# Leveraging Motion Capture System for High Accuracy AR-Assisted Assembly

**Hanning Liu, Xingjie Xie, Yujiao Li, Xiaofan Gao, Honglei Wu, Yao Zhang and Philip F. Yuan** [1]

**Abstract.** Augmented Reality (AR) allows workers to construct buildings accurately and intuitively without the need for traditional tools like 2-D drawings and rulers. However, accurately tracking worker's pose remains a significant challenge in existing experiments due to their continuous and irregular movement. This research discusses a series of methods using cameras and algorithms to achieve the 6-DoF pose tracking function and reveal the relationship between each method and corresponding tracking accuracy in order to figure out a robust approach of AR-assisted assembly. This paper begins with a consideration of the possible limitations of existing methods including the image drift associated with visual SLAM and the time-consuming nature of fiducial markers. Next, the entire hardware and software framework was introduced, which elaborates on how the motion capture system is integrated into the AR-assisted assembly system. Then, some experiments have been carried out to demonstrate the connection between the system set up and pose tracking accuracy. This research shows the possibility to easily finish assembly task based on AR technology by integrating motion capture system.

**Keywords:** Augmented Reality · Pose tracking · AR-assisted assembly · Motion capture system · Freeform steel structure

## 1 Introduction

The Architecture, Engineering, and Construction (AEC) industry is embracing augmented reality (AR) technologies to speed up the design process, improve the construction quality and ensure the safety of construction workers [1]. Specifically, in the field of AR-assisted construction processes, workers are allowed to con-

[1] Philip F. Yuan (✉)
College of Architecture and Urban Planning (CAUP), Tongji University, Shanghai, China
e-mail: philipyuan007@tongji.edu.cn

struct buildings accurately and intuitively without the need for traditional tools like 2-D drawings and rulers [2]. This not only enhances convenience for workers but also facilitates the observation of streamlined or parametrical designs that may be challenging to articulate through orthographic projection and dimensional annotation [3].

Moreover, there have already been numerous finished art installation practices worldwide utilizing AR technology, featuring variant materials and tectonic approaches such as a bent steel pipe pavilion [4], a lightweight timber structure [6], a collaboratively crafted bamboo weaving sculpture [9], etc. These practices are made possible by overlaying digital models onto real-world construction sites, facilitated by 6-DoF pose tracking advancements.

Robust 6-DoF pose tracking is a long-standing and well-established area of AR research [10]. From some of the earliest work on computer vision algorithms such as Direct Linear Transform (DLT) and Perspective-n-Point (PnP) [7] [8], to more recent work on multi-sensor data fusion approaches, integrating various types of data including those from inertial measurement units (IMUs) and simultaneous localization and mapping (SLAM) systems [11].

Despite much technical progress, however, the majority of AR-assisted construction projects nowadays still face a pose tracking issue and are only feasible for scenarios involving building materials and tectonics with significant tolerance allowances. In this paper, we investigate how to effectively improve the accuracy between digital overlays and the physical entities by integrating a 6-DoF motion capture pose tracking system into the assembly construction process. A handle-like AR device equipped with a power unit, computing unit, display unit and camera unit with highly reflective markers attached is introduced. This research shows that it outperforms more commonly used pose tracking methods in terms of accuracy, pose tracking scope, and assembly time consuming.

## 2 Related works

In the 6-DoF pose tracking research field, two primary strategies are employed to tackle this challenge: marker-less and marker-based. The marker-less strategy commonly utilizes visual SLAM and multi-sensor data fusion technology [12]. The marker-based strategy relies on highly reflective markers or fiducial markers. Fiducial markers typically involve ArUco Marker and AprilTags [13] [14].

Visual SLAM algorithm tracks 6-DoF pose by analyzing how feature points within the environment change over time [15]. However, construction sites are often open and empty, present repeating elements and dynamic changes, which pose significant challenges to accurate pose tracking. Thus, existing research indicates that marker-less tracking systems are susceptible to error accumulation over extended periods [16].

Compared to visual SLAM algorithm, fiducial marker detection and pose estimation typically require fewer computational resources, making them suitable for real-time applications on mobile devices. Besides it is much cheaper than expensive pose tracking device [18]. Highly reflective markers, once attached to the device requiring pose tracking, remain in place and require no further adjustments, offering convenience and saving considerable time [17].

The subsequent section presents a detailed discussion on the feasibility of three distinct marker-based strategies: fiducial markers only, visual SLAM combined with fiducial markers and motion capture systems combined with highly reflective markers.

## *2.1 Fiducial Markers Only*

An experiment using a designer-friendly Grasshopper plugin is conducted to test the accuracy and stability of fiducial markers-only strategy. While it exhibits impressive result when the camera is close to fiducial markers, its performance suffers under real-world conditions. Camera shake or increased distance from markers leads to frequent loss of camera pose, hindering its effectiveness. Additionally, the system's dependence on markers within the camera's view makes it unusable in common assembly scenarios where markers not always be seen. So fiducial markers only strategy isn't robust enough to be used in AR-assisted assembly scenario.

## *2.2 Visual SLAM and Fiducial markers*

A method of combining multiple fiducial markers with multi-sensor data fusion approaches has been proposed to enhance 6-DoF pose tracking accuracy [5]. However, the localized property of fiducial markers makes them can only enhance pose tracking accuracy within the vicinity of where the markers are placed. As a result, after completing construction tasks in one area, the process of setting up a new set of markers for each subsequent assembly area is necessary to achieve the desired level of accuracy. So, the time-consuming process of setting up a large number of fiducial markers needs to be optimized.

Besides, an experiment using Fologram on a mobile phone is conducted to evaluate the accuracy of the visual SLAM combined with fiducial markers strategy. The structure in this experiment consists of many round pipes and a concrete base, standing 3 meters tall. As fiducial markers can only be attached to flat planes, they are restricted to the ground and cannot be mounted on the round pipes, leading to reduced accuracy at high positions. Additionally, the rough surface of the concrete base further diminishes accuracy compared to an ideal scenar-

io. So, the stringent requirements of fiducial markers limit the scope of 6-DoF pose tracking and introduce unforeseen deviations.

## 2.3 Motion Capture System and Highly Reflective Markers

The motion capture system provides 6-DoF pose tracking with deviations of less than 1mm, establishing it as a widely accepted benchmark in research papers for ground truth data comparison [17]. This system enables pose tracking across its entire field of view, unhindered by non-flat elements. Additionally, the highly reflective markers attached to the device eliminate the need for frequent setup. However, despite its prominence, there is currently no research integrating the motion capture system with an image overlay feature to enable AR-assisted assembly work.

## 3 Methods

The following experiments were conducted in an 8.0 * 7.6m site (Fig 1), using 6-DoF motion capture system provided by NOKOV and a custom-designed AR device. The 6-DoF motion capture system comprises 8 high-speed infrared optical lenses, each equipped with numerous infrared LED lights mounted on their panels. These lights illuminate highly reflective markers, which then reflect light back to the lenses, enabling the calculation of the 2D coordinates of each marker. Finally, by integrating data from multiple lenses capturing different perspectives, the system computes the 3D coordinates of each highly reflective marker.
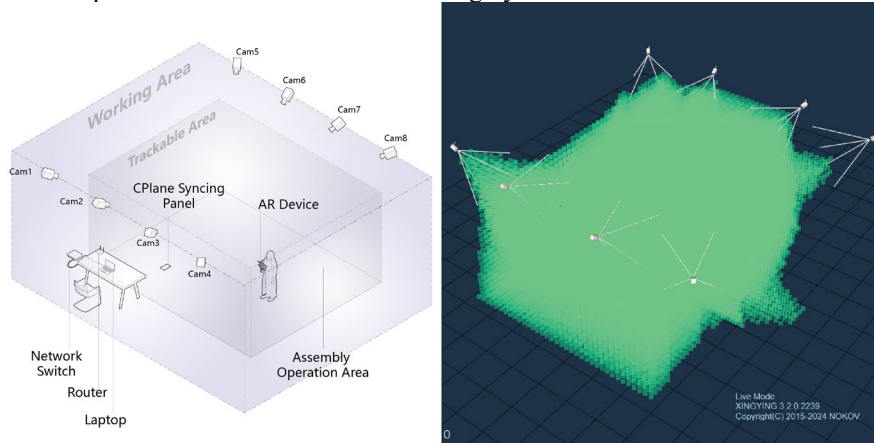


**Fig. 1** AR-assisted assembly site setup (Left) and 6-DoF pose tracking area (Right)

With three or more markers, a rigid body can be defined within the motion capture system. Rigid bodies derive their pose from the average position of all markers and are resilient to the loss of individual markers, resulting in enhanced stability and reliability. The coordinate system of the motion capture system is established during the calibration phase and remains fixed thereafter, yet it becomes invisible in the real-world space post-calibration. To ensure synchronization between real-world and digital spaces, a rigid body comprising four markers acts as the reference coordinate system throughout the tracking process.

The custom-designed AR device is comprised of five essential units (Fig 2): an image-capturing unit, a 6-DoF pose tracking unit, a computation unit, a display unit, and a power unit. The image-capturing unit features a 5-million-pixel camera tasked with capturing real-world imagery. Positioned in close proximity to the 6-DoF pose tracking unit, its captured images are overlaid with those captured by the digital camera. The 6-DoF pose tracking unit is equipped with five highly reflective markers. Three of these markers are strategically positioned near the image-capturing unit, facilitating straightforward calculation of their transform relations. The remaining two markers are placed on the left and right sides of the device to enhance visibility to the high-speed infrared optical lenses. The computation unit, a compact Intel NUC, serves as the core of the custom-designed AR device. It manages all tasks related to receiving 6-DoF pose data, rendering digital models, and overlaying images. The display unit is a screen with a resolution of 1024 * 768, perfectly matching the aspect ratio of the camera image sensor, enabling full-screen display of images. The power unit consists of a Li-Po battery with a capacity of 4000 mAh, providing continuous operation of the device for over 4 hours.
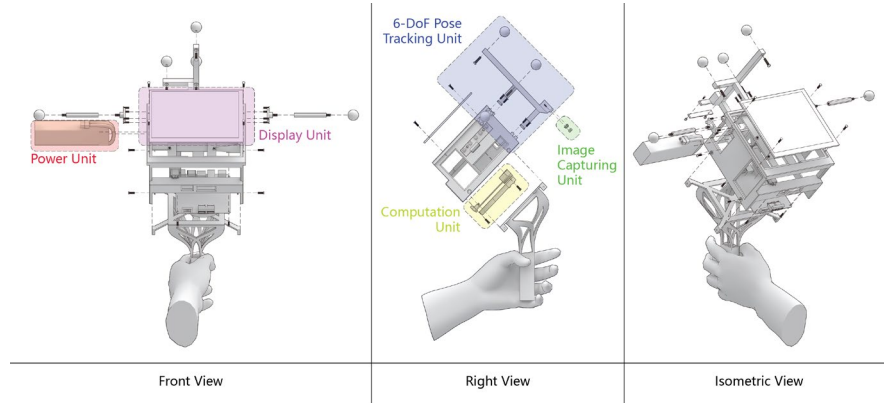


**Fig. 2** Custom-designed AR device with five units

This AR-assisted assembly system comprises two computers. The first computer connects to the network switch of the motion capture system via a cable, while the other serves as the core of the custom-designed AR device. Each computer runs a different software framework, together forming the system. The

software framework of the first computer includes Rhino & Grasshopper and necessary tools to process data from high-speed infrared optical lenses. A self-developed plugin for Grasshopper facilitates the transfer of 6-DoF pose data from the NOKOV SDK to the Rhino & Grasshopper platform. Subsequently, the data is transmitted to the AR device via UDP and local Wi-Fi. The software framework of the second computer revolves around Unity. Rhino & Grasshopper on this computer function as a plugin for Unity, with the assistance of the open-source Rhino Inside project. It receives pose data transmitted from the first computer via UDP, then utilizes the Callback method provided by Rhino Inside to establish the pose of the digital camera in Unity and synchronize digital model information.
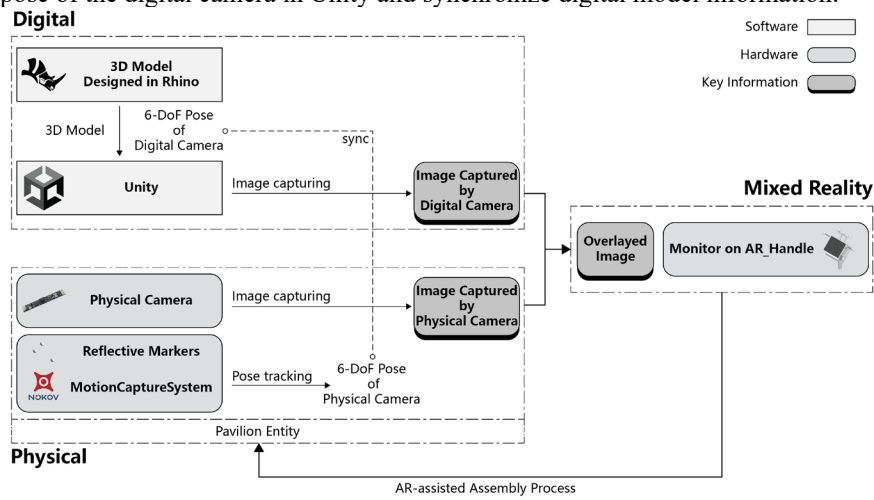


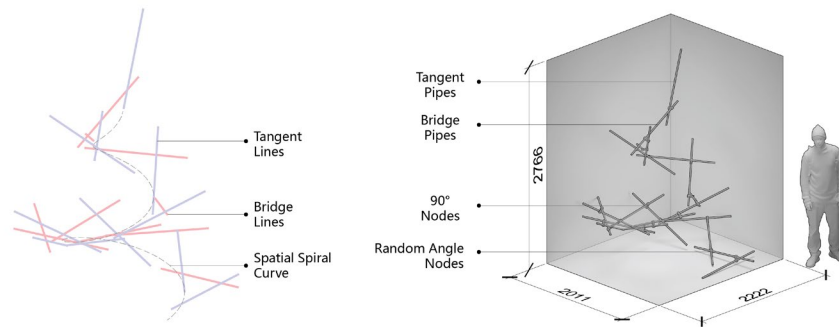**Fig. 3** AR-assisted assembly workflow



**Fig. 4** Design diagram (Left), model in digital environment (Right)

The AR-assisted assembly workflow encompasses three environments (Fig 3): the digital environment, the physical environment, and the mixed reality

environment. Designers begin by completing their design tasks in the digital environment. They then synchronize the digital and physical environments using the motion capture system to create a mixed reality environment. With the aid of the mixed reality environment, designers can achieve results with greater accuracy and intuition than was previously possible.

The assembly experiment pavilion design in this research is tailored to suit the need for testing the accuracy of the AR-assisted assembly workflow (Fig 4). It comprises 21 aluminum pipes, each 1 meter in length, and 20 connection nodes capable of rotating 360 degrees. Among the 21 pipes, 11 form tangent lines of a spatial spiral curve, while the remaining 10 act as bridges to connect with pipes ahead and behind them separately. Locating the positions of these pipes without AR assistance is challenging, allowing for the assessment of our workflow's effectiveness. The height of the pavilion is set to about 2.8m to test the 6-DoF pose tracking accuracy far from the ground. The connection nodes' angles are specially configured, with 10 set at 90 degrees and the remaining 10 at random angles, enabling the assembly process to be divided into multiple subprocesses.

# 4 Experiments

## 4.1 AR-assisted Assembly Process

The experiment comprises two preparation steps: camera parameter calibration in Unity, synchronization between the real-world and digital coordinate systems; and three assembly steps: assembly of the orthogonal pipes, fixation of the rotation angles for the nodes with random angles, and the final assembly of the pavilion.

At the beginning of the experiment series, the first step involves calibrating the parameters of the physical camera to prevent image distortion and ensure seamless overlay with the digital camera in Unity (Fig 5). This is achieved using a custom-developed Grasshopper plugin integrated with algorithms from the OpenCV library to estimate the camera's pose. Subsequently, Rhino Inside is employed to synchronize model information and camera pose with Unity. Finally, adjustments are made to the camera parameters of the digital camera in Unity to achieve perfect alignment between the images from the physical and digital cameras.

The second step is aligning the coordinate system between real-world and digital space (Fig 6). A panel equipped with four highly reflective markers is initially defined as a rigid body. The motion capture system recognizes it and provides pose information as a reference construction plane in Rhino. Subsequently, digital models are oriented with respect to the reference construction plane before being sent to Unity.
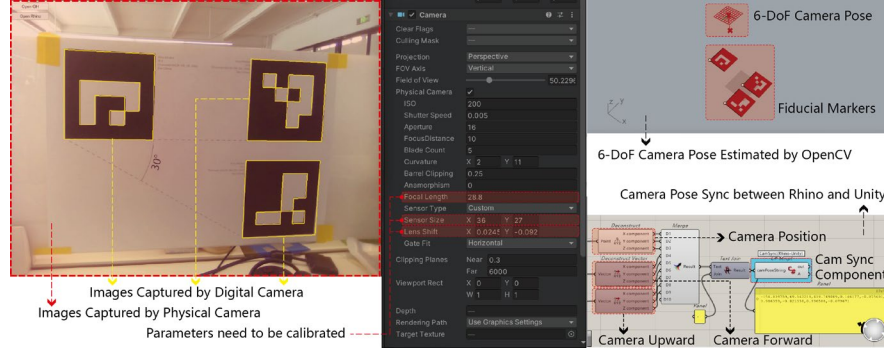
**Fig. 5** Overlaid images in Unity (Left), camera parameters in Unity (Middle), 6-DoF pose estimate and sync (Right)



**Fig. 6** Photo of the CPlane Syncing Panel (Left), panel tracked by motion capture system (Middle) and panel in Rhino space (Right)

In the first step of assembly process, the orthogonal pipes are assembled (Fig 7). The pavilion comprises 21 pipes, with 11 aligned tangentially to a spiral curve, and 10 acting as bridge pipes connecting the tangential ones. Each tangential pipe is perpendicular to its corresponding bridge pipe, collectively forming an assembly element. The position of the connection node, represented by one-dimensional data—the distance between the node and the end of each pipe—is the only variable among these elements. Initially, all the orthogonal elements are arranged in the Rhino space, then a custom-designed AR device is used to precisely locate the connection node in physical space. After confirming the overlaid images, an additional deviation check is performed solely using the motion capture system to evaluate the deviation with ground truth. The results of the check indicate deviations of less than 3.2mm, with an average value of 1.835mm (Table 1).
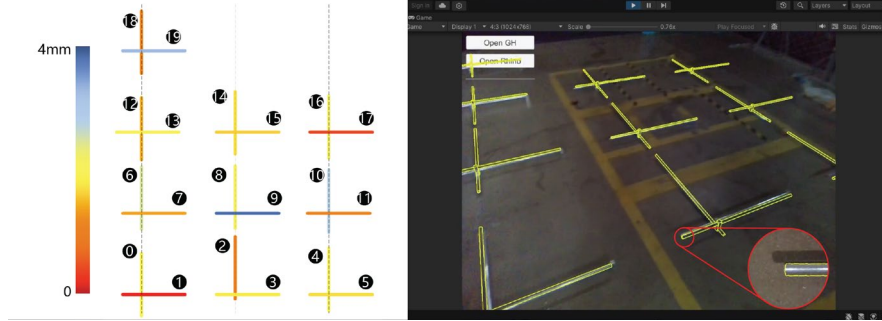
**Fig. 7** Position deviation of each pipe (Left), deviation check of overlaid image (Right)

**Table 1** Position deviation of the middle point of each pipe

| Pipe Index | X value Of Ideal Middle Point /mm | Y value Of Ideal Middle Point /mm | Z value Of Ideal Middle Point /mm | X value Of Real Middle Point /mm | Y value Of Real Middle Point /mm | Z value Of Real Middle Point /mm | Error /mm |
|---|---|---|---|---|---|---|---|
| 0 | 112.5 | 671.34 | 20 | 113.4 | 672.6 | 21.07 | 1.89 |
| 1 | 312.5 | 512.5 | 65 | 313.08 | 512.98 | 64.87 | 0.76 |
| 2 | 1612.5 | 942.05 | 20 | 1612.75 | 941.74 | 18.68 | 1.38 |
| 3 | 1812.5 | 512.5 | 65 | 1812.42 | 511.41 | 66.48 | 1.84 |
| 4 | 3112.5 | 754.54 | 20 | 3112.09 | 752.66 | 20.28 | 1.94 |
| 5 | 3312.5 | 512.5 | 65 | 3311.76 | 513.84 | 64.08 | 1.78 |
| 6 | 112.5 | 2055.11 | 20 | 111.43 | 2055.66 | 21.89 | 2.24 |
| 7 | 312.5 | 1812.5 | 65 | 311.1 | 1812.27 | 65.69 | 1.58 |
| 8 | 1612.5 | 2071.66 | 20 | 1610.77 | 2070.64 | 19.49 | 2.07 |
| 9 | 1812.5 | 1812.5 | 65 | 1814.45 | 1810.69 | 63.29 | 3.15 |
| 10 | 3112.5 | 2007.19 | 20 | 3114.12 | 2008.6 | 21.1 | 2.41 |
| 11 | 3271.5 | 1812.5 | 65 | 3272.79 | 1813.12 | 64.9 | 1.43 |
| 12 | 112.5 | 3167.43 | 20 | 113.46 | 3167.27 | 18.7 | 1.62 |
| 13 | 210.5 | 3112.5 | 65 | 211.13 | 3111.55 | 66.51 | 1.89 |
| 14 | 1612.5 | 3260.63 | 20 | 1612.8 | 3258.89 | 20.31 | 1.79 |
| 15 | 1812.5 | 3112.5 | 65 | 1812.47 | 3113.98 | 64.11 | 1.72 |
| 16 | 3112.5 | 3193.72 | 20 | 3112.14 | 3194.41 | 21.91 | 2.06 |
| 17 | 3312.5 | 3112.5 | 65 | 3311.81 | 3112.41 | 65.72 | 1.00 |
| 18 | 112.5 | 4550.93 | 20 | 111.49 | 4550.05 | 19.52 | 1.43 |
| 19 | 312.5 | 4412.5 | 65 | 311.16 | 4410.83 | 63.32 | 2.72 |
| **Avg** | / | / | / | / | / | / | **1.835** |

10

After accurately placing the orthogonal elements, the assembly process enters its second step, utilizing random rotation nodes to connect the elements sequentially (Fig 8). During this stage, two random parameters are considered: the rotation angle and the length of the node between itself and the end of each pipe. The rotation parameter can be predetermined with the help of overlaid images provided by AR technology, adding convenience to the final assembly process. After confirming the overlaid images, an additional deviation check is also performed using the motion capture system to evaluate the deviation with ground truth. The angle deviations are less than 0.2 degree, with an average value of 0.116 degree (Table 2).
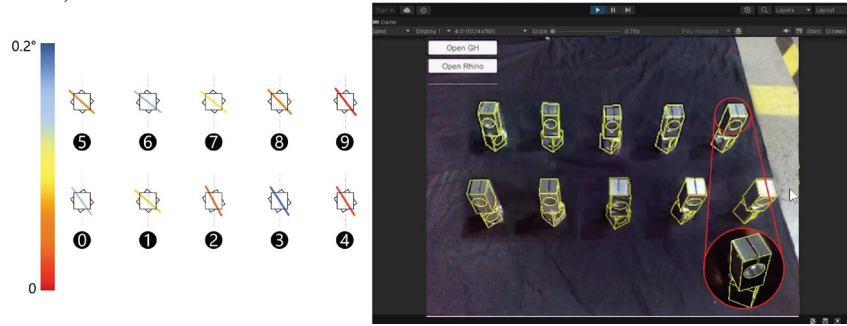


**Fig. 8** Angle deviation of each random node (Left), deviation check of overlaid image (Right)

**Table 2** Angle deviation of the random nodes

| Node Index | Ideal Angle /° | Real Angle /° | Error /° |
| --- | --- | --- | --- |
| 0 | 144.46 | 144.57 | 0.11 |
| 1 | 130.23 | 130.19 | 0.04 |
| 2 | 151.55 | 151.42 | 0.13 |
| 3 | 146.98 | 147.17 | 0.19 |
| 4 | 147.79 | 147.63 | 0.16 |
| 5 | 131.97 | 131.89 | 0.08 |
| 6 | 128.79 | 128.91 | 0.12 |
| 7 | 129.24 | 129.22 | 0.02 |
| 8 | 136.05 | 135.94 | 0.11 |
| 9 | 142.38 | 142.18 | 0.20 |
| **Avg** | / | / | **0.116** |

In the final step of the pavilion assembly (Fig 9), the process is simplified similar to the first step. The only factor to consider is the distance between the node and the end of each pipe, allowing the pavilion to be easily built with the assistance of AR technology theoretically. Due to the pipes being randomly scattered

in the air, and the lack of consideration and calculation of structural stability and deformation beforehand, the pipes are unable to remain stable in their ideal positions, even with the support of some extra vertical pipes. However, the pipes still maintain their correct orientation, experiencing a positional deviation of less than 22mm and 0.90 degrees, with average values of 10.31mm and 0.49 degrees (Table 3).
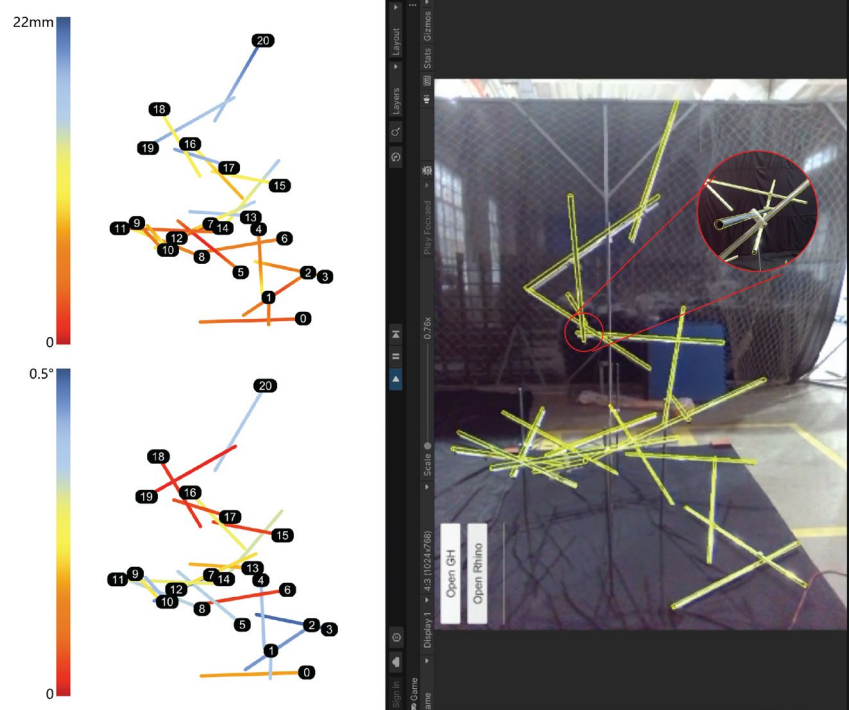


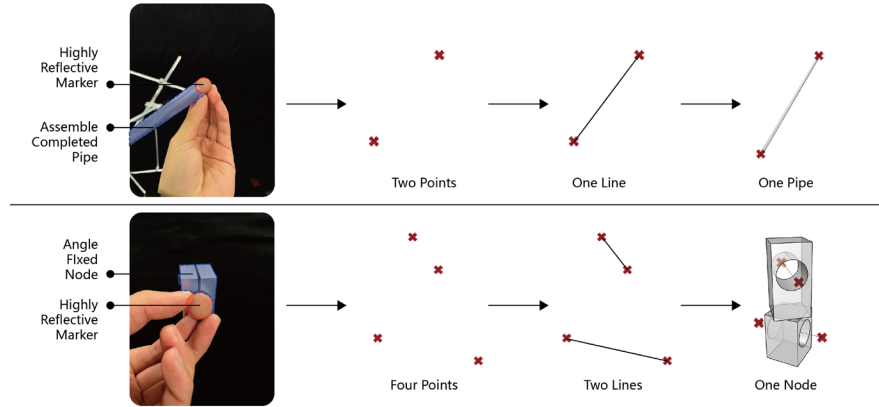**Fig. 9** Position and angle deviation of each pipe (Left), deviation check of overlaid image (Right)

**Table 3** Position and Angle deviation of all pipes

| Pipe Index | X value Of Ideal Middle Point /mm | Y value Of Ideal Middle Point /mm | Z value Of Ideal Middle Point /mm | X value Of Real Middle Point /mm | Y value Of Real Middle Point /mm | Z value Of Real Middle Point /mm | Position Error /mm | Angle Error /° |
|---|---|---|---|---|---|---|---|---|
| 0 | 1339.3 | 179.77 | 3.62 | 1341.75 | 182.69 | 182.69 | 4.71 | 0.34 |
| 1 | 1410.39 | 434.67 | 44.89 | 1406.8 | 429.84 | 429.84 | 6.09 | 0.68 |
| 2 | 1359.58 | 765.08 | 121.63 | 1360.79 | 764.52 | 764.52 | 4.10 | 0.80 |

**Table 3** (continued)

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 3 | 1308.16 | 1200.83 | 269.63 | 1309.08 | 1200.71 | 1200.71 | 6.86 | 0.90 |
| 4 | 921.51 | 1298.94 | 307.52 | 922.22 | 1299.91 | 1299.91 | 7.74 | 0.73 |
| 5 | 296.85 | 1244.92 | 419.48 | 296.41 | 1243.11 | 1243.11 | 1.86 | 0.64 |
| 6 | 649.51 | 1379.48 | 427.08 | 645.86 | 1385.71 | 1385.71 | 7.37 | 0.12 |
| 7 | 80.19 | 1061.31 | 560.54 | 78.49 | 1061.47 | 1061.47 | 5.32 | 0.82 |
| 8 | 20.34 | 670.36 | 578.94 | 27.16 | 670.01 | 670.01 | 9.37 | 0.63 |
| 9 | 24.81 | 378.29 | 724.35 | 23.84 | 373.99 | 373.99 | 5.86 | 0.43 |
| 10 | 280.21 | 184.1 | 739.41 | 278.63 | 183.6 | 183.6 | 11.12 | 0.74 |
| 11 | 469.62 | 83.7 | 858.32 | 473.08 | 82.98 | 82.98 | 5.08 | 0.54 |
| 12 | 929.74 | 191.92 | 977.18 | 928.08 | 189.5 | 189.5 | 12.38 | 0.42 |
| 13 | 961.21 | 290.38 | 1024.56 | 965.25 | 290.57 | 290.57 | 17.43 | 0.36 |
| 14 | 1066.07 | 777.56 | 1168.04 | 1062.98 | 782.49 | 782.49 | 14.81 | 0.58 |
| 15 | 908.23 | 934.47 | 1296.95 | 909.7 | 934.29 | 934.29 | 12.79 | 0.13 |
| 16 | 511.77 | 1018.17 | 1351.69 | 515.81 | 1021.19 | 1021.19 | 11.02 | 0.53 |
| 17 | 444.34 | 866.86 | 1493.02 | 446.24 | 860.12 | 860.12 | 19.92 | 0.15 |
| 18 | 424.31 | 417.73 | 1695.77 | 430.85 | 413.58 | 413.58 | 12.45 | 0.11 |
| 19 | 544.85 | 442.4 | 1921.34 | 540.22 | 445.14 | 445.14 | 19.16 | 0.05 |
| 20 | 911.84 | 802.93 | 2348.31 | 913.27 | 798.43 | 798.43 | 21.08 | 0.67 |
| **Avg** | / | / | / | / | / | / | **10.31** | **0.49** |

## *4.2 Deviation Analysis*



**Fig. 10** Method used to get the 3D coordinate data of the constructed pavilion's pipes and nodes.

A method employing the motion capture system is utilized to obtain the 3D coordinate data of the constructed pavilion's pipes and nodes (Fig. 10). Then the position and angle deviation data of the designed and constructed pavilion's pipes and nodes are compared to test the accuracy of the AR-assisted assembly workflow.

In steps 1 and 2, the pipes exhibit a position deviation of less than 3.2mm, while the nodes demonstrate an angle deviation smaller than 0.2 degrees. Furthermore, the outlines of the digital models align perfectly with the models in the physical environment. The deviation check results underscore the high accuracy of integrating the motion capture system into the AR-assisted assembly workflow.

In step 3 of the assembly process, the position deviation is less than 22mm, with the angle deviation smaller than 0.9 degrees. And some pipes are lower than the digital models in the overlaid images. The larger deviation observed in this step is primarily attributed to the deformation of the cantilever parts within the structure.
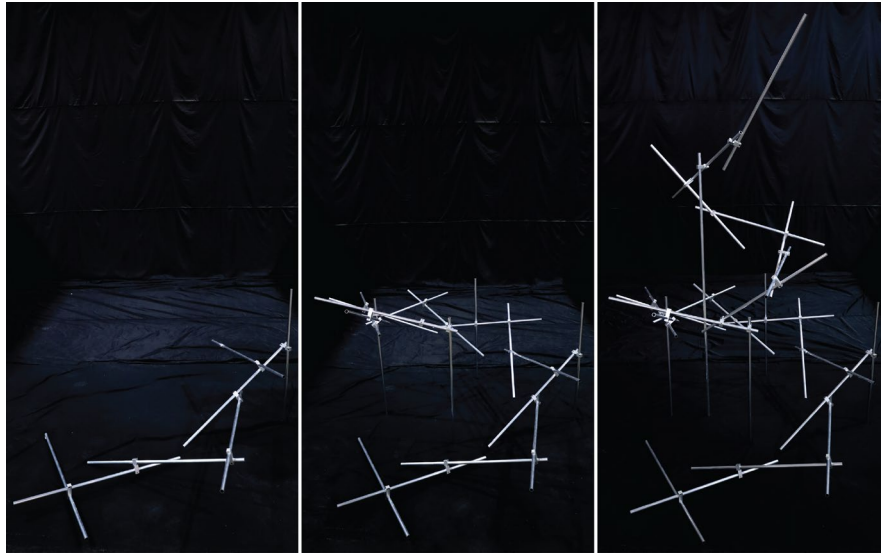
## 5 Conclusion and Discussion



**Fig. 11** Assembly progress 1/3 (Left), assembly progress 2/3 (Middle) and final assembly result (Right)

This research pioneers the integration of a 6-DoF motion capture system into the workflow of AR-assisted assembly. It develops both a software framework and a custom-designed hardware device to implement image overlaying, demonstrating the workflow's high accuracy. Moreover, it provides precise pose tracking data within a larger spatial scope beyond the typical confines of fiducial markers, thus streamlining the assembly process and saving time overall. Optimization of structural design could mitigate deviations even more. Additionally, improvement of the custom-designed device is necessary to free up the hand, thereby paving the way for more precise and intuitive assembly workflows in the future.

14

# 6 References

1. Yang S, Richard K, and Shan L. "Review and Analysis of Augmented Reality (AR) Literature for Digital Fabrication in Architecture." *Automation in Construction* 128 (August 2021): 103762.
2. Côté, S, Myriam B, Antoine G et al. "A Live Augmented Reality Tool for Facilitating Interpretation of 2D Construction Drawings." In *Augmented and Virtual Reality*, 8853:421–27. Lecture Notes in Computer Science. Cham: Springer International Publishing, 2014.
3. Gwyllim J, Cameron N, and Nick B. "Collaborative Fabrication In Mixed Reality." In *Data, Matter, Design*, 1st ed., 239–47. Routledge, 2020.
4. Gwyllim J, Cameron N, and Matthew B. "Making in Mixed Reality. Holographic Design, Fabrication, Assembly and Analysis of Woven Steel Structures," 88–97. Mexico City, Mexico, 2018.
5. Alexander HK, Arvin HX, Gwyllim J et al. "Augmented Reality for High Precision Fabrication of Glued Laminated Timber Beams," *Automation in Construction* 152 (August 2023): 104912
6. Sining W, Dandan L, and Lujie S. "Human-Cyber-Physical System for Post-Digital Design and Construction of Lightweight Timber Structures." *Automation in Construction* 154 (October 2023): 105033.
7. Vincent L, Francesc MN, and Pascal F. "EPnP: An Accurate O(n) Solution to the PnP Problem." *International Journal of Computer Vision* 81, no. 2 (February 2009): 155–66.
8. Xiao-Shan G, Xiao-Rong H, Jianliang T et al. "Complete Solution Classification for the Perspective-Three-Point Problem." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 25, no. 8 (August 2003): 930–43.
9. Garvin G and Kristof C. "Augmented Reality-Based Collaboration - ARgan, a Bamboo Art Installation Case Study," 313–22. Bangkok, Thailand, 2020.
10. Eric M, Hideaki U, and Fabien S. "Pose Estimation for Augmented Reality: A Hands-On Survey." *IEEE Transactions on Visualization and Computer Graphics* 22, no. 12 (December 1, 2016): 2633–51.
11. Timothy S and Jonas B. "Object-Based Visual-Inertial Tracking for Additive Fabrication." *IEEE Robotics and Automation Letters* 3, no. 3 (July 2018): 1370–77.
12. Schall, G, Daniel W, Gerhard R et al. "Global Pose Estimation Using Multi-Sensor Fusion for Outdoor Augmented Reality." In *2009 8th IEEE International Symposium on Mixed and Augmented Reality*, 153–62. Orlando, FL, USA: IEEE, 2009.
13. Francisco R, Rafael M, and Rafael M. "Speeded up Detection of Squared Fiducial Markers." *Image and Vision Computing* 76 (August 2018): 38–47.
14. Maximilian K, Acshi H, and Edwin O. "Flexible Layouts for Fiducial Tags." In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 1898–1903. Macau, China: IEEE, 2019.
15. Jack C, Keyu C, and Weiwei C. "Comparison of Marker-Based and Markerless AR: A Case Study of An Indoor Decoration System." In *Lean and Computing in Construction Congress - Volume 1: Proceedings of the Joint Conference on Computing in Construction*, 483–90. Heraklion, Crete, Greece: Heriot-Watt University, 2017.
16. Taihú P, Thomas F, Gastón C et al. "S-PTAM: Stereo Parallel Tracking and Mapping." *Robotics and Autonomous Systems* 93 (July 2017): 27–42.
17. Inês S, Ricardo S, Marcelo P et al. "Accuracy and Repeatability Tests on HoloLens 2 and HTC Vive." *Multimodal Technologies and Interaction* 5, no. 8 (August 23, 2021): 47.
18. Michail K, Brennan C, Sabrina C et al. "Fiducial Markers for Pose Estimation: Overview, Applications and Experimental Comparison of the ARTag, AprilTag, ArUco and STag Markers." *Journal of Intelligent & Robotic Systems* 101, no. 4 (April 2021): 71.