processes the estimated poses one frame at a time. This is done by first applying a convolutional LSTM to the input frame at time step  $t \in \mathbb{R}$ , using the preceding output of the convolutional LSTM as the hidden state. Each convolutional LSTM is followed by five 2-dimensional convolutional layers, each applying 128 filters, except for the last convolutional layer of each branch, which applies  $K$  filters. The first three convolutional layers use a kernel size of  $3 \times 3$ , whereas the following two convolutional layers use a kernel size of  $1 \times 1$ . The outputs of the two branches are then summed together element-wise.

All convolutional layers use a stride of one and zero-padding on the input, such that the output of each convolutional layer has the same dimensions as the input to the convolutional layer.

### 3.2.2 Concatenating the Sequence Directions

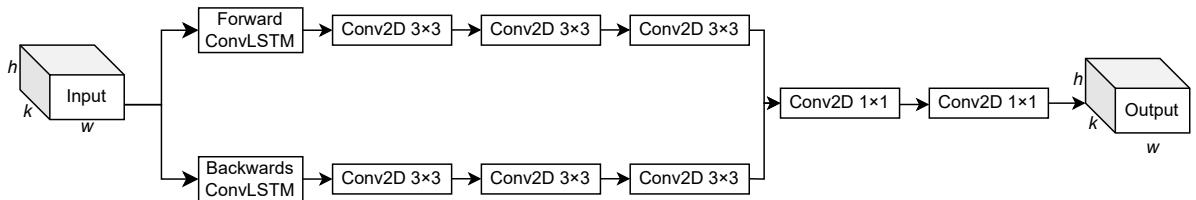


Figure 7: Illustration of the implemented bidirectional convolutional LSTM, where the two sequence orders are concatenated together and processed by two convolutional layers.

Looking at our second model, we see one major problem with the architecture behind it: the model does not have any opportunity of prioritizing one processing order over the other, as the two branches are just summed together. Our third model aims at solving this problem.

Figure 7 illustrates the architecture of our third model. This model is very similar to the previous model. It also starts by branching into two separate branches, that processes the estimated poses in opposite sequence orders, where each branch processes the estimated poses one frame at a time. Similarly to the previous model, this model also starts by applying a convolutional LSTM to the input frame, by using the preceding output of the convolutional LSTM as the hidden state. Further, each convolutional LSTM is also followed by three 2-dimensional convolutional layers, each applying 128 filters with a kernel size of  $3 \times 3$ .

Different from the previous model, this model concatenates the output of each branch to form a  $256 \times h \times w$  tensor, which is then processed by two  $1 \times 1$  2-dimensional convolutions, with 128 and  $K$  filters, respectively.

Similarly to the second model, all convolutional layers also use a stride of one and zero-padding on the input, such that the output of each convolutional layer has the same dimensions as the input to the convolutional layer.

### 3.3 DeciWatch

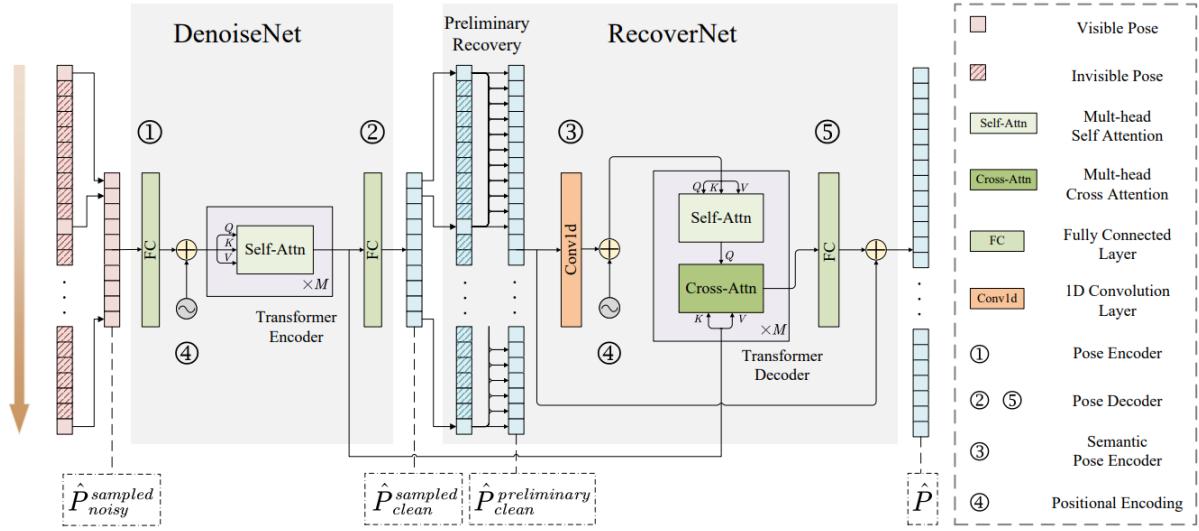


Figure 8: Illustration of the DeciWatch [26]

The last model we implement is a transformer-based model named *DeciWatch*, introduced by Zeng *Et al.*, which is illustrated in Figure 8. The model works by only processing some of the input-frames. It consists of two parts: the *DenoiseNet* and the *RecoverNet*. The aim of DenoiseNet is to denoise the estimated poses given as input to the model, whereas the aim of RecoverNet is to recover the poses of the missing frames. The following description of the model is based on an interpolation of the official paper behind DeciWatch [26].

More specifically, the model takes a sequences of  $T \in \mathbb{R}$  estimated poses  $\hat{\mathcal{P}} = \{\hat{\mathcal{P}}^t\}_{t=1}^T$  as input, where  $\hat{\mathcal{P}}^t$  is represented by 2-dimensional keypoint position. Due to redundancy in consecutive frames and continuity of human poses, the model starts by sampling every  $n$ th frame to select sparse poses  $\hat{\mathcal{P}}_{noisy}^{sampled} \in \mathbb{R}^{\frac{T}{n} \times (2K)}$ , where  $K \in \mathcal{R}$  is the number of keypoints. These sampled poses are then passed to DenoiseNet.

The goal of DenoiseNet is to denoise the sparse poses, that were estimated by a single-frame pose estimator. The denoise process can be formulated as

$$\hat{\mathcal{F}}_{clean}^{sampled} = \text{TransformerEncoder} \left( \hat{\mathcal{P}}_{noisy}^{sampled} \mathbf{W}_{DE} + \mathbf{E}_{pos} \right). \quad (24)$$

That is,  $\hat{\mathcal{P}}_{noisy}^{sampled}$  is first encoded through a linear projection matrix  $\mathbf{W}_{DE} \in \mathbb{R}^{2K \times C}$  and summed with a positional embedding  $\mathbf{E}_{pos} \in \mathbb{R}^{\frac{T}{n} \times C}$ . This is then passed to a transformer-encoder consisting of  $M \in \mathbb{R}$  multi-head Self-Attention blocks, resulting in the noisy poses being embedded into a clean feature  $\hat{\mathcal{F}}_{clean}^{sampled} \in \mathbb{R}^{\frac{T}{n} \times C}$ , where  $C \in \mathbb{R}$  is the embedding dimensions. Lastly, another linear projection matrix  $\mathbf{W}_{DD} \in \mathbb{R}^{C \times 2K}$  is used to obtain the denoised sparse poses

$$\hat{\mathcal{P}}_{clean}^{sampled} = \hat{\mathcal{F}}_{clean}^{sampled} \mathbf{W}_{DD}. \quad (25)$$

After the sparse poses has been denoised as  $\hat{\mathcal{P}}_{clean}^{sampled} \in \mathbb{R}^{\frac{T}{n} \times 2K}$ , the data is passed to the RecoverNet, whose goal is to recover the absent poses. First, a linear transformation  $\mathbf{W}_{PR} \in \mathbb{R}^{T \times \frac{T}{n}}$  is applied to perform preliminary sequence recovery to get  $\hat{\mathcal{P}}_{clean}^{preliminary} \in \mathbb{R}^{T \times 2K}$  by

$$\hat{\mathcal{P}}_{clean}^{preliminary} = \mathbf{W}_{PR} \hat{\mathcal{P}}_{clean}^{sampled}. \quad (26)$$

To improve the recovery of the absent poses a transformer-decoder and positional embedding is used together with a 1D convolutional layer to bring temporal semantics into pose encoding to encode the neighboring  $D \in \mathbb{R}$  frames' poses into pose tokens. Thus, RecoverNet, and the final prediction of DeciWatch, can be summarized by

$$\hat{\mathbf{P}} = \text{TransformerDecoder} \left( \text{Conv1d} \left( \hat{\mathbf{P}}_{\text{clean}}^{\text{preliminary}} \right) + \mathbf{E}_{\text{pos}}, \hat{\mathbf{F}}_{\text{clean}}^{\text{sampled}} \right) \mathbf{W}_{RD} + \hat{\mathbf{P}}_{\text{clean}}^{\text{preliminary}} \quad (27)$$

where  $\mathbf{W}_{RD} \in \mathbb{R}^{C \times 2K}$  is another linear transformation layer. Further, as illustrated by Figure 8, key information is drawn in the the Cross-Attention block by leveraging the denoised features  $\hat{\mathbf{F}}_{\text{clean}}^{\text{sampled}}$ .

To avoid overfitting, dropout is applied to the input of each sub-layer and sums of the embeddings of the transformer-encoder and transformer-decoder, as well as to the positional encodings.

**SKRIV OM HYPERPARAMETER-VALG VED DECIWATCH!!!!!!!!!!!!!!**

## 4 Dataset

To perform the pose estimation in Section 6.3, we need some data on which to train, validate and test our models. The following section describes the datasets that will be used.

Keypoint label	ClimbAlong	BRACE	Penn Action
Head	No	No	Yes
Nose	Yes	Yes	No
Left ear	Yes	Yes	No
Right ear	Yes	Yes	No
Left eye	No	Yes	No
Right eye	No	Yes	No
Left shoulder	Yes	Yes	Yes
Right shoulder	Yes	Yes	Yes
Left elbow	Yes	Yes	Yes
Right elbow	Yes	Yes	Yes
Left wrist	Yes	Yes	Yes
Right wrist	Yes	Yes	Yes
Left pinky	Yes	No	No
Right pinky	Yes	No	No
Left index	Yes	No	No
Right index	Yes	No	No
Left thumb	Yes	No	No
Right thumb	Yes	No	No
Left hip	Yes	Yes	Yes
Right hip	Yes	Yes	Yes
Left knee	Yes	Yes	Yes
Right knee	Yes	Yes	Yes
Left ankle	Yes	Yes	Yes
Right ankle	Yes	Yes	Yes
Left heel	Yes	No	No
Right heel	Yes	No	No
Left toes	Yes	No	No
Right toes	Yes	No	No

Table 1: Overview of the annotated keypoints of the three used datasets

### 4.1 The ClimbAlong Dataset

As the aim of our models is to perform well on climbers, we will be using some annotated data of climbers. For this, ClimbAlong ApS has developed a dataset that we will be using. The dataset consists of videos of various climbers on bouldering walls, where each video contains just a single climber. Figure 9 and 10 illustrates two windows of five consecutive frames of a single video from the ClimbAlong dataset. As shown in the figures, the videos in the dataset contains both static positions, where the climber holds a position for a while, as well as quick movements.

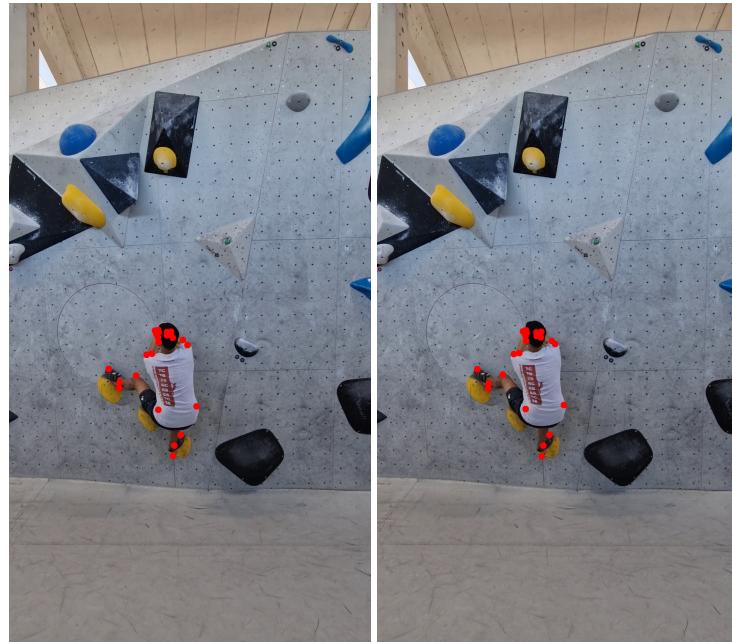
The dataset consists of 29 fully annotated videos and a total of 9,535 fully annotated frames, where each annotation consists of 25 keypoints. Table 1 gives an overview of which keypoints are annotated in the dataset. Each video is filmed in portrait mode with a resolution of



(a) Frame 17

(b) Frame 18

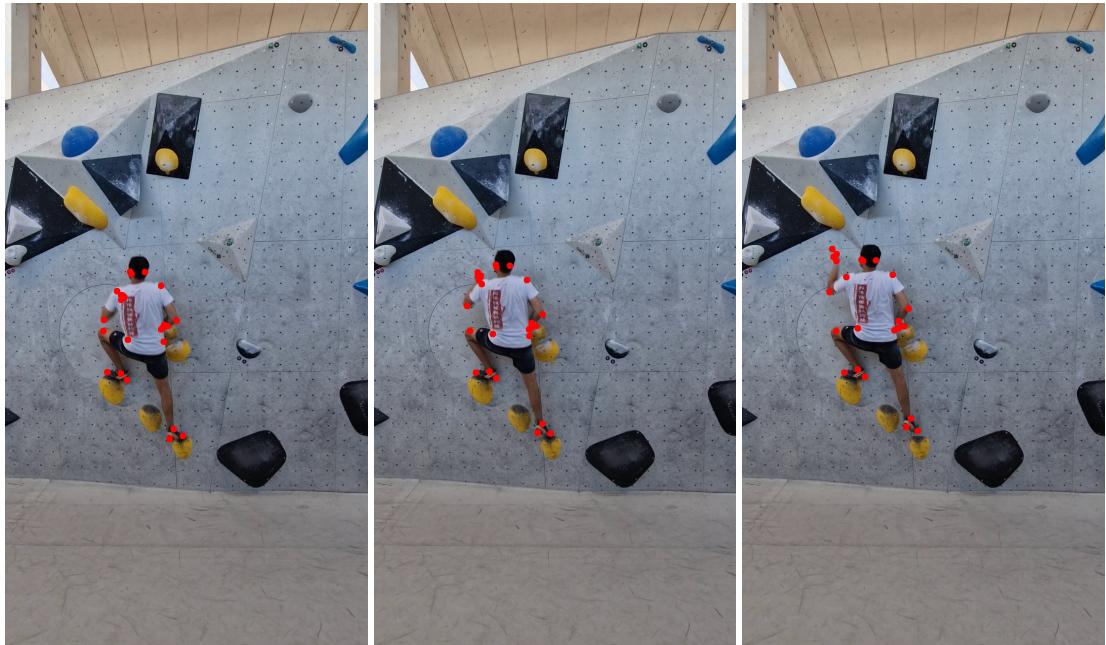
(c) Frame 19



(d) Frame 20

(e) Frame 21

Figure 9: Example of five consecutive frames of a video from the ClimbAlong dataset with the corresponding groundtruth keypoints, where the actor holds his position for a while.



(a) Frame 31

(b) Frame 32

(c) Frame 33



(d) Frame 34

(e) Frame 35

Figure 10: Example of five consecutive frames of a video from the ClimbAlong dataset with the corresponding groundtruth keypoints, where the actor performs a quick movement.

$1080 \times 1920$  and 30 frames per second.

## 4.2 The BRACE Dataset

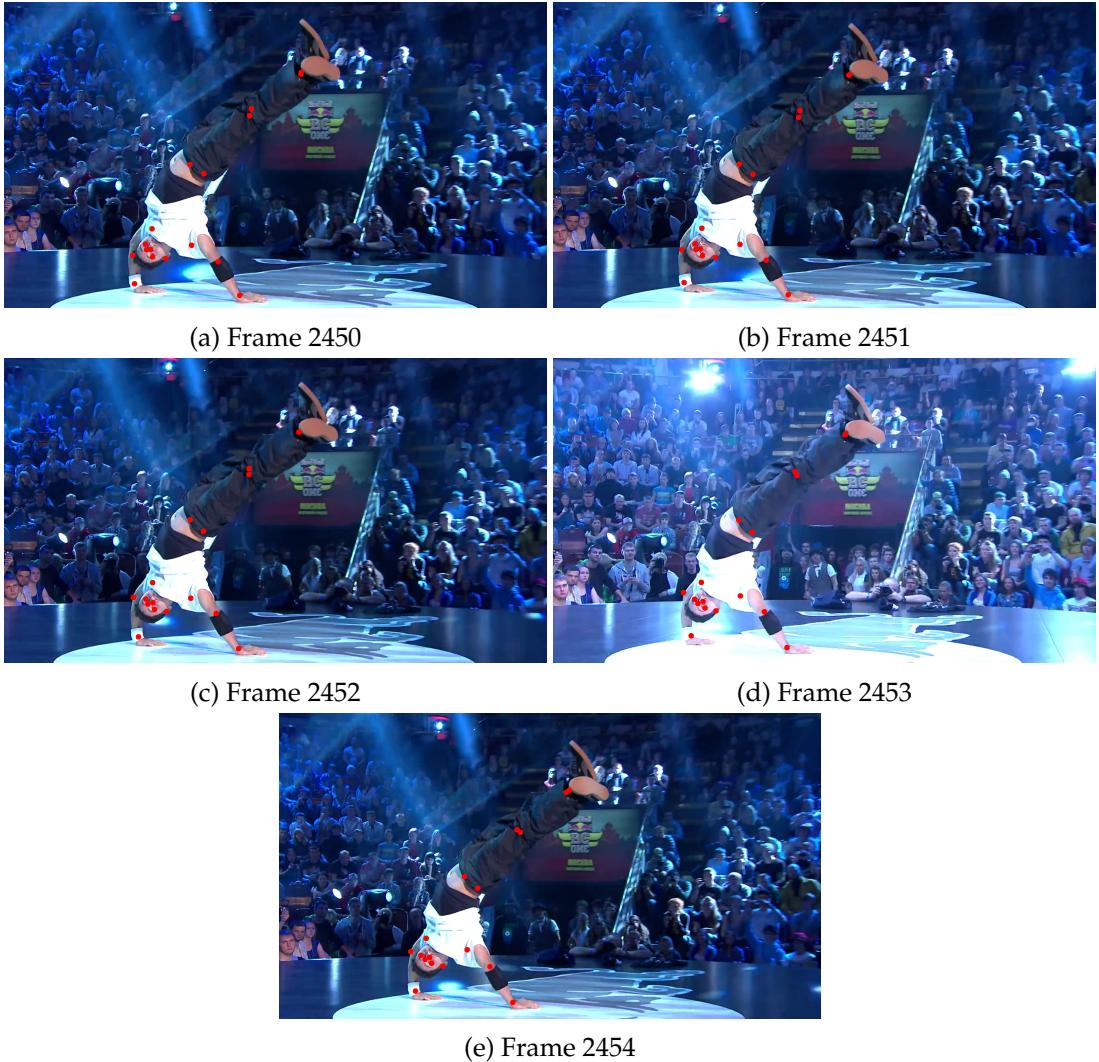


Figure 11: Example of five consecutive frames of a video from the BRACE dataset with the corresponding groundtruth keypoints, where the actor holds his position for a while.

The second dataset we will be using is the *BRACE* dataset [14]. This dataset consists of 1,352 video sequences and a total of 334,538 frames with keypoints annotations of breakdancers. The frames of the video sequences have a resolution of  $1920 \times 1080$  [14].

We chose to use this dataset as breakdancers tend to swap between static and quick poses, as well as containing some acrobatic poses, similarly to the ones seen in the ClimbAlong dataset. Generally, the movements of the BRACE dataset are quicker than the movements of the ClimbAlong dataset. The static poses of the BRACE dataset tend to occur less frequently than the static poses of the the ClimbAlong dataset, as well as the quick movements tend to be quicker than the quick movements of the ClimbAlong dataset. However, as both the actors of both datasets swap between static and quick poses, as well as both datasets containing acrobatic poses, we found the BRACE dataset relevant for our experiments in Section 6.3. Figure 11 and 12 contains two consecutive sequences, each of five frames, that illustrates these two cases.



Figure 12: Example of five consecutive frames of a video from the BRACE dataset with the corresponding groundtruth keypoints, where the actor performs a quick movement.

baseball_pitch	baseball_swing	bench_press
bowling	clean_and_jerk	golf_swing
jumping_jacks	jump_rope	pull_ups
push_ups	sit_ups	squats
strumming_guitar	tennis_forehand	tennis.Serve

Table 2: The original 15 action-types in the Penn Action dataset.

The frames of the video sequences have been annotated by initially using state-of-the-art human pose estimators to extract automatic poses. This was then followed by manually annotating bad keypoints, corresponding to difficult poses, as well as pose outliers. Finally, the automatic and manual annotations were merged by using interpolating. Each frame-annotation consists of 17 keypoints, following the COCO-format, as illustrated in Table 1 [14].

### 4.3 The Penn Action Dataset

The final dataset we will be using is the *Penn Action* dataset [28]. This dataset consists of 2,326 video sequences of 15 different action-types. Table 2 lists these 15 action-types [28].

Each sequence has been manually annotated with human joint annotation, consisting of 13 joints as well as a corresponding binary visibility-flag for each joint. The frames have a resolution within the size of  $640 \times 480$  [28].

Unlike the BRACE dataset, most of the poses in the Penn Action dataset are not very unusual and thus are not very relevant for the poses of climbers. For that reason, we have decided to focus on the action-types that may contain more unusual poses. Thus, we only keep the sequences that have `baseball_pitch`, `bench_press` or `sit_ups` as their corresponding action-type [28]. Further, the movements of the Penn Action dataset tend to be of a more similar pace to the ClimbAlong dataset than the BRACE dataset, making the Penn Action dataset relevant for our task.

In total, we use 307 video sequences from the Penn Action dataset, consisting of a total of 26,036 frames. Figure 13 illustrates five consecutive frames with its groundtruth annotations for one of these video sequences.

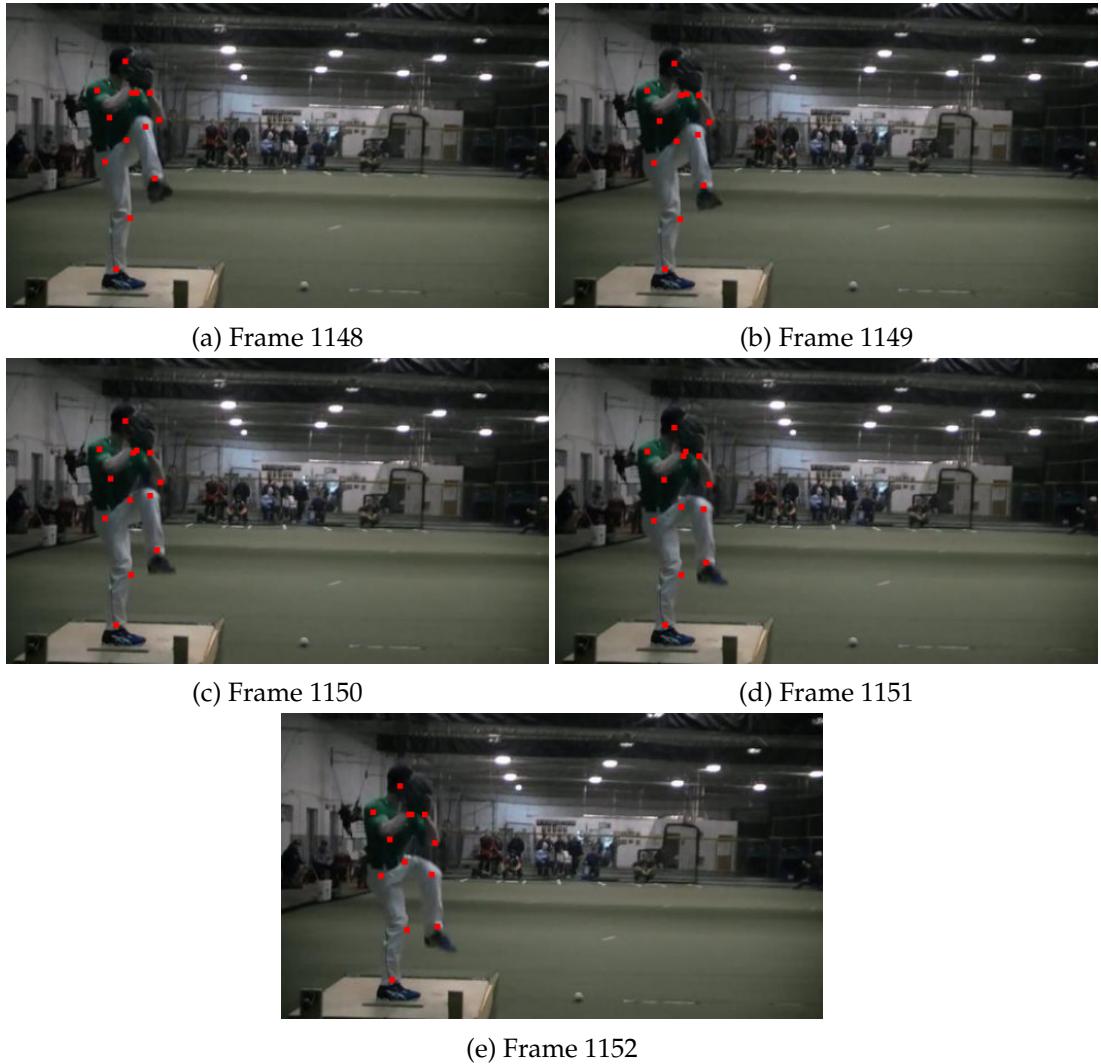


Figure 13: Example of five consecutive frames of a video from the Penn Action dataset with the corresponding groundtruth keypoints.

## 5 Pretraining

As we expect our models to require more data than what is in the ClimbAlong dataset to reach their optimal performance, we have decided to pretrain our models on the BRACE and Penn Action datasets, followed by finetuning them on the ClimbAlong dataset. By doing so we expect our models to yield better results than if we were only using the ClimbAlong dataset, as the pretraining-data will be used for adjusting the randomly initialized weights, whereas the finetuning-data will just be used for specializing the model in performing well on the ClimbAlong dataset. The following section describes the pretraining stage, including the data preprocessing, the used configuration details, as well as the obtained results.

In the pretraining stage we will not be using the already developed pose-estimator, but instead only use our temporal-inclusive models by adding noise to the data such that it simulates the output of the already developed pose-estimator on the ClimbAlong dataset. We do this, as the images of pretraining-data is very different from the ClimbAlong data, making us believe that the already developed pose-estimator will yield some very inaccurate results, as well as some predictions that will be very different from the predictions of the model on the ClimbAlong dataset.

### 5.1 Data Preprocessing

As our models take a sequence of estimated poses as input, we will not be using the images of the frames, hence why we discard the images of all frames from BRACE and Penn Action, such that we only keep the annotated poses.

We start by extracting the bounding-box of the annotated pose in each frame by using the annotated keypoints. Further, we expand each side by 10%, such that no keypoint lies on any of the boundaries of the bounding-box. To ensure that the aspect ratio of the pose is kept later on, we transform the bounding-box into a square by extending the shorter sides, such that they have the same length as the longer sides. Next, we discard everything outside the bounding-box and rescale the bounding-box to have a sidelength of 56, such that it has the same size as the output of the already developed pose-estimator.

Next, we transform each frame into twenty five heatmaps. This is done by creating twenty five  $56 \times 56$  zero-matrices for each frame, such that each zero-matrix represents a single keypoint of a single frame. Further, for each keypoint we insert a fixed value  $c \in \mathbb{R}$  at the position of the keypoint in its corresponding zero-matrix and apply a Gaussian filter with mean  $\mu_{out} = 0$  and standard deviation  $\sigma_{out} = 1$  to smear out each heatmap. For missing keypoints, we do not place the value  $c$  in the corresponding heatmap, making the heatmap consist of only zeros. Further, as Penn Action is the only dataset with the position of the head annotated, as well as the only dataset missing a annotation for the nose, we treat the head-annotation of Penn Action as if it was a nose-annotation, as the position of the two annotation would be very close to each other.

The heatmaps that we produce by following the above description will be used as the groundtruth data. However, as we will be pretraining our models detached from the already developed pose-estimator, we will also need some data as input. We acquire this data by adding some noise to the data, such that they become similar to the output of the already developed pose-estimator, essentially simulating the output of the already developed pose-estimator. The noise is introduced by randomly shifting each keypoint of each sample and by smearing our each keypoint of each sample by using a Gaussian filter with mean  $\mu_{in} = 1$  and some standard deviation  $\sigma_{in} \in \mathbb{R}_{>0}$ . For the shift-value we use  $x \cdot k$ , where the **shifting-scalar**  $k \in \mathbb{R}_{>0}$  is

some fixed positive number and  $x \in \mathbb{R}$ , is equal to 20% of the mean torso-diameter. We clip the position of the shifted keypoints between 0 and 55, such that they cannot be outside of their corresponding heatmaps.

## 5.2 Training Details

### 5.2.1 Data Configuration

For the data we use a window-size of  $k = 5$  frames, as Artacho *et al.* found this to be the optimal number of frames to use [2], making our dataset consist of 345,120 windows. Further, we randomly split our dataset into a training, validation and test set, consisting of 60%, 20% and 20% of the data, respectively, without any overlapping or repeating windows among each other. We insert the fixed value  $c = 255$  in each heatmap at the position of the corresponding keypoint. Lastly, to help the models learn the patterns of the data, we shuffle the three subsets, such that the order of which the data is given to the models does not influence the learning.

As the datasets for the pretraining stage are missing some keypoints of the ClimbAlong dataset, we have to cancel out the training of these missing keypoints, as this would otherwise result in the models learning to never predict the presence of these keypoints. There are multiple ways to do this. We handled it during training by setting the groundtruth heatmaps of the missing keypoints equal to the corresponding predicted heatmaps, making the loss of these missing heatmaps be zero and thus the weights of the model of these heatmaps would not be adjusted.

### 5.2.2 Setups

For each of the four models we use three different setups. In the first setup we uniformly at random sample the standard deviation used by the Gaussian filter at the input data  $\sigma_{in}$  from the set  $\{1, 1.5, 2, 2.5, 3\}$ . We do this, as the output heatmaps of the already developed pose-estimator does not use a fixed standard deviation, making the data a better representation of the pose-estimator output.

As we find it interesting how big of a difference the randomness of the first setup makes, we fix this standard deviation in our second setup, such that we have  $\sigma_{in} = 1$ , thus, essentially, the models only have to learn to translate the input heatmaps.

Finally, with 30 frames per second and a window-size of  $k = 5$ , we suspect that the models might be given too little context to actually be able to effectively smooth out the input data. We could fix this by increasing the window-size, however, we instead chose to make the models work at a lower frame rate, as the increased window-size would also increase the memory usage. Thus, for our last setup we still use a window-size of 5, however, with half the frame rate. For this setup we also sample  $\sigma_{in}$  from the set  $\{1, 1.5, 2, 2.5, 3\}$ .

As we find it difficult to set the optimal value of the shifting-scalar  $k$ , such that the data becomes as close as possible to the finetuning data, we run each setup twice for each model. In the first run we use a shifting-scalar of  $k = 1$ , whereas for the second run we increase this value to  $k = 2$ , making the data a lot more noisy.

### 5.2.3 Training Configuration

For optimizing the weights of the models, we use the ADAM optimizer with a batch size of 16, an initial learning rate of  $10^{-3}$  and the MSE loss-function. During training we keep track of the lowest reached validation loss of an epoch. If this lowest validation loss has not been beaten for five consecutive epochs, we reduce the learning rate by a factor of 0.1. Further, if the lowest validation loss has not been beaten for ten consecutive epochs, we terminate the training of the corresponding model. In case this never happens, we terminate the training of the corresponding model once it has been trained for fifty epochs.

To reduce the likelihood of vanishing or exploding gradients, we initialize the weights of all models by sampling from the Glorot uniform distribution, defined by

$$\mathcal{U}\left(-\frac{\sqrt{6}}{\sqrt{n_j + n_{j+1}}}, \frac{\sqrt{6}}{\sqrt{n_j + n_{j+1}}}\right), \quad (28)$$

where  $n_j$  is the size of the previous layer and  $n_{j+1}$  is the size of the layer which is having its weights initialized [6].

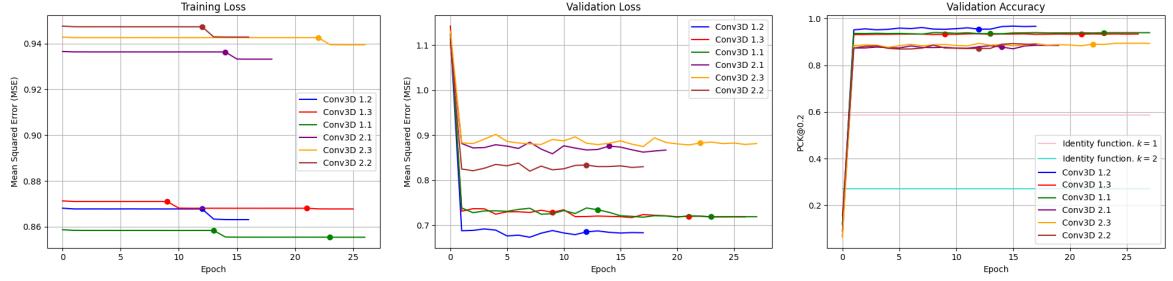
## 5.3 Training and Validation Results

We have in Figure 14 visualized the evolution of the training loss, validation loss and validation PCK@0.2 accuracy for our 24 models. Each model is delt two numbers seperated by a dot such that it follows a  $x.y$ -format (for instance, '2.3'), indicating what type of experiment and setup it was experimented on.  $x$  indicates what shifting-scalar was used when shifting the input keypoints, and  $y$  indicates which of the three setups from Section 5.2 the experiment belongs to. Each plot has at least one dot, which indicates the reduction of the learning-rate due to five consecutive epochs without beating the minimum validation loss.

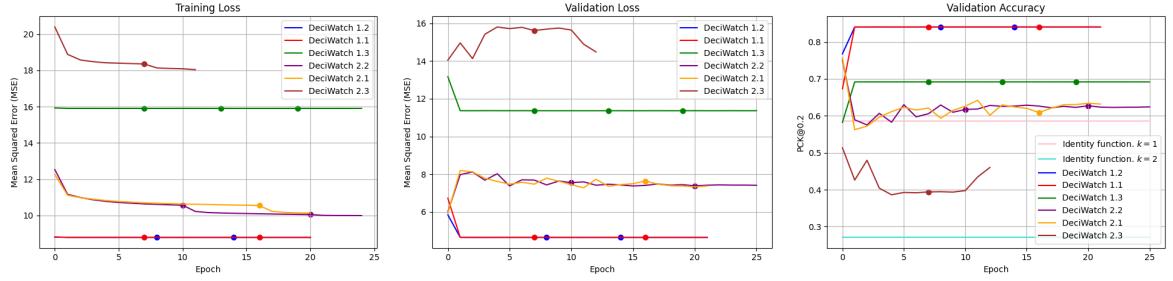
By comparing the validation accuracies of the various models against the identity function we see, that all of the models do learn to somewhat denoise the input data. The simple 3-dimensional convolutional layer seems to be the architecure that generally delivers the greatest results, whereas DeciWatch is generally the architecure that delivers the worst results. Further, the two architectures that are based on bidirectional convolutional LSTMs have close to no difference in their performance, indicating that our raised concern about one of the architectures not being able to prioritize one LSTM-branch over the other is not as great as we hypothesized it to be.

From the figure we can see, that generally the reduction of the learning rate has a pretty substantial effect on the training loss, as the training loss is immediately decreased after a reduction of the learning rate. However, the same cannot be said about the validation loss and accuracy, as these do not immediately decrease after the learning rate has been reduced. Most of the times, the training even terminates before a second reduction of the learning rate takes place, indicating the missing effect of the learning rate reduction on the validation loss and accuracy.

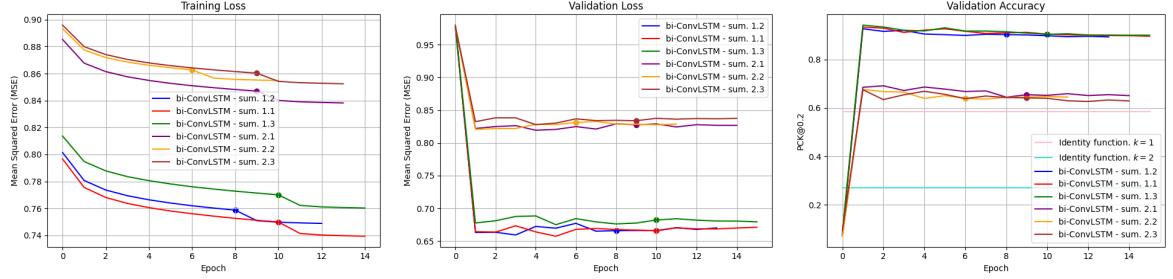
One could look at Figure 14 and argue, that the training of the models generally plateaus as soon as after the first epoch, as the evolution of the validation loss and accuracy seems rather flat after the first epoch. However, if we look at when the learning rate reduction happens, we see that more than the five required epochs have passed before the learning rate reduction is taking place, indicating that the learning has not plateaued. The evolution of the validation loss and accuracy just seem flat, because the performance difference between the first epoch and



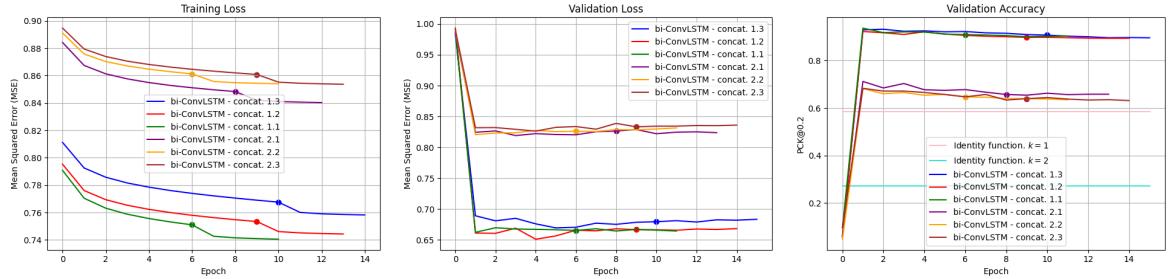
(a) Pretraining results of the 3-dimensional convolutional layer.



(b) Pretraining results of DeciWatch.



(c) Pretraining results of the bidirectional convolutional LSTM with summing.



(d) Pretraining results of the bidirectional convolutional LSTM with concatenation.

Figure 14: Evolution of the training loss, validation loss and validation PCK@0.2 accuracy of the 24 models during training, as well as the validation PCK@0.2 accuracy of the identity function of the two datasets. The dots indicates a reduction of learning-rate. First row: 3-dimensional convolutional layer. Second row: DeciWatch. Third row: Bidirectional Convolutional LSTM with summing. Fourth row: Bidirectional Convolutional LSTM with concatenation.

the models with their initial weights is so great, that we can have difficulty seeing the performance differences of the following epochs.

Further, we can clearly see the effects of the shifting-scalar on the training of the models, as all of the models performs better on the dataset with a low shifting-scalar, than on their corresponding counterpart with a higher shifting-scalar. DeciWatch 1.3 and DeciWatch 2.3 do however seems to be struggling quite a bit, as the validation accuracies of these models are much lower than the other DeciWatch models with the same noise-scalar. However, we have to keep in mind here, that DeciWatch already only considers every fifth frame. This, combined with the fact that setup 3 is only consideren every other frame, makes these two models only consider every tenth frame, which of course makes the prediction very difficult, as these models then have to recover a lot of frames.

## 5.4 Test Results

Accuracy metric	PCK@0.05			PCK@0.1			PCK@0.2		
Threshold distance	1.1	1.2	1.3	1.1	1.2	1.3	1.1	1.2	1.3
Setup	1.1	1.2	1.3	1.1	1.2	1.3	1.1	1.2	1.3
Identity function									
Conv3D	5.44	33.7	4.61	30.1	60.7	24.1	<b>94.4</b>	<b>97.0</b>	94.0
DeciWatch	<b>54.1</b>	<b>54.1</b>	<b>46.4</b>	<b>68.7</b>	<b>68.7</b>	54.8	84.0	84.0	68.9
bi-ConvLSTM - sum.	20.9	21.4	23.5	64.0	61.7	<b>65.8</b>	93.8	93.1	<b>94.6</b>
bi-ConvLSTM - concat.	21.3	20.7	20.9	62.3	61.1	61.5	94.1	92.7	93.6

Table 3: Testing accuracies of the various developed models for shifting-scalar  $k = 1$ .

Accuracy metric	PCK@0.05			PCK@0.1			PCK@0.2		
Threshold distance	1.1	1.2	1.3	1.1	1.2	1.3	1.1	1.2	1.3
Setup	1.1	1.2	1.3	1.1	1.2	1.3	1.1	1.2	1.3
Identity function									
Conv3D	0.80	2.00	0.18	13.7	18.6	9.6	<b>89.2</b>	<b>89.7</b>	<b>90.1</b>
DeciWatch	<b>8.18</b>	8.81	5.13	28.7	27.1	18.2	65.1	63.6	48.6
bi-ConvLSTM - sum.	8.05	<b>9.42</b>	<b>7.78</b>	<b>30.4</b>	<b>31.5</b>	<b>29.4</b>	69.9	68.1	68.2
bi-ConvLSTM - concat.	6.91	8.72	7.77	29.4	31.3	29.0	71.9	68.8	69.1

Table 4: Testing accuracies of the various developed models for shifting-scalar  $k = 2$ .

We have in Table 3 and Table 4 illustrated the test results of the various models on the testsets for PCK@0.05, PCK@0.1 and PCK@0.2.

For shifting-scalar  $k = 1$  we see, that DeciWatch is the architecure that yields the best results for PCK@0.05 and PCK@0.1. However, for PCK@0.2 DeciWatch actually tends to be the worst performing model, which means that most of the predicitons of DeciWatch that are considered correct, are actually very close to the groundtruth, whereas the predictions that are incorrect are pretty far off.

On the other hand, the simple 3-dimensional convolutional layer seems to be the architecture that generally performs the worst for PCK@0.05 and PCK@0.1 and actually the architecture that performs the best for PCK@0.2. Thus, the 3-dimensional convolutional layer is very good at inferencing a rough estimate of the position of the keypoints, however, the exact position of

these predictions are generally incorrect.

Generally, the bidirectional convolutional LSTMs yield some decent results, as the predictions of these models are always pretty far away from the worst-performing model for any of the three used accuracy metrics. Similarly to the validation results, the test results do not indicate any major performance differences between the two different types of bidirectional convolutional LSTMs. Likewise, similarly to the validation results, the test results indicate a major effect of a low frame rate for DeciWatch.

If we take a closer look at the four types of architectures we see, that the 3-dimensional convolutional layer generally performs much better on the second setup than it does on the two other setups. As the models of the first and third setup have to learn both translation and scaling, where for the second setup the model only has to learn translation, we can clearly see how big of an impact the task of scaling the data has for this simple model. As for comparison, the performances of the other architectures have generally a lot smaller variance, indicating these models being less impacted by the setup. This is especially visible for DeciWatch, which has no performance differences between the two first setups. However, this is due the architecture working on keypoint-coordinates instead of heatmaps, making it scaling invariant and thus making it, such that there are no differences between the two first setups.

The testing results of shifting-scalar  $k = 2$  generally follow the same pattern with some minor differences. In this case, the bidirectional convolutional LSTM with summing is actually the model that yields the best results for PCK@0.05 and PCK@0.1 instead of DeciWatch, as it was in the case for  $k = 1$ . However, DeciWatch does still yield some results that are pretty close to the results of the convolutional LSTM with summing.

Unsurprisingly, all of the models performed worse as the shifting-noise was double. Nonetheless, the 3-dimensional convolutional layer is still the architecture that yields the best results when considering PCK@0.2. However, unlike the previous case, this architecture now completely outperforms the three alternative architectures, when considering PCK@0.2, but is on the other hand completely outperformed by the other models when considering PCK@0.05. This indicates, that the architecture can still yield decent rough estimations of the keypoints when a lot of noise is added, however, it does also indicate, that the exact position of these keypoints are very much effected by the level of noise.

Which architecture is the best at denoising the data depends on one's needs. If one just need a rough estimate of the position of the keypoints, then the 3-dimensional convolutional layer would be the optimal choice, as this model has the highest PCK@0.2-values. However, if one need the model which yields the most exact predictions, then one of the bidirectional convolutional LSTMs would be the optimal choice.

We have further in Table 5 and Table 6 illustrated the keypoint-specific testing PCK@0.2 of the various models.

- $k = 1$ 
  - Wrist og ankle er de dårligste
  - Ears and Shoulders er de bedste
- $k = 2$ 
  - Elbow og knee er de dårligste
  - Ears and shoulders er de bedste

Model	Best Epoch	Ear	Shoulder	Elbow	Wrist	Hip	Knee	Ankle
Conv3D 1.1	8	94.3	94.6	94.2	93.9	95.8	94.6	93.6
Conv3D 1.2	15	95.0	<b>97.5</b>	<b>98.7</b>	<b>96.6</b>	<b>98.3</b>	<b>98.6</b>	<b>95.1</b>
Conv3D 1.3	12	94.1	93.7	95.1	92.9	97.0	93.4	92.4
DeciWatch 1.1	3	87.2	88.1	84.1	78.6	86.9	82.7	79.0
DeciWatch 1.2	11	87.2	88.1	84.1	78.6	86.9	82.8	79.0
DeciWatch 1.3	22	73.3	74.9	67.4	61.0	72.3	67.8	63.2
bi-ConvLSTM sum. 1.1	1	95.4	94.4	92.8	91.9	92.4	92.7	92.2
bi-ConvLSTM sum. 1.2	1	94.9	93.4	94.4	93.2	93.5	92.7	94.3
bi-ConvLSTM sum. 1.3	1	95.5	92.2	96.6	<b>96.6</b>	92.5	93.6	95.5
bi-ConvLSTM concat. 1.1	1	95.3	95.1	94.1	94.2	93.9	94.4	92.8
bi-ConvLSTM concat. 1.2	1	95.7	93.3	94.6	91.4	90.8	91.5	91.2
bi-ConvLSTM concat. 1.3	2	<b>95.9</b>	91.8	93.7	95.4	91.5	94.1	94.0
Total		92.0	91.4	90.8	88.7	91.0	89.9	88.5

Table 5: Keypoint-specific testing PCK@0.2-accuracies of the various models for shiting-scalar  $k = 1$ .

Model	Best Epoch	Ear	Shoulder	Elbow	Wrist	Hip	Knee	Ankle
Conv3D 1.1	17	85.1	89.5	91.0	92.0	86.0	89.9	91.9
Conv3D 1.2	15	<b>85.9</b>	89.7	91.3	<b>92.2</b>	88.5	90.8	91.0
Conv3D 1.3	26	85.2	<b>90.9</b>	<b>91.6</b>	91.9	<b>88.7</b>	<b>91.4</b>	<b>92.0</b>
DeciWatch 1.1	11	77.1	72.6	61.0	57.9	70.4	56.7	58.4
DeciWatch 1.2	5	76.1	71.2	59.6	59.7	68.8	58.3	54.0
DeciWatch 1.3	2	54.8	42.0	41.9	45.3	52.9	40.3	53.6
bi-ConvLSTM sum. 1.1	2	81.0	66.6	66.7	69.2	66.1	72.3	65.3
bi-ConvLSTM sum. 1.2	1	79.9	68.2	66.2	66.6	60.6	68.7	65.1
bi-ConvLSTM sum. 1.3	1	83.8	64.9	63.5	67.3	63.7	65.2	68.5
bi-ConvLSTM concat. 1.1	1	83.9	67.1	66.2	75.0	68.9	67.6	71.7
bi-ConvLSTM concat. 1.2	1	81.7	67.8	66.1	62.7	64.7	65.3	72.0
bi-ConvLSTM concat. 1.3	11	81.8	65.2	65.9	66.2	63.7	67.8	68.7
Total		79.7	72.2	69.2	70.5	70.2	69.5	71.0

Table 6: Keypoint-specific testing PCK@0.2-accuracies of the various models for shiting-scalar  $k = 2$ .

## 6 Finetuning

As we now have pretrained our models, we need to finetune the models, such that they are specialized to yield optimal results on the ClimbAlong dataset. The following section describes the finetuning of these models. This includes the various experiments we perform, the preprocessing of the data, the configuration details we use, as well as the obtained results.

In the finetuning stage we will be using the already developed pose-estimator to train our temporal-inclusive models. However, we will be freezing the pose-estimator, such that the weights of the model will not change during the training and we will thus only train our temporal-inclusive models. We do this for the following three reasons: (1) the training of the models will be quicker, as we just need to train the temporal-inclusive models and not the already developed pose-estimator, (2) we get a greater understanding of the effects of our models when combined with the pose-estimator, as we can clearly see how big of a difference it makes by adding our temporal-inclusive models, and (3) we lower the probability of overfitting, as we have less tuneable parameters.

### 6.1 Data Preprocessing

For the ClimbAlong dataset we perform only minor preprocessing. First, the preprocessing of each video is done by having the already developed pose-estimator process the video, such that we have the output heatmaps of the pose-estimator, containing all of the pose-estimations of each video. Next, we preprocess the heatmaps by setting all negative values to 0 and normalizing each heatmap, such that each heatmap sums up to the fixed value  $c = 255$  that we used when preprocessing the BRACE and Penn Action datasets, essentially making the heatmaps more similar to the preprocessed heatmaps of BRACE and Penn Action. These heatmaps will then be used as the input for our models.

For the groundtruth heatmaps we create twenty five heatmaps of each frame, similarly to how we did it for the BRACE and Penn Action datasets, however, in this case we use the predicted bounding-box of the pose-estimator as our bounding-box. In cases where the groundtruth keypoint is placed outside of the bounding-box, we place the groundtruth keypoint at the closest border of the bounding-box.

### 6.2 Training Details

**Data Configuration** Generally, we follow a similar approach to how we did in the pretraining stage. We again use a window-size of  $k = 5$  frames, resulting in a total of 9,419 windows. Also here are we using  $c = 255$  as a representation of the placement of each keypoint. We also split the dataset into a non-overlapping and non-repeating training, validation and test set, consisting of 60%, 20% and 20% of the data, respectively. However, we note that one incorrect frame can have a huge impact on the evaluation results, as this frame is used five times during evaluation, due to the small dataset size. For that reason, for the validation and test set we make sure that none of the windows of the same set.

**Setups** As the finetuning dataset is so small, the fitting of the models is very quick, making us fit all of the developed models from the pretraining stage. For each model we pick the epoch from the pretraining stage, that yielded the highest validation accuracy and use that for finetuning.

**Training Configuration** The optimization parameters are very similar to the ones from the pretraining stage. We again use the ADAM optimizer with a batch size of 16 and the MSE

loss-function. During training, we again keep track of the lowest reached validation loss of an epoch and use learning-rate reduction and early-stopping in a similar manner to how we did in the pretraining stage. However, unlike the pretraining stage, we here use a smaller initial learning rate of  $10^{-4}$ , as the weights only need to be fineadjusted, making us believe that greater learning rate would skew the weights too much.

### 6.3 Results

- Vi skal kun bruge én epoch - nok pga. pretræning

## Technical Details

## 7 Discussion

- Ved ikke PA's fps
- BRACE er landscape
- Der er noget usikkerhed i BRACE da den er annoteret automatisk
- BRACE of Penn Action kan semantisk være meget forskellig
- Jeg shifter efter jeg har udvidet min data med 10%. Dette kan give noget bias, da der er større sandsynlighed for, at dataen ligger på billedets kant
- Vanishing gradient for unipose - den fungerer bedre på range [0, 255] end på range [0, 1]
- Configuration details (se Camilla's)
- Future work
  - Transformer uden skipping af frames
  - Dataaugmentation til CA
  - Træn uden PA
  - Vision transformer
  - Handling af data uden for pred bbox ved CA-data
  - Hvorfor performer mine modeller så meget dårligere på pretraining end på finetuning?
    - Vi kombinerer to meget forskellige datasæt
    - PA frame rates kan ødelægge det
    - Pretrain kan være sværre, fordi vi tiløjer (for meget) støj
    - Mask RCNN kommer allerede med ret gode resultater, så vi skal kun justerer end smule
    - BRACE groundtruth kan også indeholde noget noise da det kommer fra en model, som kan være sværd at predicte.
  - Skulle have brugt en længere patience til early-stopping
  - Hvor meget hjalp det faktisk at pretraine?

## **8 Conclusion**

## 9 References

- [1] André Oskar Andersen. "2D Articulated Human Pose Estimation, Using Explainable Artificial Intelligence". Bachelor's Thesis. University of Copenhagen, Department of Computer Science, 2020.
- [2] Bruno Artacho and Andreas Savakis. *UniPose: Unified Human Pose Estimation in Single Images and Videos*. 2020. DOI: [10.48550/ARXIV.2001.08095](https://doi.org/10.48550/ARXIV.2001.08095). URL: <https://arxiv.org/abs/2001.08095>.
- [3] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. *Layer Normalization*. 2016. arXiv: [1607.06450 \[stat.ML\]](https://arxiv.org/abs/1607.06450).
- [4] Zhe Cao, Gines Hidalgo, Tomas Simon, Shih-En Wei, and Yaser Sheikh. *OpenPose: Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields*. 2018. DOI: [10.48550/ARXIV.1812.08008](https://doi.org/10.48550/ARXIV.1812.08008). URL: <https://arxiv.org/abs/1812.08008>.
- [5] Rohit Girdhar, Georgia Gkioxari, Lorenzo Torresani, Manohar Paluri, and Du Tran. *Detect-and-Track: Efficient Pose Estimation in Videos*. 2017. DOI: [10.48550/ARXIV.1712.09184](https://doi.org/10.48550/ARXIV.1712.09184). URL: <https://arxiv.org/abs/1712.09184>.
- [6] Xavier Glorot and Yoshua Bengio. "Understanding the difficulty of training deep feed-forward neural networks". In: *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings. 2010, pp. 249–256.
- [7] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. <http://www.deeplearningbook.org>. MIT Press, 2016.
- [8] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. *Mask R-CNN*. 2017. DOI: [10.48550/ARXIV.1703.06870](https://doi.org/10.48550/ARXIV.1703.06870). URL: <https://arxiv.org/abs/1703.06870>.
- [9] Sepp Hochreiter and Jürgen Schmidhuber. "Long Short-Term Memory". In: *Neural Computation* 9.8 (1997), pp. 1735–1780.
- [10] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An Introduction to Statistical Learning with Applications in R. First Edition*. Springer, 2017.
- [11] Vaclav Kosar. *Cross-Attention in Transformer Architecture*. <https://vaclavkosar.com/ml/cross-attention-in-transformer-architecture>.
- [12] Yann LeCun, Yoshua Bengio, et al. "Convolutional networks for images, speech, and time series". In: *The handbook of brain theory and neural networks* 3361.10 (1995), p. 1995.
- [13] Yue Luo, Jimmy Ren, Zhouxia Wang, Wenxiu Sun, Jinshan Pan, Jianbo Liu, Jiahao Pang, and Liang Lin. *LSTM Pose Machines*. 2017. DOI: [10.48550/ARXIV.1712.06316](https://doi.org/10.48550/ARXIV.1712.06316). URL: <https://arxiv.org/abs/1712.06316>.
- [14] Davide Moltisanti, Jinyi Wu, Bo Dai, and Chen Change Loy. "BRACE: The Breakdancing Competition Dataset for Dance Motion Synthesis". In: *European Conference on Computer Vision (ECCV)* (2022).
- [15] Alejandro Newell, Kaiyu Yang, and Jia Deng. *Stacked Hourglass Networks for Human Pose Estimation*. 2016. DOI: [10.48550/ARXIV.1603.06937](https://doi.org/10.48550/ARXIV.1603.06937). URL: <https://arxiv.org/abs/1603.06937>.
- [16] Tomas Pfister, James Charles, and Andrew Zisserman. *Flowing ConvNets for Human Pose Estimation in Videos*. 2015. DOI: [10.48550/ARXIV.1506.02897](https://doi.org/10.48550/ARXIV.1506.02897). URL: <https://arxiv.org/abs/1506.02897>.
- [17] Leonid Pishchulin, Mykhaylo Andriluka, Peter Gehler, and Bernt Schiele. "Poselet Conditioned Pictorial Structures". In: *2013 IEEE Conference on Computer Vision and Pattern Recognition*. 2013, pp. 588–595. DOI: [10.1109/CVPR.2013.82](https://doi.org/10.1109/CVPR.2013.82).

- [18] Moacir Antonelli Ponti, Leonardo Sampaio Ferraz Ribeiro, Tiago Santana Nazare, Tu Bui, and John Collomosse. "Everything You Wanted to Know about Deep Learning for Computer Vision but Were Afraid to Ask". In: *2017 30th SIBGRAPI Conference on Graphics, Patterns and Images Tutorials (SIBGRAPI-T)*. 2017, pp. 17–41. DOI: [10.1109/SIBGRAPI-T.2017.12](https://doi.org/10.1109/SIBGRAPI-T.2017.12).
- [19] Xingjian Shi, Zhourong Chen, Hao Wang, Dit-Yan Yeung, Wai-kin Wong, and Wang-chun Woo. *Convolutional LSTM Network: A Machine Learning Approach for Precipitation Nowcasting*. 2015. arXiv: [1506.04214 \[cs.CV\]](https://arxiv.org/abs/1506.04214).
- [20] Yuandong Tian, C. Lawrence Zitnick, and Srinivasa G. Narasimhan. "Exploring the Spatial Hierarchy of Mixture Models for Human Pose Estimation". In: *Computer Vision – ECCV 2012*. Ed. by Andrew Fitzgibbon, Svetlana Lazebnik, Pietro Perona, Yoichi Sato, and Cordelia Schmid. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, pp. 256–269. ISBN: 978-3-642-33715-4.
- [21] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. *Attention Is All You Need*. 2017. DOI: [10.48550/ARXIV.1706.03762](https://doi.org/10.48550/ARXIV.1706.03762). URL: <https://arxiv.org/abs/1706.03762>.
- [22] Shih-En Wei, Varun Ramakrishna, Takeo Kanade, and Yaser Sheikh. *Convolutional Pose Machines*. 2016. DOI: [10.48550/ARXIV.1602.00134](https://doi.org/10.48550/ARXIV.1602.00134). URL: <https://arxiv.org/abs/1602.00134>.
- [23] Yufei Xu, Jing Zhang, Qiming Zhang, and Dacheng Tao. *ViTPose: Simple Vision Transformer Baselines for Human Pose Estimation*. 2022. DOI: [10.48550/ARXIV.2204.12484](https://doi.org/10.48550/ARXIV.2204.12484). URL: <https://arxiv.org/abs/2204.12484>.
- [24] Sen Yang, Zhibin Quan, Mu Nie, and Wankou Yang. *TransPose: Keypoint Localization via Transformer*. 2020. DOI: [10.48550/ARXIV.2012.14214](https://doi.org/10.48550/ARXIV.2012.14214). URL: <https://arxiv.org/abs/2012.14214>.
- [25] Yi Yang and Deva Ramanan. "Articulated Human Detection with Flexible Mixtures of Parts". In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 35.12 (2013), pp. 2878–2890. DOI: [10.1109/TPAMI.2012.261](https://doi.org/10.1109/TPAMI.2012.261).
- [26] Ailing Zeng, Xuan Ju, Lei Yang, Ruiyuan Gao, Xizhou Zhu, Bo Dai, and Qiang Xu. *Deci-Watch: A Simple Baseline for 10x Efficient 2D and 3D Pose Estimation*. 2022. DOI: [10.48550/ARXIV.2203.08713](https://doi.org/10.48550/ARXIV.2203.08713). URL: <https://arxiv.org/abs/2203.08713>.
- [27] Aston Zhang, Zachary C. Lipton, Mu Li, and Alexander J. Smola. "Dive into Deep Learning". In: *arXiv preprint arXiv:2106.11342* (2021).
- [28] Weiyu Zhang, Menglong Zhu, and Konstantinos G. Derpanis. "From Actemes to Action: A Strongly-Supervised Representation for Detailed Action Understanding". In: *2013 IEEE International Conference on Computer Vision*. 2013, pp. 2248–2255. DOI: [10.1109/ICCV.2013.280](https://doi.org/10.1109/ICCV.2013.280).