NeFT-Net: N-window extended frequency transformer for rhythmic motion prediction

Adeyemi Ademola, David Sinclair, Babis Koniaris, Samantha Hannah, Kenny Mitchell



PII:	S0097-8493(25)00085-8
DOI:	https://doi.org/10.1016/j.cag.2025.104244
Reference:	CAG 104244
To appear in:	Computers & Graphics
Received date :	4 January 2025
Revised date :	14 April 2025
Accepted date :	2 May 2025

Please cite this article as: A. Ademola, D. Sinclair, B. Koniaris et al., NeFT-Net: N-window extended frequency transformer for rhythmic motion prediction. *Computers & Graphics* (2025), doi: https://doi.org/10.1016/j.cag.2025.104244.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2025 Published by Elsevier Ltd.

Graphical Abstract (for review)

Journal Click here to access/download;Graphical Abstract (for review);Neft_net_control.tiff



PDF of Manuscript, including all figures, tables etc. - must not contain any author details

NeFT-Net: N-window Extended Frequency Transformer for Rhythmic Motion Prediction

ARTICLE INFO

Keywords: Machine Learning Motion Processing Rendering Virtual reality

ABSTRACT

Advancements in prediction of human motion sequences are critical for enabling online virtual reality (VR) users to dance and move in ways that accurately mirror real-world actions, delivering a more immersive and connected experience. However, latency in networked motion tracking remains a significant challenge, disrupting engagement and necessitating predictive solutions to achieve real-time synchronization of remote motions. To address this issue, we propose a novel approach leveraging a synthetically generated dataset based on supervised foot anchor placement timings for rhythmic motions, ensuring periodicity and reducing prediction errors. Our model integrates a discrete cosine transform (DCT) to encode motion, refine high-frequency components, and smooth motion sequences, mitigating jittery artifacts. Additionally, we introduce a feed-forward attention mechanism designed to learn from N-window pairs of 3D key-point pose histories for precise future motion prediction. Quantitative and qualitative evaluations on the Human3.6M dataset highlight significant improvements in mean per joint position error (MPJPE) metrics, demonstrating the superiority of our technique over state-of-the-art approaches. We further introduce novel result pose visualizations through the use of generative AI methods.

1. Introduction

In the fields of virtual reality (VR) and computer vision, real-time tracking is crucial for recovering accurate 3D pose data. Human joint pose data is commonly captured using multi-camera or single-camera setups integrated with AI algorithms to obtain depth information and directly recover pose key points and joint orientations. Nevertheless, challenges such as limited sensor range, occlusion, and latency persist in tracking 3D pose data. In order to improve immersion and engagement in patterned motion scenarios, there is a high demand for techniques that minimize latency [21] [28] during motion tracking through motion prediction.

Deep learning techniques have significantly advanced the domain of human motion prediction [5] [10]. Among these, recurrent neural networks (RNNs) have become particularly popular for predicting sequential human pose data [19] [12]. However, when it comes to long-term horizons and periodic motions, RNNs often struggle due to their inability to effectively capture long-term history, which is essential for forecasting periodic motion actions. To address this limitation, recent approaches have incorporated encoders [22] to better represent historical information.

Our work introduces a multi-window extended frequency attention-based human motion prediction technique that utilizes synthetically generated periodic data based on re-timed foot anchor placements, as illustrated in figure 2. Our method is motivated by the observation that humans tend to repeat their motions in actions such as dancing to music beats. To validate this, we focus on the context of rhythmic motion prediction, where we demonstrate the effectiveness of our approach by re-timing *Human3.6m* [18] to match these rhythmic patterns. We present results based on analyzing relevant information from significant bones, such as the feet, over a fixed-length period.

Inspired by previous works [24], we represent each sub-sequence of foot anchors in the trajectory space using a Discrete Cosine Transform (DCT).

We then introduce our dual-windowed extended frequency motion attention as weights for DCT-encoded motion aggregation into a future motion estimate. To encode spatial dependencies between joints, we combine the motion estimate with the last observed matching period, using the result as input to a graph convolutional network (GCN) [13]. Our experiment, as shown in figure 5, demonstrates that our approach outperforms state-of-the-art methods in long-term and short-term periodic motion prediction on the Human3.6M walking and walking together datasets. Our work extends *DeFT-Net* [1] developed upon insights from Mao et al [23], specifically improving 3D pose motion prediction for known periodicity based on foot anchor placements.

^{*}Corresponding authors ORCID(s):

NeFT-Net



Figure 1: NeFT-Net stable diffusion walking sequence visualization. Note: We apply the positive prompt "A walking human" to the depth map image on the left (a) with different seeds to achieve (b), (c), and (d)

To summarize, the main contributions of this paper are:

- a re-timed motion with supervised foot anchor information of periodic cycles, such as walking, for the defined use case of rhythmic motion prediction.
- an improved overall mean per joint position error (MPJPE) results compared to state-of-the-art methods in experiments on the Human3.6M dataset for forecasting short and long-term motions by introducing *MultiWindowDCT* attention aligned on a best fit period of each motion sequence.
- a strategy for photorealistic visualization of human body motion sequences by employing the use of stable diffusion on depth maps as shown in figure 1.
- an open source version of our code implementation https://github.com/CarouselDancing/NeFT-net

2. Related Work

Human motion prediction relates to a variety of research areas like Computer Vision and Machine learning (ML), where predicting future movements is essential for applications in computer graphics and virtual reality. Section 2.1 details the various traditional techniques employed in the task of motion style synthesis and prediction. Section 2.2 describes how recurrent neural networks (RNN) has been adopted over the years for sequence-to-sequence 3D human motion prediction. Section 2.3 highlights the uniqueness of the attention-based approach compared to other approaches for motion prediction.

2.1. Traditional Approaches

Motivated by the inherent probabilistic nature of periodic human motion, early methods such as Boltzmann machines and Hidden Markov Models (HMMs) [29] [4] have been widely used to predict motion sequences. Style interpolation techniques are also frequently applied to synthesize motion, often driven by scripts, 2D video inputs, or to generate new choreography for virtual motion capture. While these methods offer robust solutions, they lack the adaptability and precision needed for capturing both short-term and long-term dependencies, particularly in dynamic contexts like dance and rhythmic walking sequences. Other advancements have introduced probabilistic models that leverage large motion databases and low-dimensional representations [27]. These methods utilize implicit empirical distributions and efficient binary tree-based search to approximate the true distribution of human motion. By structuring motion data efficiently, they allow for realistic motion synthesis and robust tracking within Bayesian frameworks, addressing both adaptability and precision challenges.

2.2. Recurrent Neural Networks (RNN) Approaches

RNNs have grown in prominence for 3D human motion prediction tasks [9]. The encoder-decoder model (ERD), first introduced by *Fragkiadaki et. al* [12], incorporates Long Short-Term Memory (LSTM) cells in the latent space for capturing motion dynamics. The work of *Jain et al.* [19] leverages a spatio-temporal graph skeleton, utilizing RNNs as

nodes to model kinematic chain joint dependencies. Aksan et al [3] replace dense output layers in the RNN architecture with structural prediction layers to explicitly model joint dependencies that follow a kinematic chain. In the works of Ghosh et al [14], a separate denoising auto-encoder is trained to correct noisy outputs. All these techniques suffer inability to capture long-range motion history trajectories.

However, RNN-based methods have historically struggled with capturing long-term motion history, leading to limitations in predicting prolonged sequences. In response, *Martinez et. al* [25] introduced a sequence-to-sequence (Seq2Seq) architecture incorporating an input-to-output skip connection, which mitigates some of the inherent bias by training the model with its own predictions. Despite improved results over earlier pose-based models [19], the discontinuity between ground truth and predicted frames persisted.

To address this, *Pavllo et al* [26] adapted the *teacher-forcing technique*, allowing the model to gradually learn from its own outputs, further enhancing prediction accuracy. Additionally, *Chiu et al* [9] introduced a hierarchical RNN model that operates across multiple time scales to better capture motion variability over different time spans. Furthermore, adversarial training methods proposed by *Gui et al.* [15] enable the generation of smoother motion sequences.

In the work of *Hernandez et al* [17], human motion forecasting was framed as a tensor imputation problem, with generative adversarial networks (GANs) adapted for long-term prediction. Although these techniques resulted in improved performance, the use of adversarial networks introduces challenges in training, such as instability due to the adversarial nature of the generator-discriminator dynamics, difficulty in achieving convergence, and sensitivity to hyperparameters, particularly when applied to periodic datasets requiring precise foot anchor encoding.

2.3. Beyond Recurrent Models

Given the drawbacks of RNNs, several works have employed the use of feed-forward networks as an alternative solution. [5] [24] The work of Butepage et. al [5] introduced a fully connected feed forward to process the recent history poses, investigating techniques to encode temporal historical information via convolution and exploiting the kinematic tree to encode spatial information. Li et. al [22] suggest a convolutional sequence-to sequence model (CNN) processing a two-dimensional pose matrix whose column represent the pose at every time step. The model was employed to extract a pose motion prior from long-term motion history of frames, which, in conjunction with more recent motion history, was used as an input to an auto regressive network for future pose prediction. While more effective than RNN-based frameworks, the manually selected size of the convolutional windows highly influences the temporal encoding of motion sequences. To address this, Aksan et. al [2] introduced a spati-temporal transformer encompassing a fully auto-regressive approach to model temporal dependencies given the recursive nature of human motion. Cai et. al [6] leverage a transformer architecture on the DCT coefficients extracted from the seed sequence and make joint predictions progressively by following a kinematic tree. Similarly, Mao et. al [24] encodes joint sequence via DCT and train a graph convolutional network (GCN) to capture/learn inter-joint dependencies. Since the GCN operates on temporal windows of poses to produce an output, the pose forecast are limited to a predetermined length. To address this they extracted DCT coefficients from shorter sub-sequences in a sliding window fashion aggregated with a 1D attention block. Guinot et al [16] introduced a stacked-attention mechanism utilizing synthetic IMU data to improve long-term dependency handling in dance motion prediction. This method addresses the limitations of traditional RNNs by transforming motion dynamics into the frequency domain using discrete cosine transform (DCT), which better encodes temporal information.

Our work is related to these approaches, but differs in two aspects. First, we introduce windowed inputs of a timebeat signal based on foot anchor pose information to the DCT windowed input so our model can learn periodic motions of short and long term history in the frequency domain. We then introduce an N-window extended frequency model with a focus on motion periodicity.

3. Method Overview

Our technique introduces a unique approach to improving human motion prediction by incorporating periodic patterns and adapting a multi-window of poses Z_i . Each Z_i consists of three concatenated slices S_i , $S_{i+p+offset}$, and $S_{i+2p+2offset}$ from the motion history $S_1 = [s_1, s_2, s_3, \dots, s_N]$. Here, *p* represents the period, and offset allows flexibility in adjusting the relative positions of these slices. This technique captures long-term temporal dependencies by analyzing different periods within human motion data, thus enhancing our model's ability to forecast future poses with improved performance. As shown in figure 2, we synthesize 3D pose data by interpolating frames containing

motion foot anchor information from natural walking sequences in the Human3.6M dataset. We apply *spherical interpolation* for pose rotations and *linear interpolation* for pose translations to ensure smooth periodic motions. Since future frame forecasting from past sequences is the main goal, our method parallels approaches that utilize Discrete Cosine Transform (DCT) to encode motion, suppress high frequencies, and smooth jittery motions as seen in prior work [24, 23]. To adapt the attention model to periodic motion cycles, we fold pose tensors to learn smooth motion transitions. Our model utilizes window slices of encoded periodic motion. For instance, if the first window captures the current motion, the second window integrates the immediate history, and the third slice looks two steps further back. This three-slice stack model enables more robust short- and long-term motion forecasting.



Figure 2: A skeleton-grid comparison of the fixed DCT motions from the HistRepeatDCT method [23] and our re-timed multi-window extended DCT motions for test subject 5 walking synchronized with right foot anchor placements. The fixed DCT motion sequence is shown as right leg purple/left leg green, and our multi-window extended re-timed DCT motions as right leg red/left leg blue skeleton. Note: The red circles represent foot placement re-timed frames and purple circles define foot placements from start to end of the original sequence.

3.1. Foot Anchor Frame Interpolation

As our goal is to learn from periodic walking sequence motions and forecast future pose motions, similar to *Cao et al.* [7], we rely on frame annotations based on the right foot placement at every n^{th} frame. For periodic actions, such as walking and walking together, *linear interpolation* is applied to the root joint for smooth transitions between frames.

In Equation 1, we compute a weighted average between the translation vectors of two key frames, p_1 and p_2 . The interpolation factor $t \in [0, 1]$ controls the degree of blending between these frames. When t = 0, the result is entirely p_1 , and when t = 1, the result is p_2 . For intermediate values of t, the linear interpolation (lerp) computes a gradual transition between the two translation vectors, creating smooth transitions in position between frames.

$$\operatorname{lerp}(p_1, p_2, t) = (1 - t)p_1 + tp_2 \tag{1}$$

In addition to translation interpolation, we also handle rotational changes between frames. Unlike translations, rotations are more complex and require spherical interpolation to compute smooth rotational transitions. Drawing from Kapoulkine's spherical linear interpolation approximation [20], we define a spherical path between the rotations and create key rotations from the rotation vectors of two consecutive frames.

In Equation 2, we perform spherical linear interpolation (slerp) between two quaternions, q_1 and q_2 , which represent rotations at two keyframes. The angle θ is the shortest angle between the two quaternions, and $t \in [0, 1]$ is the interpolation factor. The sine terms ensure that the interpolation follows the shortest path on the spherical surface, smoothly transitioning between the two rotations. When t = 0, the result is the first rotation q_1 , and when t = 1, the result is q_2 . This method provides a constant-speed rotational interpolation, crucial for preserving the natural flow of human motion.

$$\operatorname{slerp}(q_1, q_2, t) = \frac{\sin((1-t)\theta)}{\sin(\theta)}q_1 + \frac{\sin(t\theta)}{\sin(\theta)}q_2$$
(2)

We combine both interpolation techniques to achieve periodic dataset-based foot anchor frame placements and pass these sequences in an encoded DCT fashion to our multi-window frequency transformer. This method allows our model to learn and forecast future motion patterns from periodic sequences efficiently with fewer errors.

4. Multi-Window Frequency Attention

Our multi-window attention presents a novel approach to addressing the complexities of human motion forecasting, particularly in periodic actions such as walking. As natural human motion contains short-term and long-term dependencies, which can be difficult to capture using traditional forecasting models, we address these challenges by incorporating multiple temporal windows representing different segments of the motion history, an adaptive weighting mechanism, and frequency-domain transformation. Through the use of the Discrete Cosine Transform (DCT) [8] and Graph Convolutional Networks (GCNs) [11], our model is more robust to temporal and spatial dependencies present in natural human motion.



Figure 3: Overview of NeFT-Net. Our re-timed DCT input poses are shown within the solid red boxes with the multi-window extended history, and the predicted poses are shown within dotted green boxes. The last observed poses are initially used as query. For every consecutive poses in the history (key), we compute an attention score to weigh the multi-window DCT coefficients (values) of the corresponding sub-sequence. The weighted sum of such values is then concatenated with the DCT coefficients of the last observed sub-sequence to predict the future. This comprises the transformer model of *OurMultiWindowDCT*.

The core idea behind our model is the use of *N=three temporal windows*, each representing a different portion of the motion history: Current Window, Dual Window, Nth-Past temporal Window. This segmentation allows the model to better account for motion patterns over time. The introduction of learnable weights enables the model to dynamically adjust the relative importance of each window. We compute the **deltas**, or differences, between adjacent windows to capture motion changes over time:

$$\Delta_{cp} = \mathbf{X}_c - \mathbf{X}_p \tag{3}$$

$$\Delta_{pd} = \mathbf{X}_p - \mathbf{X}_d \tag{4}$$

Next, the model applies learnable weights α_c , α_p , and α_d to adaptively weight the different temporal windows:

$$\mathbf{X}_{\text{weighted}} = \alpha_c \mathbf{X}_c + \alpha_p \mathbf{X}_p + \alpha_d \mathbf{X}_d \tag{5}$$

This adaptive weighting ensures that the model remains flexible, especially when the nature of the motion changes over time.

Ne	F٦	Г-I	N	et
	•			~ ~

		Walking		Walking Together		gether
Frame No.	1	3	5	8	9	10
HistRep[23]	5.68	17.28	27.62	40.31	43.69	46.81
DeFT-Net[1]	5.45	16.78	26.50	38.41	41.69	44.78
Ours	5.31	16.23	25.48	37.04	40.75	43.54

Table 1

Following baseline setting MPJPE Batch evaluation results for test Subject 5 comparison on our re-timed interpolated vs original History Repeats Itself DCT [23] method with Human3.6m datasets for predicting human motion at various frames for activities *walking* and *walking together*.

Frequency Domain Transformation (DCT)

In addition to our N-window temporal representation, the model leverages the frequency domain through the *discrete cosine transform (DCT)* to handle periodic motion patterns. DCT transforms the motion data from the time domain to the frequency domain, which is particularly useful for periodic actions like walking, where repeating patterns occur. The DCT is defined as:

$$\mathbf{X}_{\text{DCT}}(k) = \sum_{n=0}^{N-1} \mathbf{X}(n) \cos\left[\frac{\pi}{N}\left(n + \frac{1}{2}\right)k\right]$$
(6)

where k represents the frequency index and N is the length of the sequence. Applying DCT to the weighted windows yields:

$$\mathbf{X}_{\text{DCT-weighted}} = \text{DCT}(\mathbf{X}_{\text{weighted}})$$
(7)

We apply the same principle in our Multi-windowDCT approach, where the DCT is applied to sequences from the current, dual, and *n*-past temporal windows. This transformation emphasizes the dominant frequencies in the motion while suppressing high-frequency noise, leading to smoother and more accurate predictions.

Attention Mechanism

An attention mechanism is employed to weigh the importance of different frames in the motion sequence, enabling the model to focus on the most relevant information. The general attention-weighted representation of the motion sequence is given by:

$$\mathbf{X}_{\text{attention}} = \sum_{t=1}^{T} \alpha_t \mathbf{x}_t \tag{8}$$

where α_t are the normalized attention weights for each frame \mathbf{x}_t , computed as:

$$\alpha_t = \frac{\exp(a_t)}{\sum_{i=1}^T \exp(a_i)}, \quad a_t = \operatorname{softmax}(\mathbf{q}^{\mathsf{T}} \mathbf{k}_t)$$
(9)

Here, **q** represents the query (the current motion), and \mathbf{k}_t represents the key (motion history).

To incorporate contributions from multiple temporal windows in our n-windowDCT approach, the attention mechanism is extended to weigh both individual frames within each window and the windows themselves. The n-windowed attention-weighted representation is given by:

$$\mathbf{X}_{\text{weighted}} = \sum_{i=1}^{n} \alpha_i \sum_{t=1}^{T} \beta_t^{(i)} \mathbf{x}_t^{(i)}$$
(10)

Here:

• α_i : Attention weight for the *i*th temporal window.

- $\beta_t^{(i)}$: Attention weight for the t^{th} frame in the i^{th} window, normalized over frames within that window.
- $\mathbf{x}_t^{(i)}$: The *t*th frame in the *i*th temporal window.

Our N-WindowDCT attention ensures both per-frame and per-window relevances are captured, aligning with the intuition that certain frames within each temporal window may carry more importance for the prediction task. As observed from Figure 4, transitioning from *HisRepDCT* to *OurDualwindow* yields an average 10% improvement in training loss, with a further 12% gain beyond Dual to *N-windows*. However, the observed trend suggests diminishing returns of 4 windows of observations would yield only approximately 6.7% total improvement beyond the 3-window case. Given the memory and processing overhead of tracking multiple windows, the 3-window configuration stands out as the most practical and effective choice. Similarly diminishing returns are reflected in MPJPE, reinforcing this balance between performance and efficiency.



Figure 4: Training loss (a) and MPJPE(b) over 50 epochs for HisRepFixedDCT, DualWindowDCT, and our N-WindowDCT(3 windows).

Inverse DCT and Final Prediction

After applying the GCN, we transform the output back to the time domain using the Inverse DCT (IDCT):

$$\mathbf{X}_{\text{pred}} = \text{IDCT}(\mathbf{X}_{\text{GCN}}) \tag{11}$$

This produces the final predicted motion sequence, incorporating both temporal and spatial dependencies.

NeFT-Net



Figure 5: From left to right, a plot visualization of the Mean Per Joint Position Error (MPJPE) across 72 frames for training on History Repeats Itself DCT, multi-window extended DCT, and *n*-window DCT encoded motion sequences. Note: The red vertical lines start and end of the foot placement cycle.

40

50

60

70

5. ControlNet with Depth Maps for Motion Attention Visualization

30

A powerful recent development arises where Stable Diffusion can be enhanced with ControlNet [30] to provide greater control over image generation. ControlNet allows for the incorporation of additional conditions, such as human pose and depth maps, to guide the generation process. This capability is particularly useful for visualizing motion, as depth maps can capture the spatial relationships between different body parts and their environment. ComfyUI ¹, a GUI-based Stable Diffusion interface, provides a user-friendly environment for composing images with this approach.

10

20

0

80

¹https://github.com/comfyanonymous/ComfyUI

First Author et al.: Preprint submitted to Elsevier

We focus on leveraging depth maps rendered from Blender as a guiding input to ControlNet, enabling precise and realistic depictions of both ground truth and predicted motions. By combining depth-based conditioning with the generative power of stable diffusion, this approach bridges the gap between data-driven motion prediction and its compelling visual representation, as shown in Figure 6, offering a unique perspective on how AI can translate abstract motion data into vivid, interpretable renders.

The core technique is to use depth maps as the control input for ControlNet as seen in Figure 1. By feeding a sequence of depth maps extracted from a video or generated from a simulated environment into ControlNet, we can guide the generation of a corresponding sequence of images that visualize the motion depicted in the depth maps. This approach offers several potential advantages:

- Enhanced Realism: The generative AI imagery is effortlessly realistic. Our prompting approach simply described the style of dress and context of walking, marching, etc. in a graphical depiction. Some orientation terms for example, from left to right assisted the success rate of more oriented diagrammatic results, but weren't as influential as combining all pose frames side-by-side in producing coherent outcomes.
- Precise Control: ControlNet's ability to precisely control the generation process allows for fine-tuning the visualization based on the depth information. Our experiments supplying an alternative 2D bone hierarchy stick representation directly to ControlNet proved to be far less controllable than the more information rich depth representation.
- Novel Visualizations: The combination of ControlNet and depth maps opens up possibilities for creating novel and abstract visualizations of motion.

To visualize the motion of walking of our attention mechanism, we used Blender to generate depth maps of a character model at different stages of the walking cycle. These depth maps can then be used as input to ControlNet, along with text prompts describing the desired motion, to generate images that accurately depict the character's movement. The prompting strategy used in this study was deliberately minimal, primarily to maintain consistent orientation and scene composition(e.g., "person walking forward, side view, consistent lightning"). This simplicity ensured camera alignment across various frames but limited the generative detail in body articulation, clothing variation, and scene interaction.

All frames shown in Figure 1 were generated together in a combined single diffusion pass, with depth maps concatenated to reinforce temporal coherence across the motion sequence. This batch conditioning approach helped maintain the consistency of lightening, background, and carbon appearance, which are often challenges in frame-by-frame generation.



Figure 6: From left to right: (a) NeFT-Net predicted keypoints, visualized alongside ground truth (GT) motion sequences in red and green outlines. (b) The predicted keypoints are aligned and retargeted, starting with BVH format (left) and mapped onto an SMPL-X mesh (right) using Rokoko Studio and Blender (c) Depth maps (top) are refined using stable diffusion to produce photorealistic rendered motion sequences (bottom)

In experimentation of this approach we naturally also applied a ControlNet model steered from Open Pose² derived bone hierarchy images, but found the information of such skeletal wire frame pose images led to excessive ambiguity and relatively poor posed generative image results when compared with the more information rich depth image controlled generations.

6. Conclusion

In this paper, we introduced an N-window-based motion attention model that leverages historical pose information based on the similarity between the current pose context and cyclic sub-sequences in the motion history. Our approach achieves state-of-the-art performance in predicting rhythmic motion by re-timing the Human3.6M dataset using foot anchor placements. Experimental results show strong generalization to previously unseen *walking* and *walking-together* sequences, as indicated by the training loss in Figure 4a and MPJPE in Figure 4b, demonstrating improved joint pose accuracy. Our N-Window extended frequency transformer model aligns ideally upon three historical windows, arrived at due to the observation of a clear trend of diminishing returns in both training loss and predictive accuracy. Quantitative power regression analysis of results in Figure 4b indicate that while the shift from a fixed representation to a dual-window model provides a substantial performance boost, and adding a third window slice contributes an even more meaningful improvement. Beyond this, however, predicted gains taper off progressively: estimating improvements less than 0.1% on successive windows beyond the 10th window. These reducing potential gains come at the cost of increased memory computation, storage and runtime complexity. This is coupled with the lower practical consideration of a motion pattern 10 cycles ago being as relevant to the current cycle in all but a regimented repeated march. We therefore consider three-windows both effective and efficient—avoiding unnecessary overhead while retaining strong predictive accuracy.

Whilst our analysis of varying windows of attention is a form of ablation study itself, we also compared re-timed and non-re-timed data preparations between the non-re-timed HistRepeatDCT method, and the Dual and N-Windowed approaches. Ablations of replacement or simplification of GCN and DCT elements could further indicate the relative importance of each of these measures, including further dissection of the model pipeline—such as isolating the roles of DCT encoding, the re-timing strategy, and window-based attention—as well as exploring the impact of window size for varying periodic motions and offset through hyperparameter sensitivity analysis.

We introduced the use of generative AI techniques to visualize predicted motions, which revealed the model's strong temporal consistency, particularly in sequential foot for placements—a core feature of rhythmic motion. While articulation of hands and facial expressions remain limited by the generative pipeline used, higher-fidelity synthesis approaches may offer future improvements.

Although real-time performance was not the primary target, our approach demonstrates inference speeds that are compatible with near-interactive rates and strong opportunities for optimization. Profiling and benchmarking will be key to validating deployment in time-sensitive scenarios. Moving forward, enhancing the model's real-time capabilities will be a priority—especially for interactive applications such as dance [28] and performance animation, where timing and rhythm are crucial.

7. Acknowledgment

References

- Ademola, A., Sinclair, D., Koniaris, B., Hannah, S., Mitchell, K., 2024. DeFT-Net: Dual-window extended frequency transformer for rhythmic motion prediction, in: The 42nd Eurographics UK Conference on Computer Graphics Visual Computing Conference, CVGC 2024.
- [2] Aksan, E., Kaufmann, M., Cao, P., Hilliges, O., 2021. A spatio-temporal transformer for 3d human motion prediction, in: 2021 International Conference on 3D Vision (3DV), IEEE. pp. 565–574.
- [3] Aksan, E., Kaufmann, M., Hilliges, O., 2019. Structured prediction helps 3d human motion modelling, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 7144–7153.
- [4] Brand, M., Hertzmann, A., 2000. Style machines, in: Proceedings of the 27th annual conference on Computer graphics and interactive techniques, ACM Press/Addison-Wesley Publishing Co., USA. pp. 183–192. URL: https://dl.acm.org/doi/10.1145/344779. 344865, doi:10.1145/344779.344865.
- [5] Butepage, J., Black, M.J., Kragic, D., Kjellstrom, H., 2017. Deep representation learning for human motion prediction and classification, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 6158–6166.

²https://github.com/CMU-Perceptual-Computing-Lab/openpose

- [6] Cai, Y., Huang, L., Wang, Y., Cham, T.J., Cai, J., Yuan, J., Liu, J., Yang, X., Zhu, Y., Shen, X., et al., 2020. Learning progressive joint propagation for human motion prediction, in: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VII 16, Springer. pp. 226–242.
- [7] Cao, Z., Gao, H., Mangalam, K., Cai, Q.Z., Vo, M., Malik, J., 2020. Long-term human motion prediction with scene context, in: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16, Springer. pp. 387–404.
- [8] Chen, L.H., Zhang, J., Li, Y., Pang, Y., Xia, X., Liu, T., 2023. Humanmac: Masked motion completion for human motion prediction, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 9544–9555.
- [9] Chiu, H.k., Adeli, E., Wang, B., Huang, D.A., Niebles, J.C., 2019. Action-agnostic human pose forecasting, in: 2019 IEEE winter conference on applications of computer vision (WACV), IEEE. pp. 1423–1432.
- [10] Cui, Q., Sun, H., Yang, F., 2020. Learning dynamic relationships for 3d human motion prediction, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 6519–6527.
- [11] Dang, L., Nie, Y., Long, C., Zhang, Q., Li, G., 2021. Msr-gcn: Multi-scale residual graph convolution networks for human motion prediction, in: Proceedings of the IEEE/CVF international conference on computer vision, pp. 11467–11476.
- [12] Fragkiadaki, K., Levine, S., Felsen, P., Malik, J., 2015. Recurrent network models for human dynamics, in: Proceedings of the IEEE international conference on computer vision, pp. 4346–4354.
- [13] Fu, J., Yang, F., Dang, Y., Liu, X., Yin, J., 2023. Learning constrained dynamic correlations in spatiotemporal graphs for motion prediction. IEEE Transactions on Neural Networks and Learning Systems.
- [14] Ghosh, P., Song, J., Aksan, E., Hilliges, O., 2017. Learning human motion models for long-term predictions, in: 2017 International Conference on 3D Vision (3DV), IEEE. pp. 458–466.
- [15] Gui, L.Y., Wang, Y.X., Liang, X., Moura, J.M., 2018. Adversarial geometry-aware human motion prediction, in: Proceedings of the european conference on computer vision (ECCV), pp. 786–803.
- [16] Guinot, L., Matsumoto, R., Iwata, H., 2023. Stacked dual attention for joint dependency awareness in pose reconstruction and motion prediction, in: ICAT-EGVE 2023 - International Conference on Artificial Reality and Telexistence and Eurographics Symposium on Virtual Environments.
- [17] Hernandez, A., Gall, J., Moreno-Noguer, F., 2019. Human motion prediction via spatio-temporal inpainting, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 7134–7143.
- [18] Ionescu, C., Papava, D., Olaru, V., Sminchisescu, C., 2013. Human3. 6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. IEEE transactions on pattern analysis and machine intelligence 36, 1325–1339.
- [19] Jain, A., Zamir, A.R., Savarese, S., Saxena, A., 2016. Structural-rnn: Deep learning on spatio-temporal graphs, in: Proceedings of the ieee conference on computer vision and pattern recognition, pp. 5308–5317.
- [20] Kapoulkine, A., 2015. Approximating slerp. URL: https://zeux.io/2015/07/23/approximating-slerp/.
- [21] Koniaris, B., Sinclair, D., Mitchell, K., 2024. DanceMark: An open telemetry framework for latency-sensitive real-time networked immersive experiences, in: 2024 IEEE Conference on Virtual Reality and 3D User Interfaces Abstracts and Workshops (VRW), IEEE. pp. 462–463.
- [22] Li, C., Zhang, Z., Lee, W.S., Lee, G.H., 2018. Convolutional sequence to sequence model for human dynamics, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 5226–5234.
- [23] Mao, W., Liu, M., Salzmann, M., 2020. History repeats itself: Human motion prediction via motion attention, in: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIV 16, Springer. pp. 474–489.
- [24] Mao, W., Liu, M., Salzmann, M., Li, H., 2019. Learning trajectory dependencies for human motion prediction, in: Proceedings of the IEEE/CVF international conference on computer vision, pp. 9489–9497.
- [25] Martinez, J., Black, M.J., Romero, J., 2017. On human motion prediction using recurrent neural networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 2891–2900.
- [26] Pavllo, D., Grangier, D., Auli, M., 2018. Quaternet: A quaternion-based recurrent model for human motion, in: British Machine Vision Conference (BMVC).
- [27] Sidenbladh, H., Black, M.J., Sigal, L., 2002. Implicit probabilistic models of human motion for synthesis and tracking, in: Computer Vision—ECCV 2002: 7th European Conference on Computer Vision Copenhagen, Denmark, May 28–31, 2002 Proceedings, Part I 7, Springer. pp. 784–800.
- [28] Sinclair, D., Ademola, A.V., Koniaris, B., Mitchell, K., 2023. DanceGraph: A complementary architecture for synchronous dancing online, in: 36th International Computer Animation Social Agents (CASA) 2023.
- [29] Taylor, G.W., Hinton, G.E., Roweis, S., 2006. Modeling Human Motion Using Binary Latent Variables, in: Advances in Neural Information Processing Systems, MIT Press. URL: https://proceedings.neurips.cc/paper/2006/hash/ 1091660f3dff84fd648efe31391c5524-Abstract.html.
- [30] Zhang, L., Rao, A., Agrawala, M., 2023. Adding conditional control to text-to-image diffusion models, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 3836–3847.

Title Page Template

Title: NeFT-Net: N-window Extended Frequency Transformer for Rhythmic Motion Prediction

(Article title. Article titles should be concise and informative. Please avoid abbreviations and formulae, where possible, unless they are established and widely understood, e.g., DNA).

Author Information

Author names: Adeyemi Ademola, David Sinclair, Babis Koniaris, Samantha Hannah, Kenny Mitchell

(Provide the given name(s) and family name(s) of each author. The order of authors should match the order in the submission system. Carefully check that all names are accurately spelled. If needed, you can add your name between parentheses in your own script after the English transliteration.)

Affiliations:

Adeyemi Ademola

Edinburgh Napier University School of Computing, Engineering & Built environment, Merchiston Campus, Edinburgh EH10 5DT, United Kingdom Email: <u>adeyemi.ademola@napier.ac.uk</u>

David Sinclair

Edinburgh Napier University School of Computing, Engineering & Built environment, Merchiston Campus, Edinburgh EH10 5DT, United Kingdom Email: <u>D.sinclair@napier.ac.uk</u>

Babis Koniaris

Heriot-Watt University School of Mathematical & Computer Science, Edinburgh Campus, Edinburgh EH14 4AP, United Kingdom Email: <u>B.koniaris@hw.ac.uk</u>

Samantha Hannah

Edinburgh Napier University School of Computing, Engineering & Built environment, Merchiston Campus, Edinburgh EH10 5DT, United Kingdom Email: <u>S.hannah@napier.ac.uk</u>

Kenny Mitchell

Edinburgh Napier University, School of Computing, Engineering & Built environment, Merchiston Campus, Edinburgh EH10 5DT, United Kingdom Email: <u>K.Mitchell2@napier.ac.uk</u>

Title Page Template

(Add affiliation addresses, referring to where the work was carried out, below the author names. Indicate affiliations using a lower-case superscript letter immediately after the author's name and in front of the corresponding address. Ensure that you provide the full postal address of each affiliation, including the country name and, if available, the email address of each author.)

Corresponding author:

Kenny Mitchell

Edinburgh Napier University, School of Computing, Engineering & Built environment, Merchiston Campus, Edinburgh EH10 5DT, United Kingdom Email: <u>K.Mitchell2@napier.ac.uk</u>

(Clearly indicate who will handle correspondence for your article at all stages of the refereeing and publication process and postpublication. This responsibility includes answering any future queries about your results, data, methodology and materials. It is important that the email address and contact details of your corresponding author are kept up to date during the submission and publication process).

For more information, please refer to the relevant sections under submission guidelines for the journal in the Guide for Authors.

Title Page Template

Title: NeFT-Net: N-window Extended Frequency Transformer for Rhythmic Motion Prediction

(Article title. Article titles should be concise and informative. Please avoid abbreviations and formulae, where possible, unless they are established and widely understood, e.g., DNA).

Author Information

Author names: Adeyemi Ademola, David Sinclair, Babis Koniaris, Samantha Hannah, Kenny Mitchell

(Provide the given name(s) and family name(s) of each author. The order of authors should match the order in the submission system. Carefully check that all names are accurately spelled. If needed, you can add your name between parentheses in your own script after the English transliteration.)

Affiliations:

Adeyemi Ademola

Edinburgh Napier University School of Computing, Engineering & Built environment, Merchiston Campus, Edinburgh EH10 5DT, United Kingdom Email: <u>adeyemi.ademola@napier.ac.uk</u>

David Sinclair

Edinburgh Napier University School of Computing, Engineering & Built environment, Merchiston Campus, Edinburgh EH10 5DT, United Kingdom Email: <u>D.sinclair@napier.ac.uk</u>

Babis Koniaris

Heriot-Watt University School of Mathematical & Computer Science, Edinburgh Campus, Edinburgh EH14 4AP, United Kingdom Email: <u>B.koniaris@hw.ac.uk</u>

Samantha Hannah

Edinburgh Napier University School of Computing, Engineering & Built environment, Merchiston Campus, Edinburgh EH10 5DT, United Kingdom Email: <u>S.hannah@napier.ac.uk</u>

Kenny Mitchell

Edinburgh Napier University, School of Computing, Engineering & Built environment, Merchiston Campus, Edinburgh EH10 5DT, United Kingdom Email: <u>K.Mitchell2@napier.ac.uk</u>

Title Page Template

(Add affiliation addresses, referring to where the work was carried out, below the author names. Indicate affiliations using a lower-case superscript letter immediately after the author's name and in front of the corresponding address. Ensure that you provide the full postal address of each affiliation, including the country name and, if available, the email address of each author.)

Corresponding author:

Kenny Mitchell

Edinburgh Napier University, School of Computing, Engineering & Built environment, Merchiston Campus, Edinburgh EH10 5DT, United Kingdom Email: <u>K.Mitchell2@napier.ac.uk</u>

(Clearly indicate who will handle correspondence for your article at all stages of the refereeing and publication process and postpublication. This responsibility includes answering any future queries about your results, data, methodology and materials. It is important that the email address and contact details of your corresponding author are kept up to date during the submission and publication process).

For more information, please refer to the relevant sections under submission guidelines for the journal in the Guide for Authors.

1) A re-timed motion with supervised foot anchor information of periodic cycles, such as walking, for the defined use case of rhythmic motion prediction

2) An improved overall mean per joint position error (MPJPE) results compared to state-of-the-art methods in experiments on the Human3.6M dataset for forecasting short and long-term motions by introducing MultiWindowDCT attention aligned on a best fit period of each motion sequence.

3) A strategy for photorealistic visualization of human body motion sequences by employing the use of stable diffusion on depth maps

4) An open-source version of our code implementation available at https://github.com/CarouselDancing/NeFT-net

All authors disclosed no relevant relationships.