

# Supplementary material: Pose Symmetric Network for Human Pose Regression

## A Code and Test Cases

In the supplemental materials, we have included Pytorch implementation of the proposed layers. Each layer also comes with unit-tests validating the chirality-equivariance. Please read the README.md for directory structures, usage and required dependencies. There is also a Jupyter notebook and it's HTML output visualizing the concepts introduced in the paper.

## B Additional Description for Equivariant Layers

### B.1 Equivariant fully connected layers

Recall, we achieve equivariance through parameter sharing and odd symmetry.

A fully connected layer performs the mapping  $\mathbf{y} = f_{\text{FC}}(\mathbf{x}; W, b) := W\mathbf{x} + b$ . Recall, we achieve equivariance through parameter sharing and odd symmetry:

$$W = \begin{bmatrix} \begin{bmatrix} W_{1n,1n} & W_{1n,1p} \\ W_{1p,1n} & W_{1p,1p} \end{bmatrix} & \begin{bmatrix} W_{1n,rn} & W_{1n,rp} \\ W_{1p,rn} & W_{1p,rp} \end{bmatrix} & \begin{bmatrix} W_{1n,cn} & W_{1n,cp} \\ W_{1p,cn} & W_{1p,cp} \end{bmatrix} \\ \begin{bmatrix} W_{1n,rn} & -W_{1n,rp} \\ -W_{1p,rn} & W_{1p,rp} \end{bmatrix} & \begin{bmatrix} W_{1n,1n} & -W_{1n,1p} \\ -W_{1p,1n} & W_{1p,1p} \end{bmatrix} & \begin{bmatrix} W_{1n,cn} & -W_{1n,cp} \\ -W_{1p,cn} & W_{1p,cp} \end{bmatrix} \\ \begin{bmatrix} W_{cn,1n} & W_{cn,1p} \\ \mathbf{0} & W_{cp,1p} \end{bmatrix} & \begin{bmatrix} W_{cn,1n} & -W_{cn,1p} \\ \mathbf{0} & W_{cp,1p} \end{bmatrix} & \begin{bmatrix} W_{cn,cn} & \mathbf{0} \\ \mathbf{0} & W_{cp,cp} \end{bmatrix} \end{bmatrix}, b = \begin{bmatrix} b_{1n} \\ b_{1p} \\ -b_{1n} \\ -b_{1p} \\ \mathbf{0} \\ b_{cp} \end{bmatrix}$$

Here, we prove that the design is chiral-equivariant. Through multiplying out the matrices, we can show  $W\mathcal{T}(\mathbf{x}) + b = \mathcal{T}(W\mathbf{x} + b)$ , as follows:

*Proof:*

$$\mathbf{x} = [x_{1n} \ x_{1p} \ x_{rn} \ x_{rp} \ x_{cn} \ x_{cp}]^T \text{ then } \mathcal{T}(\mathbf{x}) = [-x_{rn} \ x_{rp} \ -x_{1n} \ x_{1p} \ -x_{cn} \ x_{cp}]^T$$

With linear algebra,

$$\begin{aligned} W\mathbf{x} + b &= \begin{bmatrix} W_{1n,1n}(x_{1n}) + W_{1n,1p}(x_{1p}) + W_{1n,rn}(x_{rn}) + W_{1n,rp}(x_{rp}) + W_{1n,cn}(x_{cn}) + W_{1n,cp}(x_{cp}) + b_{1n} \\ W_{1p,1n}(x_{1n}) + W_{1p,1p}(x_{1p}) + W_{1p,rn}(x_{rn}) + W_{1p,rp}(x_{rp}) + W_{1p,cn}(x_{cn}) + W_{1p,cp}(x_{cp}) + b_{1p} \\ W_{1n,rn}(x_{1n}) - W_{1n,rp}(x_{1p}) + W_{1n,1n}(x_{rn}) - W_{1n,1p}(x_{rp}) + W_{1n,cn}(x_{cn}) - W_{1n,cp}(x_{cp}) - b_{1n} \\ -W_{1p,rn}(x_{1n}) + W_{1p,rp}(x_{1p}) - W_{1p,1n}(x_{rn}) + W_{1p,1p}(x_{rp}) - W_{1p,cn}(x_{cn}) + W_{1p,cp}(x_{cp}) + b_{1p} \\ W_{cn,1n}(x_{1n}) + W_{cn,1p}(x_{1p}) + W_{cn,1n}(x_{rn}) - W_{cn,1p}(x_{rp}) + W_{cn,cn}(x_{cn}) + \mathbf{0} \cdot (x_{cp}) + \mathbf{0} \\ \mathbf{0} \cdot (x_{1n}) + W_{cp,1p}(x_{1p}) + \mathbf{0} \cdot (x_{rn}) + W_{cp,1p}(x_{rp}) + \mathbf{0} \cdot (x_{cn}) + W_{cp,cp}(x_{cp}) + b_{cp} \end{bmatrix} \\ \mathcal{T}(W\mathbf{x} + b) &= \begin{bmatrix} -W_{1n,rn}(x_{1n}) + W_{1n,rp}(x_{1p}) - W_{1n,1n}(x_{rn}) + W_{1n,1p}(x_{rp}) - W_{1n,cn}(x_{cn}) + W_{1n,cp}(x_{cp}) + b_{1n} \\ -W_{1p,rn}(x_{1n}) + W_{1p,rp}(x_{1p}) - W_{1p,1n}(x_{rn}) + W_{1p,1p}(x_{rp}) - W_{1p,cn}(x_{cn}) + W_{1p,cp}(x_{cp}) + b_{1p} \\ -W_{1n,1n}(x_{1n}) - W_{1n,1p}(x_{1p}) - W_{1n,rn}(x_{rn}) - W_{1n,rp}(x_{rp}) - W_{1n,cn}(x_{cn}) - W_{1n,cp}(x_{cp}) - b_{1n} \\ W_{1p,1n}(x_{1n}) + W_{1p,1p}(x_{1p}) + W_{1p,rn}(x_{rn}) + W_{1p,rp}(x_{rp}) + W_{1p,cn}(x_{cn}) + W_{1p,cp}(x_{cp}) + b_{1p} \\ -W_{cn,1n}(x_{1n}) - W_{cn,1p}(x_{1p}) - W_{cn,1n}(x_{rn}) + W_{cn,1p}(x_{rp}) - W_{cn,cn}(x_{cn}) - \mathbf{0} \cdot (x_{cp}) - \mathbf{0} \\ \mathbf{0} \cdot (x_{1n}) + W_{cp,1p}(x_{1p}) + \mathbf{0} \cdot (x_{rn}) + W_{cp,1p}(x_{rp}) + \mathbf{0} \cdot (x_{cn}) + W_{cp,cp}(x_{cp}) + b_{cp} \end{bmatrix} \\ W\mathcal{T}(\mathbf{x}) + b &= \begin{bmatrix} W_{1n,1n}(-x_{rn}) + W_{1n,1p}(x_{rp}) + W_{1n,rn}(-x_{1n}) + W_{1n,rp}(x_{1p}) + W_{1n,cn}(-x_{cn}) + W_{1n,cp}(x_{cp}) + b_{1n} \\ W_{1p,1n}(-x_{rn}) + W_{1p,1p}(x_{rp}) + W_{1p,rn}(-x_{1n}) + W_{1p,rp}(x_{1p}) + W_{1p,cn}(-x_{cn}) + W_{1p,cp}(x_{cp}) + b_{1p} \\ W_{1n,rn}(-x_{rn}) - W_{1n,rp}(x_{rp}) + W_{1n,1n}(-x_{1n}) - W_{1n,1p}(x_{1p}) + W_{1n,cn}(-x_{cn}) - W_{1n,cp}(x_{cp}) - b_{1n} \\ -W_{1p,rn}(-x_{rn}) + W_{1p,rp}(x_{rp}) - W_{1p,1n}(-x_{1n}) + W_{1p,1p}(x_{1p}) - W_{1p,cn}(-x_{cn}) + W_{1p,cp}(x_{cp}) + b_{1p} \\ W_{cn,1n}(-x_{rn}) + W_{cn,1p}(x_{rp}) + W_{cn,1n}(-x_{1n}) - W_{cn,1p}(x_{1p}) + W_{cn,cn}(-x_{cn}) + \mathbf{0} \cdot (x_{cp}) + \mathbf{0} \\ \mathbf{0} \cdot (-x_{rn}) + W_{cp,1p}(x_{rp}) + \mathbf{0} \cdot (-x_{1n}) + W_{cp,1p}(x_{1p}) + \mathbf{0} \cdot (-x_{cn}) + W_{cp,cp}(x_{cp}) + b_{cp} \end{bmatrix} \end{aligned}$$

observe that  $W\mathcal{T}(\mathbf{x}) + b = \mathcal{T}(W\mathbf{x} + b)$ , which proves the claim.  $\square$

### B.2 Equivariant 1D convolution layers

**1D convolution layers [48, 24].** Pose symmetric 1D convolution layers can be based on fully connected layers. A 1D convolution is a fully connected layer with shared parameters across the time

dimension, *i.e.*, at each time step the computation is the sum of fully connected layers over a window:

$$\mathbf{y}_t = \sum_{\tau} W_{\tau} \mathbf{x}_{t-\tau} + b = \sum_{\tau} f_{\text{FC}}(\mathbf{x}_{t-\tau}; W_{\tau}, b).$$

Consequently, we enforce equivariance at each time step by employing the symmetry pattern of fully connected layers at each time slice.

$$W_{\tau} = \begin{bmatrix} \begin{bmatrix} W_{1n,1n,\tau} & W_{1n,1p,\tau} \\ W_{1p,1n,\tau} & W_{1p,1p,\tau} \end{bmatrix} & \begin{bmatrix} W_{1n,rn,\tau} & W_{1n,rp,\tau} \\ W_{1p,rn,\tau} & W_{1p,rp,\tau} \end{bmatrix} & \begin{bmatrix} W_{1n,cn,\tau} & W_{1n,cp,\tau} \\ W_{1p,cn,\tau} & W_{1p,cp,\tau} \end{bmatrix} \\ \begin{bmatrix} W_{1n,rn,\tau} & -W_{1n,rp,\tau} \\ -W_{1p,rn,\tau} & W_{1p,rp,\tau} \end{bmatrix} & \begin{bmatrix} W_{1n,1n,\tau} & -W_{1n,1p,\tau} \\ -W_{1p,1n,\tau} & W_{1p,1p,\tau} \end{bmatrix} & \begin{bmatrix} W_{1n,cn,\tau} & -W_{1n,cp,\tau} \\ -W_{1p,cn,\tau} & W_{1p,cp,\tau} \end{bmatrix} \\ \begin{bmatrix} W_{cn,1n,\tau} & W_{cn,1p,\tau} \\ \mathbf{0} & W_{cp,1p,\tau} \end{bmatrix} & \begin{bmatrix} W_{cn,1n,\tau} & -W_{cn,1p,\tau} \\ \mathbf{0} & W_{cp,1p,\tau} \end{bmatrix} & \begin{bmatrix} W_{cn,cn,\tau} & \mathbf{0} \\ \mathbf{0} & W_{cp,cp,\tau} \end{bmatrix} \end{bmatrix},$$

for all  $\tau$ . The bias of a 1D convolution is identical to that of a fully connected layer, *i.e.*, the same bias is added for each time step. Hence the same parameter sharing is used.

### B.3 Equivariant LSTM and GRU layers

LSTM and GRU modules which satisfy chirality can be obtained from fully connected layers. However, naïvely setting all matrix multiplies within an LSTM to satisfy the equivariance property will not lead to an equivariant LSTM because gates are elementwise *multiplied* with the cell state. If both gate and cell preserve the negation then the product will not. Therefore, we change the weight sharing scheme for the gates. We set  $D_n^{\text{out}}$  for the gates to be the empty set, *i.e.*, the gates will be invariant to negation at the input,  $T_{\text{neg}}^{\text{in}}$ , but still equivariant to the switch operation,  $T_{\text{swi}}^{\text{in}}$ . With this setup, the product of the gates and the cell's output will preserve the sign, as the gates are invariant to negation and passed through a Sigmoid to be within the range of  $(0, 1)$ . GRU modules are modified in the same manner.

More formally, the computation in an LSTM module are as follows:

$$\begin{aligned} i_t &= \sigma(W^{\text{ii}}x_t + b^{\text{ii}} + W^{\text{hi}}h_{(t-1)} + b^{\text{hi}}) && \text{(Input Gate)} \\ o_t &= \sigma(W^{\text{io}}x_t + b^{\text{io}} + W^{\text{ho}}h_{(t-1)} + b^{\text{ho}}) && \text{(Output Gate)} \\ f_t &= \sigma(W^{\text{if}}x_t + b^{\text{if}} + W^{\text{hf}}h_{(t-1)} + b^{\text{hf}}) && \text{(Forget Gate)} \\ g_t &= \tanh(W^{\text{ig}}x_t + b^{\text{ig}} + W^{\text{hg}}h_{(t-1)} + b^{\text{hg}}) && \text{(Cell State)} \\ c_t &= f_t \cdot c_{(t-1)} + i_t \cdot g_t \\ h_t &= o_t \cdot \tanh(c_t) && \text{(Recurrent State)} \end{aligned},$$

where  $\sigma$  denotes an element-wise sigmoid non-linearity.

Observe that the LSTM operations consist of fully connected layers. For the cell state's parameters, *e.g.*,  $W^{\text{ig}}$ ,  $W^{\text{hg}}$ ,  $b^{\text{ig}}$ ,  $b^{\text{hg}}$ , we follow the weight sharing scheme discussed for fully connected layers.

Due to multiplication in the cell state, we redesigned the parameter sharing for the input, output and forget gate, to be invariant to  $T_{\text{neg}}^{\text{in}}$ , by setting  $D_n^{\text{out}}$  to be the empty set: no negation is needed for all dimension. This results in the following parameter sharing scheme for the parameters  $W^{\text{ii}}$ ,  $b^{\text{ii}}$ ,  $W^{\text{hi}}$ ,  $b^{\text{hi}}$ ,  $W^{\text{io}}$ ,  $b^{\text{io}}$ ,  $W^{\text{ho}}$ ,  $b^{\text{ho}}$ ,  $W^{\text{if}}$ ,  $b^{\text{if}}$ ,  $W^{\text{hf}}$ ,  $b^{\text{hf}}$ :

$$W = \begin{bmatrix} \begin{bmatrix} W_{1p,1n} & W_{1p,1p} \\ -W_{1p,rn} & W_{1p,rp} \end{bmatrix} & \begin{bmatrix} W_{1p,rn} & W_{1p,rp} \\ -W_{1p,1n} & W_{1p,1p} \end{bmatrix} & \begin{bmatrix} W_{1p,cn} & W_{1p,cp} \\ -W_{1p,cn} & W_{1p,cp} \end{bmatrix} \\ \begin{bmatrix} \mathbf{0} & W_{cp,1p} \end{bmatrix} & \begin{bmatrix} \mathbf{0} & W_{cp,1p} \end{bmatrix} & \begin{bmatrix} \mathbf{0} & W_{cp,cp} \end{bmatrix} \end{bmatrix}, b = \begin{bmatrix} b_{1p} \\ b_{1p} \\ b_{cp} \end{bmatrix}.$$

This LSTM is chirality equivariant, as the computation of the cell state is equivariant. Other computations are linear combinations of chirality equivariant operations, which remains equivariant. We note that the chirality equivariant GRU module is modified by following the same sharing scheme for the gates.

### B.4 Equivariant batch-norm layers

A batch normalization layer performs an element-wise standardization, followed by an element-wise affine layer (with learnable parameters  $\gamma$  and  $\beta$ ):

$$\mathbf{y} = f_{\text{BN}}(\mathbf{x}) := \gamma \cdot \frac{\mathbf{x} - \mu}{\sqrt{\sigma^2 + \epsilon}} + \beta.$$

| App.                      | Walk        |             |      | Jog         |             |             | Box         |             |             | Avg.<br>-   |
|---------------------------|-------------|-------------|------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
|                           | S1          | S2          | S3   | S1          | S2          | S3          | S1          | S2          | S3          |             |
| Pavlo [36]                | 17.6        | 12.5        | 37.6 | 28.1        | 19.1        | 19.2        | 29.5        | 44.0        | 43.1        | 33.3        |
| Pavlo [36] ( $\ddagger$ ) | <b>17.5</b> | 12.3        | 37.4 | <b>27.7</b> | 19.0        | 19.0        | 27.7        | 43.4        | 42.5        | 33.0        |
| Ours                      | 18.9        | <b>12.3</b> | 38.1 | 28.5        | <b>18.1</b> | <b>18.2</b> | <b>27.1</b> | <b>40.9</b> | <b>40.2</b> | <b>32.2</b> |

Table A1: Results on HumanEva-I for multi-action (MA) models reported in Protocol 1 (MPJPE), lower the better.  $\ddagger$  indicates test time augmentation.

451 *Equivariance for  $\gamma$ , and  $\beta$*  is obtained by following the principle applied to fully connected layers:  
452 we achieve equivariance via parameter sharing and odd symmetry:

$$453 \quad \gamma = \begin{bmatrix} [\gamma_{1n} & \gamma_{1p}] & [\gamma_{1n} & \gamma_{1p}] & [\gamma_{cn} & \gamma_{cp}] \end{bmatrix}^T \text{ and } \beta = \begin{bmatrix} [\beta_{1n} & \beta_{1p}] & [-\beta_{1n} & \beta_{1p}] & [0 & \beta_{cp}] \end{bmatrix}^T.$$

454 *Equivariance for  $\mu$ , and  $\sigma$*  is obtained by computing the mean and standard deviation on the ‘‘aug-  
455 mented batch’’ and by keeping track of its running average. Formally, given a batch  $\mathcal{B}$  of data,

$$456 \quad \mu = \frac{1}{2|\mathcal{B}|} \sum_{\mathbf{x} \in \mathcal{B}} \mathbf{x} + \mathcal{T}^{\text{in}}(\mathbf{x}), \quad \sigma = \sqrt{\frac{\sum_{\mathbf{x} \in \mathcal{B}} (\mathbf{x} - \mu)^2 + (\mathcal{T}^{\text{in}}(\mathbf{x}) - \mu)^2}{2|\mathcal{B}|}}.$$

## 457 B.5 Dropout.

At test time, dropout scales the input by  $p$ , where  $p$  is the dropout probability. The equivariance property is satisfied because of the associativity property of a scalar multiplication. The input and output dimension and symmetry of a dropout layer are identical. Therefore,  $\mathcal{T}^{\text{out}}$  and  $\mathcal{T}^{\text{in}}$  are identical. From the definition:

$$\mathcal{T}^{\text{out}}(p \cdot \mathbf{x}) = \mathcal{T}^{\text{in}}(p \cdot \mathbf{x}) = T_{\text{neg}}^{\text{in}} T_{\text{swi}}^{\text{in}}(p \cdot \mathbf{x}) = p \cdot (T_{\text{neg}}^{\text{in}} T_{\text{swi}}^{\text{in}} \mathbf{x}) = p \cdot (\mathcal{T}^{\text{in}}(\mathbf{x})) \quad \forall \mathbf{x} \in \mathbb{R}^{|J^{\text{in}}||D^{\text{in}}|}.$$

458 Hence, a dropout layer naturally satisfies the equivariance property. At training-time, we do not  
459 enforce equivariance for the dropped units, *i.e.*, we do not jointly drop symmetric units as we found  
460 this to prevent overfitting. This is likely application dependent.

## 461 C Additional Results

### 462 C.1 3D pose estimation

463 In Tab. A1, we report the HumanEva-I for multi-action models evaluated on Protocol 1 (MPJPE).  
464 Our approach have benefits the most from the Boxing action while maintaing the performance on  
465 other actions. We also provide qualitative evaluation in Fig. A1 and Fig. A2. We observe that our  
466 model successfully estimates 3D poses from 2D key-points. We have also attached animations in the  
467 supplemental.

### 468 C.2 Skeleton based action recognition

469 In Fig. A3, we show the visualization of the input skeleton sequences computed by OpenPose [2] and  
470 the predicted action class by our chiral invariant skeleton based action recognition model.

## 471 D Implementation Details

### 472 D.1 3D pose estimation

473 **Implementation details.** Our model follows the temporal convolutional architecture proposed  
474 by Pavlo et al. [36], and replaced all layers with their chiral versions; code for the layers are attached  
475 in the supplemental as well. We also changed ReLU to tanh to achieve chiral equivariance. For the  
476 temporal models, we follow their 4 blocks design which has the receptive field of 243. For the single  
477 frame model, we follow their 3 blocks design. These models all contains 1020 hidden dimensions so  
478 it is a factor the number of joints, 17, this is slightly smaller than the 1024 used in [36]. We also use  
479 their data processing and batching stragety as described in Section 5 and Appendix A.5 of [36]. For

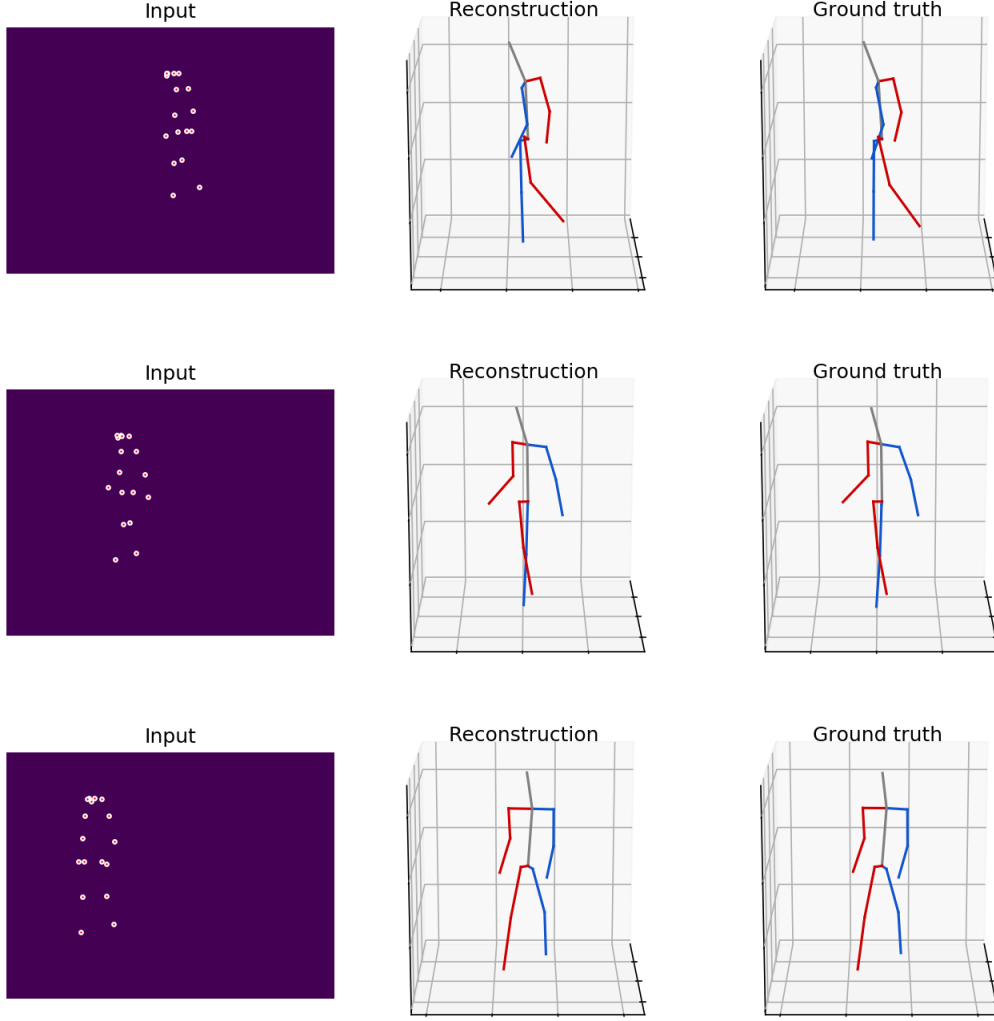


Figure A1: Qualitative visualization of 2D to 3D pose estimation for the action “Walking” on HumanEva-I dataset.

480 training the model, we utilized the Adam optimizer with  $\beta_1=0.9$  and  $\beta_2=0.9999$ . We decay the  
481 batch-normalizations’ momentum as suggested in [36]. Other details follows the publicly available  
482 implementation by Pavlo et al. [36]. We enforced chiral equivariance by choosing the  $|D_n^{\text{out}}|$  to be  $\frac{1}{3}$   
483 of the hidden dimension. The  $|D_n^{\text{in}}|$  for the input layer is 17 and the  $|D_n|^{\text{out}}$  for the output layer is  
484 17, as one for each joint.

## 485 D.2 2D pose forecasting

486 **Implementation details.** The non-chiral equivariant baseline is a seq2seq model consisting of an  
487 encoder and decoder, which are stacked-LSTMs with hidden size of 1040 and 2 stacked layers. We  
488 trained using teacher forcing with the Adam optimizer. The batch-size is 256, and we trained for 30  
489 epochs. Dropout is applied to the LSTMs’ hidden layer with drop probability of 0.5. Following prior  
490 works, we use max norm gradient clipping of 5, a learning rate of 0.005 with a decay of 0.95 every 2  
491 epochs. The data processing and evaluation setting follows [5]. Other details follows the publicly  
492 available implementation by Chiu et al. [5]. We enforced chiral equivariance by choosing the  $|D_n^{\text{out}}|$   
493 to be  $\frac{1}{2}$  of the hidden dimension, as the output is two dimensional per joint.

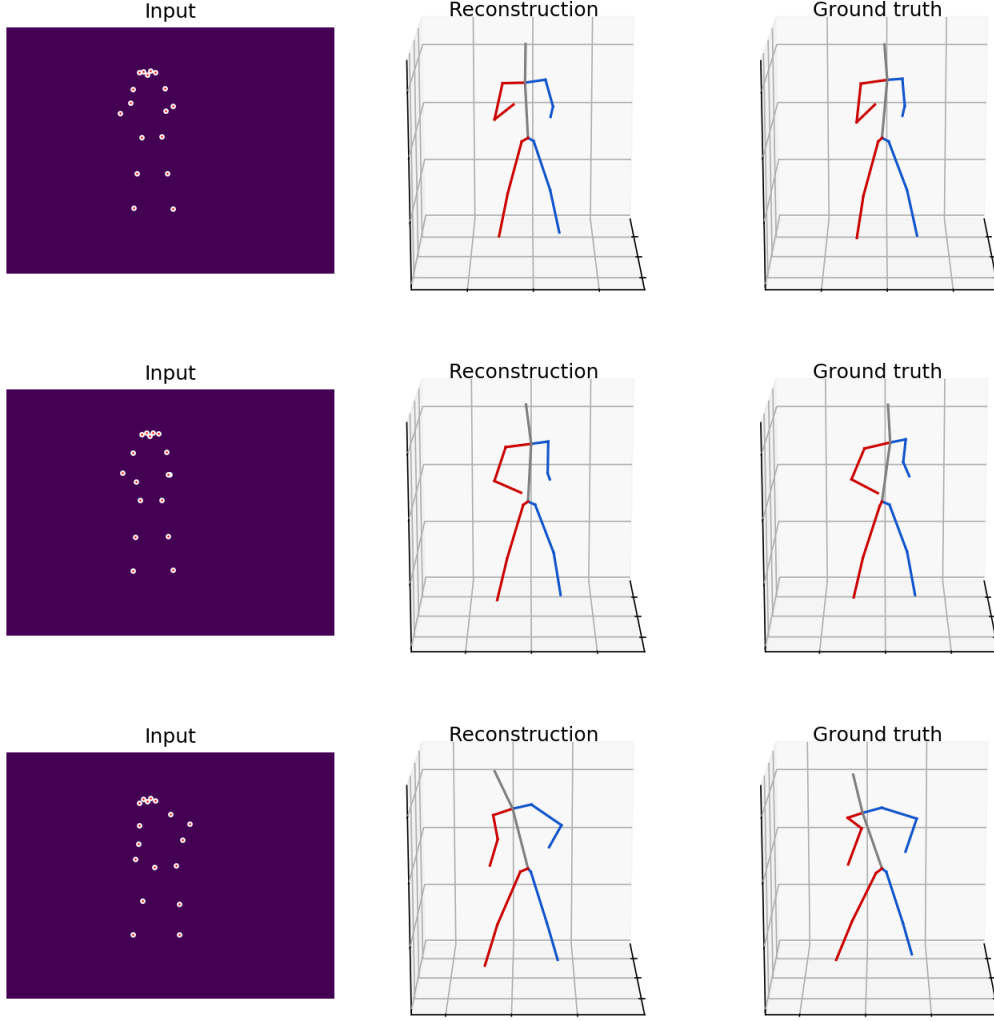


Figure A2: Qualitative visualization of 2D to 3D pose estimation for the action “Boxing” on HumanEva-I dataset.

### 494 D.3 Skeleton-based action recognition

495 **Implementation details.** The non-chiral version of the model, Ours-Conv, follows Temporal-  
496 Conv [21] while we modified the model to have not only temporal convolution but also spatial  
497 convolution. There are ten spatial-temporal convolution blocks and each block we first perform  
498 spatial convolution and then temporal convolution. The temporal convolution considers the intra-  
499 frame information while the spatial convolution considers the inter-frame information. For the  
500 recognition task, we need chiral invariance, *i.e.*, a chiral pair should be classified as the same action  
501 class. To this end, we use a chiral invariance layer where we let both  $J_r^{\text{out}}$ ,  $J_l^{\text{out}}$  as well as  $D_n^{\text{out}}$   
502 to be empty sets, which means there are no left and right joints but only center joints and there is  
503 no dimension that will be negated in the output of the layer after applying the chirality transform.  
504 Note that the chiral transformation exchange the left and right joints and negate the dimension in the  
505 index set  $D_n^{\text{out}}$ . Given  $J_r^{\text{out}}$ ,  $J_l^{\text{out}}$  and  $D_n^{\text{out}}$  are all empty, it's obvious that the output will be chiral  
506 invariance. For the chiral invariance model, Ours-Conv-Chiral, we replace the all the non-symmetric  
507 layers before the chiral invariance layer with their corresponding chiral equivariance version. All the  
508 layers after the chiral invariance layer remains the same as in the Ours-Conv model. Similar to [21],  
509 there are in total 10 convolution blocks in Ours-Conv and we put the chiral invariance layer at the  
510 fourth layer. Also, we gradually reduce the ratio of the dimension to be negated ( $|D_n^{\text{out}}|/|D^{\text{out}}|$ ) from

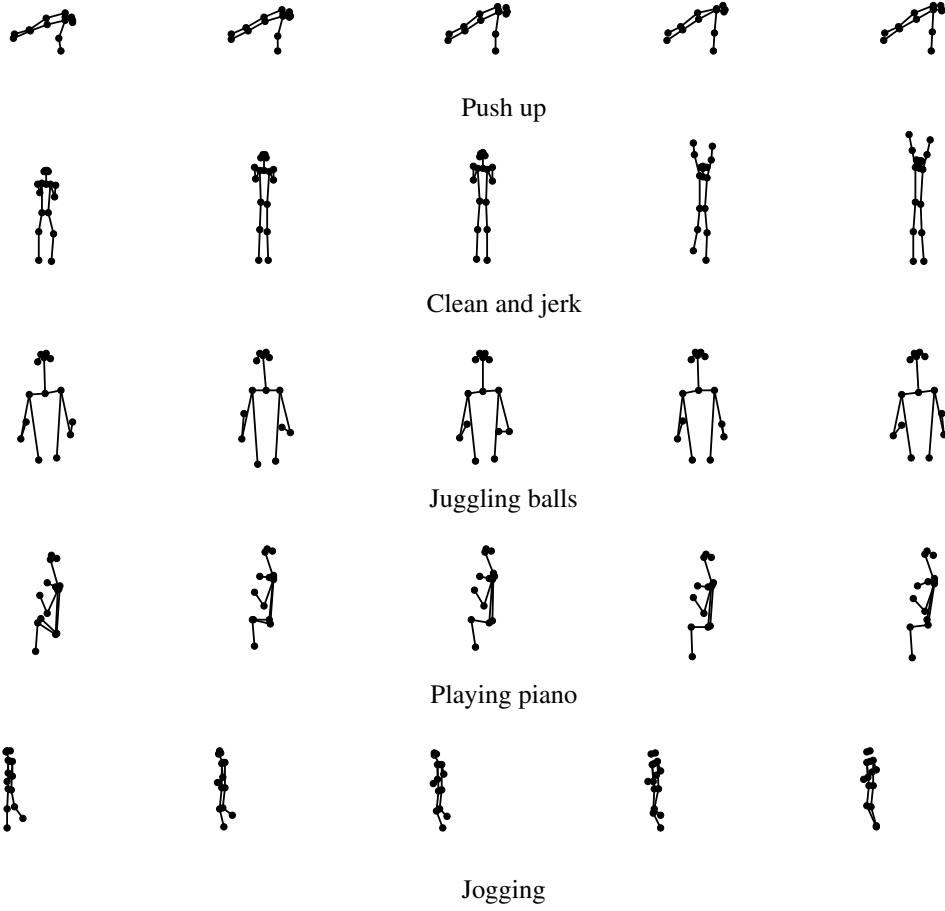


Figure A3: Visualization of the input skeleton sequences and the corresponding predicted action classes of our method on the Kinetics-400 dataset [20].

511  $\frac{1}{3}$  to  $\frac{1}{6}$  at the first layer, from  $\frac{1}{6}$  to  $\frac{1}{12}$  at the second layer and from  $\frac{1}{12}$  to 0 at the third layer. We use  
512 the SGD optimizer with a momentum of 0.9 as in [51] with a batch size of 256. We train the model  
513 for 90 epochs.