

Distinguishing Between Head and Phone Gestures On a Smartphone With Front-Facing Camera and IMU

James Whiffing

jw204@bath.ac.uk

University of Bath - Department of Computer Science
Bath, England

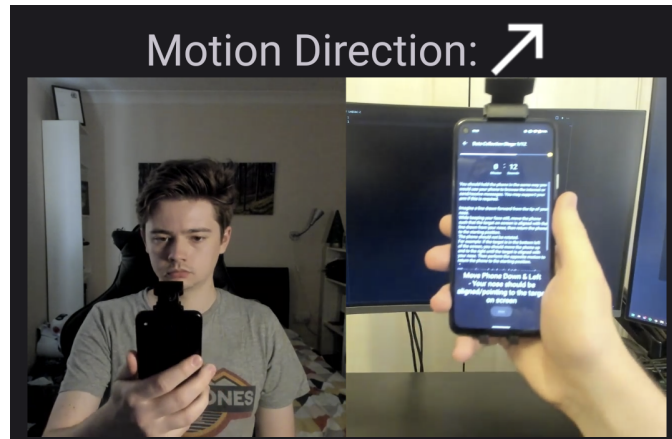


Figure 1: *TODO: Replace with Image showing diff between head vs phone moving resulting in similar photo (at least head pose)*

ABSTRACT

TODO

ACM Reference Format:

James Whiffing. 2022. Distinguishing Between Head and Phone Gestures On a Smartphone With Front-Facing Camera and IMU. In *Proceedings of (Conference acronym 'XX)*. ACM, New York, NY, USA, 5 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

There are attempts to introduce an additional modal of interaction with smart devices utilising the user's face. These exist on a spectrum with regards to interaction techniques: Using the face as a pointer, typically based on the movement/position of the user's nose; detecting gestures based on the movement/pose of the user's face; and a combination of the two.

A common issue that afflicts many of these systems/approaches is that they don't distinguish between the movement of the phone or the movement of the user's head. For example, the user moving

their head to the left, will be treated the same as the user moving the phone to the right, since from the front-facing camera's perspective it looks like the head is moving in the same way. This reduces the number of 'recognisable' gestures.

We look to explore whether such a system could distinguish between the user moving their head vs the phone being moved.

In order to develop such a system, a data driven approach was taken. As such a study was undertaken to collect the camera feed and IMU/Gyro of the smart device, an IMU within an earbud worn by the user, and 3D positioning of the user's head and the smart device via a motion capture stage. With the motion capture data being synced to the IMU/Gyro data and photos, a system could be trained to recognise several gestures and learn to distinguish between whether an observed gesture was due to the phone or the user's head moving.

2 LITERATURE REVIEW

In this section we will review existing literature to build an understanding of: the gestures we can expect to process, how they may be used and what they mean; Methods with which to obtain data pertaining to the pose of a head; and finally the means with which we can track movement and determine the gesture being performed.

2.1 Types of Head Gestures

Gestures can be classified into 5 classes[7]:

Dietic These are gestures that involve pointing, and mapping this to either a specific object or location within the interface.

Permission to make digital or hard copies of all or part 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 components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference acronym 'XX, N/A, N/A

© 2022 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

<https://doi.org/XXXXXXX.XXXXXXX>

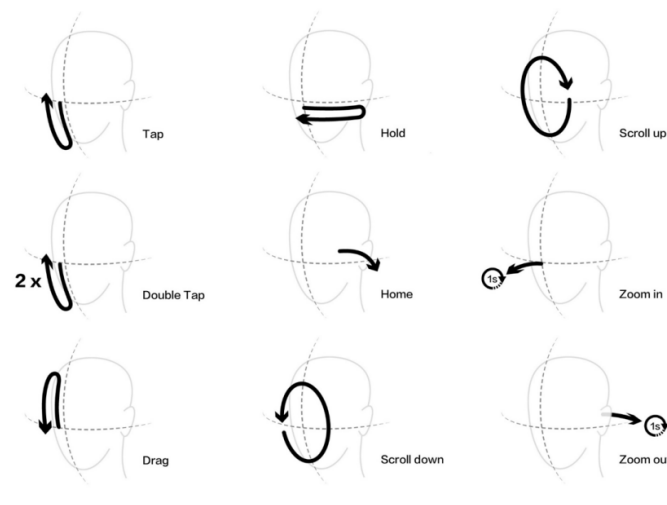


Figure 2: Proposed head gestures and their corresponding actions[12].

Manipulative A gesture which indicates intent to manipulate some object.

Semaphoric Gestures which map to a specific action or intent.

Gesticulative Typically accompany speech, but not required to transfer meaning.

Language A substitute to written or verbal language.

Of these gesture classes only Dietic, Manipulative, and Semaphoric can be applied to head gestures.

2.2 Head Tracking

2.2.1 Head Mounted IMU. One of the simpler ways to track the movement of a user's head is with an Inertial Measurement Unit (IMU), as this is a physical device that can be used to measure rotational and linear acceleration¹. An example of this can be seen in the work of Yan et al. who propose 9 gestures (Figure 2) and utilise the IMU embedded within the HoloLens[12].

2.2.2 HAAR Cascades. (Viola Jones Algorithm), find ref [8] Utilise HAAR-like features with an Adaboost classifier to detect 4 features: left eye, right eye, nose-tip, and mouth. Use Greyscale image, and use integral image (summed area table) to apply the features. Determine head pitch and yaw (not roll).

[9] Also utilise HAAR-like features applied to an integral image, however they utilise a histogram of pixel intensity (amount of variation in a given slice of the image) to identify the eyes. This is done once during calibration to then extract the eyes (left eye 30% along the line, right is 70%) and nose (if the eyes are d far apart, the nose is $d * 0.45$ below the midpoint between the eyes), assuming head is upright and user is facing camera.

2.2.3 Convolutional Neural Network (CNN).

¹Linear acceleration is typically less accurately tracked compared to angular acceleration (SOURCE!!!!?)

2.3 Head Gesture Recognition

2.3.1 Regression.

2.3.2 Recurrent Neural Network (RNN).

2.3.3 Markov Model. One way that a gesture can be described is via a set of possible states (e.g. head poses or movement) and a set of rules which describe how these states can change. However you may not be able to directly observe the Elmezain et al. utilise such a method through the use of a Hidden Markov Model to [3]

OLD

One of the simpler ways to track the movement of a user's head is with an Inertial Measurement Unit (IMU), as this is a physical device that can be used to measure rotational and linear acceleration².

The gestures they propose were derived from a study wherein they asked participants to suggest head movements that they believed corresponded to the action taken. These were then collated by manually into 80 gestures, which were then effectively voted upon by the participants for their respective actions. The gestures with the most votes for a given action were selected, with some minor adjustments to ensure there were no clashes between actions.

To extract the gestures from the HoloLens, the IMU output was segmented via detection of acceleration (20 degrees per second) and deceleration (4 degrees per second), not exceeding 2 seconds.

Feature extraction is performed with Dynamic Time Warping (DTW)[1], followed by a Support Vector Machine (SVM) classifier to classify the observed gesture into one of 9 the categories, or unintentional movement. With just the DTW they were able to achieve 90% accuracy, but with the SVM they were able to boost this to 97%.

To evaluate the head gestures, they compared them with existing hand gestures. They found that head gestures caused more fatigue and generally felt less natural, while being equivalent or better with regards to learnability.

Using a mobile device we won't have access to an IMU on the user's head, however we will be able to try and utilise the same set of gestures, and to use a similar approach to track the phone's movement, which could allow us to try and differentiate between the phone or head being moved.

An alternative approach is to try and extract the face using an RGB camera.

One way to do this was developed by Gorodnichy, which 'finds' the nose under the assumption that it should have the greatest intensity gradient since it should always be closest to the camera, and given it is convex in nature it should be the 'brightest' feature[4].

The tracked point isn't a specific point on the nose (e.g. the tip), but a point that can move across the surface of the nose, based on what is closest to the camera.

They go on to extend this work with the usage of the user's nose to control a pointer on their screen[5]. They extract just the nose since it meets their two requirements for a trackable feature:

- (1) It is always visible, presuming the user is facing within 180 degrees towards the camera.

²Linear acceleration is typically less accurately tracked compared to angular acceleration

- (2) Only one feature should be used to define the cursor, to reduce/eliminate potential jitter.

To click with the 'Nouse' the user blinks twice within short succession. Blinking is determined by reviewing the change in the sequence of 3 frames.

To ensure the system is realtime they use a reduced resolution, and could find that a resolution of 160x120 was robust enough to accurately track nose, and map the cursor to an accurate location on the screen.

They make claims about accuracy and enjoyment, but relevant data not provided. They only seemed to make statements suggesting Nouse was as good as, if not better than, typical mouse control. They claim mouse usage caused wrist ache, but movement of entire head doesn't present neck ache, which was reported in the IMU head tracking describe4d above[12].

Another nose controlled cursor is presented by Varona et al., however they use Haar cascades to extract the region containing the face, within which they use a similar technique as above[4] to extract points for the corners of the nose, or the nostrils[11].

To detect eyes, the system determines the user's skin colour by sampling the pixels within the detected face region. They then presume the eyes will be a different colour, and as such filter based on the extracted skin-tone. They then select the features closest to the nose, that are symmetrical.

A UI is provided with possible actions as buttons. The user moves the cursor to the action they wish to perform, then wink (with either eye) to select it. When they then fixate on part of the screen (move the cursor to a point and keep it stationary), the action will be performed.

Only evaluated for click recognition and accuracy of where the click was performed within a grid of points. However >80% accuracy even for users with no training time, just instructions.

Moving closer to a tool for mobile devices we have the work of Roig-Maimó et al., who use the front faced camera to scroll, using the head angle w/r/t the device as direction of scrolling.[10]

They extend upon the work of Varona et al.[11] the nose of the user and to correlate it's motion to a virtual cursor on the screen.

Selection/tapping is performed via tapping anywhere on the screen, however the tap will actually occur under the virtual cursor.

They evaluated the system by asking users to select elements of varying sizes, phone held in different orientations (portrait vs landscape), and with varying gain applied to the velocity of the cursor in response to head movement.

Elements below 88x88pt³ were found to be less successfully selected, this is primarily due to the low resolution used for the gesture tracking being unable to be mapped to a finer resolution on the device screen. Potentially increasing the resolution used for the tracking could permit finer accuracy.

They do not distinguish between the user moving the phone, or the user moving their head. It could be seen as a feature, either move the phone or head to scroll, but this would be interesting to try and distinguish, to potentially support additional actions/gestures.

The above systems describe the ability to identify where on a screen the user is looking, or at the very least intends to perform some action, through providing them with a virtual cursor they

can manipulate via moving their head, either through rotation or physically moving the head. We can look to extend upon this to understand where the user's attention may be focused.

Some smartphones now also include front-facing depth cameras, a technique for tracking a user's head, and detecting head/facial features is provided by Deepateep and Vichitvejpaisal, wherein they utilised the ARKit Framework for iOS[2].

Objective similar to works described above to control a virtual cursor, however instead of specifically using the nose, they are using the perceived pose of the user's entire head.

Cursor motion is tracked based on head pitch and yaw in the Y and X directions respectively. Requires user to directly face the camera for zero movement of the cursor.

Additionally they combine this control with facial gestures, which perform specific actions, or permit the beginning of specific actions, such as zoom, drag, and tapping. These utilised poses obtained from the eyebrows and mouth. Timings specific to each action, some gestures overloaded based on timing.

A depth camera affords greater accuracy for tracking (particularly for understanding the distance from the screen), and in case of iPhone there is consistency with specific hardware. However for the general smartphone population, specifically android, depth cameras aren't standard, and when present can have different hardware. For our project we will presume 3D depth cameras aren't available.

3 METHODOLOGY

This section details the process undertaken to develop the system which can meet the goal (outlined in the Introduction): distinguishing between head and phone based gestures on a smartphone.

In order to develop the aforementioned system we opted to take a data-driven approach. The benefit of taking a data-driven approach is that we can leverage Machine Learning, and train the system with exemplar data, rather than needing to manually determine the features and derive the algorithm needed to accomplish our goal.

3.1 Data Collection Study

For a data-driven approach to work we need to first determine what data we need to collect in order to train our system. We have identified the following types of data:

Images From The Front Facing Camera Given the majority of papers we reviewed *Citations again?*

utilise a camera to track the user's face, from which they can derive the gestures, we feel it necessary to do the same.

Smartphone Acceleration Data To understand whether the smartphone's PoV is changing we need to know how it is moving.

Head Acceleration Data If we can determine the movement of the smartphone via acceleration data, it is reasonable to see if we can also do the same with the user's head.

Actual Head and Phone Pose (Ground-Truth) In order to accurately train one models we propose below, we will need some Ground-Truth data.

³Apple Point, effectively 2 pixels on a retina display, so 88pt == 176px

With these data-types we can then build-up a dataset by recording each of the data-types during the performance of a series of gestures. Knowing the gesture associated to the recorded data will allow us to then train our system to recognise the gestures based on the data.

To create the required dataset we decided upon 11 gestures, each with 2-8 variations (effectively directions the gesture could be performed in), resulting in a total of 44 distinct motions to obtain samples of. A table of the gestures and variations can be found under Appendix A.

3.1.1 Apparatus and Techniques. Given the data types listed above, we decided to use the following tools:

- (1) **Pixel 4a** - An Android Smartphone with Bluetooth, a front-facing camera, and an IMU
- (2) **eSense Earable** - A Bluetooth Earbud with an IMU
- (3) **CAMERA Motion Capture Studio** - A Motion Capture (MoCap) studio found on campus within the University of Bath

The smartphone and earbud (when paired via bluetooth) will be able to provide the first three data-types defined above. While the Ground-Truth data can then be supplied by the MoCap studio.

To collect the data we developed an application to run on the smartphone. This was developed in Kotlin⁴ and the Android SDK. The application was designed to show participants a motion (a gesture and direction/variation) to perform. This is detailed in text, images, and a video. *Include figures to show this (spread across columns at top of page?)*

The participant would then be asked to perform this motion after pressing a record button. While recording the app would do the following:

- Capture images as frequently as possible from the front-facing camera, saving them as raw YUV bytes, with the UTC timestamp as the name.
- Record the smartphone IMU data (linear and angular accelerations), saving them to a csv with the UTC timestamp.
- Record the earbud IMU data (linear and angular accelerations), saving them to a csv with the UTC timestamp.

Once the participant has finished with the motion they could press the same button to stop the recording. Otherwise the recording will automatically terminate after 10 seconds, since the gestures shouldn't take more than a couple seconds to perform and the phone has limited RAM and storage with which to save data. To prevent accidentally stopping the recording too soon, say by accidentally double-tapping the screen, we disable the button for 2 seconds.

Once a motion has been recorded, the app shows the participant the next motion to perform. When the participant completes the final motion to perform, the app returns to the first motion. This repeats two times, such that each motion is captured 3 times. This is to collect variance in each motion for each participant.

In order to collect the Ground-Truth data, the study was performed within the MoCap studio. The participant was asked to wear hat that

had a motion-tracker attached, such that the tracker was placed around the middle of the back of their head. An exact position wasn't important as we only needed to determine the relative movement of their head, rather than the exact position. The smartphone was then tracked via a motion tracker attached to a 3D-printed mount, such that the tracker would not affect the participant's grip on the phone, or interfere with the images captured from the front-facing camera. *Include figure to show this*

Each tracker was composed of 5 points. 3 were positioned such that they formed a right-angle triangle, allowing the orientation of the tracker to be derived. The other 2 points were there to improve tracking accuracy, and help make the trackers unique and distinguishable. The MoCap system would track each of the 10 points at 60 fps and export the data as an fbx file.

In order to later synchronise the data collected we required the user to shake the smartphone, with a force of at least 2G, prior to beginning each round of 44 motions. The app would record the shake magnitude (in the X, Y, and Z axis) and the UTC timestamp of when it happened.

The full study protocol can be found under Appendix C

3.1.2 Study Results. Our study was run with 8 participants.

Unfortunately due to an issue with the application, the earbud IMU data was not recorded, despite the earbud being on and paired with the phone. This was not caught until after the study was completed.

Due to a late start, and overrunning into the next participant's slot, participant 0 was unable to complete their 3rd round of motions.

Some participants didn't initially stop recording upon completing a motion, as such their initial motions have superfluous frames that don't contain data relevant to the motion they're recorded for.

Stats from the data, including figs and tables? Range of motion per gesture. Time taken by gesture and participant, Average sample rate for IMU and Images, Average Number of frames containing face by participant and gesture

3.2 Data Post-Processing

Before being able to use the data for training, we needed to synchronise the data recorded from the smartphone, and the fbx data from the MoCap studio. To do this we derived the acceleration of the phone based on the MoCap data to find where it meets/exceeds the magnitude of the shake recorded by the app. From this we can determine the frame of the fbx data that corresponds to the recorded timestamp. We can determine the frame for any subsequent timestamp based on the known frame-rate of 60fps. To verify the data didn't drift we resync the data based on the other 2 recorded shakes, verifying that they're within 10 frames of the expected frame. In doing this verification, we did not come across any recording wherein subsequent shakes were not found to be at the expected frame.

Synchronisation and post-processing of the fbx data was performed with Blender and Python. Blender was used as permitted viewing the fbx data and verify derived location, roll, pitch, and yaw were correct. Also only way I was able to access the fbx data

⁴A programming language that runs on the JVM and is used to develop applications for Android.

programmatically. Synchronised data was exported to CSVs for each motion recording, containing a path to the image, the raw IMU data, and the derived MoCap data.

To increase the amount of effective data we have for training we shall do some fps scaling, such that we copy the data, but assuming we're only capturing images every X fps. We will find the photo closest to the new frame where it would have been captured, and average appropriate data (such as acceleration).

We will also slice the data in overlapping chunks (at least for the RNN model).

3.3 The Proposed Models

To achieve our goal we opted to train 2 models with which we could evaluate and compare performance. The first is a cascading classifier which predicts the motion being performed given a sequence of data. It first identifies if a face is present in an image via a CNN which returns a bounding box of the face[13]. If a bounding box is present the pixels within the box are passed to another CNN which extracts the position of 68 landmarks of the face[6]. These landmarks, the bounding box location, the average IMU data since the last image, and the last X frames are then passed into an RNN we trained to classify the gesture. If no bounding box or landmarks are found for a given image, zeros are provided **or previous data? . is input going to be padded with zeros for first frames / last frames, or require certain number of frames before attempting classification?**

The second model is 2 models which will be trained to predict the direction of movement in each of the 6 DoF for the head and phone (the head model will also take the landmarks and bounding box as input). It will output as a 2d one-hot encoded array, each row being the Degree of Freedom, the column being the direction (0 = stationary, -1 = negative, +1 = positive). The output of the 2 can then be fed into a HMM trained to predict the gesture performed based on the derived motion. (possibly an RNN if easier)

3.3.1 Training.

3.3.2 Model Evaluation.

3.4 Model Deployment

4 SYSTEM EVALUATION

4.1 Results

4.2 Discussion

5 LIMITATIONS AND FURTHER WORK

Collect more data (to improve accuracy)

Collect the earbud data (if model unsuccessful)

Depending on the model limited by only recognising gestures, not pointing

6 CONCLUSION

REFERENCES

- [1] Donald J Berndt and James Clifford. 1994. Using dynamic time warping to find patterns in time series.. In *KDD workshop*, Vol. 10. Seattle, WA, USA, 359–370.
- [2] Chatson Deepateep and Pongsagon Vichitvejpaisal. 2020. Facial Movement Interface for Mobile Devices Using Depth-sensing Camera. In *2020 12th International Conference on Knowledge and Smart Technology (KST)*. IEEE, 115–120.
- [3] Mahmoud Elmezain, Ayoub Al-Hamadi, Jorg Appenrodt, and Bernd Michaelis. 2008. A hidden markov model-based continuous gesture recognition system for hand motion trajectory. In *2008 19th international conference on pattern recognition*. IEEE, 1–4.
- [4] Dmitry O Gorodnichy. 2002. On importance of nose for face tracking. In *Proceedings of Fifth IEEE International Conference on Automatic Face Gesture Recognition*. IEEE, 188–193.
- [5] Dmitry O Gorodnichy and Gerhard Roth. 2004. Nouse 'use your nose as a mouse' perceptual vision technology for hands-free games and interfaces. *Image and Vision Computing* 22, 12 (2004), 931–942.
- [6] Yin Guobing. 2021. Head Pose Estimation By TensorFlow and OpenCV. <https://github.com/yinguobing/head-pose-estimation>
- [7] Maria Karam et al. 2005. A taxonomy of gestures in human computer interactions. (2005).
- [8] Jin Kim, Gyun Hyuk Lee, Jason J Jung, and Kwang Nam Choi. 2017. Real-time head pose estimation framework for mobile devices. *Mobile Networks and Applications* 22, 4 (2017), 634–641.
- [9] Euclides N Arcoverde Neto, Rafael M Barreto, Rafael M Duarte, Joao Paulo Magalhaes, Carlos ACM Bastos, Tsang Ing Ren, and George DC Cavalcanti. 2012. Real-time head pose estimation for mobile devices. In *International Conference on Intelligent Data Engineering and Automated Learning*. Springer, 467–474.
- [10] Maria Francesca Roig-Maimó, Javier Varona Gómez, and Cristina Manresa-Yee. 2015. Face Me! Head-tracker interface evaluation on mobile devices. In *Proceedings of the 33rd Annual ACM Conference Extended Abstracts on Human Factors in Computing Systems*. 1573–1578.
- [11] Javier Varona, Cristina Manresa-Yee, and Francisco J Perales. 2008. Hands-free vision-based interface for computer accessibility. *Journal of Network and Computer Applications* 31, 4 (2008), 357–374.
- [12] Yukang Yan, Chun Yu, Xin Yi, and Yuanchun Shi. 2018. Headgesture: hands-free input approach leveraging head movements for hmd devices. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 2, 4 (2018), 1–23.
- [13] Shiqi Yu and Yuntao Feng. 2022. YuNet. <https://github.com/ShiqiYu/libfacedetection>

A TABLE OF GESTURES AND VARIATIONS

B DATA COLLECTION APP CLASS DIAGRAM

C STUDY PROTOCOL DIAGRAM

D STUDY DATA ANALYSIS