# Supplementary Material for ProtoRes: Proto-Residual Network for Pose Authoring via Learned Inverse Kinematics
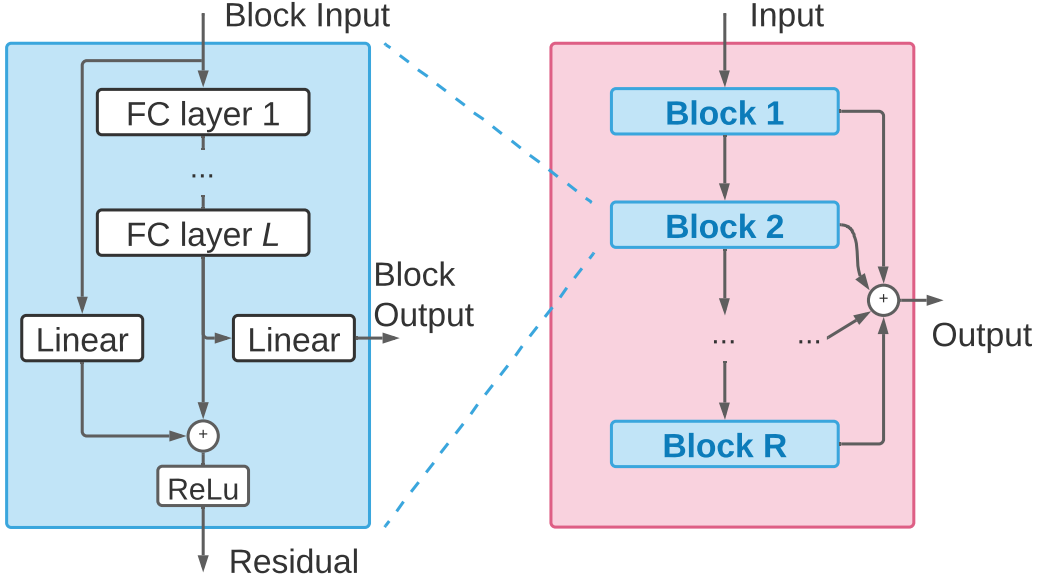
## Table of Contents

Figure 4: Block diagram of the fully-connected residual (FCR) decoder architecture. Left: the diagram of one residual block of the FCR decoder. Note that the basic residual block of the encoder architecture is exactly the same. Right: residual blocks connected in the FCR architecture.

## A    ARCHITECTURE DETAILS

**Residual block** depicted in Fig. 4 (left) is used as the basic building block of the ProtoRes encoder and decoder.

**Decoders** The block diagram of the global position and the inverse kinematics decoders used in the main architecture (see Fig. 2) is presented in Fig. 4. The architecture has fully connected residual topology consisting of multiple fully connected blocks connected using residual connections. Each block has residual and forward outputs. The forward output contributes to the final output of the decoder. The residual connection sums the hidden state of the block with the linear projection of the input and applies a ReLU non-linearity.

In the main text we use a convention that the number of layers and blocks in the encoder, as well as in GPD and IKD decoders is the same and is given by $L$ and $R$ respectively. Obviously, using a different number of layers and residual blocks in each of the blocks might be more optimal.

**Forward Kinematics** pass is applied to the output of the IKD, transforming local joint rotations and global root position into the global joint rotations and positions using skeleton kinematic equations. The FK pass relies on the offset vector $\mathbf{o}_j = [o_{x,j}, o_{y,j}, o_{z,j}]^\mathsf{T}$ and the rotation matrix $\mathbf{R}_j$ for each joint $j$. The offset vector is a fixed non-learnable vector representing bone length constraint for joint $j$. It provides the displacement of this joint with respect to its parent joint when joint $j$ rotation is zero. $\mathbf{R}_j$ can be naïvely represented using local Euler rotation angles $\alpha_j, \beta_j, \gamma_j$:

$$\mathbf{o}_j = \begin{bmatrix} o_{x,j} \\ o_{y,j} \\ o_{z,j} \end{bmatrix}; \quad \mathbf{R}_j = \begin{bmatrix} \cos\alpha_j & -\sin\alpha_j & 0 \\ \sin\alpha_j & \cos\alpha_j & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \cos\beta_j & 0 & \sin\beta_j \\ 0 & 1 & 0 \\ -\sin\beta_j & 0 & \cos\beta_j \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos\gamma_j & -\sin\gamma_j \\ 0 & \sin\gamma_j & \cos\gamma_j \end{bmatrix}. \quad (13)$$

However, we use a more robust representation proposed by (Zhou et al., 2019), relying on vector norm $\overrightarrow{\mathbf{u}} \equiv \mathbf{u}/\|\mathbf{u}\|_2$ and vector cross product $\mathbf{u} \times \mathbf{v} = \|\mathbf{u}\|\|\mathbf{v}\|\cos(\gamma)\overrightarrow{\mathbf{n}}$ ($\gamma$ is the angle between $\mathbf{u}$ and $\mathbf{v}$ in the plane containing them and $\overrightarrow{\mathbf{n}}$ is the normal to the plane):

$$\widehat{\mathbf{r}}_{j,x} = \overrightarrow{\widehat{\mathbf{f}}_{R,j}[1:3]}, \quad \widehat{\mathbf{r}}_{j,z} = \overrightarrow{\widehat{\mathbf{r}}_{j,x} \times \widehat{\mathbf{f}}_{R,j}[4:6]}, \quad \widehat{\mathbf{r}}_{j,y} = \widehat{\mathbf{r}}_{j,z} \times \widehat{\mathbf{r}}_{j,x}, \quad \widehat{\mathbf{R}}_j = [\widehat{\mathbf{r}}_{j,x}\, \widehat{\mathbf{r}}_{j,y}\, \widehat{\mathbf{r}}_{j,z}]. \quad (14)$$

Provided with the local offset vectors and rotation matrices of all joints, the global rigid transform of any joint $j$ is predicted following the tree recursion from the parent joint $p(j)$ of joint $j$:

$$\widehat{\mathbf{G}}_j = \widehat{\mathbf{G}}_{p(j)} \begin{bmatrix} \widehat{\mathbf{R}}_j & \mathbf{o}_j \\ \mathbf{0} & 1 \end{bmatrix}. \tag{15}$$

The global transform matrix $\widehat{\mathbf{G}}_j$ of joint $j$ contains its global rotation matrix, $\widehat{\mathbf{G}}_j^{13} \equiv \widehat{\mathbf{G}}_j[1:3, 1:3]$, and its 3D global position, $\widehat{\mathbf{g}}_j = \widehat{\mathbf{G}}_j[1:3, 4]$.

## B   EFFECTOR NOISE MODEL

This section describes the details of the NOISEMODEL that is used in Algorithm 1 to corrupt model effector input $\mathbf{x}[i,:]$ based on appropriate noise level $\sigma(\Lambda_i)$.

### B.1   POSITION EFFECTOR NOISE MODEL

If effector type is positional ($T_i = 0$), *i.e.* effector $i$ is a coordinate in 3D space, typically corresponding to the desired position of joint $I_i$ in 3D space, we employ Gaussian white noise model:

$$\mathbf{x}[i, 1:3] = \mathbf{g}_{I_i} + \sigma(\Lambda_i)\varepsilon_i; \quad \mathbf{x}[i, 4:6] = 0. \tag{16}$$

Here $\mathbf{x}[i,:]$ is the $i$-th model input, $\mathbf{g}_{I_i}$ is the ground truth location of joint $I_i$, $\sigma(\Lambda_i)$ is the noise standard deviation computed based on eq. (20) and $\varepsilon_i$ is a 3D vector sampled from the zero-mean Normal distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

### B.2   ROTATION EFFECTOR NOISE MODEL

If effector type is angular ($T_i = 1$), *i.e.* effector $i$ is a 6DoF rotation matrix representation, we employ random rotation model that is implemented in the following stages. First, suppose $\mathbf{f}_{I_i}$ is the ground truth 6DoF representation of the global rotation of joint $I_i$ corresponding to effector $i$. We transform it to the rotation matrix representation $\mathbf{G}_{I_i}^{13}$ using equation (14). Second, we generate the random 3D Euler angles vector $\varepsilon_i$ from the zero-mean Gaussian distribution $\mathcal{N}(\mathbf{0}, \sigma(\Lambda_i)\mathbf{I})^1$ and convert it to the random rotation matrix $\Psi_i$ using eq. (13):

$$\Psi_i = \begin{bmatrix} \cos\varepsilon_i[1] & -\sin\varepsilon_i[1] & 0 \\ \sin\varepsilon_i[1] & \cos\varepsilon_i[1] & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \cos\varepsilon_i[2] & 0 & \sin\varepsilon_i[2] \\ 0 & 1 & 0 \\ -\sin\varepsilon_i[2] & 0 & \cos\varepsilon_i[2] \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos\varepsilon_i[3] & -\sin\varepsilon_i[3] \\ 0 & \sin\varepsilon_i[3] & \cos\varepsilon_i[3] \end{bmatrix}. \tag{17}$$

Third, we apply random rotation to the ground truth matrix, $\mathbf{G}_{I_i}^{13\prime} = \Psi_i \mathbf{G}_{I_i}^{13}$. Finally, we convert the randomly perturbed rotation matrix back to the 6DoF representation:

$$\mathbf{x}[i, 1:3] = \mathbf{G}_{I_i}^{13\prime}[:, 1], \quad \mathbf{x}[i, 4:6] = \mathbf{G}_{I_i}^{13\prime}[:, 2]. \tag{18}$$

### B.3   LOOK-AT EFFECTOR NOISE MODEL

If effector type is look-at ($T_i = 2$), *i.e.* effector $i$ is a position of the target at which a given joint is supposed to look, we employ random sampling of the target point along the ray cast in the direction formed by the global rotation of a given joint.

First, we sample the local direction vector $\mathbf{d}_i$ from the zero-mean normal 3D distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$ and normalize it to unit length. Second, we sample the distance between the joint and the target object, $d_{\mathbf{t}}$, from the normal distribution $\mathcal{N}(0, 5)$ folded over at 0 by taking the absolute value. The location of the target object is then determined as $\mathbf{t}_i = \mathbf{g}_{I_i} + d_{\mathbf{t}}\mathbf{d}_i + \sigma(\Lambda_i)\varepsilon_i$. Finally, the output is constructed as follows:

$$\mathbf{x}[i, 1:3] = \mathbf{t}_i, \quad \mathbf{x}[i, 4:6] = \mathbf{d}_i. \tag{19}$$

As previously, $\varepsilon_i$ is a 3D vector sampled from the zero-mean Normal distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

---

**Algorithm 1** Loss calculation for a single item in the training batch of ProtoRes.

---

**Require:** $\mathbf{R}_j, \mathbf{G}_j; N \sim \text{UNIFORM}[3, 16]$  ▷ Ground truth for all joints $j \in [0, J)$; number of effectors
**Ensure: x**  ▷ Sample inputs
  $I_1, \ldots, I_N \leftarrow \text{MULTINOMIAL}(\{0, \ldots, J-1\}, N)$  ▷ Effector IDs
  $T_1, \ldots, T_N \leftarrow \text{MULTINOMIAL}(\{0, 1, 2\}, N)$  ▷ Effector type
  **for** $i$ in $1 \ldots N$ **do**
    $\Lambda_i \leftarrow \text{UNIFORM}[0, 1]$  ▷ Effector tolerance
    $\sigma(\Lambda_i); W(\Lambda_i) \leftarrow \sigma_M \Lambda_i^{\eta}; \min(W_M, 1/\sigma(\Lambda_i))$  ▷ Effector noise std and weight
    $\mathbf{x}[i, :] \leftarrow \text{NOISEMODEL}(\mathbf{G}_{I_i}, \sigma(\Lambda_i), T_i)$  ▷ Generate noisy effector
  **end for**
**Predict:** $\widetilde{\mathbf{f}}_{R,j}, \widehat{\mathbf{R}}_j, \widehat{\mathbf{G}}_j$   $\forall j$ based on $\mathbf{x}$
  $\mathcal{L}_{gpd-L2}^{rnd} \leftarrow \frac{1}{\sum_{i=1}^{N} \mathbb{1}_{T_i=0} W(\Lambda_i)} \sum_{i=1}^{N} \mathbb{1}_{T_i=0} W(\Lambda_i) \text{MSE}(\mathbf{g}_{I_i}, \widetilde{\mathbf{f}}_{R, I_i})$  ▷ Randomized GPD position loss
  $\mathcal{L}_{ikd-L2}^{rnd} \leftarrow \frac{1}{\sum_{i=1}^{N} \mathbb{1}_{T_i=0} W(\Lambda_i)} \sum_{i=1}^{N} \mathbb{1}_{T_i=0} W(\Lambda_i) \text{MSE}(\mathbf{g}_{I_i}, \widehat{\mathbf{g}}_{I_i})$  ▷ Randomized IKD position loss
  $\mathcal{L}_{gpd-L2}^{det} \leftarrow \sum_{j=1}^{J} \text{MSE}(\mathbf{g}_j, \widetilde{\mathbf{f}}_{R,j})$  ▷ Deterministic GPD position loss
  $\mathcal{L}_{ikd-L2}^{det} \leftarrow \sum_{j=1}^{J} \text{MSE}(\mathbf{g}_j, \widehat{\mathbf{g}}_j)$  ▷ Deterministic IKD position loss
  $\mathcal{L}_{loc-geo}^{det} \leftarrow \sum_{j=1}^{J} \text{GEO}(\mathbf{R}_j, \widehat{\mathbf{R}}_j)$  ▷ Deterministic local rotation loss
  $\mathcal{L}_{glob-geo}^{rnd} \leftarrow \frac{1}{\sum_{i=1}^{N} \mathbb{1}_{T_i=1} W(\Lambda_i)} \sum_{i=1}^{N} \mathbb{1}_{T_i=1} W(\Lambda_i) \text{GEO}(\mathbf{G}_{I_i}^{13}, \widehat{\mathbf{G}}_{I_i}^{13})$  ▷ Randomized global rotation loss
  $\mathcal{L}_{lat}^{det} \leftarrow \frac{1}{\sum_{i=1}^{N} \mathbb{1}_{T_i=2}} \sum_{i=1}^{N} \mathbb{1}_{T_i=2} \text{LAT}(\mathbf{x}[i, 1:3], \mathbf{x}[i, 4:6], \widehat{\mathbf{G}}_{I_i}^{13})$  ▷ Randomized Look-at loss
  $\mathcal{L} \leftarrow \frac{W_{pos}}{J}(\mathcal{L}_{gpd-L2}^{rnd} + \mathcal{L}_{ikd-L2}^{rnd} + \mathcal{L}_{gpd-L2}^{det} + \mathcal{L}_{ikd-L2}^{det}) + \frac{1}{J}(\mathcal{L}_{lat}^{det} + \mathcal{L}_{glob-geo}^{rnd} + \mathcal{L}_{loc-geo}^{det})$  ▷ Total loss

---

## C   TRAINING AND EVALUATION METHODOLOGY: DETAILS

The training methodology involves techniques targeting to (i) regularize model via data augmentation, (ii) learn handling of sparse inputs and (iii) effectively combine multi-task loss terms.

**Data augmentation** is based on the rotation and mirror augmentations. The former rotates the skeleton around the vertical $Y$ axis by a random angle in $[0, 2\pi]$. Rotation w.r.t. ground $XZ$ plane is not applied to avoid creating poses implausible according to the gravity direction. Mirror augmentation removes any implicit left- or right-handedness biases by flipping the skeleton w.r.t. the $YZ$ plane.

**Sparse inputs** modeling relies on effector sampling. First, the total number of effectors is sampled uniformly at random in the range [3, 16]. Given the total number of effectors, the effector IDs (one of 64 joints) and types (one of 3 types: position, rotation, or look-at) are sampled from the Multinomial without replacement. This sampling scheme produces an exponentially large number of different permutations of effector types and joints, resulting in strong regularizing effects.

**Effector tolerance and randomized loss weighting.** The motivation behind the randomized loss weighting is two-fold. First, the randomized loss weighting was originally introduced as a binary indicator to force the model to better respect constraints provided as effectors, compared to the joints predicted by the model. Afterwards, we realized that this can be made more flexible by generating a continuous variable representing the tolerance level. This variable can be provided as an input to the network and it can be exposed as a user interface feature to let the user control the degree of responsiveness of the model to different effectors. We also discovered that the latter feature only works when a noise is added to effector value and the standard deviation of the noise is appropriately synchronised with the tolerance. The noise teaches the model to disregard the effector completely if the tolerance input value corresponds to the high noise variance regime.

Second, we observed that the use of the randomized weighting improves multi-task training and generalization performance. Initially, we noticed that increasing the weight of position loss would drive the generalization on the position metric to a better spot, while the rotation metric generalization would be compromised, which is not surprising. This was especially evident when the position loss weight was increased by one or two orders of magnitude. This is a well-known phenomenon when

---

[1]Note that in the case of angles, sampling from the Tikhonov (a.k.a. circular normal or von Mises) distribution might be a better idea, but Gaussian worked well in our case.

dealing with multiple loss terms, which we informally call "fighting" between losses (related to the Pareto front, more formally). This effect can be observed when comparing two bottom rows in Table 2. Introducing the randomized loss weighting scheme we observed two things. "Fighting" disappeared, i.e. the randomly generated weights of position effectors varied in a wide range between 1e-1 and 1e5 within a batch, but the fact that some of the weight values are one or two orders of magnitude greater than the baseline position weight of 100, did not lead to the deterioration of the rotation loss. Moreover, the introduction of the randomized loss weighting positively affected the generalization on both position and rotation metrics, which can be assessed by comparing the first row of Table 2 with its bottom rows. This leads us to believe that the randomized loss weighting introduces a sinergy in the multi-task training that is not achievable by simple adjustment of static loss weights. We believe this technique could be more generally applicable to multi-task training, but a more detailed investigation of this is outside of the current scope.

We now describe the technical details behind randomized loss weighting implementation. For each sampled effector, we further uniformly sample $\Lambda \in [0,1]$ treated as effector tolerance. Given an effector tolerance $\Lambda$, noise (noise models used for different effector types are described in detail in Appendix B) with variance proportional to $\Lambda$ is added to effector data before feeding them to the neural network:

$$\sigma(\Lambda) = \sigma_M \Lambda^\eta. \tag{20}$$

We use $\eta > 10$ to shape the distribution of $\sigma$ to smaller values. Furthermore, to each effector is attached a randomized loss weight reciprocal to $\sigma(\Lambda)$, capped at $W_M$ if $\sigma(\Lambda) < 1/W_M$:

$$W(\Lambda) = \min(W_M, 1/\sigma(\Lambda)). \tag{21}$$

$\Lambda$ drives network inputs and is simultaneously used to weigh losses by $W(\Lambda)$. Thus ProtoRes learns to respect effector tolerance, leading to two positive outcomes. First, ProtoRes provides a tool allowing one to emphasize small tolerance effectors ($\Lambda \approx 0$) and relax the large tolerance ones ($\Lambda \approx 1$). Second, randomized loss weighting improves the overall accuracy in the multi-task training scenario.

The detailed procedure to compute the ProtoRes loss based on one batch item is presented in Algorithm 1 and the summary is provided below. First, we sample (i) the number of effectors and (ii) their associated type and ID. For each effector, we randomly sample the tolerance level and compute the associated noise std and loss weight. Given noise std, an appropriate noise model is applied to generate input data based on effector type as described in Appendix B. Then ProtoRes predicts draft joint positions $\widetilde{\mathbf{f}}_{R,j}$, local joint rotations $\widehat{\mathbf{R}}_j$, as well as world-space rotations and positions $\widehat{\mathbf{G}}_j$ for all joints $j \in [0,J)$. We conclude by calculating the individual deterministic and randomized loss terms, whose weighted sum is used for backpropagation.

# D DATASETS: DETAILS

## D.1 DATASETS DESCRIPTIONS

**miniMixamo** We use the following procedure to create our first dataset from the publicly available MOCAP data available from `mixamo.com`, generously provided by Adobe Inc. (2020). We download a total of 1598 clips and retarget them on our custom 64-joint skeleton using the Mixamo online tool. This skeleton definition is used in Unity to extract the global positions as well as global and local rotations of each joint at the rate of 60 frames per second (total 356,545 frames). The resulting dataset is partitioned at the clip level into train/validation/test splits (with proportion 0.8/0.1/0.1, respectively) by sampling clip IDs uniformly at random. Splitting by clip makes the evaluation framework more realistic and less prone to overfitting: frames belonging to the same clip are often similar. At last, the final splits retain only 10% of randomly sampled frames (miniMixamo has 33,676 frames total after subsampling) and all the clip identification information (clip ID, meta-data/description, character information, etc.) is discarded. This anonymization guarantees that the original sequences from `mixamo.com` cannot be reconstructed from our dataset, allowing us to release the dataset for reproducibility purposes without violating the original dataset license (Adobe Inc., 2020).

For miniMixamo our contribution is as follows. Mixamo data is not available as a single file. Therefore, anyone who wants to use the data for academic purposes needs to go through a lengthy process of downloading individual files. Importantly, this step creates additional risks for the reproducibility

of results. We have gone through this step and assembled all files in one place. Furthermore, Mixamo data cannot be redistributed, according to Adobe licensing, which is again a reproducibility risk. However, we do not need the entire dataset for benchmarking on the task we defined. Therefore, we defined a suitable subsampling and anonymization procedure that allowed us to obtain (i) a high quality reproducible benchmark dataset for our task and (ii) a legal permission from Mixamo/Adobe to redistribute this benchmark for academic research purposes. We are extremely grateful to the representatives from Mixamo and Adobe who approved it to facilitate the democratization of character animation. The entire process of creating the benchmark took us a few months of work, which we consider a significant contribution to the research community.

**miniAnonymous**    To collect our second dataset we predefine a wide range of human motion scenarios and hire a qualified MOCAP studio to record 1776 clips (967,258 total frames @60 fps). Then we create a dataset of a total of 96,666 subsampled frames following exactly the same methodology that was employed for miniMixamo.

The following action scenarios were used to collect MOCAP sequences in miniAnonymous. The detailed hierarchy of motion scenarios is presented in Table 6. First, the following locomotion types: compass crouch, compass jogs, compass runs, compass walks were collected for female and male subjects under high energy, low energy and injured scenarios. The same locomotion types were collected under neutral energy feminine and neutral energy macho scenarios. Furthermore, under neutral energy generic scenario, for both female and male subjects, we collected following action, object and environment interaction types: archery, bokken fighting, calisthetics, door interactions, fist fighting, food, handgun, hands, knife fighting, locomotion, longsword fighting, phone, place, railing interactions, rifle, seated interactions, shotgun, standings, sword, wall. Among the latter categories, locomotion and handgun had following more detailed subdivisions. Locomotion: compass crawls, compass crouch, compass jogs, compass rifle aim walks, compass rifle crawl, compass rifle crouch, compass rifle jogs, compass rifle runs, compass rifle walks, compass runs, compass walks, rifle idles, walk carry heavy backpack, walk carry heavy sack, walk carry ladder, walk dragging heavy object. Handgun: downwards, level, upwards, verticals. The motion categories existing in miniAnonymous and miniMixamo can be compared by looking at Tables 6 and 7 respectively.

The key differentiators of the datasets that we release that make them significant contributions toward AI driven artistic pose development are as follows:

- Both miniMixamo (derived from the Mixamo, which is generously provided by Adobe) and miniAnonymous are collected by professional studio contractors relying on the service of professional actors using high-end MOCAP studio equipment.

- Both datasets are clean and contain data of very high quality. For our dataset, we specifically had to go through multiple cleaning iterations to make sure all the data collection and conversion artifacts are removed. We are very grateful to our contractor for being diligent, detail oriented, and determined to provide the high quality data.

- Both datasets provide data in the industry standard skeleton format compatible with multiple existing animation rigs and therefore making it easy to experiment with the ML assisted pose authoring results in 3D development environments such as Unity. This is in contrast to CMU and AMASS datasets that are collected in heterogeneous environments using non-standard sensor placements.

- Both our datasets provide 64 joint skeletons and contain fine grain hands and feet data, unlike other publicly available datasets.

## D.2    DATASETS DESCRIPTIVE STATISTICS

The standard deviations of joint positions (in the coordinate system relative to hips joint) and joint local quaternions are presented in Tables 3 and 4.

The distributions of 30 most popular tags across the frames of the original dataset used to build miniMixamo and miniAnonymous are shown in Figure 5. It is clear that the two dataset cover some common activities such as walking and idling, for example. Additionally, there are numerous categories the two datasets emphasize separately. For example miniMixamo focuses a lot on fighting

and handling guns, whereas miniAnonymous has more neutral energy activities, provides extensive labeling of male/female poses as well as covers scenarios of handling objects and food.

The distribution of PCA of flattened rotations of all joints of miniMixamo (red) and miniAnonymous (blue) datasets is shown in Figure 6. Each point corresponds to a pose. It is clear that there are both overlapping areas corresponding to clusters with similar poses and areas with disjoint clusters, showing distinct pose configurations characteristic of the two different datasets.

Table 3: Per-joint hip-local positions standard deviations computed over two proposed datasets

| | miniMixamo | | | miniAnonymous | | |
|---|---|---|---|---|---|---|
| Joint | X | Y | Z | X | Y | Z |
| Hips | 0 | 0 | 0 | 0 | 0 | 0 |
| Spine0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Spine1 | 0.0049 | 0.0034 | 0.0125 | 0.0019 | 0.0003 | 0.0045 |
| Chest | 0.0153 | 0.0114 | 0.0353 | 0.008 | 0.0032 | 0.0238 |
| Neck | 0.0395 | 0.0338 | 0.0859 | 0.0257 | 0.0245 | 0.0805 |
| Head | 0.066 | 0.0609 | 0.1288 | 0.0484 | 0.0608 | 0.1175 |
| ClavicleLeft | 0.0409 | 0.0481 | 0.0798 | 0.0254 | 0.0414 | 0.0743 |
| ClavicleRight | 0.0412 | 0.0475 | 0.0794 | 0.0255 | 0.0412 | 0.0741 |
| BicepLeft | 0.0469 | 0.0723 | 0.0971 | 0.0297 | 0.0483 | 0.093 |
| ForarmLeft | 0.0941 | 0.1608 | 0.1577 | 0.0825 | 0.1293 | 0.1563 |
| HandLeft | 0.174 | 0.2697 | 0.1948 | 0.1564 | 0.2381 | 0.1964 |
| Index0Left | 0.2026 | 0.3167 | 0.21 | 0.1871 | 0.2878 | 0.2116 |
| Index1Left | 0.2121 | 0.3312 | 0.2164 | 0.1978 | 0.3022 | 0.2183 |
| Index2Left | 0.2184 | 0.3415 | 0.2208 | 0.2068 | 0.3133 | 0.2238 |
| Index2LeftEnd | 0.2247 | 0.3519 | 0.2259 | 0.2151 | 0.324 | 0.2292 |
| Middle0Left | 0.2049 | 0.317 | 0.2147 | 0.1862 | 0.2873 | 0.2161 |
| Middle1Left | 0.2155 | 0.3333 | 0.2216 | 0.1979 | 0.3028 | 0.2235 |
| Middle2Left | 0.2218 | 0.3437 | 0.226 | 0.2069 | 0.3144 | 0.2294 |
| Middle2LeftEnd | 0.2278 | 0.354 | 0.2314 | 0.2151 | 0.3258 | 0.2351 |
| Ring0Left | 0.2081 | 0.318 | 0.2207 | 0.1866 | 0.2872 | 0.2221 |
| Ring1Left | 0.2175 | 0.3324 | 0.2268 | 0.1974 | 0.3017 | 0.2296 |
| Ring2Left | 0.2225 | 0.341 | 0.2304 | 0.2053 | 0.3126 | 0.2352 |
| Ring2LeftEnd | 0.2273 | 0.3493 | 0.2349 | 0.2129 | 0.3231 | 0.2408 |
| Pinky0Left | 0.2106 | 0.3178 | 0.2262 | 0.187 | 0.2859 | 0.2276 |
| Pinky1Left | 0.2187 | 0.3301 | 0.2315 | 0.1969 | 0.2982 | 0.235 |
| Pinky2Left | 0.2232 | 0.3379 | 0.2347 | 0.2041 | 0.3071 | 0.2402 |
| Pinky2LeftEnd | 0.2272 | 0.3449 | 0.2382 | 0.2095 | 0.3142 | 0.2449 |
| Thumb0Left | 0.1884 | 0.2971 | 0.1982 | 0.1757 | 0.2681 | 0.2003 |
| Thumb1Left | 0.197 | 0.3108 | 0.2039 | 0.1849 | 0.2791 | 0.2021 |
| Thumb2Left | 0.2073 | 0.3258 | 0.2121 | 0.1961 | 0.2935 | 0.2075 |
| Thumb2LeftEnd | 0.2162 | 0.3374 | 0.2195 | 0.2054 | 0.3054 | 0.2124 |
| BicepRight | 0.05 | 0.0725 | 0.0969 | 0.0296 | 0.0543 | 0.086 |
| ForarmRight | 0.098 | 0.17 | 0.1524 | 0.0897 | 0.138 | 0.1603 |
| HandRight | 0.1703 | 0.2842 | 0.185 | 0.169 | 0.2372 | 0.1899 |
| Index0Right | 0.196 | 0.3329 | 0.2045 | 0.2037 | 0.2843 | 0.2093 |
| Index1Right | 0.2048 | 0.348 | 0.2118 | 0.215 | 0.2988 | 0.2168 |
| Index2Right | 0.2099 | 0.3586 | 0.2173 | 0.2251 | 0.3114 | 0.2223 |
| Index2RightEnd | 0.2145 | 0.3691 | 0.2234 | 0.2352 | 0.3235 | 0.227 |
| Middle0Right | 0.1984 | 0.3335 | 0.2077 | 0.2025 | 0.2821 | 0.2123 |
| Middle1Right | 0.2078 | 0.35 | 0.2154 | 0.2147 | 0.2968 | 0.2194 |
| Middle2Right | 0.2119 | 0.3586 | 0.2187 | 0.2221 | 0.3061 | 0.2232 |
| Middle2RightEnd | 0.215 | 0.3649 | 0.2207 | 0.2274 | 0.3133 | 0.2258 |
| Ring0Right | 0.2017 | 0.3347 | 0.212 | 0.2021 | 0.2803 | 0.2165 |
| Ring1Right | 0.2099 | 0.3495 | 0.2187 | 0.2142 | 0.2946 | 0.224 |
| Ring2Right | 0.2133 | 0.3566 | 0.2218 | 0.223 | 0.3047 | 0.2284 |
| Ring2RightEnd | 0.2153 | 0.3612 | 0.2234 | 0.2304 | 0.3136 | 0.2322 |
| Pinky0Right | 0.2043 | 0.3346 | 0.2158 | 0.2014 | 0.278 | 0.2201 |
| Pinky1Right | 0.2112 | 0.3472 | 0.2215 | 0.2118 | 0.2888 | 0.2261 |
| Pinky2Right | 0.2144 | 0.3544 | 0.2251 | 0.2192 | 0.2969 | 0.2296 |
| Pinky2RightEnd | 0.2165 | 0.3591 | 0.2276 | 0.2242 | 0.3039 | 0.2316 |
| Thumb0Right | 0.1828 | 0.3126 | 0.1926 | 0.1911 | 0.2676 | 0.198 |
| Thumb1Right | 0.1907 | 0.3265 | 0.1985 | 0.2013 | 0.2801 | 0.2021 |
| Thumb2Right | 0.1996 | 0.3399 | 0.205 | 0.2123 | 0.2949 | 0.2091 |
| Thumb2RightEnd | 0.2064 | 0.3494 | 0.2096 | 0.2205 | 0.3061 | 0.2148 |
| ThighLeft | 0 | 0 | 0 | 0 | 0 | 0 |
| CalfLeft | 0.083 | 0.1416 | 0.1341 | 0.0673 | 0.152 | 0.1385 |
| FootLeft | 0.1263 | 0.1611 | 0.1942 | 0.1185 | 0.1861 | 0.2062 |
| ToeLeft | 0.1453 | 0.1705 | 0.2202 | 0.1338 | 0.1943 | 0.2308 |
| ToeLeftEnd | 0.1632 | 0.1789 | 0.2254 | 0.1494 | 0.1985 | 0.2349 |
| ThighRight | 0 | 0 | 0 | 0 | 0 | 0 |
| CalfRight | 0.0781 | 0.1352 | 0.1381 | 0.0695 | 0.1495 | 0.1436 |
| FootRight | 0.1258 | 0.1659 | 0.2047 | 0.1181 | 0.1908 | 0.209 |
| ToeRight | 0.1435 | 0.1743 | 0.2342 | 0.1319 | 0.1937 | 0.2353 |
| ToeRightEnd | 0.158 | 0.1788 | 0.2414 | 0.1471 | 0.1924 | 0.2402 |

Table 4: Per-joint quaternion component standard deviations computed over two proposed datasets

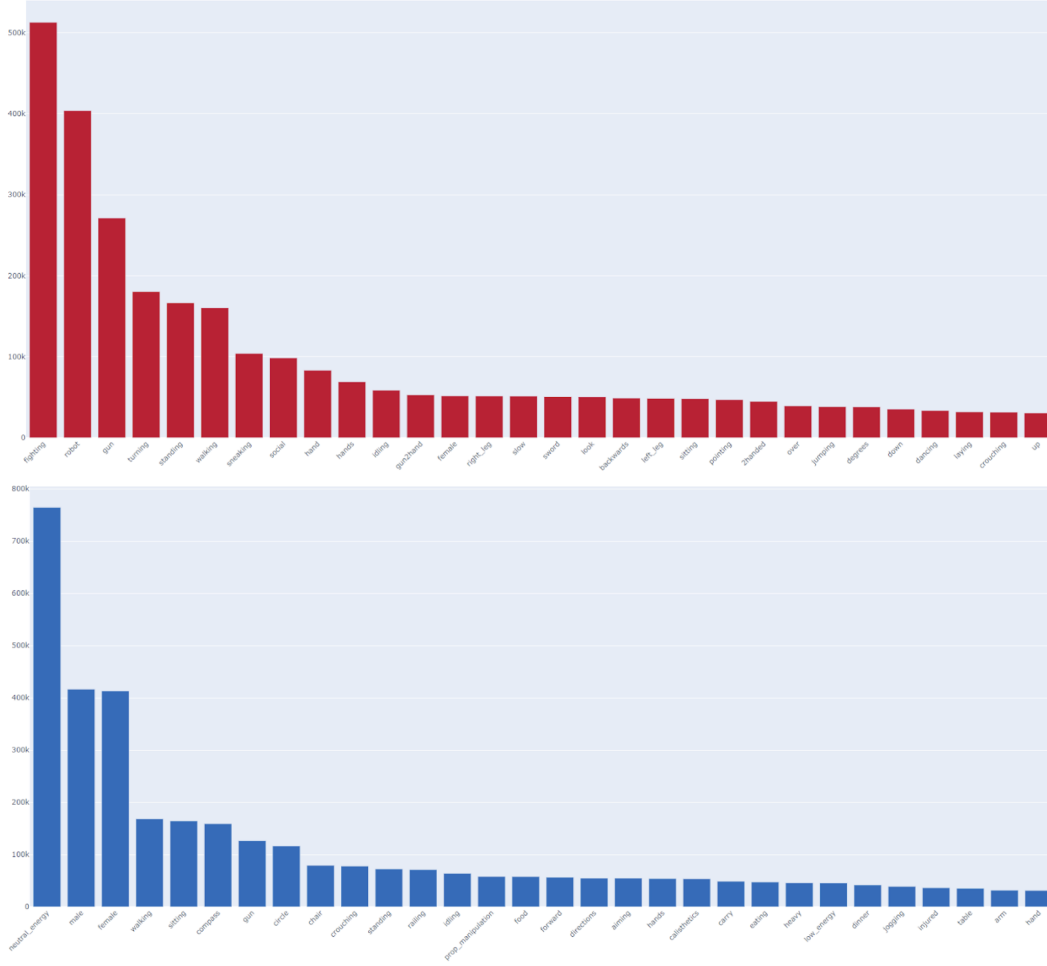| Joint | miniMixamo | | | | miniAnonymous | | | |
|---|---|---|---|---|---|---|---|---|
| | X | Y | Z | W | X | Y | Z | W |
| Hips | 0.187 | 0.2239 | 0.0952 | 0.5381 | 0.1614 | 0.4735 | 0.0951 | 0.4673 |
| Spine0 | 0.0928 | 0.0343 | 0.0358 | 0.0133 | 0.0325 | 0.019 | 0.0139 | 0.0015 |
| Spine1 | 0.0579 | 0.038 | 0.0344 | 0.0076 | 0.0694 | 0.0308 | 0.021 | 0.0065 |
| Chest | 0.0612 | 0.0292 | 0.0227 | 0.0065 | 0.0695 | 0.0319 | 0.0204 | 0.0075 |
| Neck | 0.1241 | 0.0682 | 0.0572 | 0.0236 | 0.1368 | 0.0887 | 0.0691 | 0.0347 |
| Head | 0.1183 | 0.1206 | 0.0703 | 0.0374 | 0.1349 | 0.0948 | 0.0746 | 0.0337 |
| ClavicleLeft | 0.2139 | 0.2156 | 0.2234 | 0.2039 | 0.4424 | 0.3554 | 0.4045 | 0.3964 |
| ClavicleRight | 0.4307 | 0.4887 | 0.4985 | 0.4309 | 0.4785 | 0.459 | 0.475 | 0.4696 |
| BicepLeft | 0.188 | 0.2132 | 0.1935 | 0.0744 | 0.2015 | 0.2523 | 0.2225 | 0.1059 |
| ForarmLeft | 0.1131 | 0.1535 | 0.259 | 0.1768 | 0.1958 | 0.2931 | 0.2892 | 0.2372 |
| HandLeft | 0.1732 | 0.2283 | 0.1276 | 0.2091 | 0.0034 | 0.1179 | 0.0603 | |
| Index0Left | 0.188 | 0.0278 | 0.047 | 0.0623 | 0.2013 | 0.0117 | 0.063 | 0.055 |
| Index1Left | 0.222 | 0.0055 | 0.0225 | 0.1145 | 0.1915 | 0.0007 | 0.0014 | 0.0787 |
| Index2Left | 0.178 | 0.0056 | 0.0157 | 0.0675 | 0.1743 | 0.0006 | 0.0014 | 0.0587 |
| Index2LeftEnd | 0 | 0 | 0 | 0 | 0.0007 | 0.0006 | 0.0014 | 0.0013 |
| Middle0Left | 0.2031 | 0.0203 | 0.0395 | 0.0757 | 0.2614 | 0.0175 | 0.0166 | 0.0933 |
| Middle1Left | 0.2265 | 0.0056 | 0.0236 | 0.1177 | 0.1648 | 0.0033 | 0.0017 | 0.0594 |
| Middle2Left | 0.186 | 0.0052 | 0.0153 | 0.0714 | 0.1758 | 0.0012 | 0.0017 | 0.0522 |
| Middle2LeftEnd | 0 | 0 | 0 | 0 | 0.0007 | 0.0005 | 0.0014 | 0.0012 |
| Ring0Left | 0.2151 | 0.0242 | 0.046 | 0.0908 | 0.2561 | 0.0249 | 0.0297 | 0.087 |
| Ring1Left | 0.2306 | 0.0075 | 0.0235 | 0.1218 | 0.1867 | 0.0004 | 0.0014 | 0.0773 |
| Ring2Left | 0.186 | 0.0058 | 0.0179 | 0.0706 | 0.1498 | 0.0005 | 0.0014 | 0.048 |
| Ring2LeftEnd | 0 | 0 | 0 | 0 | 0.0007 | 0.0005 | 0.0014 | 0.0012 |
| Pinky0Left | 0.2179 | 0.0377 | 0.0701 | 0.094 | 0.2533 | 0.0585 | 0.0901 | 0.0874 |
| Pinky1Left | 0.2133 | 0.0104 | 0.0238 | 0.0929 | 0.1978 | 0.0004 | 0.0014 | 0.0853 |
| Pinky2Left | 0.1999 | 0.009 | 0.0187 | 0.1088 | 0.1798 | 0.0008 | 0.0012 | 0.0701 |
| Pinky2LeftEnd | 0 | 0 | 0 | 0 | 0.0009 | 0.0008 | 0.0012 | 0.0012 |
| Thumb0Left | 0.1266 | 0.0871 | 0.2072 | 0.0443 | 0.0957 | 0.1374 | 0.071 | 0.1124 |
| Thumb1Left | 0.0447 | 0.0613 | 0.0984 | 0.0422 | 0.0635 | 0.0006 | 0.0006 | 0.0127 |
| Thumb2Left | 0.0625 | 0.0502 | 0.1524 | 0.0531 | 0.1265 | 0.0004 | 0.0007 | 0.0174 |
| Thumb2LeftEnd | 0 | 0 | 0 | 0 | 0.0012 | 0.0004 | 0.0007 | 0.001 |
| BicepRight | 0.1912 | 0.206 | 0.2099 | 0.0897 | 0.2043 | 0.2555 | 0.2305 | 0.1083 |
| ForarmRight | 0.1039 | 0.1349 | 0.3054 | 0.2408 | 0.2122 | 0.2554 | 0.4638 | 0.4425 |
| HandRight | 0.1609 | 0.1848 | 0.1284 | 0.0902 | 0.2169 | 0.0082 | 0.15 | 0.0648 |
| Index0Right | 0.2545 | 0.0292 | 0.0448 | 0.9321 | 0.197 | 0.0148 | 0.0535 | 0.0535 |
| Index1Right | 0.3376 | 0.0053 | 0.0201 | 0.9063 | 0.1763 | 0.0013 | 0.0009 | 0.0662 |
| Index2Right | 0.2512 | 0.0055 | 0.0155 | 0.9351 | 0.158 | 0.0008 | 0.0013 | 0.0449 |
| Index2RightEnd | 0 | 0 | 0 | 0 | 0.0008 | 0.0009 | 0.0013 | 0.0014 |
| Middle0Right | 0.2965 | 0.0163 | 0.0516 | 0.9213 | 0.261 | 0.0177 | 0.0176 | 0.0915 |
| Middle1Right | 0.4421 | 0.0074 | 0.0226 | 0.8616 | 0.1968 | 0.0048 | 0.0017 | 0.0765 |
| Middle2Right | 0.2967 | 0.0079 | 0.0154 | 0.9217 | 0.2133 | 0.0014 | 0.0014 | 0.0762 |
| Middle2RightEnd | 0 | 0 | 0 | 0 | 0.0036 | 0.0009 | 0.0012 | 0.0013 |
| Ring0Right | 0.3277 | 0.026 | 0.0569 | 0.911 | 0.2494 | 0.0263 | 0.0405 | 0.0754 |
| Ring1Right | 0.4355 | 0.0079 | 0.0238 | 0.8645 | 0.2002 | 0.001 | 0.0011 | 0.073 |
| Ring2Right | 0.3114 | 0.0099 | 0.0167 | 0.9168 | 0.1543 | 0.0009 | 0.0012 | 0.0415 |
| Ring2RightEnd | 0 | 0 | 0 | 0 | 0.0009 | 0.0009 | 0.0012 | 0.0012 |
| Pinky0Right | 0.3372 | 0.0452 | 0.076 | 0.9051 | 0.2767 | 0.0574 | 0.0737 | 0.1006 |
| Pinky1Right | 0.3847 | 0.0131 | 0.0251 | 0.8887 | 0.1991 | 0.001 | 0.001 | 0.0786 |
| Pinky2Right | 0.314 | 0.0109 | 0.0171 | 0.9143 | 0.1973 | 0.001 | 0.0011 | 0.0668 |
| Pinky2RightEnd | 0 | 0 | 0 | 0 | 0.001 | 0.001 | 0.0011 | 0.0012 |
| Thumb0Right | 0.1855 | 0.0962 | 0.1144 | 0.9211 | 0.2739 | 0.6352 | 0.1285 | 0.4566 |
| Thumb1Right | 0.0484 | 0.0561 | 0.265 | 0.9293 | 0.0747 | 0.0007 | 0.0012 | 0.0179 |
| Thumb2Right | 0.0652 | 0.0359 | 0.2321 | 0.9364 | 0.169 | 0.0009 | 0.001 | 0.04 |
| Thumb2RightEnd | 0 | 0 | 0 | 0 | 0.0012 | 0.0009 | 0.001 | 0.0014 |
| ThighLeft | 0.107 | 0.0975 | 0.0818 | 0.2296 | 0.1164 | 0.0884 | 0.0704 | 0.2397 |
| CalfLeft | 0.2419 | 0.075 | 0.0448 | 0.145 | 0.2698 | 0.049 | 0.0427 | 0.1751 |
| FootLeft | 0.1158 | 0.0685 | 0.073 | 0.0554 | 0.1258 | 0.091 | 0.0894 | 0.0697 |
| ToeLeft | 0.0658 | 0.026 | 0.018 | 0.0552 | 0.052 | 0.045 | 0.0365 | 0.0328 |
| ToeLeftEnd | 0 | 0 | 0 | 0 | 0.0009 | 0.0002 | 0.0009 | 0.0008 |
| ThighRight | 0.1083 | 0.0959 | 0.0783 | 0.2298 | 0.1234 | 0.0895 | 0.0783 | 0.2442 |
| CalfRight | 0.2458 | 0.0738 | 0.0437 | 0.1525 | 0.2835 | 0.0426 | 0.03 | 0.1927 |
| FootRight | 0.1237 | 0.0685 | 0.077 | 0.059 | 0.1215 | 0.1012 | 0.0912 | 0.0627 |
| ToeRight | 0.0786 | 0.0257 | 0.0154 | 0.0649 | 0.0595 | 0.0431 | 0.0286 | 0.0483 |
| ToeRightEnd | 0 | 0 | 0 | 0 | 0.001 | 0.0002 | 0.0009 | 0.0008 |

Figure 5: The distribution of tags across frames in the original Mixamo (top) and Anonymous (bottom) datasets.

## E    TRAINING AND EVALUATION SETUP: DETAILS

We use Algorithm 1 of Appendix C to sample batches of size 2048 from the training subset. The number of effectors is sampled once per batch and is fixed for all batch items to maximize data throughput. The training loop is implemented in PyTorch (Paszke et al., 2019) using Adam optimizer (Kingma & Ba, 2015) with a learning rate of 0.0002. Hyperparameter values are adjusted on the validation set (see Appendix E for hyperparameter settings).

We report $\mathcal{L}_{gpd-L2}^{det}$, $\mathcal{L}_{ikd-L2}^{det}$, $\mathcal{L}_{loc-geo}^{det}$ metrics calculated on the test set, using models trained on the training set. $\mathcal{L}_{gpd-L2}^{det}$ is computed only on the root joint. These metrics characterise both the 3D position accuracy ($\mathcal{L}_{gpd-L2}^{det}$, $\mathcal{L}_{ikd-L2}^{det}$) and the bone rotation accuracy ($\mathcal{L}_{loc-geo}^{det}$). They are defined as follows:

$$\mathcal{L}_{gpd-L2}^{det} = \text{MSE}(\mathbf{g}_0, \widetilde{\mathbf{f}}_{R,0}) \tag{22}$$

$$\mathcal{L}_{ikd-L2}^{det} = \sum_{j=1}^{J} \text{MSE}(\mathbf{g}_j, \widehat{\mathbf{g}}_j) \tag{23}$$

$$\mathcal{L}_{loc-geo}^{det} = \sum_{j=1}^{J} \text{GEO}(\mathbf{R}_j, \widehat{\mathbf{R}}_j) \tag{24}$$
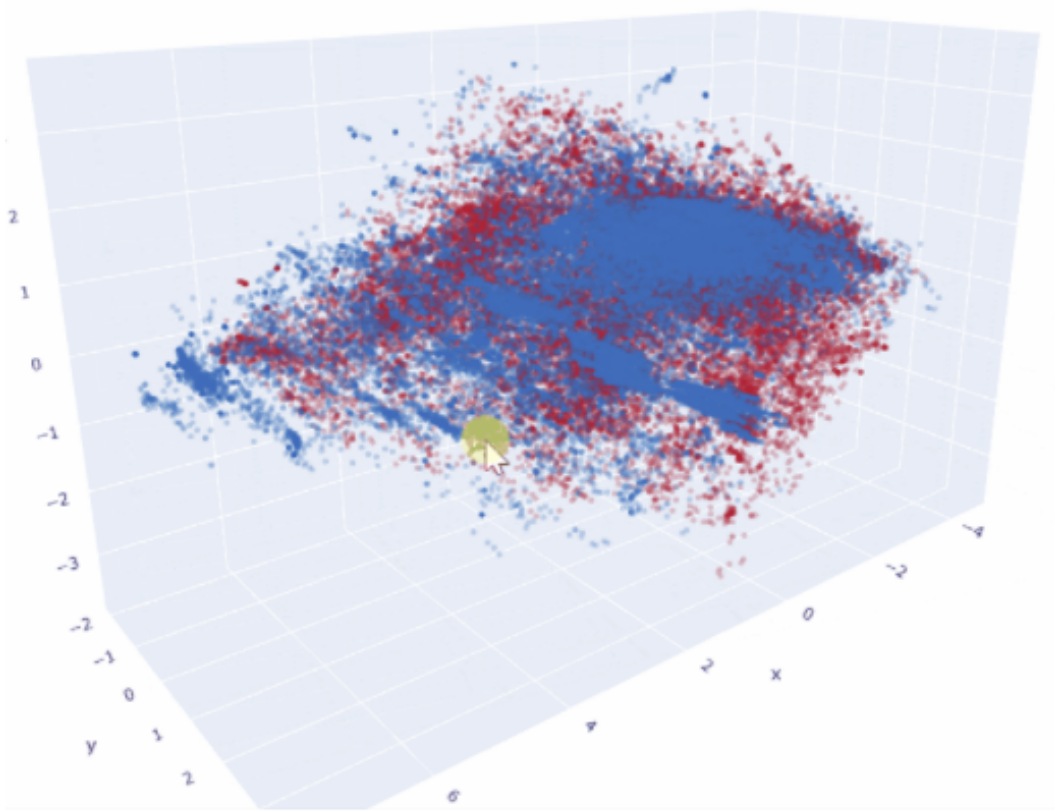
21

Figure 6: The distribution of PCA of flattened rotations of all joints of miniMixamo (red) and miniAnonymous (blue) datasets. Each point corresponds to a pose. It is clear that there are both overlapping areas corresponding to clusters with similar poses and areas with disjoint clusters, showing distinct pose configurations characteristic of the two different datasets.

| Hyperparameter | Value | Grid |
|---|---|---|
| Epochs, miniMixamo/ miniAnonymous | 40k/15k | [20k, 40k, 80k] / [10k, 15k, 40k] |
| Losses | MSE, GEO, LAT | MSE, GEO, LAT |
| Width ($d_h$) | 1024 | [256, 512, 1024, 2048] |
| Blocks ($R$) | 3 | [1, 2, 3] |
| Layers ($L$) | 3 | [2, 3, 4] |
| Batch size | 2048 | [512, 1024, 2048, 4096] |
| Optimizer | Adam | [Adam, SGD] |
| Learning rate | 2e-4 | [1e-4, 2e-4, 5e-4, 1e-3] |
| Base L2 loss scale ($W_{pos}$) | 1e2 | [1, 10, 1e2, 1e3, 1e4] |
| Max noise scale ($\sigma_{M,0}, \sigma_{M,1}$) | 0.1 | [0.01, 0.1, 1] |
| Max effector weight ($W_M$) | 1e3 | [10, 1e2, 1e3, 1e4] |
| Noise exponent, $\eta$ | 13 | 13 |
| Dropout | 0.01 | [0.0, 0.01, 0.05, 0.1, 0.2] |
| Embedding dimensionality | 32 | [16, 32, 64, 128] |
| Augmentattion | mirror, rotation | [mirror, rotation, translation] |

Table 5: Settings of ProtoRes hyperparameters and the hyperparameter search grid.

Here $\widetilde{\mathbf{f}}_{R,0}$ and $\mathbf{g}_0$ are ground truth global location of the root joint and its prediction from the GPD, respectively; $\mathbf{g}_j$ and $\widehat{\mathbf{g}}_j$ ground truth location of joint $j$ obtained by subjecting the rotation prediction

| High Energy | Low Energy | Injured | Neutral Energy | | |
|---|---|---|---|---|---|
| Generic | Generic | Generic | Feminine | Macho | Generic |
| *Female & Male* | *Female & Male* | *Female & Male* | *Female* | *Male* | *Female & Male* |
| Locomotion | Locomotion | Locomotion | Locomotion | Locomotion | Locomotion |
| - Crouch | - Crouch | - Crouch | Crouch | Crouch | - Crawls |
| - Jogs | - Jogs | - Jogs | - Jogs | - Jogs | - Crouch |
| - Runs | - Runs | - Runs | - Runs | - Runs | - Jogs |
| - Walk | - Walk | - Walk | - Walk | - Walk | - Rifle Aim Walk |
| | | | | | - Rifle Crawl |
| | | | | | - Rifle Crouch |
| | | | | | - Rifle Jogs |
| | | | | | - Rifle Walks |
| | | | | | - Runs |
| | | | | | - Walks |
| | | | | | - Idles |
| | | | | | - Rifle Idles |
| | | | | | - Walk Heavy Backpack |
| | | | | | - Walk Heavy Sack |
| | | | | | - Walk Ladder |
| | | | | | - Walk Dragging Object |
| | | | | | Archery |
| | | | | | Bokken Fighting |
| | | | | | Calisthetics |
| | | | | | Door Interactions |
| | | | | | Fist Fighting |
| | | | | | Food |
| | | | | | Handgun |
| | | | | | - Downwards |
| | | | | | - Level |
| | | | | | - Misc |
| | | | | | - Upwards |
| | | | | | - Verticals |
| | | | | | Hands |
| | | | | | Knife Fighting |
| | | | | | Longsword Fighting |
| | | | | | Misc |
| | | | | | Phone |
| | | | | | Place |
| | | | | | Railing Interactions |
| | | | | | Rifle |
| | | | | | Seated Interactions |
| | | | | | Shotgun |
| | | | | | Standing |
| | | | | | Sword |
| | | | | | Wall |

Table 6: The hierarchy of motion styles and categories in miniAnonymous. *Female & Male* indicate that every clip of every underlying category has been captured with both a female and a male actor.

from IKD to the forward kinematics process; $\mathbf{R}_j$ and $\widehat{\mathbf{R}}_j$ are ground truth local rotation matrix of joint $j$ and its prediction obtained from the IKD, respectively.

The evaluation framework tests model performance on a pre-generated set of seven files containing 6, 7, ..., 12 effectors respectively. Skeleton is split in six zones, with four main zones including each limb, the hip zone and the head zone. In each file, we first sample one positional effector from each main zone. Remaining effectors are sampled randomly from all zones and effector types, mimicking pose authoring scenarios observed in practice. Metrics are averaged over all samples in all files, assessing the overall quality of pose reconstruction in scenario with sparse and variable inputs. All tables present results averaged over 4 random seed retries and metric values computed every 10 epochs over last 1000 epochs, rounded to the last statistically significant digit.

**Hyperparameter settings**   The training loop is implemented in PyTorch (Paszke et al., 2019) using Adam optimizer (Kingma & Ba, 2015) with a learning rate of 0.0002. We tried to use SGD optimizer

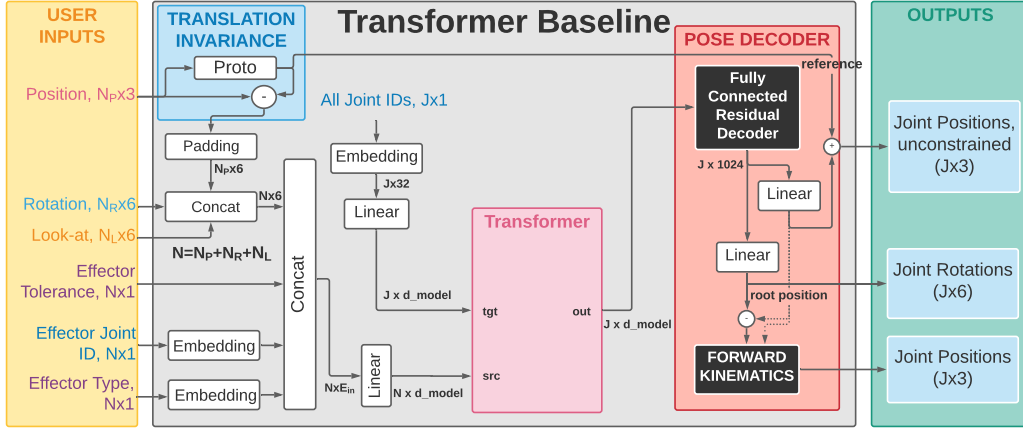| Motion Category |
| --- |
| Combat |
| Adventure |
| Sport |
| Dance |
| Fantasy |
| Superhero |

Table 7: Motion categories in miniMixamo



Figure 7: Block diagram of the Transformer baseline architecture.

to train the architecture, but it was very difficult to obtain stable results with it. Adam optimizer turned out to be much more suitable for our problem. The learning rate was selected to be 0.0002, which is lower than Adam's default. Obtaining stable training results with higher learning rates was not feasible. Batch size is selected to be 2048 to accelerate training speed. In practice we observed slightly better generalization results with smaller batch size (1024 and 512). The detailed settings of ProtoRes hyperparameters are presented in Table 5.

## F  MASKED-FCR BASELINE ARCHITECTURE

Masked-FCR is a brute-force unstructured baseline that uses a very wide $J \cdot 3 \cdot 7$ input layer ($J$ joints, 3 effector types, 6D effector plus one tolerance value) to handle all possible effector permutations. Each missing effector is masked with one of $3 \cdot J$ learnable 7D placeholders. Masked-FCR has 3 encoder and 6 decoder blocks to match ProtoRes.

## G  TRANSFORMER BASELINE ARCHITECTURE

We implement Transformer baseline using the default Transformer module available from PyTorch `https://pytorch.org/docs/stable/generated/torch.nn.Transformer.html`. The block diagram of the Transformer baseline architecture is shown in Fig. 7.

We use a standard transformer application scenario in which the transformer source input is fed with the variable length input and the required outputs are queried via target input. In our case the variable length input corresponds to the effector data concatenated with embedded effector categorical variables. The query for the output consists of the embeddings of all joints. Note that the joint embedding is reused both for source and target inputs and both inputs are projected to the internal d_model dimensionality of transformer.

| Hyperparameter | Value | Grid |
|---|---|---|
| FCR decoder parameteres | | |
| FCR Blocks ($R$) | 6 | N/A |
| FCR Width ($d_h$) | 1024 | N/A |
| Transformer parameteres | | |
| d_model | 128 | [64, 128, 256] |
| nhead | 8 | [1, 2, 4, 8] |
| num_encoder_layers | 2 | [1,2,3] |
| num_decoder_layers | 2 | [1,2,3] |
| dim_feedforward | 1024 | [256, 512, 1024] |
| dropout | 0.01 | 0.01 |
| activation | relu | relu |
| Base L2 loss scale ($W_{pos}$) | 100 | [10, 100] |

Table 8: Settings of Transformer baseline hyperparameters and the hyperparameter search grid.

The embeddings of joint IDs are used to query the Transformer output, producing one encoder embedding for each of $J$ joints. The $J$ encodings are fed into the 6-block FCR decoder (to match the total number of decoder blocks in ProtoRes) with two heads: one predicting rotation and one predicting unconstrained position. This is similar to the use of Transformer to predict bounding box class IDs and sizes for object detection (Carion et al., 2020). Predictions of rotations and of the root joint are used in the forward kinematics pass, just as in ProtoRes.

Internally, Transformer processes both source and target inputs via self-attention first and then applies the multi-head attention between source and target after self-attention. This results in the output embedding for each skeleton joint that depends on all the input information as well as the learned interactions across all output skeleton joints. The output embedding of the transformer is then decoded to unconstrained position and rotation outputs using a two-headed Fully-Connected residual stack (this is the same architecture as the one used in ProtoRes decoders). Note that this is a well-known Transformer application scheme that has recently been used to achieve SOTA results in object detection, for example (Carion et al., 2020). Table 8 lists the hyperparameter settings for the Transformer baseline. Note that only hyperparameters that are unique to this baseline or different from the ProtoRes defaults appearing in Table 5 are listed.

## H   ABLATION OF THE DECODER: DETAILS

The detailed results of the decoder ablation are shown in Table 9. One block of decoder is much less computationally expensive than one block of encoder. One block of encoder processes $N$ effectors, whereas the decoder deals with the partially defined pose representation collapsed to a vector. Hence the decoder block is $N$ times less expensive. To demonstrate the effectiveness of decoder, we keep all hyperparameters at defaults described in Section 3.2 and vary the number of encoder and decoder blocks in ProtoRes (0 decoder blocks corresponds to a simple linear projection of encoder output). We measure the train time of each configuration on NVIDIA M40 24GB GPU installed on Dell PowerEdge R720 server with two Intel Xeon E5-2667 2.90GHz CPU. The table reveals a few things. First, adding more decoder blocks significantly increases accuracy when the number of encoder blocks is lower (e.g. 3 or 5). When the number of encoder blocks is high (e.g. 7) linear projection provides similar accuracy. Second, using non-trivial decoder is computationally more efficient. For example, the 3 encoder and 3 decoder blocks configuration has comparable accuracy with 5 encoder and 1 decoder blocks, however it is noticeably more compute efficient (5+1 configuration uses 40% more compute time than 3+3).

Table 9: Ablation of the decoder. Random benchmark, lower values are better. Train time is measured on a 2GPU Dell PowerEdge R720 server with two NVIDIA M40 24GB GPUs and two Intel Xeon E5-2667 2.90GHz processors, each GPU running one ProtoRes training session.

| | | miniMixamo | | | | miniAnonymous | | | |
|---|---|---|---|---|---|---|---|---|---|
| Encoder blocks | Decoder blocks | $\mathcal{L}_{gpd-L2}^{det}$ | $\mathcal{L}_{ikd-L2}^{det}$ | $\mathcal{L}_{loc-geo}^{det}$ | Train time, h | $\mathcal{L}_{gpd-L2}^{det}$ | $\mathcal{L}_{ikd-L2}^{det}$ | $\mathcal{L}_{loc-geo}^{det}$ | Train time, h |
| 3 | 0 | 1.54e-3 | 4.59e-3 | 0.2485 | 91 | 1.05e-3 | 3.65e-3 | 0.1939 | 105 |
| 3 | 1 | 1.35e-3 | 4.34e-3 | 0.2433 | 95 | 0.93e-3 | 3.52e-3 | 0.1895 | 110 |
| 3 | 2 | 1.34e-3 | 4.24e-3 | 0.2397 | 102 | 0.93e-3 | 3.34e-3 | 0.1840 | 116 |
| 3 | 3 | 1.36e-3 | 4.16e-3 | 0.2381 | 106 | 0.93e-3 | 3.28e-3 | 0.1817 | 121 |
| 5 | 0 | 1.27e-3 | 4.20e-3 | 0.2399 | 144 | 0.84e-3 | 3.27e-3 | 0.1824 | 166 |
| 5 | 1 | 1.28e-3 | 4.30e-3 | 0.2390 | 148 | 0.81e-3 | 3.24e-3 | 0.1818 | 171 |
| 5 | 2 | 1.15e-3 | 4.02e-3 | 0.2345 | 153 | 0.77e-3 | 3.10e-3 | 0.1791 | 176 |
| 5 | 3 | 1.18e-3 | 4.03e-3 | 0.2351 | 157 | 0.82e-3 | 3.07e-3 | 0.1785 | 182 |
| 7 | 0 | 1.13e-3 | 4.00e-3 | 0.2355 | 196 | 0.74e-3 | 2.98e-3 | 0.1762 | 226 |
| 7 | 1 | 1.23e-3 | 4.16e-3 | 0.2356 | 200 | 0.79e-3 | 3.07e-3 | 0.1780 | 230 |
| 7 | 2 | 1.13e-3 | 4.51e-3 | 0.2383 | 205 | 0.78e-3 | 3.10e-3 | 0.1780 | 236 |
| 7 | 3 | 1.15e-3 | 3.98e-3 | 0.2352 | 209 | 0.82e-3 | 3.14e-3 | 0.1783 | 241 |

## I  ABLATION OF THE PROTOTYPE-SUBTRACT-ACCUMULATE STACKING PRINCIPLE

Here we compare the proposed Prototype-Subtract-Accumulate (PSA) residual stacking scheme described in equations (2)-(5) against the ResPointNet (Niemeyer et al., 2019) stacking scheme: Maxpool-Concat Daisy Chain (MCDC). To implement the MCDC stacking proposed by Niemeyer et al. (2019) we setup the experiment as follows.

- Equation (2) is replaced with the concatenation of the output of the previous block, $\mathbf{b}_r$, with the maxpool of the previous block output along axis 1 (we use batch, effector, channels convention for tensor axes 0,1,2; respectively)

- In equation (4) we only compute $\mathbf{b}_r$, since $\mathbf{f}_r$ is not used

- Equation (5) is removed

- The final pose embedding is created using the maxpool along axis 1 of $\mathbf{b}_r$ at the last encoder block

- We use decoder with 0 blocks, *i.e.* only the linear projection at the end to derive both PSA and MCDC results. This is because (i) residual decoder is not part of the original design by Niemeyer et al. (2019), (ii) in this study we focus exclusively on the effects of stacking within the encoder, equalizing all other experimental conditions.

- Hyperparameters of both architectures are taken from Table 5 in Appendix C

We show that our stacking scheme is more accurate and allows stacking deeper networks more effectively. Quantitative results are shown in Table 10. It is clear that the proposed stacking approach provides gain in setups with varying number of encoder blocks. For example, in the case of small number of 3 blocks (computationally efficient setup) our PSA approach is clearly more accurate than MCDC approach of (Niemeyer et al., 2019). Similarly, with 7 blocks (training time is more than two times longer) PSA is again more accurate than MCDC. Moreover, in the case of PSA we see noticeable accuracy improvement while increasing the number of encoder blocks from 5 to 7, whereas MCDC provides almost no additional gain beyond 5 blocks. We conclude that the proposed PSA stacking mechanism is better than MCDC as it provides better accuracy in the computationally efficient configuration and it is significantly more effective at supporting deeper architectures that provide better accuracy in the case of PSA, whereas MCDC accuracy saturates at 5 blocks providing little to no accuracy gain with deeper architectures.

Table 10: Ablation of the Prototype-Subtract-Accumulate. Random benchmark, lower values are better.

| | | miniMixamo | | | miniAnonymous | | |
|---|---|---|---|---|---|---|---|
| | | $\mathcal{L}_{gpd-L2}^{det}$ | $\mathcal{L}_{ikd-L2}^{det}$ | $\mathcal{L}_{loc-geo}^{det}$ | $\mathcal{L}_{gpd-L2}^{det}$ | $\mathcal{L}_{ikd-L2}^{det}$ | $\mathcal{L}_{loc-geo}^{det}$ |
| Stacking scheme | Encoder blocks | | | | | | |
| MCDC | 3 | 1.52e-3 | 4.65e-3 | 0.2539 | 1.12e-3 | 3.93e-3 | 0.1984 |
| MCDC | 5 | 1.28e-3 | 4.25e-3 | 0.2395 | 0.85e-3 | 3.25e-3 | 0.1831 |
| MCDC | 7 | 1.28e-3 | 4.24e-3 | 0.2384 | 0.85e-3 | 3.26e-3 | 0.1830 |
| PSA (**ours**) | 3 | 1.54e-3 | 4.59e-3 | 0.2485 | 1.05e-3 | 3.65e-3 | 0.1939 |
| PSA (**ours**) | 5 | 1.27e-3 | 4.20e-3 | 0.2399 | 0.84e-3 | 3.27e-3 | 0.1824 |
| PSA (**ours**) | 7 | 1.13e-3 | 4.00e-3 | 0.2355 | 0.74e-3 | 2.98e-3 | 0.1762 |

## J  QUALITATIVE COMPARISON TO TRANSFORMER AND FINALIK BASELINES

In this section we provide additional qualitative comparison between ProtoRes and two baselines, the non-learnable baseline FinalIK and the machine learning baseline Transformer. The FinalIK comparison results appear in the top row of Figure 8. The Transformer comparison results appear in the bottom row of Figure 8.

The characteristic feature of FinalIK is that it lacks inductive bias towards realistic poses. Therefore it is very easy to create effector configurations that result in unnatural poses as can be seen in Figure 8. Conversely, creating realistic poses requires significant amount of tuning of the individual joints as most joints of the body act very independently and locally. With FinalIK, it is relatively hard to produce a believable pose with only a few effectors. ProtoRes fills this gap by providing a learned prior for realistic poses. As a result, most of the poses created by ProtoRes look natural, no matter how many effectors are used to steer the pose. In addition, the effector space can be augmented dynamically to capture missing accents, if necessary.

When comparing ProtoRes against Transformer, we can see that Transformer has a less global approach to forming a pose than ProtoRes. This manifests itself in the Tranformer based model having a tendency to create poses in which limbs penetrate each other or the rest of the body (see Figure 8 bottom row, pictures 1,2,3 from the left). Also, Transformer is very good at following the effector inputs and reproducing them in the reconstructed pose — so much that it is willing to sacrifice the overall final pose plausibility at the expense of sticking to the effector, taking the effector guidance very "literally" and missing its re-interpretation given global context. A good example of this behaviour is presented in Figure 8 bottom row, rightmost picture. We see that ProtoRes brings the entire body to the floor and prepares the arms to touch the floor to support the body, all in the attempt to produce a plausible pause of someone looking at the look-at target that is placed relatively low. On the other hand, Transformer takes the look-at effector guidance very literally and locally, willing to break the neck of the pose to look at the provided target and not realizing that the entire body position needs to be changed to accommodate for this rather peculiar configuration of effectors to create a believable pose. In general, our observation is that the Transformer based baseline response to effector inputs tends to be localized. As a consequence, its behaviour has a very pronounced pure IK flavour to it, at times reminding the behaviour of FinalIK, especially when facing difficult effector configurations. A lot of the time, however, the outputs of ProtoRes and Transformer are comparable and consistent, since both of them develop strong inductive bias towards realistic poses.

## K  VIDEOS AND DEMONSTRATIONS

We provide here descriptions related to each video found in the supplementary materials. The Unity tool that is used to produce demo videos and qualitative pictures shows the raw neural network outputs without applying any post-processing. This is done to ensure fair comparison of different methods and avoid any misleading results due to post-processing. Note that most of the demonstrations are done using a ProtoRes model trained on a large internal dataset not evaluated in this work. One of

Figure 8: Qualitative comparison results. ProtoRes vs. FinalIK (top row). ProtoRes vs. Transformer (bottom row).

our demonstrations, described below in Appendix K.5 qualitatively shows some of the most severe impacts of training on ablated datasets.

## K.1 PROTORES DEMO

The video presents an overview of the integration of ProtoRes as a posing tool inside the Unity game engine, showcasing how different effector types can be manipulated and how they can influence the resulting pose. In this demo, a user-defined configuration is presented in the UI, allowing one to choose which effectors are enabled within that configuration. Note that this configuration could contain more or less effectors of each type, and can be built for specific posing needs. The most generic configuration would present all possible effectors inside each effector type sub-menu.

## K.2 POSING FROM IMAGES

This video presents screen recordings of a novice user using ProtoRes to quickly prototype poses taken from 2D silhouette images. Note that one can reach satisfactory results in less than a minute in each case, with a relatively low number of manipulations. Note also that fine-tuning the resulting poses can always be achieved by adding more effectors and applying more manipulations.

## K.3 LOSS ABLATION

This recording shows a setup where different models are used with identical effector setup. Both models use the ProtoRes architecture. On the left, the model uses the total loss presented in Algorithm 1, whereas the right-hand side model uses positional losses only (GPD and IKD positional losses) and both local and global rotational losses, as well as look-at losses are disabled. This demonstration clearly shows how positional constraints, even when respected, do not suffice to produce realistic human poses. Joint rotations have to be modeled as well.

## K.4 FINALIK COMPARISON

This recording shows a setup where ProtoRes is compared to a full body biped IK system provided by FinalIK RootMotion (2020), with an identical effector setup. Note that FinalIK solves constraints by modifying the current pose, often resulting in smaller changes in the output, when compared to ProtoRes that predicts a full pose at each update. The lack of a learned model of human poses in FinalIK becomes quickly noticeable when manipulating effectors significantly.

## K.5 DATASETS COMPARISON

In these demonstrations, we showcase how training ProtoRes on different datasets can impact the results. We showcase models trained on the two ablated datasets presented in this work, i.e. miniAnonymous (left) and miniMixamo (right). We also show performance of a model trained on the full Mixamo Adobe Inc. (2020) dataset (center) to qualitatively show how performance can be improved with more training data. In all of these recordings, one can notice differences in the resulting poses, emphasizing the fact that human posing from few effectors, when no extra conditioning signals are used, is an ambiguous task that will be influenced by the training data.

The first sequence makes this fact especially obvious on the finger joints and the head's look-at direction. The second sequence shows how good data coverage in the training set can significantly impact performance in special or rare effector configurations. Finally, the third sequence shows a similar pattern for look-at targets, where the difference in training data can be noticed in the general posture of the character and the varying levels of robustness with respect to those targets.

## K.6 RETARGETING

This video shows that ProtoRes trained on one dataset can be successfully applied to multiple characters defined on skeletons on which our model was never trained. Currently, the model is trained on a dataset with a specific skeleton layout with specific bone offsets. In a naïve application scenario, a new skeleton would require retraining the model on a dataset for this specific skeleton, which could be a prohibitively expensive exercise. To overcome this limitation, we show in this video that the model learned using our approach can be retargeted to very different skeletons using conventional retargeting methods via a process that is fully automatic, without changing the model in any way. Thus, our model can be seamlessly used without any retraining with a wide variety of skeletons. In practice this solution allows the use of the same model on characters having similar skeleton in terms of topology and/or morphology and very different in other aspects such as proportions (bone lengths) or bone count. For instance, as our demo shows, it can be used with any humanoid character and still produces valuable results. Transfer of the model across skeletons using more sophisticated ML-based techniques such as fine-tuning, conditioning, domain adaptation is deemed to a be a fruitful future research area, for which our current results form a strong baseline.

## K.7 APPLICATION TO THE QUADRUPED SKELETON

In this video we demonstrate that ProtoRes can be applied to a completely different type of skeleton (quadruped) and it does not require any additional coding as opposed to FinalIK (no hand-crafted skeleton, effectors, bone chains, pulling, etc). It does require data and re-training however, but does not contain anything specific to humanoids. In fact, we have successfully trained the proposed model on a quadruped (dog) dataset without changing the model or hyperparameters, which demonstrates the generality of the method.

## K.8 COMPARISON TO TRANSFORMER

In this video we demonstrate qualitative examples demonstrating the basic difference between Transformer and ProtoRes. First, generally speaking, in many cases with challenging effector settings ProtoRes shows much greater robustness to outliers with respect to pose plausibility, whereas in many cases Transformer's behaviour feels like that of a pure non-learnable IK model. Second, Transformer has a tendency to generate self-intersections between the limbs to satisfy some of the constraints more closely at the expense of the overall pose plausibility. Finally, in response to an extreme look-at constraint Transformer can bend the neck behind the body, which in reality never happens in the data and would break the neck of a real human. In contrast, ProtoRes keeps the neck naturally oriented when the look-at target is unattainable.

From the theoretical perspective, the fundamental difference between the Transformer framework and the ProtoRes is as follows. ProtoRes creates one global representation of the partially specified pose and then reconstructs the full pose in one shot through a global IKD. On the other hand, Transformer generates embeddings of individual joints for the full output pose, from which a shared IKD reconstructs the local rotation angles for each joint. Transformer has all the ingredients to reproduce processing similar to our approach. However, in reality the Transformer takes a different learning route and (i) learns more localized joint predictions and (ii) learns to more strictly respect local input constraints, even at the expense of creating poses that are not statistically plausible. Combining the global (ProtoRes) and local (Transformer) approaches to derive better hybrid models seems like a promising direction for future work.

## K.9 LIMITATIONS

The final video shows examples of some specific limitations of the approach that are listed in Appendix L. Namely, we first expose specific consequences of the lack of temporal consistency in

our problem formulation, where smoothly moving an effector can cause flickering on some joints, such as the fingers. We also show how between some effector configurations, the character must be flipped completely to stay in a plausible pose, and how it's possible to place some effectors to reach an invalid pose coming from that *flip region* of the latent manifold. Finally, we showcase some problematic behaviors that can be caused by extreme look-at targets. In some cases, especially with many other constraints, ProtoRes will tend to produce a plausible pose that will not respect the look-at constraint. In other cases, the extreme look-at target may cause an unrealistic pose, e.g. by causing the character to have an impossible neck rotation. It is interesting to note how invalid poses from look-at effectors tend to happen more often than from other effector types with novice users. We hypothesize that the plausible region of a look-at target, given a current character pose, is less intuitive to grasp than for other effector types. Indeed, the current pose of the character seems to guide more precisely the placement of positional and rotational effectors than look-at effectors, leading more often to configurations outside of the training distribution for look-at effectors.

## L  LIMITATIONS

The limitations of our work can be summarized as follows:

- Constraints are not satisfied exactly, as opposed to the conventional systems. The limitation can be observed in the demo video 4_Final_IK_comparison.mp4, which compares ProtoRes and FinalIK side-by-side. For example, at seconds 18-21 we can see the chest position is not satisfied exactly by ProtoRes and the generated left foot position traverses the vicinity of the constraint as the user changes the position of the left foot. Another example like this can be seen at seconds 35-41, in which the left foot effector is not satisfied exactly by ProtoRes. It is interesting that for the rather peculiar configuration of effectors, FinalIK produces outputs that tend to severely twist the joints, perhaps beyond the capabilities of an average human. At the same time ProtoRes tends to trade the precision of following individual effectors with the plausibility of the overall pose. This is the price to pay for the ability of the model to inject the data-driven inductive bias that can be used to reconstruct pose from very sparse inputs. This could be mitigated using a conventional solver on top of the trained model. In this case, the model will produce a globally plausible pose, whereas the solver will only do the final pass to strictly satisfy certain constraints. Also, to provide additional flexibility in solving some of the constraints more strictly than the others, our model provides an effector tolerance mechanism that can help the user trade off the strictness of satisfying certain effectors vs. some others.

- Lack of temporal consistency. Our work solves the problem of creating a discrete pose. Therefore, it is limited in how it can be applied to modify an underlying smooth animation clip. For example, we can see flickering of joints (especially fingers) when effectors follow an underlying smooth animation (the finger embedding space is not smooth and has a high ambiguity). This happens to a smaller degree with the head when it is not constrained with look-at or rotation inputs.

- Exotic poses significantly deviating from the the training data distribution (a common ML/DL problem) may be hard to achieve. For example, the Lotus yoga pose is very hard to achieve with small number of effectors. Extreme or rare effector configurations may not be respected. Extreme look-at targets may not be followed or can cause artifacts in the resulting pose.

- Some effector displacement can cause a complete flip in the final pose as it makes more sense to be e.g. left-oriented or right-oriented to reach a hand position. This is normal, but we can sometimes reach "in-between" poses on the boundary of the hand effector that causes the flip, leading to weird poses

- We also noted a limitation as "aiming" poses (holding something in the hands). For example, Finger poses are generally wrong w.r.t. to a gun without additional finger constraints. It may be cumbersome to place hands + look-at for each aiming pose/angle?

- Runtime. In its current state, the model allows interactive real-time rate (about 100 FPS). This is very good for the primary application area of the model in the interactive pose design. However, the current model cannot be used for runtime applications such as driving

charachters directly in real-time games, because it would consume too high of a time budget (about 10ms, which is too much to be usable in the game runtime context).

- No contextual input is supported (text description or environment awareness), in particular for finger posing and feet collisions

- Good for realism, but might limit creativity. In particular, no bone stretching support, which is sometimes used by animators to add more expressiveness to non-realistic characters.

- The current model struggles when a large number of fine-grain controls, especially fingers are used simultaneously. Perhaps, a more structured hierarchical approach can be used to enable this functionality.