

PROTORES: PROTO-RESIDUAL NETWORK FOR POSE AUTHORING VIA LEARNED INVERSE KINEMATICS

Anonymous authors

Paper under double-blind review

ABSTRACT

Our work focuses on the development of a learnable neural representation of human pose for advanced AI assisted animation tooling. Specifically, we tackle the problem of constructing a full static human pose based on sparse and variable user inputs (*e.g.* locations and/or orientations of a subset of body joints). To solve this problem, we propose a novel neural architecture that combines residual connections with prototype encoding of a partially specified pose to create a new complete pose from the learned latent space. We show that our architecture outperforms a baseline based on Transformer, both in terms of accuracy and computational efficiency. Additionally, we develop a user interface to integrate our neural model in Unity, a real-time 3D development platform. Furthermore, we introduce two new datasets representing the static human pose modeling problem, based on high-quality human motion capture data, which will be released publicly along with model code.

1 INTRODUCTION

Modeling human pose and learning pose representations have received increasing attention recently due to their prominence in applications, such as computer graphics and animation (Harvey et al., 2020; Xu et al., 2020); immersive augmented reality (Facebook Reality Labs, 2021; Capece et al., 2018; Lin & O’Brien, 2019; Yang et al., 2021); entertainment (McDonald, 2018; Xpire, 2019); sports and wellness (Rosenhahn et al., 2008; Kim et al., 2021) as well as human machine interaction (Heindl et al., 2019; Casillas-Perez et al., 2016; Schwarz et al., 2014) and autonomous driving (Kumar et al., 2021). In the gaming industry, state-of-the-art real-time pose manipulation tools, such as CCD (Kenwright, 2012), FABRIK (Aristidou & Lasenby, 2011) or FinalIK (RootMotion, 2020), are popular for rapid execution and rely on forward and inverse kinematics models defined via non-learnable kinematic equations. *Inverse kinematics* (IK) is the process of computing the internal geometric parameters of a kinematic system resulting in the desired configuration (*e.g.* global positions) of system’s joints (Paul, 1992). *Forward kinematics* (FK) refers to the use of the kinematic equations to compute the positions of joints from specified values of internal geometric parameters. While mathematically accurate, non-learnable IK models do not guarantee that the underconstrained solutions derived from sparse constraints (*e.g.* positions of a small subset of joints) result in plausible human poses.

In this paper we develop a neural modeling approach to reconstruct full human pose from a sparse set of constraints supplied by a user, in the context of pose authoring and game development. We bridge the gap between skeleton-aware human pose representation based on IK/FK ideas and the neural embedding of human pose. Our approach effectively implements a learnable model for skeleton IK, mapping desired joint configuration into predictions of skeleton internal parameters (local rotations), learning the statistics of natural poses using datasets derived from high-quality motion capture (MOCAP) sequences. The approach, which we call ProtoRes, models the semantics of joints and their interactions using a novel prototypical residual neural network architecture. Inspired by prototypical networks, which showed that one semantic class can be represented by the prototype (mean) of a few examples (Snell et al., 2017), we extend it using a multi-block residual approach: the final pose embedding is a mean across embeddings of sparse constraints and across partial pose predictions produced in each block. We show that in terms of the pose reconstruction accuracy, ProtoRes outperforms existing gaming industry tools such as FinalIK, as well as out-of-the-box machine-learning solution based on Transformer (Vaswani et al., 2017), which also happens to be 10 times less effective in terms of training speed than the proposed architecture.

Finally, we develop user-facing tools that integrate learned ProtoRes pose representation in the Unity game engine, providing impressive qualitative examples of solutions to the problem of the AI assisted human pose authoring. The examples reveal that at the qualitative level, getting traditional workflows to behave the way ProtoRes does would require one to use many techniques in tandem, including IK, FK, layered animation pose libraries, along with procedural rigs encoding explicit heuristics. The process would be highly labor-intensive even for an experienced user while the results would still be of variable fidelity depending on the skill of the user. This is because traditional rigs have no inductive bias towards realistic poses and only allow the user to explore a limited linear latent space defined by uniform interpolation of a heuristic constraint system. ProtoRes forms a foundation that allows any junior or indie/studio user to bypass these existing complexities and create entirely new workflows for meaningfully exploring learned latent space using a familiar yet far more powerful way. We believe that our model and tools will help speed up the animation process and alleviate game artist animation skill requirements thus simplifying and democratizing game development.

1.1 BACKGROUND

We consider the full-body pose authoring animation task depicted in Fig. 1. The animator provides a few inputs, which we call *effectors*, that the target pose has to respect. For example, in Fig. 1, the look-at effector specifies that the head should be facing the orange dot, the positional effectors constrain the right foot and the right hand to be pinned to the pink dots and the rotational effector, shown in cyan, constrains the world-space rotation of the pelvis. We assume that the animator can generate arbitrary number of such effectors placed on any skeletal joint (one joint can be driven by more than one effector). The task of the model is to combine all the information provided via effectors and generate a plausible full-body pose respecting provided effector constraints. We define the full-body pose as the set of all kinematic parameters necessary to recreate the appearance of the body in 3D.

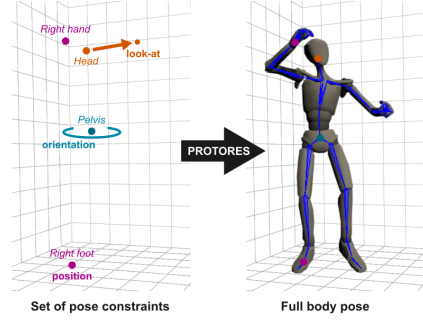


Figure 1: ProtoRes completes full human pose using a combination of 3D coordinates, look-at targets and world-space rotations specified by a user for an arbitrary subset of body joints.

Mathematically, we assume that each effector can be represented in the space $\mathbb{R}^{d_{\text{eff}}}$, where d_{eff} is taken to be maximum over all effector types. Suppose we have 3D position and 6D rotation effectors: d_{eff} is 6. In position effectors, 3 extra values are 0. We formulate the pose authoring problem as learning the mapping $\Upsilon_{\theta} : \mathbb{R}^{N \times d_{\text{eff}}} \rightarrow \mathbb{R}^{d_{\text{kin}}}$ with learnable parameters $\theta \in \Theta$. Υ_{θ} maps the input space $\mathbb{R}^{N \times d_{\text{eff}}}$ of variable effector dimensionality N (the number of effectors is not known in advance) to the space $\mathbb{R}^{d_{\text{kin}}}$, containing all kinematic parameters to reconstruct full-body pose. For example, a body with J joints can be fully defined using a tree model with 6D local rotation per joint and 3D coordinate for the root joint, in which case $d_{\text{kin}} = 6J + 3$, assuming fixed bone lengths. Given a sufficiently representative dataset $\mathcal{D} = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^M$ of poses containing pairs of inputs $\mathbf{x}_i \in \mathbb{R}^{N \times d_{\text{eff}}}$ and outputs $\mathbf{y}_i \in \mathbb{R}^{d_{\text{kin}}}$ it is viable to define the empirical risk minimization problem to learn Υ_{θ} :

$$\Upsilon_{\theta} = \arg \min_{\theta \in \Theta} \frac{1}{M} \sum_{\mathbf{x}_i, \mathbf{y}_i \in \mathcal{D}} L(\Upsilon_{\theta}(\mathbf{x}_i), \mathbf{y}_i) \quad (1)$$

1.2 RELATED WORK

Joint representations. 3D position joint representations can be used to specify a pose (Cheng et al., 2021; Cai et al., 2019; Khapugin & Grishanin, 2019). However, this approach is sub-optimal as it does not enforce fixed-length bones, nor specifies joint rotations. Bone length constraints play important role in modeling realistic poses (Pavlo et al., 2018). Joint rotations are crucial in downstream applications, such as deforming a 3D mesh on top of the skeleton, to avoid unrealistic twisting. A practical solution is to predict joint rotations, automatically satisfying bone lengths and adequately modelling rotations (Pavlo et al., 2018). This is viable via skeleton representations based on Euler angles (Han et al., 2017), rotation matrices (Zhang et al., 2018) and quaternions (Pavlo et al., 2018). In this work, we use the two-row 6D rotation matrix representation that addresses the continuity issues reminiscent of quaternion and Euler representations (Zhou et al., 2019).

Pose modeling architectures. Multi-Layer Perceptrons (MLPs) (Cho & Chen, 2014; Khapugin & Grishanin, 2019; Mirzaei et al., 2020) and kernel methods (Grochow et al., 2004; Holden et al., 2015) have been used to learn single pose representations. Beyond single pose, skeleton moving through time can be modeled as a spatio-temporal graph (Jain et al., 2016) or as a graph convolution (Yan et al., 2018; Mirzaei et al., 2020). A common limitation of these approaches is their reliance on a fixed set of inputs, whereas our architecture is specifically designed to handle sparse variable inputs.

Pose prediction from sparse constraints. Real-time methods based on nearest-neighbor search, local dynamics and motion matching have been used on sparse marker position and accelerometer data (Tautges et al., 2011; Riaz et al., 2015; Chai & Hodgins, 2005; Büttner & Clavet, 2015). MLPs and RNNs have been used for real-time processing of sparse signals such as accelerometers (Huang et al., 2018; Holden et al., 2017; Starke et al., 2020; Lee et al., 2018) and VR constraints (Lin & O’Brien, 2019; Yang et al., 2021). These approaches rely on the past pose information to disambiguate next frame prediction and as such are not applicable to our problem, in which only current pose constraints are available. Iterative IK algorithms such as FinalIK (RootMotion, 2020) have been popular in real-time applications. FinalIK works by setting up multiple IK chains for each limb of the body of a predefined human skeleton and for a fixed set of effectors. Several iterations are executed to solve each of these chains using a conventional bone chain IK method, e.g. CCD. In FinalIK, the end effector (hands and feet) can be positioned and rotated, while mid-effectors (shoulders and thighs) can only be positioned. Effectors can have a widespread effect on the body via a hand-crafted pulling mechanism that gives a different weight to each chain. This and similar tools suffer from limited realism when used for human full-body IK, as they are not data-driven. Learning-based methods strive to alleviate this by providing learned model of human pose. Grochow et al. (2004) proposed a kernel based method for learning a pose latent space in order to produce the most likely pose satisfying sparse effector constraints via online constrained optimization. The more recent commercial tool Cascadeur uses a cascade of several MLPs (each dealing with fixed set of positional effectors: 6, 16, 28) to progressively produce all joint positions without respecting bone constraints (Khapugin & Grishanin, 2019; Cascadeur, 2019). Unlike our approach, Cascadeur cannot handle arbitrary effector combinations, rotation or look-at constraints and requires post processing to respect bone constraints.

Permutation invariant architectures. Models for encoding unstructured variable inputs have been proposed in various contexts. Attention models (Bahdanau et al., 2015) and Transformer (Vaswani et al., 2017) have been proposed in the context of natural language processing. Prototypical networks (Snell et al., 2017) used average pooled embedding to encode semantic classes via a few support images in the context of few-shot image classification. Maxpool representations over variable input dimension were proposed by Qi et al. (2017) as PointNet and Zaheer et al. (2017) as DeepSets for segmentation and classification of 3D point clouds, image tagging, set anomaly detection and text concept retrieval. Niemeyer et al. (2019) further generalized the PointNet by chaining the basic maxpool/concat PointNet blocks resulting in ResPointNet architecture.

1.3 SUMMARY OF CONTRIBUTIONS

The contributions of our paper can be summarized as follows.

- We define the 3D character posing task and publicly release two associated benchmarks.
- We show that learned inverse kinematics solution can construct better poses, qualitatively and quantitatively, compared to a non-learned approach.
- We extend existing architectures with (i) semantic conditioning of joint ID and type at the input, (ii) novel residual scheme involving prototype subtraction and accumulation across blocks, as opposed to maxpool/concat daisy chain of ResPointNet, (iii) two-stage architecture with computationally efficient residual decoder that improves accuracy at smaller computational cost, as opposed to the naive final linear projection approach of PointNet and ResPointNet, and (iv) two-stage decoder design.
- We propose a novel look-at loss function.
- We propose a novel randomized weighting scheme combining randomly generated effector tolerance levels and effector noise to increase the effectiveness of multi-task training.

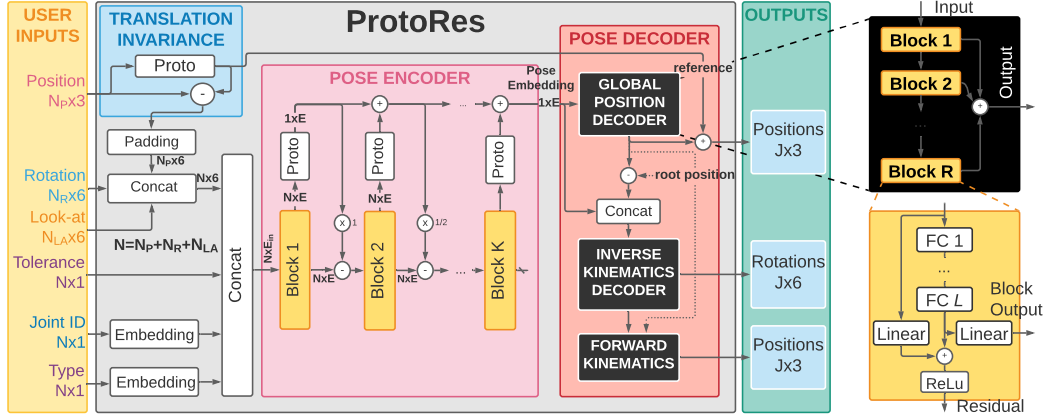


Figure 2: ProtoRes follows the encoder-decoder pattern and produces predictions in three steps. First, the variable number and type of user supplied inputs are processed for translation invariance and embedded. Second, proto-residual encoder transforms the pose specified via effectors into a pose embedding. Finally, the pose decoder expands the pose embedding into the full-body pose representation including local rotation and global position of each joint.

2 PROTORES

ProtoRes, shown in Fig. 2, follows the encoder-decoder pattern, unlike PointNet and ResPointNet that use a linear layer as a decoding mechanism (Qi et al., 2017; Niemeyer et al., 2019). The encoder has to deal with N effectors, whereas the decoder processes the collapsed representation of the pose and is therefore N times more compute efficient. Adding more decoder blocks thus results in accuracy gains at a fraction of compute cost. Below, we describe the rest of architecture in more detail.

User inputs. ProtoRes accepts position (3D coordinates), rotation (6D representation of Zhou et al. (2019)) and look-at (3D target position and 3D local facing direction) effectors. All positions are re-referenced relative to the centroid of the positional effectors to achieve translation invariance. Translation invariance simplifies the handling of poses in global space removing the need in universal reference frame, which is tricky to define. Each effector is further characterized by a positive tolerance value. Smaller tolerance implies that the effector has to be more strictly respected in the reconstructed pose. Moreover, the input includes effector semantics encoded via joint ID, an integer in $[0, J]$, and effector type, indicating a positional (0), rotational (1) or look-at (2) effector. Unlike e.g. PointNet (Qi et al., 2017) acting on dense extensive point clouds, ProtoRes acts on sparse inputs that barely provide enough information about the full pose. Therefore, it is critical to provide the semantic information on whether the given effector affects hands or feet and whether this is a positional or a rotational effector, for example. Type and joint ID variables are embedded into continuous vectors and concatenated with effector data, resulting in the encoder input, a matrix $\mathbf{x}_{in} \in \mathbb{R}^{N \times E_{in}}$ with E_{in} corresponding to the combined dimension of all embeddings plus 7D effector data and tolerance.

Pose Encoder is a two-loop residual network. The first residual loop is implemented inside each block depicted in Fig. 2 (bottom right). The second residual loop shown in Fig. 2 (left) implements the proposed Prototype-Subtract-Accumulate (PSA) residual stacking principle, which we empirically found to outperform ResPointNet’s Maxpool-Concat daisy chain proposed by Niemeyer et al. (2019). Next, we first lay out the encoder equations and then describe the motivation behind them in detail. We assume the encoder input to be $\mathbf{x}_1 = \mathbf{x}_{in} \in \mathbb{R}^{N \times E_{in}}$, omitting the batch dimension for brevity, in which case the fully-connected layer $FC_{r,\ell}$, with $\ell = 1 \dots L$, in the residual block r , $r = 1 \dots R$, with weights $\mathbf{W}_{r,\ell}$ and biases $\mathbf{a}_{r,\ell}$ can be conveniently described as $FC_{r,\ell}(\mathbf{h}_{r,\ell-1}) \equiv \text{RELU}(\mathbf{W}_{r,\ell} \mathbf{h}_{r,\ell-1} + \mathbf{a}_{r,\ell})$. The prototype layer is defined as $\text{PROTOTYPE}(\mathbf{x}) \equiv \frac{1}{N} \sum_{i=1}^N \mathbf{x}[i, :]$. The pose encoder is then described as:

$$\mathbf{x}_r = \text{RELU}[\mathbf{b}_{r-1} - 1/(r-1) \cdot \mathbf{p}_{r-1}], \quad (2)$$

$$\mathbf{h}_{r,1} = FC_{r,1}[\mathbf{x}_r], \dots, \mathbf{h}_{r,L} = FC_{r,L}[\mathbf{h}_{r,L-1}], \quad (3)$$

$$\mathbf{b}_r = \text{RELU}[\mathbf{L}_r \mathbf{x}_r + \mathbf{h}_{r,L}], \mathbf{f}_r = \mathbf{F}_r \mathbf{h}_{r,L}, \quad (4)$$

$$\mathbf{p}_r = \mathbf{p}_{r-1} + \text{PROTOTYPE}[\mathbf{f}_r]. \quad (5)$$

Equations (3) and (4) implement the MLP and the first residual loop. The proposed PSA residual mechanism, described in equations (2) and (5), is motivated by the following. First, it implements the inductive bias that the information in individual effectors is only valuable when it is different from what is already stored in the embedding of entire pose. Equation (2) implements this logic by forcing delta-mode in effectors w.r.t. to the pose embedding, \mathbf{p}_{r-1} , from the previous block, which additionally creates another residual loop that should facilitate gradient flow. Second, equation (5) collapses the forward encoding of individual effectors into the representation of the entire pose via prototype, which is known to be very effective at representing information from sparse examples (see e.g. Snell et al. (2017)). Finally, the representation of pose is accumulated across residual blocks in (5), effectively implementing skip connections from very early layers. Distant skip connections implemented via concatenation were shown to be effective in DenseNet (Huang et al., 2017). Concatenation based skipping requires additional computation and is most efficient with convolutional networks, in which kernel size can be traded for feature width while implementing distant skip connections, whereas accumulation is more compute efficient in our context.

Pose Decoder has two blocks: global position decoder (GPD) and inverse kinematics decoder (IKD). Both rely on the fully-connected residual (FCR) architecture depicted in Fig. 2 (right). GPD unrolls the pose embedding generated by encoder into the unconstrained predictions of 3D joint positions. IKD generates the internal geometric parameters (joint rotations) of the skeleton that are guaranteed to generate feasible joint positions after forward kinematics pass.

GPD accepts the encoded pose embedding, $\tilde{\mathbf{b}}_0 = \mathbf{p}_R \in \mathbb{R}^E$, and produces 3D position predictions $\tilde{\mathbf{f}}_R \in \mathbb{R}^{3J}$ of all skeletal joints using the FCR whose r -th block is described as follows:

$$\begin{aligned} \mathbf{h}_{r,1} &= \text{FC}_{r,1}^{gpd}[\tilde{\mathbf{b}}_{r-1}], \dots, \mathbf{h}_{r,L} = \text{FC}_{r,L}^{gpd}[\mathbf{h}_{r,L-1}], \\ \tilde{\mathbf{b}}_r &= \text{RELU}[\mathbf{L}_r^{gpd}\tilde{\mathbf{b}}_{r-1} + \mathbf{h}_{r,L}], \tilde{\mathbf{f}}_r = \tilde{\mathbf{f}}_{r-1} + \mathbf{F}_r^{gpd}\mathbf{h}_{r,L}. \end{aligned} \quad (6)$$

Since GPD produces predictions with no regard to skeleton constraints, its predictions do not respect bone lengths. For the IKD to provide correct rotations, the origin of the kinematic chain in world space must be given, and GPD conveniently provides the prediction of the reference (root) joint.

IKD accepts input $\hat{\mathbf{b}}_0 \in \mathbb{R}^{E+3J}$, consisting of the concatenation of the encoder-generated pose embedding, $\mathbf{p}_R \in \mathbb{R}^E$, and the output of GPD, $\tilde{\mathbf{f}}_R \in \mathbb{R}^{3J}$. Effectively, the draft pose generated by GPD, is used to condition IKD. We show in Section 3 that this additional conditioning improves accuracy. IKD predicts the 6DoF angle for each joint, $\hat{\mathbf{f}}_R \in \mathbb{R}^{6J}$, and its r -th block operates as follows:

$$\begin{aligned} \mathbf{h}_{r,1} &= \text{FC}_{r,1}^{ikd}[\hat{\mathbf{b}}_{r-1}], \dots, \mathbf{h}_{r,L} = \text{FC}_{r,L}^{ikd}[\mathbf{h}_{r,L-1}], \\ \hat{\mathbf{b}}_r &= \text{RELU}[\mathbf{L}_r^{ikd}\hat{\mathbf{b}}_{r-1} + \mathbf{h}_{r,L}], \hat{\mathbf{f}}_r = \hat{\mathbf{f}}_{r-1} + \mathbf{F}_r^{ikd}\mathbf{h}_{r,L}. \end{aligned} \quad (7)$$

Forward Kinematics (FK) pass, described in detail in Appendix A, applies skeleton kinematic equations to the local joint rotations and global root position produced by IKD. For each joint j , it produces the global transform matrix $\hat{\mathbf{G}}_j$ containing the global rotation matrix, $\hat{\mathbf{G}}_j^{13} \equiv \hat{\mathbf{G}}_j[1:3, 1:3]$, and the 3D global position, $\hat{\mathbf{g}}_j = \hat{\mathbf{G}}_j[1:3, 4]$, of the joint.

2.1 LOSSES

We use three losses to train the architecture in a multi-task fashion. The total loss combines loss terms additively with weights chosen to equalize their magnitude orders.

L2 loss penalizes the mean squared error between the prediction $\hat{\mathbf{y}}$ and the ground truth \mathbf{y} :

$$\text{MSE}(\mathbf{y}, \hat{\mathbf{y}}) = \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2. \quad (8)$$

L2 loss is used to supervise the GPD as well as the IKD output after FK pass. In the latter case it drives IKD to learn to predict rotations that lead to small position errors after FK.

Geodesic loss penalizes the errors of the IKD’s rotational outputs. It represents the smallest arc (in radians) to go from one rotation to another over the surface of a sphere. The geodesic loss is defined for the ground truth rotation matrix \mathbf{R} and its prediction $\hat{\mathbf{R}}$ as (see e.g. Salehi et al. (2018)):

$$\text{GEO}(\mathbf{R}, \hat{\mathbf{R}}) = \arccos \left[(\text{tr}(\hat{\mathbf{R}}^T \mathbf{R}) - 1)/2 \right]. \quad (9)$$

We believe that using a combination of L2 and Geodesic losses is necessary to learn a high-quality pose representation. This is especially important when the task is to reconstruct a sparsely specified pose, giving rise to multiple plausible reconstructions. We argue that a model trained to reconstruct both plausible joint positions and rotations is better equipped to solve the task accurately. Empirical evidence presented in Section 3.4 supports this intuition: a model trained on both L2 and Geodesic generalizes better on both losses than models trained only on one of those terms.

Look-at loss, proposed in this paper, enables the “look-at” feature, *i.e.* the ability to orient a joint to face a particular global position (e.g. having the head looking at a given object). It allows the model to align any direction vector $\mathbf{d}_j \in \mathbb{R}^3$ of a joint, expressed in its local frame of reference, towards a global target location \mathbf{t} . Given the predicted global transform matrix $\hat{\mathbf{G}}_j$, look-at loss is defined as:

$$\text{LAT}(\mathbf{t}, \mathbf{d}_j, \hat{\mathbf{G}}_j) = \arccos \left[\overrightarrow{(\mathbf{t} - \hat{\mathbf{g}}_j)} \cdot \hat{\mathbf{G}}_j^{13} \mathbf{d}_j \right]. \quad (10)$$

$\overrightarrow{(\mathbf{t} - \hat{\mathbf{g}}_j)}$ is a unit-length vector pointing at the target object in world space. $\hat{\mathbf{G}}_j^{13}$, when multiplied by \mathbf{d}_j , represents the global predicted look-at direction. The look-at loss trains the IKD to produce $\hat{\mathbf{G}}_j^{13}$ consistent with the look-at direction defined by \mathbf{t} and \mathbf{d}_j , both provided as network inputs.

2.2 TRAINING METHODOLOGY

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

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 (position, rotation, or look-at) are sampled from the Multinomial without replacement. This induces exponentially large number of effector type and joint permutations, resulting in strong regularizing effects and teaching the network to deal with variable inputs.

Effector tolerance and randomized loss weighting. The motivation behind randomized loss weighting is two-fold. First, we empirically find that when a random weight is multiplicatively applied to the respective loss term and its reciprocal is used as one of the network inputs for the corresponding effector, the network learns to respect the tolerance level. When exposed as a user interface feature, it lets the user control the degree of responsiveness of the model to different effectors. We also discovered that this only works when noise is added to the effector value and the standard deviation of the noise is appropriately modulated by tolerance. For example, the noise teaches the model to disregard the effector completely if the tolerance input value corresponds to the high noise variance regime. Second, we notice that the randomized weighting improves multi-task training and generalization performance. In particular, we observe significant competition between rotation and position losses on our task. The introduction of the randomized loss weighting seems to turn the competition into cooperation as our empirical results suggest in Section 3. We implement the randomized loss weighting scheme via the following steps. For each sampled effector, we uniformly sample $\Lambda \in [0, 1]$ treated as the effector tolerance. Given an effector tolerance Λ , noise with the maximum standard deviation σ_M modulated by Λ (noise models used for different effector types are described in detail in Appendix B) is added to effector data before feeding them to the neural network:

$$\sigma(\Lambda) = \sigma_M \Lambda^\eta, \quad (11)$$

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

$$W(\Lambda) = \min(W_M, 1/\sigma(\Lambda)). \quad (12)$$

Λ drives network input and $W(\Lambda)$ weighs the loss term affected by the effector.

The detailed procedure to compute the ProtoRes loss based on one batch item is presented in Algorithm 1 of Appendix C and the summary is provided below. First, we sample (i) the number of effectors and (ii) their associated types and IDs. 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 $\hat{\mathbf{f}}_{R,j}$, local joint rotations $\hat{\mathbf{R}}_j$, as well as world-space rotations and positions $\hat{\mathbf{G}}_j$ 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.

Table 1: Key quantitative results: ProtoRes vs. baselines. Lower values are better.

	miniMixamo			miniAnonymous		
	$\mathcal{L}_{gpd-L2}^{det}$	$\mathcal{L}_{ikd-L2}^{det}$	$\mathcal{L}_{loc-geo}^{det}$	$\mathcal{L}_{gpd-L2}^{det}$	$\mathcal{L}_{ikd-L2}^{det}$	$\mathcal{L}_{loc-geo}^{det}$
5-point benchmark						
FinalIK (RootMotion, 2020)	5.53e-3	8.54e-3	0.5287	3.76e-3	7.83e-3	0.5164
Masked-FCR	1.30e-3	2.49e-3	0.2607	1.11e-3	2.38e-3	0.2124
Transformer	1.10e-3	2.06e-3	0.2698	0.92e-3	1.79e-3	0.2138
ProtoRes	1.00e-3	2.02e-3	0.2534	0.76e-3	1.74e-3	0.2037
Random benchmark						
Masked-FCR	15.02e-3	35.21e-3	0.3136	1.57e-2	3.23e-2	0.2694
Transformer	1.63e-3	4.32e-3	0.2599	1.27e-3	3.49e-3	0.2006
ProtoRes	1.36e-3	4.16e-3	0.2381	0.93e-3	3.28e-3	0.1817

3 EMPIRICAL RESULTS

Our results demonstrate that (i) ProtoRes reconstructs sparsely defined pose more accurately than existing non-ML IK solution, (ii) ProtoRes is more accurate than two ML baselines, (iii) the proposed encoder-decoder design is accurate and efficient, (iv) our Prototype-Subtract-Accumulate residual scheme is more effective than the Maxpool-Concat daisy chain of Niemeyer et al. (2019), (v) two-stage GPD+IKD decoding is more effective than IKD-only decoding, (vi) the proposed randomized loss weighting improves multi-task training, (vii) joint Geodesic/L2 loss training is synergetic.

3.1 DATASETS

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).

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.

3.2 TRAINING AND EVALUATION SETUP

We use Algorithm 1 of Appendix C to sample batches of size 2048 from the training subset, hyperparameters are adjusted on the validation set. The number of effectors is sampled once per batch and is fixed for all batch items to maximize data throughput. 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 and the bone rotation accuracy. The evaluation framework tests model performance on a pre-generated set of seven files containing 6, 7...12 effectors respectively. Metrics are averaged over all files, assessing the overall quality of pose reconstruction in scenario with sparse variable inputs. Tables present results averaged over 4 random seed retries and metrics computed every 10 epochs over last 1000 epochs, rounded to the last statistically significant digit. Additional details and hyperparameter settings appear in Appendix E.



Figure 3: Qualitative posing results. Left 4 poses: adding position, look-at and rotation effectors to specify pose. Right 2 poses: achieving interesting poses with sparse constraints (4 and 7 effectors).

3.3 KEY RESULTS

To demonstrate the advantage of the proposed architecture, we perform two evaluations. First, we compare ProtoRes against two ML baselines in the random effector evaluation setup described in Section 3.2. The first baseline, 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 data and 1D tolerance) handling all effector permutations. Missing effectors are masked with $3 \cdot J$ learnable 7D placeholders. Masked-FCR has 3 encoder and 6 decoder blocks to match ProtoRes. The second baseline is based on the Transformer encoder (Vaswani et al. (2017), see Appendix G for architecture and hyperparameter settings). The bottom of Table 1, summarizing this study, shows clear advantage of ProtoRes w.r.t. both baselines. Additionally, training Transformer on NVIDIA M40 24GB GPU for 40k epochs of miniMixamo takes 1055 hours (batch size 1024 to fit Transformer in GPU memory), whereas training ProtoRes takes 106 hours. ProtoRes is clearly more compute efficient.

Second, Table 1 (top) compares ProtoRes against a non-ML IK solution FinalIK (RootMotion, 2020), as well as Transformer and Masked-FCR, on a 5-point evaluation benchmark. The 5-point benchmark tests the reconstruction of the full pose from five position effectors: chest, left and right hands, left and right feet. It is chosen, because generating the exponentially large number of FinalIK configurations to process all heterogeneous effector combinations in the random benchmark is not feasible. Note that the 5-point benchmark and random benchmark results are not directly comparable. We can see that all ML methods significantly outperform FinalIK in reconstruction accuracy, ProtoRes being the best overall. Clearly, ML methods learn the right inductive biases from the data to solve the ill-posed sparse input pose reconstruction problem, unlike the pure non-learnable IK method FinalIK.

Third, qualitative posing results are shown in Fig. 3, demonstrating that visually appealing poses can be obtained with small number of effectors (4 and 7 effectors for the two poses on the right). The left 4 poses demonstrate how pose can be refined successively by adding more effectors. Please refer to Appendix J and supplementary videos for more demonstration examples.

3.4 ABLATION STUDIES

Decoder ablation is shown in Table 2 (top). We keep all hyperparameters at defaults described in Section 3.2 and vary the number of decoder blocks (0 blocks corresponds to a simple linear projection) and compare to the ProtoRes baseline with 3 decoder blocks. We see consistent gain adding more decoder blocks at small compute cost (91, 95, 102, 106 hours of train time on miniMixamo dataset and NVIDIA M40 GPU for 0,1,2,3 blocks). Please see more detailed results in Appendix H.

Prototype-Subtract-Accumulate ablation is presented in Appendix I, in which we compare it against the ResPointNet stacking scheme (Maxpool-Concat daisy chain by Niemeyer et al. (2019)). We show that our stacking scheme is more accurate and allows stacking deeper networks gaining more accuracy, whereas stacking by Niemeyer et al. (2019) saturates at 3 encoder blocks.

GPD ablation is shown in Table 2 (middle). We remove GPD and increase IKD depth to 6 blocks to match the capacity of IKD+GPD. Comparing to the baseline, we see that GPD creates consistent gain across metrics and datasets by conditioning IKD with a draft pose.

Table 2: Ablation studies on the random benchmark. Lower values are better.

		miniMixamo			miniAnonymous		
		$\mathcal{L}_{gpd-L2}^{det}$	$\mathcal{L}_{ikd-L2}^{det}$	$\mathcal{L}_{loc-geo}^{det}$	$\mathcal{L}_{gpd-L2}^{det}$	$\mathcal{L}_{ikd-L2}^{det}$	$\mathcal{L}_{loc-geo}^{det}$
		ProtoRes baseline					
		1.36e-3	4.16e-3	0.2381	0.93e-3	3.28e-3	0.1817
Decoder	blocks	Ablation of decoder					
	0	1.54e-3	4.59e-3	0.2485	1.05e-3	3.65e-3	0.1939
	1	1.35e-3	4.34e-3	0.2433	0.93e-3	3.52e-3	0.1895
	2	1.34e-3	4.24e-3	0.2397	0.93e-3	3.34e-3	0.1840
GPD	blocks, GPD/IKD	Ablation of GPD					
\times	0/6	1.43e-3	4.39e-3	0.2413	0.93e-3	3.34e-3	0.1830
\mathcal{L}_{ikd-L2}^*	$\mathcal{L}_{loc-geo}^*$	Ablation of rotation and position loss terms					
\checkmark	\times	1.60e-3	4.49e-3	0.2742	1.12e-3	3.65e-3	0.2392
\times	\checkmark	2.04e-3	6.19e-3	0.2442	1.33e-3	4.63e-3	0.1862
W_{pos}	Randomized Loss	Ablation of randomized loss weighting					
100	\times	1.77e-3	4.93e-3	0.2549	1.15e-3	3.58e-3	0.1905
1000	\times	1.66e-3	4.75e-3	0.2668	1.09e-3	3.40e-3	0.2029

Ablation of loss terms, shown in Table 2 (middle), studies the effect of (i) removing all L2 loss terms from the output of the FK pass and (ii) removing all Geodesic loss terms from the rotation output of IKD. Interestingly, removing either of the loss terms results in the degradation of all monitored metrics on both datasets. We conclude that jointly penalizing positions with L2 and rotations with Geodesic results in positive synergetic effects and improves the overall quality of pose model.

Randomized loss weighting ablation is shown in Table 2 (bottom). The randomized loss weighting scheme (see Algorithm 1 in Appendix C) is disabled by replacing all randomized loss terms with their deterministic counterparts. For example, $\mathcal{L}_{ikd-L2}^{rnd}$ is replaced with $\sum_{j=1}^J \text{MSE}(\mathbf{g}_j, \hat{\mathbf{g}}_j)$. The inclusion of randomized weighting significantly improves generalization performance on all datasets and metrics. Additionally, when L2 weight W_{pos} increases with disabled randomized weighting, position L2 metrics improve, but at the expense of declining rotation metrics. Therefore, randomized weighting scheme contributes positive effect that cannot be achieved by tweaking the deterministic loss weights.

The limitations of the current work, discussed in detail in Appendix K, include (i) the lack of temporal consistency as we focus on the problem of authoring a discrete pose, (ii) constraints are satisfied approximately, as opposed to the more conventional systems, (iii) exotic poses significantly deviating from the training data distribution may be hard to achieve.

4 CONCLUSIONS

We define and solve the discrete full-body pose authoring task using sparse and variable user inputs. We define and release two datasets to support the development of ML models for discrete pose authoring and animation. We propose ProtoRes, a novel ML architecture which processes a variable number of heterogeneous user inputs (position, angle, direction) to reconstruct a full-body pose. We compare ProtoRes against two strong ML baselines, Masked-FCR and Transformer, showing superior results for ProtoRes, both in terms of accuracy and computational efficiency. We also show that ML models reconstruct full-body poses from sparse user inputs more accurately than existing non-learnable inverse kinematics models. We develop a suite of UI tools for the integration of our model in Unity and provide demos showing how our model can be used effectively to solve the discrete pose authoring problem by the end user. Our results have a few implications. First, our ML based tools will have positive impacts on the simplification and democratization of the game development process by helping a wide audience materialize their creative animation ideas. Second, our novel approach to neural pose representation could be applied in a variety of tasks where efficient and accurate reconstruction of full-body poses from noisy intermittent measurements is important.

REFERENCES

- Adobe Inc. Adobe general terms of use. <https://www.adobe.com/legal/terms.html>, 2020. Accessed: 2021-05-06.
- Andreas Aristidou and Joan Lasenby. Fabrik: A fast, iterative solver for the inverse kinematics problem. *Graphical Models*, 73(5):243–260, 2011.
- Dzmitry Bahdanau, Kyung Hyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. In *ICLR*, January 2015.
- Michael Büttner and Simon Clavet. Motion matching - the road to next gen animation. In *Proc. of Nucl.ai 2015*, 2015. URL https://www.youtube.com/watch?v=z_wpgHFSWss&t=658s.
- Yujun Cai, Lihao Ge, Jun Liu, Jianfei Cai, Tat-Jen Cham, Junsong Yuan, and Nadia Magnenat Thalmann. Exploiting spatial-temporal relationships for 3d pose estimation via graph convolutional networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.
- Nicola Capece, Ugo Erra, and Giuseppe Romaniello. A low-cost full body tracking system in virtual reality based on microsoft kinect. In Lucio Tommaso De Paolis and Patrick Bourdot (eds.), *Augmented Reality, Virtual Reality, and Computer Graphics*, pp. 623–635. Springer International Publishing, 2018.
- Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm (eds.), *ECCV 2020*, pp. 213–229, 2020.
- Cascadeur. How to use deep learning in character posing. <https://cascadeur.com/ru/blog/general/how-to-use-deep-learning-in-character-posing>, 2019. Accessed: 2021-05-06.
- David Casillas-Perez, Javier Macias-Guarasa, Marta Marron-Romera, David Fuentes-Jimenez, and Alvaro Fernandez-Rincon. Full body gesture recognition for human-machine interaction in intelligent spaces. In Francisco Ortuño and Ignacio Rojas (eds.), *Bioinformatics and Biomedical Engineering*, pp. 664–676. Springer International Publishing, 2016.
- Jinxiang Chai and Jessica K Hodgins. Performance animation from low-dimensional control signals. In *ACM Transactions on Graphics (ToG)*, volume 24, pp. 686–696. ACM, 2005.
- Yu Cheng, Bo Wang, Bo Yang, and Robby T Tan. Graph and temporal convolutional networks for 3d multi-person pose estimation in monocular videos. *AAAI*, 2021.
- Kyunghyun Cho and Xi Chen. Classifying and visualizing motion capture sequences using deep neural networks. In *2014 International Conference on Computer Vision Theory and Applications (VISAPP)*, volume 2, pp. 122–130. IEEE, 2014.
- Facebook Reality Labs. Inside facebook reality labs: Research updates and the future of social connection. <https://tech.fb.com/inside-facebook-reality-labs-research-updates-and-the-future-of-social-connection/>, 2021. Accessed: 2021-04-30.
- Keith Grochow, Steven L Martin, Aaron Hertzmann, and Zoran Popović. Style-based inverse kinematics. In *ACM SIGGRAPH 2004 Papers*, pp. 522–531. 2004.
- Fei Han, Brian Reily, William Hoff, and Hao Zhang. Space-time representation of people based on 3d skeletal data: A review. *Computer Vision and Image Understanding*, 158:85–105, 2017.
- Félix G. Harvey, Mike Yurick, Derek Nowrouzezahrai, and Christopher Pal. Robust motion in-betweening. 39(4), 2020.
- Christoph Heindl, Markus Ikeda, Gernot Stübl, Andreas Pichler, and Josef Scharinger. Metric pose estimation for human-machine interaction using monocular vision. *ArXiv*, 2019.

- Daniel Holden, Jun Saito, and Taku Komura. Learning an inverse rig mapping for character animation. In *Proceedings of the 14th ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, pp. 165–173, 2015.
- Daniel Holden, Taku Komura, and Jun Saito. Phase-functioned neural networks for character control. *ACM Transactions on Graphics (TOG)*, 36(4):1–13, 2017.
- Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks. In *CVPR*. IEEE Computer Society, 2017.
- Yinghao Huang, Manuel Kaufmann, Emre Aksan, Michael J Black, Otmar Hilliges, and Gerard Pons-Moll. Deep inertial poser: Learning to reconstruct human pose from sparse inertial measurements in real time. *ACM Transactions on Graphics (TOG)*, 37(6):1–15, 2018.
- Ashesh Jain, Amir R Zamir, Silvio Savarese, and Ashutosh Saxena. Structural-rnn: Deep learning on spatio-temporal graphs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5308–5317, 2016.
- Ben Kenwright. Inverse kinematics–cyclic coordinate descent (ccd). *Journal of Graphics Tools*, 16(4):177–217, 2012.
- Evgeniy Khapugin and Alexander Grishanin. Physics-based character animation with cascadeur. In *ACM SIGGRAPH 2019 Studio*, SIGGRAPH ’19, New York, NY, USA, 2019. Association for Computing Machinery.
- Hyoungun Kim, Abhaysinh Zala, Graham Burri, and M. Bansal. Fixmypose: Pose correctional captioning and retrieval. In *AAAI*, 2021.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2015.
- Chandan Kumar, Jayanth Ramesh, Bodhisattwa Chakraborty, Renjith Raman, Christoph Weinrich, Anurag Mundhada, Jain. Arjun, and Fabian B Flohr. VRU Pose-SSD: Multiperson pose estimation for automated driving. In *AAAI 2021*, 2021.
- Kyungho Lee, Seyoung Lee, and Jehee Lee. Interactive character animation by learning multi-objective control. In *SIGGRAPH Asia 2018 Technical Papers*, pp. 180. ACM, 2018.
- James Lin and James O’Brien. Temporal ik: Data-driven pose estimation for virtual reality. 2019.
- Kyle McDonald. Dance x machine learning: First steps. <https://medium.com/@kcimc/discrete-figures-7d9e9c275c47>, 2018. Accessed: 2021-05-03.
- Maryam Sadat Mirzaei, Kourosh Meshgi, Etienne Frigo, and Toyoaki Nishida. Animgan: A spatiotemporally-conditioned generative adversarial network for character animation. In *2020 IEEE International Conference on Image Processing (ICIP)*, pp. 2286–2290. IEEE, 2020.
- Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Occupancy flow: 4d reconstruction by learning particle dynamics. In *ICCV*, October 2019.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS 2019*, pp. 8024–8035, 2019.
- R.P. Paul. *Robot Manipulators: Mathematics, Programming, and Control : The Computer Control of Robot Manipulators*. The MIT Press Series in Artificial Intelligence. MIT Press, 1992.
- Dario Pavllo, David Grangier, and Michael Auli. QuaterNet: A quaternion-based recurrent model for human motion. In *British Machine Vision Conference (BMVC)*, 2018.
- C. Qi, Hao Su, Kaichun Mo, and L. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. *CVPR*, pp. 77–85, 2017.

- Kaiser Riaz, Guan hong Tao, Björn Krüger, and Andreas Weber. Motion reconstruction using very few accelerometers and ground contacts. *Graphical Models*, 79:23–38, 2015.
- RootMotion. Advanced character animation systems for Unity. <http://root-motion.com>, 2020. Accessed: 2021-04-30.
- Bodo Rosenhahn, Christian Schmaltz, Thomas Brox, Joachim Weickert, Daniel Cremers, and Hans-Peter Seidel. Markerless motion capture of man-machine interaction. In *CVPR*, pp. 1–8, 2008.
- Seyed Sadegh Mohseni Salehi, Shadab Khan, Deniz Erdogmus, and Ali Gholipour. Real-time deep pose estimation with geodesic loss for image-to-template rigid registration. *IEEE transactions on medical imaging*, 38(2):470–481, 2018.
- Julia Schwarz, Charles Claudius Marais, Tommer Leyvand, Scott E. Hudson, and Jennifer Mankoff. Combining body pose, gaze, and gesture to determine intention to interact in vision-based interfaces. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI ’14, pp. 3443–3452, New York, NY, USA, 2014. Association for Computing Machinery.
- Jake Snell, Kevin Swersky, and Richard S. Zemel. Prototypical networks for few-shot learning. In *NIPS*, pp. 4080–4090, 2017.
- Sebastian Starke, Yiwei Zhao, Taku Komura, and Kazi Zaman. Local motion phases for learning multi-contact character movements. *ACM Transactions on Graphics (TOG)*, 39(4):54–1, 2020.
- Jochen Tautges, Arno Zinke, Björn Krüger, Jan Baumann, Andreas Weber, Thomas Helten, Meinard Müller, Hans-Peter Seidel, and Bernd Eberhardt. Motion reconstruction using sparse accelerometer data. *ACM Transactions on Graphics (ToG)*, 30(3):18, 2011.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *NeurIPS*, volume 30, 2017.
- Xpire. Using ai to make nba players dance. <https://tinyurl.com/y3bdj5p5>, 2019. Accessed: 2021-05-03.
- Zhan Xu, Yang Zhou, Evangelos Kalogerakis, Chris Landreth, and Karan Singh. Rignet: Neural rigging for articulated characters. *ACM Trans. on Graphics*, 39, 2020.
- Sijie Yan, Yuanjun Xiong, and Dahua Lin. Spatial temporal graph convolutional networks for skeleton-based action recognition. In *AAAI 2018*, volume 32, 2018.
- Dongseok Yang, Doyeon Kim, and Sung-Hee Lee. Real-time lower-body pose prediction from sparse upper-body tracking signals. *arXiv preprint arXiv:2103.01500*, 2021.
- Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov, and Alexander J Smola. Deep sets. In *NeurIPS*, volume 30. Curran Associates, Inc., 2017.
- He Zhang, Sebastian Starke, Taku Komura, and Jun Saito. Mode-adaptive neural networks for quadruped motion control. *ACM Transactions on Graphics (TOG)*, 37(4):1–11, 2018.
- Yi Zhou, Connelly Barnes, Lu Jingwan, Yang Jimei, and Li Hao. On the continuity of rotation representations in neural networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

Supplementary Material for ProtoRes: Proto-Residual Network for Pose Authoring via Learned Inverse Kinematics

Table of Contents

A	Architecture Details	14
B	Effector noise model	15
B.1	Position effector noise model	15
B.2	Rotation effector noise model	15
B.3	Look-at effector noise model	15
C	Training Methodology: Details	16
D	Datasets: Details	17
E	Training Setup: Details	18
F	Masked-FCR baseline architecture	19
G	Transformer baseline architecture	19
H	Ablation of the Decoder: Details	20
I	Ablation of the Prototype-Subtract-Accumulate stacking principle	21
J	Videos and demonstrations	22
J.1	ProtoRes Demo	22
J.2	Posing from Images	22
J.3	Loss Ablation	22
J.4	FinalIK Comparison	22
J.5	Datasets Comparison	23
J.6	Limitations	23
K	Limitations	23

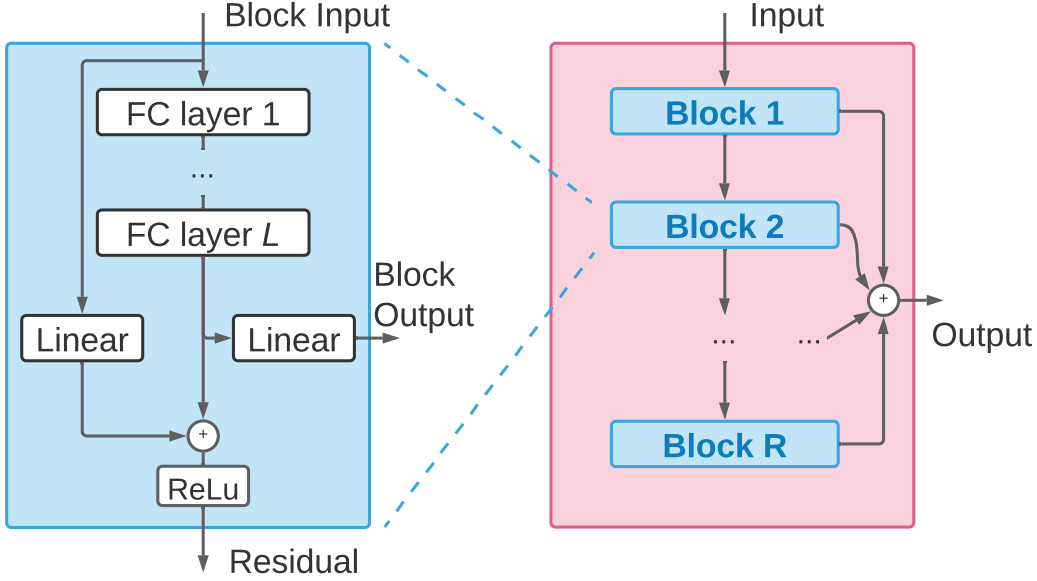


Figure 4: Block diagram of the fully-connected residual (FCR) decoder architecture. Left: the diagram of one residual block of the FCR decoder. Note that the basic residual block of the encoder architecture is exactly the same. Right: residual blocks connected in the FCR architecture.

A ARCHITECTURE DETAILS

Encoder The residual block depicted in Fig. 4 (left) is used as the basic building block of the ProtoRes encoder.

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 $\mathbf{o}_j = [o_{x,j}, o_{y,j}, o_{z,j}]^T$ 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 $\alpha_j, \beta_j, \gamma_j$:

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

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

$$\hat{\mathbf{r}}_{j,x} = \overrightarrow{\hat{\mathbf{f}}_{R,j}[1:3]}, \quad \hat{\mathbf{r}}_{j,z} = \overrightarrow{\hat{\mathbf{r}}_{j,x} \times \hat{\mathbf{f}}_{R,j}[4:6]}, \quad \hat{\mathbf{r}}_{j,y} = \hat{\mathbf{r}}_{j,z} \times \hat{\mathbf{r}}_{j,x}, \quad \hat{\mathbf{R}}_j = [\hat{\mathbf{r}}_{j,x} \ \hat{\mathbf{r}}_{j,y} \ \hat{\mathbf{r}}_{j,z}]. \quad (13)$$

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 :

$$\widehat{\mathbf{G}}_j = \widehat{\mathbf{G}}_{p(j)} \begin{bmatrix} \widehat{\mathbf{R}}_j & \mathbf{o}_j \\ \mathbf{0} & 1 \end{bmatrix}. \quad (14)$$

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 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:

$$\mathbf{x}[i, 1:3] = \mathbf{g}_{I_i} + \sigma(\Lambda_i)\boldsymbol{\varepsilon}_i; \quad \mathbf{x}[i, 4:6] = 0. \quad (15)$$

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. equation 19 and $\boldsymbol{\varepsilon}_i$ 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 13. Second, we generate the random 3D Euler angles vector $\boldsymbol{\varepsilon}_i$ from the zero-mean Gaussian distribution $\mathcal{N}(\mathbf{0}, \sigma(\Lambda_i)\mathbf{I})^1$ and convert it to the random rotation matrix Ψ_i using eq. equation 16:

$$\Psi_i = \begin{bmatrix} \cos \boldsymbol{\varepsilon}_i[1] & -\sin \boldsymbol{\varepsilon}_i[1] & 0 \\ \sin \boldsymbol{\varepsilon}_i[1] & \cos \boldsymbol{\varepsilon}_i[1] & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \cos \boldsymbol{\varepsilon}_i[2] & 0 & \sin \boldsymbol{\varepsilon}_i[2] \\ 0 & 1 & 0 \\ -\sin \boldsymbol{\varepsilon}_i[2] & 0 & \cos \boldsymbol{\varepsilon}_i[2] \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos \boldsymbol{\varepsilon}_i[3] & -\sin \boldsymbol{\varepsilon}_i[3] \\ 0 & \sin \boldsymbol{\varepsilon}_i[3] & \cos \boldsymbol{\varepsilon}_i[3] \end{bmatrix}. \quad (16)$$

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

$$\mathbf{x}[i, 1:3] = \mathbf{G}_{I_i}^{13'}[:, 1], \quad \mathbf{x}[i, 4:6] = \mathbf{G}_{I_i}^{13'}[:, 2]. \quad (17)$$

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 formed by the ground truth 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_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 $\mathbf{t}_i = \mathbf{g}_{I_i} + d_t \mathbf{d}_i + \sigma(\Lambda_i)\boldsymbol{\varepsilon}_i$. Finally, the output is constructed as follows:

$$\mathbf{x}[i, 1:3] = \mathbf{t}_i, \quad \mathbf{x}[i, 4:6] = \mathbf{d}_i. \quad (18)$$

As previously, $\boldsymbol{\varepsilon}_i$ 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: $\mathbf{R}_j, \mathbf{G}_j; N \sim \text{UNIFORM}[3, 16]$ \triangleright Ground truth for all joints $j \in [0, J]$; number of effectors

Ensure: \mathbf{x} \triangleright Sample inputs

$I_1, \dots, I_N \leftarrow \text{MULTINOMIAL}(\{0, \dots, J-1\}, N)$ \triangleright Effector IDs

$T_1, \dots, T_N \leftarrow \text{MULTINOMIAL}(\{0, 1, 2\}, N)$ \triangleright Effector type

for i in $1 \dots N$ **do**

$\Lambda_i \leftarrow \text{UNIFORM}[0, 1]$ \triangleright Effector tolerance

$\sigma(\Lambda_i); W(\Lambda_i) \leftarrow \sigma_M \Lambda_i^\eta; \min(W_M, 1/\sigma(\Lambda_i))$ \triangleright Effector noise std and weight

$\mathbf{x}[i, :] \leftarrow \text{NOISEMODEL}(\mathbf{G}_{I_i}, \sigma(\Lambda_i), T_i)$ \triangleright Generate noisy effector

end for

Predict: $\hat{\mathbf{f}}_{R,j}, \hat{\mathbf{R}}_j, \hat{\mathbf{G}}_j \quad \forall j$

$\mathcal{L}_{gpd-L2}^{rnd} \leftarrow \frac{1}{\sum_{i=1}^N \mathbb{1}_{T_i=0} W(\Lambda_i)} \sum_{i=1}^N \mathbb{1}_{T_i=0} W(\Lambda_i) \text{MSE}(\mathbf{g}_{I_i}, \hat{\mathbf{f}}_{R,I_i})$ \triangleright Randomized GPD position loss

$\mathcal{L}_{ikd-L2}^{rnd} \leftarrow \frac{1}{\sum_{i=1}^N \mathbb{1}_{T_i=0} W(\Lambda_i)} \sum_{i=1}^N \mathbb{1}_{T_i=0} W(\Lambda_i) \text{MSE}(\mathbf{g}_{I_i}, \hat{\mathbf{g}}_{I_i})$ \triangleright Randomized IKD position loss

$\mathcal{L}_{gpd-L2}^{det} \leftarrow \sum_{j=1}^J \text{MSE}(\mathbf{g}_j, \hat{\mathbf{f}}_{R,j})$ \triangleright Deterministic GPD position loss

$\mathcal{L}_{ikd-L2}^{det} \leftarrow \sum_{j=1}^J \text{MSE}(\mathbf{g}_j, \hat{\mathbf{g}}_j)$ \triangleright Deterministic IKD position loss

$\mathcal{L}_{loc-geo}^{det} \leftarrow \sum_{j=1}^J \text{GEO}(\mathbf{R}_j, \hat{\mathbf{R}}_j)$ \triangleright Deterministic local rotation loss

$\mathcal{L}_{glob-geo}^{rnd} \leftarrow \frac{1}{\sum_{i=1}^N \mathbb{1}_{T_i=1} W(\Lambda_i)} \sum_{i=1}^N \mathbb{1}_{T_i=1} W(\Lambda_i) \text{GEO}(\mathbf{G}_{I_i}^{13}, \hat{\mathbf{G}}_{I_i}^{13})$ \triangleright Randomized global rotation loss

$\mathcal{L}_{lat}^{det} \leftarrow \frac{1}{\sum_{i=1}^N \mathbb{1}_{T_i=2}} \sum_{i=1}^N \mathbb{1}_{T_i=2} \text{LAT}(\mathbf{x}[i, 1:3], \mathbf{x}[i, 4:6], \hat{\mathbf{G}}_{I_i}^{13})$ \triangleright Randomized Look-at loss

$\mathcal{L} \leftarrow \frac{W_{pos}}{J} (\mathcal{L}_{gpd-L2}^{rnd} + \mathcal{L}_{ikd-L2}^{rnd} + \mathcal{L}_{gpd-L2}^{det} + \mathcal{L}_{ikd-L2}^{det}) + \frac{1}{J} (\mathcal{L}_{lat}^{det} + \mathcal{L}_{glob-geo}^{rnd} + \mathcal{L}_{loc-geo}^{det})$ \triangleright Total loss

C TRAINING 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

¹Note 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 synergy 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 Λ , noise (noise models used for different effector types are described in detail in Appendix B) with variance proportional to Λ is added to effector data before feeding them to the neural network:

$$\sigma(\Lambda) = \sigma_M \Lambda^\eta. \quad (19)$$

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

$$W(\Lambda) = \min(W_M, 1/\sigma(\Lambda)). \quad (20)$$

Λ 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 $\tilde{\mathbf{f}}_{R,j}$, local joint rotations $\hat{\mathbf{R}}_j$, as well as world-space rotations and positions $\hat{\mathbf{G}}_j$ 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

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.

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, η	13	13
Dropout	0.01	[0.0, 0.01, 0.05, 0.1, 0.2]
Embedding dimensionality	32	[16, 32, 64, 128]
Augmentation	mirror, rotation	[mirror, rotation, translation]

Table 3: Settings of ProtoRes hyperparameters and the hyperparameter search grid.

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 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.

E TRAINING 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.

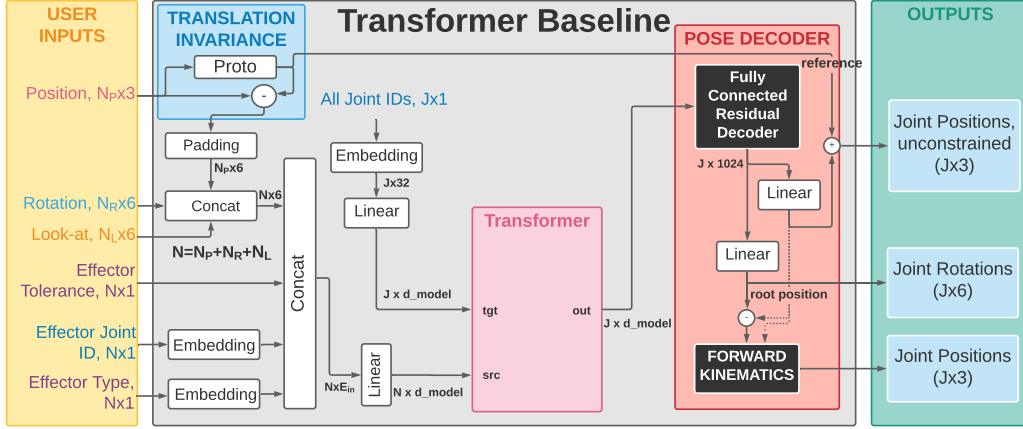


Figure 5: Block diagram of the Transformer baseline architecture.

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 and the bone rotation accuracy. The evaluation framework tests model performance on a pre-generated set of seven files containing 6 to 12 effectors each. 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 500 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 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 3.

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. 5.

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

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

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.

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 transformer is then decoded to unconstrained position and rotation outputs using two-headed Fully-Connected residual stack (this is the same architecture as the one used in ProtoRes decoders). Note that this is a classical Transformer application scheme that has recently been used to achieve SOTA results in object detection, for example Carion et al. (2020). Table 4 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 3 are listed.

H ABLATION OF THE DECODER: DETAILS

The detailed results of the decoder ablation are shown in Table 5. 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 5: 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.

Encoder blocks	Decoder blocks	miniMixamo			Train time, h	miniAnonymous			Train time, h
		$\mathcal{L}_{gpd-L2}^{det}$	$\mathcal{L}_{ikd-L2}^{det}$	$\mathcal{L}_{loc-geo}^{det}$		$\mathcal{L}_{gpd-L2}^{det}$	$\mathcal{L}_{ikd-L2}^{det}$	$\mathcal{L}_{loc-geo}^{det}$	
3	0	1.54e-3	4.59e-3	0.2485	91	1.05e-3	3.65e-3	0.1939	
3	1	1.35e-3	4.34e-3	0.2433	95	0.93e-3	3.52e-3	0.1895	
3	2	1.34e-3	4.24e-3	0.2397	102	0.93e-3	3.34e-3	0.1840	
3	3	1.36e-3	4.16e-3	0.2381	106	0.93e-3	3.28e-3	0.1817	
5	0	1.27e-6	4.20e-3	0.2399	144	0.84e-3	3.27e-3	0.1824	
5	1	1.28e-3	4.30e-3	0.2390	148	0.81e-3	3.24e-3	0.1818	
5	2	1.15e-3	4.02e-3	0.2345	153	0.77e-3	3.10e-3	0.1791	
5	3	1.18e-3	4.03e-3	0.2351	157	0.82e-3	3.07e-3	0.1785	
7	0	1.13e-3	4.00e-3	0.2355	196	0.74e-3	2.98e-3	0.1762	
7	1	1.23e-3	4.16e-3	0.2356	200	0.79e-3	3.07e-3	0.1780	
7	2	1.13e-3	4.51e-3	0.2383	205	0.78e-3	3.10e-3	0.1780	
7	3	1.15e-3	3.98e-3	0.2352	209	0.82e-3	3.14e-3	0.1783	

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 modify our encoder 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 the same decoder with the MCDC encoder to make sure we exactly match the overall architecture capacity with that of ProtoRes. However, note that decoder is not part of the original design by Niemeyer et al. (2019) and without it, the overall performance degrades further. Here we focus exclusively on the effects of stacking within the encoder.
- Hyperparameters of both architectures are taken from Table 3 in Appendix C

We show that our stacking scheme is more accurate and allows stacking deeper networks effectively. Quantitative results are shown in Table 6. It is clear that the proposed stacking approach provides gain in setups with varying number of encoder blocks. More importantly, when stacking more blocks our approach provides more gain, whereas for the Maxpool-Concat daisy chain approach of (Niemeyer et al., 2019), the gain saturates at 3 blocks and adding more blocks beyond this does not help. Moreover, we can see on the second dataset, miniAnonymous, that become worse, this is the result of less stable convergence of the MCDC stacking. In some cases with large number of encoder blocks, we observe that the MCDC approach converges to a noticeably worse generalization performance. We attribute this to convergence issues, since for larger number of blocks, MCDC stacking approach results in both higher generalization errors and training losses. The proposed PSA approach is significantly more robust when used with deeper architectures as well as offers better generalization performance in general. We conclude that the proposed approach is more effective,

Table 6: Ablation of the Prototype-Subtract-Accumulate. Random benchmark, lower values are better.

Stacking scheme	Encoder blocks	miniMixamo			miniAnonymous		
		$\mathcal{L}_{gpd-L2}^{det}$	$\mathcal{L}_{ikd-L2}^{det}$	$\mathcal{L}_{loc-geo}^{det}$	$\mathcal{L}_{gpd-L2}^{det}$	$\mathcal{L}_{ikd-L2}^{det}$	$\mathcal{L}_{loc-geo}^{det}$
MCDC	3	1.39e-3	4.27e-3	0.2403	0.92e-3	3.23e-3	0.1818
MCDC	5	1.30e-3	4.19e-3	0.2394	0.88e-3	3.22e-3	0.1809
MCDC	7	1.32e-3	4.26e-3	0.2399	0.99e-3	3.49e-3	0.1864
PSA (ours)	3	1.36e-3	4.16e-3	0.2381	0.91e-3	3.19e-3	0.1810
PSA (ours)	5	1.18e-3	4.03e-3	0.2351	0.82e-3	3.07e-3	0.1785
PSA (ours)	7	1.15e-3	3.98e-3	0.2352	0.82e-3	3.14e-3	0.1783

because it has higher accuracy for computationally efficient configuration with 3 blocks and it results in even more accuracy gain when more encoder blocks are stacked together.

J VIDEOS AND DEMONSTRATIONS

We provide here descriptions related to each video found in the supplementary materials. 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 our demonstrations, described below in Appendix J.5 qualitatively shows some of the most severe impacts of training on ablated datasets.

J.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.

J.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.

J.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.

J.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.

J.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.

J.6 LIMITATIONS

The final video shows examples of some specific limitations of the approach that are listed in Appendix K. 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.

K LIMITATIONS

The limitations of our work can be summarized as follows:

- Constraints are not satisfied exactly, as opposed to the conventional systems. 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 characters 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
- Single skeleton layout with specific bone offsets is currently supported. A new skeleton requires either a re-trained model, or a retargeting pass, which may be expensive and has the potential to reduce the realism of the reconstructed pose.
- The approach relies on MOCAP data. This type of data may be hard to obtain for certain characters, for example an octopus.
- 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.