

Use Case

Our aim is to create a boxing game played using computer vision to detect human movements as a control scheme. Handheld controllers and VR headsets are often bulky and unwieldy, so we aim to incorporate computer vision to create a less restrictive and more immersive gaming experience, while providing exercise for users. Our target audience will be gamers and people who enjoy boxing/exercise.



3 components:

Computer Vision Pipeline (Mediapipe/Python/Unity scripts), Taiming
Video Game (Unity), Shithe
Haptic Feedback Device (Arduino/PCB wristband), Eric



Quantitative Design Requirements

Hardware

- Overall design must not be too large/heavy to not impede movement ($<100\text{g}$)
- Haptic feedback produced must be noticeable and have different, clearly distinguishable levels
- Battery must sustain arduino/motor setup for reasonable amount of time ($>20\text{ min}$)
- Wireless communication within reasonable distance (3 meters)

Computer Vision Pipeline

- Low latency($\leq 50\text{ ms}$ response time between user movement and movement being displayed on screen)
- High accuracy(80% accuracy in matching the real-world gestures with Unity avatars over 5 frames.)
- Good stability(smooth real-time gesture mapping with minimal fluctuation)

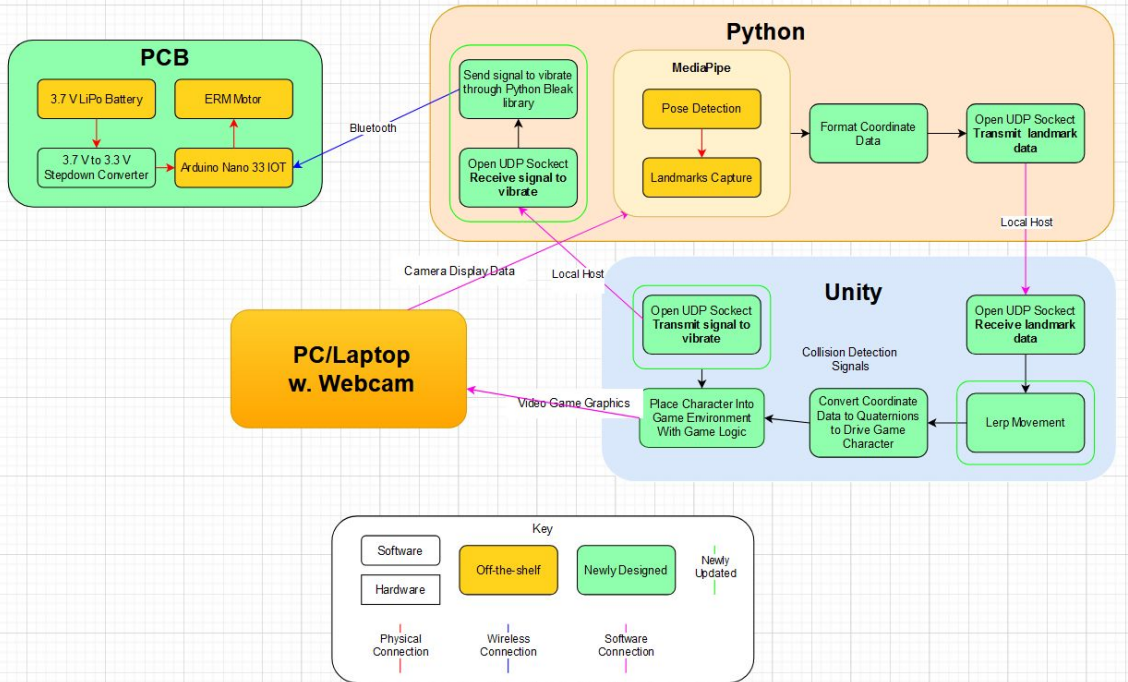
Video Game

- High-frame rate, $\geq 40\text{ FPS}$ on average, no sudden drops in frame rate
- Low Latency, $\leq 50\text{ ms}$ between human movement and movement displayed on screen
- Loading time of at most 45 seconds for any required load screens

Overall

- Low latency for entire system ($\leq 100\text{ ms}$ between human movement and haptic feedback)

Solution Approach



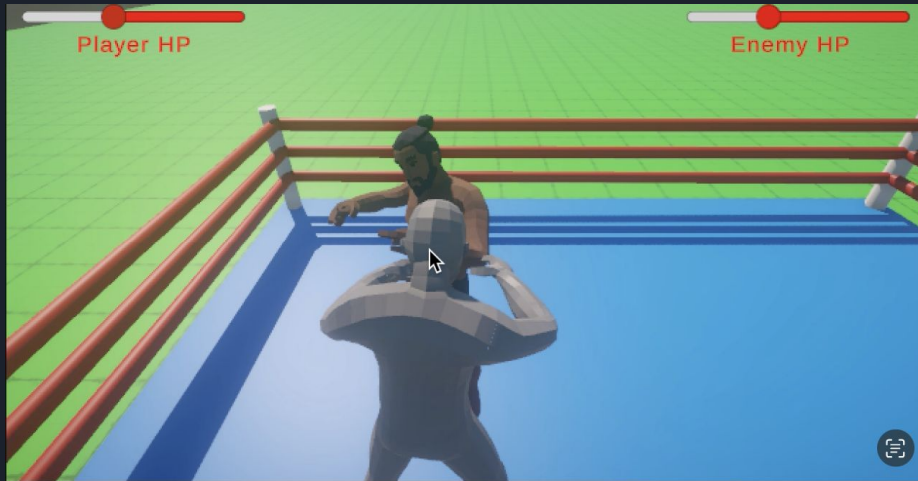
Important Update:

Quaternion Conversion replaced the naive 3D Coordinate Mapping to ensure smooth and accurate mapping free of distortions



Complete Solution


- Video game demo where user is able to move around in environment and throw punches, while receiving haptic feedback
- “Fight” against enemy AI, hitting each other until either player wins or loses, depending on whose health goes to zero first





Test, Verification, Validation - (Watch)

- **Weight Test** - Weigh watch to determine weight, with and without battery
 - Target: 100 g, on par with average watch
- **Vibration Test** - test vibration capabilities, determine if noticeable and if distinct levels can be felt
 - Target: 2 distinct levels of vibration, for player punch connecting and enemy punch connecting
- **Power Test** - Measure power consumption of watch in two ways, while idle and while vibrating to determine minimum and maximum power consumption
 - Target: 1.2 A, this would provide us with 20 minutes of continuous operation for a 400 mAh battery
- **Distance Test** - Measure maximum connectivity distance of bluetooth
 - Target: 3 meters, which is the amount of space we aim to use for our game
- **Ping Test** - Measure delay between bluetooth signal being sent by Python and after received by Arduino, using Python time library
 - Target: 30 ms



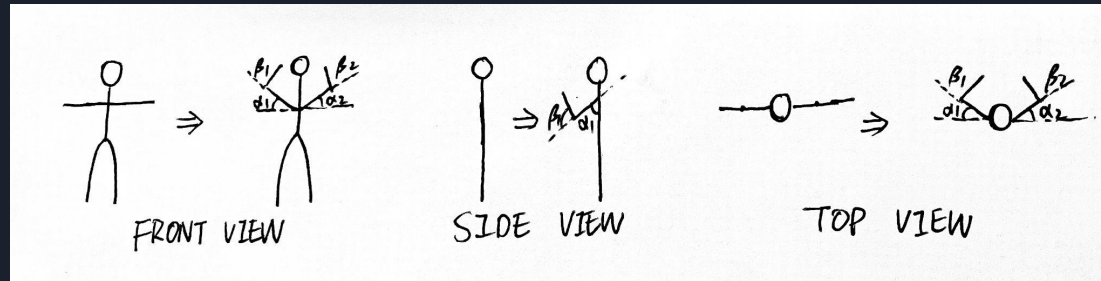
Test, Verification, Validation - (Watch) Results

Test	Target	Results
Weight	100 g	53 g
Vibration	2 distinct levels	2 distinct levels
Power	1.2 A	Min: 50 mA Max: 77 mA
Distance	3 m	20+ m
Ping	30 ms	22 ms

Test, Verification, Validation - (Computer Vision)

- **Accuracy Test**

- Target Accuracy is 80%, which means in five consecutive $n \cdot 20\text{ms}$ timestamps, 80 % of the ground truth-to-game avatar mappings should have less than 10 degrees difference in each joint angles.
- Ground truth data are obtained by measuring arm rotation angles from three views (front, side, top).
- Chosen over 3D inference due to hardware limitations (The latter requires a depth camera.)



The views above cover 4 D.O.Fs of the human arm and 2 perspectives.

Test, Verification, Validation - (Computer Vision)

- **Accuracy Test**

- Game Avatar data are obtained by projecting 3D transformations to x, z, and x plane (corresponding to each view) and extraction rotational values.
- The test is repeated five times for User 1 (height: 1.73m (5'8""); arm length: 1.74m) and User 2 (height: 1.62m (5'4""); arm length: 1.65m)

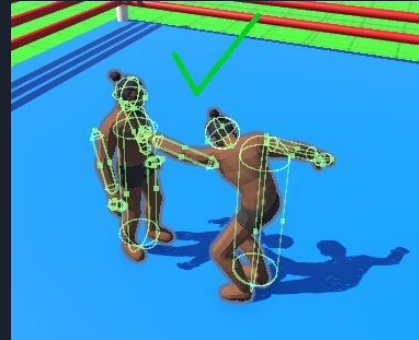
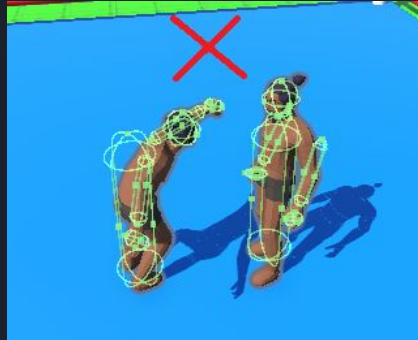
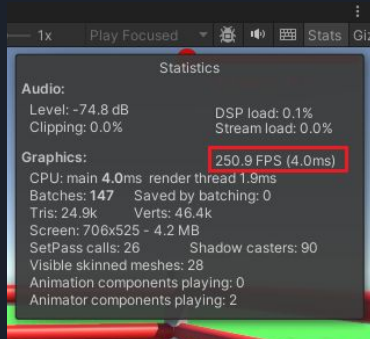
	Front View	Side View	Top View
$\alpha 1$	92%	86%	73%
$\alpha 2$	93%	N/A	76%
$\beta 1$	89%	83%	67%
$\beta 2$	91%	N/A	65%
Total	92.5%	84.5%	70.25%

Only **Top View**
underperforms due to
MediaPipe's limitation
on depth estimation

Mapping Accuracy Result (Target: 80%)

Test, Verification, Validation - (Video Game)

1. Passing = Higher than 40 FPS average with no sudden frame drops far below that target.
Tested FPS through the use of Unity's built in statistics.
Result: Passed with min average FPS of 50 on lower end laptop, Max 250 on desktop
2. Passing = Low Latency with ≤ 50 ms between human movement and movement in video game.
Result: Passed with 35 ms ping between computer vision and unity scripts
3. Passing = 100% Accurate hit detection in game with dynamic damage calculation depending on where is hit
Result: Passed with damage detection working properly



Design Trade-offs

- **Smoothing: Lerp (Linear Interpolation) VS. Calman Filter**

- **Lerp** transitions between two values (in our case, the target (next UDP received set of coordinates) and the current set of coordinates) over time by blending a fraction of the target value into the current value, ensuring smooth transitions between poses rather than snapping to the target positions immediately.

$$\text{Lerp}(a, b, t) = (1 - t)a + tb$$

- The **Kalman filter** is a mathematical algorithm that estimates the true state of a system by combining noisy measurements and predictions based on system dynamics.
- Lerp Pros: Simple to implement, Computationally Efficient;
- Cons: No noise handling, Constant smoothing rate (Introducing lag for fast-changing data)
- Calman Filter Pros: Noise Robustness, Predictive Power;
- Cons: Computationally intensive, Higher Latency

Calman Filtered:
Smoother at the cost
of a 600ms latency!



Lerped:
Fast,
but less smoother.





Design Trade-offs

- Lerp (Linear Interpolation) Internal Trade-off

```
for (int i = 0; i < pred3d.Count; i++)  
{  
    pred3dCurrent[i] = Vector3.Lerp(pred3dCurrent[i], pred3d[i], Time.deltaTime * MULTIPLIER);  
}
```

- `pred3dCurrent[i]` starts at the current position of the joint.
- `pred3d[i]` is the new position detected by the pose detection system for that joint.
- `Vector3.Lerp` moves `pred3dCurrent[i]` closer to `pred3d[i]` based on a smoothing factor (`Time.deltaTime * MULTIPLIER`).

High `Time.deltaTime` reduces latency by speeding up the lerping process, but it compromises smoothness by making transitions abrupt.

Large `MULTIPLIER` -> Smaller latency **BUT** worse smoothing performance

Small `MULTIPLIER` -> Better smoothing performance **BUT** larger latency

