

1 Dataset

To perform the pose estimation in section ?? and section ??, we need some data on which to train, validate and test our models. The following section describes the datasets that will be used. This starts off in section 1.1, where we describe the dataset provided by ClimbAlong. This is then followed by section 1.2, where we describe the datasets that we will be using during pretraining of our models.

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

1.1 The Finetuning Dataset

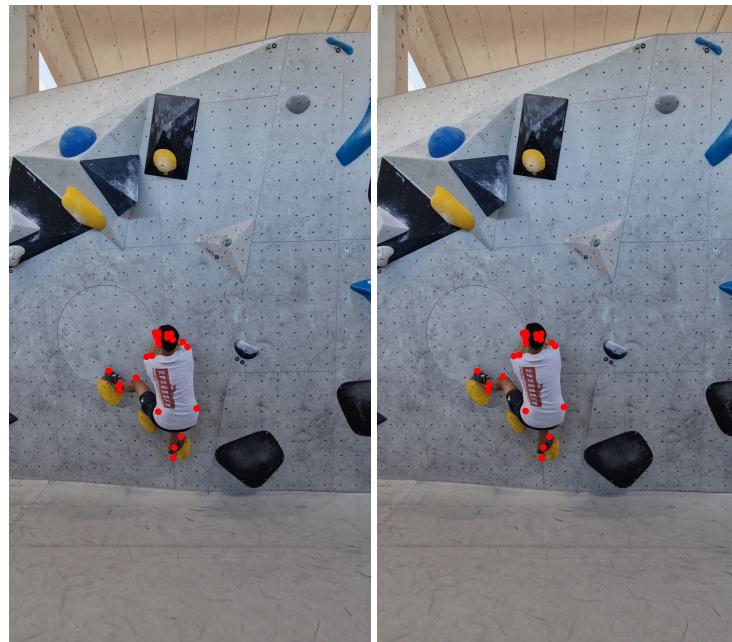
As the aim of our models is to perform well on climbers, we will be using some annotated data of climbers. For this, ClimbAlong has developed a dataset that we will be using. The dataset consists of videos of various climbers on bouldering walls, where each video contains just a single climber. Figure 1 and 2 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 movements, where the climber holds a position for a while, as well as quick movements.



(a) Frame 17

(b) Frame 18

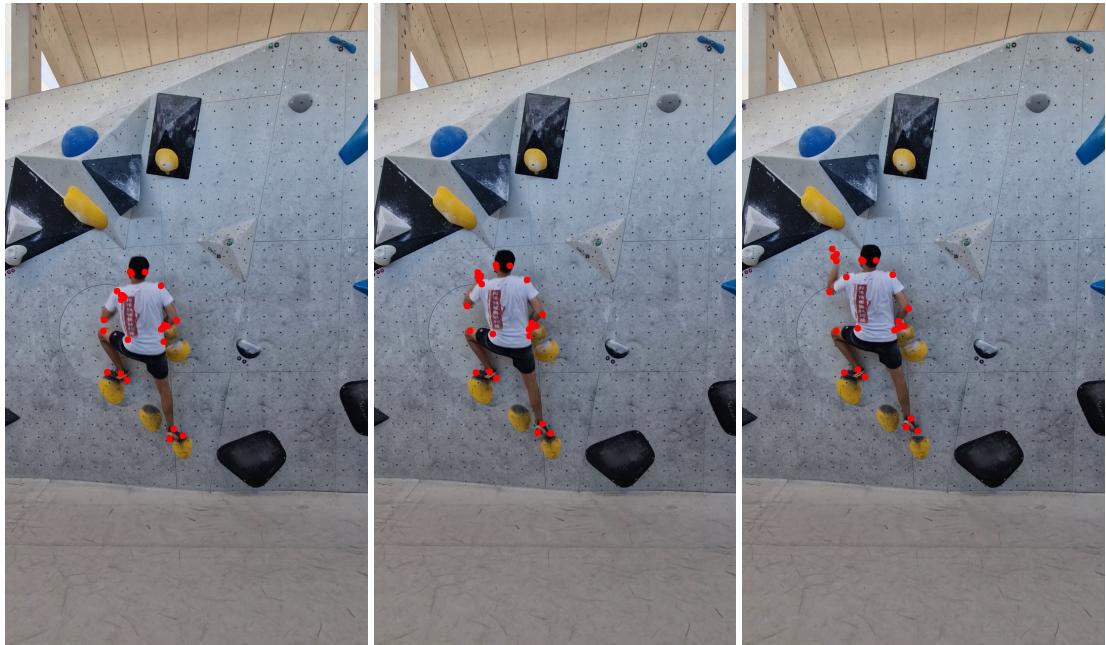
(c) Frame 19



(d) Frame 20

(e) Frame 21

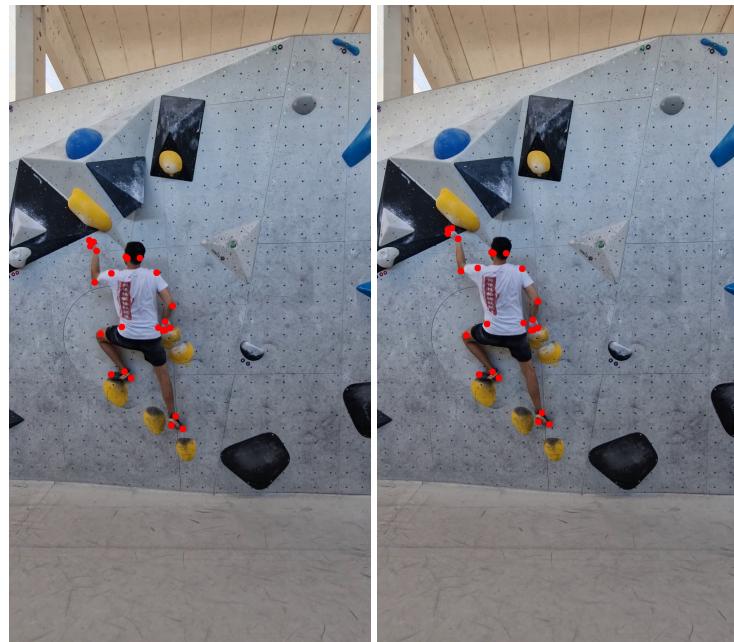
Figure 1: Example of five consecutive frames of a video from the ClimbAlong dataset with the corresponding ground truth keypoints, where the actor holds his position for a while.



(a) Frame 31

(b) Frame 32

(c) Frame 33



(d) Frame 34

(e) Frame 35

Figure 2: Example of five consecutive frames of a video from the ClimbAlong dataset with the corresponding ground truth 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 [2].

The dataset consists of 30 fully annotated videos and a total of 10,293 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 1080×1920 and 30 frames per second. We will throughout this project be referring to this dataset as either the **ClimbAlong dataset** or the **finetuning dataset**.

1.2 The Pretraining Dataset

As we will in section ?? be pretraining our models before specializing our models on bouldering, we will be requiring some pretraining data. For this we use the BRACE dataset [1] and parts of the Penn Action dataset [2]. We chose these datasets, as they are rather similar to the ClimbAlong dataset, as they also mostly consist of unusual movements and only very few usual movements, such as walking. Generally, the movements of BRACE tend to be quicker than the movements of ClimbAlong, whereas the movements of Penn Action tend to be of a more similar pace. Thus, by incorporating both BRACE and Penn Action, we should capture both the fast and slow movements of the ClimbAlong dataset. The following section introduces these datasets for pretraining.

1.2.1 The BRACE Dataset

The first pretraining dataset we will be using is the BRACE dataset [1]. This dataset consists of 1,352 video sequences and a total of 334,538 frames with keypoints annotations of break-dancers. The frames of the video sequences have a resolution of 1920×1080 [1].

Generally, the movements of the BRACE dataset are quicker than the movements of the ClimbAlong dataset. Similarly, the annotated person of BRACE tend to swap between static and quick movements just like the annotated person 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. Figure 3 and 4 contains two consecutive sequences, each of five frames, that illustrates these two cases.

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 [1].

1.2.2 The Penn Action Dataset

The second pretraining dataset we will be using is the Penn Action dataset [2]. This dataset consists of 2,326 video sequences of 15 different action-types. Table 2 lists these 15 action-types [2]. Each frame have a resolution within the size of 640×480 [2].

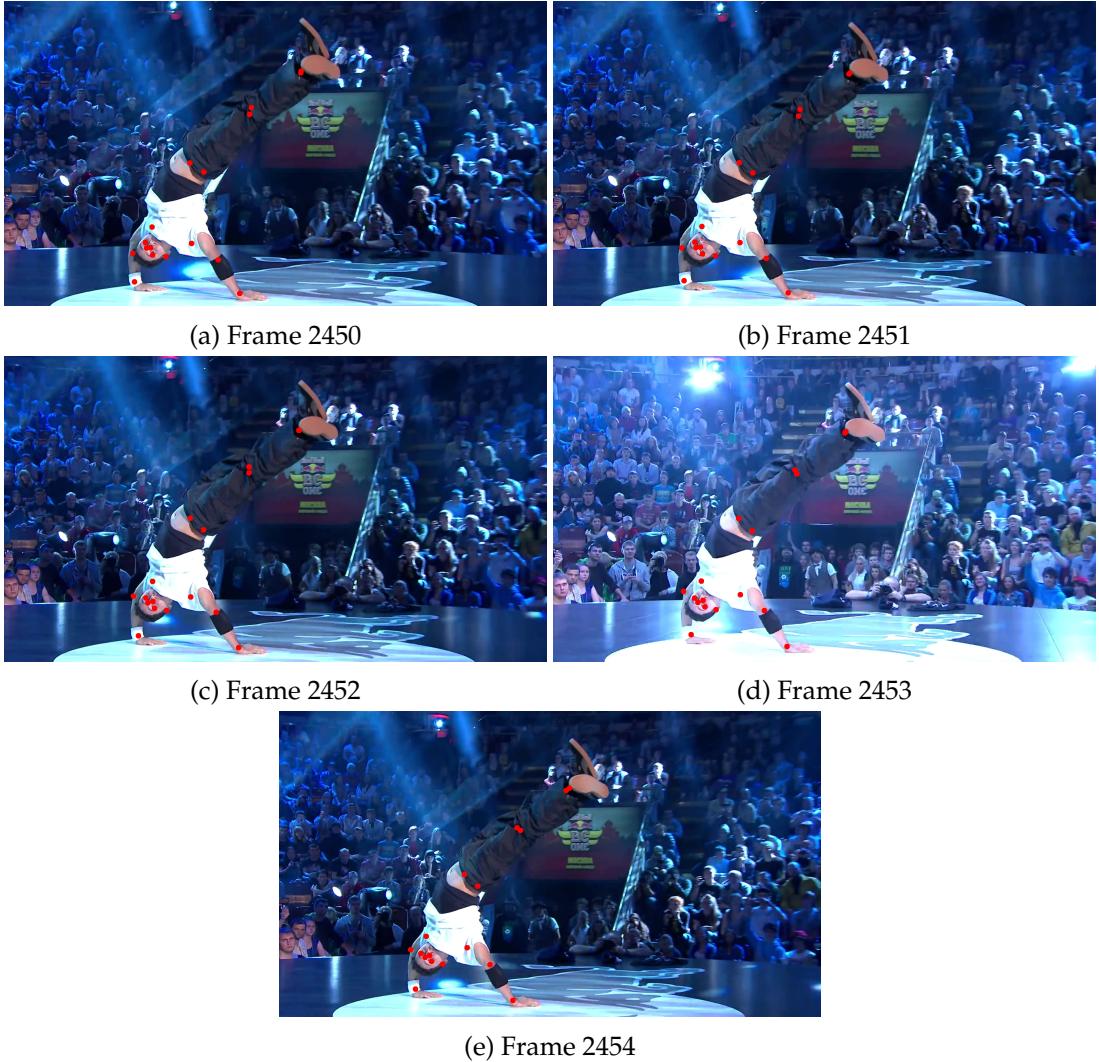


Figure 3: Example of five consecutive frames of a video from the BRACE dataset with the corresponding ground truth keypoints, where the actor holds his position for a while [1].



Figure 4: Example of five consecutive frames of a video from the BRACE dataset with the corresponding ground truth keypoints, where the actor performs a quick movement [1].

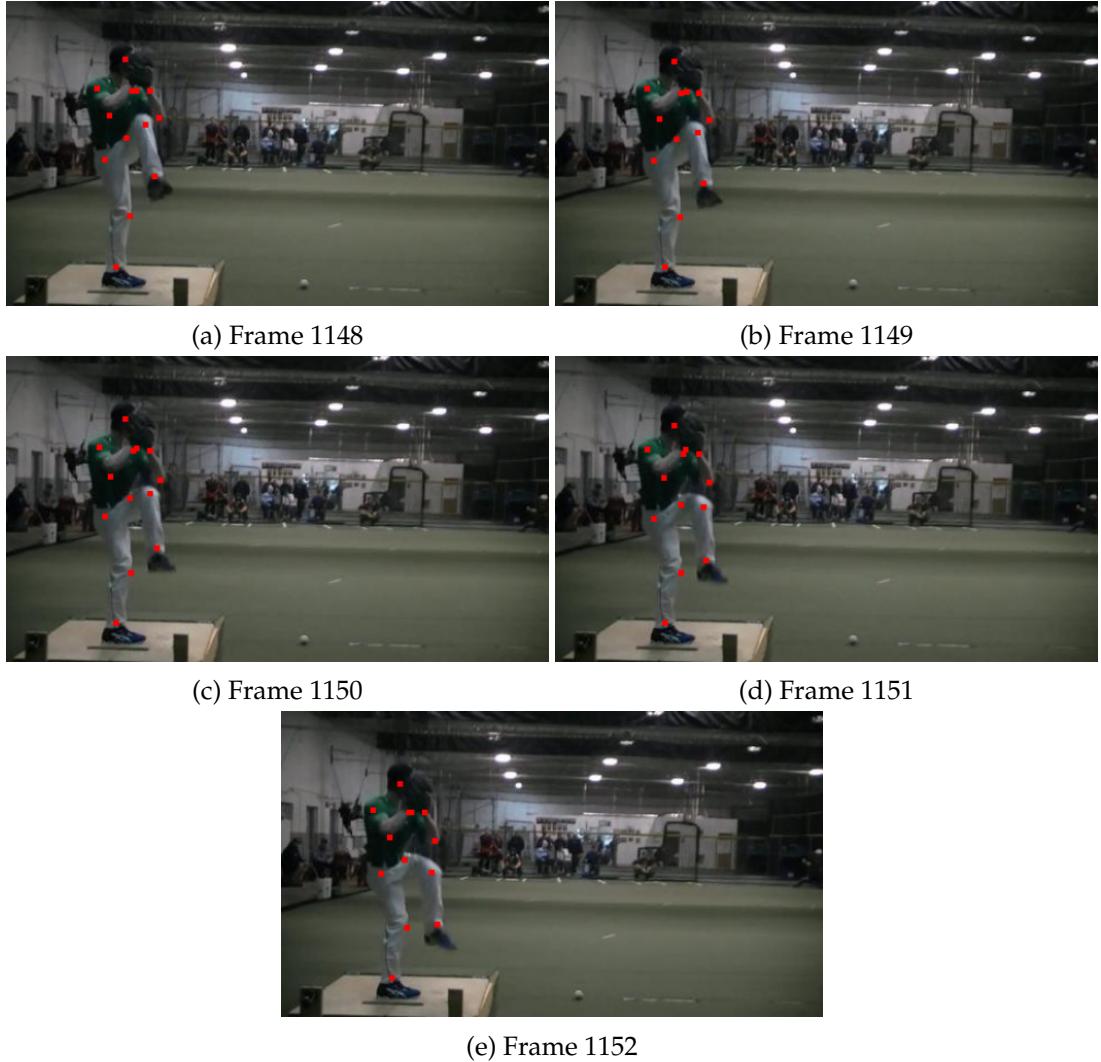


Figure 5: Example of five consecutive frames of a video from the Penn Action dataset with the corresponding ground truth keypoints [2].

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 handling of the non-visible keypoints is very inconsistent, as the position of some of the keypoints have been estimated, whereas for others it have been placed near one of the corners of the frame.

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 keypoints that have `baseball_pitch`, `bench_press` or `sit_ups` as their corresponding action-type [2].

In total, we use 307 video sequences from the Penn Action dataset, consisting of a total of 26,036 frames. Figure 5 illustrates five consecutive frames with its ground truth annotations for one of these video sequences.