

Rapid-Motion-Track: Markerless Tracking of Fast Human Motion with Deeper Learning

Renjie Li, Chun Yu Lao, Rebecca St. George,
Katherine Lawler, Saurabh Garg, Son N. Tran, Quan Bai, Jane Alty

February 17, 2023

Abstract

Objective: The coordination of human movement directly reflects function of the central nervous system. Small deficits in movement are often the first sign of an underlying neurological problem. The objective of this research is to develop a new end-to-end, deep learning-based system, Rapid-Motion-Track (RMT) that can track the fastest human movement accurately when webcams or laptop cameras are used.

Materials and Methods: We applied RMT to finger tapping, a well-validated test of motor control that is one of the most challenging human motions to track with computer vision due to the small keypoints of digits and the high velocities that are generated. We recorded 160 finger tapping assessments simultaneously with a standard 2D laptop camera (30 frames/sec) and a high-speed wearable sensor-based 3D motion tracking system (250 frames/sec). RMT and a range of DLC models were applied to the video data with tapping frequencies up to 8Hz to extract movement features.

Results: The movement features (e.g. speed, rhythm, variance) identified with the new RMT system exhibited very high concurrent validity with the gold-standard measurements (97.3% of RMT measures were within +/-0.5Hz of the Optotrak measures), and outperformed DLC and other advanced computer vision tools (around 88.2% of DLC measures were within +/-0.5Hz of the Optotrak measures). RMT also accurately tracked a range of other rapid human movements such as foot tapping, head turning and sit-to-stand movements.

Conclusion: With the ubiquity of video technology in smart devices, the RMT method holds potential to transform access and accuracy of human movement assessment.

1 Background and Significance

Assessment of human movement is critical to neuroscience as well as a broad range of other fields including sport science and rehabilitation medicine as it allows interrogation of brain networks, and measurement of the response to interventions [1, 2, 3]. Currently, movement features are evaluated by a clinician or researcher, using validated rating scales that provide ordinal scores for individual components such as speed and amplitude, or composite scores of these combined together [4, 5, 6]. These methods are subjective, expensive (as it requires at least one rater, often a medical specialist, per participant) and imprecise, with considerable inter-rater variability [7]. To reduce bias, multiple raters are recommended to rate one participant's movement [7] but this is labor dependent, time consuming, and lacks the granularity of a continuous measure. Wearable sensors provide accuracy, but they require specialist equipment and tend to be expensive. Furthermore they cannot be applied at a population level for epidemiological or health screening purposes, and are largely inaccessible for people in remote areas. No wearable sensor methods have yet been routinely incorporated into standard clinical practice.

Due to the ubiquity of cameras in personal smart devices worldwide, deep learning-based computer vision methods applied to video recordings of human movement is a promising solution to overcome these limitations. However, current computer vision methods require cameras with a high sampling frequency that produce high resolution videos, and this limits the use of computer vision techniques with standard cameras in laptops or desktop webcams that have relatively low sampling rates. This is problematic as laptops and desktops are the most common method used in telemedicine [8] which has boomed in use since the COVID-19 pandemic.

The Finger Tapping Test (FTT) is a well validated test of human movement function. It is used clinically and in neuroscience research studies as a measure of fine motor control. Participants are required to tap their index finger against their thumb repeatedly and usually instructed to tap ‘as big and fast as possible’. The FTT is used to assess motor function across neuroscience research and an array of neurological disorders including the two most common neurodegenerative disorders: Parkinson’s disease [9] and Alzheimer’s disease [10]. With ageing populations, these two conditions already affect 61 million people worldwide and are predicted to affect 165 million by 2050 [11], so there is growing need for objective methods to evaluate the FTT.

Our recent work, was the first to assess the performance of currently available cutting-edge computer vision models for extracting FTT movement features from relatively low frame per second (fps) web-cameras. This demonstrated that DeepLabCut (DLC) and other similar computer vision models (8 were assessed in total) [8] were able to reliably track finger tapping movements up to 4Hz. However, when FTT frequencies were above 4Hz, the motion blur on videos prevents accurate keypoint detection and renders DLC and other computer vision methods inaccurate. This means that whilst DLC techniques can track slower tapping speeds of patients with obvious motor impairment (e.g. in more advanced Parkinson’s), it would not be able to measure subtle changes in fast movements that are generated by young adult participants or by those in the pre-motor, or early, stages of Parkinson’s and other neurological disorders, where movements are generally above 4Hz. This severely limits the application of computer vision methods in real world settings. Thus, there remains an urgent unmet need for a computer vision method that can accurately measure human hand movements, and other fast human motions, using standard webcam technologies. This would transform our ability to remotely assess humans in their own homes for both clinical and research purposes. With the global COVID-19 pandemic, this need is even more pressing.

To overcome these shortcomings of DLC and other advanced computer vision techniques, we have developed a new system, called Rapid-Motion-Track (RMT) that can extract accurate features of fast human movements from standard (relatively low) 30fps laptop cameras. We designed and conducted a number of experiments whereby participants completed the FTT at a range of metronome- and self-paced frequencies. Tapping movements were simultaneously recorded with a standard laptop camera and high speed 3D wearable sensors (Optotrak, 250 fps). Features determined through RMT applied to the video recordings were validated against the gold standard, wearable sensor method and compared to DLC and other deep learning methods.

Our main contribution is that the new RMT system can extract valid and accurate fast human movement features using relatively low fps cameras that are superior to current DLC and other computer vision models. Sub-contributions are:

- We show RMT outperforms DLC and other state-of-the-art computer vision methods in two challenging cases: 1) tapping frequencies from 0.5Hz to 6Hz and 2) videos in low resolution (256×256) and low frame rates (30 fps).
- We demonstrate that RMT has robust tracking across a range of human movements used in motor control assessments including sit-to-stand, head turning, foot tapping and leg agility tests.

2 Materials and Methods

2.1 Data Collection

The data collection process has been described in detail previously [8]. In brief, participants sat facing a laptop camera with an Optotrak high speed 3D camera behind them. The Optotrak system (Northern Digital Inc.) sensors (infrared Light Emitting Diodes) were attached on participant’s index fingertip and thumb-tip (both on the back of the hand). Sixteen participants performed FTT for 20 seconds under 5 conditions (0.5 Hz, 1Hz, 2Hz, 3Hz and maximal speed (which was typically in the range 5-8 Hz)). The laptop camera recorded 2D videos and the Optotrak system recorded the real life positions of the two sensors. Further details can be found in our previous work [8].

Each video contains around 600 frames (20×30) and 20 frames were extracted from each video based on K-means clustering algorithm (K=10) for manually labelling the index fingertip and thumb-tip position. Total 4,400 (with finger digit position being labelled) frames were used for evaluating the

digit tracking result. The hand movement features calculated based on Optotrak sensors were used as the feature ground truths for evaluating the feature extraction result.

2.2 RapidMotionTrack System

RapidMotionTrack takes 2D FTT videos as the input and the outputs are hand kinematic features. This system consists of three modules, i.e. Fingertip Tracking, Adaptive Vertex Recognition and Feature Extraction, shown in Figure 1. The Fingertip Tracking module adapts our previous work P-MSDSNet [12] to track thumb-tip and index fingertip positions on the 2D finger tapping videos. The role of the Adaptive Vertex Recognition module is to smooth the distance-versus-time curve adaptively based on the dominant tapping frequency of a participant to accurately localize the time point of tapping peaks and tapping valleys. The Feature Extraction module takes charge of extracting different hand movement features that would be useful for neuroscience research.

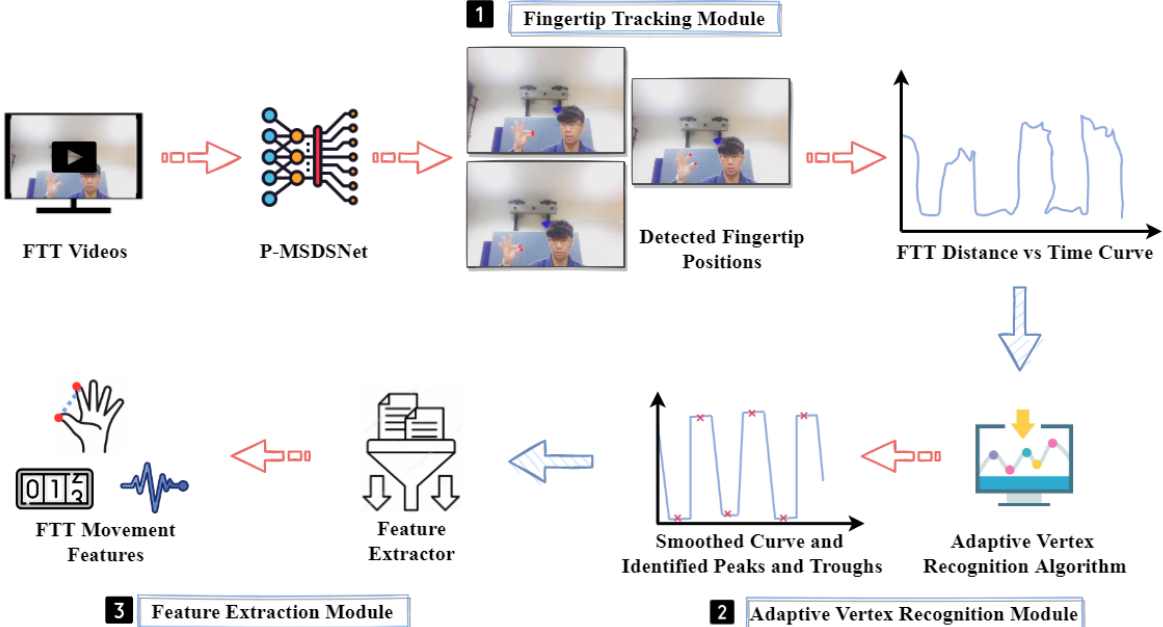


Figure 1: An illustration of RapidMotionTrack system with three modules, including Fingertip Tracking Module, Adaptive Vertex Recognition Module and Feature Extraction Module.

2.2.1 Fingertip Tracking Module

The Fingertip Tracking module takes 2D finger tapping videos as input and outputs a distance-versus-time curve graph. The distance is measured between the thumb-tip and the index fingertip. We applied and adjusted our previous work P-MSDSNet [12] to track thumb-tip and index fingertip positions on each frame of the video. P-MSDSNet is a multi-scale neural network which can learn abundant features from high resolution feature maps to low resolution feature maps in a cyclical and cascade pattern. Due to P-MSDSNet’s deep supervision-based spatial attention mechanism on different scale levels of the input images, the network is able to extract discriminable features to effectively predict the locations of the thumb tip and the index fingertip.

In the Fingertip Tracking module, we employed a P-MSDSNet with a stack of neural network blocks, each consisting of 5 different scaling operators. The multi-scale features are fused together in an up-down manner and are propagated to deeper levels under a deep supervision mechanism to increase the prediction accuracy. For FTT, some frames may have motion blur areas due to the very quick finger tapping. To mitigate the impact of the motion blur on the fingertips detection, we improve the final prediction stage network by taking into account the learned feature maps from different scales at the last neural network block in the stack (Figure 2).

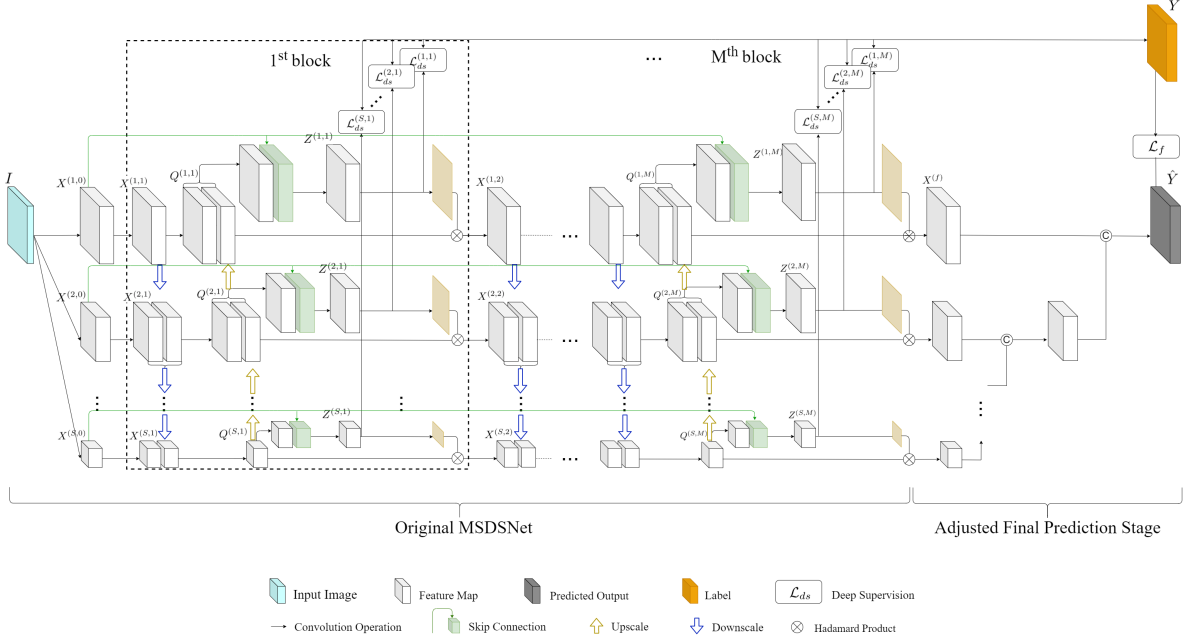


Figure 2: An illustration of adjusted P-MSDSNet used in RapidMotionTrack. RapidMotionTrack applies 3 up-downscale blocks ($M = 3$) and 5 different scale sizes ($S = 5$). It also adds a final prediction stage to fuse the outputs from different scale feature maps together to make the final prediction.

In the network forward propagation process, the original frame is transformed to 5 ($S = 5$) different scale feature maps (from $X^{(1,0)}$ to $X^{(5,0)}$) through some convolutional operators. Then, along each scale, feature maps propagate information forward in a parallel manner, meanwhile, fuse different information through down-up scale connection blocks between adjacent scales. For example, $X^{(1,1)}$ feature map not only propagates forward along its own scale path, but also does the downscale operation and fuse with $X^{(2,1)}$ feature map. Then, $X^{(2,1)}$ feature map propagates forward after it fuses the downsampled information from $X^{(1,1)}$. In this case, information is spread not only along its self scale path but also across different scales paths in a downscale-upscale fusing mechanism. We apply 3 downscale-upscale blocks ($M = 3$) and add deep supervision module at the end of each block [12]. In the final prediction stage, we fuse the final feature map from the smallest scale to largest scale step by step using convolutional operator and concatenating operator. At last, the final prediction of heatmap is inferred at the largest scale \hat{Y} (the same scale with the original frame).

In the network training process, twenty frames were selected from each of the 220 videos by K-means ($K = 10$) clustering algorithm [13] for manually annotating the positions of thumb-tip and index fingertip. Then, 4,400 annotated frames ($20 \text{ frames} \times 220 \text{ videos}$) were randomly split into training dataset (95%) and testing dataset (5%). We trained different deep learning models on training dataset and compared the performance on testing dataset. The annotated fingertip positions on each frame is transposed to a $H \times W \times 2$ heatmap, where H and W refer to the height and width of the frame and 2 refers to the number of fingertips (in this case, thumb-tip and index fingertip). The heatmap is generated by a Gaussian function with $\mu = \begin{bmatrix} x \\ y \end{bmatrix}$ and $\Sigma = \begin{bmatrix} 3 & 3 \\ 3 & 3 \end{bmatrix}$ for fingertip, (x, y) is the position of the fingertip. For training, Mean Square Error (MSE) between the true fingertip position heatmap and predicted heatmap is used as loss function and Adam algorithm [14] is used as learning optimizer. The trained model will be used on each frame of each video to track the positions of thumb-tip and index fingertip.

After obtaining the positions of thumb-tip and index fingertip on each frame, we calculate the Euclidean distance between thumb-tip and index fingertip (in pixels) on each frame and draw the distance-versus-time curve graph.

2.2.2 Adaptive Vertex Recognition Module

Since the predicted fingertip position covers only one pixel area of the frame, even when the performance of fingertip tracking is good, it is inevitable that the distance vs time curve will not be smooth (this can be caused by little tremor when participant does the finger tapping). Unsmooth curve increases the difficulty in identifying peaks and troughs, which further leads to inaccurate feature extraction. To solve this issue, we developed an Adaptive Vertex Recognition algorithm to identify peaks and troughs along the distance-versus-time curve graph adaptively by filtering out rough areas.

Firstly, we use the distance difference of the adjacent frames ΔS as a criteria to remove the fluctuation of the signal.

$$\Delta S' = \begin{cases} 0 & \text{when } \Delta S < \gamma_{flatness}R \\ \Delta S & \text{when } \Delta S \geq \gamma_{flatness}R \end{cases} \quad (1)$$

where S is the distance signal after time-averaged mean removal, S' is the curve after fluctuation removal, Δ represents the difference between adjacent frames of specified signal, $\gamma_{flatness}$ is the threshold which set as 0.1 in the present study, and R is the range of the interested signal. The distance signal before (ΔS) and after fluctuation removal ($\Delta S'$) are plotted side-by-side in Figure 3 (left).

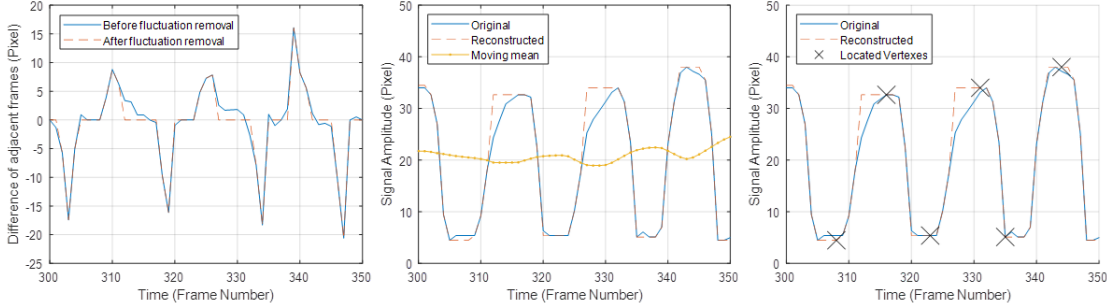


Figure 3: Adaptive vertex recognition process.

Secondly, we reconstructed the signal based on $\Delta S'$ and S . The former obtains sharp change between sections (where $\Delta S' = 0$, i.e., finger open or close), the latter guides the magnitude of the reconstructed signal. Moreover, the existence of peaks and troughs can be determined by calculating the moving average value μ over the signal. Peaks are determined when the local section are above μ_i , in contrast, troughs are recognised when the local sections are below.

$$\mu_i = \frac{1}{n} \sum_{i-n/2}^{i+n/2} S'_i \quad (2)$$

$$\text{w.r.t } n = \gamma_{window} \times N$$

where γ_{window} is the relative window size (set as 0.1 in the present study) and N is the total number of frames in the signal. Please note that the duration of the first and last platforms are omitted, since it often is meaningless waiting period. Figure 3 (center) shows the original signal S , reconstructed signal S' , and the moving mean μ .

Due to the fact that the features, such as the time between each tap and open duration of the fingers, of every signal varies significantly. Thus, the cases of short and long opening are both need to be considered and processed separately to recognise the vertexes in a reliable manner throughout all videos, following are the formulas applied.

$$t_{vertex} = \begin{cases} t_{localvertex} & \text{when } l_{section} \leq \gamma_{platform}N \\ t_{median} & \text{when } l_{section} > \gamma_{platform}N \end{cases} \quad (3)$$

where $t_{localvertex}$ is the moment of the local maxima or minima of the section, t_{median} is the central time of the local section, $l_{section}$ is the duration of the local section, and $\gamma_{platform}$ is the constant which set as 0.01 in the present study.

Once the time of local vertexes t_{vertex} are located, the local height S_{vertex} can be found by below:

$$S_{vertex} = \begin{cases} S_{localsection} & \text{when } \Delta S'(t_{vertex}) = 0 \\ S(t_{vertex}) & \text{when } \Delta S'(t_{vertex}) > 0 \end{cases} \quad (4)$$

where $S_{localsection}$ is the maxima or minima of the local section. The sample signal, processed signal, and recognised vertexes are demonstrated in the Figure 3 (right). The curve of recognised vertexes (t_{vertex} , S_{vertex}) are ready for the next step, feature extraction in the following section.

2.2.3 Feature Extraction

The Feature Extraction module takes reconstructed distance (between thumb-tip and index fingertip)-versus-time graph data as the input and different hand kinematic features are the outputs. Based on key neuroscience measures, we summarized speed, amplitude and rhythm related features of the FTT and their calculations. These features are known to deteriorate in Parkinson's and other neurological disorders.

Speed related features include Mean Tapping Frequency (M-TF), Total Tapping Count (TTC), Maximum Speed (MS), Mean Inter Tap Interval (M-ITI) and Decrements on Speed (DoS). Amplitude related features include Coefficient of Variance of Amplitude (COV-A) and Decrements on Amplitude (DoA). Rhythm related features include Coefficient of Variance of Tapping Frequency (COV-TF) and Intra Individual Variability (IIV). Detailed validation results on all 9 features are presented in supplementary document.

$$\text{M-TF} = \frac{1}{K_p - 1} \sum_{k=2}^{K_p} \frac{1}{t_{(k)} - t_{(k-1)}} \quad (5)$$

$$\text{TTC} = \min(K_p, K_v) \quad (6)$$

$$\text{MS} = 1 / \min(t_{(k)} - t_{(k-1)}) \text{ for } k = 2, 3, \dots, K_p \quad (7)$$

$$\text{M-ITI} = \frac{1}{K_v - 1} \sum_{k=2}^{K_v} (t_{(k)} - t_{(k-1)}) \quad (8)$$

$$\text{DoS} = \frac{1}{K_p - 1} \ln\left(\frac{1}{t_{(2)} - t_{(1)}} / \frac{1}{t_{(K_p)} - t_{(K_p-1)}}\right) \quad (9)$$

$$\text{COV-A} = \frac{\sqrt{(a_{(k)} - \frac{1}{K_p} \sum_{k=1}^{K_p} a_{(k)})^2}}{\frac{1}{K_p} \sum_{k=1}^{K_p} a_{(k)}} \quad (10)$$

$$\text{DoA} = \frac{1}{K_p} \ln \frac{a_{(1)}}{a_{(K_p)}} \quad (11)$$

$$\text{COV-TF} = \sqrt{\frac{\sum_{k=2}^{K_p} \left(\frac{1}{t_{(k)} - t_{(k-1)}} - \text{M-TF}\right)^2}{K_p - 1}} / \text{M-TF} \quad (12)$$

$$\text{IIV} = \sqrt{\frac{1}{K_p - 1} \sum_{k=2}^{K_p} [(t_{(k)} - t_{(k-1)}) - \text{M-ITI}]^2} \quad (13)$$

where K_p refers to the number of peaks, K_v refers to the number of valleys, $t_{(k)}$ refers to the time point at k^{th} peak and a_k refers to the normalized amplitude.

2.3 Evaluation Method

We compare the performance of RMT with other state-of-the-art methods, including DLC, to precisely track the thumb-tip and index fingertip during the FTT. The Percentage of Correct Keypoints (PCK) measures the percentage of correctly localized thumb-tip and index fingertip. A correctly localized keypoint is confirmed if the distance (in pixels) between the predicted position and true position is within a pre-set threshold (see equation 14). The Mean of Per Joint Position Error (MPJPE) measures how precisely the predicted position of fingertip can reach compared with the true position (see equation 15). Figure 4 shows the PCK of different methods and Table 1 shows the MPJPE of different methods. RMT achieves the best performance in both two metrics.

$$\text{PCK}_{@T} = \frac{1}{N \times J} \sum_{j=1}^J \sum_{n=1}^N (\|P_n^{(j)} - Y_n^{(j)}\|_2 < T) \quad (14)$$

$$\text{MPJPE} = \frac{1}{N \times J} \sum_{j=1}^J \sum_{n=1}^N \|P_n^{(j)} - Y_n^{(j)}\|_2 \quad (15)$$

where N is the number of frames and J is the number of fingertips to be detected (in this case, $J = 2$), $P_n^{(j)}$ is the predicted position of the j^{th} fingertip on the n^{th} frame, $Y_n^{(j)}$ is the true position of the j^{th} fingertip on the n^{th} frame and T is the pre-set threshold.

3 Results

In this section, we show that RMT achieves accurate digit tracking from 2D videos and accurate movement feature extraction.

3.1 Digit Tracking Results

Digit tracking by RMT is more accurate than other available state-of-the-art methods. This is reflected in a higher Percentage of Correct Keypoints (PCK) and a lower Mean of Per Joint Position Error (MPJPE) on the testing dataset shown in Figure 4 and Table 1. The calculations of PCK and MPJPE are described in the Online Method section.

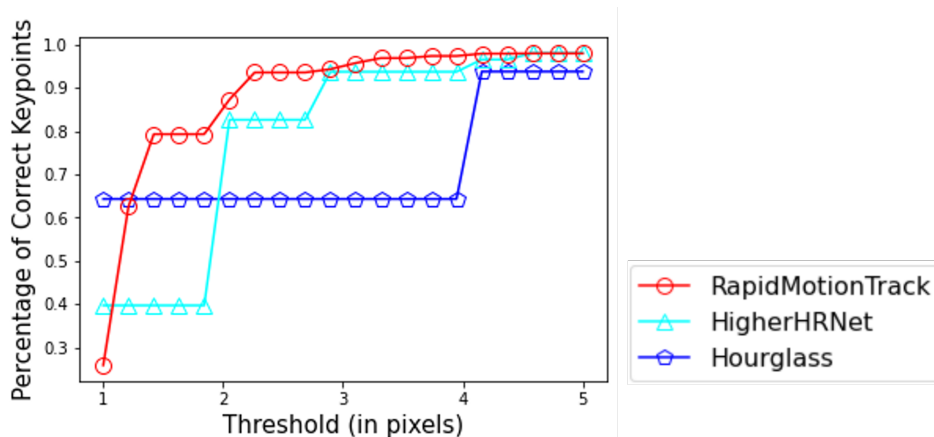


Figure 4: Different models’ PCK on FTT testing dataset (DLC software does not report PCK). This shows that RMT has the highest percentage of detected correct keypoints at different thresholds when compared to the wearable sensor data (ground truth).

Table 1: Different computer vision models’ Mean of Per Joint Position Error on FTT testing dataset

Methods	Thumb-tip	Index Fingertip	Average
RapidMotionTrack	1.10	1.32	1.21
DLC-ResNet50 [9]	-	-	3.00
DLC-MobileNet [9]	-	-	2.30
HigherHRNet [15]	1.47	1.71	1.59
Hourglass [16]	1.37	1.82	1.60

The distance between detected index fingertip and thumb-tip positions normalized to maximum aperture at each frame over time is plotted in Figure 2 (Figure 5). A visual comparison to the wearable sensor method demonstrates that when applied to the unlabelled videos, RMT provides good tracking performance in the ‘maximal speed’ condition (approximately 6Hz).

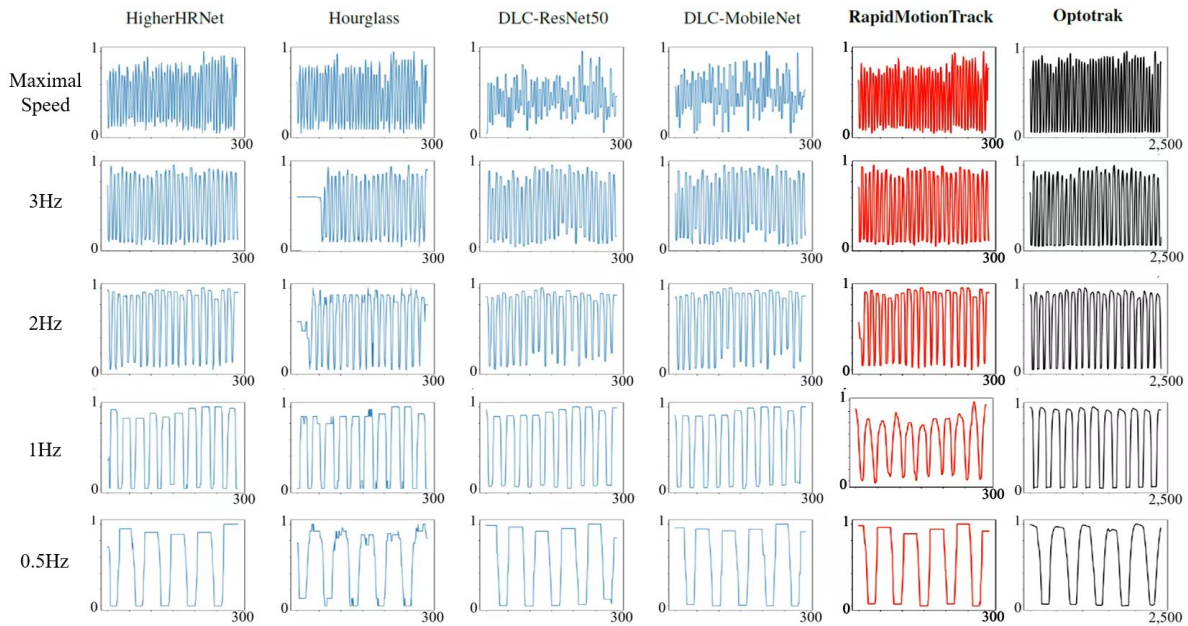


Figure 5: Normalized sample distance vs time graphs calculated from computer vision methods in one individual performing the FTT. The blue graphs show the results from each computer vision method compared to the wearable sensor gold standard method, Optotrak (black graph). The X-axis shows the time (frame by frame where 300 frames for computer vision methods and 2,500 frames for the Optotrak method) and the Y-axis shows normalized distance (measured in pixels for computer vision methods and in mm for the Optotrak method). Each row denotes a different FTT condition, where Row 1 is the self paced maximal speed condition and Rows 2, 3, 4 and 5 are the 3Hz, 2Hz, 1Hz and 0.5Hz metronome-paced conditions respectively.

3.2 Feature Extraction Results

The RMT system, HigherHRNet and Hourglass models can extract valid hand movement features even when participants tap in a high frequency range ($>4\text{Hz}$). In comparison, hand movement speed, amplitude, and rhythm related features calculated based on the DLC system’s digit tracking result are invalid when participants tap more quickly than 4Hz. These are reflected in the comparison between different computer vision methods and the Optotrak in terms of Bland Altman plots, X-Y plots (Figure 6) and statistical tests (Table 2). Details of movement features are explained on the Online Method section.

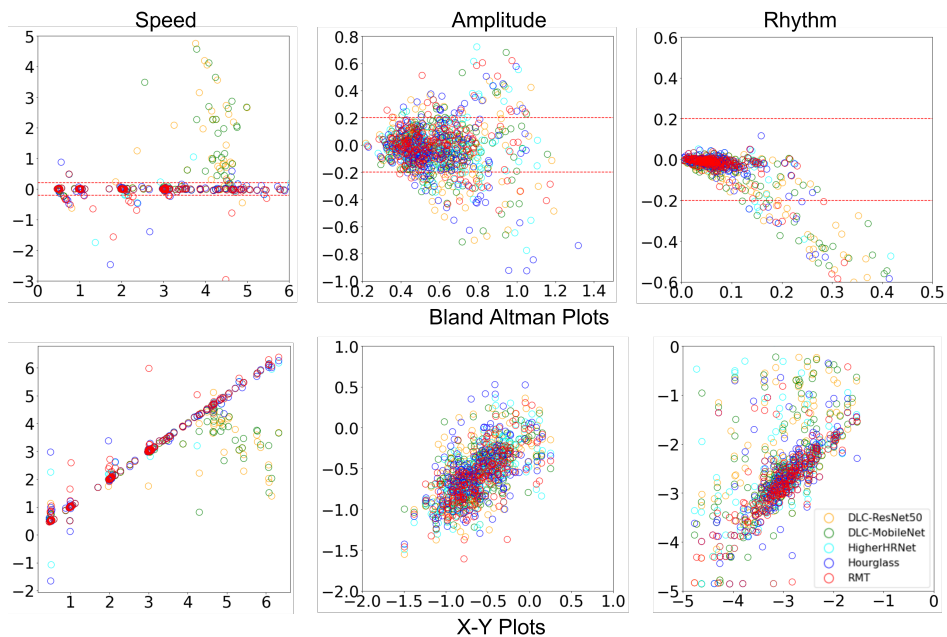


Figure 6: Features comparison: computer vision method (Y-axis) vs Optotrak method (X-axis). The top row presents Bland Altman Plots of different features calculated by computer vision methods (Y-axis) vs the Optotrak method (X-axis). The bottom row presents X-Y Plots of different features calculated by computer vision method (Y-axis) vs the Optotrak method (X-axis). Speed is calculated from the mean of tapping frequency, amplitude is calculated from by the logarithm of the coefficient of variance of normalized tapping amplitude and rhythm is calculated by the logarithm of the coefficient of variance of tapping frequency. Detailed comparisons of other features are shown in the supplementary material.

Table 2: Welch’s t test (at 0.05 significant level) results of different computer vision methods compared with the Optotrak measures on speed, amplitude and rhythm. Results from two methods have no significant difference when $P > 0.05$ (accept null hypothesis), and have significant difference when $P < 0.05$ (reject null hypothesis).

Condition	Computer Vision Method	Speed	Amplitude	Rhythm
Maximal speed (<4Hz)	RMT	accept	accept	accept
	DLC-ResNet50	accept	reject	accept
	DLC-MobileNet	accept	reject	accept
Maximal speed (>4Hz)	RMT	accept	accept	accept
	DLC-ResNet50	reject	reject	accept
	DLC-MobileNet	reject	reject	accept
3Hz	RMT	accept	accept	accept
	DLC-ResNet50	accept	accept	accept
	DLC-MobileNet	accept	accept	accept
2Hz	RMT	accept	accept	accept
	DLC-ResNet50	accept	accept	accept
	DLC-MobileNet	accept	accept	accept
1Hz	RMT	accept	accept	accept
	DLC-ResNet50	accept	accept	accept
	DLC-MobileNet	accept	accept	accept
0.5Hz	RMT	accept	accept	accept
	DLC-ResNet50	accept	accept	accept
	DLC-MobileNet	accept	accept	accept

3.3 Generalisability of RMT to other fast human movements

RMT was also able to accurately track a range of other human movements when applied to webcam videos of head turning, toe tapping [17], leg agility, and sit-to-stand movements [18, 19] as shown in Figure 7. The tracking videos are provided in the supplementary materials.

4 Discussion and Conclusion

4.1 Precise Digit Tracking from 2D Laptop Camera Videos

Digit tracking by RMT is more accurate than available state-of-the-art methods. When participants tap in a fast range (e.g. higher than 4Hz), DLC cannot accurately and reliably track the digits [8]. Additionally, this will further lead to incorrect digit distance vs time graph. Figure 5 demonstrated an individual participant’s distance (between index fingertip and thumb-tip) vs time graphs under 5 conditions of finger tapping calculated from different computer vision methods and the Optotrak method. In the maximal speed condition, motion blur frames were very common. Compared with other deep learning networks to detect keypoints, the distance vs time graph obtained from the RMT system is more stable and closer to the gold standard Optotrak methods. This is due to the P-MSDSNet fusing multi-scale features from both same depth levels and different depth levels in a parallel manner, which can gradually refine the features and detected fingertip positions layer by layer and scale by scale even on a blurred frame. An accurate digit tracking lays the foundation for extracting hand movement features.

4.2 Feature Extraction Discussion

Despite the accurate digit tracking results, the unsmooth distance vs time curve may also affect the accuracy of final feature extraction. To overcome this issue, the proposed Adaptive Vertex Recognition algorithm can help RMT system correctly identify the peaks and troughs from a particular distance vs time graph without being affected by the noise or unsmooth part.



Figure 7: RMT accurately tracks a range of rapid human motions including: head turning, arising from a chair, toe tapping and leg agility. The keypoints are denoted by red dots and the video in supplementary materials demonstrates these keypoints track throughout the motion.

4.3 Strengths, Limitations and Future Direction

A particular strength of RMT is that it takes original finger tapping videos as input, through the Fingertip Tracking Module and Adaptive Vertex Recognition Module, and outputs accurate and valid hand movement features for both assessment and research. This is an additional output compared with DLC which only tracks digits (RMT not only tracks digits but also produces valid movement feature data). In addition, both digit tracking and the feature extraction achieve the best results out of all currently available state of the art computer vision deep learning models and this can serve for not only FTT but also other movement test analysis including sit-to-stand movement test [18, 19] and foot-tapping test [17] from recorded 2D videos. We showed a comparison of RMT with the DLC and the Optotrak system in Table 3; this shows that the features extracted using RMT are equivalent to wearable sensor methods but, in addition, have the potential to be performed remotely (e.g., for telemedicine) and at the population level (due to the wide reach of webcams). Rather than the Optotrak system and DLC whose output is fingertip position, the RMT system can directly output validated finger tapping features.

Table 3: Output data from RMT , DLC and wearable sensor method

	RMT system	DLC	Wearable sensor methods
Finger Digits Tracking	Yes	Yes	Yes
Feature Extraction	Yes	No	Yes
Feature Validity	Yes	No	Yes

Currently RMT is offline but future work will develop this into an online system. Additionally, we will extend RMT that can be used for general human movement feature extraction. First, we will build a model zoo consisting of different models for different movements, and keep updating the zoo. Second, we will make RMT as a windows-based system that can be used in a more efficient manner for neuroscience. RMT will also be used in the TAS Test project [20], which aims to detect the earliest stages of Alzheimer’s disease using hand movement analysis.

References

[1] Ji-Won Kim, Jae-Ho Lee, Yuri Kwon, Chul-Seung Kim, Gwang-Moon Eom, Seong-Beom Koh, Do-Young Kwon, and Kun-Woo Park. Quantification of bradykinesia during clinical finger taps using a gyrosensor in patients with parkinson’s disease. *Medical & biological engineering & computing*, 49(3):365–371, 2011.

- [2] Stefan Schneider, Vera Abeln, Julia Popova, Elena Fomina, Amrei Jacubowski, Romain Meeusen, and Heiko K Strüder. The influence of exercise on prefrontal cortex activity and cognitive performance during a simulated space flight to mars (mars500). *Behavioural brain research*, 236:1–7, 2013.
- [3] Renjie Li, Xinyi Wang, Katherine Lawler, Saurabh Garg, Quan Bai, and Jane Alty. Applications of artificial intelligence to aid detection of dementia: a scoping review on current capabilities and future directions. *Journal of Biomedical Informatics*, page 104030, 2022.
- [4] CG Goetz, BC Tilley, SR Shaftman, G Stebbins, S Fahn, P Martinez-Martin, W Poewe, C Sampaio, M Stern, R Dodel, et al. Updrs revision task force movement disorder society-sponsored revision of the unified parkinson’s disease rating scale (mds-updrs): scale presentation and clinimetric testing results. *Mov. Disord*, 23:2129–2170, 2008.
- [5] Dustin A Heldman, Joseph P Giuffrida, Robert Chen, Megan Payne, Filomena Mazzella, Andrew P Duker, Alok Sahay, Sang Jin Kim, Fredy J Revilla, and Alberto J Espay. The modified bradykinesia rating scale for parkinson’s disease: reliability and comparison with kinematic measures. *Movement Disorders*, 26(10):1859–1863, 2011.
- [6] Tiago A Mestre, Maria João Forjaz, Philipp Mahlknecht, Francisco Cardoso, Joaquim J Ferreira, Ralf Reilmann, Cristina Sampaio, Christopher G Goetz, Esther Cubo, Pablo Martinez-Martin, et al. Rating scales for motor symptoms and signs in huntington’s disease: critique and recommendations. *Movement disorders clinical practice*, 5(2):111–117, 2018.
- [7] Afsaneh Shirani, Braeden D Newton, and Darin T Okuda. Finger tapping impairments are highly sensitive for evaluating upper motor neuron lesions. *BMC neurology*, 17(1):1–8, 2017.
- [8] Renjie Li, Rebecca J St George, Xinyi Wang, Katherine Lawler, Edward Hill, Saurabh Garg, Stefan Williams, Samuel Relton, David Hogg, Quan Bai, et al. Moving towards intelligent telemedicine: Computer vision measurement of human movement. *Computers in Biology and Medicine*, page 105776, 2022.
- [9] Pablo Arias, Verónica Robles-García, Nelson Espinosa, Yoanna Corral, and Javier Cudeiro. Validity of the finger tapping test in parkinson’s disease, elderly and young healthy subjects: Is there a role for central fatigue? *Clinical Neurophysiology*, 123(10):2034–2041, 2012.
- [10] Shota Suzumura, Aiko Osawa, Taishi Nagahama, Izumi Kondo, Yuko Sano, and Akihiko Kandori. Assessment of finger motor skills in individuals with mild cognitive impairment and patients with alzheimer’s disease: Relationship between finger-to-thumb tapping and cognitive function. *Japanese Journal of Comprehensive Rehabilitation Science*, 7:19–28, 2016.
- [11] ER Dorsey, T Sherer, MS Okun, and BR Bloem. The emerging evidence of the parkinson pandemic. *J parkinsons dis*. 2018; 8 (s1): S3-8.
- [12] Renjie Li, Son Tran, Saurabh Garg, Katherine Lawler, Jane Alty, and Quan Bai. Parallel multi-scale networks with deep supervision for hand keypoint detection. *arXiv preprint arXiv:2112.10275*, 2021.
- [13] Alexander Mathis, Pranav Mamidanna, Kevin M Cury, Taiga Abe, Venkatesh N Murthy, Mackenzie Weygandt Mathis, and Matthias Bethge. Deeplabcut: markerless pose estimation of user-defined body parts with deep learning. *Nature neuroscience*, 21(9):1281–1289, 2018.
- [14] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [15] Bowen Cheng, Bin Xiao, Jingdong Wang, Honghui Shi, Thomas S Huang, and Lei Zhang. Higherhrnet: Scale-aware representation learning for bottom-up human pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5386–5395, 2020.
- [16] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *European conference on computer vision*, pages 483–499. Springer, 2016.

- [17] Hayato Enoki, Toshikazu Tani, and Kenji Ishida. Foot tapping test as part of routine neurologic examination in degenerative compression myelopathies: a significant correlation between 10-sec foot-tapping speed and 30-m walking speed. *Spine surgery and related research*, pages 2018–0033, 2019.
- [18] Susan L Whitney, Diane M Wrisley, Gregory F Marchetti, Michael A Gee, Mark S Redfern, and Joseph M Furman. Clinical measurement of sit-to-stand performance in people with balance disorders: validity of data for the five-times-sit-to-stand test. *Physical therapy*, 85(10):1034–1045, 2005.
- [19] Ryan P Duncan, Abigail L Leddy, and Gammon M Earhart. Five times sit-to-stand test performance in parkinson’s disease. *Archives of physical medicine and rehabilitation*, 92(9):1431–1436, 2011.
- [20] Jane Alty, Quan Bai, Renjie Li, Katherine Lawler, Rebecca J St George, Edward Hill, Aidan Bindoff, Saurabh Garg, Xinyi Wang, Guan Huang, et al. The tas test project: a prospective longitudinal validation of new online motor-cognitive tests to detect preclinical alzheimer’s disease and estimate 5-year risks of cognitive decline and dementia. *BMC neurology*, 22(1):1–13, 2022.