# VirtualGrasp: Leveraging Experience of Interacting with Physical Objects to Facilitate Digital Object Retrieval

Yukang Yan, Chun Yu, Xiaojuan Ma, Xin Yi, Ke Sun, Yuanchun Shi

# Thor's Hammer

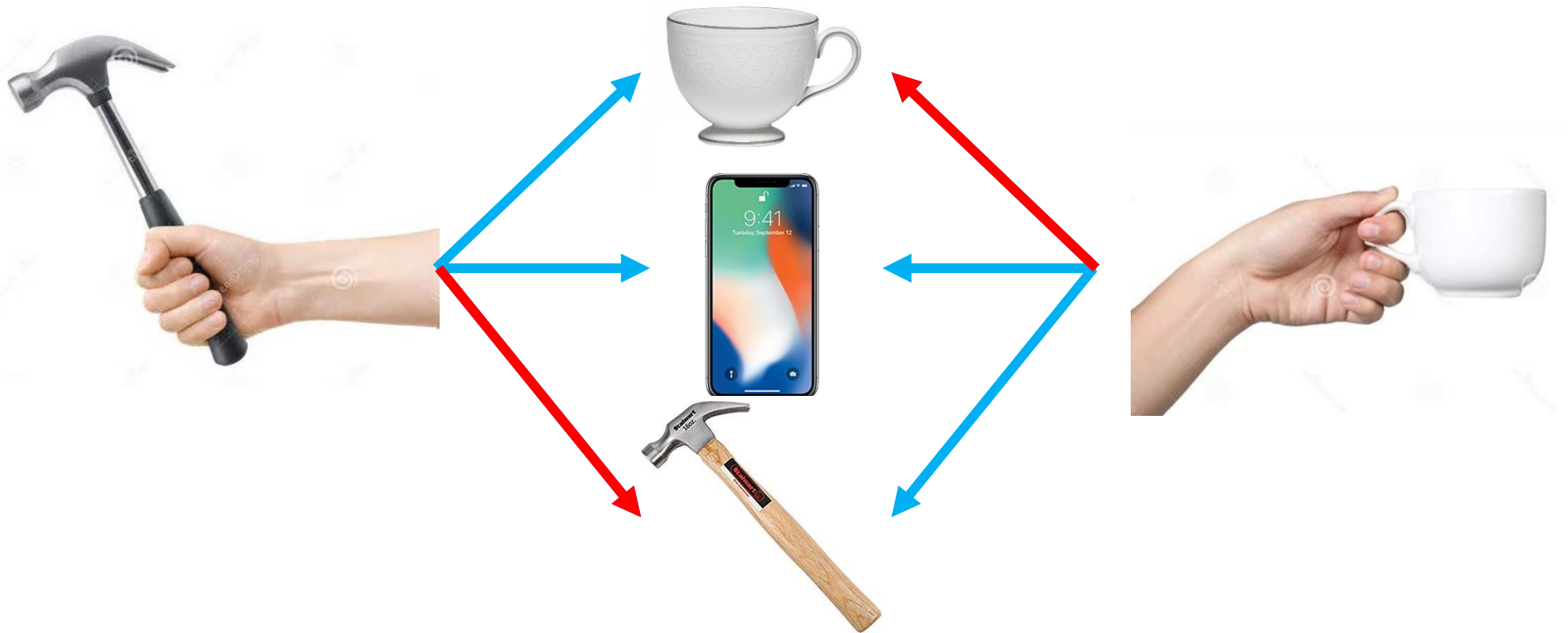To retrieve a **Virtual** object in VR, users perform the gesture of **Grasping** it in physical world.

Users perform **Swipe** gesture to select among the objects, and **Dwell** to confirm

To provide a set of **Self-Revealing Gestures** for object retrieval

To provide a set of **Self-Revealing Gestures** for object retrieval

1. Will users **consistently** perform the same grasping gesture for each object?

2. Can the grasping gestures of objects be **distinguished** by algorithms?

# Gesture Interaction

## Advantages

- Intuitive

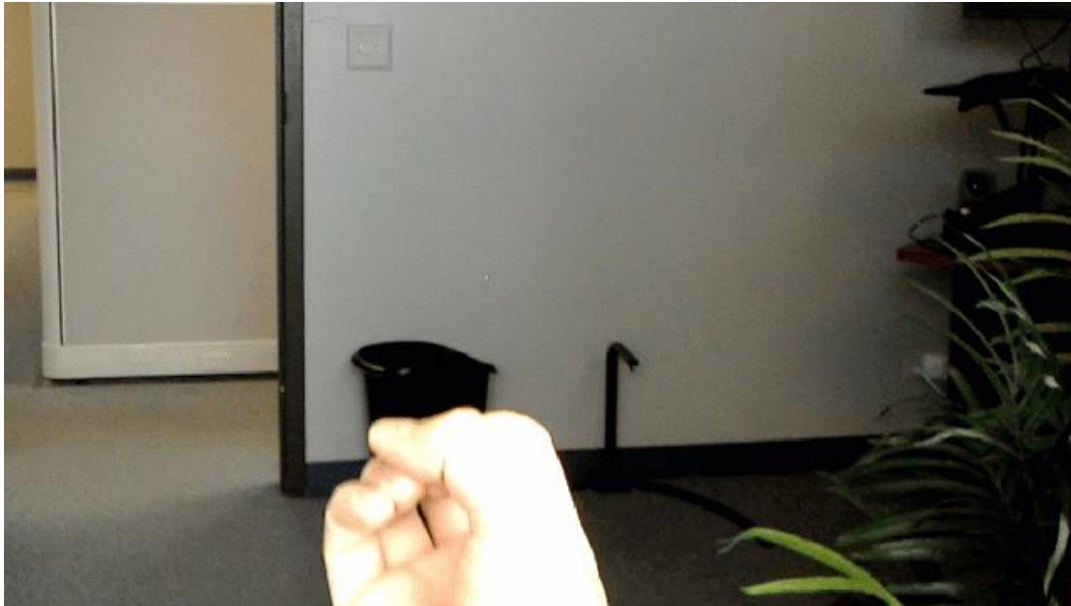- Direct Interaction

- Semantic Meaning

- Eyes-Free Interaction

## Disadvantages

- **Non Self-Revealing**

- Fatigue

- Fuzzy Input

# Gesture Interaction

## Hard to Discover



Bloom gesture to open the menu



Draw circle to open the camera

# Gesture Interaction

## Hard to Learn

# Gesture Interaction

## Hard to Remember

# Gesture Interaction

**Mappings** from **Targets** to **Gestures**

- Simple and easy to understand

- Consistent with acquired experience

- Consensus across different users

## Approaches for Mapping Problems

**Look-and-Feel** **Design of the Targets**



Yatani et al. CHI 08          Bragdon et al. CHI 11          Wagner et al. CHI 14



Esteves et al. UIST 15  Carter et al. CHI 16     Clarke et al. UIST 17        Esteves et al. UIST 17

## Gesture Interaction

**Mapping from Targets to Gestures**

- Simple and easy to understand

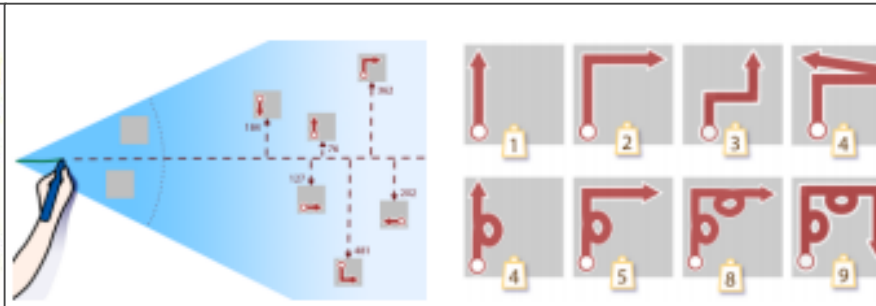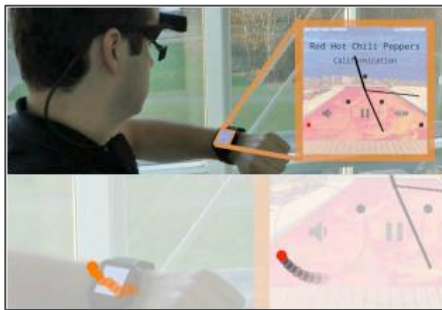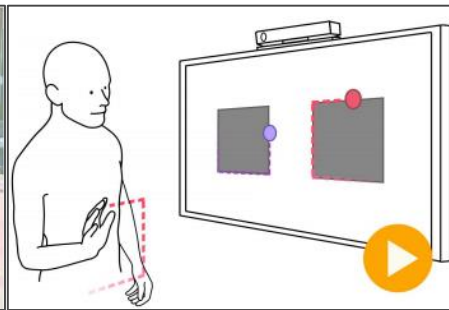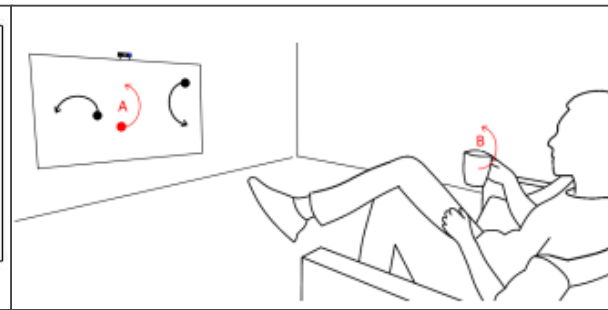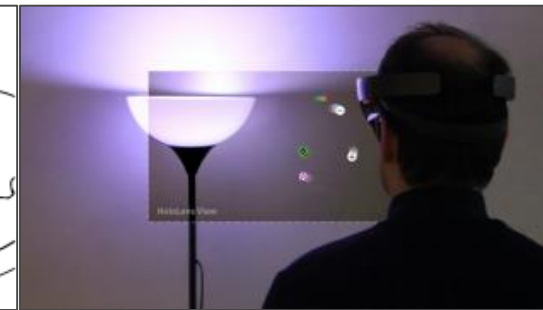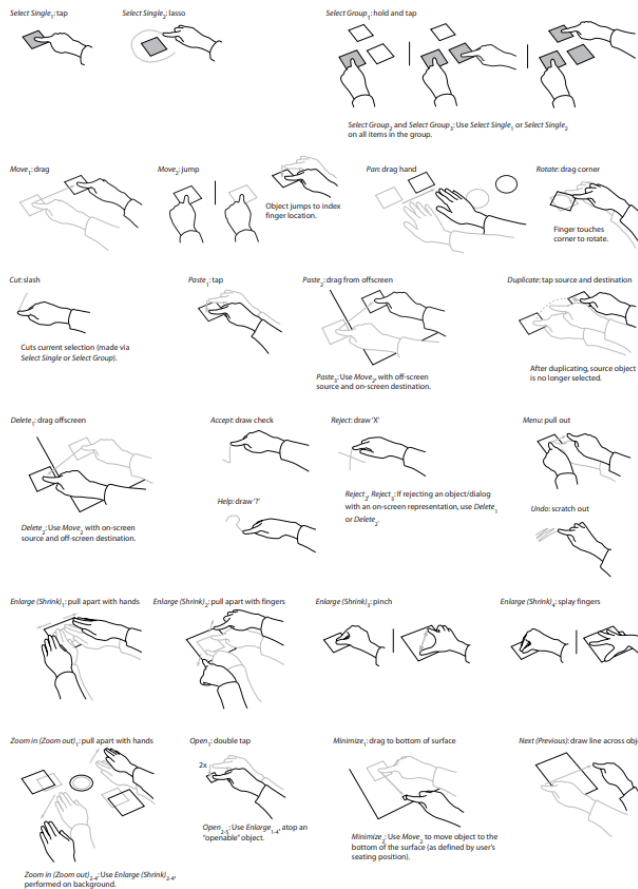- Consistent with acquired experience

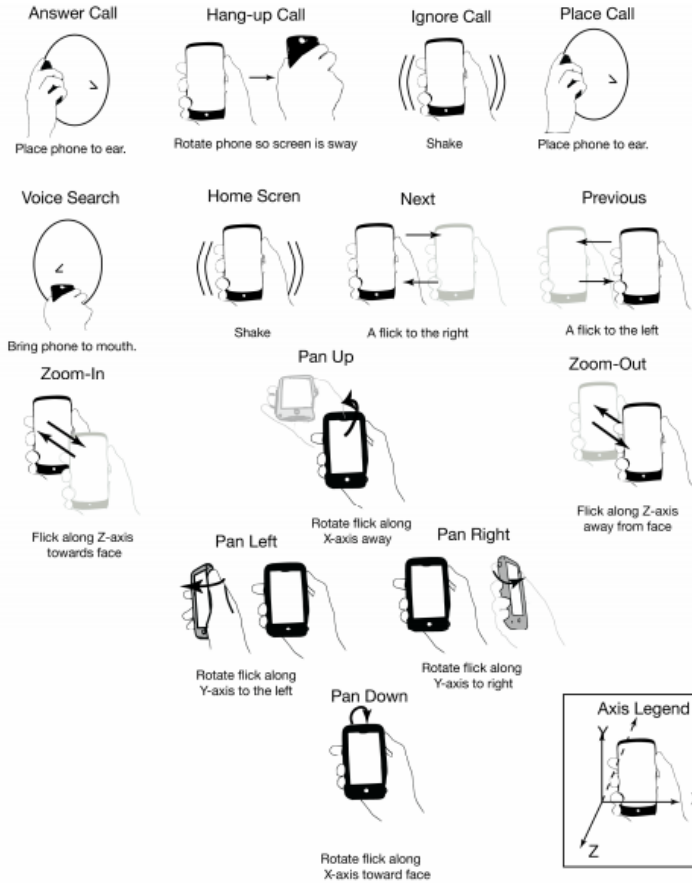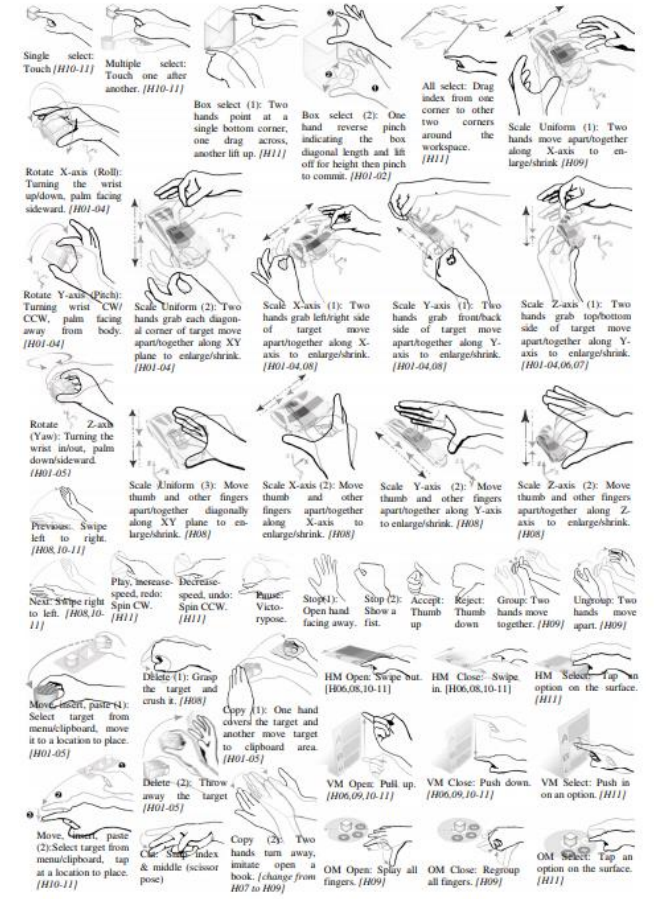- Consensus across different users

# Approaches for Mapping Problems

## **User Defined** Gestures (Participatory Design)



Wobbrock et al. CHI 09               Ruiz et al. CHI 11               Piumsomboon et al. INTERACT 13

## Gesture Interaction

**Mapping** from **Targets** to **Gestures**
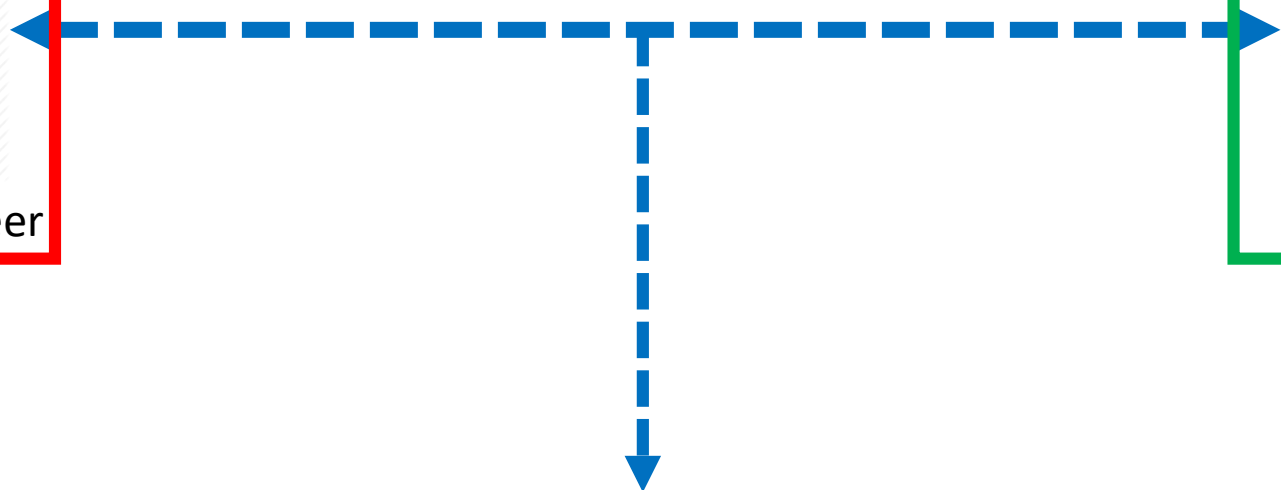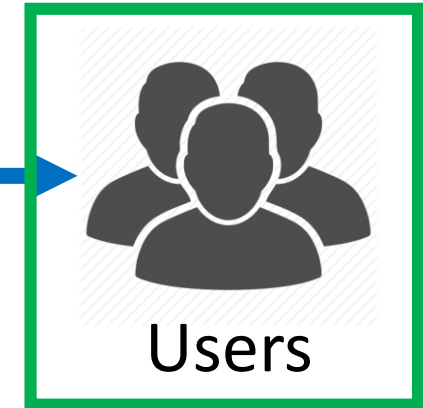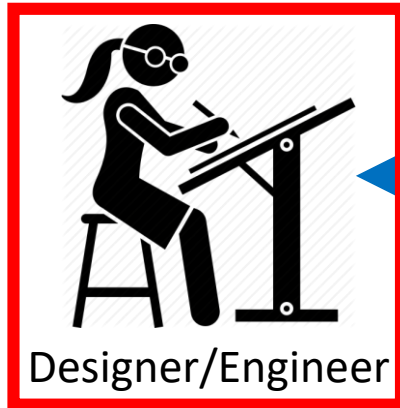
- Simple and easy to understand ➡ **One Simple Metaphor**

- Consistent with acquired experience

- Consensus across different users

| *Nature* | symbolic | Gesture visually depicts a symbol. |
| --- | --- | --- |
| | physical | Gesture acts physically on objects. |
| | metaphorical | Gesture indicates a metaphor. |
| | abstract | Gesture-referent mapping is arbitrary. |

# **Trade-Off** between Mapping and Recognition

**Designer/Engineer**

**Users**

**Current Gestures**
- ✓ Robust to Recognize
- ✓ No Conflicts within the Set
- ✗ Non Self-Revealing to Users

**Balanced Gestures**
- ✓ Self-Revealing to Users
- ✓ Consistent across Users
- ✓ Robust to Recognize
- ✓ Large Vocabulary

**User Defined Gestures**
- ✓ Intuitive to Users
- ✓ Consistent across Users
- ✗ No Concerns of Recognition

# Object Retrieval with **VirtualGrasp**

**1. Consistency**: Can users achieve high agreement on the mappings between the objects and their grasping gestures?

**2. Recognition:** Can grasping gestures of different objects be correctly distinguished by algorithms?

**3. Self-Revealing:** Can users discover the object-gesture mappings themselves? If not, can they learn and remember them easily?

- User Study: Gesture Elicitation -> **Consistency**

- Experiment: Gesture Recognition -> **Recognition**

- User Study: Object Retrieval -> **Self-Revealing**

- Summary

- Discussion

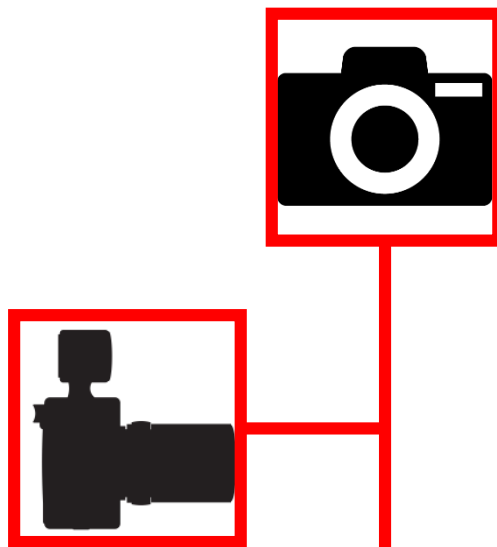- **User Study: Gesture Elicitation -> Consistency**

- Experiment: Gesture Recognition -> **Recognition**

- User Study: Object Retrieval -> **Self-Revealing**

- Summary

- Discussion

Names of **49** different objects were shown on the front screen.

Two cameras recorded the gestures from the front and the side view.

We recruited 20 participants (14M/6F) to perform the grasping gestures for each object.

**Object Set** (49 Objects)
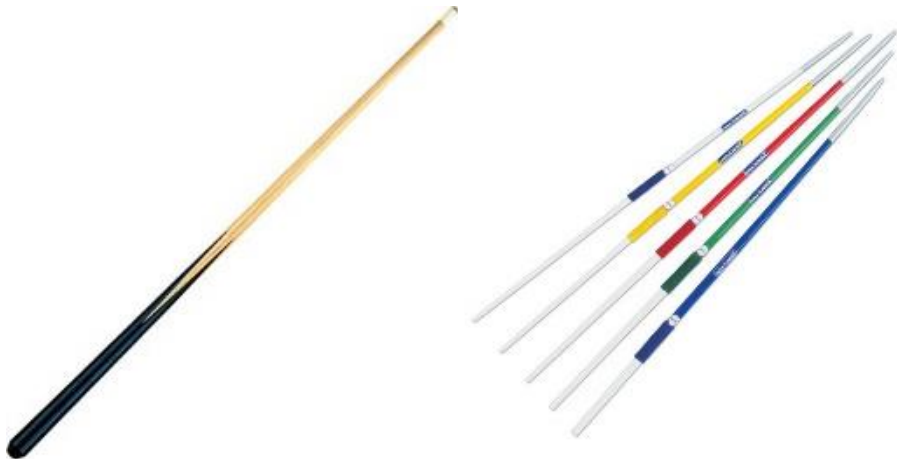
- **Consensus across Users**
  - $20 \times 49 = 980$ gestures, 140 gesture-object pairs.
  - 18/49 objects mapped to one unique gesture
  - 49/49 objects mapped to no more than five gestures
  - Agreement Score: AVG = 0.68, SD = 0.27
- **Key Properties of Objects**
  - Shapes: 41.3/49  Usages: 40.8/49  Sizes: 29.8/49

- **Breakdown and Distribution of the Gestures**
    - The taxonomy: "Single/Double Hands", "Hand Position", "Palm Orientation" and "Hand Shape".
    - One-to-one V.S. N-to-one mapping.
    - Infrequent gestures to be leveraged.

**Discussion**

- Half Open-Ended Elicitation Study
  - The power of the metaphor: high consistency across users.
  - The using experience of the objects are required (39/980).

# OUTLINE

- User Study: Gesture Elicitation -> **Consistency**

- Experiment: Gesture Recognition -> **Recognition**

- User Study: Object Retrieval -> **Self-Revealing**

- Summary

- Discussion

- **Participants**
  - **12** participants, with an average of 24.3 (SD = 1.5).  Four of them had experience of mid air gesture interaction. All were familiar with touchscreen gesture interaction.
- **Apparatus**
  - Perception Neuron, which was a MEMS (Micro-Electro-Mechanical System)  based tracking device, with a resolution of **0.02** degrees.
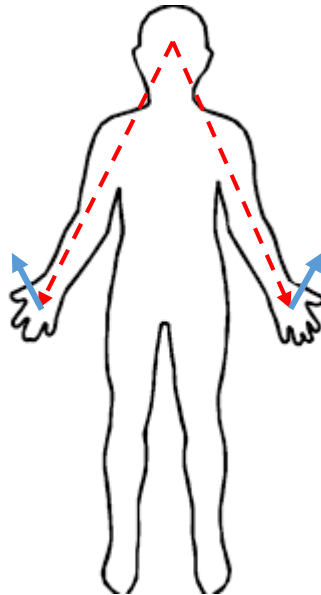


The sensors that participants put on

- **Data Collection**
  - The positions and orientations of the hand palms relative to the head.
  - The positions of the 14 joints relative to the hand palms.
  - 40 frames for each gesture that participants performed.
  - $2\ hands\ \times 16\ vectors\ \times 3\ values = 96\ values\ per\ frame$
  - $12\ participants\ \times 101\ gestures\ \times 2\ rounds\ \times 40\ frames = 96960\ frames$

- **Leave-Two-Out Validation**
  - Data of two participants as test set and the left as training set. ($C_{12}^2 = 66\ rounds$)
  - Top-N accuracy: N most possible objects contain the target. (Top-1, Top-3, Top-5)
  - Average accuracy:

|  | Top-1 | Top-3 | Top-5 |
|---|---|---|---|
| Mean | 70.96% | 89.65% | 95.05% |
| SD | 9.25% | 6.39% | 4.56% |



**Too small objects**

- **Leave-Two-Out Validation**
  - Data of two participants as test set and the left as training set. ($C_{12}^{2} = 66\ rounds$)
  - Top-N accuracy: N most possible objects contain the target. (Top-1, Top-3, Top-5)
  - Average accuracy:

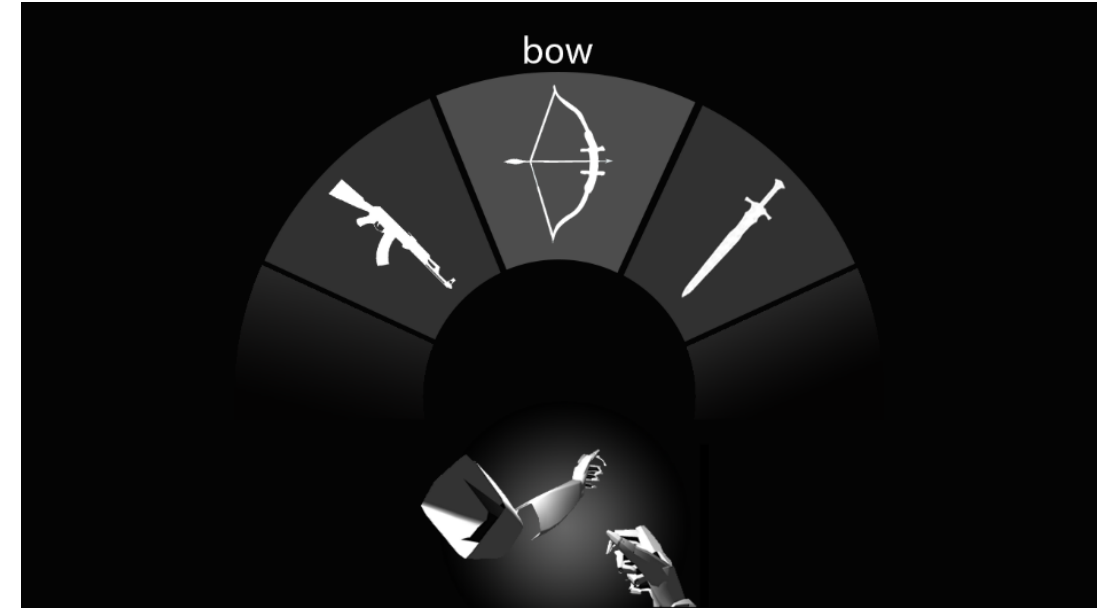|  | Top-1 | Top-3 | Top-5 |
|---|---|---|---|
| Mean | 70.96% | 89.65% | 95.05% |
| SD | 9.25% | 6.39% | 4.56% |

**Strong connection to usages**

- User Study: Gesture Elicitation -> **Consistency**

- Experiment: Gesture Recognition -> **Recognition**

- User Study: Object Retrieval -> **Self-Revealing**

- Summary

- Discussion

- **Participants**
  - 12 new participants, who never participated in STUDY1 or STUDY2.
- **Apparatus**
  - We showed the name of the target object on the top, visualized the current gesture of the participants, and showed the recognition result of top three possible objects in the center.



User Interface

- **Discovery** Session
  - ***Without learning*** the gesture-object mappings in the system, we asked participants to perform their own grasping gestures.
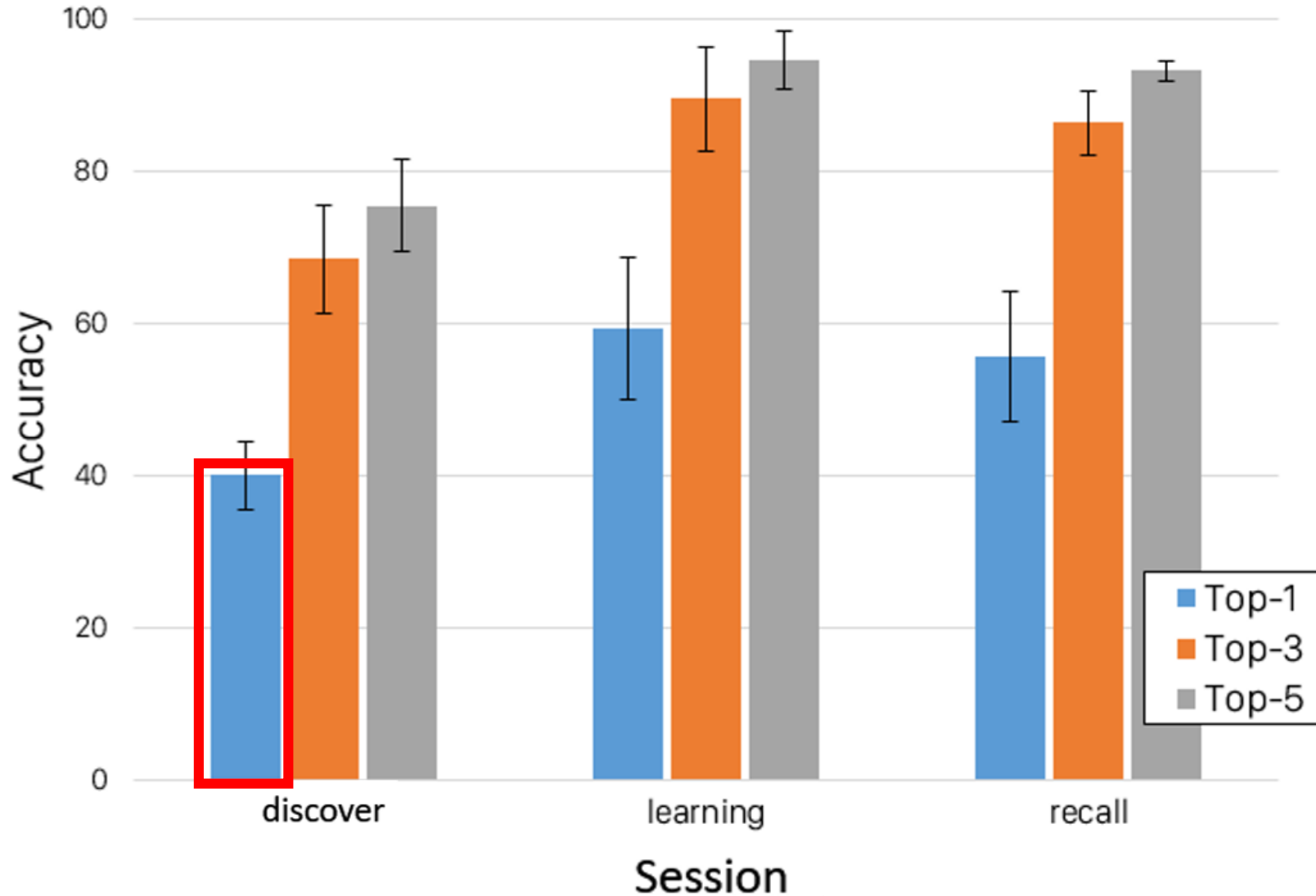- **Learning** Session
  - Before test, we let participants learn the standard gestures. They were free to ***practice*** the gestures until they confirm to be ready.
- **Recall** Session
  - ***A week later***, participants came back to lab and perform 49 object retrieval tasks again. During the week, they were not exposed to the standard gestures again.
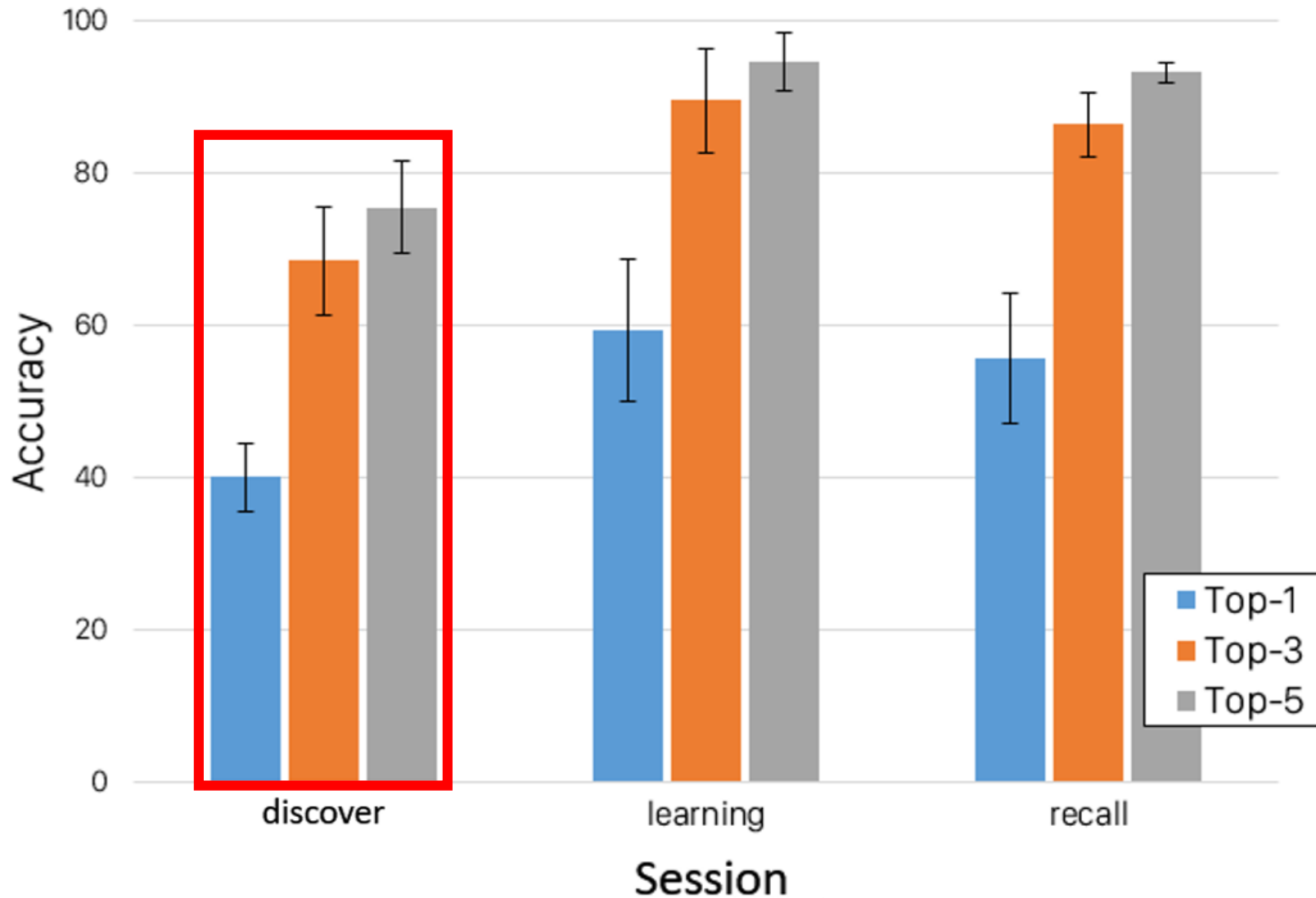
**40%** of the gestures were triggered without training

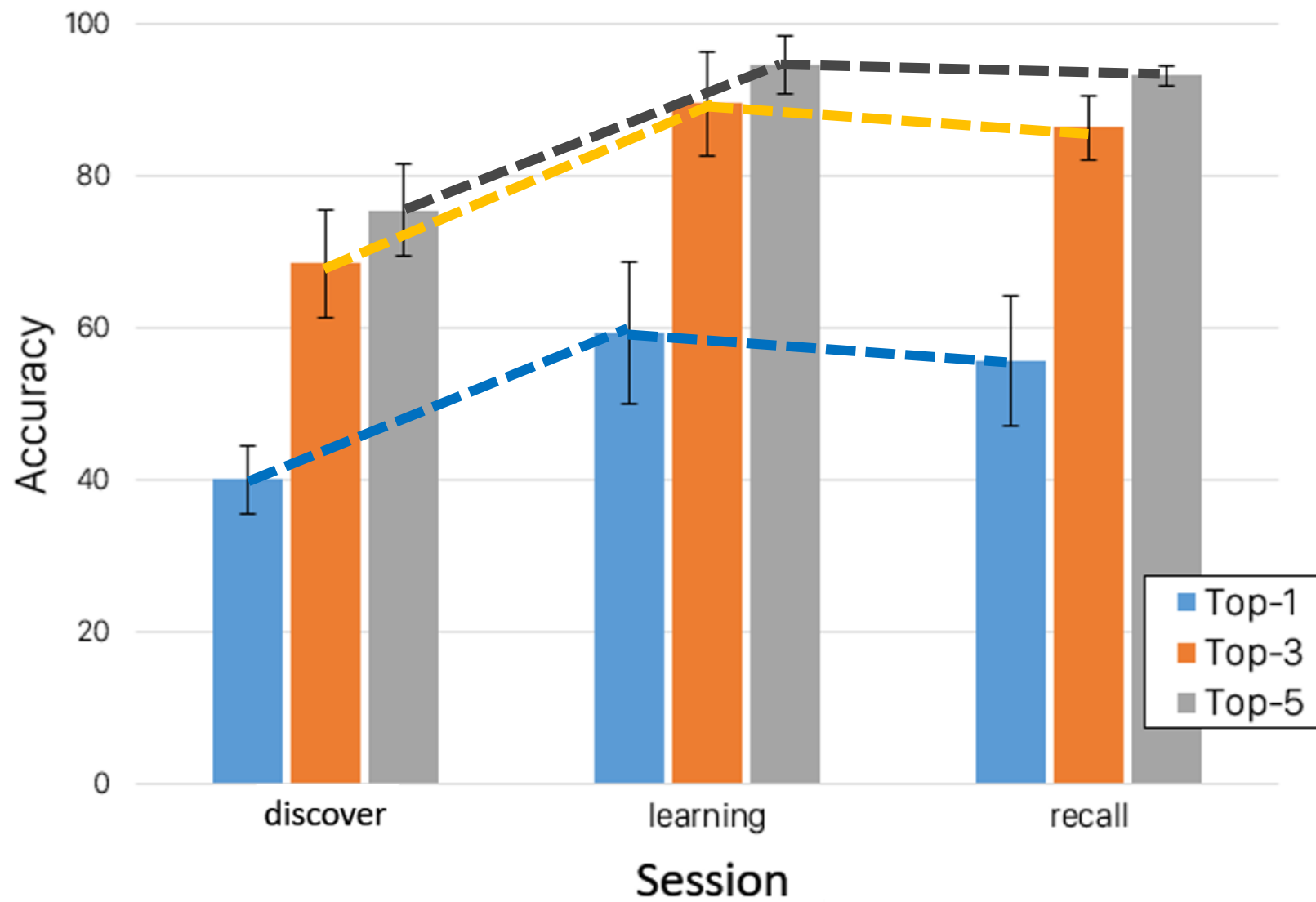**76%** of the objects were successfully retrieved

- **Discoverability**: Without any training, participants could discover **40%** of the exact mappings by themselves, and could directly use VirtualGrasp to retrieve **76%** of the objects with top five candidates.

- **Memorability**: A week after the learning session, participants could still recall the mappings well, and could successfully retrieve **93%** of the objects with top five candidates.

**Subjective Feedback**

- **The system is intelligent**
- *"Two different gestures came to me for grasping the camera and it was intelligent that the system correctly recognized the one I performed." [P4]*
- **The gestures make sense**
- *"I never used a grenade before, but I agreed with Gesture 3 which was grasping it over the shoulder to throw it." [P6]*
- **New tricks under the concept**
- *"For 'Stapler', I chose to perform the gesture of pressing it instead of holding it, because few other objects require pressing." [P8]*

| 5-Point Scale | Discoverability | Fatigue | Memorability | Fun |
|---|---|---|---|---|
| Mean | 4.2 | 4.4 | 4.5 | 4.4 |
| SD | 0.78 | 0.70 | 0.53 | 0.52 |

**High Consistency**

**Study 1
Gesture Elicitation**

an elicitation study ⟶

1. Consistency of Mapping
2. Object Properties
3. Gesture Taxonomy and Distribution

Gesture Set

**Good Accuracy**

**Study 2
Object-Gesture
Recognition**

an offline evaluation ⟶

1. Cross Validation
2. Effect of Scenarios
3. Consistency along Time
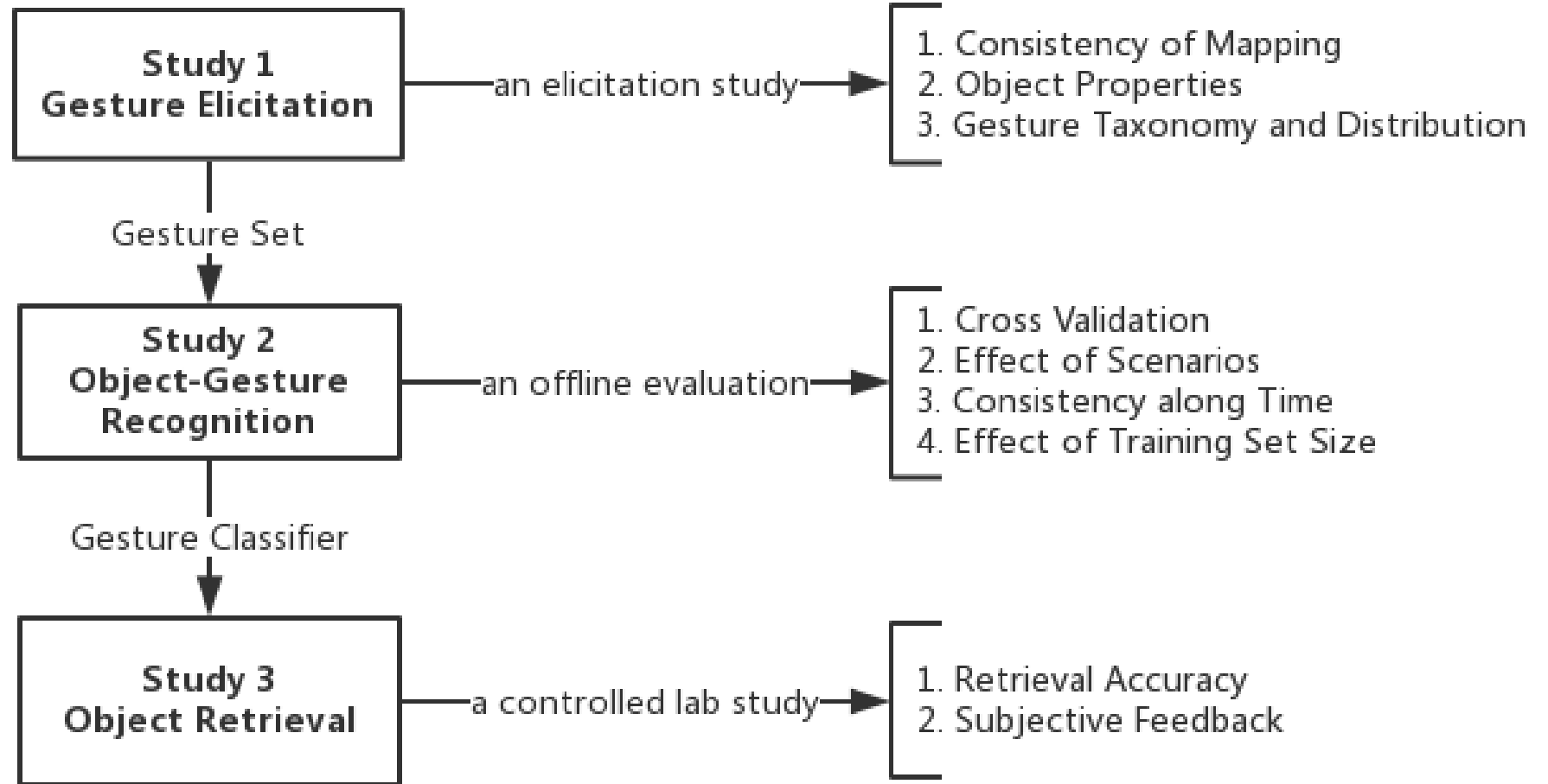4. Effect of Training Set Size

Gesture Classifier

**Little Effort**

**Study 3
Object Retrieval**

a controlled lab study ⟶

1. Retrieval Accuracy
2. Subjective Feedback

- **Object-Gesture Mappings**
    - Objects with different property values.
        - Not from objects of the same type.

- **Object-Gesture Mappings**
  - Objects with different property values.
    - Not from objects of the same type.
  - Grasping gestures reflect different properties of objects.
    - Difficult to distinguish grasping gestures of too small objects.

- **Object-Gesture Mappings**
  - Objects with different property values.
    - Not from objects of the same type.
  - Grasping gestures reflect different properties of objects.
    - Difficult to distinguish grasping gestures of too small objects.
- **Sensing Technique**
  - Hand gesture, hand position and hand orientation.
    - Vision-based sensing techniques.
  - Hand gesture.
    - Data gloves, EMG sensors, Vision-based.
  - Hand position and hand orientation.
    - VR controllers.

Thanks