

MI-Poser: Human Body Pose Tracking Using Magnetic and Inertial Sensor Fusion with Metal Interference Mitigation

RIKU ARAKAWA*, Carnegie Mellon University, United States

BING ZHOU[†], Snap Research, United States

GURUNANDAN KRISHNAN, Snap Research, United States

MAYANK GOEL, Carnegie Mellon University, United States

SHREE K. NAYAR, Snap Research, United States

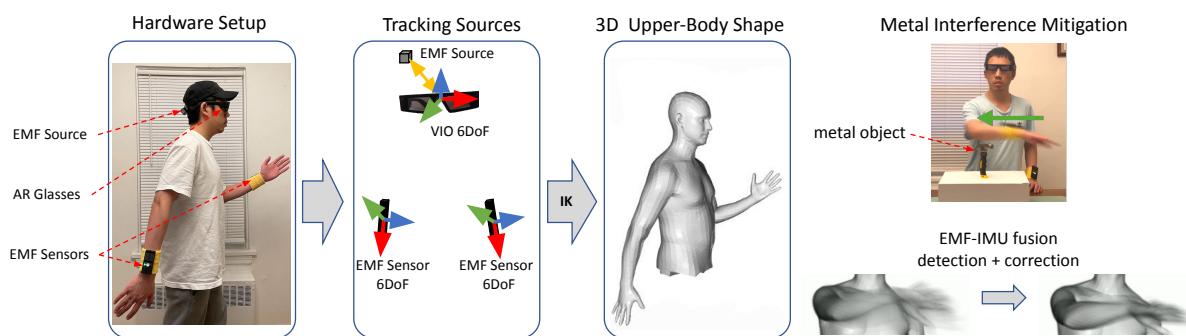


Fig. 1. MI-Poser takes the Visual-Inertial Odometry (VIO) tracking data from AR glasses and two EMF sensors on the wrists as input and generates 3D body shapes through a machine learning model for IK. To tackle the well-known magnetic metal interference issue, we propose Metal Interference Mitigation (MIM) which actively detects and corrects metal interference with EMF-IMU sensor fusion. As a result, the output body movements become steady and have more fidelity.

Inside-out tracking of human body poses using wearable sensors holds significant potential for AR/VR applications, such as remote communication through 3D avatars with expressive body language. Current inside-out systems often rely on vision-based methods utilizing handheld controllers or incorporating densely distributed body-worn IMU sensors. The former limits hands-free and occlusion-robust interactions, while the latter is plagued by inadequate accuracy and jittering. We introduce a novel body tracking system, *MI-Poser*, which employs AR glasses and two wrist-worn electromagnetic field (EMF) sensors to achieve high-fidelity upper-body pose estimation while mitigating metal interference. Our lightweight system demonstrates a minimal error (6.6 cm mean joint position error) with real-world data collected from 10 participants. It remains robust against various upper-body movements and operates efficiently at 60 Hz. Furthermore, by incorporating an IMU sensor co-located with the EMF sensor, MI-Poser presents solutions to counteract the effects of metal interference, which inherently

*The majority of work was done during an internship at Snap Research.

[†]Corresponding author.

Authors' addresses: Riku Arakawa, rarakawa@cs.cmu.edu, Carnegie Mellon University, Pittsburgh, United States; Bing Zhou, bzhou@snapchat.com, Snap Research, New York, United States; Gurunandan Krishnan, gkrishnan@snap.com, Snap Research, New York, United States; Mayank Goel, mayankgoel@cmu.edu, Carnegie Mellon University, Pittsburgh, United States; Shree K. Nayar, snayar@snap.com, Snap Research, New York, United States.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

© 2023 Copyright held by the owner/author(s).

2474-9567/2023/9-ART85

<https://doi.org/10.1145/3610891>

disrupts the EMF signal during tracking. Our evaluation effectively showcases the successful detection and correction of interference using our EMF-IMU fusion approach across environments with diverse metal profiles. Ultimately, MI-Poser offers a practical pose tracking system, particularly suited for body-centric AR applications.

CCS Concepts: • **Human-centered computing** → **Mobile devices**; • **Computing methodologies** → *Computer vision*.

Additional Key Words and Phrases: body pose tracking, inverse kinematics, sensor fusion

ACM Reference Format:

Riku Arakawa, Bing Zhou, Gurunandan Krishnan, Mayank Goel, and Shree K. Nayar. 2023. MI-Poser: Human Body Pose Tracking Using Magnetic and Inertial Sensor Fusion with Metal Interference Mitigation. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 7, 3, Article 85 (September 2023), 24 pages. <https://doi.org/10.1145/3610891>

1 INTRODUCTION

Human motion tracking is integral to augmented reality (AR) and virtual reality (VR). Existing VR devices [20, 35, 54] primarily use cameras in head-mounted displays to track head pose and hand-held controllers for spatial input. However, cameras are power-intensive and impractical for AR glasses. Moreover, their limited field of view can lose track of controllers or hands, constraining user interaction. To achieve occlusion-robust, hands-free tracking, researchers have explored wearable-IMU-based solutions [21, 68–70]. However, these approaches typically require more sensors attached to the body (e.g., waist to track torso movement), and the error can accumulate over time as the IMU sensor does not provide 3D position information directly. As a result, the solutions often suffer from imprecise tracking results.

In this paper, we introduce *MI-Poser*, an upper-body pose-tracking system utilizing magnetic tracking in wristbands and AR glasses, as depicted in Figure 1. MI-Poser incorporates an electromagnetic field (EMF) source in AR glasses and two wrist-worn EMF sensors, enabling 6-DoF wrist tracking relative to the head. Combined with AR glasses' Visual-Inertial Odometry (VIO) tracking, MI-Poser can track 6-DoF poses for head and wrists. We trained deep neural networks for human pose inverse kinematics (IK) on a large dataset (AMASS [34]) to reconstruct upper-body pose from sparse signals. Importantly, EMF sensing often suffers from interference from nearby metallic objects. Thus, we propose *metal interference mitigation* (MIM) to enhance the input data of pose reconstruction by utilizing a collocated IMU sensor. MIM detects metal interference on EMF sensors and actively corrects the measured values. MI-Poser's IK and MIM incur little latency and the pipeline runs efficiently at 60 Hz.

We first evaluated MI-Poser's feasibility in body pose tracking by comparing its output with Microsoft Kinect [36] tracking. The study involved ten participants performing various upper-body movements, including out-of-sight hand motions. Results indicated that MI-Poser tracks upper-body pose with a mean joint position error of 6.6 cm. Additionally, we found that MI-Poser outperforms prior IMU-based work with the same sensor placement, thanks to the precise 6-DoF information EMF tracking provides.

Next, we assessed our approach to addressing the inherent metal interference issue by collecting a dataset of synchronized EMF sensor data, IMU sensor data, and ground truth pose data. Data were collected under three conditions with varying metal profiles: open-space, standard, and extreme cases. We first quantified errors between EMF and ground truth pose in each condition, finding significant tracking errors exist only briefly when the EMF sensor passes by a metal object. Our EMF-IMU fusion approach successfully detects such short-period metal interference (0.62 MCC in standard and 0.56 in extreme cases). Moreover, our correction approach reduces tracking error under interference by 11.6° max rotation error per session and 3.3 cm max position error per session in standard cases.

In sum, this work offers the following contributions:

- (1) We devised a hands-free wearable upper-body pose tracking system with a natural form factor, using wrist-worn EMF sensors and AR glasses.
- (2) We proposed metal interference mitigation (MIM) to address the inherent issue in utilizing EMF tracking in dynamic user environments.
- (3) In User Study 1, we demonstrated our system's ability to reconstruct 3D body pose in various movements with a small error in real-time, even when hands are out of view from the AR glasses.
- (4) In User Study 2, we showed that MIM significantly reduces error in EMF tracking in environments with different metal profiles. We also qualitatively demonstrated its efficacy in improving body pose tracking with smoother and more accurate output.

Although prior work used EMF sensors for tracking VR controllers [62], MI-Poser is the first hands-free, occlusion-robust body tracking system employing sparse on-body EMF sensors with metal interference solutions. The results suggest MI-Poser provides a practical body pose tracking system for everyday body-centered AR applications, such as remote communication through avatars with upper-body expression while walking outdoors.

2 RELATED WORK

To situate our work, we first review existing wearable systems for body pose tracking and discuss the need to address the limitations of current vision- and IMU-based systems. Then, we examine existing research employing magnetic tracking to provide a background for our EMF-based pose tracking system.

2.1 Wearable Systems for Body Pose Tracking

Researchers have explored wearable body pose estimation systems as a means to achieve portable and flexible interactions. Prior work has involved attaching an array of sensors on the body [51, 60] or using exoskeletons [71]. Recent developments in machine learning approaches have enabled systems with lower costs using sparse sensors, reducing the burden of wearing numerous sensors on the body. Vision-based approaches are the most popular in this context [2–4, 22, 39, 50, 58, 63, 66]. For example, Ahuja *et al.* [4] attached additional cameras to the Meta Quest 2 VR controllers to create an inside-out body capture system. Magic Leap 2 [33] employs a similar technique. Some researchers opted for wrist-worn vision sensors instead of VR controllers, using a spherical camera [8] or an array of small cameras [30]. Furthermore, several IK models have been developed to achieve high-fidelity body reconstruction from sparse sensor inputs. These models use the poses of the head and two hands as inputs to estimate the full body [5, 15, 22]. For instance, Jiang *et al.* [22] proposed a full-body pose tracking system using HTC VIVE (note it necessitates an additional base station in the environment). To train and evaluate these IK models, researchers utilized the extensive human motion database AMASS [34], which comprises a collection of high-precision MoCap datasets. Although these vision-based systems minimize error within the dataset, they unavoidably face challenges such as heavy computation (e.g., model inference) and sensing costs¹, which can be critical in resource-constrained devices like AR glasses. Moreover, they depend on line of sight; therefore, trained IK models may not function correctly if a user's hand moves out of the camera's view, limiting the possible tracking range of human body movement. Unlike existing VR headsets (e.g., Meta Quest Pro [35]) that use cameras for tracking hand-held controllers in 3D space, AR glasses lack sufficient space to accommodate multiple cameras, resulting in a limited field of view.

As alternatives, researchers have investigated IMU-based pose tracking [21, 37, 49, 56, 61, 64, 68–70]. Shen *et al.* [49] proposed a method for reconstructing an arm movement from a single smartwatch by using its embedded IMU sensors. Similarly, Tautges *et al.* [56] proposed an approach to reconstructing full-body animation from four acceleration sensors attached to wrists and ankles. In the context of IMU-based pose tracking, recent works

¹For reference, the HTC Vive Lighthouse consumes approximately 5W for the source [62].

leverage more sophisticated machine/deep learning models using the AMASS dataset [34]. For example, Sparse Inertial Poser [61] allows 3D human pose estimation using six IMU sensors attached to wrists, lower legs, back, and head. Deep Inertial Poser [21] improves the approach by incorporating temporal pose priors through deep learning. TransPose [70], LoBSTR [68], and Physical Inertial Poser [69] build upon those works and further advance the performance of IMU-based body pose tracking. However, these approaches require not a small number of sensor-instrumented joints (typically six) and exhibit a certain degree of errors.

Considering previous research, there is a demand for hands-free, occlusion-robust pose tracking systems with minimal errors, particularly in AR contexts. As a result, we developed a body pose tracking system utilizing a practically sparse sensor input (AR glasses and two wrist-worn sensors) based on a different sensing modality: EMF sensing.

2.2 Magnetic Tracking for Interaction

Magnetic field sensing-based tracking has a long history [46, 47] with a comprehensive review available in [44]. Several HCI applications have been proposed [9, 11, 18, 43, 52], leveraging occlusion-free tracking. For example, Abracadabra [18] tracks finger radial position relative to a watch using an attached magnet, and Nenya [9] measures magnetic field changes with a magnetometer-equipped smartwatch and a ring with two permanent magnets. Chen *et al.* [11] proposed uTrack, which tracks finger movements using a pair of magnetometers on the back of the fingers and a permanent magnet to the back of the thumb. While precise in short-range tracking like centimeters, these approaches can not be extended to long-range (e.g., body-scale) since the Earth's geomagnetic field easily influences the tracking. Razer Hydra [52] extends the sensing range of the controller by using a base station that generates a weak magnetic field.

Electromagnetic field (EMF) tracking involves oscillating magnetic fields [25, 29] and has gained attention for its precision in medium-range 6-DoF tracking. For a detailed review, see [17]. Several HCI applications have been proposed using EMF tracking, such as Finexus [12], which advances uTrack [11] by tracking multiple fingertips in real time. AuraRing [42] offers 5-DoF finger tracking for VR/AR applications with low power consumption (*i.e.*, around 2.3mW for a sensor in a ring and 73.3mW for a transmitter in a wristband). Whitmire *et al.* [62] extended the tracking range to body scale, enabling VR controller tracking by embedding three coils in HMD and a set of orthogonal receivers in hand-held devices with reasonable power consumption (*i.e.*, around 45mW for a sensor in a controller and 224mW for a transmitter in HMD without wireless communication module). Similar work to ours is EM-Pose [23], a wearable EMF-based body tracking system that uses 6 or 12 on-body sensors, which requires users to wear an additional EMF source on their back.

A known drawback of EMF tracking is its susceptibility to magnetic field distortion by environmental metals [42, 62], particularly in dynamic environments like VR/AR. Interference becomes more significant as the tracking range increases, such as body-scale [62]. Some work has attempted offline calibration to account for magnetic interference [26, 27], but online calibration is desirable. Although previous work [23, 62] recognized the issue, no work has been proposed to date to quantify the interference effect in different user environments and to devise solutions to it. Therefore, we propose approaches to addressing the metal interference issue in our EMF-based body pose tracking system.

3 PROPOSED METHOD

We design a hands-free, wearable upper-body pose tracking system with an extensive tracking range, MI-Poser. The system pipeline is illustrated in Figure 2. Unlike existing VR tracking systems that employ cameras and handheld controllers, MI-Poser utilizes wrist-worn EMF sensors for tracking. As depicted in Figure 1, users wear an EMF receiver on each wrist while an EMF source is mounted to the AR glasses. The EMF tracking captures the wrist poses relative to the EMF source. Simultaneously, the AR glasses track the user's head position in the



Fig. 2. The overview of MI-Poser’s pipeline.

world coordinate using Visual-Inertial Odometry [16]. These sparse measurements are input into our IK model to reconstruct a high-fidelity upper-body pose. We prioritized upper-body pose tracking due to its applicability to various AR applications. While our sensor setup could estimate full-body pose in a data-driven manner for specific motions through hallucination, like walking as shown in [37], we focus on the fidelity of the reconstructed pose in general conditions.

Metal interference is inevitable when using EMF sensors for body tracking in dynamic user environments. Therefore, we propose Metal Interference Mitigation (MIM) methods that operate online with minimal latency. While previous work [23, 62] acknowledged the issue, no concrete solutions have been proposed, as discussed in Section 2.2. Our solution incorporates an IMU sensor embedded with the EMF receiver. The fusion of EMF and IMU sensors has been studied to enhance EMF tracking performance in ideal, metal-free environments using a static filter like Kalman Filter [47]. However, in practical situations, the tracking algorithm should actively detect interference presence in real time and dynamically correct the trajectory.

3.1 MIM Overivew

Initially, we examined the behavior of our EMF tracking under metal interference. Within a 1.5 m range from the EMF source, metal effects are seldom present in open-space environments (e.g., outside), leading to accurate tracking performance. However, in typical spaces with some metal objects (e.g., a desk with a laptop), interference emerges when the EMF receiver approaches a metal object, confirming previous observations [62]. Even over a short period, interference significantly impacts the position and rotation of the EMF sensor, potentially degrading IK performance. We identified two types of metal interference based on how the sensor moves around metal objects. On one hand, when the sensor passes by a static metal object, a spike-like short-period error occurs in the tracking. On the other hand, when a metal object and the EMF receiver move together, an error persists as long as they remain in close range. Both cases can occur in end-user scenarios, such as users moving their arm near metal objects or holding a metal can while interacting with AR content.

To address this, we divided the problem into two parts: *interference detection* and *interference correction*. The overview of our MIM approach is presented in Figure 3. The detection part aims to identify moments when tracking errors arise due to metal interference, while the correction part seeks to mitigate errors by adjusting the interfered EMF sensor values. In this paper, we focus on correcting the first type of interference, the short-period error that occurs when metal objects are placed statically in an environment and the sensor encounters them occasionally (e.g., users swinging their arms). This is due to the challenge of tracking pose over a long period (more than a few seconds) using EMF and IMU sensors under metal interference. However, our detection method can also address the second type of interference, informing users of the reasons for degraded body tracking performance and improving user experience [6]. For instance, if a user holds a smartphone that causes interference to the tracking of the corresponding wrist, MI-Poser can notify the user via AR glasses that the tracking performance is low because the metal object is close to the hand.

3.2 Interference Detection

The first part is interference detection. There are two values regarding the rotation of the sensor based on different principles: angular momentum from the gyro sensor and orientation from the EMF sensor. Let’s assume a rotation in axis-angle representation $\Phi_s^{EMF}(t) \in \mathbb{R}^{1 \times 3}$ given time t when there is no metal interference. We use $I(t)$ as

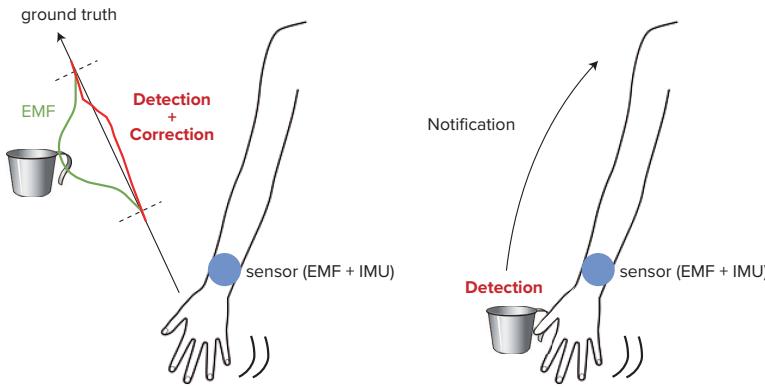


Fig. 3. Overview of the proposed metal interference mitigation (MIM) according to the two types of errors. When a user encounters a metal object for a short period, it is detected and the trajectory is corrected (left). When the interference is longer, e.g., when the user holds a metal object, it is detected for notifying the user (right).

a binary index to represent the presence or absence of interference; $I(t) = 0$ in this case. Simultaneously, an angular momentum $\Delta\Phi_s^{IMU}(t)$ is measured in the same coordinate as $\Phi_s^{EMF}(t)$. At time $t + \Delta t$, there is rotation information from the EMF sensor as $\Phi_s^{EMF}(t + \Delta t)$. If no interference occurs at $t + \Delta t$, an approximation holds: $\Phi_s^{EMF}(t + \Delta t) \sim \Phi_s^{EMF}(t) + \Delta\Phi_s^{IMU}(t) \times \Delta t$.

We can then introduce an error threshold $e_{th}^{\Delta\Phi}$ to estimate the interference state $I(t + \Delta t)$ by comparing $e_{th}^{\Delta\Phi}$ with the distance between $\Phi_s^{EMF}(t + \Delta t)$ and $\Phi_s^{EMF}(t) + \Delta\Phi_s^{IMU}(t) \times \Delta t$. The distance is the intrinsic geodesic distance between two angles. If the distance is larger than the threshold, the system predicts $\hat{I}(t + \Delta t) = 1$. Otherwise, it predicts $\hat{I}(t + \Delta t) = 0$.

3.3 Interference Correction

The second part is interference correction, which dynamically adjusts the measured value from the EMF sensor based on the detection result. It is essential to perform this correction online, as MI-Poser aims to be a real-time body pose tracking system. This means that if $\hat{I}(t) = 1$, we need to correct the current position $P_s^{EMF}(t)$ and rotation $\Phi_s^{EMF}(t)$ using past tracking and sensor data up to time t . If interference persists in $\hat{I}(t + \Delta t) = 1$, we must correct them using the past data up to $t + \Delta t$. As noted, we apply this correction as long as the detected interference is of short duration to avoid drift error.

For the rotation correction, we use $\Phi_s^{EMF}(t) + \Delta\Phi_s^{IMU}(t) \times \Delta t$ instead of $\Phi_s^{EMF}(t + \Delta t)$. Given that IMU-based rotation tracking is fairly feasible, we expected this simple solution to work well. Meanwhile, IMU-based position tracking is known to be a challenging problem, and we have prepared three methods.

3.3.1 IMU Odometry Model. This physics-based method uses initial velocity and a time series of acceleration from the IMU sensor to calculate position through dual integration.

$$x(t_0 + \Delta t) = x(t_0) + v(t_0) \times \Delta t + \int_{t_0}^{t_0 + \Delta t} \int_{t_0}^{\tau} a(\tau') d\tau' dt \quad (1)$$

, where $x(t)$, $v(t)$, and $a(t)$ represent position, velocity, and acceleration at time t , and t_0 represents the initial reference time. While straightforward, this method performs poorly when there is noise in $a(t)$ and $v(t_0)$ [57]. Since we use this approach around the interference moments, the EMF position tracking can contain noises

during the moments, resulting in noisy $v(t_0)$. Another limitation is that successful short-time arm tracking based on an IMU sensor [32] requires a high sampling frequency like 2000 Hz, indicating that lower frequency leads to larger errors.

3.3.2 Trajectory Forecasting Model. The IMU odometry method may not account for trajectory trends and seasonality, which are often used in practical time-series forecasting methods [14]. Human body movements, particularly arm movements, include short-duration trends and can be forecasted based on past trajectories [19, 55, 65]. We anticipated that a short-period future trajectory could be forecasted using previous tracking history, which could then be used to correct EMF position data under metal interference. We adopted the N-BEATS method [41], a state-of-the-art deep learning approach using backward and forward residual links and a deep stack of fully-connected layers. The model's input and output can be written as:

$$\{x(t_0), \dots, x(t_0 + \Delta t_{output})\} = \text{N-BEATS}(\{x(t_0 - \Delta t_{input}), \dots, x(t_0)\}) \quad (2)$$

, where Δt_{output} and Δt_{input} correspond to the amount of future data the model outputs and the amount of previous data the model takes as inputs, respectively.

In testing this model, we found that significant prediction errors tended to occur when there were large position changes right after the moment the model predicted. This can be understood as the time-series forecasting model estimating the trajectory based on past data but not reflecting future acceleration information. This observation led us to introduce the following model.

3.3.3 Fusion Model. To consider future acceleration while avoiding error due to noisy $v(t_0)$, we approximate the trajectory as follows:

$$x(t_0 + \Delta t) = \text{N-BEATS}(\{x(t_0 - \Delta t_{input}), \dots, x(t_0)\})|_{t=t_0+\Delta t} + \int_{t_0}^{t_0+\Delta t} \int_{t_0}^{\tau} a(\tau') d\tau' d\tau \quad (3)$$

, where Δt is small enough (at least, $\Delta t < \Delta t_{output}$). We iteratively use the same N-BEATS model. In detail, while the N-BEATS model outputs estimation for Δt_{output} seconds, we use the single prediction value corresponding to time $t_0 + \Delta t$. After adding the acceleration component through integration, we use this value as the input for the next N-BEATS inference for the next frame ($t_0 + 2\Delta t$) if metal interference still exists (i.e., $\hat{I}(t_0 + 2\Delta t) = 1$). In this way, we can adjust the N-BEATS prediction by adding the acceleration component, which further influences the subsequent trajectory forecasting.

4 IMPLEMENTATION AND SYSTEM PERFORMANCE

4.1 Hardware

Our EMF tracking system has an EMF transmitter (source) and two EMF receivers (sensors) using off-the-shelf 3D coils (See Figure 4). These components communicate with AR glasses through Bluetooth Low Energy (BLE) using the Enhanced ShockBurst protocol [40] to minimize latency. Designing an EMF tracking system for human body tracking involves multiple trade-offs, such as source coil size, electric current, sensor coil size, and tracking range/error requirements. Our final sensor configuration meets our 1.5 meter range requirement (typical arm reach) with a position RMS error of 0.9 mm and angle RMS error of 0.5° within a 1 meter range in an *ideal metal interference-free* lab environment. Additionally, an IMU sensor is integrated into the sensor, which we leverage for a sensor-fusion approach to MIM. The EMF and IMU data streams are synchronized and accessible via a BLE connection. We run the tracking at 120 fps with a latency of around 15 ms, and the tracking algorithm runs locally on the sensor. The algorithm incorporates Kalman Filter to stabilize long-range tracking and reduce jitter. We use Spectacles [53], commercial AR glasses, featuring a precise VIO algorithm running at 60 Hz, which is accessible via a BLE connection.

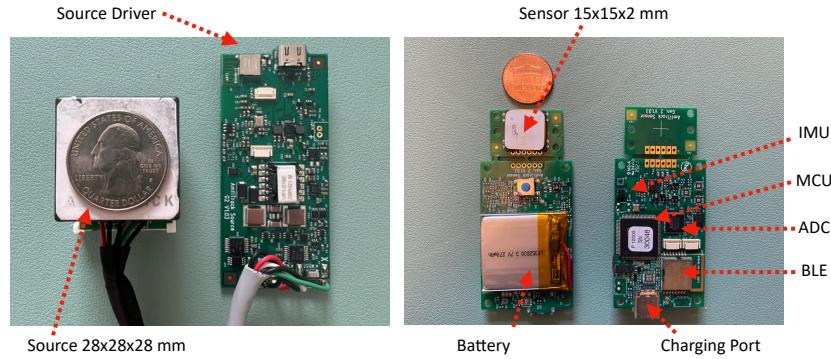


Fig. 4. EMF tracking hardware. The source (left) is integrated into AR glasses and the sensors (right) are attached to the user’s wrists.

4.2 IK Model for Pose Estimation from Sparse Sensor Inputs

We used the SMPL model [31] to represent and animate the human body pose. We trained our IK model to reconstruct the upper body from a sparse sensor set using the AMASS dataset [34], similarly to prior work discussed in Section 2.1. For the IK model, we adopted the state-of-the-art model from AvatarPoser [22]. The key difference is that our proposed system has EMF sensors on the wrists, while AvatarPoser assumes hand-held controllers. We expected that this different sensor placement would improve upper-body tracking by avoiding rotation noises from hand movements and helps the model infer more plausible arm poses.

In AMASS, we used a subset combination of the CMU, Eyes_Japan, KIT, MPI_HDM05, and TotalCapture datasets as the training set, and MPI_Limits as the validation set. We down-sampled the MoCap dataset from 120 Hz to 60 Hz and generated windowed segments of 40 frames (*i.e.*, 2/3 second window) with a stride length of 0.1 seconds to match with the original work [22]. We used the Adam optimizer [28] with a batch size of 32, and a starting learning rate was 0.001, which decays by a factor of 0.8 every 20 epochs. We performed the training with PyTorch on Google Cloud Platforms with NVIDIA Tesla V100 GPU.

To account for variations in body size and sensor-wearing positions, we calibrated the sensor outputs before inputting them into the IK model. For sensor-position calibration, the user simply maintained the default T-pose (Figure 13 in Appendix) for a few seconds. We used sensor measurements taken during this period for calibration, similar to prior work [21, 70]. We first estimated a scaling factor by comparing the arm span between actual sensor measurements and the SMPL model definition. To compensate for minor sensor offsets, we applied spatial transformations to the sensor output, ensuring alignment with the SMPL model definitions in the default pose. We applied the scaling factor and transformations to each frame throughout the entire body pose tracking session.

4.3 Trajectory Forecasting Model for MIM

To train the N-BEATS model for MIM (interference correction), we first collected data by moving the EMF sensor freely in open space without metal objects. Approximately 40 minutes of data were used for training the N-BEATS model with two stacks where there are two blocks per stack with 512 hidden layer units. The model takes 120 samples (corresponding to 1 second of Δt_{input}) as input and outputs 60 samples (corresponding to 0.5 seconds of Δt_{output}) position data (Recall the EMF tracking runs at 120 Hz). We used the Adam optimizer [28] with batch size 16 up to 20 epochs. Within the dataset, the best model performed a 1.26 cm mean absolute position error in the validation dataset. This means that the model can forecast a position of 0.5 seconds ahead with small errors.

4.4 Real-Time System Performance

Currently, the IK model and MIM process run on a laptop (MacBook Pro with a 2.6 GHz 6-Core CPU and 16 GB memory) written in Python while streaming data in real time. Future work will involve transferring the process to AR glasses using JavaScript. Spectacles are equipped with an Octa-core CPU (2×2.52 GHz + 6×1.7 GHz). Given the limited computational resources, it is crucial to consider power consumption and inference speed, and we report them in our current prototype in this section.

4.4.1 Power Consumption. The power consumption at the source and the sensor is 1.4W and 0.68W, involving the communication modules, respectively. Further optimization, as demonstrated in [62], is advisable. For instance, replacing the currently used microcontroller (F7), which has more capabilities than necessary, with a lower power consumption alternative (H7) could reduce power consumption. Nonetheless, the sensor-level power consumption is significantly lower than existing camera-based research work. For example, ControllerPose [4] attaches a camera to each controller to capture upper-body movements, and a single camera's power consumption is approximately 3.3W, which does not involve the hand tracking algorithm. While we must consider the power required to run the model on the device, our prototype is suggested to operate with reasonable power consumption.

4.4.2 Inference Speed. The IK model's current latency is 4.2 ms on the laptop, from captured EMF values to the output. Likewise, MIM's detection and correction models incur average latencies of 0.09 ms and 0.50 ms, respectively. These latencies do not significantly impact the body pose tracking pipeline, making MIM a suitable complement to the EMF-based upper-body pose tracking system. Together with the MI-Poser's IK model for reconstructing body pose, our pipeline takes approximately 5 ms to process one frame on a laptop. The current MI-Poser pipeline operates efficiently at 60 Hz (recall the EMF sensor runs at 120 Hz and the VIO tracking runs at 60 Hz). Notably, commercial on-device tracking speeds, such as Meta Quest 2 [35], are also 60 Hz. For further comparison, we ran ControllerPose [4] and IMUPoset [37] systems using the same laptop, with pipeline speeds of approximately 4 Hz and 48 Hz, respectively. Please refer to the Video Figure for a real-time demonstration.

5 USER STUDY 1: UPPER-BODY POSE TRACKING IN THE ABSENCE OF METAL INTERFERENCE

Since MI-Poser is the first setup for an upper-body tracking system with two wrist-worn EMF sensors and AR glasses, we first examined its tracking performance in an open space (without visible metal objects) and compared it with similar setups using IMU sensors. We trained an IK model using the existing AMASS dataset [34] and tested it with real sensor data.

5.1 Data Collection

We collected sensor data from AR glasses and EMF sensors to demonstrate system performance using the setup shown in Figure 1. Ground truth data were obtained using Microsoft Kinect [36], following prior work [4]. To ensure the reliability of the ground truth data from Kinect, we filtered out unreliable inferred tracking frames in post-processing, based on the tracking state Kinect logged, which constituted approximately 10% of the data.

To inspect the fine-grained performance of MI-Poser across various upper-body movements, we designed an obstacle course-style setup inspired by previous works [3, 21]. We selected motions that encompassed a diverse range of upper-body movements:

- *Punch*: the participant alternately punches with both arms in front of their body.
- *Wave*: the participant randomly raises their arms and waves in the air.
- *Swing*: the participant alternately swings their arms from side to side.
- *Rotate*: the participant rotates both arms against each other in front of their chest.
- *Walk*: the participant walks randomly with their arms swinging naturally.
- *Basketball*: the participant jumps and performs basketball shooting gestures over their head.

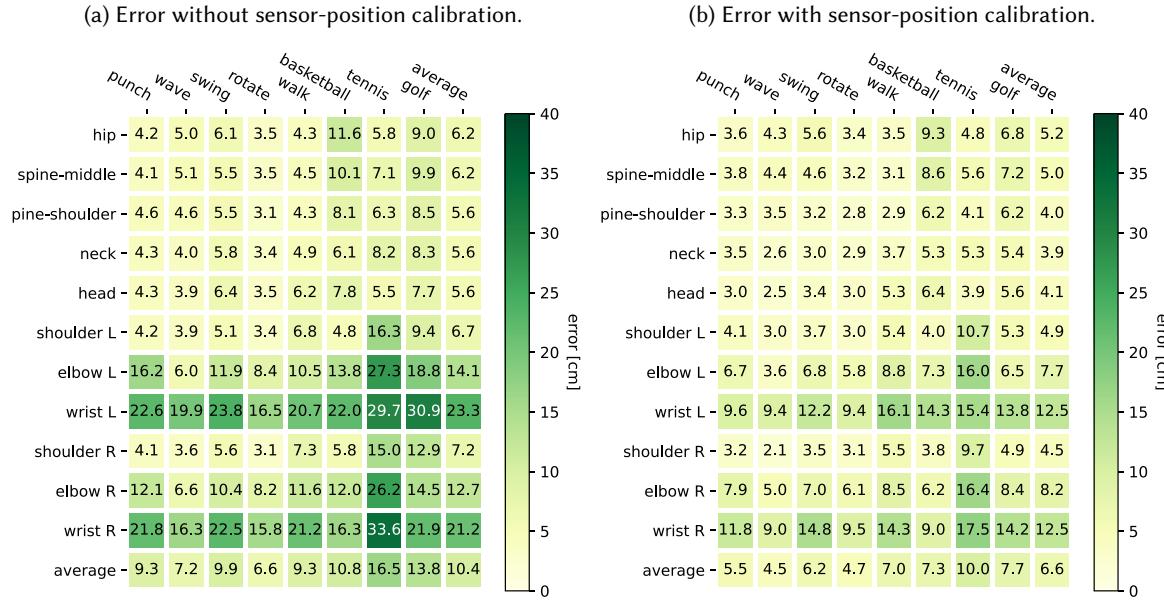


Fig. 5. MI-Poser’s performance (joint position error) on real sensor data across body region and motion. The error is overall small thanks to the precise EMF tracking. The sensor-position calibration significantly reduces the error. Error accumulates from the hip (alignment root) to the end-effectors. The largest error is observed during the tennis motion, which involves rapid and extensive torso movements.

- *Tennis*: the participant swings their arms from behind their body to the front, with torso rotation.
- *Golf*: the participant swings both arms together, with torso rotation.

Several of these motions, such as swinging, walking, basketball, tennis, and golf, include moments when hands move outside the field of view of the cameras on AR glasses. These moments are often challenging to track in conventional camera-based AR systems.

We recruited 10 participants from our institution with diverse genders, ages, weights, and body shapes for data collection. Participants performed each motion for 50 seconds, with 10-second rest periods in between. The entire data collection process, including the initial calibration, took approximately 10 minutes per participant. We obtained approval from our institution to conduct the study.

5.2 Results

5.2.1 Fine-Grained Error Metric. The results across different joints and motions are presented in Figure 5. The overall error (a) without and (b) with sensor-position calibration is 10.4 cm and 6.6 cm, respectively, demonstrating a significant improvement due to the sensor-position calibration. As we align the root (hip) in calculating the error metric following prior work such as [21, 37], the hip generally has the smallest error, and the error propagates to the end-effectors like wrists, leading to the largest error. Still, the overall error is reasonably small after the sensor-position calibration, including when hands are out of view from the AR glasses.

However, there are several performance limitations to consider. First, by examining the error by motions in Figure 5, larger errors arise from those involving fast and extensive torso movements, such as tennis. Next, the EMF tracking method has a minimum working distance of about 10 cm to prevent signal saturation. As a result,

gestures involving close proximity to the source, like both hands being close in golf, can cause significant signal jitters. Lastly, while our sensor-position calibration accounts for body skeleton scale and minor sensor position shifts, it does not adjust for variations in body shape. This can lead to inaccurate hand-body contact, such as when the hand penetrates or hovers over the body mesh despite physical contact.

Additionally, we tested a scenario where a user wears a single EMF receiver on their wrist, assuming the sensor is embedded in a smartwatch. We trained a different model with this configuration using the AMASS dataset and evaluated its performance with our dataset, using data corresponding to one EMF receiver and AR glasses. The overall error without and with sensor-position calibration is 22.0 cm and 15.9 cm. The largest error comes from the hand without the EMF receiver, while similar performance is maintained for the hand with the receiver. Anecdotally, the error decreases in some movements (*e.g.*, rotating, walking), which can be attributed to hallucination from the training dataset, as indicated by [37]. Thus, while depending on applications, utilizing a pervasive configuration of on-body sensors could be valuable when users employ tracking in their everyday lives, such as in AR applications while walking.

5.2.2 Comparison to Prior Work. Our system is the first of its kind to use two wrist-worn EMF sensors for upper-body pose tracking. To the best of our knowledge, no prior research has reported the performance of upper-body pose tracking using such a sparse set of real sensors (head and two wrists). Although most IMU-based approaches were tested using the AMASS dataset, several works also reported performance on real sensor data. However, these studies used different datasets and involved more sensors (including leg-worn IMUs) in reporting full-body performance, making direct comparisons difficult. Moreover, the current Spectacles’ SDK² does not provide raw IMU signals, preventing us from obtaining the IMU signal for the head in our setup, which makes a fair comparison with IMU-based solutions using our dataset impossible. This remains a limitation of the current study.

As a remedy, we adopted IMUPoser’s model [37] (two-layer bi-directional LSTM) as a baseline for the IMU-based upper-body pose tracking and tested its error on the DIP-IMU dataset and IMUPoser dataset. Both datasets contain human body poses across different motions similar to ours, such as arm raise, arm swing, and walking, involving 10 participants. The DIP-IMU dataset [21] includes 17 IMUs (X-Sense sensors), while IMUPoser uses common wearables as sensors, such as smartwatches and earbuds. We trained their model with the AMASS dataset (the same subsets we used for MI-Poser) using the configuration of three IMU sensors corresponding to our setup (*i.e.*, the head and two wrists) and evaluated its performance on the datasets. As a result, the joint position error for the upper body is 8.3 cm and 10.4 cm with sensor-position calibration for the DIP-IMU and IMUPoser datasets, respectively. Although not a direct comparison in terms of evaluating with different datasets, it is suggested that MI-Poser has a better performance compared to IMU-based systems, thanks to the high precision EMF/VIO tracking.

6 USER STUDY 2: MEASURING AND MITIGATING METAL INTERFERENCE

User Study 1 demonstrated that our EMF-based tracking setup achieves accurate upper-body pose tracking. In this section, we quantify the metal interference on EMF tracking in various environments and examine the effectiveness of MIM in enhancing the input data for our pose-tracking pipeline.

6.1 Data Collection

6.1.1 Apparatus. We used the same sensor described in Section 4.1, which streams EMF tracking and IMU tracking data synchronously. To obtain ground truth tracking data, we used Apple ARKit 4 [7], which enables a state-of-the-art self-localization in world coordinates based on VIO tracking [24]. We developed a custom

²<https://docs.spectacles.dev/app/reference/api>

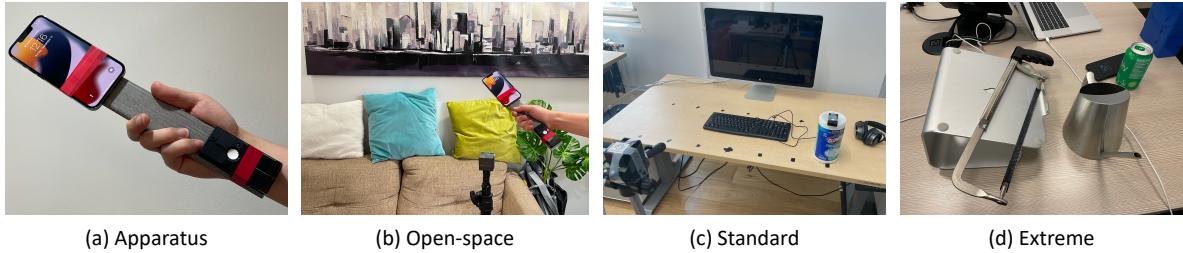


Fig. 6. Setup for data collection in User Study 2. (a) An EMF sensor and iPhone 13 Pro are tightly attached to a non-metal object. The iPhone tracks its pose based on ARKit's world tracking function. (b) Open-space with few metal objects. We intentionally chose locations with visually striking objects to aid ARKit's VIO tracking. (c) Standard condition with common metal objects such as a monitor. (d) Extreme condition where we intentionally move the sensor close to common metal objects such as a can. Note that the apparatus is not visible in (c) and (d).

iOS application using Xcode 13.3 that records the position, orientation, and UNIX timestamp. We attach the EMF sensor and iPhone 13 Pro, which can utilize built-in LiDAR for enhanced ARKit tracking performance, to a non-metal rigid body with a space of 20 cm between them. This decision was based on the observation by Whitmire *et al.* [62] that significant interference occurs if an EMF receiver and a smartphone are closer than approximately 5 cm. The apparatus used for data collection is shown in Figure 6 (a). When the apparatus is moved, the sensor streams its EMF pose tracking data along with the IMU data at 120 Hz relative to the EMF source, while the iPhone captures its pose at 60 Hz relative to the reference world coordinate. Since these two sensors are attached to a rigid body and their transformation is constant, we can align their coordinates using a calibration process, which we elaborate on in Appendix B.

We did not use high-end MoCap systems such as OptiTrack [38] because we wanted to conduct data collection in multiple and actual environments. Our apparatus offers a convenient method for data collection. We tested whether ARKit could work properly when we moved the apparatus naturally to represent arm movements in VR/AR scenarios. ARKit reports its tracking state, and based on that, we observed that the tracking performance worsens if we rotate the apparatus too quickly. Consequently, we could not include such movements in the data collection.

6.1.2 Condition. We selected three representative cases for data collection concerning the level of metal interference: open-space, standard, and extreme conditions. The open-space condition represents locations with minimal metal interference (See Figure 6 (b)). In contrast, the standard condition includes typical places like desks and rooms with a few metal objects (*e.g.*, a desk with metal support, Figure 6 (c)). For the open-space condition, we chose two locations in a building where no visible metal objects were present within a 2-meter range, except for the floor³. For the standard condition, we selected three locations in the same building: a desk with a few everyday metal objects, such as a laptop and monitor, a meeting space with a metal door, and a crafting room with some metal objects like a hammer. In addition to these two conditions, we added an extreme condition, where we intentionally moved the apparatus closer to metal objects for an extended period, such as touching a laptop or holding a metal can (Figure 6 (d)). This condition was introduced to examine the extent of errors that could occur in potential end-user environments.

6.1.3 Procedure. We collected four session data for each environment, resulting in 8 sessions for the open-space condition, 12 sessions for the standard condition, and 4 sessions for the extreme condition. Five people from our

³As a typical building, the floor has steel support.

Table 1. Summary of the dataset we collected in different conditions.

	Open-space	Standard	Extreme
Number of Sessions	8	16	4
Mean Position Error Per Session (cm)	0.72 ± 0.18	2.09 ± 1.10	9.24 ± 3.66
Max Position Error Per Session (cm)	6.64 ± 4.53	25.5 ± 18.2	99.3 ± 33.3
Mean Rotation Error Per Session (°)	1.27 ± 0.09	2.08 ± 0.68	4.73 ± 1.35
Max Rotation Error Per Session (°)	4.91 ± 2.32	18.2 ± 22.7	93.6 ± 12.1

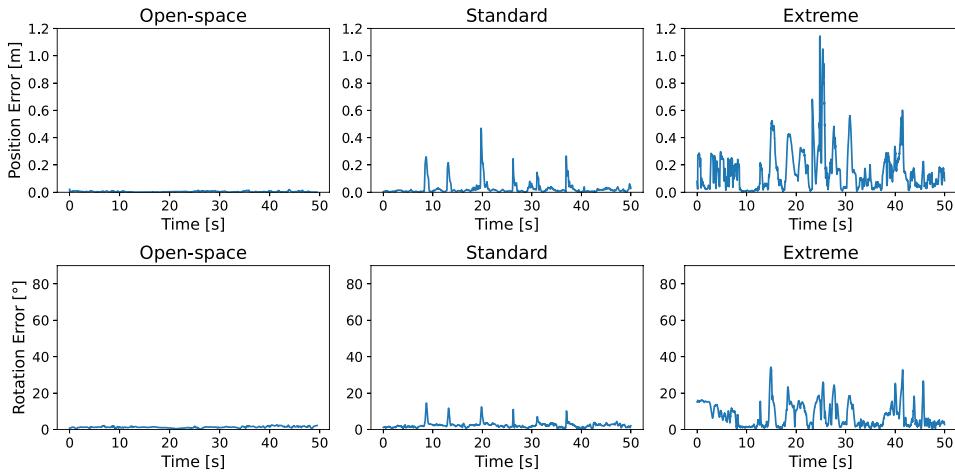


Fig. 7. Example plots of the position and rotation error between the ground truth and EMF sensor values in the three different conditions.

institution were asked to hold the apparatus and move around freely within a roughly 1.5-meter range from the EMF source for approximately one minute, emulating natural body movements. We informed them that this data would be used for tracking body pose in VR/AR applications and asked them to move freely as if they were playing with VR/AR content, such as exercising or manipulating imaginary objects. Note that the EMF source is fixed to the environment; as detailed in Appendix A, we can compute the sensor values with respect to the source using the HMD’s VIO tracking, so we considered the metal interference problem in the EMF coordinate.

Before each session, we initialized the ARKit world tracking coordinate and scanned the environment, preventing it from entering the “extending map” status during the recording to avoid degraded tracking performance. In this step, we also placed some non-metal objects in the environment if there were not many visually striking objects (e.g., white wall) to aid ARKit’s VIO tracking. In every condition, we used the first 20 seconds for calibration; during this time, we were careful not to move the apparatus close to metal objects in the environment, regardless of the condition. Furthermore, the apparatus was placed statically within a range of 1 meter from the EMF transmitter for about 3 seconds within the same 20 seconds. We used the static part to calibrate the IMU sensor to obtain linear acceleration, while the entire 20-second data was used to calibrate the sensor and iPhone coordinates. The calibration was conducted as an offline process after collection. Each session lasted approximately 80 seconds, resulting in 1-minute data after excluding the calibration part.

Table 2. Correlation between position error between ground truth and EMF, rotation error between ground truth and EMF, and gyro error between EMF and IMU.

	Open-space	Standard	Extreme
Position error v.s. Rotation error	0.49	0.61	0.72
Gyro error v.s. Position error	0.65	0.68	0.58
Gyro error v.s. Rotation error	0.50	0.56	0.59

6.1.4 Dataset Overview. Table 1 summarizes the dataset we collected. The position and rotation error was calculated as the distance between the ground truth and the EMF position and rotation, respectively. This table shows that the open-space condition has the smallest error, while the extreme condition has the largest. Figure 7 shows sample error plots within a session in the three different conditions. From this plot, we can observe a trend corresponding to the two types of errors we discussed in Section 3.1. There is almost no error in the open-space condition, supporting the use of MI-Poser in such environments. However, there are spike-like significant errors in the standard condition, indicating that MI-Poser’s pose tracking can deteriorate for those short periods. On the other hand, significant errors persist for several seconds in the extreme condition corresponding to when users intentionally hold or get close to metal objects. This result confirmed the two types of interference we discussed in Section 3.1.

6.1.5 Preliminary Analysis. Figure 7 shows a high correlation between the rotation and position error. As discussed in Section 3.2, we rely on the gyro error between the IMU and EMF sensor to detect metal interference. To verify the validity of the approach, we analyzed the correlation between these values. Table 2 summarizes the frame-by-frame correlation across the three conditions. The position and rotation errors between ground truth and EMF are highly correlated, indicating that the interference causes both errors. This trend is more obvious when there is more metal interference, *i.e.*, in the extreme condition. Moreover, the gyro error correlates with these two errors with high coefficients, suggesting that the gyro error can serve as an indicator to identify moments with interference, supporting our detection approach.

6.2 Detection Result

We first evaluated our method’s efficacy in detecting interference. Although interference is inherently continuous, simplifying it into a binary state streamlines user notification (*e.g.*, through an HMD) and subsequent trajectory correction, as detailed in Section 3.3. Consequently, we implemented a threshold-based approach to determine frame interference. As previously noted, position and rotation errors correlate and impact IK models. Thus, we used position error as the interference criterion in the following analyses, though similar results can be achieved with rotation error.

We first regarded a frame at time t where the error e_t^P between the ground truth and EMF position is larger than a given threshold e_{th}^P as a frame with interference. By varying e_{th}^P from 0 cm to 50 cm in 0.1 cm increments, we identified the optimal threshold $e_{th}^{\Delta\Phi}$ (refer to Section 3.2) using a naïve full search. Due to imbalanced label distribution (*i.e.*, few frames with interference), we employed Matthews correlation coefficient (MCC) [13] as the performance metric.

Figure 8 depicts the evaluation outcomes as e_{th}^P was altered. The blue line denotes the MCC value, while the orange line signifies the positive (interference) sample ratio. For instance, using $e_{th}^P = 0.1$ as an empirically determined threshold, positive frame ratios for open-space, standard, and extreme conditions are 0.13%, 3.14%, and 25.8%, respectively. Table 3 shows an in-depth sample result with the same threshold. The results indicate that our approach, comparing IMU and EMF gyro data, can effectively detect rare interference instances with

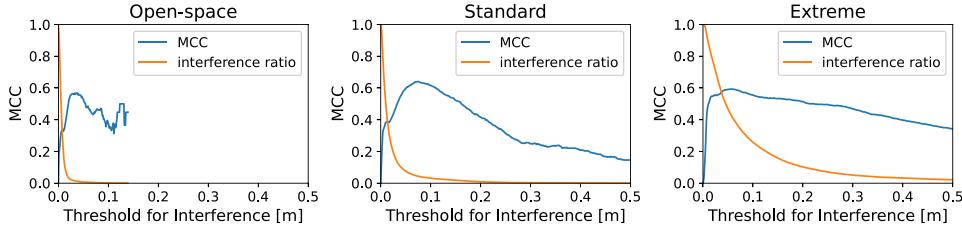


Fig. 8. Result of the interference detection by condition. Our approach reasonably detects the interference ($MCC \sim 0.6$).

Table 3. Example result of our interference detection when we regard $e_{th}^P = 0.1$ as a criteria for interference.

	Open-space	Standard	Extreme
Interference Ratio	0.13%	3.14%	25.8%
MCC	0.38	0.62	0.56
F1	0.33	0.63	0.68
Precision	0.22	0.56	0.59
Recall	0.63	0.71	0.80

reasonable accuracy. By adjusting $e_{th}^{\Delta\Phi}$, we can develop alternative models, such as high-precision models, to minimize false-positive user notifications. Our current online frame-by-frame method does not utilize future data; however, if developers permit slight delays, employing data from several frames ahead could reduce false positives and enhance detection accuracy.

6.3 Correction Result

Our correction approach's hyperparameter, $e_{th}^{\Delta\Phi}$, determines the frequency of system corrections. A smaller $e_{th}^{\Delta\Phi}$ leads to prolonged IMU-based trajectory corrections, which can cause drift problems and result in incorrect adjustments. Consequently, we varied $e_{th}^{\Delta\Phi}$ to optimize performance within this trade-off. In extreme conditions, trajectory correction proves challenging due to extended interference periods (See Figure 7). Therefore, we focused on enhancing tracking performance in open-space and standard conditions, where spike-like interference is prevalent. We recommend using the detection model to inform users of interference presence for the extreme condition, as outlined in Section 3.1.

6.3.1 Rotation Error. Figure 9 shows the best performance of our correction approach when we varied $e_{th}^{\Delta\Phi}$. The used $e_{th}^{\Delta\Phi}$ is fixed across sessions. In the open-space condition, the lowest mean and maximum rotation errors per session, $1.17^\circ \pm 0.13^\circ$ and $3.35^\circ \pm 0.70^\circ$, respectively, are reduced with the correction compared to the raw errors. Similarly, the standard condition exhibits lower mean and maximum rotation errors per session, $1.78^\circ \pm 0.52^\circ$ and $6.59^\circ \pm 4.21^\circ$, respectively. It is worth noting that the improvement in the mean error per session appears small since the correction only occurred during a short moment in each session. Here, in the same manner as the detection result, we used $e_{th}^P = 0.1$ and calculated the improvement within the interference period instead of per session. As a result, the improvements in the open-space and standard conditions are 71.0% ($6.92^\circ \rightarrow 2.01^\circ$) and 57.1% ($8.88^\circ \rightarrow 3.81^\circ$), respectively. This result confirms that the large rotation error is significantly reduced thanks to our EMF-IMU fusion approach.

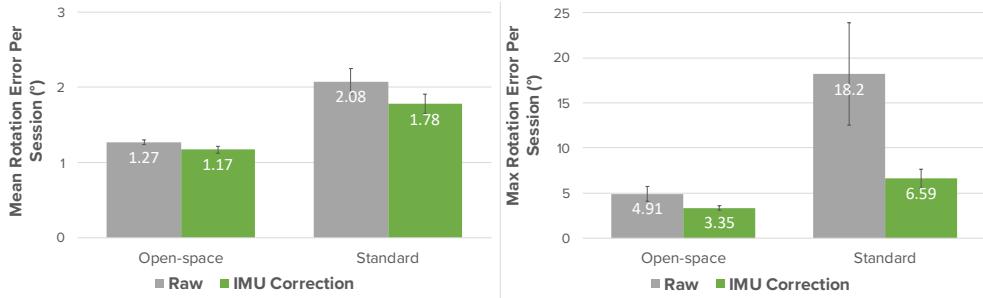


Fig. 9. Result of the interference correction in rotation error by condition. Mean rotation error per session (left). Max rotation error per session (right). Our rotation correction approach reduces the error, especially suppressing large errors. The error bars are standard error.

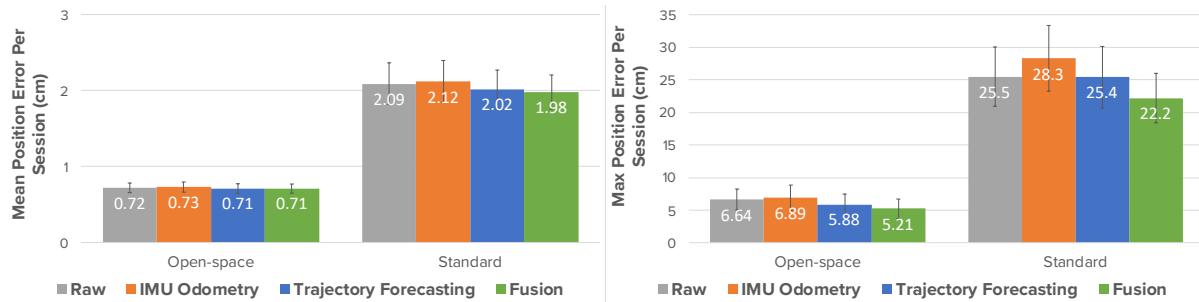


Fig. 10. Result of the interference correction in position error by condition and method. Mean position error per session (left). Max position error per session (right). Our position correction approach with the fusion model reduces the error but the improvement is relatively smaller than in the rotation error. The error bars are standard error.

6.3.2 Position Error. Similar to the rotation correction, the performance of our position correction depends on $e_{th}^{\Delta\Phi}$, so we varied it and computed the total improvement achieved by the model. Figure 10 compares the result across different conditions and methods. Without correction, the position error in the open-space condition was 0.72 cm ($SD=0.18$ cm). Notably, vision-based marker-less fingertip tracking in Meta Quest 2 achieved approximately 1.0 cm static positional error in a similar setting [1], indicating that EMF tracking in open-space conditions performs reasonably accurately.

Comparing different correction approaches, the IMU odometry model produced inferior results due to the drift error and sensor noise, as anticipated in Section 3.3. Conversely, our fusion model effectively leverages both IMU odometry and trajectory forecasting. Specifically, the lowest mean and maximum position errors per session in the open-space condition are 0.71 cm ($SD=0.18$ cm) and 5.21 cm ($SD=4.22$ cm), respectively. In the standard condition, the lowest mean and maximum position errors per session are 1.98 cm ($SD=0.90$ cm) and 22.2 cm ($SD=15.2$ cm). It is important to note that the mean position error improvement is not readily apparent since only a small portion of a session undergoes correction, while the maximum position error displays a significant reduction.

Despite this, the improvement (error reduction rate) is not as substantial as the rotation correction for several reasons. First, the noise profile in the measured acceleration and gyro values is different in the IMU sensor.



Fig. 11. Demonstration of the interference detection. The interference level is estimated on a frame-by-frame basis and displayed on the laptop (red line). When the sensor passes near the metal object (left figure), the interference level increases, while it remains unchanged in other situations (right figure).



Fig. 12. Demonstration of the interference correction. A user moves their hand horizontally. The figure is generated using 1/3 seconds of the rendered video. MIM contributes to improved output, resulting in a smoother and more accurate hand pose.

Secondly, while a double integral is applied to the measured acceleration for correcting position, only a single integral is applied to the measured gyro for correcting rotation, resulting in less drift error.

7 END-TO-END SYSTEM DEMONSTRATION

We now demonstrate how MIM works in the MI-Poser’s end-to-end pipeline. Figure 11 shows a user moving their hands along with the detected interference level. Only when the hand (*i.e.*, the EMF sensor) goes close to the metal door is there an increased value in the interference level, providing users with a way to identify potential tracking degradation due to metal objects in the environment. Moreover, Figure 12 presents a user moving their hand horizontally over a metal object and illustrates two outputs of the proposed method with and without MIM. Without MIM, the hand pose exhibits significant variance due to interference. However, with MIM, the hand pose becomes more stable, and the orientation appears more reasonable (*i.e.*, more horizontal). Please refer to the Video Figure for demonstrations.

8 DISCUSSION

We have demonstrated that MI-Poser can achieve real-time pose reconstruction with minimal errors while maintaining natural form factors for AR scenarios, such as hands-free operation. Furthermore, its robustness

against metal interference due to interference detection and correction makes it a promising solution for AR/VR applications. Several limitations and areas for future work remain, which we discuss in this section.

8.1 Limitations

First, the ground truth upper-body pose data in User Study 1 was obtained using Kinect rather than high-end MoCap systems with external cameras, such as OptiTrack [38] and Vicon [59]. Consequently, there is uncertainty in the ground truth data. For instance, Kinect is known to be slower than OptiTrack by 50 ms [10], which may have led to overestimated errors for fast motions (e.g., tennis) in our results. While our setup is sufficient to demonstrate the proof-of-concept (as done in [4]), further investigation with high-end MoCap systems is necessary to assess more precise performance.

Similarly, in User Study 2, we used ARKit on an iPhone with LiDAR to collect ground truth pose data, aiming for convenient data collection in multiple environments with different metal profiles. Although state-of-the-art in VIO, this approach is less accurate compared to MoCap systems. Nevertheless, the difference between EMF and ARKit tracking is small (0.72 cm) in the open-space condition (See Table 1). Considering the error of EMF tracking in an *ideal metal-free environment* is 0.9 mm, we conjecture that ARKit’s tracking performance has been reasonable, at least for analyzing significantly larger errors due to interference. We should note, as discussed in Section 6.1.1, our data collection setup could not include excessively fast arm movements, which necessitates an alternative setup.

In addition, our system relies on VIO tracking on AR glasses, causing it to suffer from degraded performance under low light conditions. Therefore, applications should allow users to disable global body tracking in such situations, as implemented in conventional VR headsets.

8.2 Future Work

8.2.1 Further Study about the Efficacy of MIM. While MIM effectively mitigated metal interference at the sensor tracking level, and we demonstrated its qualitative improvement in body pose tracking, its quantitative improvement remains to be examined. In the future, we plan to emulate a variety of user environments (e.g., a living room with some metal objects) and collect in-situ ground truth body pose tracking using a system such as OptiTrack, and evaluate MI-Poser. Such a study will reveal when metal interference occurs in actual scenes and how much MIM contributes to pose estimation and remains an important future work for this research. Furthermore, as outlined in Section 3.1, MI-Poser can notify users about potentially degraded tracking performance based on interference detection. Examining how users appreciate such feedback from the perspective of maintaining trust between users and systems [6] would be insightful.

8.2.2 Algorithm Refinement for MIM. The solutions we tested for MIM are simple yet effective. While simple solutions are preferred for their ease of deployment and maintenance, more sophisticated algorithms can be explored. For instance, we employed an algorithm to detect interference by analyzing the data of a single frame. However, other unsupervised anomaly detection methods [67] can leverage time-series trends to improve precision. Furthermore, collecting more data in various environments to increase samples with interference makes it possible to train models in a supervised manner for both detection and correction.

8.2.3 Evaluation of Full-Body Tracking Performance. As we mentioned in Section 3, prior work showed the possibility of full-body tracking with a similar setup as ours (i.e., sensors on the head and wrists) by training a model with a large dataset. The IK model we used [22] can also estimate the pose of joints from the lower-body. However, this is essentially a hallucination from the data and is known not to be robust against motions not in the training dataset [37]. Thus, we focused on upper-body tracking in our proof-of-concept. It is worth testing the current performance of full-body tracking with various motions.

8.2.4 Integration of the EMF Sensor into Smartwatch. Given the form factor of MI-Poser, integrating an EMF sensor into a smartwatch is a promising direction. In this way, we can use MI-Poser with an even sparser sensor setup, namely, AR glasses and a smartwatch, enabling more ubiquitous scenarios such as video calling with a rich upper-body expression of a 3D avatar from outside. Undoubtedly, this setup loses information about the wrist without the watch, leading to lower fidelity as demonstrated in Section 5.2; hence, the benefit of applications needs to be tested. Note that a shield must be added to the electronics inside the watch to prevent interference with the EMF sensor.

9 CONCLUSION

We proposed MI-Poser, a body pose tracking system using an EMF transmitter attached to AR glasses and two wrist-worn EMF sensors, which offers a hands-free, wide-range solution. In User Study 1, we demonstrated that MI-Poser can reconstruct upper-body pose with small errors across various movements, including cases where hands are out of sight from AR glasses, highlighting the benefit of combining EMF tracking and IK. As critical in utilizing body-scale EMF sensing with a sparse sensor setup, we also dealt with the metal interference issue to improve MI-Poser’s robustness in end-user environments, namely, metal interference mitigation (MIM). In User Study 2, using a newly collected dataset, we quantified errors due to interference in different metal conditions and proposed solutions based on the EMF-IMU fusion approach. The effectiveness of the proposed interference detection and correction was demonstrated, which is the first of its kind in using EMF tracking. While future work remains, the results suggest that MI-Poser offers developers a practical body pose tracking system, especially for enabling many interesting everyday AR applications.

REFERENCES

- [1] Diar Abdulkarim, Massimiliano Di Luca, Poppy Aves, Sang-Hoon Yeo, R Chris Miall, Peter Holland, and Joseph M Galea. 2022. A methodological framework to assess the accuracy of virtual reality hand-tracking systems: A case study with the oculus quest 2. *BioRxiv* (2022), 2022–02.
- [2] Karan Ahuja, Chris Harrison, Mayank Goel, and Robert Xiao. 2019. MeCap: Whole-Body Digitization for Low-Cost VR/AR Headsets. In *Proceedings of the 32nd Annual ACM Symposium on User Interface Software and Technology, UIST 2019, New Orleans, LA, USA, October 20-23, 2019*. ACM, New York, 453–462. <https://doi.org/10.1145/3332165.3347889>
- [3] Karan Ahuja, Sven Mayer, Mayank Goel, and Chris Harrison. 2021. Pose-on-the-Go: Approximating User Pose with Smartphone Sensor Fusion and Inverse Kinematics. In *CHI ’21: CHI Conference on Human Factors in Computing Systems, Virtual Event / Yokohama, Japan, May 8-13, 2021*. ACM, New York, 9:1–9:12. <https://doi.org/10.1145/3411764.3445582>
- [4] Karan Ahuja, Vivian Shen, Cathy Mengying Fang, Nathan Riopelle, Andy Kong, and Chris Harrison. 2022. ControllerPose: Inside-Out Body Capture with VR Controller Cameras. In *CHI ’22: CHI Conference on Human Factors in Computing Systems, New Orleans, LA, USA, 29 April 2022 - 5 May 2022*. ACM, New York, 108:1–108:13. <https://doi.org/10.1145/3491102.3502105>
- [5] Sadegh Aliakbarian, Pashmina Cameron, Federica Bogo, Andrew W. Fitzgibbon, and Thomas J. Cashman. 2022. FLAG: Flow-based 3D Avatar Generation from Sparse Observations. *CoRR* abs/2203.05789 (2022). <https://doi.org/10.48550/arXiv.2203.05789>
- [6] Saleema Amershi, Daniel S. Weld, Mihaela Vorvoreanu, Adam Fournier, Besmira Nushi, Penny Collisson, Jina Suh, Shamsi T. Iqbal, Paul N. Bennett, Kori Inkpen, Jaime Teevan, Ruth Kikin-Gil, and Eric Horvitz. 2019. Guidelines for Human-AI Interaction. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems, CHI 2019, Glasgow, Scotland, UK, May 04-09, 2019*. ACM, New York, NY, 3. <https://doi.org/10.1145/3290605.3300233>
- [7] Apple. 2020. ARKit Framework. <https://developer.apple.com/documentation/arkit>
- [8] Riku Arakawa, Azumi Maekawa, Zendai Kashino, and Masahiko Inami. 2020. Hand with Sensing Sphere: Body-Centered Spatial Interactions with a Hand-Worn Spherical Camera. In *SUI ’20: Symposium on Spatial User Interaction, Virtual Event, Canada, October 31 - November 1, 2020*. ACM, New York, 1:1–1:10. <https://doi.org/10.1145/3385959.3418450>
- [9] Daniel Ashbrook, Patrick Baudisch, and Sean White. 2011. Nenya: subtle and eyes-free mobile input with a magnetically-tracked finger ring. In *Proceedings of the International Conference on Human Factors in Computing Systems, CHI 2011, Vancouver, BC, Canada, May 7-12, 2011*. ACM, New York, NY, 2043–2046. <https://doi.org/10.1145/1978942.1979238>
- [10] Chien-Yen Chang, Belinda Lange, Mi Zhang, Sebastian Koenig, Phil Requejo, Noom Somboon, Alexander A. Sawchuk, and Albert A. Rizzo. 2012. Towards pervasive physical rehabilitation using Microsoft Kinect. In *6th International Conference on Pervasive Computing Technologies for Healthcare, PervasiveHealth 2012 and Workshops, San Diego, CA, USA, May 21-24, 2012*. IEEE, New York, 159–162.

- <https://doi.org/10.4108/icst.pervasivehealth.2012.248714>
- [11] Ke-Yu Chen, Kent Lyons, Sean White, and Shwetak N. Patel. 2013. uTrack: 3D input using two magnetic sensors. In *The 26th Annual ACM Symposium on User Interface Software and Technology, UIST'13, St. Andrews, United Kingdom, October 8-11, 2013*. ACM, New York, NY, 237–244. <https://doi.org/10.1145/2501988.2502035>
 - [12] Ke-Yu Chen, Shwetak N. Patel, and Sean J. Keller. 2016. Finexus: Tracking Precise Motions of Multiple Fingertips Using Magnetic Sensing. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems, San Jose, CA, USA, May 7-12, 2016*. ACM, New York, NY, 1504–1514. <https://doi.org/10.1145/2858036.2858125>
 - [13] Davide Chicco and Giuseppe Jurman. 2020. The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation. *BMC genomics* 21, 1 (2020), 1–13.
 - [14] Robert B Cleveland, William S Cleveland, Jean E McRae, and Irma Terpenning. 1990. STL: A seasonal-trend decomposition. *J. Off. Stat.* 6, 1 (1990), 3–73.
 - [15] Andrea Dittadi, Sebastian Dziadzio, Darren Cosker, Ben Lundell, Thomas J. Cashman, and Jamie Shotton. 2021. Full-Body Motion from a Single Head-Mounted Device: Generating SMPL Poses from Partial Observations. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*. IEEE, New York, NY, 11667–11677. <https://doi.org/10.1109/ICCV48922.2021.01148>
 - [16] Christian Forster, Luca Carlone, Frank Dellaert, and Davide Scaramuzza. 2017. On-Manifold Preintegration for Real-Time Visual-Inertial Odometry. *IEEE Trans. Robotics* 33, 1 (2017), 1–21. <https://doi.org/10.1109/TRO.2016.2597321>
 - [17] Alfred M. Franz, Tamás Haidegger, Wolfgang Birkfellner, Kevin Cleary, Terry M. Peters, and Lena Maier-Hein. 2014. Electromagnetic Tracking in Medicine - A Review of Technology, Validation, and Applications. *IEEE Trans. Medical Imaging* 33, 8 (2014), 1702–1725. <https://doi.org/10.1109/TMI.2014.2321777>
 - [18] Chris Harrison and Scott E. Hudson. 2009. Abracadabra: wireless, high-precision, and unpowered finger input for very small mobile devices. In *Proceedings of the 22nd Annual ACM Symposium on User Interface Software and Technology, Victoria, BC, Canada, October 4-7, 2009*. ACM, New York, NY, 121–124. <https://doi.org/10.1145/1622176.1622199>
 - [19] Phuong Hoang, Michael Hamilton, Joseph Murray, Corey Stafford, and Hien Tran. 2015. A Dynamic Feature Selection Based LDA Approach to Baseball Pitch Prediction. In *Trends and Applications in Knowledge Discovery and Data Mining - PAKDD 2015 Workshops: BigPMA, VLSP, QIMIE, DAEBH, Ho Chi Minh City, Vietnam, May 19-21, 2015. Revised Selected Papers (Lecture Notes in Computer Science, Vol. 9441)*. Springer, 125–137. https://doi.org/10.1007/978-3-319-25660-3_11
 - [20] HTC. 2020. VIVE. <https://www.vive.com>
 - [21] Yinghao Huang, Manuel Kaufmann, Emre Aksan, Michael J. Black, Otmar Hilliges, and Gerard Pons-Moll. 2018. Deep inertial poser: learning to reconstruct human pose from sparse inertial measurements in real time. *ACM Trans. Graph.* 37, 6 (2018), 185. <https://doi.org/10.1145/3272127.3275108>
 - [22] Jiaxi Jiang, Paul Strelci, Huajian Qiu, Andreas Fender, Larissa Laich, Patrick Snape, and Christian Holz. 2022. AvatarPoser: Articulated Full-Body Pose Tracking from Sparse Motion Sensing. *CoRR* abs/2207.13784 (2022). <https://doi.org/10.48550/arXiv.2207.13784>
 - [23] Manuel Kaufmann, Yi Zhao, Chengcheng Tang, Lingling Tao, Christopher D. Twigg, Jie Song, Robert Wang, and Otmar Hilliges. 2021. EM-POSE: 3D Human Pose Estimation from Sparse Electromagnetic Trackers. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*. IEEE, New York, NY, 11490–11500. <https://doi.org/10.1109/ICCV48922.2021.01131>
 - [24] Jungha Kim, Minkyeong Song, Yeoeun Lee, Moonkyeong Jung, and Pyojin Kim. 2022. An Empirical Evaluation of Four Off-the-Shelf Proprietary Visual-Inertial Odometry Systems. *CoRR* abs/2207.06780 (2022). <https://doi.org/10.48550/arXiv.2207.06780>
 - [25] Volodymyr V. Kindratenko. 1999. Calibration of electromagnetic tracking devices. *Virtual Real.* 4, 2 (1999), 139–150. <https://doi.org/10.1007/BF01408592>
 - [26] Volodymyr V. Kindratenko. 2000. A survey of electromagnetic position tracker calibration techniques. *Virtual Real.* 5, 3 (2000), 169–182. <https://doi.org/10.1007/BF01409422>
 - [27] Volodymyr V. Kindratenko and William R. Sherman. 2005. Neural network-based calibration of electromagnetic tracking systems. *Virtual Real.* 9, 1 (2005), 70–78. <https://doi.org/10.1007/s10055-005-0005-3>
 - [28] Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
 - [29] Jack B. Kuipers. 1980. SPASYN-an electromagnetic relative position and orientation tracking system. *IEEE Transactions on Instrumentation and Measurement* 29, 4 (Dec. 1980), 462–466. <https://doi.org/10.1109/tim.1980.4314980>
 - [30] Hyunchul Lim, Yaxuan Li, Matthew Dressa, Fang Hu, Jae Kim, Ruidong Zhang, and Cheng Zhang. [n.d.]. ([n. d.]). ([n. d.]).
 - [31] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. 2015. SMPL: a skinned multi-person linear model. *ACM Trans. Graph.* 34, 6 (2015), 248:1–248:16. <https://doi.org/10.1145/2816795.2818013>
 - [32] Azumi Maekawa, Seito Matsubara, Sohei Wakisaka, Daisuke Uriu, Atsushi Hiyama, and Masahiko Inami. 2020. Dynamic Motor Skill Synthesis with Human-Machine Mutual Actuation. In *CHI '20: CHI Conference on Human Factors in Computing Systems, Honolulu, HI, USA, April 25-30, 2020*. ACM, New York, NY, 1–12. <https://doi.org/10.1145/3313831.3376705>
 - [33] Magic Leap. 2022. Magic Leap 2. <https://www.magicleap.com/ml2-devices>

- [34] Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael J. Black. 2019. AMASS: Archive of Motion Capture As Surface Shapes. In *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*. IEEE, New York, NY, 5441–5450. <https://doi.org/10.1109/ICCV.2019.00554>
- [35] Meta Technologies LLC. 2022. Meta Quest Pro. <https://www.oculus.com/quest>
- [36] Microsoft. 2010. Kinect. <https://en.wikipedia.org/wiki/Kinect>
- [37] Vimal Mallyn, Riku Arakawa, Mayank Goel, Chris Harrison, and Karan Ahuja. 2023. IMUPoser: Full-Body Pose Estimation using IMUs in Phones, Watches, and Earbuds. In *CHI '23: CHI Conference on Human Factors in Computing Systems, Hamburg, Germany, 22 April 2022 - 29 April 2022*. ACM, in press. <https://doi.org/10.1145/3544548.3581392>
- [38] NaturalPoint Inc. 2020. OptiTrack. <http://optitrack.com>
- [39] Evonne Ng, Donglai Xiang, Hanbyul Joo, and Kristen Grauman. 2020. You2Me: Inferring Body Pose in Egocentric Video via First and Second Person Interactions. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*. Computer Vision Foundation / IEEE, New York, 9887–9897. <https://doi.org/10.1109/CVPR42600.2020.00991>
- [40] Nordic. 2019. Enhanced ShockBurst (ESB). https://developer.nordicsemi.com/nRF_Connect_SDK/doc/1.0.0/nrf/ug_esb.html
- [41] Boris N. Oreshkin, Dmitri Carpow, Nicolas Chapados, and Yoshua Bengio. 2020. N-BEATS: Neural basis expansion analysis for interpretable time series forecasting. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- [42] Farshid Salemi Parizi, Eric Whitmire, and Shwetak N. Patel. 2019. AuraRing: Precise Electromagnetic Finger Tracking. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 3, 4 (2019), 150:1–150:28. <https://doi.org/10.1145/3369831>
- [43] Keun-Woo Park, Daehwa Kim, Seongkook Heo, and Geehyuk Lee. 2020. MagTouch: Robust Finger Identification for a Smartwatch Using a Magnet Ring and a Built-in Magnetometer. In *CHI '20: CHI Conference on Human Factors in Computing Systems, Honolulu, HI, USA, April 25-30, 2020*. ACM, New York, 1–13. <https://doi.org/10.1145/3313831.3376234>
- [44] Valter Pasku, Alessio De Angelis, Guido De Angelis, Darmindra D. Arumugam, Marco Dionigi, Paolo Carbone, Antonio Moschitta, and David S. Ricketts. 2017. Magnetic Field-Based Positioning Systems. *IEEE Commun. Surv. Tutorials* 19, 3 (2017), 2003–2017. <https://doi.org/10.1109/COMST.2017.2684087>
- [45] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and Michael J. Black. 2019. Expressive Body Capture: 3D Hands, Face, and Body from a Single Image. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*.
- [46] Frederick H Raab, Ernest B Blood, Terry O Steiner, and Herbert R Jones. 1979. Magnetic position and orientation tracking system. *IEEE Transactions on Aerospace and Electronic systems* 5 (1979), 709–718.
- [47] Daniel Roetenberg, Per J. Slycke, and Peter H. Veltink. 2007. Ambulatory Position and Orientation Tracking Fusing Magnetic and Inertial Sensing. *IEEE Trans. Biomed. Eng.* 54, 5 (2007), 883–890. <https://doi.org/10.1109/TBME.2006.889184>
- [48] Mili Shah. 2013. Solving the Robot-World/Hand-Eye Calibration Problem Using the Kronecker Product. *Journal of Mechanisms and Robotics* 5, 3 (June 2013). <https://doi.org/10.1115/1.4024473>
- [49] Sheng Shen, He Wang, and Romit Roy Choudhury. 2016. I am a Smartwatch and I can Track my User’s Arm. In *Proceedings of the 14th Annual International Conference on Mobile Systems, Applications, and Services, MobiSys 2016, Singapore, June 26-30, 2016*. ACM, New York, NY, 85–96. <https://doi.org/10.1145/2906388.2906407>
- [50] Hideki Shimobayashi, Tomoya Sasaki, Arata Horie, Riku Arakawa, Zendai Kashino, and Masahiko Inami. 2021. Independent Control of Supernumerary Appendages Exploiting Upper Limb Redundancy. In *AHs ’21: Augmented Humans Conference 2021, Rovaniemi, Finland, February 22-24, 2021*, Jonna Häkkilä and Paul Strohmeier (Eds.). ACM, New York, 19–30. <https://doi.org/10.1145/3458709.3458980>
- [51] Takaaki Shiratori, Hyun Soo Park, Leonid Sigal, Yaser Sheikh, and Jessica K. Hodgins. 2011. Motion capture from body-mounted cameras. *ACM Trans. Graph.* 30, 4 (2011), 31. <https://doi.org/10.1145/2010324.1964926>
- [52] Sixense. 201. Razer Hydra. <https://dl.razerzone.com/master-guides/Hydra/HydraOMG-ENG.pdf>
- [53] Snap. 2018. Spectacles. <https://www.spectacles.com/>
- [54] Sony Interactive Entertainment LLC. 2020. PlayStation VR. <https://www.playstation.com/en-us/ps-vr/>
- [55] Shuya Suda, Yasutoshi Makino, and Hiroyuki Shinoda. 2019. Prediction of Volleyball Trajectory Using Skeletal Motions of Setter Player. In *Proceedings of the 10th Augmented Human International Conference 2019, Reims, France, March 11-12, 2019*. ACM, New York, NY, 16:1–16:8. <https://doi.org/10.1145/3311823.3311844>
- [56] Jochen Tautges, Arno Zinke, Björn Krüger, Jan Baumann, Andreas Weber, Thomas Helten, Meinard Müller, Hans-Peter Seidel, and Bernd Eberhardt. 2011. Motion reconstruction using sparse accelerometer data. *ACM Trans. Graph.* 30, 3 (2011), 18:1–18:12. <https://doi.org/10.1145/1966394.1966397>
- [57] Yushuang Tian, Xiaoli Meng, Dapeng Tao, Dongquan Liu, and Chen Feng. 2015. Upper limb motion tracking with the integration of IMU and Kinect. *Neurocomputing* 159 (2015), 207–218. <https://doi.org/10.1016/j.neucom.2015.01.071>
- [58] Denis Tomè, Patrick Peluse, Lourdes Agapito, and Hernán Badino. 2019. xR-EgoPose: Egocentric 3D Human Pose From an HMD Camera. In *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*. IEEE, New York, 7727–7737. <https://doi.org/10.1109/ICCV.2019.00782>

- [59] Vicon Motion Systems Ltd. 2020. Vicon. <https://vicon.com>
- [60] Daniel Vlasic, Rolf Adelsberger, Giovanni Vannucci, John Barnwell, Markus H. Gross, Wojciech Matusik, and Jovan Popovic. 2007. Practical motion capture in everyday surroundings. *ACM Trans. Graph.* 26, 3 (2007), 35. <https://doi.org/10.1145/1276377.1276421>
- [61] Timo von Marcard, Bodo Rosenhahn, Michael J. Black, and Gerard Pons-Moll. 2017. Sparse Inertial Poser: Automatic 3D Human Pose Estimation from Sparse IMUs. *Comput. Graph. Forum* 36, 2 (2017), 349–360. <https://doi.org/10.1111/cgf.13131>
- [62] Eric Whitmire, Farshid Salemi Parizi, and Shwetak N. Patel. 2019. Aura: Inside-out Electromagnetic Controller Tracking. In *Proceedings of the 17th Annual International Conference on Mobile Systems, Applications, and Services, MobiSys 2019, Seoul, Republic of Korea, June 17-21, 2019*. ACM, New York, 300–312. <https://doi.org/10.1145/3307334.3326090>
- [63] Alexander Winkler, Jungdam Won, and Yuting Ye. 2022. QuestSim: Human Motion Tracking from Sparse Sensors with Simulated Avatars. In *SIGGRAPH Asia 2022 Conference Papers*. ACM, New York. <https://doi.org/10.1145/3550469.3555411>
- [64] Chengshuo Xia, Xinrui Fang, Riku Arakawa, and Yuta Sugiura. 2022. VoLearn: A Cross-Modal Operable Motion-Learning System Combined with Virtual Avatar and Auditory Feedback. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 6, 2 (2022), 81:1–81:26. <https://doi.org/10.1145/3534576>
- [65] Fanglu Xie, Xina Cheng, and Takeshi Ikenaga. 2017. Motion State Detection Based Prediction Model for Body Parts Tracking of Volleyball Players. In *Advances in Multimedia Information Processing - PCM 2017 - 18th Pacific-Rim Conference on Multimedia, Harbin, China, September 28-29, 2017, Revised Selected Papers, Part I (Lecture Notes in Computer Science, Vol. 10735)*. Springer, 280–289. https://doi.org/10.1007/978-3-319-77380-3_27
- [66] Weipeng Xu, Avishek Chatterjee, Michael Zollhöfer, Helge Rhodin, Pascal Fua, Hans-Peter Seidel, and Christian Theobalt. 2019. Mo²Cap²: Real-time Mobile 3D Motion Capture with a Cap-mounted Fisheye Camera. *IEEE Trans. Vis. Comput. Graph.* 25, 5 (2019), 2093–2101. <https://doi.org/10.1109/TVCG.2019.2898650>
- [67] Kenji Yamanishi, Jun’ichi Takeuchi, Graham J. Williams, and Peter Milne. 2004. On-Line Unsupervised Outlier Detection Using Finite Mixtures with Discounting Learning Algorithms. *Data Min. Knowl. Discov.* 8, 3 (2004), 275–300. <https://doi.org/10.1023/B:DAMI.0000023676.72185.7c>
- [68] Dongseok Yang, Doyeon Kim, and Sung-Hee Lee. 2021. LoBSTR: Real-time Lower-body Pose Prediction from Sparse Upper-body Tracking Signals. *Comput. Graph. Forum* 40, 2 (2021), 265–275. <https://doi.org/10.1111/cgf.142631>
- [69] Xinyu Yi, Yuxiao Zhou, Marc Habermann, Soshi Shimada, Vladislav Golyanik, Christian Theobalt, and Feng Xu. 2022. Physical Inertial Poser (PIP): Physics-aware Real-time Human Motion Tracking from Sparse Inertial Sensors. *CoRR* abs/2203.08528 (2022). <https://doi.org/10.48550/arXiv.2203.08528>
- [70] Xinyu Yi, Yuxiao Zhou, and Feng Xu. 2021. TransPose: real-time 3D human translation and pose estimation with six inertial sensors. *ACM Trans. Graph.* 40, 4 (2021), 86:1–86:13. <https://doi.org/10.1145/3450626.3459786>
- [71] Zheng Zhao, Weihai Chen, Yang Li, Jiahua Wang, and Zhongcai Pei. 2020. A Wearable Body Motion Capture System and Its Application in Assistive Exoskeleton Control. In *2020 15th IEEE Conference on Industrial Electronics and Applications (ICIEA)*. IEEE. <https://doi.org/10.1109/iciea48937.2020.9248141>

A DETAILED IMPLEMENTATION OF IK

Our hardware setup has two subsystems for measuring the pose of certain body parts: AR glasses that use VIO tracking and wrist-worn EMF sensors. Accordingly, we define the coordinates for our system as follows:

- **World Global Coordinate:** The coordinate where AR glasses provide the absolute positions $P_g^W \in \mathbb{R}^{1 \times 3}$ and orientations in axis-angle representation $\Phi_g^W \in \mathbb{R}^{1 \times 3}$, where g denotes glasses and W denotes the world coordinate.
- **HMD Local Coordinate:** This is the local coordinate of the AR glasses with the origin in the center of the AR glasses. We represent sensor positions and rotations in HMD local coordinate as $P_s^{HMD} \in \mathbb{R}^{1 \times 3}$ and $\Phi_s^{HMD} \in \mathbb{R}^{1 \times 3}$, where s denotes sensor and HMD denotes the HMD coordinate.
- **EMF Local Coordinate:** The coordinate where two wrist-worn EMF sensors are tracked relative to the EMF source on the head. We represent their positions and rotations as $P_s^{EMF} \in \mathbb{R}^{1 \times 3}$ and $\Phi_s^{EMF} \in \mathbb{R}^{1 \times 3}$, where s denotes sensor and EMF denotes the EMF coordinate.
- **Body Local Coordinate:** Human body pose can be represented by each joint position $P_j^B \in \mathbb{R}^{1 \times 3}$ and orientation $\Phi_j^B \in \mathbb{R}^{1 \times 3}$ in the body local coordinate, where j denotes body joints and B denotes the body coordinate.

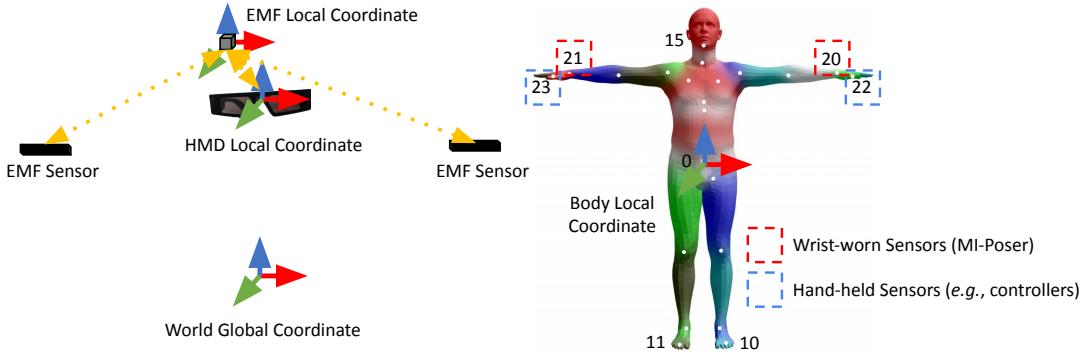


Fig. 13. Coordinate definitions of MI-Poser and SMPL model joints. MI-Poser maps the HMD and two EMF sensor positions and rotations to the joints 15, 20, and 21 and outputs positions and rotations of 22 joints (from joint 0 to joint 21) for generating 3D body mesh. Typical existing systems (e.g., [4, 35]) use hand-held controllers at joint 22 and joint 23.

The joints defined in the kinematic tree of the SMPL model are shown in Figure 13, where the root is the pelvis (joint 0). We used 22 joints, from 0 to 21, for pose reconstruction, while existing work using hand-held controllers uses 24 joints from 0 to 23. This difference comes from our setting of a wrist-worn setup; previous work uses hand-held devices for hand tracking. This setting can help the model estimate arm pose better (excluding hands) by reducing the forearms' total degree of freedom.

The SMPL model takes the relative rotations of all the joints as input and outputs the 3D body mesh through a learned rigged template mesh with 6,890 vertices. The body shape S is defined by identity-dependent shape parameters β , pose parameters θ , soft deformation ϕ as the following equation:

$$S(\beta, \theta, \phi) = G(T(\beta, \theta, \phi), J(\beta), \theta, w) \quad (4)$$

, where $T(\beta, \theta, \phi)$ represents the body vertices in rest pose as shown in Figure 13, β are the parameters to define the body shape, $\theta = \{\Phi_j^B\}_{j=0,\dots,21} \in \mathbb{R}^{3 \times 22}$ denotes body pose (*i.e.*, relative rotation angles between adjacent body joints), and ϕ are parameters to define soft tissue dynamic deformations. $J \in \mathbb{R}^{3 \times 22}$ denotes joint locations, and $w \in \mathbb{R}^{4 \times 3 \times 22}$ are the blend weights. G is a learned machine learning human pose prior model which maps the joints to 3D body mesh [45]. In MI-Poser, we assume the parameters β and ϕ are constant, which means we use a standard body shape in our IK model. According to Equation 4, we need to estimate rotations $\theta \in \mathbb{R}^{3 \times 22}$ for generating the body shape of a standard body shape.

In our settings (See Figure 1), the user wears AR glasses, two EMF sensors on the wrist, and the EMF source on the back of the head which has a fixed relative pose to AR glasses. By coordinate transformations, we can get the absolute pose of the AR glasses and two EMF sensors in the body coordinate. Specifically, with calibration, we map the AR glasses to joint 15 and two EMF sensors to joints 20 and 21 in Figure 13. Then, our IK problem can be formulated as,

$$\theta = \text{IK-MODEL}(\{P_j^B, \Phi_j^B\}_{j \in \{15, 20, 21\}}) \quad (5)$$

, where P_j^B and Φ_j^B are the positions and orientations of joint j . Our problem becomes learning the *IK-MODEL* function.

To examine how our setup difference (*i.e.*, wrist-worn v.s. hand-held sensors) affects the pose reconstruction performance, we trained a model adapted from the state-of-the-art, which is AvatarPoser [22]. We trained the same model architecture while providing a different configuration of sensor locations (*i.e.*, joints 15, 20, and 21 instead of joints 15, 22, and 23, as shown in Figure 13). Multiple frames are used to predict θ_t in the model.

In detail, the model takes 40 previous frames (equivalent to 2/3 seconds at 60 FPS), so the IK can be written as $\theta_t = \text{IK-MODEL}(\{P_{\tau,j}^B, \Phi_{\tau,j}^B\}_{j \in \{15, 20, 21\}, \tau=t-39, t-38, \dots, t})$. This model considers global transition, and thus it is suitable for application scenarios involving body movements, such as room-scale games.

B DETAILED PREPROCESSING IN USER STUDY 2

B.1 Ground Truth Calibration

We conducted a preprocessing as a calibration for each session. The data collection setup has two coordinates: the EMF coordinate and iPhone coordinate. We denote the homogeneous transformation matrix of the iPhone with respect to its origin as $X_t \in \mathbb{R}^{4 \times 4}$, given the time t . Similarly, we denote the homogeneous transformation matrix of the EMF sensor with respect to its origin as $Y_t \in \mathbb{R}^{4 \times 4}$. In addition, since the origin of the iPhone coordinate and the EMF coordinate are both fixed, their transformation is constant. We denote the homogeneous transformation matrix of the iPhone coordinate's origin to the EMF coordinate's origin as $A \in \mathbb{R}^{4 \times 4}$. Similarly, since both the iPhone and the EMF sensor are attached to a rigid body, their transformation is also constant. We denote the homogeneous transformation matrix of the iPhone to the EMF sensor as $B \in \mathbb{R}^{4 \times 4}$. Using these notations, as to the EMF sensor transformation with respect to the origin of the iPhone coordinate, we have the following equation,

$$X_t B = A Y_t \quad (6)$$

This equation is known to be solvable through a data-driven numerical approach [48]. We used the first 20 seconds of data as a calibration (recall that we paid attention that there was no metal close to the sensor during the time) to obtain \hat{A} and \hat{B} , an estimate of A and B , respectively. Once we calculate them, we can get the ground truth pose data of the EMF sensor with respect to the origin of the EMF coordinate; its homogeneous transformation matrix can be written as $\hat{A}^{-1} X_t \hat{B}$.

B.2 IMU Calibration

In addition, since the IMU values are measured in their local coordinates at each frame, we needed to convert them into the EMF coordinate. Here, we can assume the position and orientation are identical between the IMU and EMF sensors, embedded in the same PCB. As we discussed in Section 3.2, there are two ways to represent rotation at time $t + \Delta t$: $\Phi_s^{EMF}(t + \Delta t)$, and $\Phi_s^{EMF}(t) + \Delta \Phi_s^{IMU}(t) \times \Delta t$. We use the former if $\hat{I}(t + \Delta t) = 0$ and otherwise use the latter. We denote the acceleration of the IMU sensor with respect to the origin of the EMF coordinate as $a(t) \in \mathbb{R}^{1 \times 3}$, while denoting the measured acceleration in the IMU's local coordinate as $a^{local}(t) \in \mathbb{R}^{1 \times 3}$. By rotating $a^{local}(t)$ using $\Phi_s^{EMF}(t)$, we can obtain $a(t)$. Moreover, since $a^{local}(t)$ contains the gravity, we need a calibration to cancel it from a_t . Specifically, we identified the static moments in the calibration phase based on the EMF position data and compute the average $a(t)$ during the moments, as a_{base} , which represents the gravity component. Then, we used $a(t) - a_{base}$ as the linear acceleration of the sensor at time t in the EMF sensor coordinate.