# Supplementary Material of Diff3DHPE: A Diffusion Model for 3D Human Pose Estimation

Jieming Zhou[1], Tong Zhang[2], Zeeshan Hayder[3], Lars Petersson[3], Mehrtash Harandi[4]
[1]Australian National University, [2]EPFL, [3]CSIRO, [4]Monash University

jieming.zhou@anu.edu.au, tong.zhang@epfl.ch,
{zeeshan.hayder, Lars.Petersson}@data61.csiro.au, mehrtash.harandi@monash.edu

## 1. Proof of Iteration Steps Required by DDIM

The reverse diffusion process proposed by DDIM [2] is:

$$\hat{\boldsymbol{y}}_{\tau_{i-1}} = \sqrt{\bar{\alpha}_{\tau_{i-1}}}\Big(\frac{\hat{\boldsymbol{y}}_{\tau_i} - \sqrt{1 - \bar{\alpha}_{\tau_i}}\hat{\boldsymbol{\epsilon}}_{\tau_i}}{\sqrt{\bar{\alpha}_{\tau_i}}}\Big) + \sqrt{1 - \bar{\alpha}_{\tau_{i-1}}}\hat{\boldsymbol{\epsilon}}_{\tau_i}, \tag{1}$$

$$\hat{\boldsymbol{y}}_0 = \frac{\hat{\boldsymbol{y}}_{\tau_1} - \sqrt{1 - \bar{\alpha}_{\tau_1}}\hat{\boldsymbol{\epsilon}}_{\tau_1}}{\sqrt{\bar{\alpha}_{\tau_1}}}, \tag{2}$$

where $\tau_i$ is sampled every $\lceil T/S \rceil$ steps from $\{t_1, t_2, ..., t_T\}$, $\tau_1 < \tau_2 < ... < \tau_S \in [1, T]$, $S < T$, $\hat{\boldsymbol{y}}_t$ is the estimated 3D coordinates at step $t$, $\hat{\boldsymbol{y}}_{\tau_S} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})$, and $\bar{\alpha}_t$ is a predefined noise schedule. In this paper, we select $cos$ schedule for $\bar{\alpha}_t$ proposed by [1]:

$$\bar{\alpha}_t = \frac{f(t)}{f(0)}, f(t) = cos\Big(\frac{t/T + s}{1 + s} \cdot \frac{\pi}{2}\Big)^2, s = 0.008. \tag{3}$$

We assume the 3D coordinate value of a human joint is between $[-1000, 1000]$ mm after centralizing the body. Then, we normalize the coordinate value to $[-1, 1]$, which is required by the diffusion model. Since $\epsilon \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})$, we have $95\%$ probability that $|\epsilon| < 2\sigma = 2$. $\sigma$ is the standard deviation of $\epsilon$. Therefore, we shall have

$$\frac{\sqrt{1 - \bar{\alpha}_{\tau_1}}}{\sqrt{\bar{\alpha}_{\tau_1}}} < 10^{-3} \cdot \frac{1}{|\epsilon|} = 5 \times 10^{-4} \tag{4}$$

in Eq. 2, which ensures impact introduced by noise value to the final prediction has $95\%$ probability smaller than 1 mm. To achieve this, the minimum $\tau_1 = 1$. Then, we derive

$$\bar{\alpha}_1 = \frac{cos\Big(\frac{1/T + s}{1+s} \cdot \frac{\pi}{2}\Big)^2}{cos\Big(\frac{s}{1+s} \cdot \frac{\pi}{2}\Big)^2} > \frac{1}{1 + (5 \times 10^{-4})^2}, \tag{5}$$

$$\frac{cos\Big(\frac{1/T + s}{1+s} \cdot \frac{\pi}{2}\Big)}{cos\Big(\frac{s}{1+s} \cdot \frac{\pi}{2}\Big)} > \sqrt{\frac{1}{1 + (5 \times 10^{-4})^2}} \tag{6}$$

from Eq. 4. According to small-angle approximations, we can have

$$\frac{1 - \frac{(\frac{1/T + s}{1+s} \cdot \frac{\pi}{2})^2}{2}}{1 - \frac{(\frac{s}{1+s} \cdot \frac{\pi}{2})^2}{2}} > \sqrt{\frac{1}{1 + (5 \times 10^{-4})^2}}, \tag{7}$$

when $T \gg 1$ and $s = 0.008$. Thus, we obtain:

$$\begin{aligned} T &> \frac{1}{\frac{2(1+s)}{\pi}\sqrt{2 - \sqrt{\frac{1}{1+(5\times 10^{-4})^2}}\big(2 - (\frac{s}{1+s} \cdot \frac{\pi}{2})^2\big)} - s} \\ &\approx 1.55 \times 10^5, \end{aligned} \tag{8}$$

which can meet the target.

## 2. Hyper-parameter settings

The hyper-parameter search space and the final choice in our experiments are listed in Table 1 and 2.

Table 1. Hyper-parameter search space.$lr$: learning rate. $StepEmb$: whether or not using step embedding. $S$: the number of reverse diffusion steps.

| Param. | Search Space |
|---|---|
| $lr$ | 1E-4,4E-4,1E-3,4E-3 |
| $StepEmb$ | T, F |
| $S$ | 1,3,4,5,6,7,8,10,15,20,40,80,160,320 |

## References

[1] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021. 1

[2] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021. 1

Table 2. Final hyper-parameters of each model. Diff3DHPE-M: Diff3DHPE with MixSTE backbone. Diff3DHPE-P: Diff3DHPE with PoseFormer backbone. DDIM-M: Diffusion model with DDIM reverse diffusion method and MixSTE backbone. $*$: we train the baselines with only L2 loss of 3D pose prediction error and normalize the training target 3D pose ground truth to $[-1, 1]$. $F$: the number of frames. $bs$: batch size. $lr$: learning rate. $dim$: embedding dimension. $depth$: the number of Transformer blocks. $StepEmb$: whether or not using step embedding. $S$: the number of reverse diffusion steps. $dr$: dropout rate. $wd$: weight decay. $lrd$: learning rate decay factor.

| Model | Dataset | $F$ | $bs$ | $lr$ | $dim$ | $depth$ | $StepEmb$ | $S$ | $dr$ | $wd$ | $lrd$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Diff3DHPE-P | H3.6M CPN | 81 | 1024 | 4E-3 | | | T | 5 | | | |
| Diff3DHPE-P | H3.6M GT | 81 | 1024 | 4E-3 | 32 | 8 | | 5 | | | |
| PoseFormer | H3.6M CPN | 81 | 1024 | 1e-4 | | | N/A | N/A | | | |
| PoseFormer | H3.6M GT | 81 | 1024 | 1e-4 | | | | | | | |
| Diff3DHPE-M | H3.6M CPN | 81 | 64 | 4E-4 | | | | 9 | | | |
| Diff3DHPE-M | H3.6M CPN | 243 | 24 | 4E-4 | | | | 5 | | | |
| Diff3DHPE-M | H3.6M GT | 81 | 64 | 4E-4 | | | | 5 | | | |
| Diff3DHPE-M | H3.6M GT | 243 | 24 | 4E-4 | | | T | 6 | | | |
| Diff3DHPE-M w/o PDE | H3.6M CPN | 81 | 64 | 1E-4 | | | | 6 | 0.1 | 0.1 | 0.1 |
| Diff3DHPE-M w/o PDE | H3.6M CPN | 243 | 24 | 1E-4 | | | | 5 | | | |
| DDIM-M | H3.6M CPN | 81 | 64 | 4E-4 | 512 | 16 | | 40 | | | |
| DDIM-M | H3.6M CPN | 243 | 24 | 4E-4 | | | | 80 | | | |
| MixSTE | H3.6M CPN | 81 | 64 | 1E-4 | | | | | | | |
| MixSTE | H3.6M CPN | 243 | 24 | 1E-4 | | | N/A | N/A | | | |
| MixSTE | H3.6M GT | 81 | 64 | 1E-4 | | | | | | | |
| MixSTE | H3.6M GT | 243 | 24 | 1E-4 | | | | | | | |
| Diff3DHPE-M | 3DHP GT | 27 | 64 | 4E-4 | | | F | 7 | | | |