

Supplementary Material:

Residual Force Control for Agile Human Behavior Imitation and Extended Motion Synthesis

1 Symbol Tables

Table 1: Important notations used in the main paper.

| Notation | Description |
|--|--|
| \mathbf{s} | state |
| \mathbf{a} | action |
| r | reward |
| r^{im} | motion imitation reward |
| r^{reg} | regularizing reward of residual forces |
| γ | discount factor |
| \mathbf{x} | humanoid state, $\mathbf{x} = (\mathbf{q}, \dot{\mathbf{q}})$ |
| $\hat{\mathbf{x}}$ | reference humanoid state, $\mathbf{x} = (\hat{\mathbf{q}}, \hat{\dot{\mathbf{q}}})$ |
| $\mathbf{x}_{0:T}$ | humanoid motion $\mathbf{x}_{0:T} = (\mathbf{x}_0, \dots, \mathbf{x}_{T-1})$ |
| $\hat{\mathbf{x}}_{0:T}$ | reference motion $\hat{\mathbf{x}}_{0:T} = (\hat{\mathbf{x}}_0, \dots, \hat{\mathbf{x}}_{T-1})$ |
| \mathbf{q} | humanoid DoFs, $\mathbf{q} = (\mathbf{q}_r, \mathbf{q}_{\text{nr}})$ |
| \mathbf{q}_r | humanoid root DoFs (global position and orientation) |
| \mathbf{q}_{nr} | humanoid non-root DoFs |
| $\dot{\mathbf{q}}$ | joint velocities |
| $\ddot{\mathbf{q}}$ | joint accelerations |
| \mathbf{u} | target joint angles of PD controllers |
| $\boldsymbol{\tau}$ | joint torques computed from PD control |
| $\mathbf{k}_p, \mathbf{k}_d$ | PD controller gains |
| $\tilde{\mathbf{a}}$ | corrective action (residual forces) |
| $\mathcal{T}(\mathbf{s}_{t+1} \mathbf{s}_t, \mathbf{a}_t)$ | original humanoid dynamics |
| $\tilde{\mathcal{T}}(\mathbf{s}_{t+1} \mathbf{s}_t, \mathbf{a}_t, \tilde{\mathbf{a}}_t)$ | RFC-based humanoid dynamics |
| $\pi_\theta(\mathbf{a}_t \mathbf{s}_t)$ | original humanoid control policy |
| $\tilde{\pi}_\theta(\mathbf{a}_t, \tilde{\mathbf{a}}_t \mathbf{s}_t)$ | RFC-based composite policy |
| $\tilde{\pi}_{\theta_1}(\mathbf{a}_t \mathbf{s}_t)$ | humanoid control policy same as $\pi_\theta(\mathbf{a}_t \mathbf{s}_t)$ |
| $\tilde{\pi}_{\theta_2}(\tilde{\mathbf{a}}_t \mathbf{s}_t)$ | residual force policy |
| $\boldsymbol{\xi}_j$ | residual force vector |
| \mathbf{e}_j | the contact point of $\boldsymbol{\xi}_j$ |
| \mathbf{h}_i | contact force vector determined by simulation |
| \mathbf{v}_i | the contact point of \mathbf{h}_i |
| $\mathbf{J}_{\mathbf{e}_j}$ | Jacobian matrix $d\mathbf{e}_j/d\mathbf{q}$ |
| $\mathbf{J}_{\mathbf{v}_i}$ | Jacobian matrix $d\mathbf{v}_i/d\mathbf{q}$ |
| $\mathbf{B}(\mathbf{q})$ | inertial matrix |
| $\mathbf{C}(\mathbf{q}, \dot{\mathbf{q}})$ | matrix of Coriolis and centrifugal terms |
| $\mathbf{g}(\mathbf{q})$ | gravity vector |
| $\boldsymbol{\eta}$ | total joint torques $\sum \mathbf{J}_{\mathbf{e}_j}^T \boldsymbol{\xi}_j$ of residual forces, $\boldsymbol{\eta} = (\boldsymbol{\eta}_r, \boldsymbol{\eta}_{\text{nr}})$ |
| $\boldsymbol{\eta}_r$ | torques of root DoFs from residual forces |
| $\boldsymbol{\eta}_{\text{nr}}$ | torques of non-root DoFs from residual forces |
| \mathbf{z} | latent variable for human intent |
| $\kappa_\psi(\mathbf{x}_{t:t+f} \mathbf{x}_{t-p:t}, \mathbf{z})$ | kinematic policy (decoder distribution) in CVAE |
| $q_\phi(\mathbf{z} \mathbf{x}_{t-p:t}, \mathbf{x}_{t:t+f})$ | approximate posterior (encoder distribution) in CVAE |

2 Motion Imitation Reward

In the following, we give a detailed definition of the motion imitation reward r_t^{im} introduced in Sec. 3 of the main paper, which is used to encourage the humanoid motion $\mathbf{x}_{0:T}$ generated by the policy to match the reference motion $\hat{\mathbf{x}}_{0:T}$. We use two types of imitation reward r_t^{im} depending on the length of the motion:

- (1) World coordinate reward r_t^{world} , which is the reward used in DeepMimic [3]. It has proven to be effective for imitating short clips ($< 5\text{s}$) of locomotions, but not suitable for long-term motions due to global position drifts as pointed out in [6]. We still use this reward for imitating short clips of locomotions (e.g., acrobatics) to have a fair comparison with DeepMimic.
- (2) Local coordinate reward r_t^{local} , which is the reward used in EgoPose [6]. It is more robust to global position drifts and we use it to imitate longer motion clips (e.g., ballet dance). For extended motion synthesis, we also use r_t^{local} to imitate the output motion of the kinematic policy κ_ψ .

Note that for the same reference motion, we always use the same motion imitation reward for both our RFC models and DeepMimic for a fair comparison.

2.1 World Coordinate Reward

As defined in DeepMimic [3], the world coordinate reward r_t^{world} consists of four sub-rewards:

$$r_t^{\text{world}} = w_p r_t^p + w_v r_t^v + w_e r_t^e + w_c r_t^c, \quad (1)$$

where w_p, w_v, w_e, w_c are weighting factors, which we set to (0.3, 0.1, 0.5, 0.1).

The pose reward r_t^p measures the mismatch between joint DoFs \mathbf{q}_t and the reference $\hat{\mathbf{q}}_t$ for non-root joints. We use \mathbf{b}_t^j and $\hat{\mathbf{b}}_t^j$ to denote the local orientation quaternion of joint j computed from \mathbf{q}_t and $\hat{\mathbf{q}}_t$ respectively. We use $\mathbf{b}_1 \ominus \mathbf{b}_2$ to denote the relative quaternion from \mathbf{b}_2 to \mathbf{b}_1 , and $\|\mathbf{b}\|$ to compute the rotation angle of \mathbf{b} .

$$r_t^p = \exp \left[-\alpha_p \left(\sum_j \|\mathbf{b}_t^j \ominus \hat{\mathbf{b}}_t^j\|^2 \right) \right]. \quad (2)$$

The velocity reward r_t^v measures the difference between joint velocities $\dot{\mathbf{q}}_t$ and the reference $\hat{\dot{\mathbf{q}}}_t$. The reference velocity $\hat{\dot{\mathbf{q}}}_t$ is computed using finite difference:

$$r_t^v = \exp \left[-\alpha_v \|\dot{\mathbf{q}}_t - \hat{\dot{\mathbf{q}}}_t\|^2 \right]. \quad (3)$$

The end-effector reward r_t^e measures the difference between end-effector position \mathbf{g}_t^e and the reference position $\hat{\mathbf{g}}_t^e$ in the world coordinate:

$$r_t^e = \exp \left[-\alpha_e \left(\sum_e \|\mathbf{g}_t^e - \hat{\mathbf{g}}_t^e\|^2 \right) \right]. \quad (4)$$

The center-of-mass reward r_t^c encourages the humanoid's center of mass \mathbf{c}_t to match the reference $\hat{\mathbf{c}}_t$:

$$r_t^c = \exp \left[-\alpha_c \|\mathbf{c}_t - \hat{\mathbf{c}}_t\|^2 \right]. \quad (5)$$

The weighting factors $\alpha_p, \alpha_v, \alpha_e, \alpha_c$ are set to (2, 0.005, 5, 100). All the hyperparameters of r_t^{world} are tuned to achieve best performance for the DeepMimic baseline.

2.2 Local Coordinate Reward

The local coordinate reward r_t^{local} also consists of four sub-rewards:

$$r_t^{\text{local}} = w_p r_t^p + w_e r_t^e + w_{\text{rp}} r_t^{\text{rp}} + w_{\text{rv}} r_t^{\text{rv}}, \quad (6)$$

where $w_p, w_e, w_{\text{rp}}, w_{\text{rv}}$ are weighting factors set to (0.5, 0.3, 0.1, 0.1) same as EgoPose [6].

The pose reward r_t^p is the same as that in r_t^{world} as defined in Eq. (2). The end-effector reward r_t^e takes the same form as Eq. (4) but the end-effector positions \mathbf{g}_t^e and $\hat{\mathbf{g}}_t^e$ are computed in the humanoid's

local heading coordinate. The root pose reward r_t^{rp} encourages the humanoid’s root joint to have the same height y_t and orientation quaternion \mathbf{o}_t as the reference \hat{y}_t and $\hat{\mathbf{o}}_t$:

$$r_t^{\text{rp}} = \exp \left[-\alpha_{\text{rp}} \left((y_t - \hat{y}_t)^2 + \|\mathbf{o}_t \ominus \hat{\mathbf{o}}_t\|^2 \right) \right]. \quad (7)$$

The root velocity reward r_t^{rv} penalizes the deviation of the root’s linear velocity \mathbf{l}_t and angular velocity $\boldsymbol{\omega}_t$ from the reference $\hat{\mathbf{l}}_t$ and $\hat{\boldsymbol{\omega}}_t$:

$$r_t^{\text{rv}} = \exp \left[-\|\mathbf{l}_t - \hat{\mathbf{l}}_t\|^2 - \alpha_{\text{rv}} \|\boldsymbol{\omega}_t - \hat{\boldsymbol{\omega}}_t\|^2 \right]. \quad (8)$$

Note that all features are computed in the local heading coordinate of the humanoid instead of the world coordinate. The weighting factors $\alpha_p, \alpha_e, \alpha_{\text{rp}}, \alpha_{\text{rv}}$ are set to (2, 20, 300, 0.1) same as EgoPose [6].

3 Additional Implementation Details

Residual Forces. In RFC-Explicit, each $\boldsymbol{\xi}_j$ is a 6 dimensional vector including both the force and torque applied at the contact point e_j . The torque is needed since it along with the force can model the total effect of multiple forces applied at different contact points for a single rigid body of the humanoid. We scale $\boldsymbol{\xi}_j$ by 100 after it is output by the policy. In RFC-Implicit, we similarly scale the total root torques $\boldsymbol{\eta}_r$ by 100. The weight w_{reg} for the regularizing reward r_t^{reg} defined in Eq. (3) and Eq. (5) of the main paper is set to 0.1 for both RFC-Explicit and RFC-Implicit. For RFC-Explicit, k_f and k_{cp} are set to 1 and 4 respectively. For RFC-Implicit, k_r is set to 1.

Time Efficiency. Without the residual forces, the total time of a single policy step with simulation is 3.9ms. After adding the residual forces, the time is 4.0ms for RFC-Explicit and 4.3ms for RFC-Implicit. The slight increase in time is due to the larger action space of RFC as well as the Jacobian computation in RFC-Explicit. The processing time 4.0ms translates to 250 FPS, which is well above the interactive frame rate, and the performance gain from RFC justifies the small sacrifice in speed.

Humanoid Model. As mentioned in the main paper, we have different humanoid models for different datasets, and the number of DoFs and rigid bodies varies for different humanoid models. Each DoF except for the root is implemented as a hinge joint in MuJoCo. Most joints have 3 DoFs meaning 3 consecutive hinge joints that form a 3D rotation parametrized by Euler angles. Only the elbow and knee joints have just one DoF. In MuJoCo, there are many parameters one can specify for each joint (e.g., stiffness, damping, armature). We only set the armature inertia to 0.01 to further stabilize the simulation. We leave the stiffness and damping to 0 since they are already modeled in the gains k_p and k_d of the PD controllers. The gains k_p range from 200 to 1000 where stronger joints like spine and legs have larger gains while weaker joints like arms and head have smaller gains. The gains k_d are set to $0.2k_p$. We also set proper torque limits ranging from 50 to 200 based on the gains to prevent instability. In our experiments, we find that the learned motion policies are not sensitive to the gains and scaling the gains by reasonable amount yields similar results. We believe the reason is that the policy can learn to adjust to different scales of gains.

3.1 Motion Imitation

The reference motions we use are from the following motion clips in the CMU MoCap database¹: 05_06 (ballet1), 05_07 (ballet2), 05_13 (ballet3), 88_01 (backflip), 90_02 (cartwheel), 90_05 (jump kick), 90_08 (side flip), 90_11 (handspring). The motion clips are downsampled to 30Hz. We do not apply any filters to smooth the motions and directly use the downsampled motions for training.

Table 2: Training hyperparameters for motion imitation.

| Parameter | γ | GAE(λ) | Batch Size | Minibatch Size | Policy Stepsize | Value Stepsize | PPO clip ϵ |
|-----------|----------|------------------|------------|----------------|--------------------|--------------------|---------------------|
| Value | 0.95 | 0.95 | 50000 | 2048 | 5×10^{-5} | 3×10^{-4} | 0.2 |

Training. In each RL episode, the state of the humanoid agent is initialized to the state of a random frame in the reference motion. The episode is terminated when then end of the reference motion

¹<http://mocap.cs.cmu.edu/>

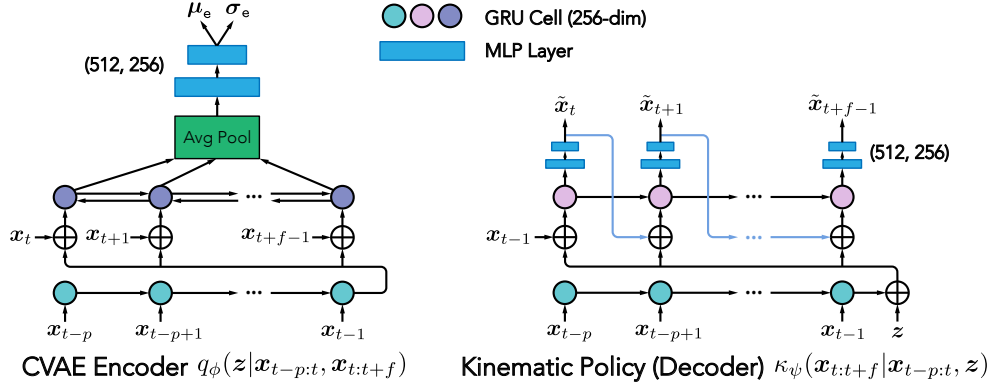


Figure 1: Network architectures for the CVAE encoder q_ϕ and kinematic policy (decoder) κ_ψ .

is reached or the humanoid’s root height is 0.1 below the minimum root height in the reference motion. As discussed in the main paper, the control policy π_θ is a Gaussian policy whose mean μ_θ is generated by a multi-layer perceptron (MLP). The MLP has two hidden layers (512, 256) with ReLU activations. The diagonal elements of the policy’s covarian matrix Σ are set to 0.1. We use the proximal policy optimization (PPO [5]) to learn the policy π_θ . We use the generalized advantage estimator GAE(λ) [4] to compute the advantage for policy gradient. The policy is updated with Adam [2] for 2000 epochs. The hyperparameter settings are available in Table 2.

3.2 Extended Motion Synthesis

Network Architectures. The CVAE Encoder distribution $q_\phi(z|x_{t-p:t}, x_{t:t+f}) = \mathcal{N}(\mu_e, \text{Diag}(\sigma_e^2))$ is a Gaussian distribution whose parameters μ_e and σ_e are generated by a recurrent neural network (RNN) as shown in Fig. 1 (Left). Similarly, the kinematic policy (decoder) $\kappa_\psi(x_{t:t+f}|x_{t-p:t}, z) = \mathcal{N}(\tilde{x}_{t:t+f}, \beta I)$ is also a Gaussian distribution whose mean $\tilde{x}_{t:t+f}$ is generated by another RNN as illustrated in Fig. 1 (Right), and β is a hyperparameter which we set to 10. We use GRUs [1] as the recurrent units for both q_ϕ and κ_ψ . For the RFC-based control policy $\tilde{\pi}_\theta(a, \tilde{a}|x, \hat{x}, z)$, we concatenate (x, \hat{x}, z) together and input them to an MLP to produce the mean of the actions (a, \tilde{a}) . The MLP has two hidden layers (512, 256) with ReLU activations.

Table 3: CVAE Training hyperparameters for the kinematic policy κ_ψ .

| Parameter | Dim(z) | Batch Size | Minibatch Size | Initial Learning Rate | KL Tolerance |
|-----------|------------|------------|----------------|-----------------------|--------------|
| Value | 128 | 10000 | 256 | 1×10^{-3} | 10 |

Training. The kinematic policy κ_ψ is trained with the CVAE objective in Eq. (6) of the main paper. We jointly optimize q_ϕ and κ_ψ for 200 epochs using Adam [2] with a fixed learning rate. We then continue to optimize the model for another 800 epochs while linearly decreasing the learning rate to 0. We list the hyperparameters for training the CVAE in Table 3. The training procedure for the RFC-based control policy $\tilde{\pi}_\theta(a, \tilde{a}|x, \hat{x}, z)$ is outlined in Alg. 1 of the main paper. Similar to motion imitation, we optimize the policy $\tilde{\pi}_\theta$ with PPO for 3000 epochs using Adam. The hyperparameter settings are the same as motion imitation as shown in Table 2.

References

- [1] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.
- [2] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [3] X. B. Peng, P. Abbeel, S. Levine, and M. van de Panne. Deepmimic: Example-guided deep reinforcement learning of physics-based character skills. *ACM Transactions on Graphics (TOG)*, 37(4):1–14, 2018.
- [4] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.

- [5] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [6] Y. Yuan and K. Kitani. Ego-pose estimation and forecasting as real-time pd control. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 10082–10092, 2019.