# GazePointAR

A Context-Aware Multimodal Voice Assistant for Pronoun
Disambiguation in Wearable Augmented Reality

Jaewook Lee[1], Jun Wang[1], Elizabeth Brown[1], Liam Chu[1],
Sebastian S. Rodriguez[2], and Jon E. Froehlich[1]
University of Washington[1], University of Illinois at Urbana-Champaign[2]

MAKEABILITY LAB

PAUL G. ALLEN SCHOOL
OF COMPUTER SCIENCE & ENGINEERING

UNIVERSITY OF ILLINOIS
URBANA-CHAMPAIGN

But voice assistants don't understand "**this**"

# GazePointAR Timeline

**01** — Designing GazePointAR v1

**02** — Three-part lab study with 12 participants

**03** — Designing GazePointAR v2

**04** — 5-days first-person diary study

**05** — Discussing the future of context-aware VAs

# GazePointAR Timeline



**01**

Designing
GazePointAR v1

# GazePointAR v1

# GazePointAR v1

# GazePointAR v1

# GazePointAR v1

# GazePointAR v1

# GazePointAR v1

# GazePointAR v1

# GazePointAR v1



Hey Glass?

Hi, I'm listening.

unity
Microsoft HoloLens

How much is **this**?

Object Recognition
OCR
Celebrity Recognition

Eye Gaze
Pointing Gesture
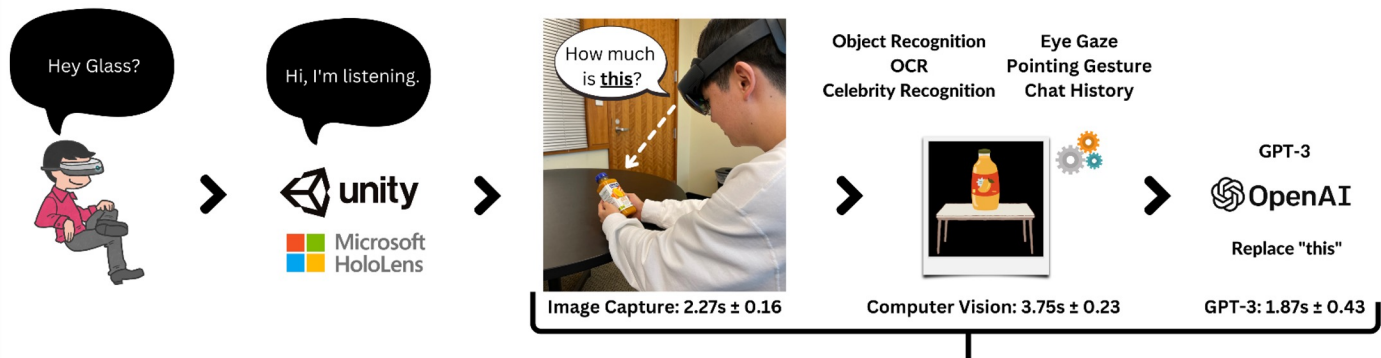Chat History

GPT-3
OpenAI

Replace "this"

Image Capture: 2.27s ± 0.16      Computer Vision: 3.75s ± 0.23      GPT-3: 1.87s ± 0.43

How much is **this**?

User's Gaze

User's Field of View

Bottle
Person

ML Results

Parent:
Bottle

Children:
Naked, Mighty, Mango, 290, calories

How much is <u>Bottle with text that says Naked Mighty Mango 290 calories</u>?

# GazePointAR v1

LANDFILL

RECYCLING

COMPOST

# GazePointAR Timeline

**01**

Designing
GazePointAR v1

**02**

Three-part lab study
with 12 participants

# Part 1 - Comparing VAs

In Part 1, participants used Google VA, Google Lens, and GazePointAR to find a recipe that uses a specific pasta sauce.

# Part 2 - Ambiguous Queries with GazePointAR

In Part 2, participants interacted with GazePointAR to complete three additional query tasks: math, price comparison, and celebrity search.



**Math Task**

**Price Comparison Task**

**Celebrity Task**

# Part 3 - Design Probe & Co-Design

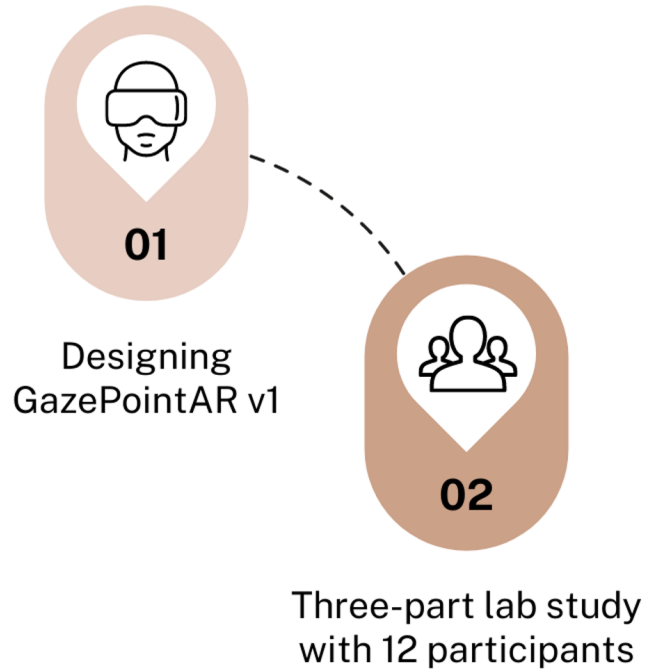In Part 3, participants first watched five design probe videos, then brainstormed and tried their own context-sensitive queries.



Cooking



Math



Language Translation



Recycling



Accessibility

# Key Findings

- GazePointAR is simple, fast, natural, and human-like.

- Participants preferred to speak pronouns, but not always.

- Pronouns are often for difficult-to-pronounce, long, or unknown object names.

- Gaze-only is preferred to keep interactions hands-free.

# Answerable Queries

# Unanswerable Queries

# Key Limitations

- Support multiple pronouns (e.g., "Which is healthier, **this** or **that**?").

- Support queries with no pronoun (e.g., "What would be good for dinner?").

- Provide explanations to its answers.

- Reduce having to dwell.

# GazePointAR Timeline



**01**

Designing
GazePointAR v1

**02**

Three-part lab study
with 12 participants

**03**

Designing
GazePointAR v2

GazePointAR v1

GazePointAR v2

GazePointAR v1

GPT-3

→ GPT-3.5

GazePointAR v2

GazePointAR v1

GPT-3

→ GPT-3.5

Google Cloud Vision API

→ YOLOv8

GazePointAR v2

GazePointAR v1

Google Cloud Vision API
→ YOLOv8

GPT-3
→ GPT-3.5

Manual Pronoun Replacement
→ Prompt Engineering

GazePointAR v2

# GazePointAR v2's Prompt Method

**Prompt**

The user asked, "<user-spoken query>"

To help you answer this question, here is what the user looked at: <gaze data>

The user also pointed at the following objects: <pointing data>

Finally, here are all other objects in the user's view: <all objects not gazed or pointed at>

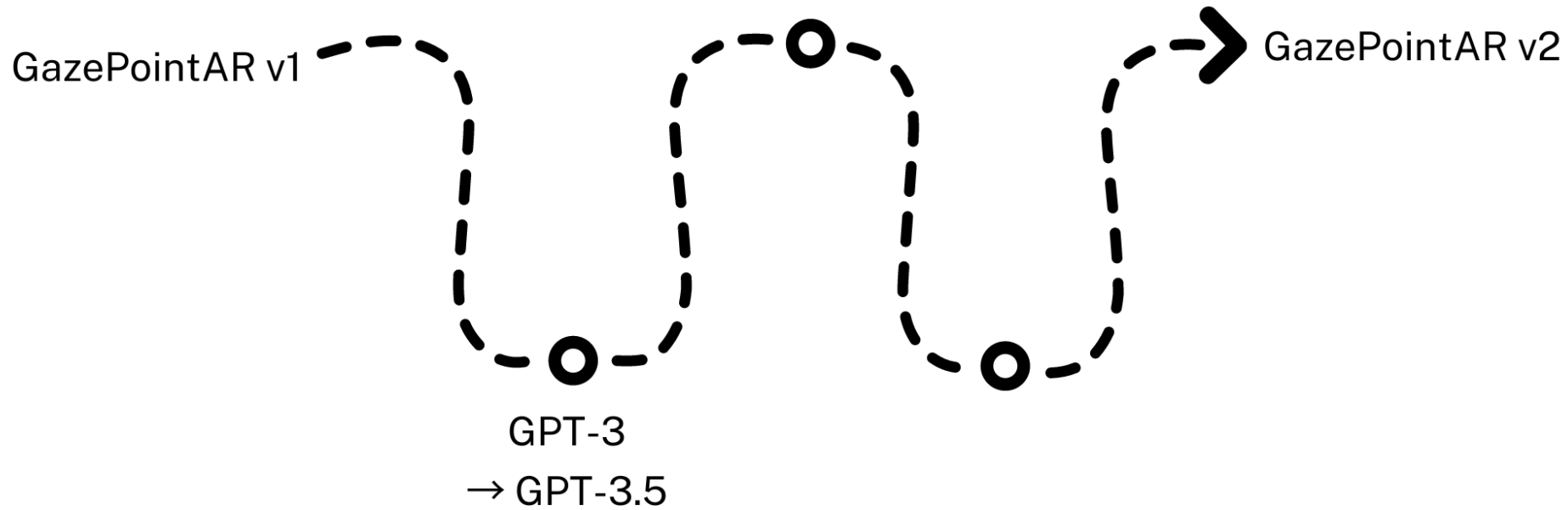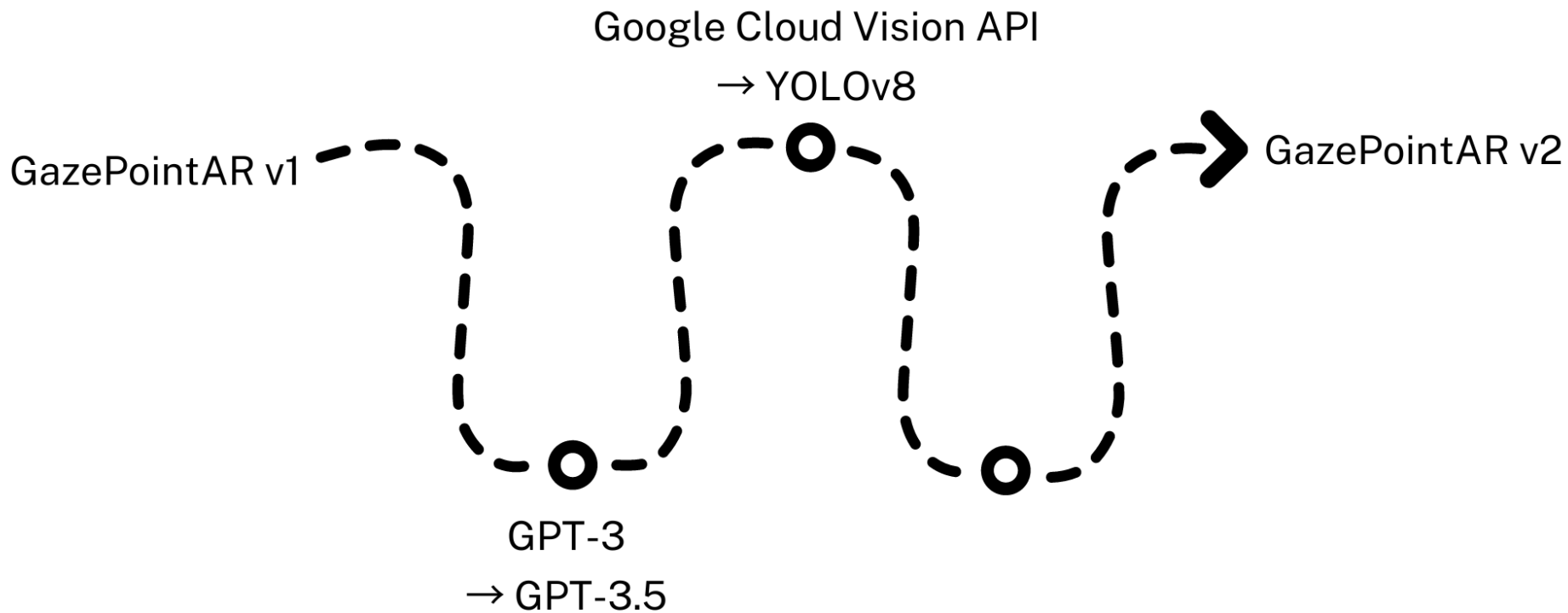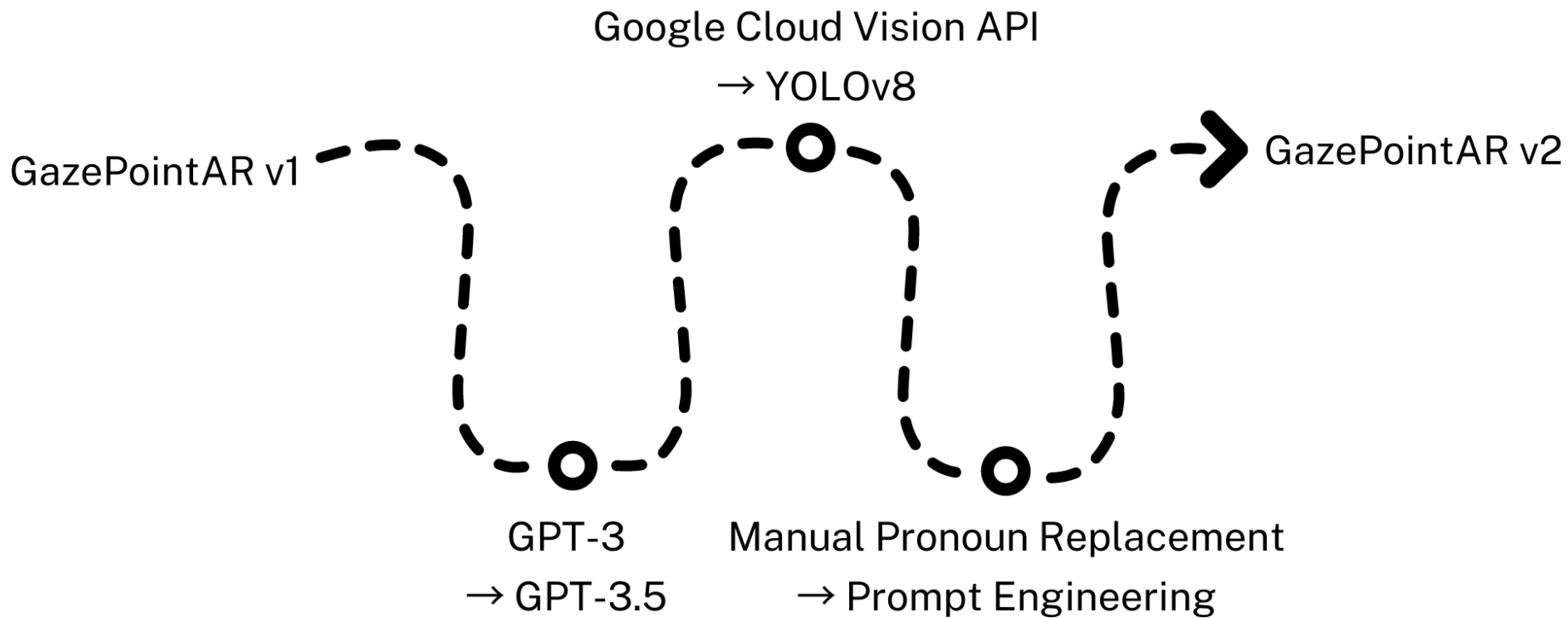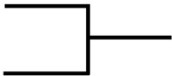Use the information above when answering the user's question, "<user-spoken query>". You should answer this question in one sentence. As part of your answer, include a short explanation. Even If you do not have enough information or an exact answer is unknown, you should do your best to provide an estimate or a range of possible answers.

**Explanations**

Insert original user-spoken query

Gaze data is still formatted as "<object or person name> with text that says <text 1> <text 2> <text 3> ..." However, child layer is no longer limited to 5 largest bounding box.

This line is included only if the user pointed at something. Pointing data is formatted the same as gaze data.

Insert semi-colon separated list of phrases describing objects not gazed or pointed at.

Output formatting to return a result that is exactly one sentence long with brief explanation.

# GazePointAR v2's Prompt Method

**Prompt**

The user asked, "<user-spoken query>"

**Explanations**

Insert original user-spoken query

# GazePointAR v2's Prompt Method

**Prompt**

The user asked, "<user-spoken query>"

To help you answer this question, here is what the user looked at: <gaze data>
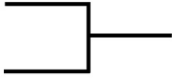
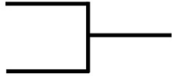**Explanations**

Insert original user-spoken query

Gaze data is still formatted as "<object or person name> with text that says <text 1> <text 2> <text 3> …" However, child layer is no longer limited to 5 largest bounding box.

# GazePointAR v2's Prompt Method

**Prompt**

The user asked, "<user-spoken query>"

To help you answer this question, here is what the user looked at: <gaze data>

The user also pointed at the following objects: <pointing data>

**Explanations**

Insert original user-spoken query
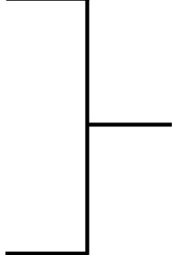
Gaze data is still formatted as "<object or person name> with text that says <text 1> <text 2> <text 3> ..." However, child layer is no longer limited to 5 largest bounding box.

This line is included only if the user pointed at something. Pointing data is formatted the same as gaze data.
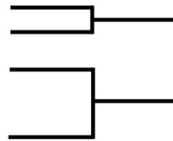
# GazePointAR v2's Prompt Method

**Prompt**

The user asked, "<user-spoken query>"

To help you answer this question, here is what the user looked at: <gaze data>

The user also pointed at the following objects: <pointing data>

Finally, here are all other objects in the user's view: <all objects not gazed or pointed at>

**Explanations**

Insert original user-spoken query

Gaze data is still formatted as "<object or person name> with text that says <text 1> <text 2> <text 3> ..." However, child layer is no longer limited to 5 largest bounding box.

This line is included only if the user pointed at something. Pointing data is formatted the same as gaze data.

Insert semi-colon separated list of phrases describing objects not gazed or pointed at.

# GazePointAR v2's Prompt Method

**Prompt**

The user asked, "<user-spoken query>"

To help you answer this question, here is what the user looked at: <gaze data>
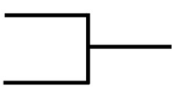
The user also pointed at the following objects: <pointing data>

Finally, here are all other objects in the user's view: <all objects not gazed or pointed at>

Use the information above when answering the user's question, "<user-spoken query>". You should answer this question in one sentence. As part of your answer, include a short explanation. Even If you do not have enough information or an exact answer is unknown, you should do your best to provide an estimate or a range of possible answers.

**Explanations**

Insert original user-spoken query

Gaze data is still formatted as "<object or person name> with text that says <text 1> <text 2> <text 3> ..." However, child layer is no longer limited to 5 largest bounding box.

This line is included only if the user pointed at something. Pointing data is formatted the same as gaze data.

Insert semi-colon separated list of phrases describing objects not gazed or pointed at.

Output formatting to return a result that is exactly one sentence long with brief explanation.

# GazePointAR Timeline



**01** Designing GazePointAR v1

**02** Three-part lab study with 12 participants

**03** Designing GazePointAR v2

**04** 5-days first-person diary study

# Study 2 - First-Person Diary Study

With GazePointAR v2, the first author used GazePointAR in their day-to-day activities for **five** days, recording their observations.



A: …you may consider brand like Hugo Boss, Brooks Brothers, or J.Crew, which offer quality men's cloth at relatively lower costs.

A: Based on your gaze, you are interested in the cappuccino, so I would recommend you to try the cappuccino.

A: …The Brave New World Revisited by Aldous Huxley, as it is written by the same author and is often considered a companion piece.

A: …consider making stuffed zucchini boats, as they are a delicious and healthy way to use both zucchini and peppers.

# ARMI STICE

| | 8oz | 12oz | 16oz | 20oz* |
|---|---|---|---|---|
| Drip Coffee | 3.75 | all sizes | | |
| Chai | 4.5 | 5 | 5.5 | 6 |
| Tumeric Latte | 4.5 | 5 | 5.5 | 6 |
| Matcha | 4.5 | 5 | 5.5 | 6 |
| Hot Chocolate | 3.75 | 4.25 | 4.75 | 5.25 |
| Steamer | 3 | 3.50 | 4 | 4.50 |
| Brewed Tea | 3.75 | all sizes | | |
| Pot of Tea | 5.75 | | | |
| | | | | |
| Extra Espresso | 1.00 | | | |
| *4 Shots | | | | |

## HOT OR ICED DRINKS

| | 8oz | 12oz | 16oz | 20oz* |
|---|---|---|---|---|
| Americano | | | | |
| Latte | 4.25 | 4.25 | 4.25 | 4.75 |
| Mocha | 4.75 | 5.25 | 4.25 | 4.75 |
| Cold Brew | 4.75 | 5.25 | 5.75 | 6.25 |
| Espresso | | 4.5 | 5 | 6.25 |
| Macchiato | 4 | one size | | |
| Cappuccino | 4.75 | one size | | |
| Cortado | 4.75 | one size | | |
| Turkish Coffee | 4.75 | | | |
| | 4.75 | | | |

### MILKS

Whole
Nonfat

Soy*
Oat*
Almond*
Coconut*

### FLAVORS

+0.50
Vanilla
Sugar Free Vanilla
Hazelnut
House Lavender
Caramel
Chocolate
White Chocolate

Hey glass?
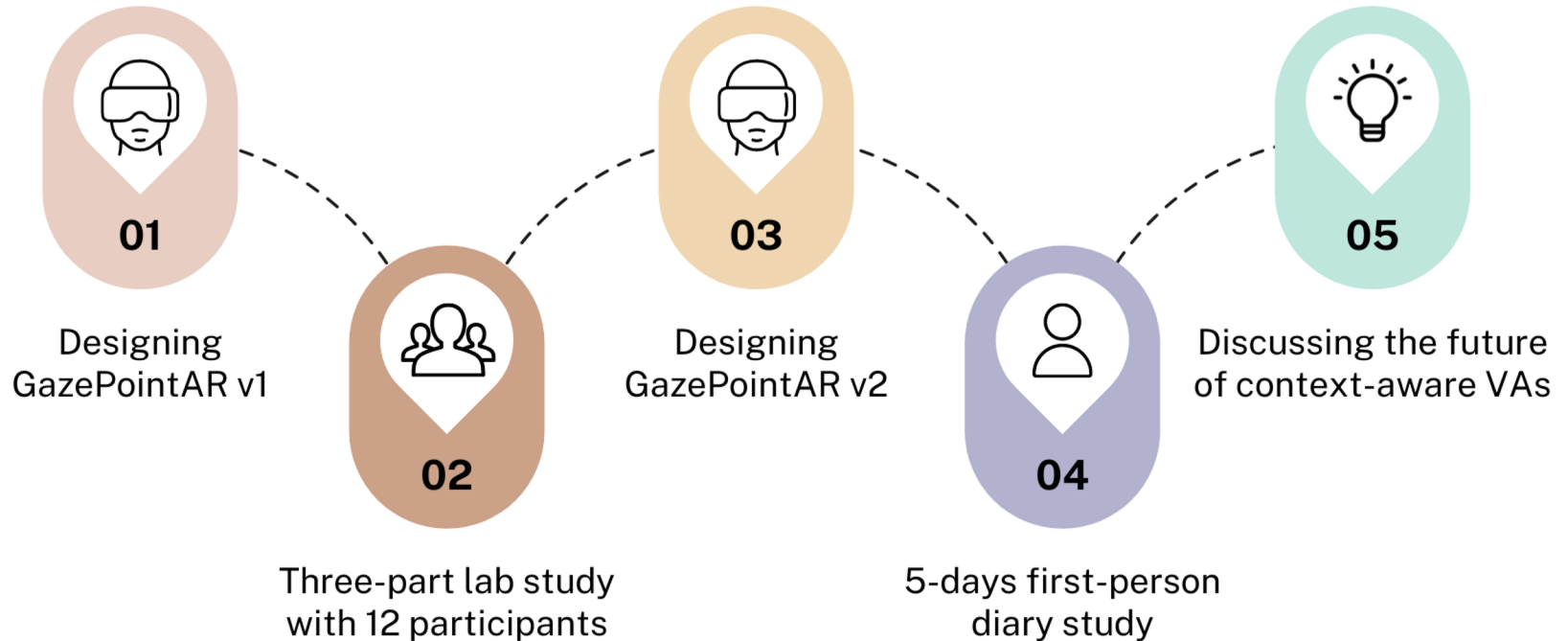
huh

# Study 2 Findings

- GazePointAR is natural and companion-like.

- GazePointAR still has several limitations:

  - Referents in the past require gaze history (e.g., "Where did I leave my keys?").

  - Multiple referents require gaze shift (e.g., "Which is the healthier, **this** or **that**?") .

  - Pointing is impractical in public.

  - Extended dwelling causes eye fatigue.

# GazePointAR Timeline

**01** Designing GazePointAR v1

**02** Three-part lab study with 12 participants

**03** Designing GazePointAR v2

**04** 5-days first-person diary study

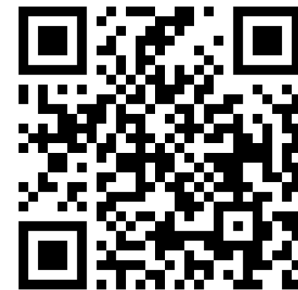**05** Discussing the future of context-aware VAs

# So, What's Next?



- GazePointAR is simple, natural, and human-like to speak to.

- Future context-aware VAs should leverage longitudinal natural eye gaze input.

- Future research should further study LLM-driven query disambiguation.

"I want to say queries with and without pronouns, because whichever comes to mind first, that's the one I want to say. GazePointAR should adapt to me and my natural eye gaze" [P12].

# Thanks for listening!

## Let's Connect!

**Email**         jaewook4@cs.washington.edu

**Website**       https://jaewook-lee.com

**Twitter**        https://twitter.com/jaewook_jae