

---

# TINKER@HOME 2021 TEAM DESCRIPTION PAPER

**Bingbai Li, Cuijie Xu, Dik Hin Leung, Guan Wang, Robin Ananda,  
Tengfei Zhang, Tiantian Yuan, Xinyu Han, Yanjun Jiang**  
Tinker Robotics Laboratory  
Tsinghua University  
xcj20@mails.tsinghua.edu.cnimonenext@gmail.com

## ABSTRACT

In this paper, we provide our robot system design details, including software and hardware, for RoboCup@HOME 2021 competition. We also show our improvements since last year. Our system design includes a framework of behavior modeling and human robot interface. Our main contributions are a novel logic model, and a framework to integrate various open source algorithms.

## 1 INTRODUCTION

The RoboCup@Home competition (Matamoros et al., 2020) aims at bringing robotic platforms to use in realistic domestic environments. In contrast to other leagues like soccer – which predefine and standardize the playing field – in Robocup@HOME, robots need to deal with different apartment layouts, continuously changing surroundings, unknown sites, unstructured public spaces, and interacting or cooperating with humans who are only very briefly – or not at all, in many situations – instructed about how to interact with the robot. The environment the competition provides varies from year to year. And thus it is typically required that the robots must be capable for multiple tasks, including navigation, human recognition and tracking, speech understanding and simple dialogues, object recognition and manipulation. How could a robot manage all these requirements, and in the meantime, deals with any possible exception or interference during the tasks is the main problem that teams participating in RoboCup@HOME have to deal with.

Future Robotics Club, or FuRoC team in short, is a student club of undergraduates from Tsinghua University, focusing on domestic robots, artificial intelligence and related fields. Robocup@HOME 2021 will be our 7th participation in the @home League of World RoboCup. Our team made a lot technical improvements since last year, and thus upgraded Tinker to a new stage.

## 2 ROBOCUP@HOME

The RoboCup@Home league (Matamoros et al., 2020) aims at bringing robots into domestic environments and support people in daily scenarios. The leagues are divided into three sub-leagues. Two of them are Standard Platform Leagues for which all competitors use the same robot, and the rest is Open Platform League that grants complete freedom to all competitors. The official leagues and their names are:

- RoboCup@Home Domestic Standard Platform League (DSPL)
- RoboCup@Home Social Standard Platform League (SSPL)
- RoboCup@Home Open Platform League (OPL)

Each league points out to a different aspect of service robotics, reason for which they target specific abilities. The robots used in the first two Standard Platform Leagues competition are fixed. They aim at benchmarking teams on their integrated software on a common hardware platform, and the standard platform provided to them is mature ones provided by commercial enterprises. Given the target population and anticipate results, the tasks are also nichetargeting. On the contrary, teams competing in the Open Platform League build the robot themselves. They are provided with the possibility to customize their robots and also the possibility to be portable within the limitation

defined by the rulebook. The competition consists of 2 Stages and the Finals. Each stage consists of a series of Tests that are being held in a daily life environment. The best teams from Stage I advance to Stage II which consists of more difficult tests. The competition ends with the Finals where only the two highest ranked teams of each league compete to select the winner. In Stage I, robots are tested for their basic functionality, such as finding, object detecting, speech recognition and navigating. In Stage II, they are challenged for their capability of integrating the functionality in Stage I to deal with more complex environment and tasks. In Final, the top - two teams will display a demonstration that arouses interest in the aspect of future research or application.

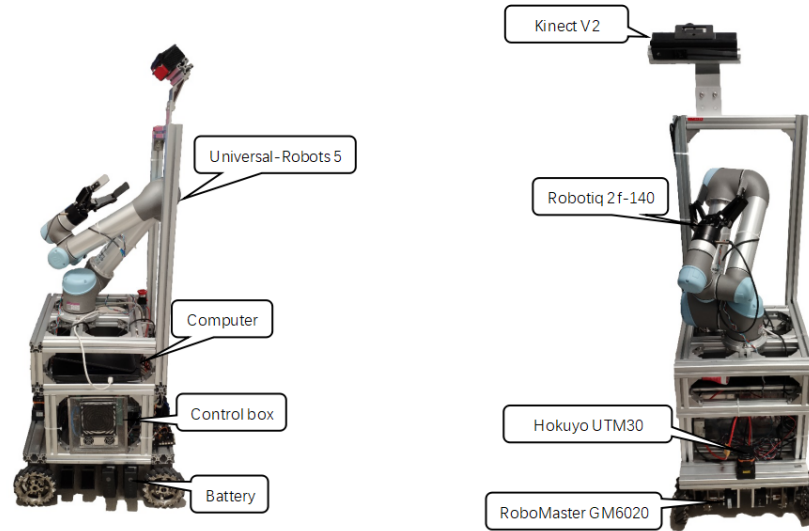


Figure 1: The robot Tinker in the front view (left) and Tinker in the side view (right). The overall height of Tinker is  $\approx 140$  cm. The overall weight of Tinker is  $\approx 70$  kg.

### 3 LAYER

Tinker is designed to be an autonomous robot mainly for domestic service. It is equipped with a agile mobile chassis using 4 mecanum wheels, a UR5 robotic arm, a powerful gripper (Robotiq-2F140), and various types of sensor. Depth cameras (Kinect v2) are used for imaging and recognizing environment, objects and people. Furthermore, an additional camera (realsense d435i) is assembled to the gripper, which assists the head camera while grasping. A single-line laser (UTM-30LX) is used in observing the environment and avoiding obstacles. From hardware interference to artificial intelligence decision making, Tinker has to deal with challenges at different levels. Thus, a multi-level, distributed structure based on the ROS platform is employed to meet such need. The hardware layer contains an embedded board driving motors and preprocessing odometry data. The hardware-communication layer is responsible for the overall control of the motors and the collection of data from the sensors. The output of the hardware layer is the ROS-compatible-sensor images, including camera images, point clouds and multiple other topics. The logic layer is in charge of providing basic functions of Tinker such as manipulation, navigation, person tracking, object recognition, speech understanding and synthesis etc. The decision layer serves as the collector of topics published by the logic layer and decision maker for the next integrated action to accomplish the task.

#### 3.1 THE LAYERS

Power management modules:

- Dji battery with DCDC transformer for the other equipments.

Actuators:

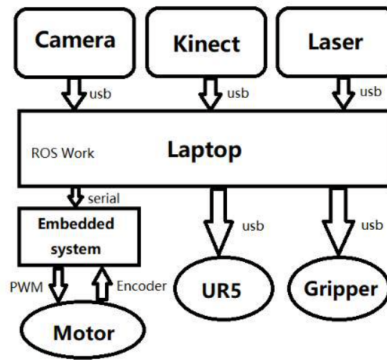


Figure 2: Structure of the hardware- communication layer

- DJI GM6020 motors for driving the chassis and the platform.
- Universal-Robots UR5 robot arm for accessing objects.
- Robotiq-2f140 mechanical gripper.

Sensors:

- Hokuyo UTM30 laser scanner for navigation.
- Kinect v2 depth camera for navigation and object detection.
- Realsense D435I camera for object recognition and detection.
- Encoder on motors for motor controlling.

### 3.2 THE HARDWARE-COMMUNICATION LAYER

The hardware-communication layer must be highly scalable to quickly install and uninstall different executors of the sensors. All control commands of the robot are sent to ROS nodes that are running on the laptop. The laptop also gathers the data collected by the sensors and give orders to the mechanical sections.

### 3.3 THE LOGIC LAYER

The most important functions of Tinker are implemented in this layer. The main components in this layer include:

- Navigation: Mapping, localization, route-planning and collision avoidance.
- Vision: Human recognition, object recognition and their tracking.
- Speech: Speech recognition and synthesis.
- Manipulation: Robot arm planing with feedback from vision and laser.

### 3.4 THE DECISION MAKING LAYER

Task planning is done in decision layer. The modules in the decision layer are run as state machines. They integrate the data uploaded from the lower layers to decide which state they are in, and then give different orders or make different responses. Each module deal with just one single task, and share the same information from the lower layers.

---

## 4 HARDWARE

### 4.1 CHASSIS

Owing to the Mecanum wheels, Tinker is capable of moving in all directions without having to turn around. For safety reason, the speed is limited to no more than 0.25m/s. Compared with the locomotion system in the past, the most important difference is the actuators. The new actuators make far less noise than the former ones because they don't have gearbox inside. Thus Tinker can operate more quietly, which makes it more user-friendly. The chassis consists of 4 separated Mecanum wheel systems, each of which is made up of a Mecanum wheel attached with DJI M3058P19 motor. The laptop sends order to the embedded board to control the chassis as they move along trajectory planned by the ROS Manipulation logic layer. The chassis has a size of 800mm × 500mm × 200mm.

### 4.2 THE ARM AND HAND

The arm and hand is the most major part of the robot. We choose Universal Robotic 5 (UR5) as its arm for its length and power meets the need to hold most of the objects at home, and Robotiq-2f140 as its hand for it gives a reliable grasping. A Realsense attached at the end of arm is used to object relocalization.

### 4.3 HEAD & HAND CAMERA

We choose Kinectv2 as the head camera, providing the central processor with point source and depth images to navigate as well as detect targets, and Realsense D435i as the hand camera, assisting in object detection and recognition.

## 5 APPROACH

### 5.1 HUMAN TRACKING

For human tracking and following, we implemented the TLD (Track Learning Detection) algorithm (Kalal et al., 2012). TLD was originally proposed by Zdenek Kalal in 2010. Currently it has developed to be one of the frontier real time tracking algorithms. Apart from combining the traditional tracking and detecting algorithms, it is more robust when taking into consideration the distortion and partial occlusions. The algorithm consists of three modules as its name indicated. The Tracking Module estimates the moving direction of the object according to the difference between two adjacent frames. The Detection Module detect the object in each frame independently. The Learning Module integrates the results of the Tracking Module and Detection Module to correct the detection errors and update the features of the target object. The algorithm is applied to human tracking and following tasks. Before the robot starts tracking, the human partner that is to be followed will be asked to stand in front of Tinker so that it can record his/her features. Once the record is done, the robot will track and keep up with him. Tinker also uses the depth information to keep itself a safe distance away from the instructor. Moreover, Kinect camera is used to analyze human skeleton, make simple recognition and understand of body language.

### 5.2 FACE RECOGNITION

In order to support human-robot interaction, the robot is required to recognize different masters or guests in domestic service. We established a face recognition system with two steps: enrollment and recognition. During the enrollment section, a person will be asked to stand in front of the RGB camera. The face detector based on Haar feature from OpenCV is applied and the detected feature will be stored. For a single being, the system stores 3-5 pictures. We implemented the face recognition algorithm based on sparse representation. A redundant dictionary is trained offline with a set of pre-imported faces. The algorithm seeks the most sparse representation coefficient by solving a L1 optimization problem. The residual errors for different classes (persons) will tell who is the unknown person: if the residual error for a specific class, for example, person A, is smaller than a specified threshold, and the errors for other classes are larger than another specified threshold,

---

the newcoming person is identified as person A. More details of this face recognition pipeline can be found in (Wright et al., 2009).

### 5.3 OBJECT RECOGNITION

Tinker uses a two-phase approach to recognize objects and precisely manipulate them. In the first phase, a point cloud is built according to the features collected by the Kinect depth camera, and we use Fast Plane Extraction in Organized Point Clouds inside. In simple terms, it is to extract the object from the two-dimensional picture, use the ransac and least square method to fit its shape parameters, and use the known three-dimensional spatial position information to reproject it into the three-dimensional space. The real-time plane extraction in 3-dimension point clouds is crucial to many robotics applications. We present an innovative algorithm to reliably detect multiple planes in organized point clouds obtained from devices - such as Kinect sensors, in real time. By uniformly dividing such a point cloud into non-overlapping groups of points in the image space, we are able to construct a graph in which the nodes and edges represent a group of points and their neighborhood respectively. We then perform an hierarchical clustering on this graph to systematically merge nodes that belong to the same plane until the squared error of the plane fitting mean exceeds a threshold. Finally we refine the extracted planes using pixel-wise region growing. Our experiments demonstrate that the proposed algorithm can reliably detect all major planes in the scene at a frame rate of more than 15Hz (for point clouds generated by 640\*480 depth images), which is much faster than many other algorithms we know. For object classification, another image processing method is implemented. We use fast YOLO (Redmon & Farhadi, 2016) for general object type detection, which is a precise while light-weight neural network, designed for general object detection and classification. Then we implemented the Word Bag Model (Csurka et al., 2004) to pair the object image with those collected in the data base.

### 5.4 SLAM

SLAM is an important algorithms that enables a robot to navigate and explore in an unknown environment (Grisetti et al.). The Mapping task requires the robot to record, integrate and update the information it collected about the surroundings while the Localization task requires the robot to identify its own location, referring to the estimated environment. Using a laser range finders (LRFs), we adopted the SLAM package to produce 2D occupancy grid map of the space. The raw data from LRFs are collected as the input stream. Features, or landmarks, are also extracted from the environment. When Tinker moves around, these features are used to estimate its location. It is called the Laser-Scan-Matcher process. However, the estimation of this process is imprecise and the error accumulates. The GMapping process is adopted (Grisetti et al., 2005), assisting by the EKF (Extended Kalman Filter) to correct the estimated result and draw the map (Grisetti et al.).

### 5.5 NAVIGATION

Navigation is one of the basic capability that a domestic robot must acquire. Based on the generated map, the robot needs to plan the route from its current position to the target one. Considering both distance and collision avoidance, the A\* algorithm is adopted to find the route. Moreover, the robot must be able to handle unexpected obstacles when moving around. Thus, the navigation package is applied and modified for Tinker so that parameters in the move\_base package are tuned and the navigation task can be achieved functionally. However, the behavior and speed is far from satisfactory. We now extend a local package which subscribes the global plan from the origin and linearizes the curve. In this way, the whole process becomes much more fluent. To avoid small objects and non-cylinder-like objects like chairs and cups on the floor, the robot is equipped with depth cameras - including a kinect2 and - to build another local obstacle layer. Since the pointcloud tend to be noisy, we introduce a filter to this obstacle layer to achieve more stable navigation performance, and a social layer is added to classify Bayesian data. Once a person enters the sight of the camera, he will be identified something alive and tagged. Even if he leaves the sight of the camera, he will be marked in the clustering model formed by radar to provide better effect for obstacle avoidance.

---

## 5.6 SPEECH RECOGNITION & UNDERSTANDING

For 2021 competition, we implement speech interaction system based on Google TTS(Text-To-Speech) and STT(Speech-To-Text) Cloud API (Lamere et al., 2003). To overcome the delay of communication between the robot and cloud platform, we import a stream based on audio transform method for speech recognition. Moreover, Tinker caches a huge amount of audio response template locally to speed up the robot reaction. Apart from the Speech-to-text layer, the dialogue system also contains a simple keyword parser, which takes keywords in certain patterns to operate task switches. Once the software recognizes a sequence with one of the predefined patterns, the robot will interpret one’s intention and makes responses.

## 6 CONCLUSION

In this paper we introduced our team and the robot, gave a brief outlook of RoboCup@Home, described the layer of Tinker, and its hardware and software system. A major focus was set on the description of the algorithm approaches. We proposed an novel approach for detecting multiple planes in organized point clouds which we proved to be more efficient. We now focus on two topics: computer-human interaction through gestures and facial expressions and the visual SLAM loop closure detection based on neural network. For the former one, the majority of the researches and applications now put emphasis on how the robot can better understand human, while actually, communication is something bi-directional. Thus, how human can understand what the robot want to “express” is also crucial. Moreover, apart from words, gesture is a very important method to exchange ideas. For the latter one, compared with the traditional Bag of Words algorithm, the design of the neural network-based relocation algorithm is greatly simplified. The idea is more natural, and it is hoped to obtain better performance than the traditional relocation algorithm.

## ACKNOWLEDGEMENT

The authors of this paper would like to thank previous team members of Tinker@Home2014, Tinker@Home2015, Tinker@Home2017, Tinker@Home2019, for their help and support through out building the robot and writing this manuscript. Particular thanks to old members Jingsong Peng, Jiacheng Guo and Yilin Zhu. The authors would like to thank teachers Yanxiong E, Xiaoliang Zhu in Youth League Committee of Tsinghua University.

## REFERENCES

- G. Csurka, C. Dance, L. Fan, J. Willamowski, and C. Bray. Visual categorization with bags of key-points. *Workshop on statistical learning in computer vision, ECCV*, 1, 2004.
- Giorgio Grisetti, Cyrill Stachniss, and Wolfram Burgard. Improved techniques for grid mapping with rao-blackwellized particle filters, *ieee transactions on robotics*, volume 23. 2007:34–46.
- Giorgio Grisetti, Cyrill Stachniss, and Wolfram Burgard. Improving grid-based slam with rao-blackwellized particle filters by adaptive proposals and selective resampling. *IEEE International Conference on Robotics and Automation (ICRA)*, 2005.
- Z. Kalal, K. Mikolajczyk, and J. Matas. Tracking-learning- detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(7):1409–1422, 2012.
- P. Lamere, P. Kwok, E. Gouvea, B. Raj, R. Singh, and W. Walker. The cmu sphinx-4 speech recognition system. *IEEE Intl. Conf. on Acoustics, Speech and Signal Processing (ICASSP 2003)*, 4, 2003.
- Mauricio Matamoros, Alexander Moriarty, Justin Hart, and Hiroyuki Okada. Robocup@home 2020: Rules and regulations (draft). [https://athome.robocup.org/wp-content/uploads/2020\\_rulebook.pdf](https://athome.robocup.org/wp-content/uploads/2020_rulebook.pdf), 2020.
- J. Redmon and A. Farhadi. Yolo9000: Better, faster, stronger. 2016.

---

J. Wright, A. Y. Yang, A. Ganesh, S. S. Sastry, and Y. Ma. Robust face recognition via sparse representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(2):210–227, 2009.

## A TEAM REPOSITORY

Our team repository can be found at <https://github.com/tinkerfuroc>. The repository may be of help to other teams by providing:

1. Implementation of all the algorithms and needed parameters described in the paper
2. Robot setup scripts and tools
3. Code for RoboCup@Home tasks

## B 3RD PARTY DEPENDENCIES

- ROS
- Tensorflow
- Face++
- Google Cloud API
- OpenCV
- PCL