

## 1. Supplementary Material

### 1.1. Discrete cosine transformation

First, the full gaze-motion sequence  $X \in \mathbb{R}^{3 \times (j+1) \times (T+F)}$  is encoded to transformed sequence  $Y \in \mathbb{R}^{3 \times (j+1) \times L}$  temporally by discrete cosine transform (DCT) [ANR74] as follows:

$$Y = \text{DCT}(X) = DX \quad (1)$$

where  $D \in \mathbb{R}^{(T+F) \times (T+F)}$  is the DCT matrix. The reason we applied DCT first is that extensive work [MLSL19, CZL\*23] has demonstrated such a transformer can extract both current and periodic temporal features, which leads to smoother predictions and a more compact representation.

Due to the orthogonality of DCT, the original gaze-motion sequence can be recovered from the transformed sequence:

$$X = \text{iDCT}(Y) = D^T Y \quad (2)$$

To reduce the computational complexity and weaken the effects of high-frequency noise, we select the first  $L$  row of  $D$ ,  $D_L \in \mathbb{R}^{L \times (T+F)}$  to replace  $D$ . In a subsequent paper, we will replace  $D$  with  $D_L$  for simplicity.

### 1.2. Details spatial graph attention network

Given features  $H' = [h'_1, h'_2, \dots, h'_L] \in \mathbb{R}^{f \times j \times L}$  extracted by the temporal GAT, the spatial GAT aggregated features across spatial dimensions as follows:

$$h''_i = \text{LeakyReLU} \left( \frac{1}{\bar{N}_{head}} \sum_{n=1}^{\bar{N}_{head}} \sum_{k=1}^j \bar{\alpha}_{ik}^n \bar{h}'_k \right), \quad (3)$$

where  $h''_i \in \mathbb{R}^{f \times L}$  is the output feature of node  $i$  and  $\bar{N}_{head}$  denotes the number of heads for attention. We also applied an average function to fuse different output features from each head in spatial GAT. For each head, the attention matrix  $\bar{\alpha}_{ik}^n$  represents interactions between each timestamp, calculated as follows:

$$\bar{\alpha}_{ik}^n = \frac{\exp(\text{LeakyReLU}(\bar{\mathbf{a}}^n [\bar{h}'_i \oplus \bar{h}'_k]))}{\sum_{l=1}^j \exp(\text{LeakyReLU}(\bar{\mathbf{a}}^n [\bar{h}'_i \oplus \bar{h}'_l]))}, \quad (4)$$

where  $\bar{\mathbf{a}}^n$  is a parameter vector  $\in \mathbb{R}^{2f \times L \times 1}$

### 1.3. Details of self-attention block

Given noised sequence  $Y^t$  at timestep  $t$ , We then apply an efficient self-attention block [ZCP\*22] to further model temporal correlations between each frame:

$$\mathbf{Y} = \text{Dropout} \left( \text{softmax}(\mathbf{Q}) \text{softmax}(\mathbf{K}^\top) \right) \text{LN}(\mathbf{V}) + Y^t, \quad (5)$$

where LN is layer normalisation,  $\mathbf{Q} \in \mathbb{R}^{L \times d}$ ,  $\mathbf{K} \in \mathbb{R}^{L \times d}$ , and  $\mathbf{V} \in \mathbb{R}^{L \times d}$  are calculated using the original self-attention mechanism [VSP\*17]:

$$\mathbf{Q} = \mathbf{W}_q Y^t, \mathbf{K} = \mathbf{W}_k Y^t, \mathbf{V} = \mathbf{W}_v Y^t, \quad (6)$$

where  $\mathbf{W}_q$ ,  $\mathbf{W}_k$  and  $\mathbf{W}_v$  are learnable parameter matrices.

### 1.4. Details of step hint module

In this paper, we employed feature-wise linear modulation to inject the fused embedding  $\mathbf{e}$  into the output  $\mathbf{Y}$  of the self-attention block:

$$w = \phi_w(\psi(\mathbf{e}_t)), b = \phi_b(\psi(\mathbf{e}_t)), \mathbf{Y}' = \text{SiLU}(w \cdot \mathbf{Y} + b), \quad (7)$$

where  $(\cdot)$  denotes element-wise multiplication,  $\phi_w$  and  $\phi_b$  are linear projections, and  $\psi$  is a single layer MLP with SiLU activation function. This modulation allowed the step hint embedding to influence the self-attention features. We also applied this block after each cross-attention and MLP block, enabling the timestep embedding to provide hints throughout the network.

### 1.5. Detailed training process

The detailed training procedure is shown in [algorithm 1](#).

---

#### Algorithm 1: Training procedure of proposed method

---

**Input:** observed gaze-motion sequence  $X_{obs}$ , noising steps  $T$ , the initialized gaze-motion fusion network  $GazePoseFuse$ , the initialized gaze encoder  $GazeEncoder$ , the initialized gaze encoder  $PoseEncoder$ , the initialized noise prediction network  $\epsilon_\theta$ , full gaze-motion sequence  $X_{full}$ , max iterations  $I_{max}$

**Output:** trained  $GazePoseFuse$  and  $\epsilon_\theta$

$Y = \text{DCT}(\text{Pad}(X));$

$H'' = GazePoseFuse(Y);$

$Y_{full} = \text{DCT}(X_{full});$

**for**  $t = 1, 2, \dots, I_{max}$  **do**

$Y_{full}^0 = Y_{full};$

$t = \text{Uniform}(\{1, 2, \dots, T\});$

$\epsilon \sim N(0, I);$

$\theta = \theta - \nabla_\theta \|\epsilon - \epsilon_\theta(Y_{full}, H'', t)\|^2$

**return**  $GazePoseFuse$  and  $\epsilon_\theta$ .

---

### 1.6. Detailed Completion denoising process

We employed an ingenious prediction mask mechanism [CZL\*23] to denoise progressively:

$$Y^{t-1} = \text{DCT} \left( \mathbf{M} \odot \text{iDCT}(Y_{orig}^{t-1}) + (\mathbf{1} - \mathbf{M}) \odot \text{iDCT}(Y_{pred}^{t-1}) \right), \quad (8)$$

where  $\mathbf{M} = \underbrace{[1, 1, \dots, 1]}_H, \underbrace{[0, 0, \dots, 0]}_F \in \mathbb{R}^{(H+F) \times 1}$  is a mask vector

indicating which frames are observed.  $Y_{orig}^{t-1}$  is obtained by adding  $t-1$  iterations of Gaussian noise, and  $Y_{pred}^{t-1}$  is obtained by denoising the output  $Y^t$  from the previous iteration:

$$Y_{orig}^{t-1} = \sqrt{\bar{\beta}_{t-1}} Y + \sqrt{1 - \bar{\beta}_{t-1}} \epsilon, \bar{\beta}_t = \prod_{i=1}^t \beta_i, \beta_i \in [0, 1]. \quad (9)$$

$$Y_{pred}^{t-1} = \frac{1}{\sqrt{\bar{\beta}_t}} \left( Y^t - \frac{1 - \bar{\beta}_t}{\sqrt{1 - \bar{\beta}_t}} \epsilon_\theta(Y^t, H'', t) \right) + (1 - \beta) \epsilon, \quad (10)$$

where  $t$  denotes the  $t$ -th noise iteration,  $\beta_t$  controls noise level, and  $\epsilon \sim N(0, I)$ . At the start,  $Y^t$  is sampled from Gaussian noise.

Through this iteration process, we can obtain the full generated sequence  $\bar{\mathbf{p}} = \text{iDCT}(Y^0)$ . The last  $F$  frames were predicted motions.

The detailed denoising procedure is shown in [algorithm 2](#).

---

**Algorithm 2:** Inference procedure of proposed method
 

---

**Input:** observed gaze-motion sequence  $X$ , noising steps  $T$ , the trained gaze-motion fusion network  $\text{GazePoseFuse}$ , the trained noise prediction network  $\epsilon_\theta$ , noising steps  $T$

**Output:** future motions  $\mathbf{p}$

$Y^T \sim N(0, I)$ ;

$Y = \text{DCT}(\text{Pad}(X_{full}))$ ;

$H'' = \text{GazePoseFuse}(Y_{full})$ ;

**for**  $t \in T, T-1, \dots, 1$  **do**

$\epsilon \sim N(0, I)$  if  $t > 1$ , else  $\mathbf{z} = 0$ ;

$Y_{orig}^{t-1} = \sqrt{\beta_t} Y + \sqrt{1 - \beta_t} \epsilon$ ;

$Y_{pred}^{t-1} = \frac{1}{\sqrt{\beta_t}} \left( Y^t - \frac{1 - \beta_t}{\sqrt{1 - \beta_t}} \epsilon_\theta(Y^t, H'', t) \right) + (1 - \beta) \epsilon$ ;

$Y^{t-1} = \text{DCT}[\mathbf{M} \odot \text{iDCT}(Y_{orig}^{t-1}) + (\mathbf{1} - \mathbf{M})$

$\odot \text{iDCT}(Y_{pred}^{t-1})$ ];

$\mathbf{p} = \text{iDCT}(Y^0)$ ;

$\bar{\mathbf{p}} = \mathbf{p}_{H:H+F}$ ;

**return**  $\bar{\mathbf{p}}$

---

### 1.7. User Study Detail

[Figure 1](#) shows an example of the instructions and definition of realism shown to the user before starting the test. In our user study, we randomly selected 18 samples for comparison. [Figure 2](#) shows an example of our interface to enable participants to rank them. Note that in each question, the order of each method is randomly shuffled.

### 1.8. Further Visualisation Results

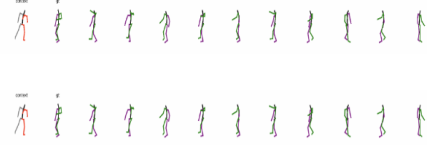
We illustrated more visualisation samples in [Figure 3](#), [Figure 4](#), [Figure 5](#), [Figure 6](#), and [Figure 7](#)

### References

- [ANR74] AHMED N., NATARAJAN T., RAO K. R.: Discrete cosine transform. *IEEE transactions on Computers* 100, 1 (1974), 90–93. 1
- [CZL\*23] CHEN L.-H., ZHANG J., LI Y., PANG Y., XIA X., LIU T.: Humanmac: Masked motion completion for human motion prediction. *arXiv preprint arXiv:2302.03665* (2023). 1
- [KBM\*20] KRATZER P., BIHLMAIER S., MIDLAGAJNI N. B., PRAKASH R., TOUSSAINT M., MAINPRICE J.: Mogaze: A dataset of full-body motions that includes workspace geometry and eye-gaze. *IEEE Robotics and Automation Letters* 6, 2 (2020), 367–373. 4
- [MLSL19] MAO W., LIU M., SALZMANN M., LI H.: Learning trajectory dependencies for human motion prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2019), pp. 9489–9497. 1

**IMPORTANT:** Please read it carefully before the test.

In this test, you will find several sets of **human motions** represented by a moving skeleton. The first column "context" means the input of the motion prediction model, and the second column "gt" means the actual movements in the real world. For the rest of the 10 different moving skeletons, they are predicted from a model. Different rows denote the results from different models. For example:



We aim to compare different generative models. Your mission is to rate them according to the **realism** and **precision** of their predictions.

**Realism:** If these poses are **plausible**. You can check if there are any angle distortions, too short/long limbs, or implausible poses, any sudden or unreasonable changes during the whole motion.

**Precision:** If these motions **align with the 'gt'**. You can measure the similarity between each motion and gt.

PS: It will take some time for the webpage to fully load the GIFs, so please be patient and wait for all the animations to load smoothly before making your selection.

Back

Next

**Figure 1:** The instructions and definition of realism of our questionnaire in user study

Sequence 1



1.1 Select the row with more realistic predictions (S1 and S2 means the 1st and 2nd row respectively)

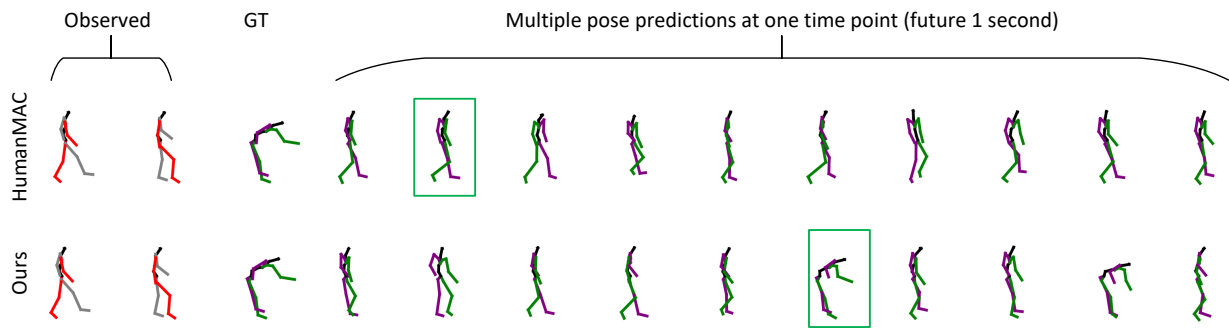
S1  
 S2

1.2 Select the row with more precise predictions (S1 and S2 means the 1st and 2nd row respectively)

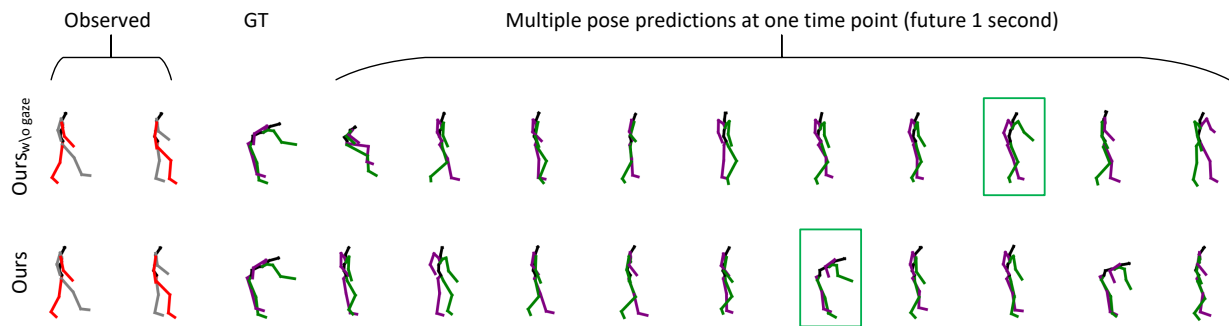
S1  
 S2

**Figure 2:** The interface for ranking the predictions from different models according to realism and precision.

- [VSP\*17] VASWANI A., SHAZEER N., PARMAR N., USZKOREIT J., JONES L., GOMEZ A. N., KAISER Ł., POLOSUKHIN I.: Attention is all you need. *Advances in neural information processing systems* 30 (2017). 1
- [ZCP\*22] ZHANG M., CAI Z., PAN L., HONG F., GUO X., YANG L., LIU Z.: Motiondiffuse: Text-driven human motion generation with diffusion model. *arXiv preprint arXiv:2208.15001* (2022). 1

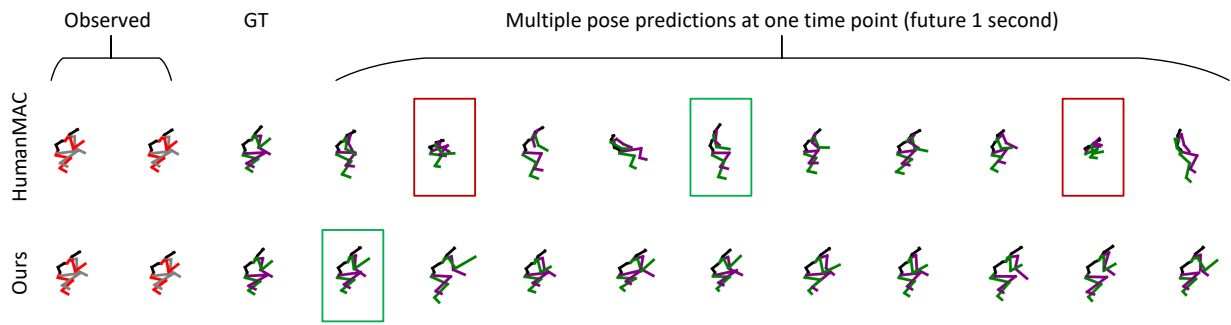


**Figure 3:** Ground truth (GT) human pose at future one second and multiple pose predictions generated by different methods on the GIMO dataset [ZYM\*22] with the best prediction (lowest  $l_2$  distance to GT) boxed in green. The ground truth poses show a motion an action to bend down and take an object. Predictions from HumanMAC generally fail to forecast this intention, while our gaze-guided method is better able to predict future motion aligned better with this intention.

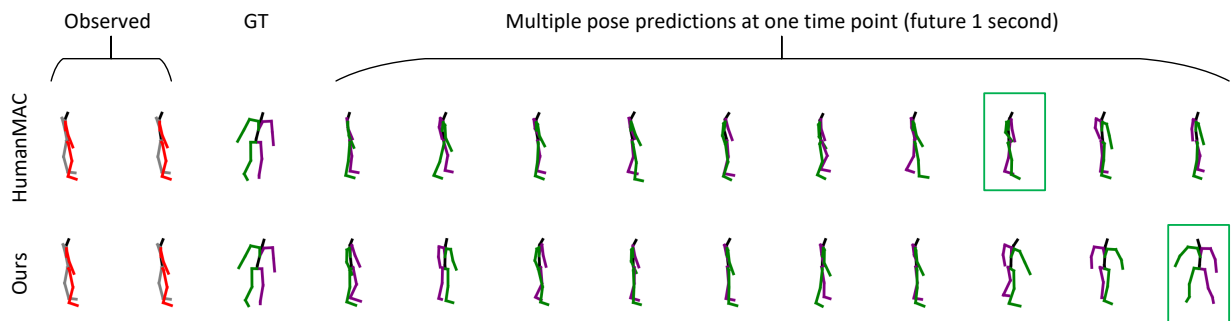


**Figure 4:** Ground truth (GT) human pose at future one second and multiple pose predictions generated by different methods on the GIMO dataset [ZYM\*22] with the best prediction (lowest  $l_2$  distance to GT) boxed in green. The ground truth poses show a motion an action to bend down and take an object. Predictions from  $Ours_{w/o\ gaze}$  generally fail to forecast this intention, while our full model is better able to predict future motion aligned better with this intention.

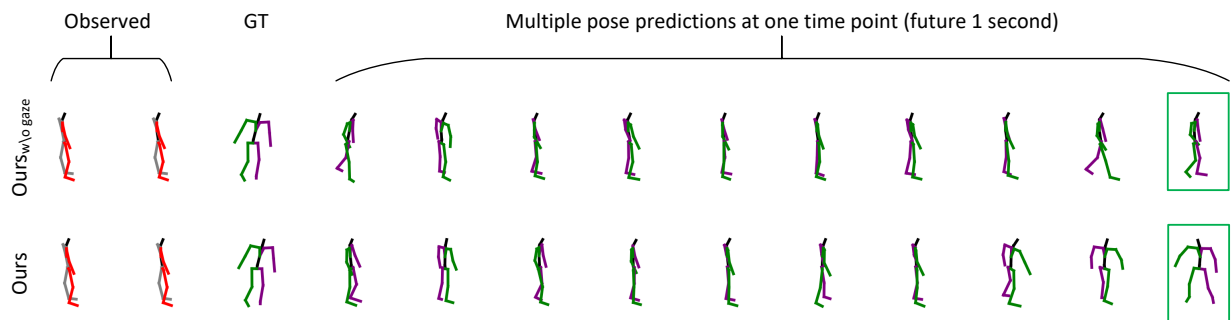
[ZYM\*22] ZHENG Y., YANG Y., MO K., LI J., YU T., LIU Y., LIU C. K., GUIBAS L. J.: Gimo: Gaze-informed human motion prediction in context. In *European Conference on Computer Vision* (2022), Springer, pp. 676–694. 3, 4



**Figure 5:** Ground truth (GT) human pose at future one second and multiple pose predictions generated by different methods on the GIMO dataset [ZYM\*22] with the best prediction (lowest  $l_2$  distance to GT) boxed in green and implausible cases marked in red. Our method can generate multiple reasonable and precise predictions while some poses predicted by HumanMAC are physically implausible.



**Figure 6:** Ground truth (GT) human pose at future one second and multiple pose predictions generated by different methods on the MoGaze dataset [KBM\*20] with the best prediction (lowest  $l_2$  distance to GT) boxed in green. The observed motion sequence is in place, and the groundtruth pose a second later is a sudden turn to the right at roughly 90 degrees. Predictions from HumanMAC are all in place still, while our gaze-guided method can recognise this possible intention and one of our generation accurately predicts the future of this pose.



**Figure 7:** Ground truth (GT) human pose at future one second and multiple pose predictions generated by different methods on the MoGaze dataset [KBM\*20] with the best prediction (lowest  $l_2$  distance to GT) boxed in green. The observed motion sequence is in place, and the groundtruth pose a second later is a sudden turn to the right at roughly 90 degrees. Predictions from  $Ours_{w/o\ gaze}$  are all in place still, while our full model can recognise this possible intention and one of our generation accurately predicts the future of this pose.