

# Multimodal AI

## Lecture 14.2 – Human-AI Interaction

**Paul Liang**

Assistant Professor

MIT Media Lab & MIT EECS



<https://pliang279.github.io>

[ppliang@mit.edu](mailto:ppliang@mit.edu)

 [@pliang279](https://twitter.com/pliang279)



# Assignments for This Coming Week

Please fill out course evaluations and give us feedback!

HW5 extension -- due **next Wednesday May 13**.

For project:

- Make sure to meet with myself and TAs this week if you need advice.
- Presentations next Tuesday May 12 – in-class slides presentation
- Final report due Tuesday May 19.

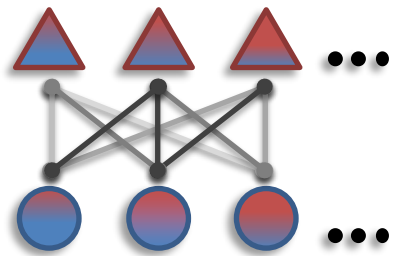
# Today's lecture

- 1 AI for smell, taste, temperature
- 2 Human-AI interaction
- 3 Ethics and safety

# Multimodal Human-AI Interaction

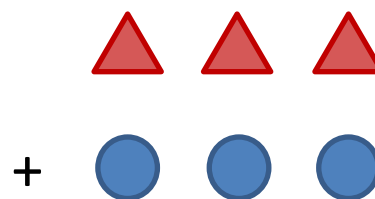
Solving hard problems by breaking them down into step-by-step reasoning steps in multiple modalities

*It's just a privilege to watch your mind at work.*



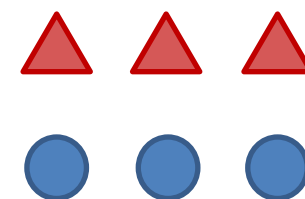
Multimodal representation

*This person is being sarcastic.  
They seem to be close friends.*



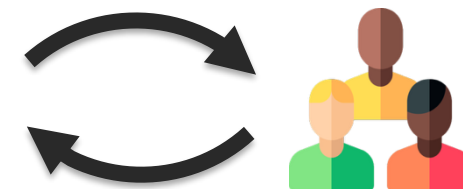
*(quote previous episodes)  
(highlight multimodal information)*

*Here's a story of them in a different culture...*



*(generate future episodes)*

**Models: Multimodal fusion and generation**  
**Data: Hard challenges + human reasoning steps**  
**Training: Reinforcement learning for emergent reasoning**  
**Human: Trustworthy, safe, controllable**



# AI for Smell

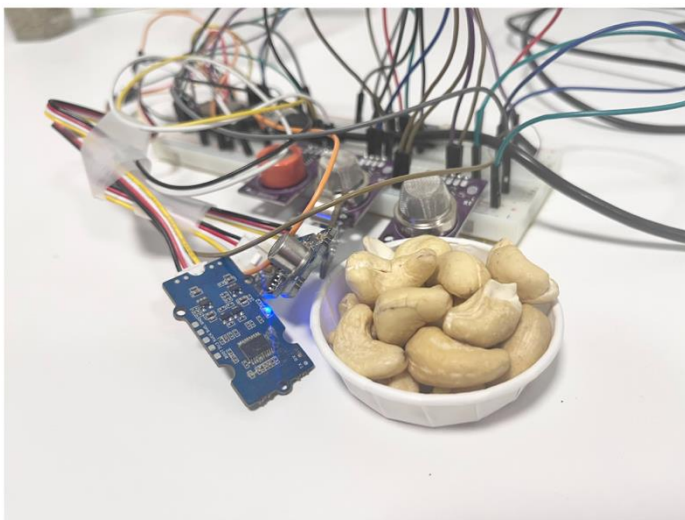
<https://github.com/MIT-MI/SmellNet>

## Smell Sensor Detection (Oregano)

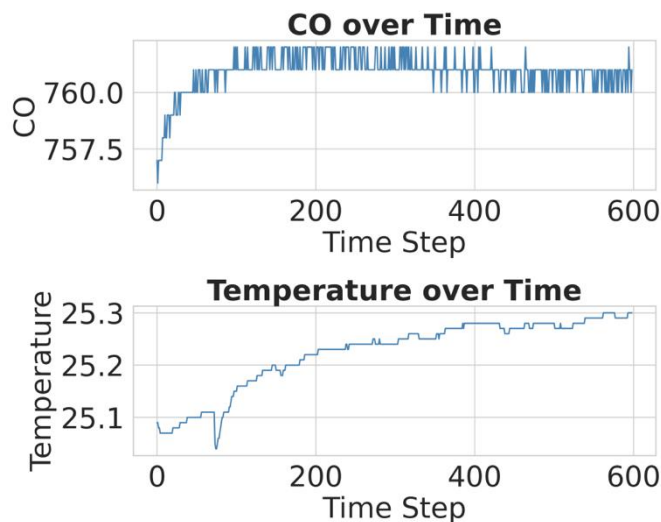


# Creating SmellNet

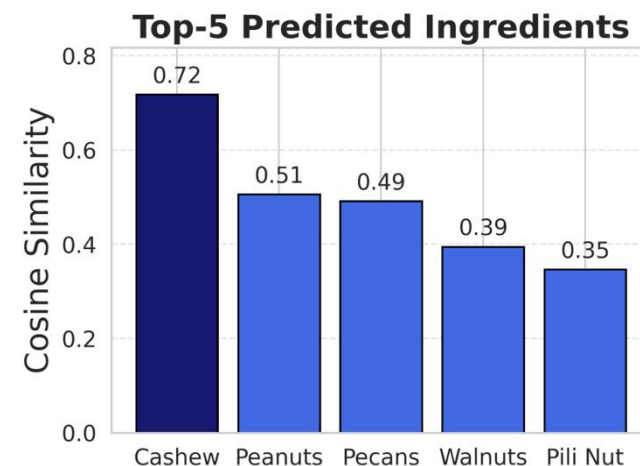
<https://github.com/MIT-MI/SmellNet>



(a) Sensor setup detecting cashew.



(b) Time-series signals from CO and temperature sensors.

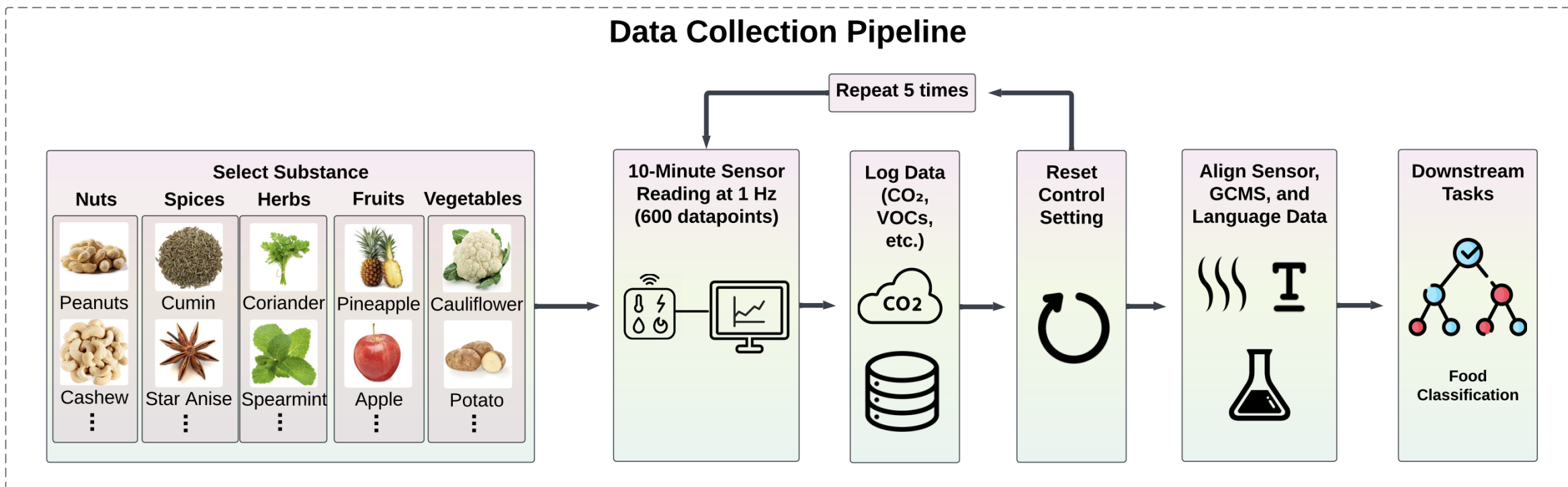


(c) Top-5 model predictions using cosine similarity.

# Creating SmellNet

<https://github.com/MIT-MI/SmellNet>

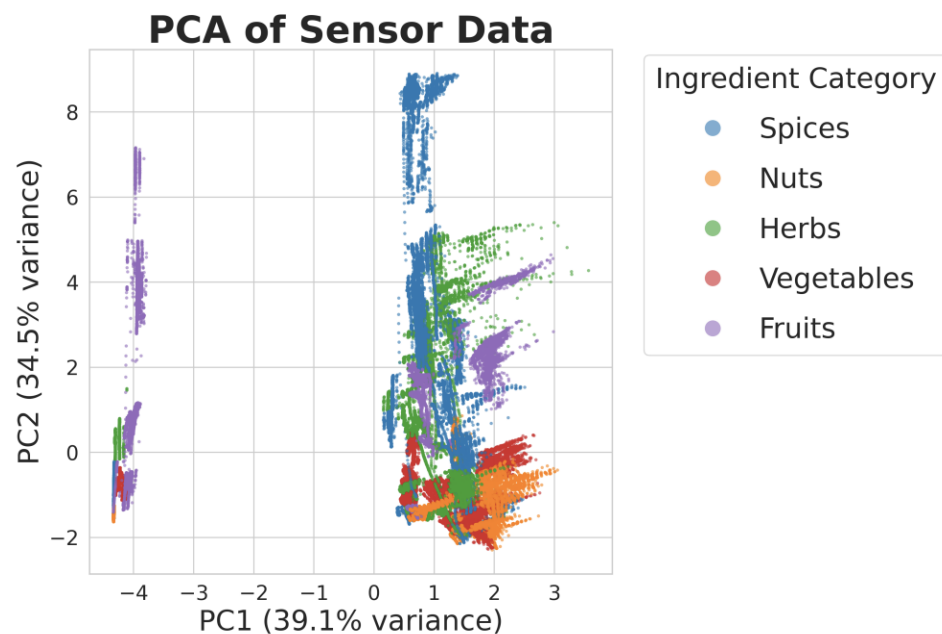
Data is key! More than 50 hours across 50 substances. >300,000 datapoints



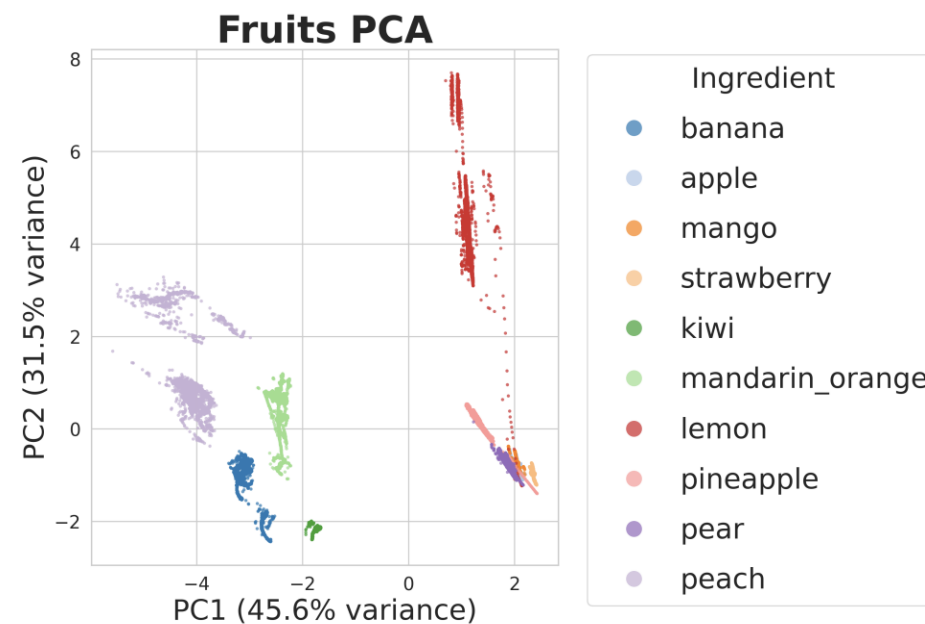
# AI for Smell Recognition

<https://github.com/MIT-MI/SmellNet>

PCA: a way to visualize >2-dimensional data in 2 dimensions



(a) All Substances PCA

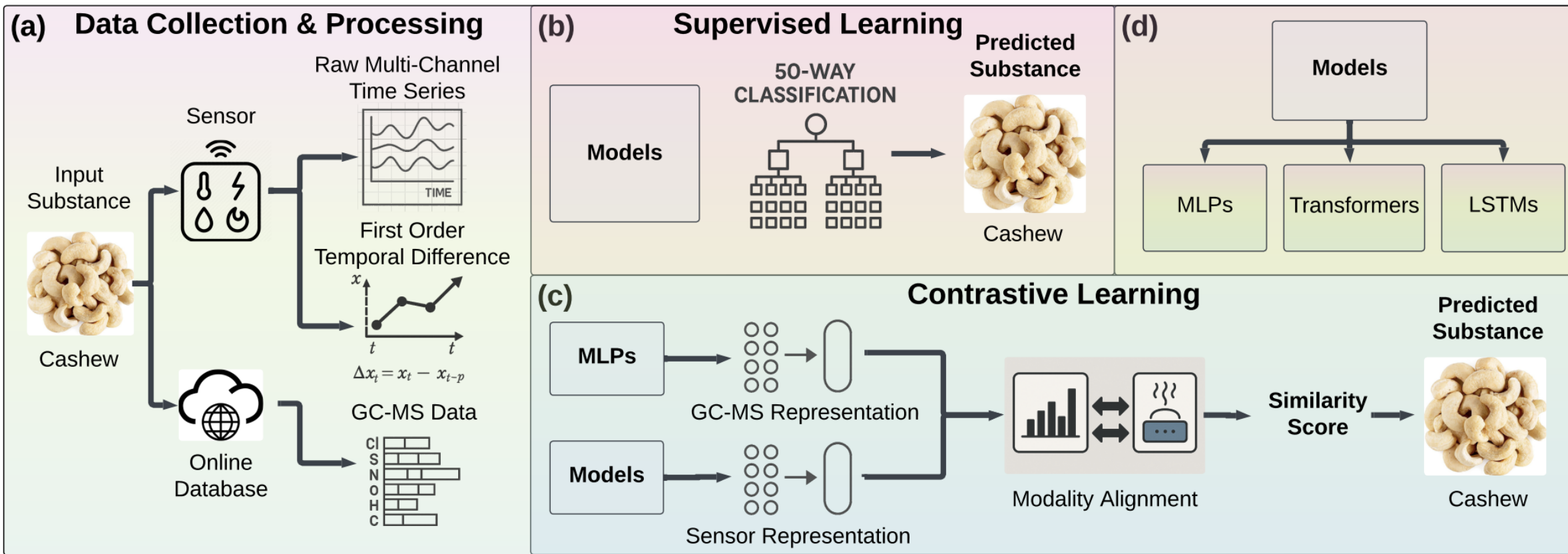


(b) Fruits PCA

# AI for Smell Recognition

<https://github.com/MIT-MI/SmellNet>

Modeling techniques: sequence models, temporal difference models, alignment with online databases



# Smell Recognition Results

<https://github.com/MIT-MI/SmellNet>

ScentFormer applies temporal difference and sliding windows to capture relative changes

Per-Category Accuracy @1 ( $p=25$ ) — Regular vs Contrastive

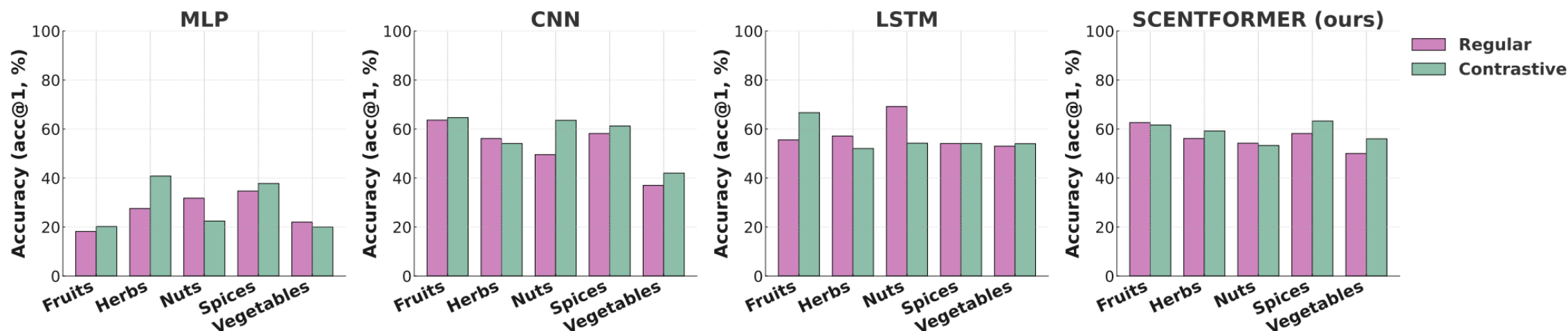
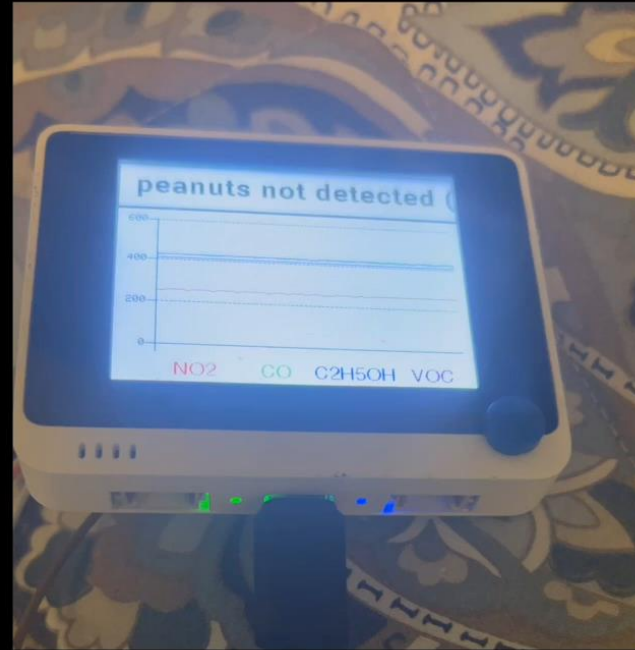
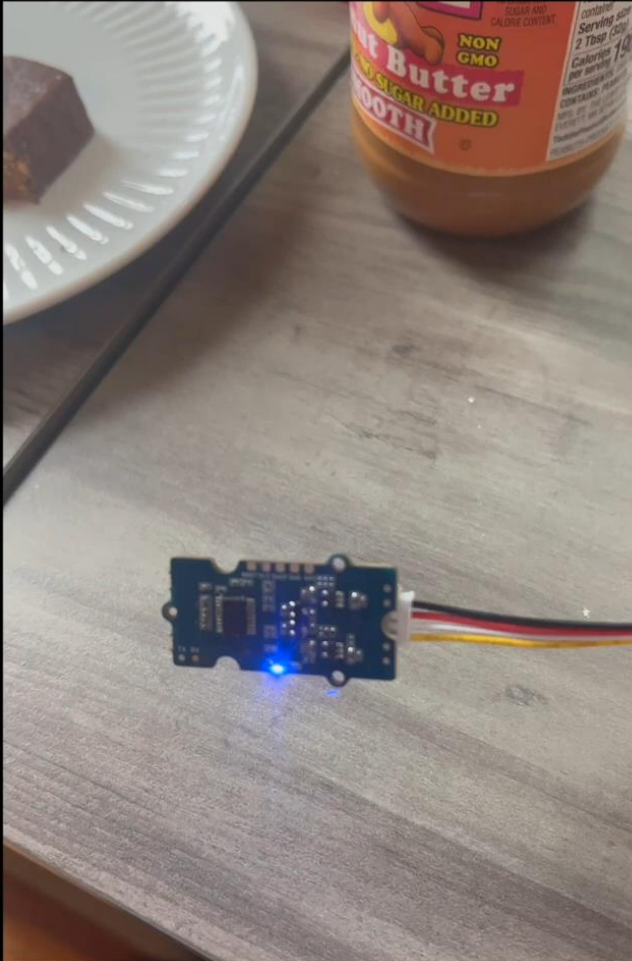


Figure 5: **Per-category accuracy** ( $\text{acc}@1$ ) for four models at lag  $p = 25$ . Bars are paired per category (Regular vs. Contrastive). Figure shows the non temporal models suffer from categories like vegetables, but temporal architectures demonstrate stronger robustness across categories.

# Applications in Allergen Detection

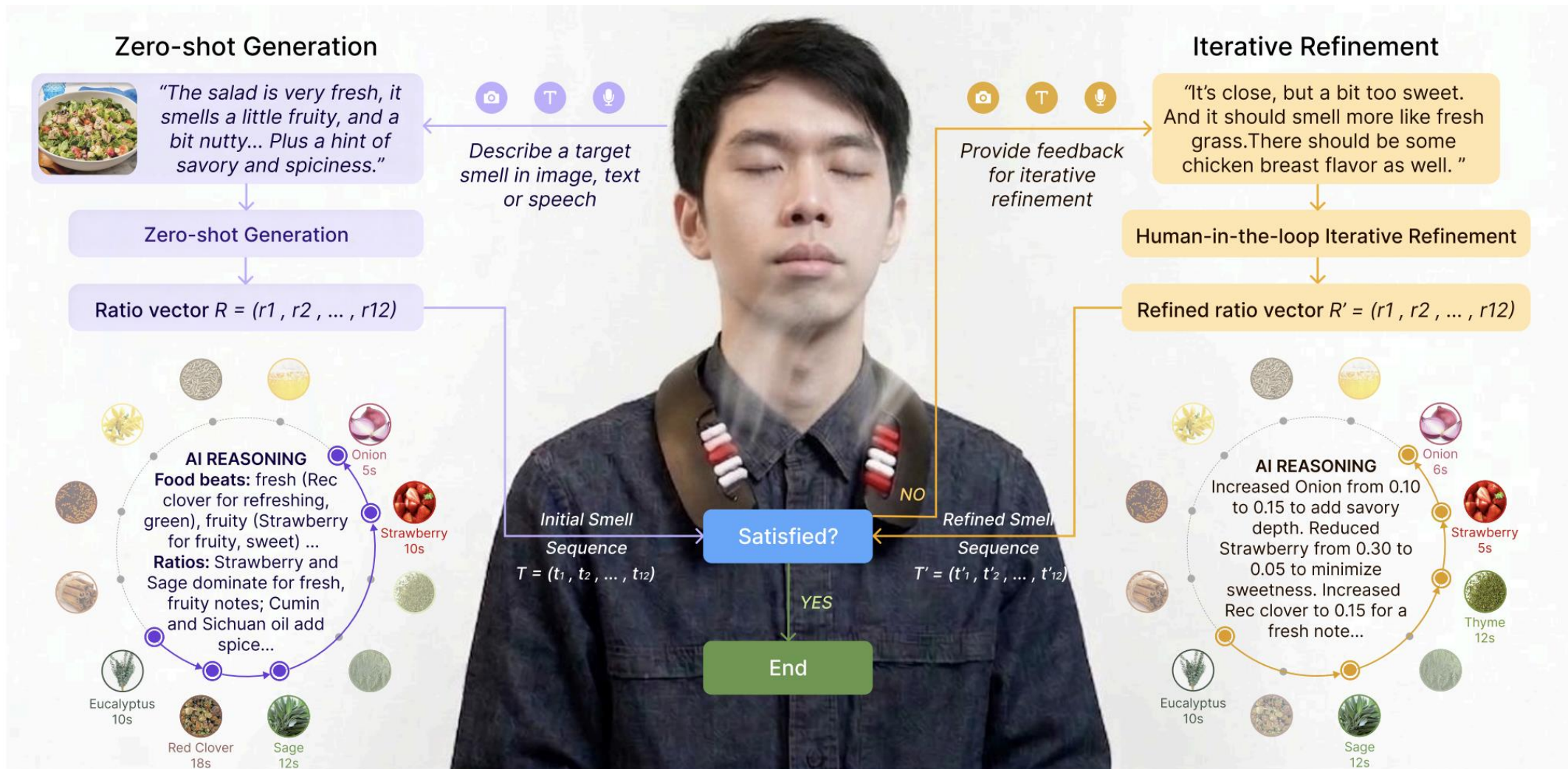
<https://github.com/MIT-MI/SmellNet>



# AromaGen: Towards Smell Transmission



# AromaGen: Towards Smell Transmission



# A Language for Smell?

Participants mostly described smell based on taste, texture, and memory, but a shared vocabulary emerges.

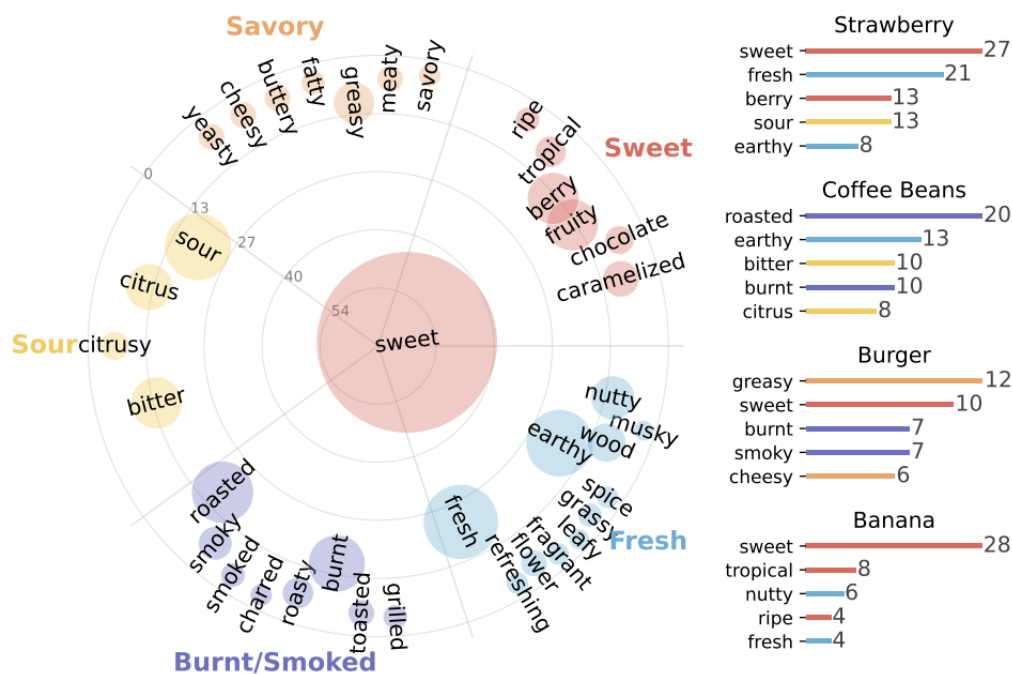


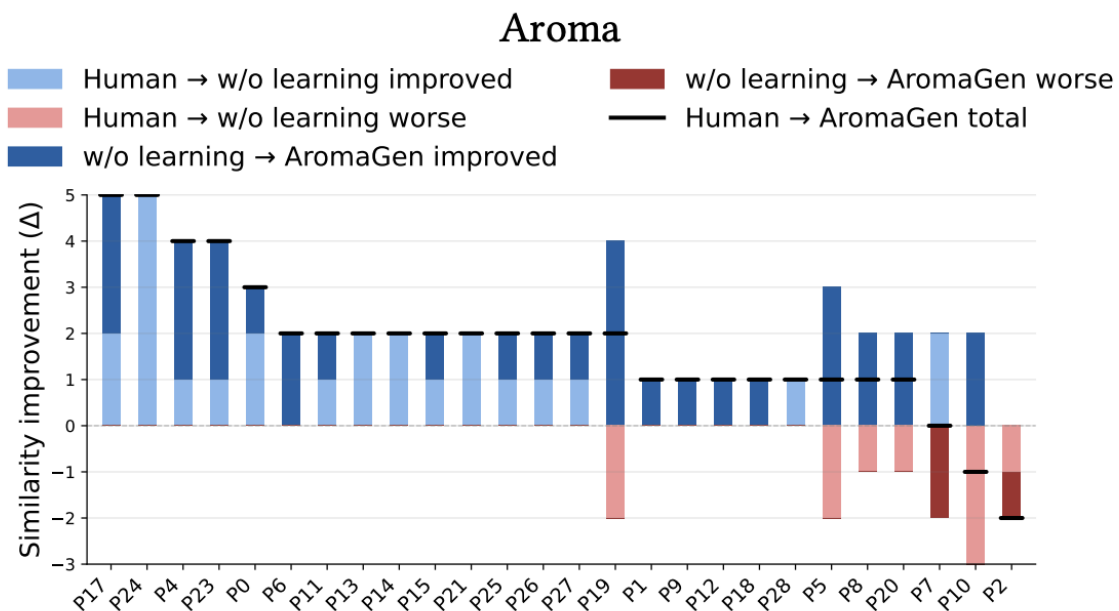
Figure 4: The 12 base odorants used in AROMAGEN's palette.

Odorant	Vol.	Notes
Cumin	6	Smoky, spice
Ylang Ylang	6	Warm, light spice
Sichuan Oil	3	Light, chai, spice
Cinnamon	5	Sweet, spice, coffee, warm
Eucalyptus	5	Refreshing, spa
Red Clover	5	Mint, clover, green, refreshing
Sage	6	Refreshing
Cypress	5	Woody stability
Thyme	5	Bitter, green, vegetable
Strawberry	5	Elegant clarity, fruity, sweet
Onion	6	Umami, onion, chips, savory
Isovaleric Acid	8	Cheesy, sweat, sour

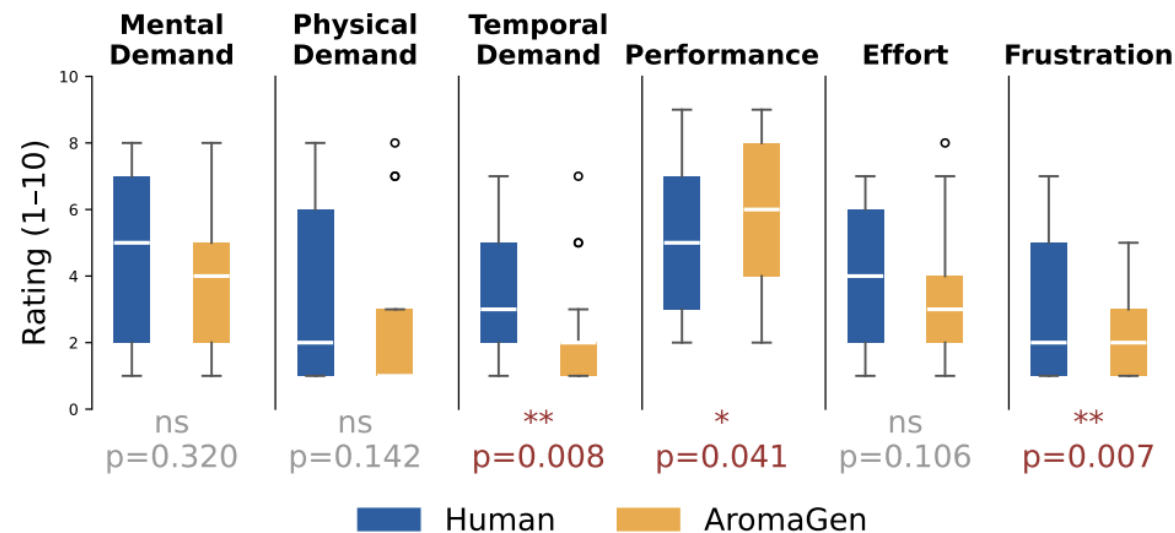
Table 3: The 12 base odorants in AROMAGEN's palette, with volatility score (1–10) and characteristic notes. Colored squares indicate perceptual categories: ■ Sweet, ■ Savory, ■ Sour, ■ Burnt/Smoked, ■ Fresh.

# AromaGen Key Results

Learning procedure in AromaGen improves similarity.



AromaGen saves effort of human mixing with similar performance.

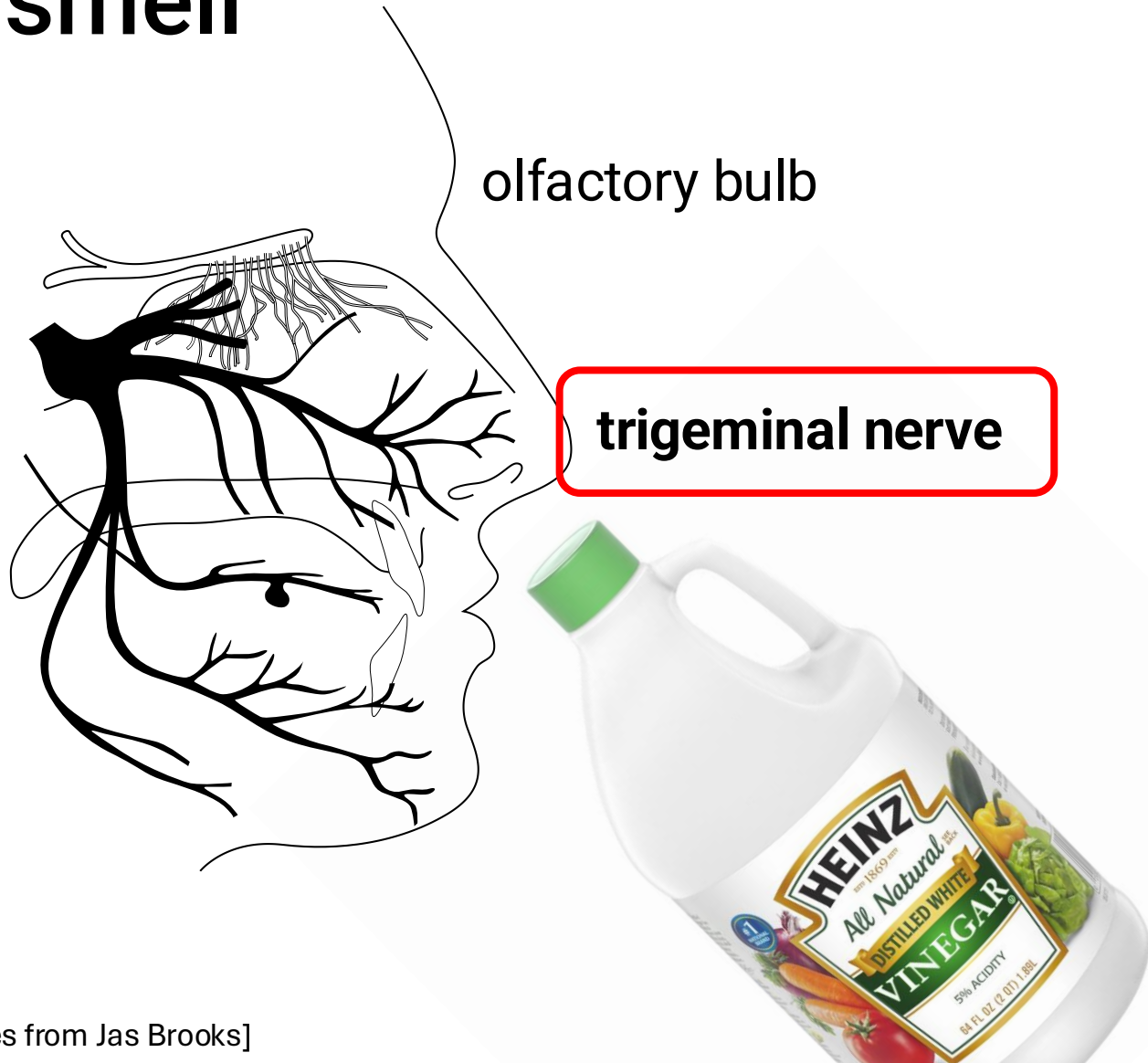


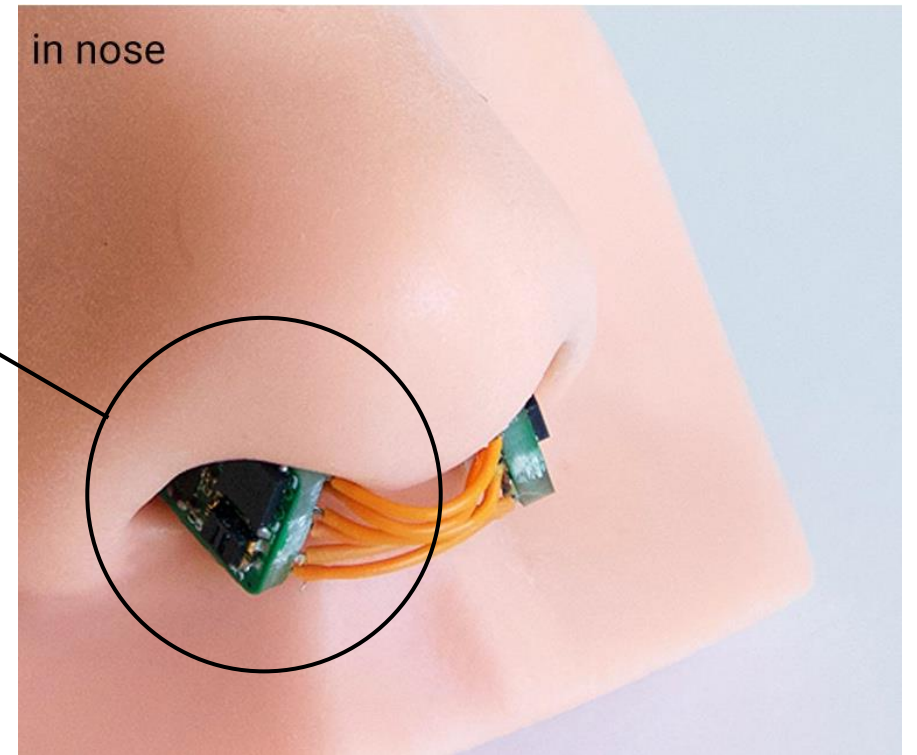
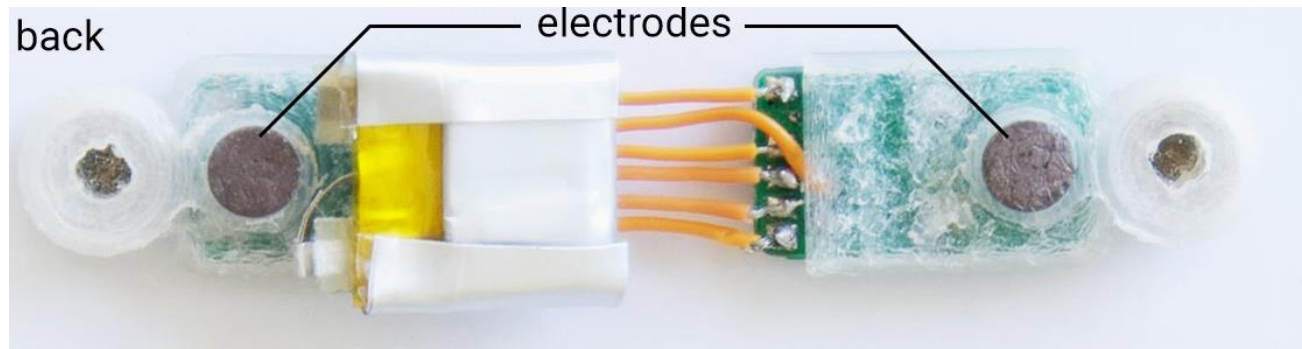
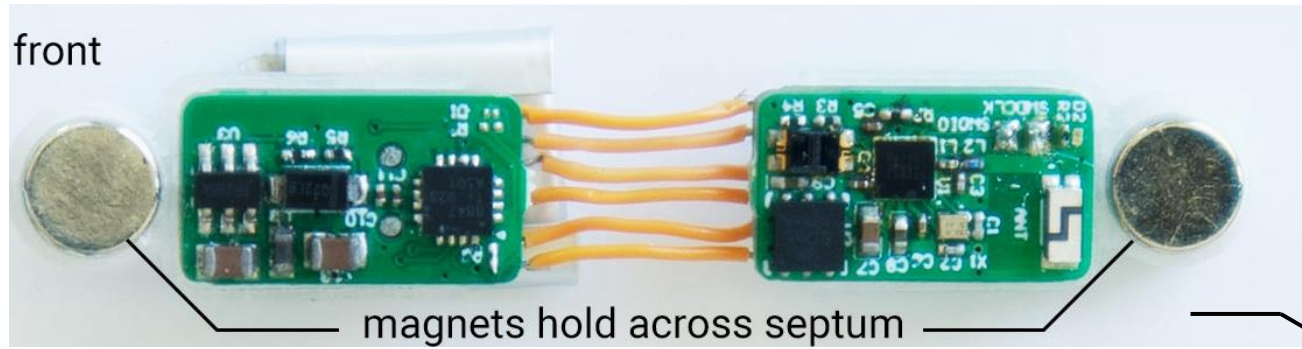


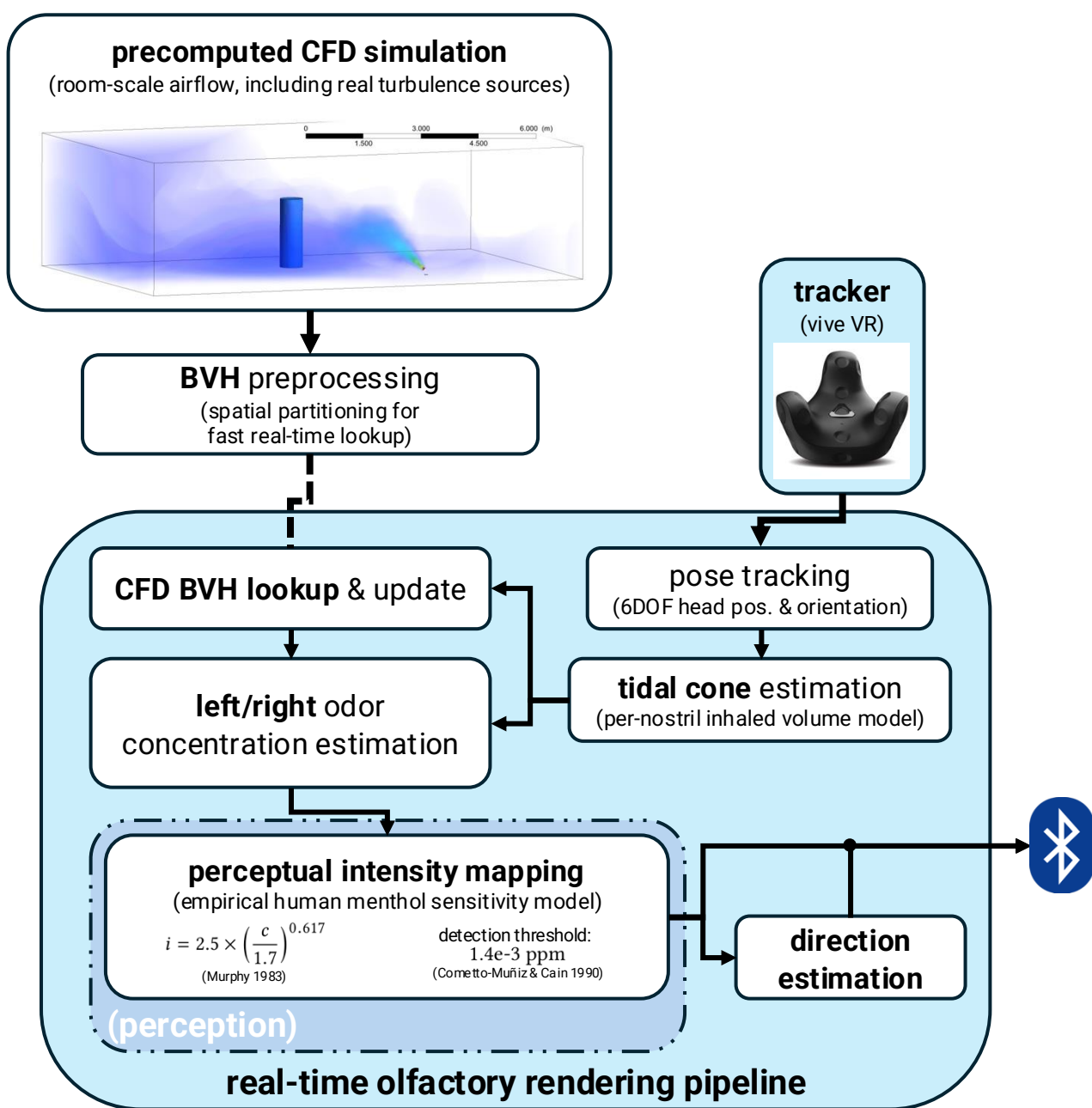
**our goal:**  
can we deliver smell sensations  
**without chemicals?**

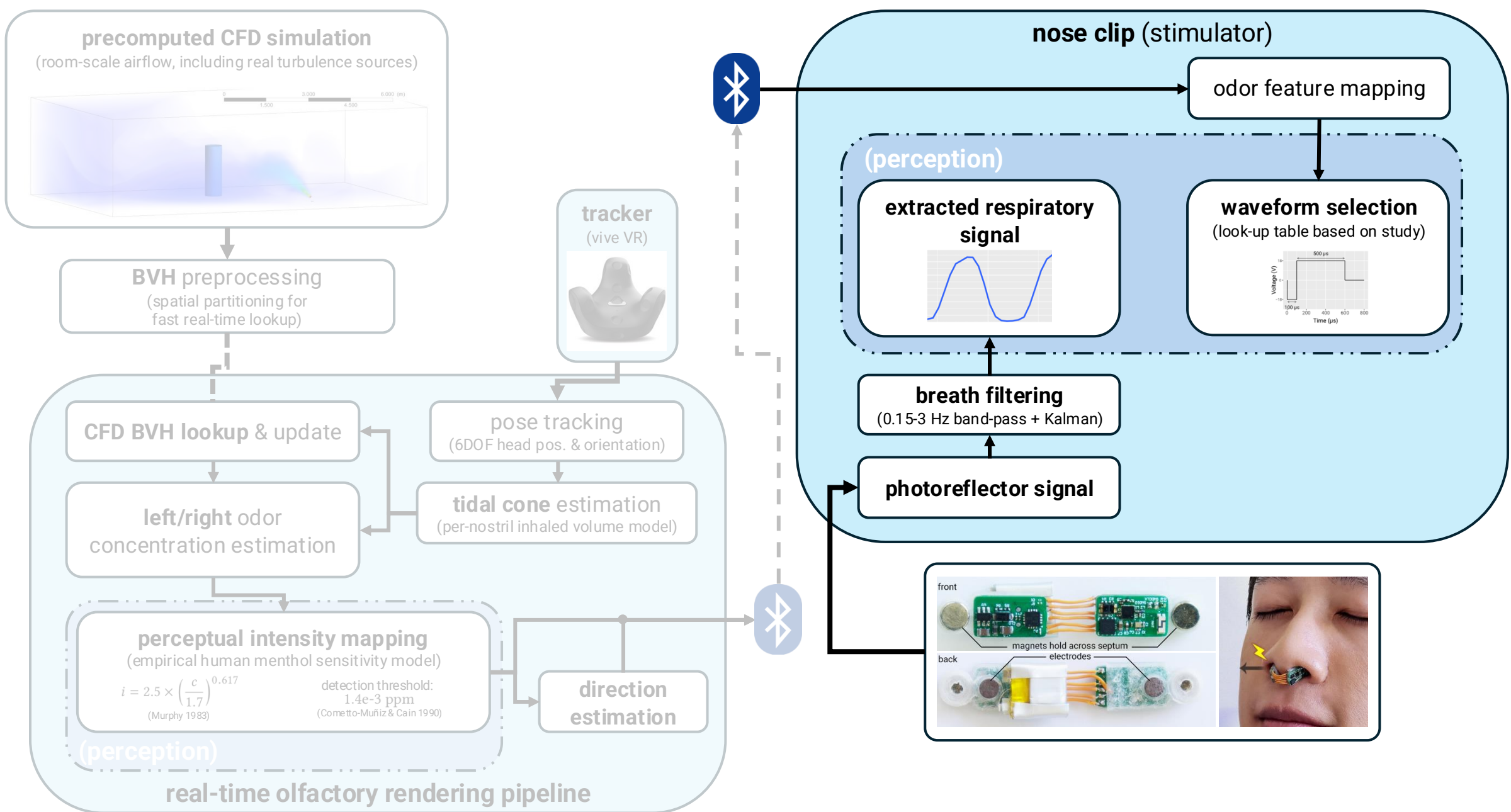


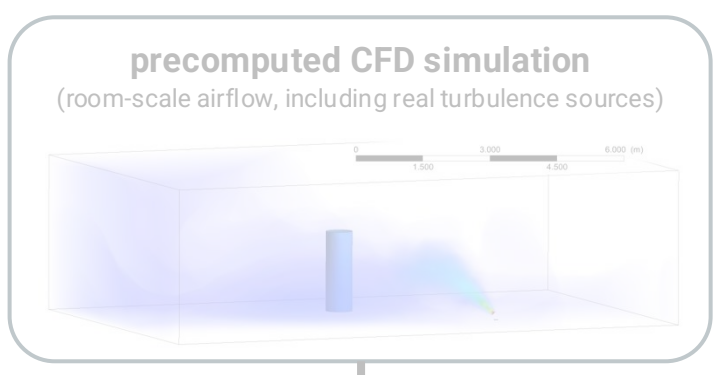
# smell



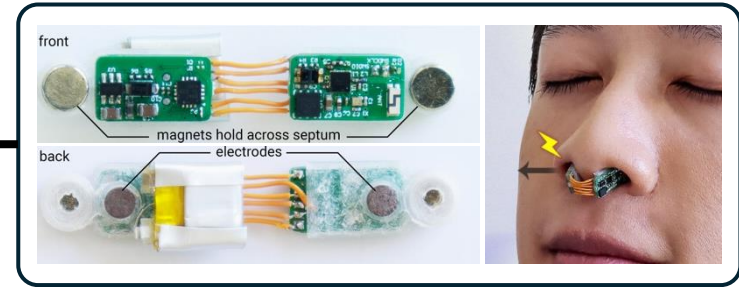
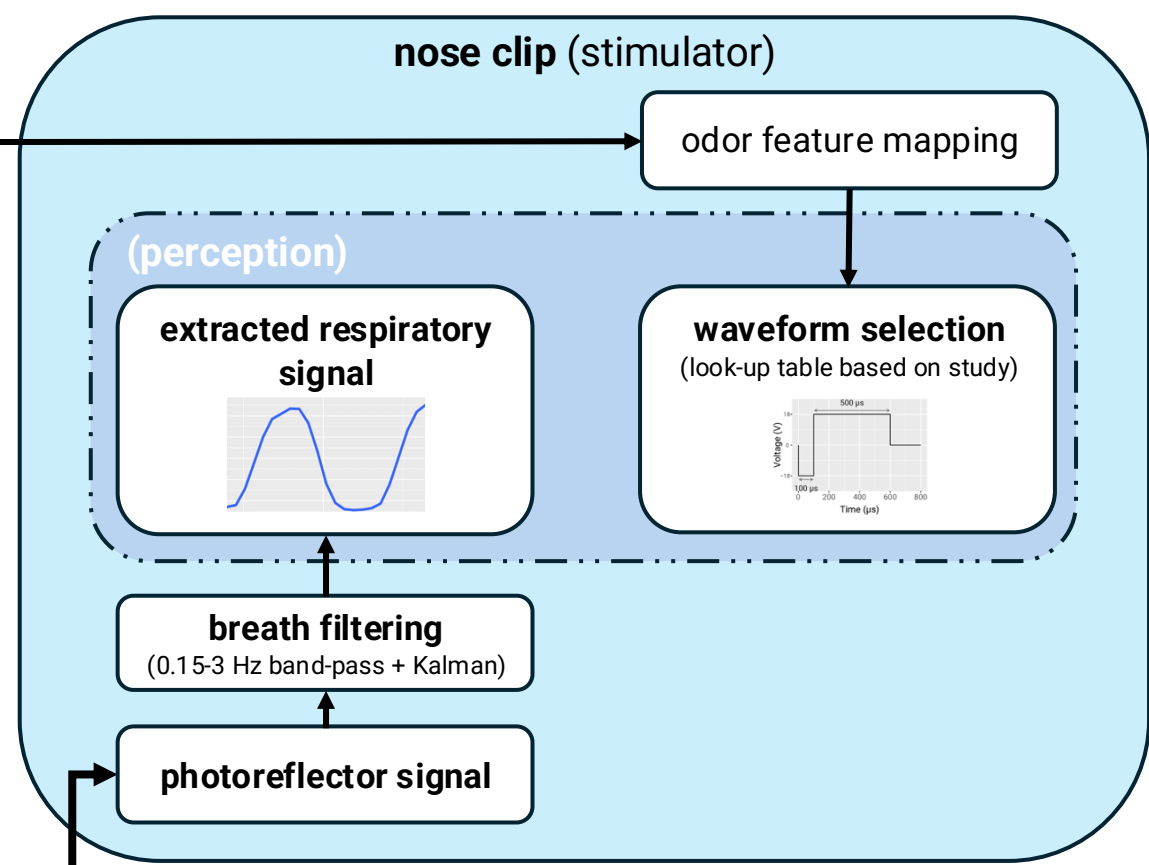
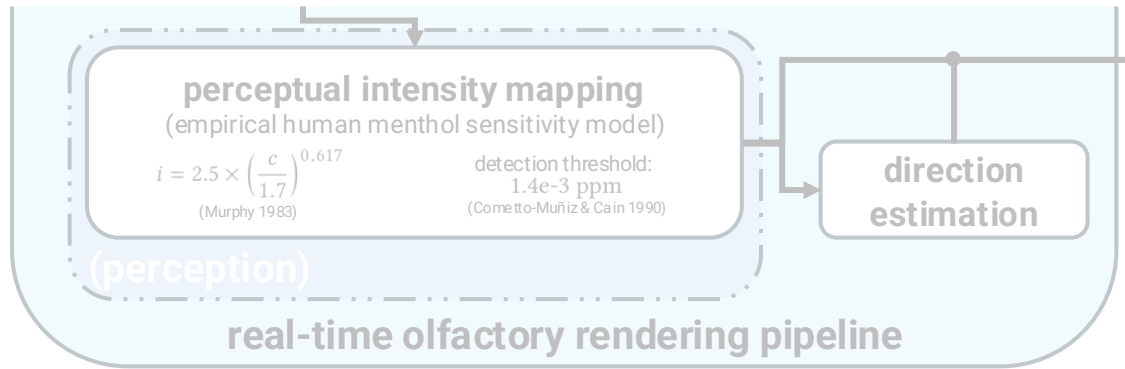
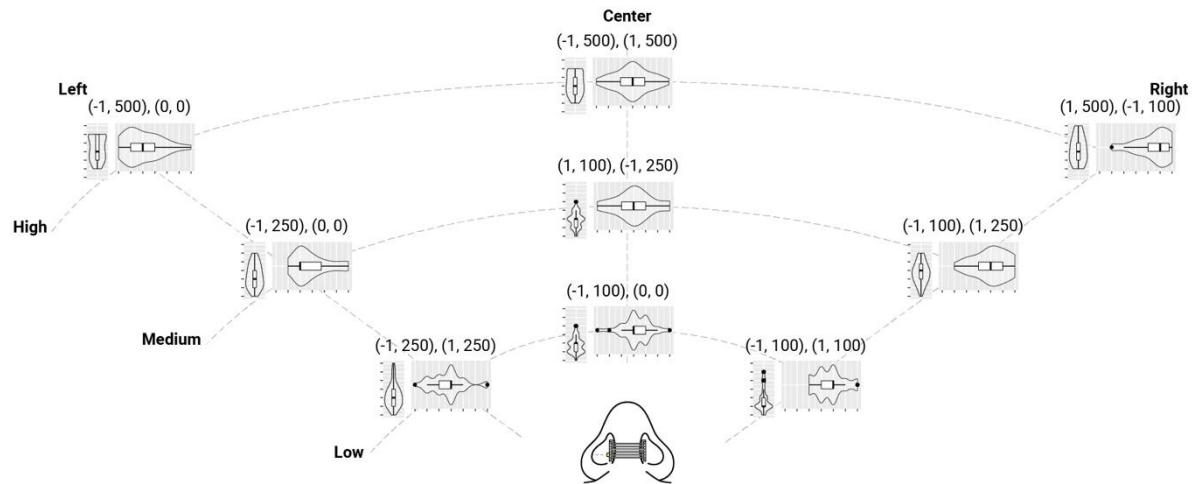


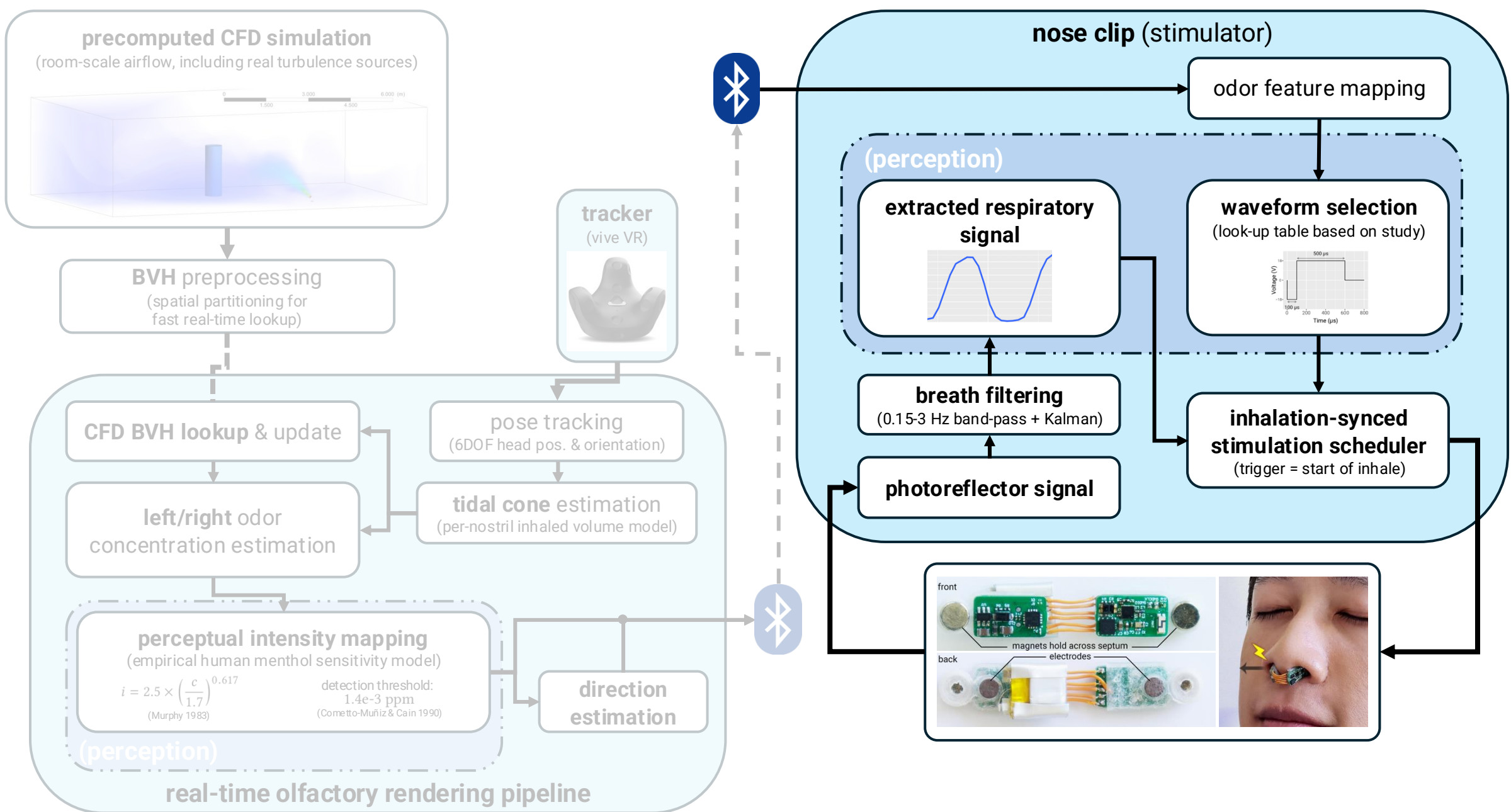


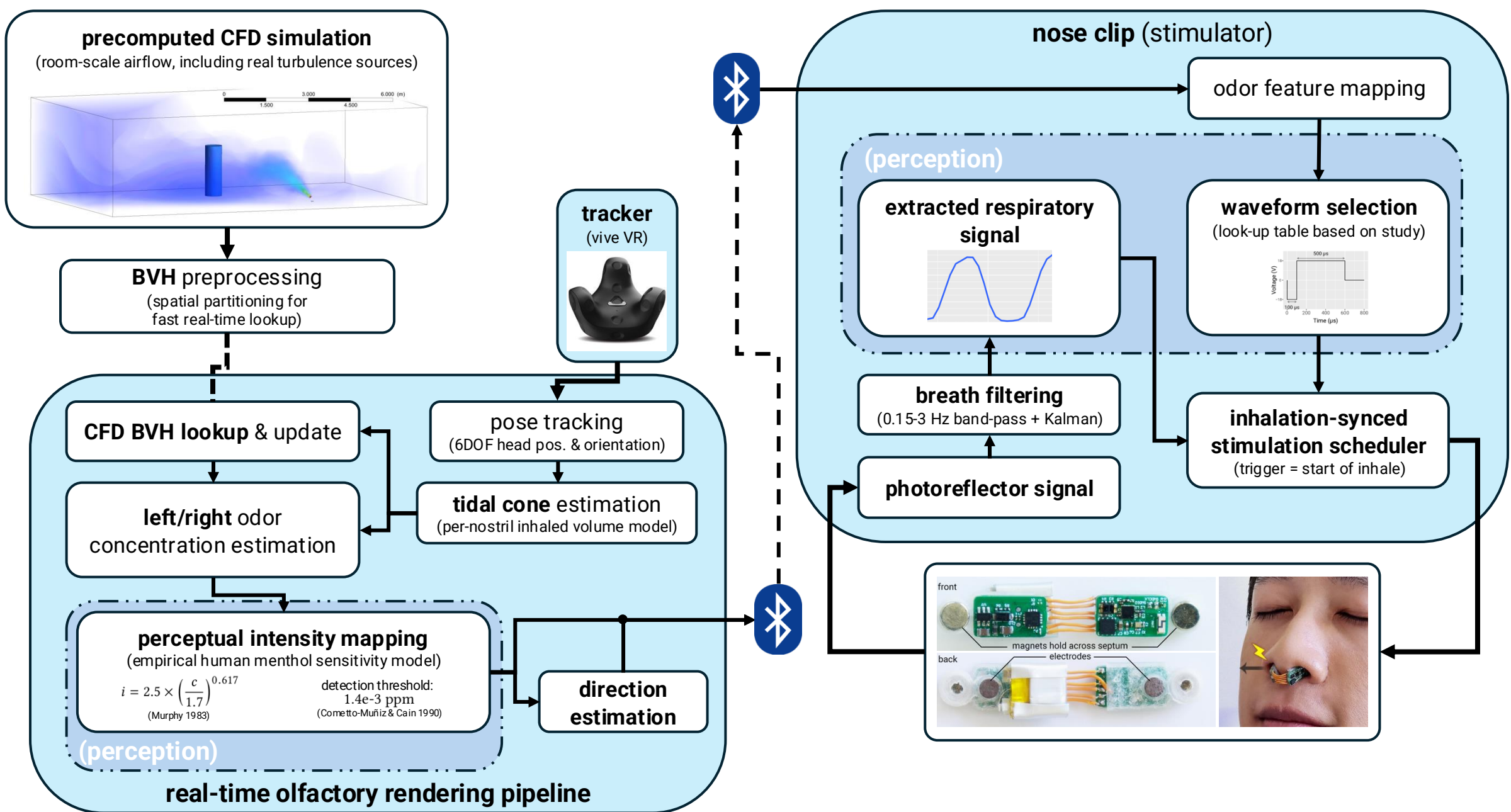


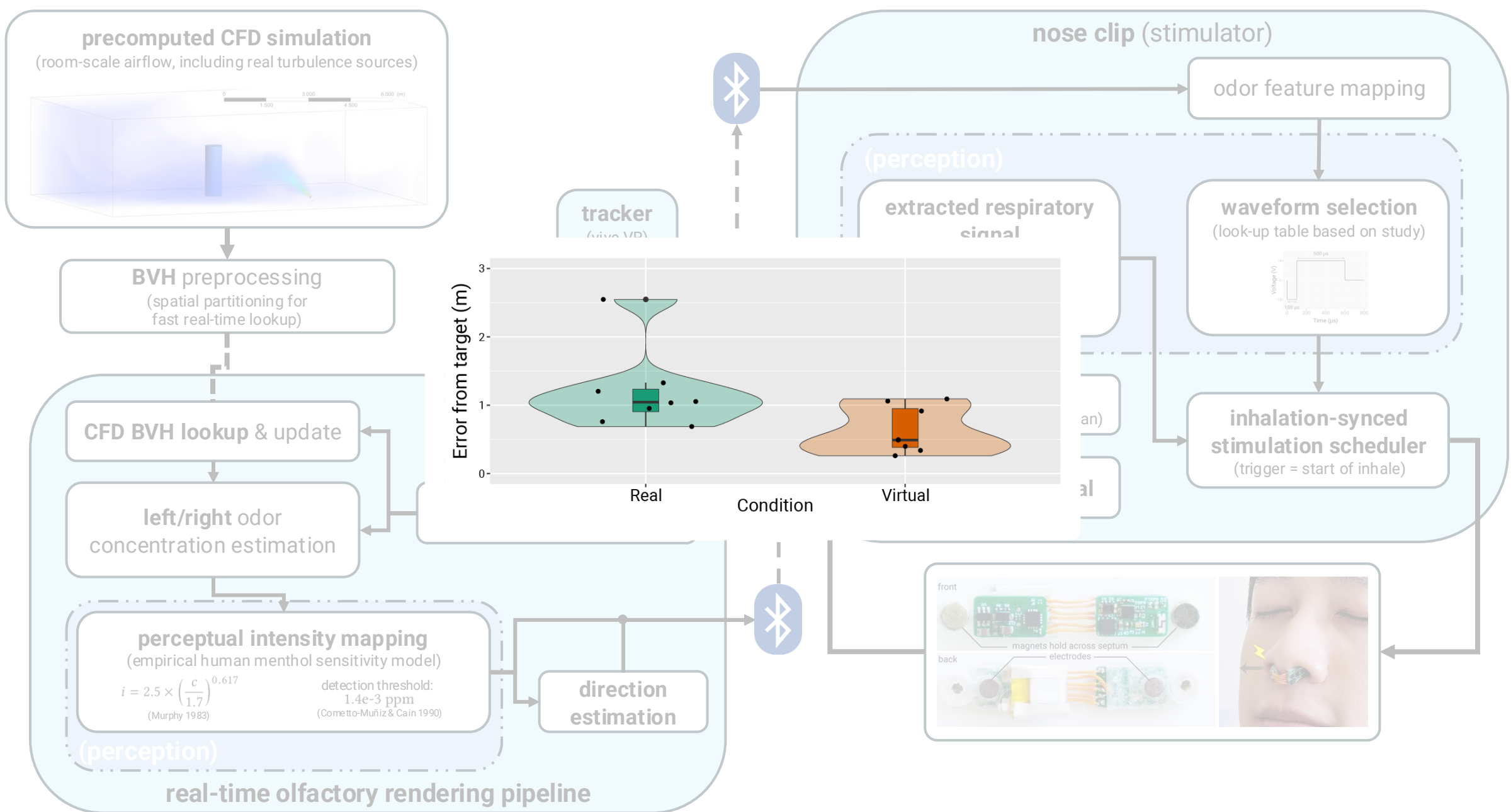


**tracker**  
(live VD)

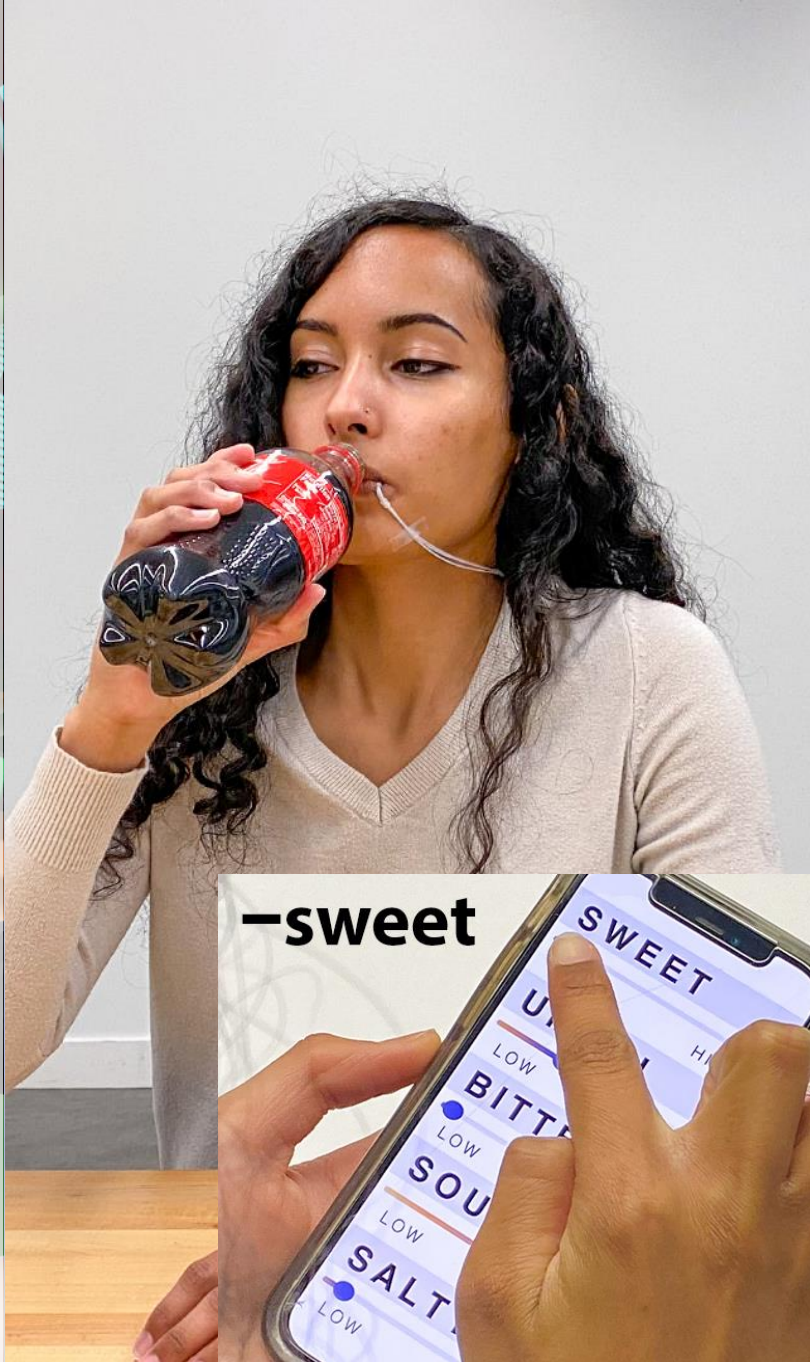
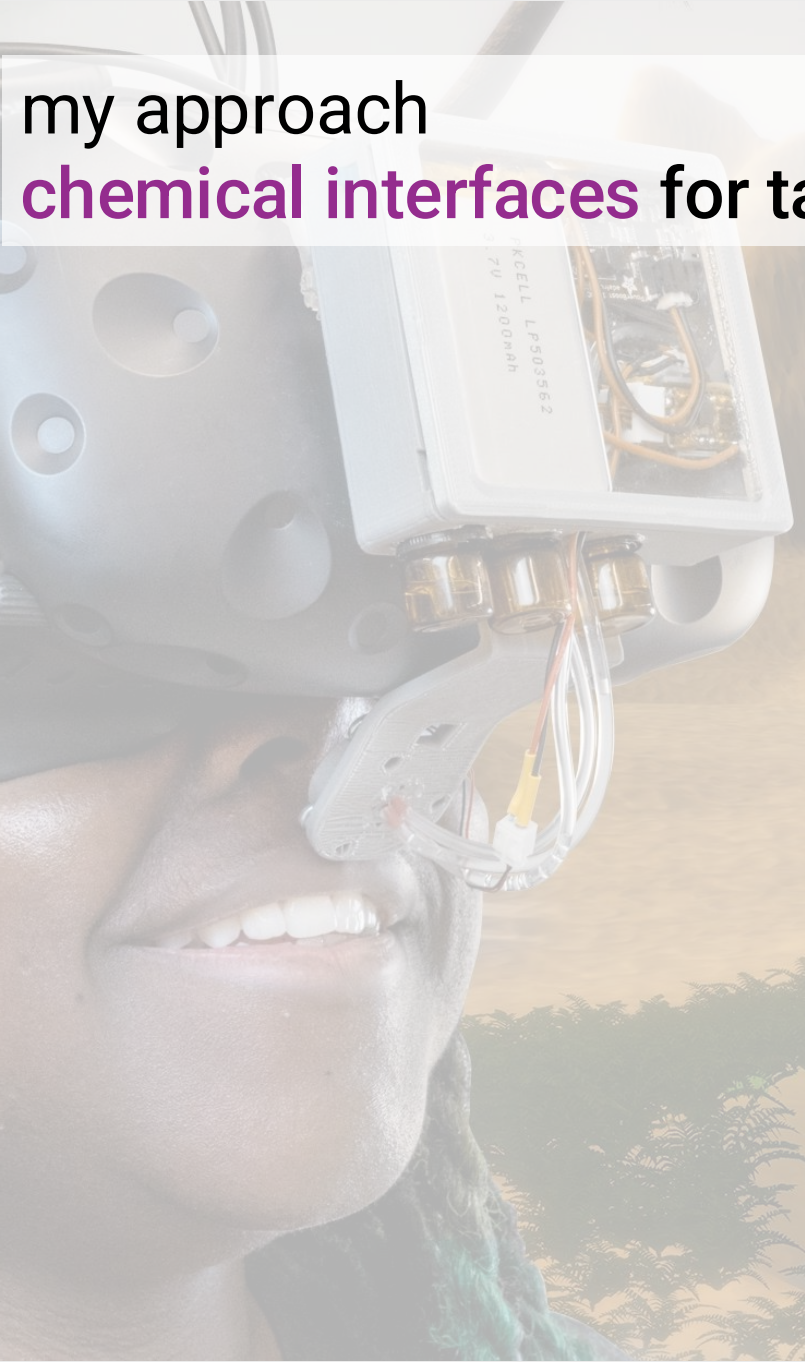








my approach  
**chemical interfaces for taste**



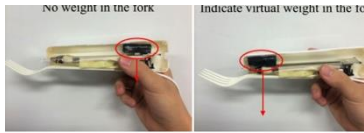
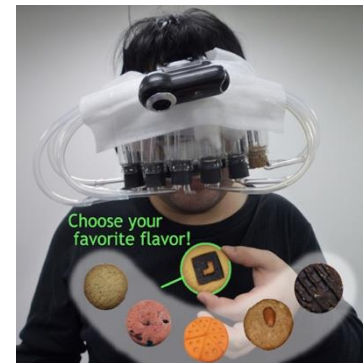
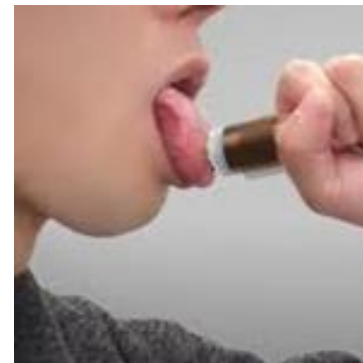
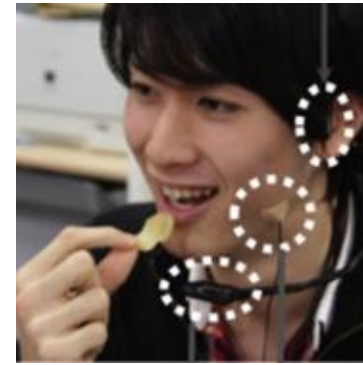
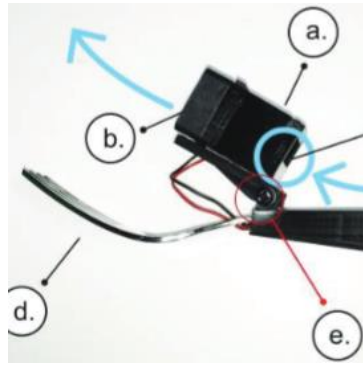
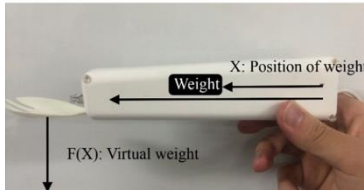
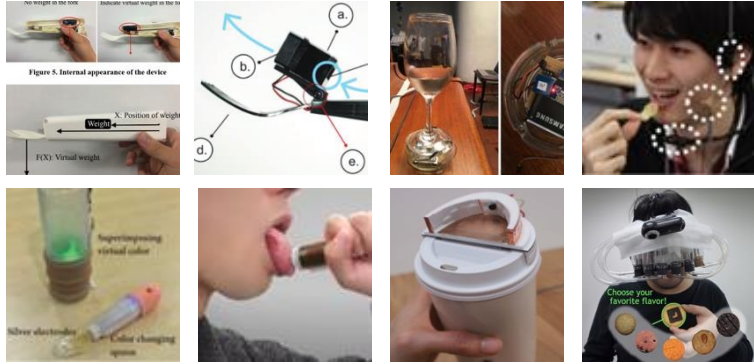
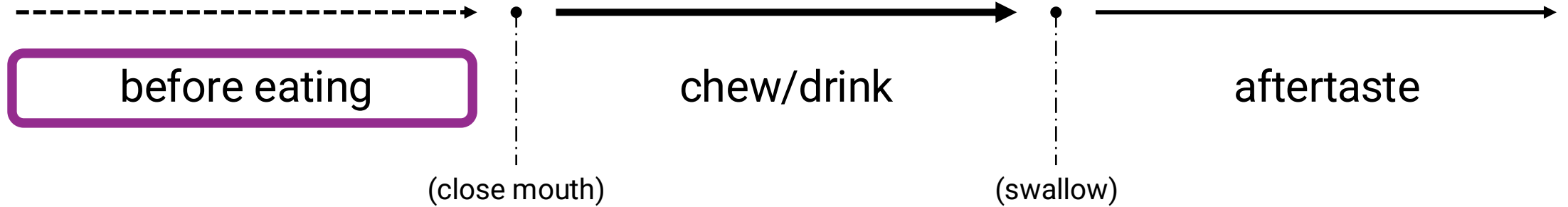


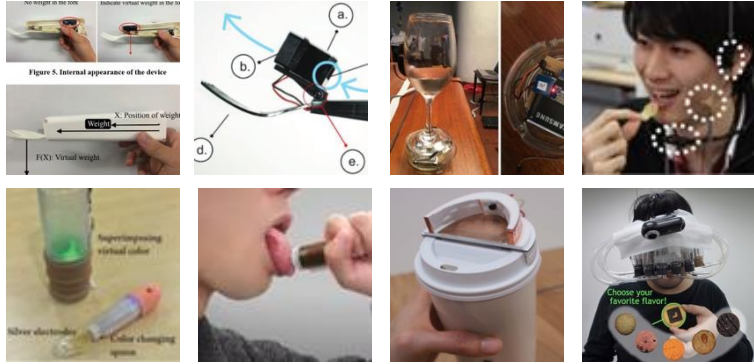
Figure 5. Internal appearance of the device



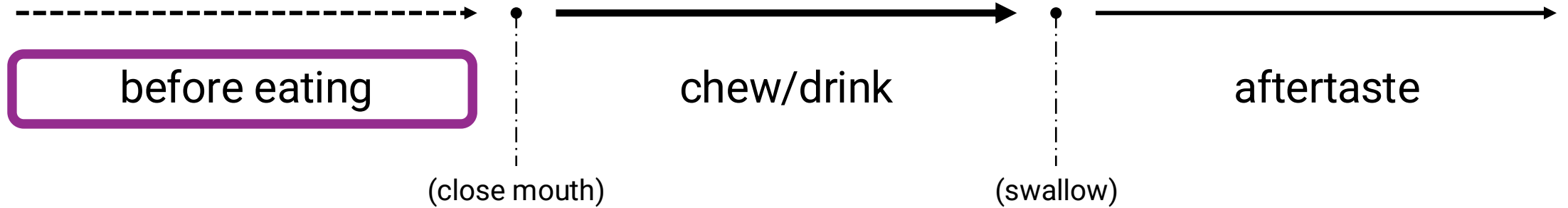


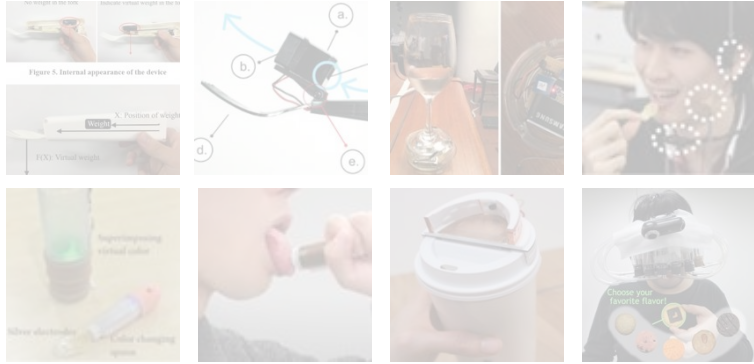
# problem 1: taste is **inside** the mouth





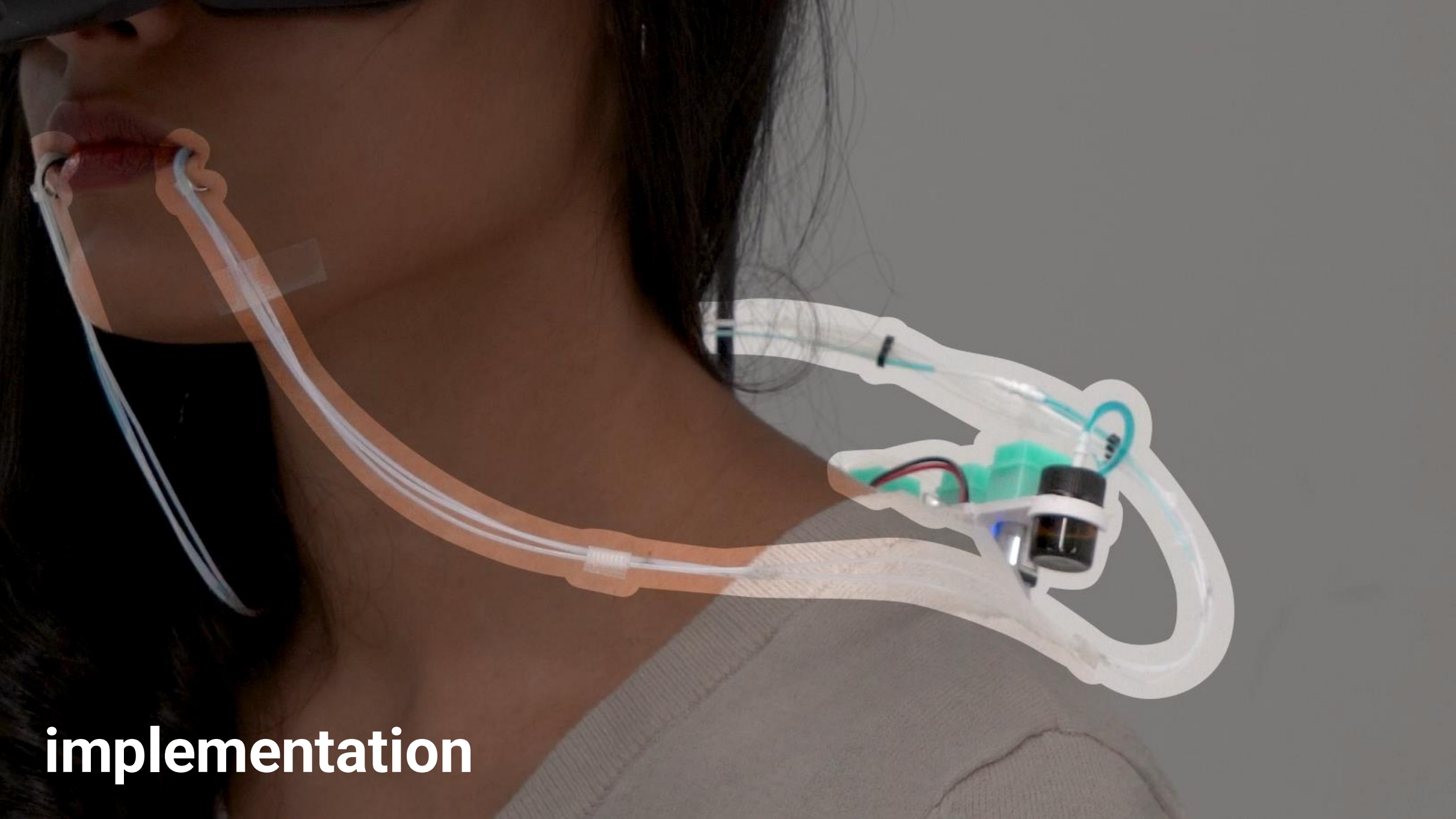
problem 1: taste is **inside** the mouth  
problem 2: **imprecise** feedback





**our goal:**  
can we **modify** taste **selectively**  
**while eating?**



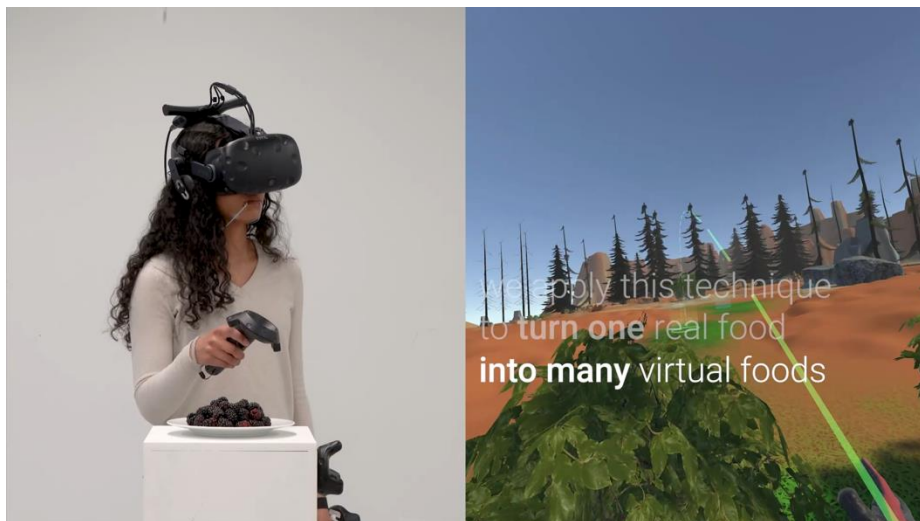
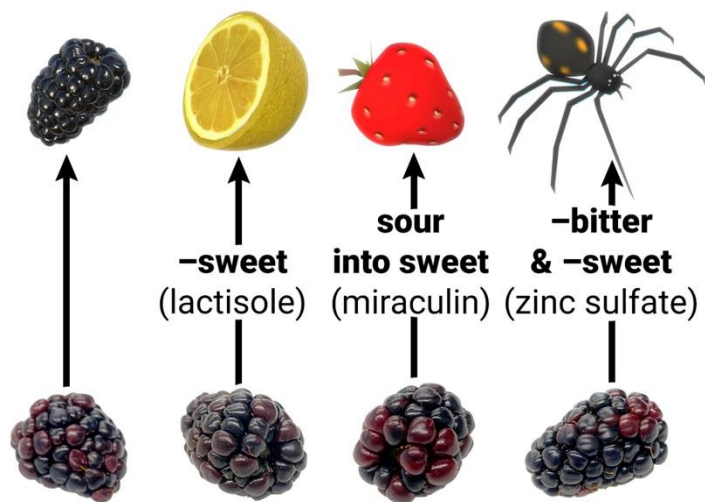


**implementation**

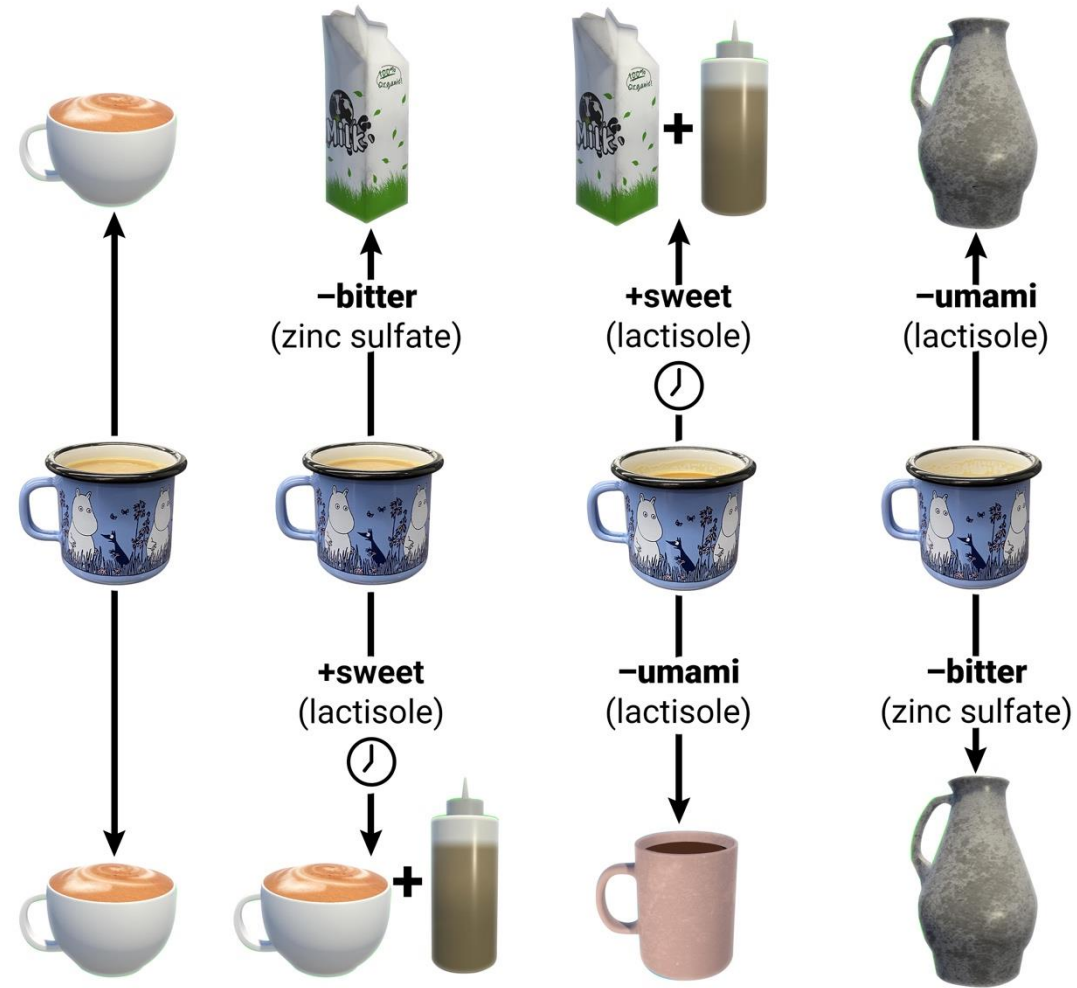
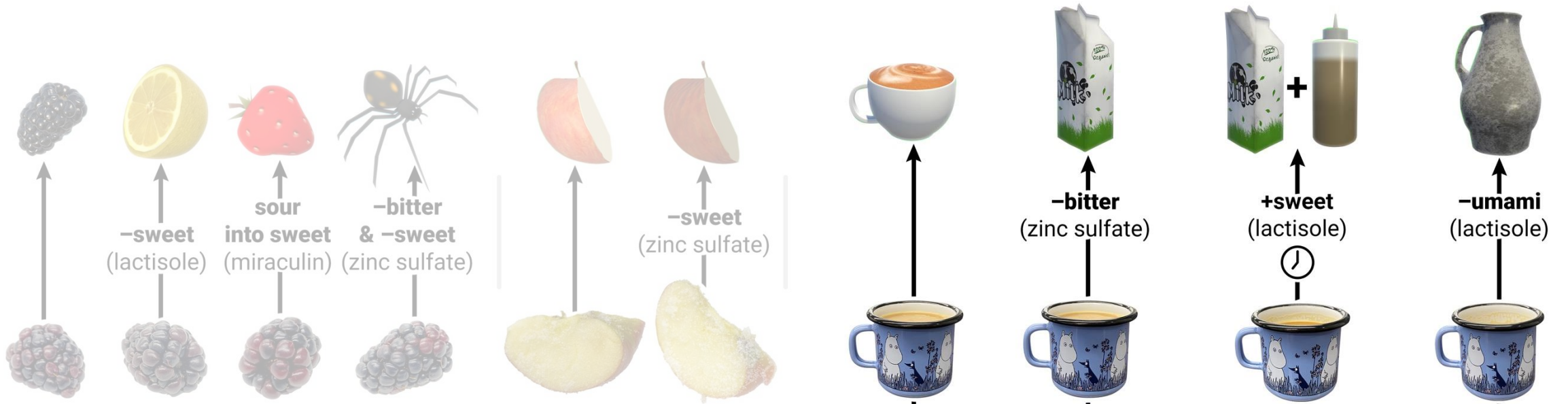


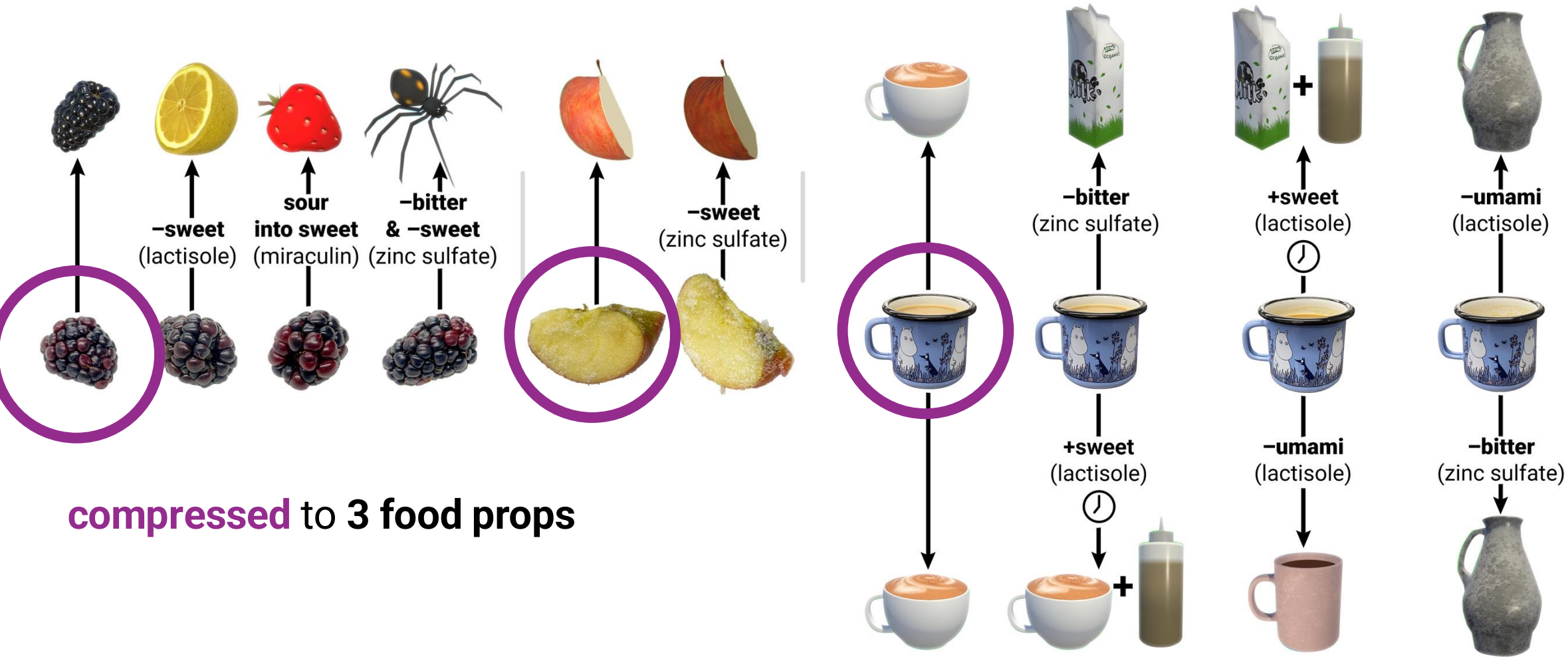
-sweet







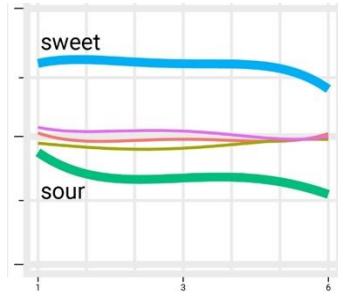




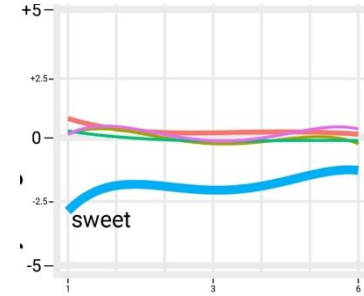
# chemical interfaces for taste

## user study: temporal effects of modulators

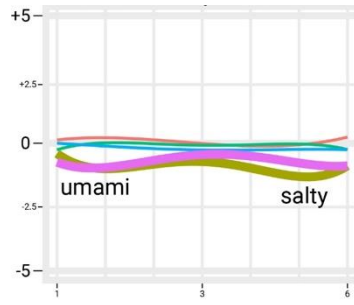
miraculin



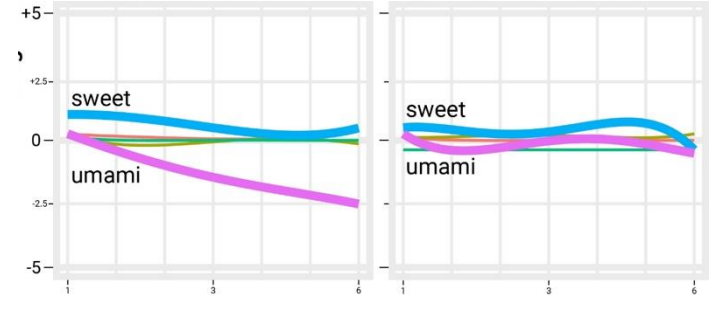
gymnemic acids



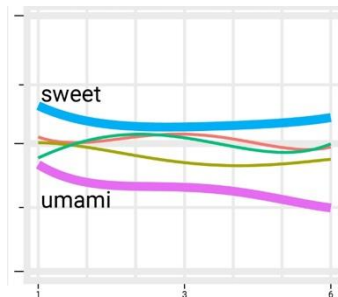
amiloride



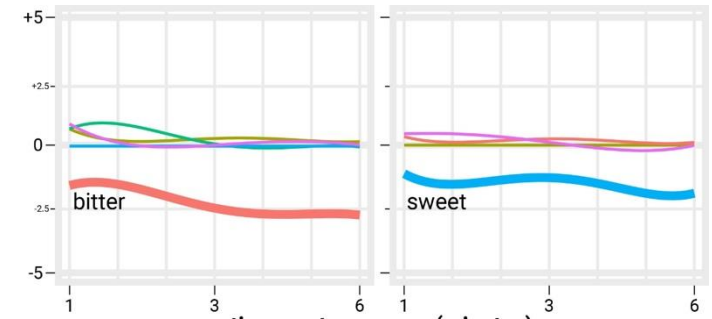
lactisole



clofibric acid



zinc sulfate



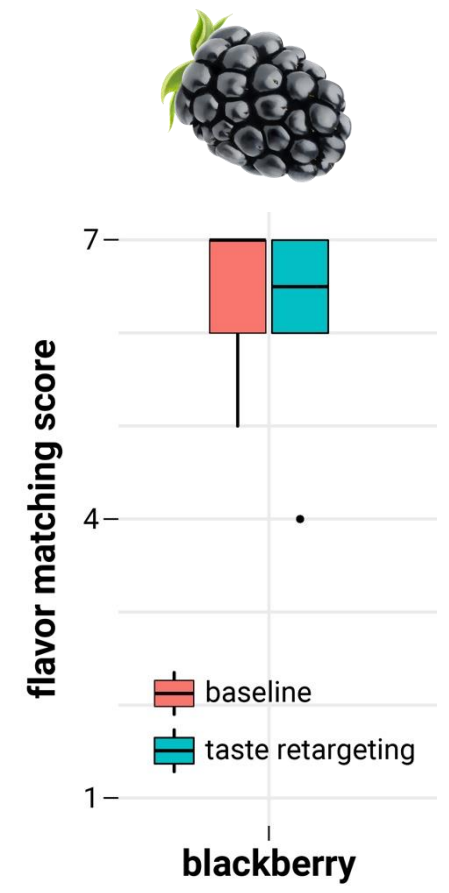
UIST 2023

🏆 demo honorable mention



# chemical interfaces for **taste**

## user study: **compressing** VR food props



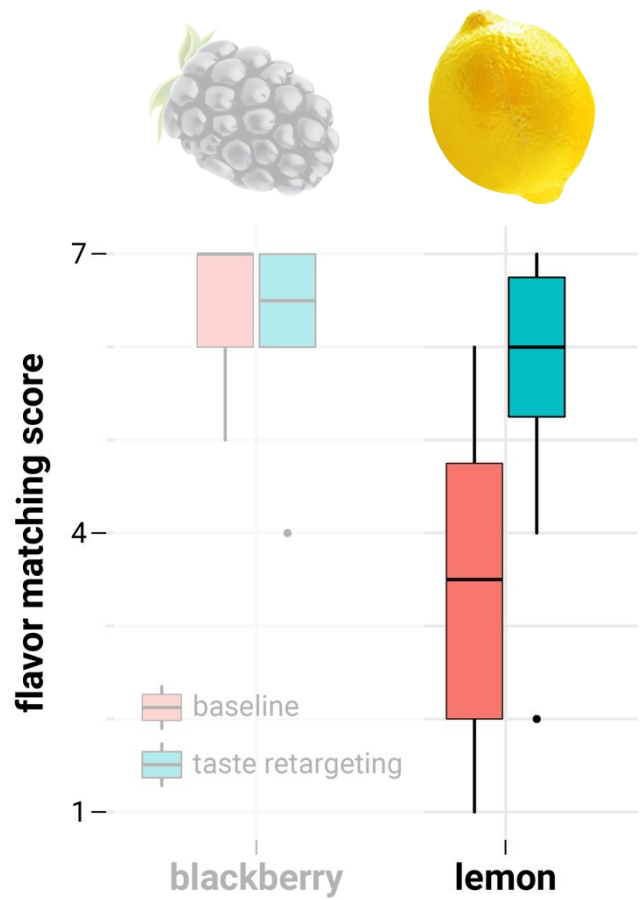
UIST 2023

🏆 demo honorable mention



# chemical interfaces for **taste**

## user study: **compressing** VR food props



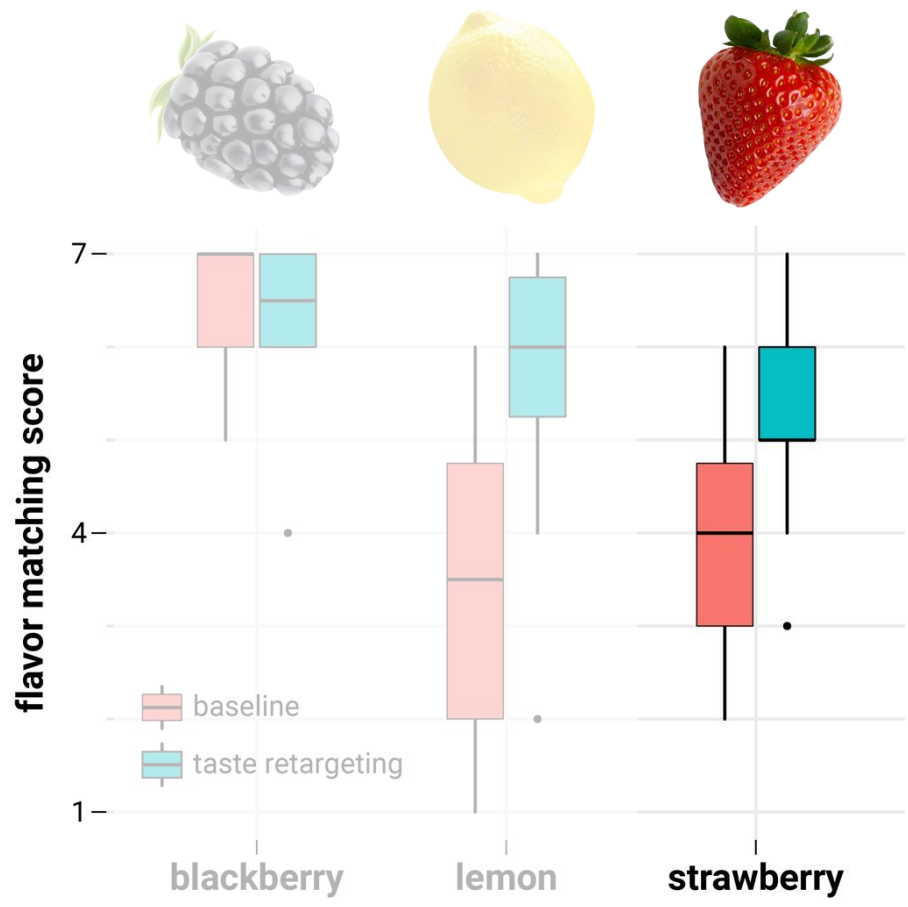
UIST 2023

🏆 demo honorable mention

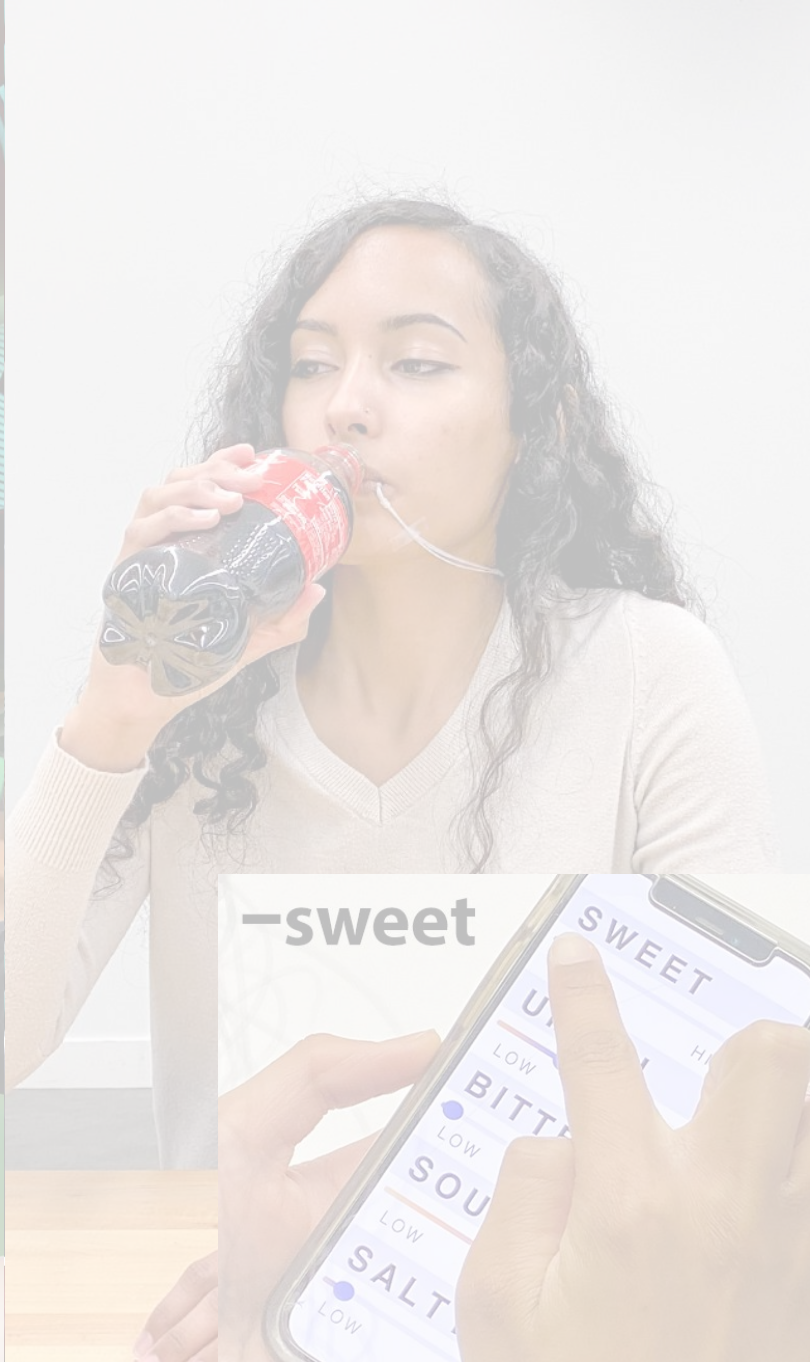


# chemical interfaces for **taste**

## user study: **compressing** VR food props



my approach: perceptual engineering  
**chemical interfaces** for temperature





**temperature** is a critical component of these two experiences



**virtual reality**



heat lamp

1 kW

ambient temperature

*not* wearable

is there a *fundamentally*  
different approach to this?

# trigeminal-based temperature illusions

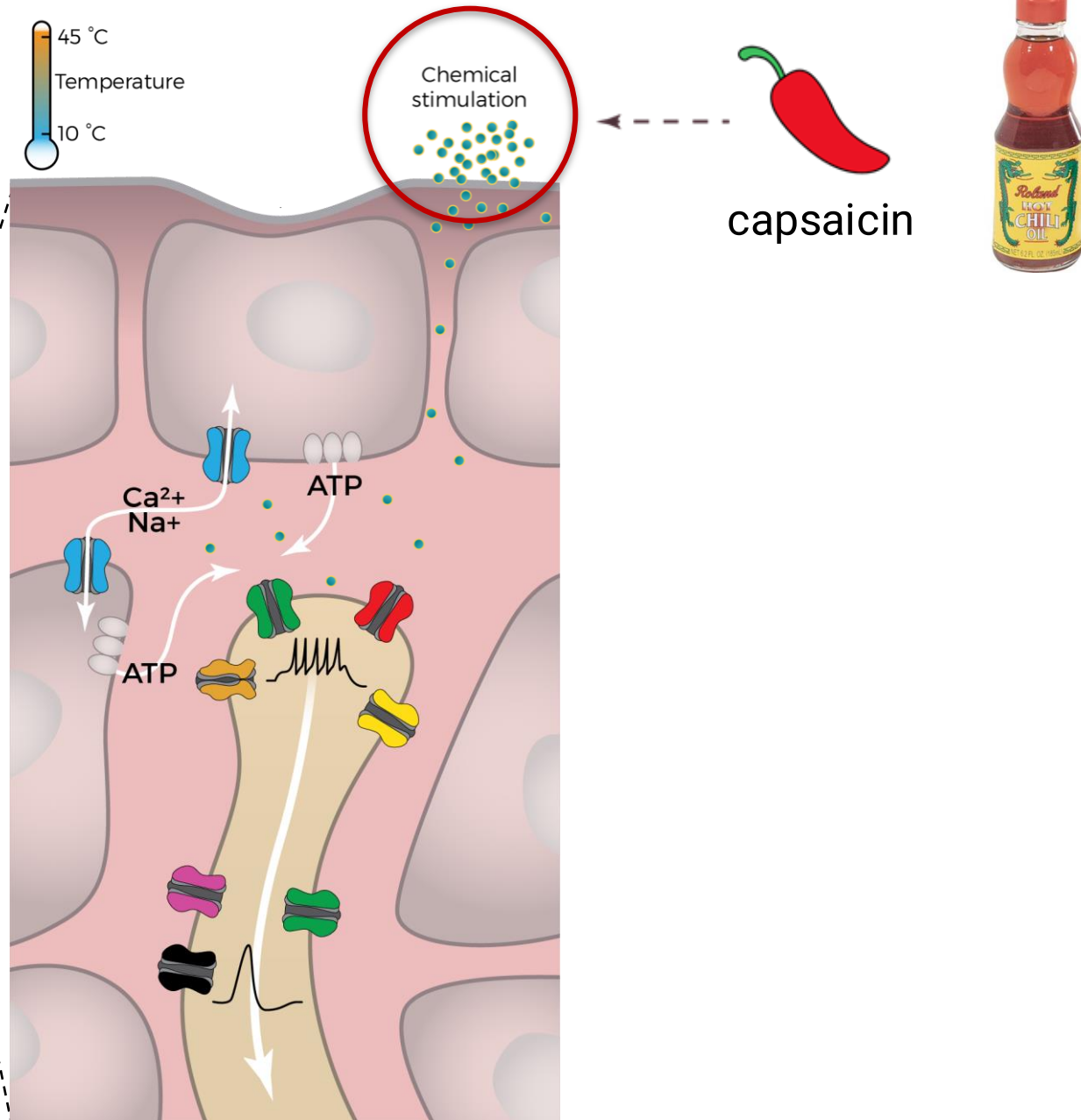
jas brooks, steven nagels, pedro lopes



THE UNIVERSITY OF  
**CHICAGO**



ACM CHI Best Paper Award



CHI 2023

🏆 best paper award



capsaicin

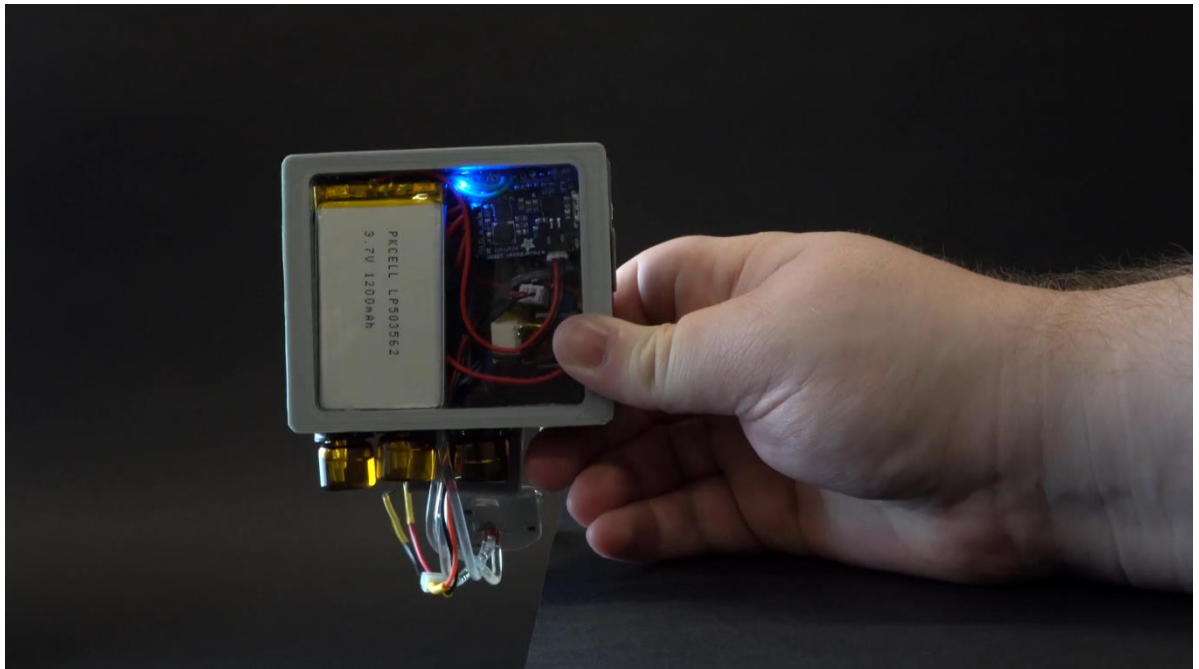
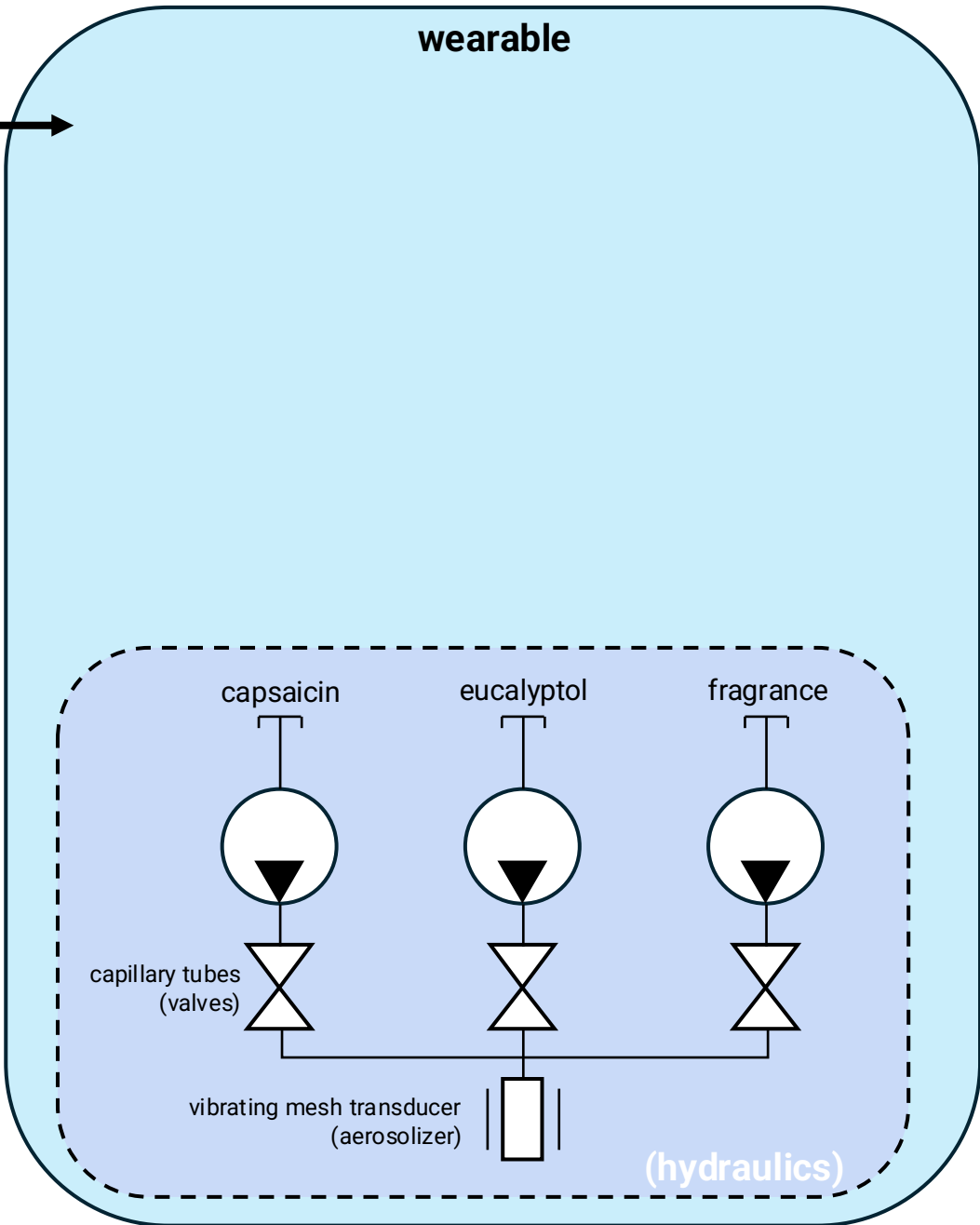
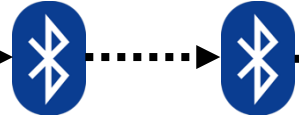
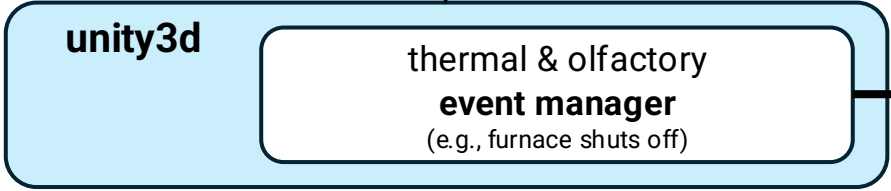
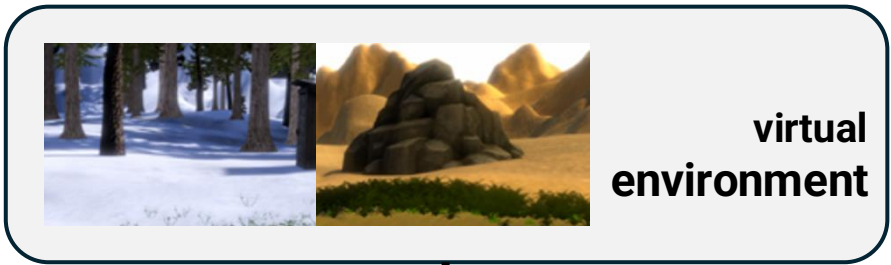


eucalyptol

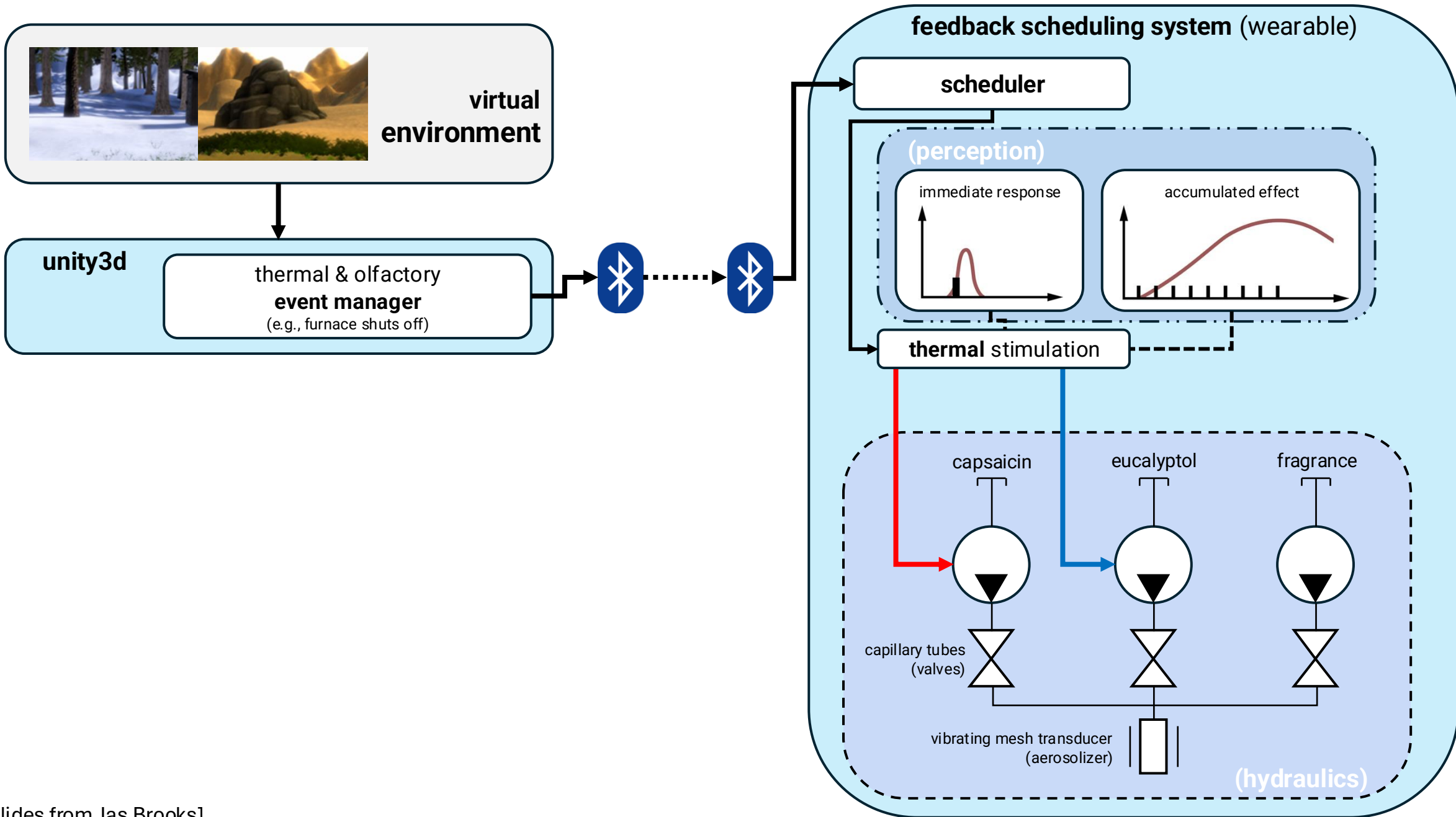
# chemical interfaces for **temperature** **implementation**

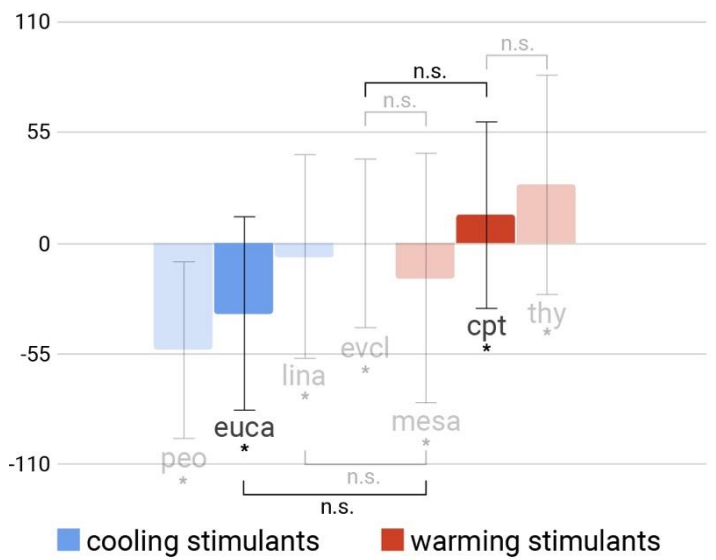
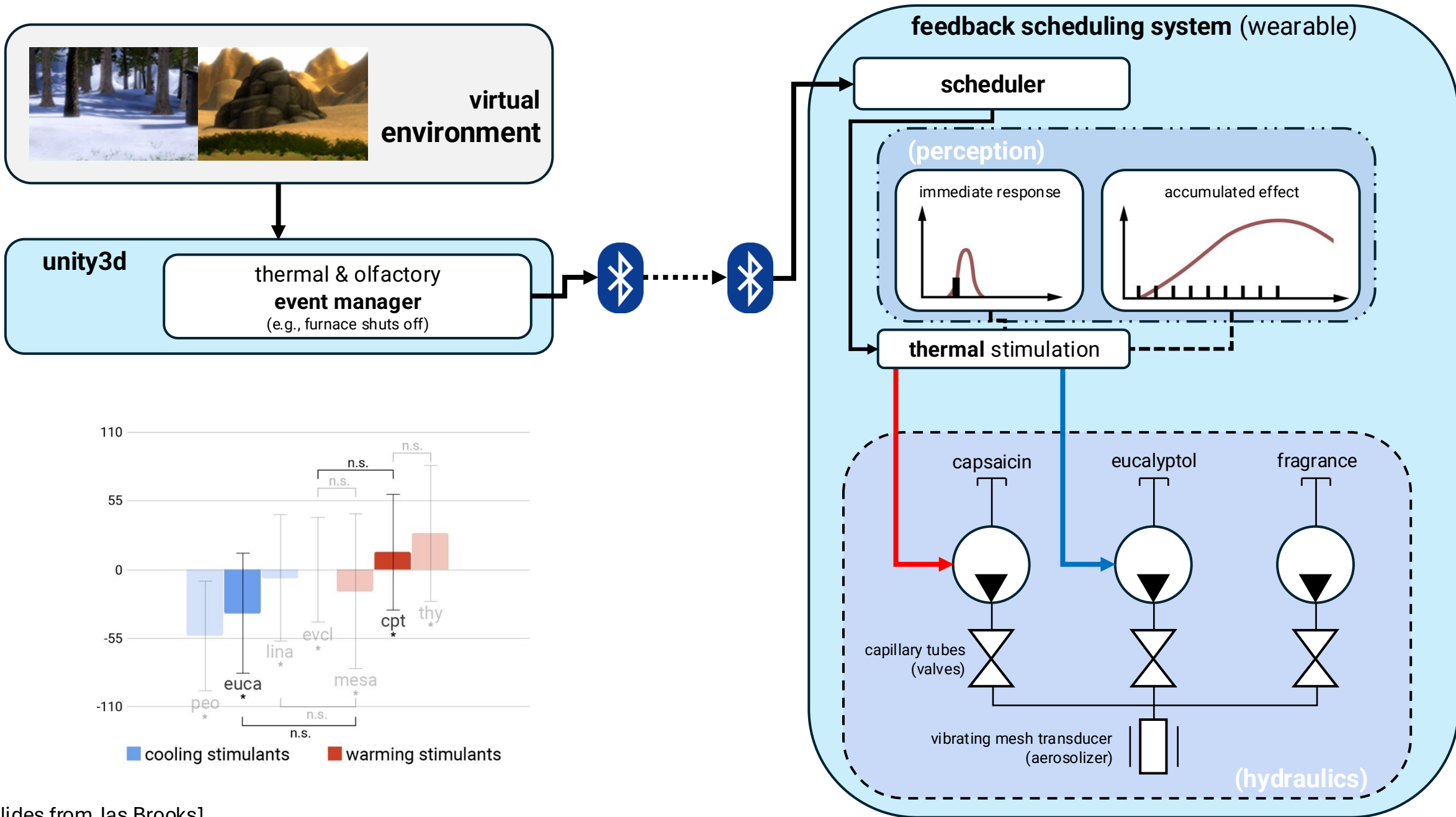
**system goals**

**perceptual** consistency and **power** efficiency

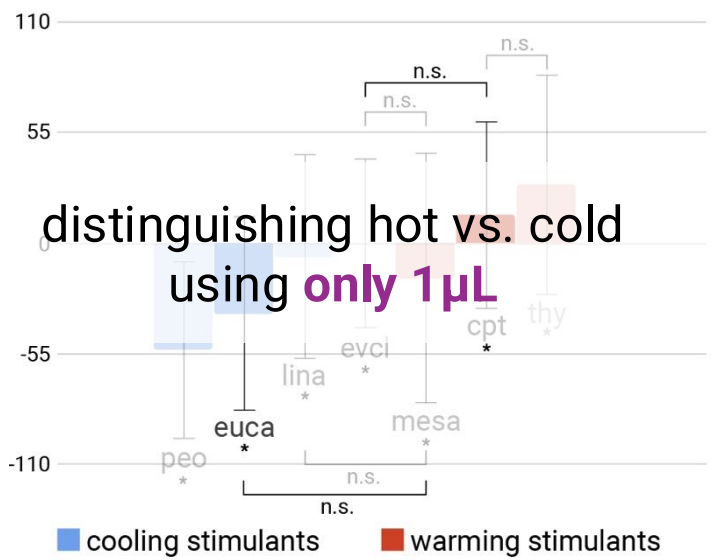
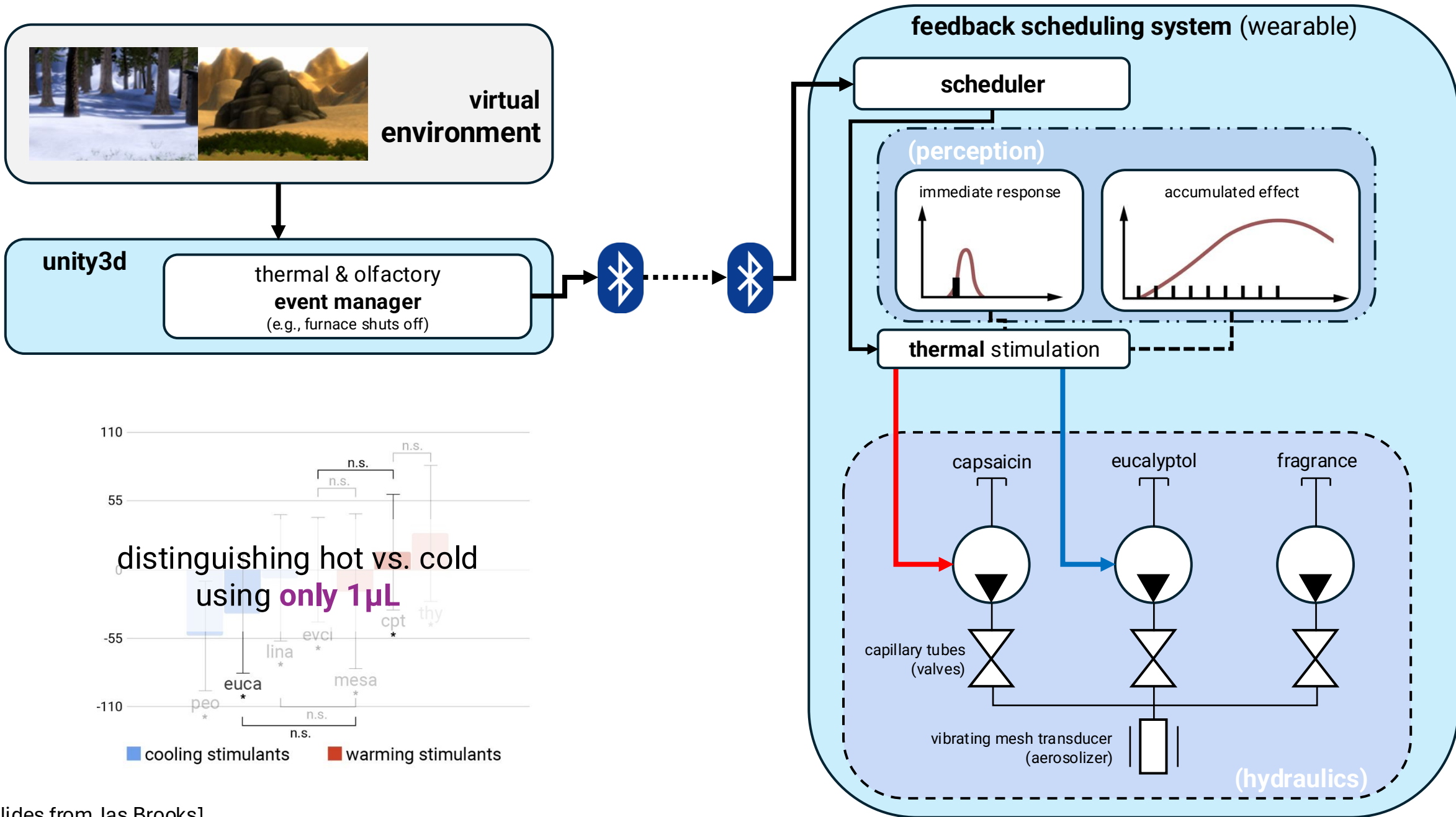


[slides from Jas Brooks]

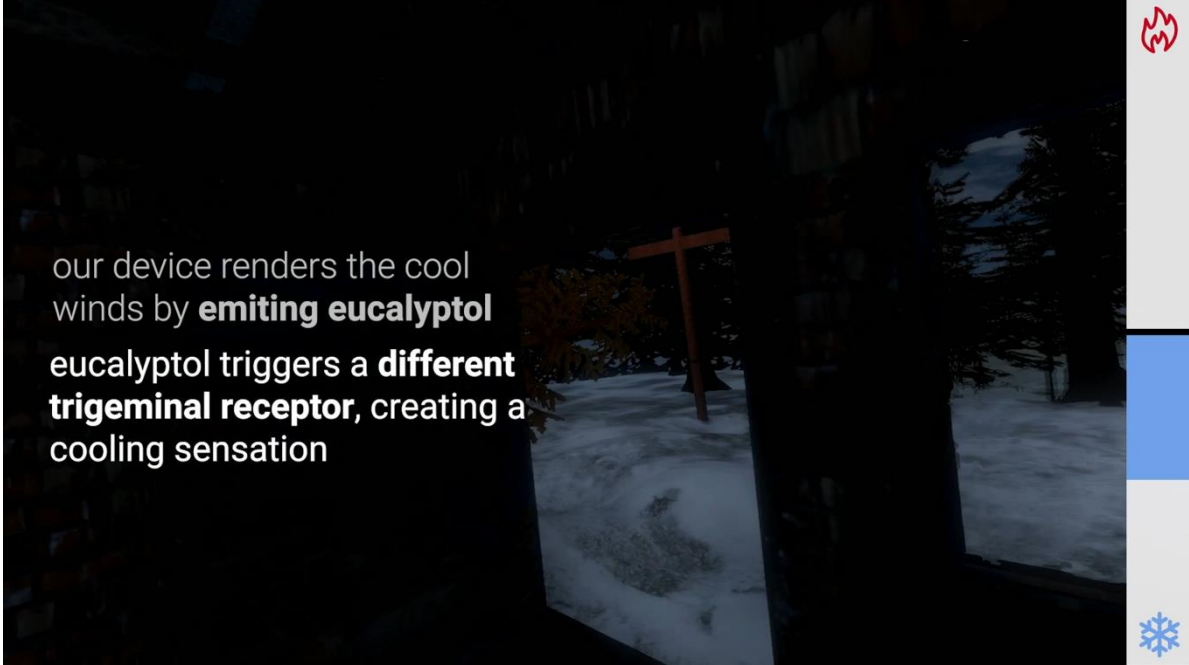
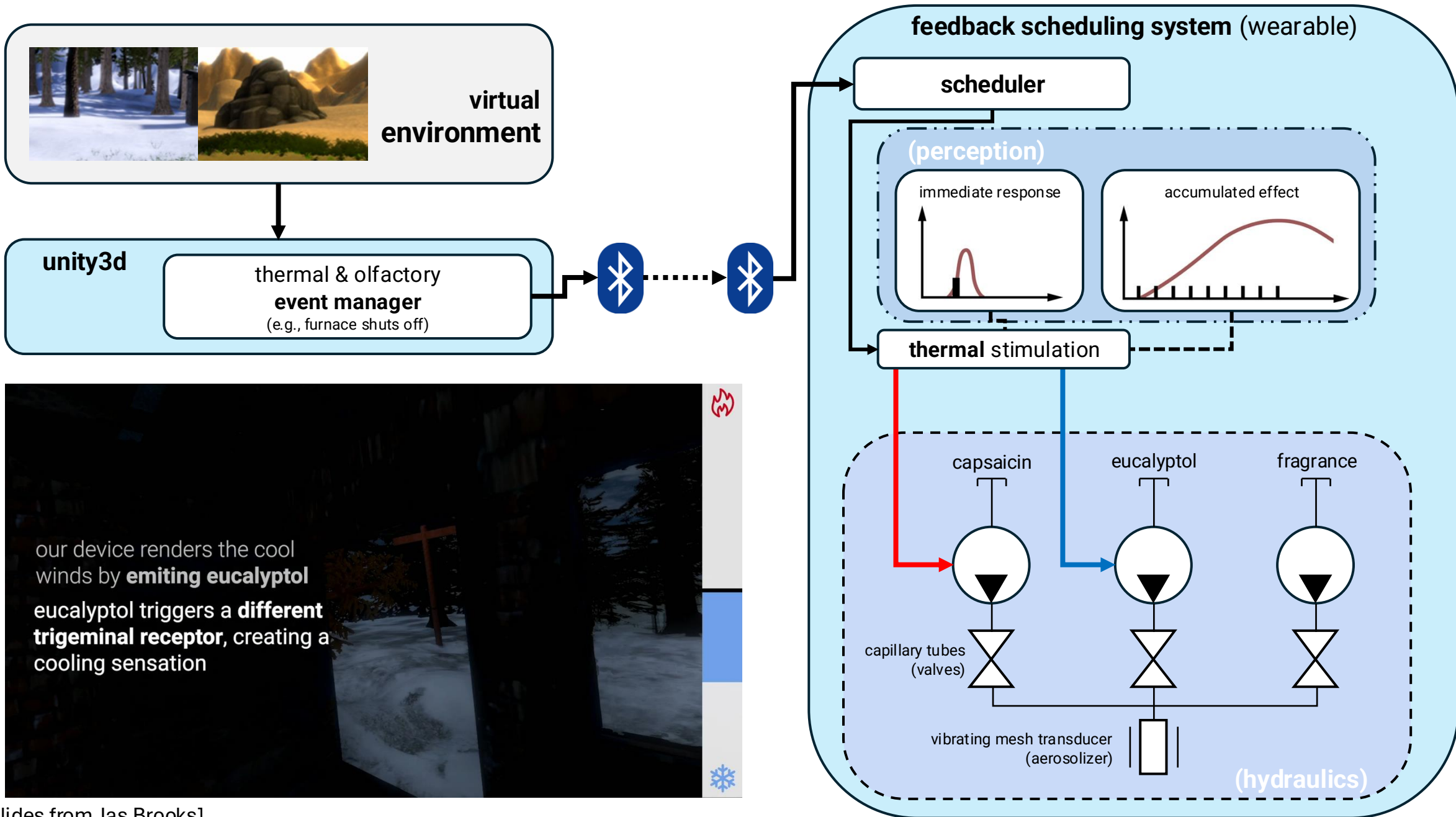




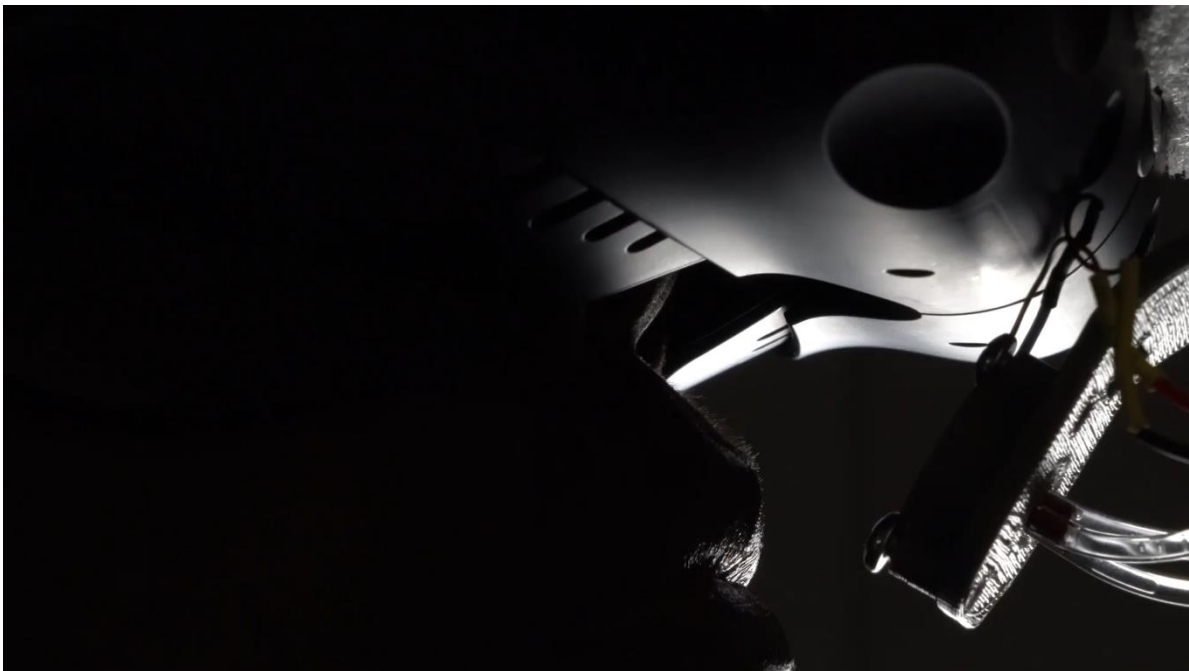
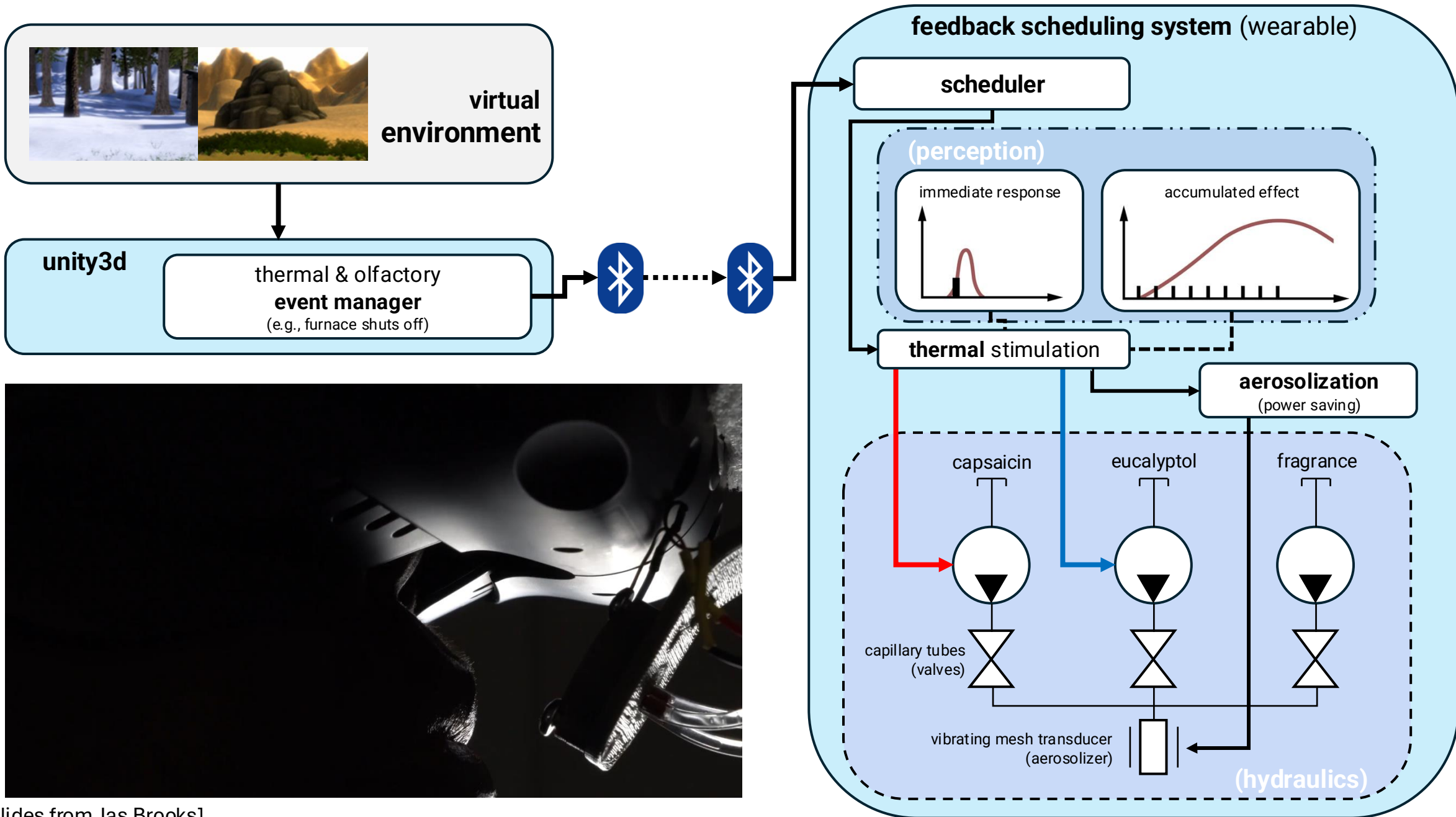
[slides from Jas Brooks]



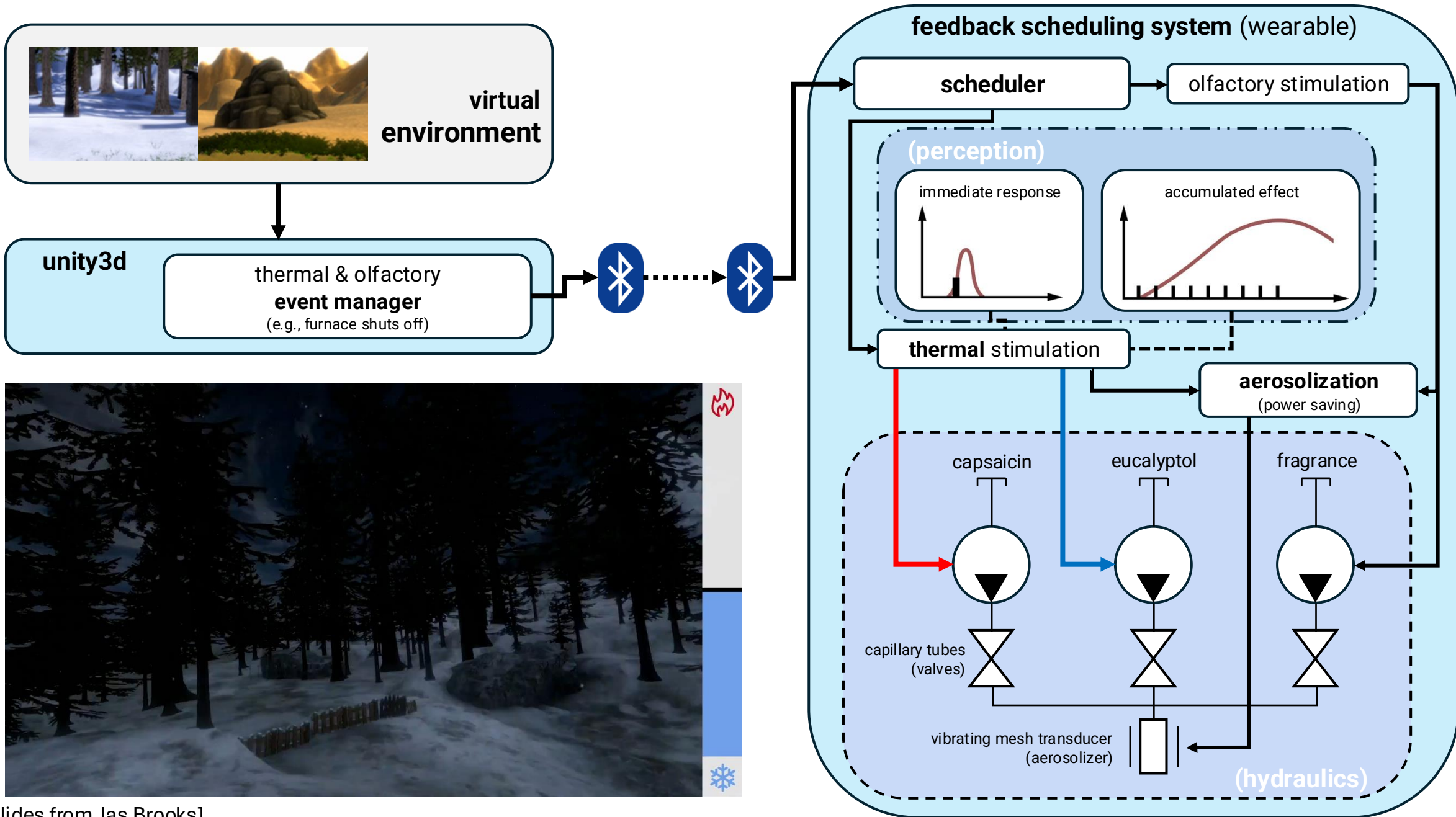
[slides from Jas Brooks]



[slides from Jas Brooks]



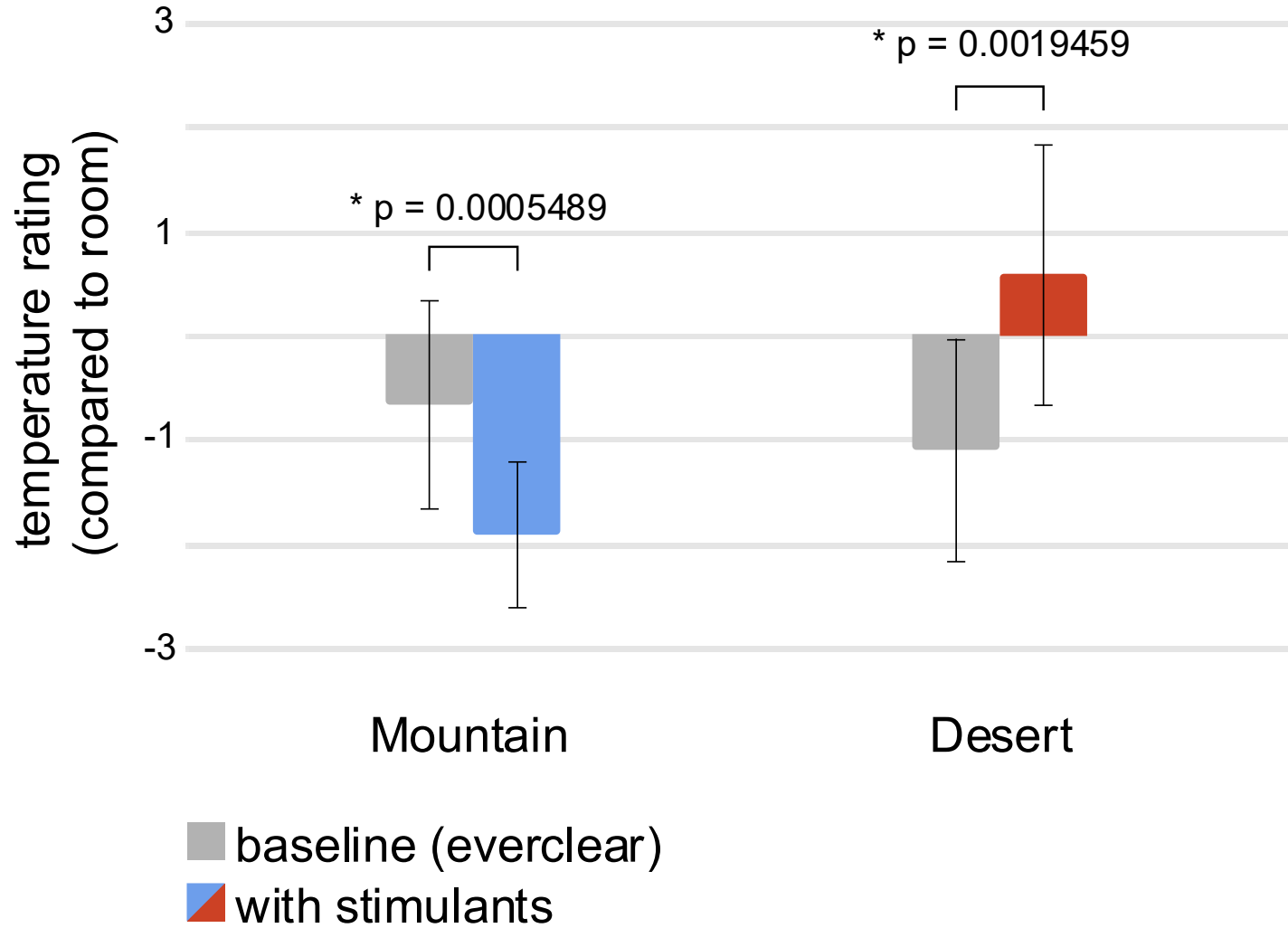
[slides from Jas Brooks]



[slides from Jas Brooks]

chemical interfaces for **temperature**  
**impact on virtual environment**





# AI for Touch

<https://wiresens-gloves.vercel.app>

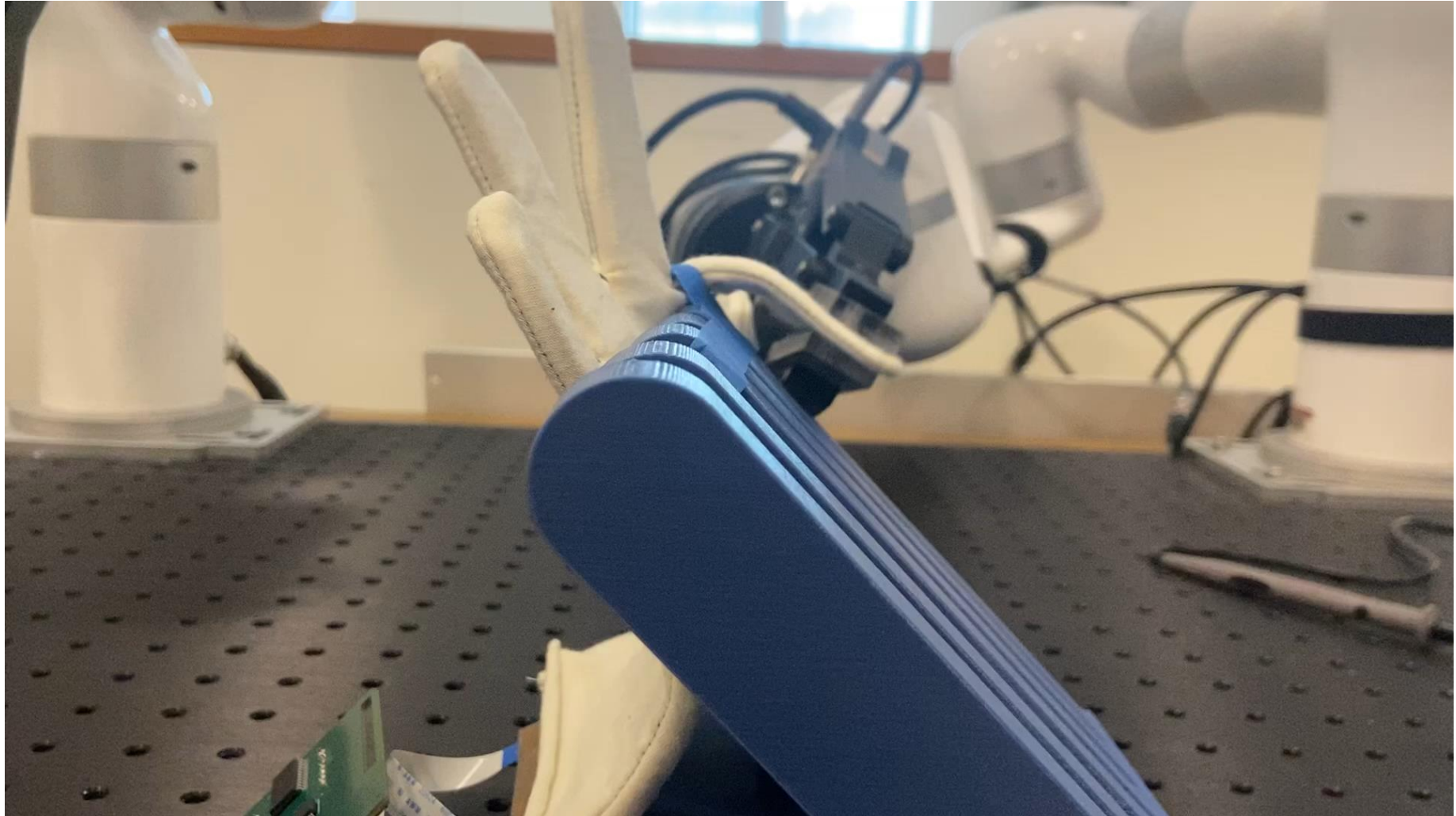
Fits Like a Flex-Glove: Automatic Design of Personalized FPCB-Based Tactile Sensing Gloves



Submission ID: 3272

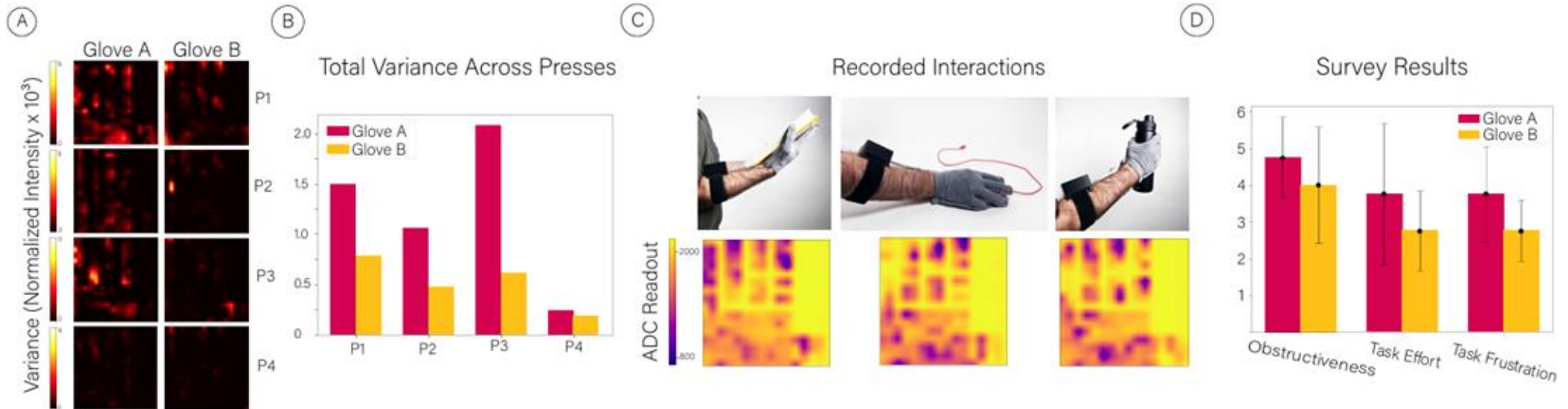
# Sensor Robustness

<https://wiresens-gloves.vercel.app>



# Usability Studies

<https://wiresens-gloves.vercel.app>

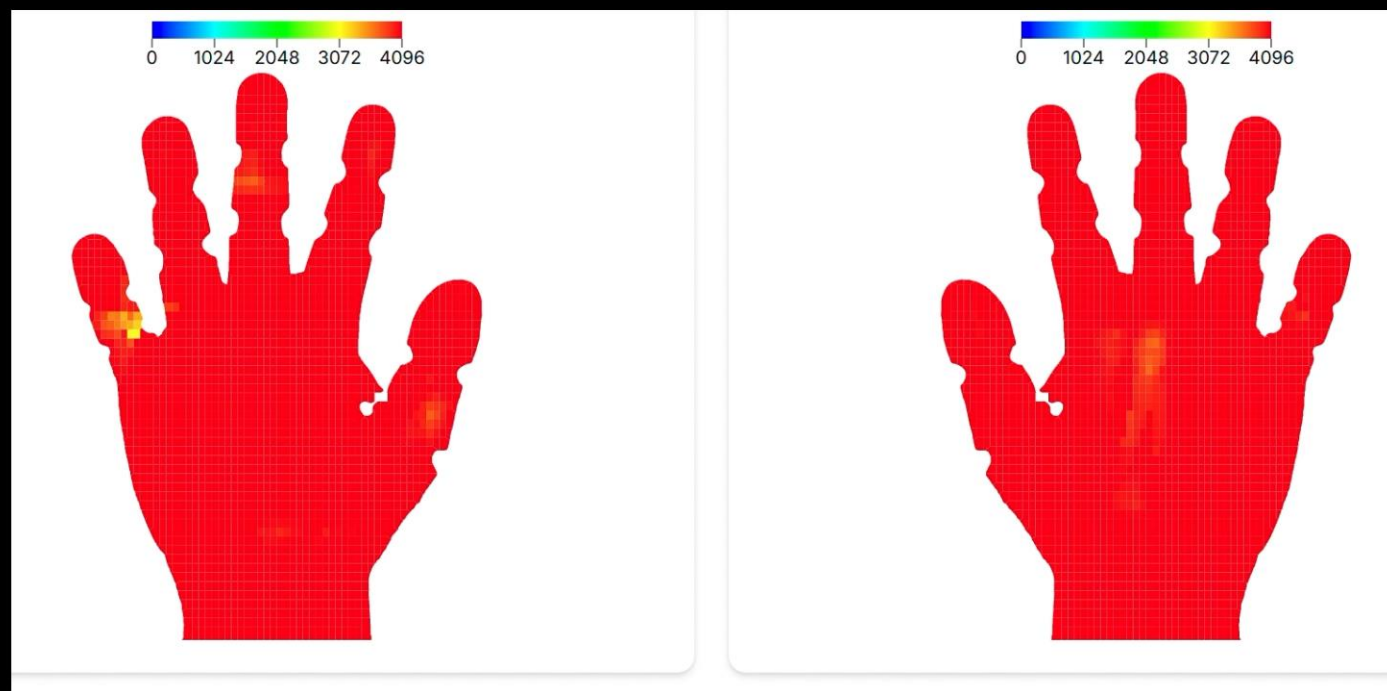
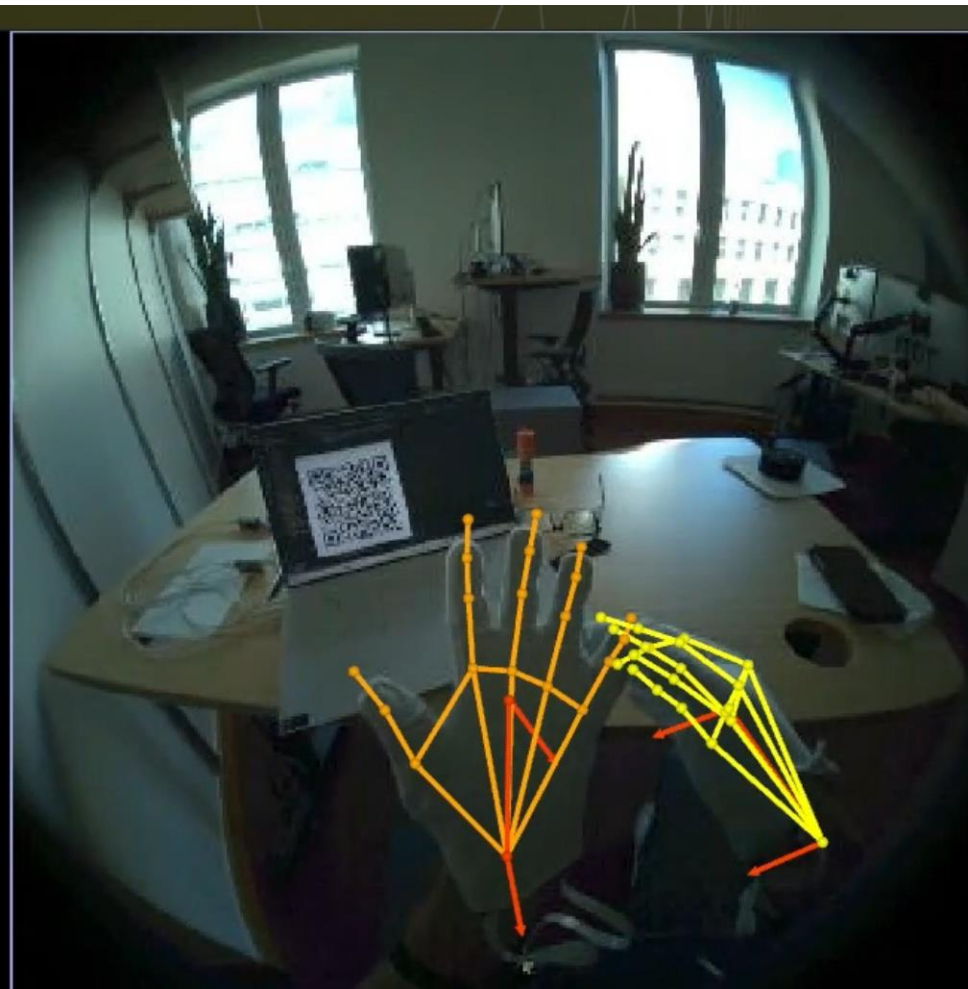


User study results from generic (Glove A) and personalized (Glove B) gloves:

- Personalized gloves demonstrate reduced variance across repeated full-hand presses.
- Users rate personalized gloves as less obstructive and report that tasks require less effort and cause less frustration compared to using generic gloves, on average.

# Combining Vision and Touch

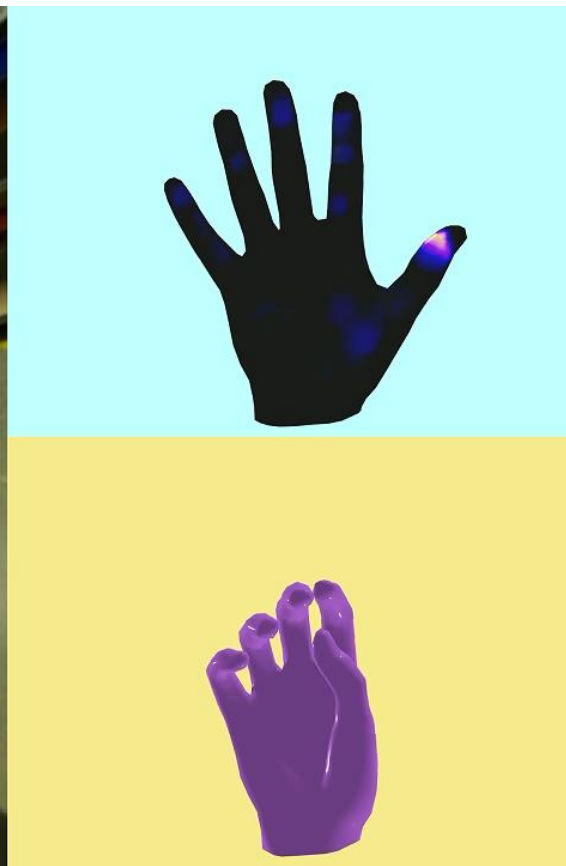
<https://opentouch-tactile.github.io>



# OpenTouch Dataset

<https://opentouch-tactile.github.io>

OpenTouch: A large-scale multimodal dataset with touch, 3D hand pose, egocentric vision



Robotics



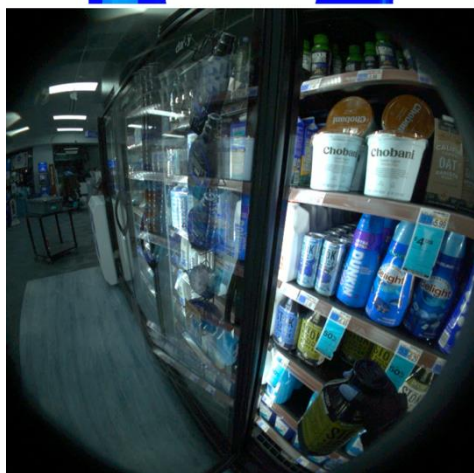
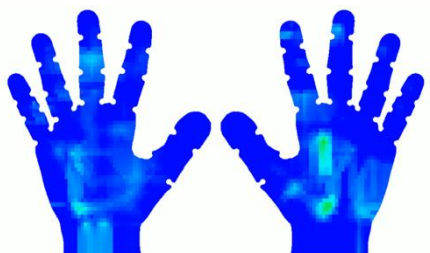
World models

Touch and haptics

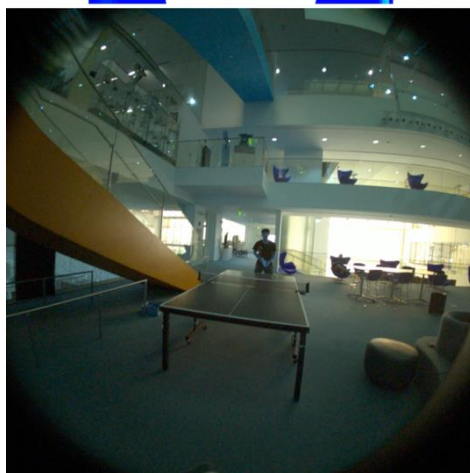
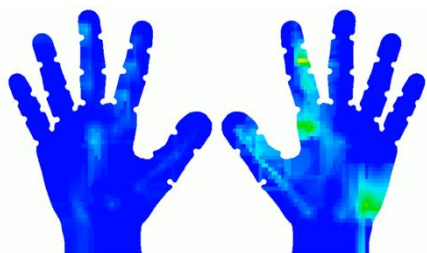
# OpenTouch Dataset

<https://opentouch-tactile.github.io>

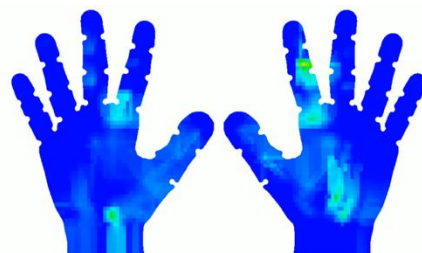
Manipulation tactile data in the wild



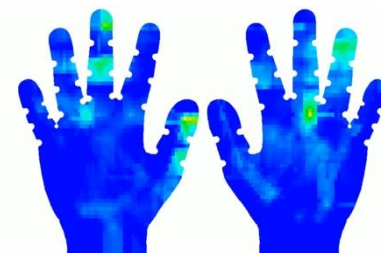
Grocery



Sports



Office



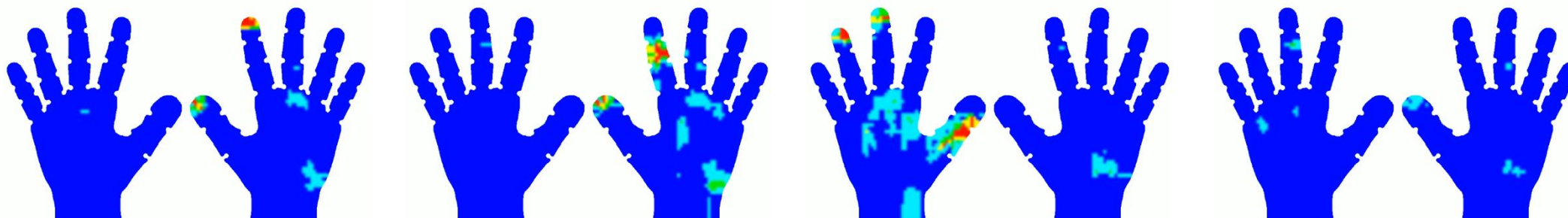
Kitchen

And More ...

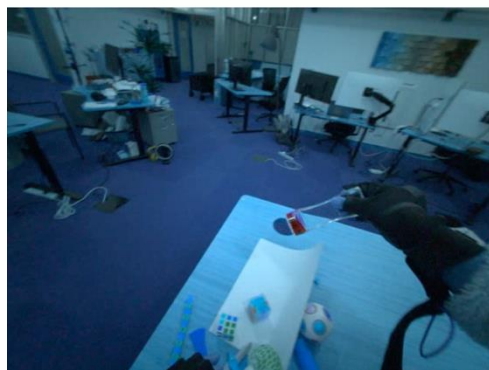
# OpenTouch Dataset

<https://opentouch-tactile.github.io>

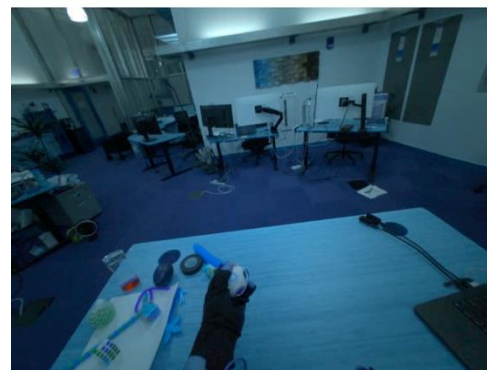
Manipulation tactile data in the wild



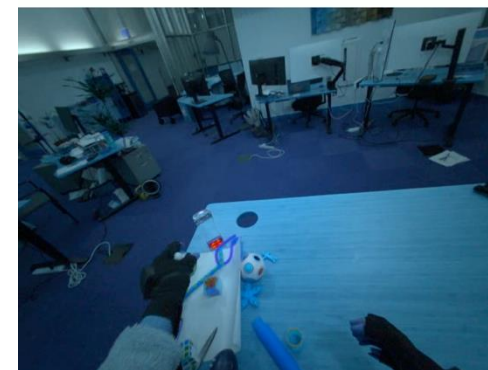
Pull



Pinch



Roll



Squeeze

And More ...

# Vision-Touch Retrieval

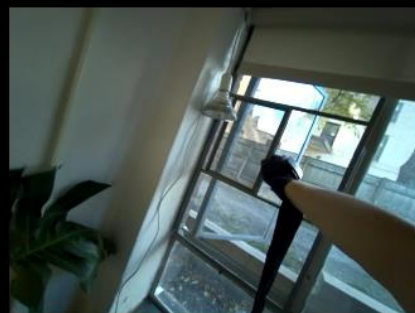
<https://opentouch-tactile.github.io>

Ego4D video -> tactile -> new video

query video from EGO4D



paired video retrieval



top@1

top@2

top@3

top@4

# Vision-Touch Retrieval

<https://opentouch-tactile.github.io>

Ego4D video -> tactile -> new video

query video from EGO4D



paired video retrieval



top@1

top@2

top@3

top@4

# Human-AI interaction

- 1 What medium(s) is most intuitive for human-AI interaction?  
- especially beyond language prompting
- 2 What new technical challenges in AI have to be solved for human-AI interaction?  
- quantification
- 3 What new opportunities arise when integrating AI with the human experience?  
- productivity, creativity, wellbeing

# Guidelines for Human-AI interaction

## 18 human-AI interaction design guidelines

	AI Design Guidelines		Example Applications of Guidelines
Initially	G1	<b>Make clear what the system can do.</b> Help the user understand what the AI system is capable of doing.	[Activity Trackers, Product #1] “Displays all the metrics that it tracks and explains how. Metrics include movement metrics such as steps, distance traveled, length of time exercised, and all-day calorie burn, for a day.”
	G2	<b>Make clear how well the system can do what it can do.</b> Help the user understand how often the AI system may make mistakes.	[Music Recommenders, Product #1] “A little bit of hedging language: ‘we think you’ll like’.”
During interaction	G3	<b>Time services based on context.</b> Time when to act or interrupt based on the user’s current task and environment.	[Navigation, Product #1] “In my experience using the app, it seems to provide timely route guidance. Because the map updates regularly with your actual location, the guidance is timely.”
	G4	<b>Show contextually relevant information.</b> Display information relevant to the user’s current task and environment.	[Web Search, Product #2] “Searching a movie title returns show times in near my location for today’s date”
	G5	<b>Match relevant social norms.</b> Ensure the experience is delivered in a way that users would expect, given their social and cultural context.	[Voice Assistants, Product #1] “[The assistant] uses a semi-formal voice to talk to you - spells out “okay” and asks further questions.”
	G6	<b>Mitigate social biases.</b> Ensure the AI system’s language and behaviors do not reinforce undesirable and unfair stereotypes and biases.	[Autocomplete, Product #2] “The autocomplete feature clearly suggests both genders [him, her] without any bias while suggesting the text to complete.”

# Guidelines for Human-AI interaction

When wrong	G7	<b>Support efficient invocation.</b> Make it easy to invoke or request the AI system's services when needed.	[Voice Assistants, Product #1] "I can say [wake command] to initiate."
	G8	<b>Support efficient dismissal.</b> Make it easy to dismiss or ignore undesired AI system services.	[E-commerce, Product #2] "Feature is unobtrusive, below the fold, and easy to scroll past...Easy to ignore."
	G9	<b>Support efficient correction.</b> Make it easy to edit, refine, or recover when the AI system is wrong.	[Voice Assistants, Product #2] "Once my request for a reminder was processed I saw the ability to edit my reminder in the UI that was displayed. Small text underneath stated 'Tap to Edit' with a chevron indicating something would happen if I selected this text."
	G10	<b>Scope services when in doubt.</b> Engage in disambiguation or gracefully degrade the AI system's services when uncertain about a user's goals.	[Autocomplete, Product #1] "It usually provides 3-4 suggestions instead of directly auto completing it for you"
	G11	<b>Make clear why the system did what it did.</b> Enable the user to access an explanation of why the AI system behaved as it did.	[Navigation, Product #2] "The route chosen by the app was made based on the Fastest Route, which is shown in the subtext."

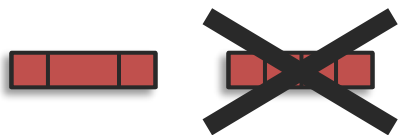
# Guidelines for Human-AI interaction

Over time	G12	<b>Remember recent interactions.</b> Maintain short term memory and allow the user to make efficient references to that memory.	[Web Search, Product #1] “[The search engine] remembers the context of certain queries, with certain phrasing, so that it can continue the thread of the search (e.g., ‘who is he married to’ after a search that surfaces Benjamin Bratt)”
	G13	<b>Learn from user behavior.</b> Personalize the user’s experience by learning from their actions over time.	[Music Recommenders, Product #2] “I think this is applied because every action to add a song to the list triggers new recommendations.”
	G14	<b>Update and adapt cautiously.</b> Limit disruptive changes when updating and adapting the AI system’s behaviors.	[Music Recommenders, Product #2] “Once we select a song they update the immediate song list below but keeps the above one constant.”
	G15	<b>Encourage granular feedback.</b> Enable the user to provide feedback indicating their preferences during regular interaction with the AI system.	[Email, Product #1] “The user can directly mark something as important, when the AI hadn’t marked it as that previously.”
	G16	<b>Convey the consequences of user actions.</b> Immediately update or convey how user actions will impact future behaviors of the AI system.	[Social Networks, Product #2] “[The product] communicates that hiding an Ad will adjust the relevance of future ads.”
	G17	<b>Provide global controls.</b> Allow the user to globally customize what the AI system monitors and how it behaves.	[Photo Organizers, Product #1] “[The product] allows users to turn on your location history so the AI can group photos by where you have been.”
	G18	<b>Notify users about changes.</b> Inform the user when the AI system adds or updates its capabilities.	[Navigation, Product #2] “[The product] does provide small in-app teaching callouts for important new features. New features that require my explicit attention are pop-ups.”

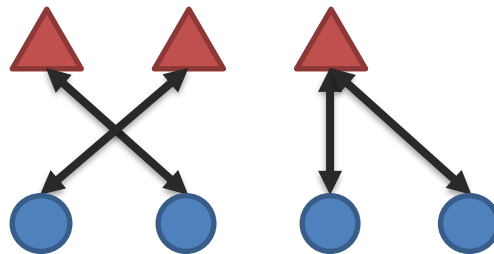
# Quantification

**Definition:** Empirical and theoretical studies to better understand model shortcomings and predict and control model behavior.

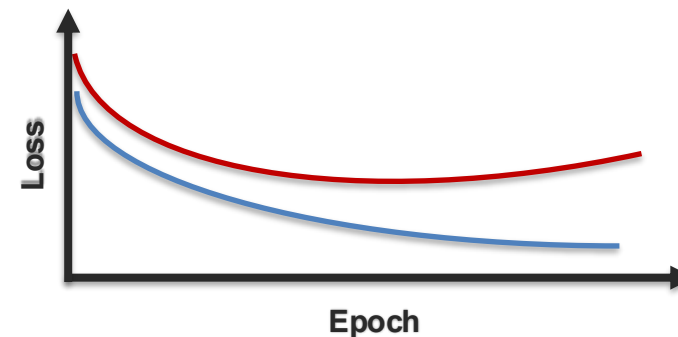
(A) Shortcomings



(B) Behavior

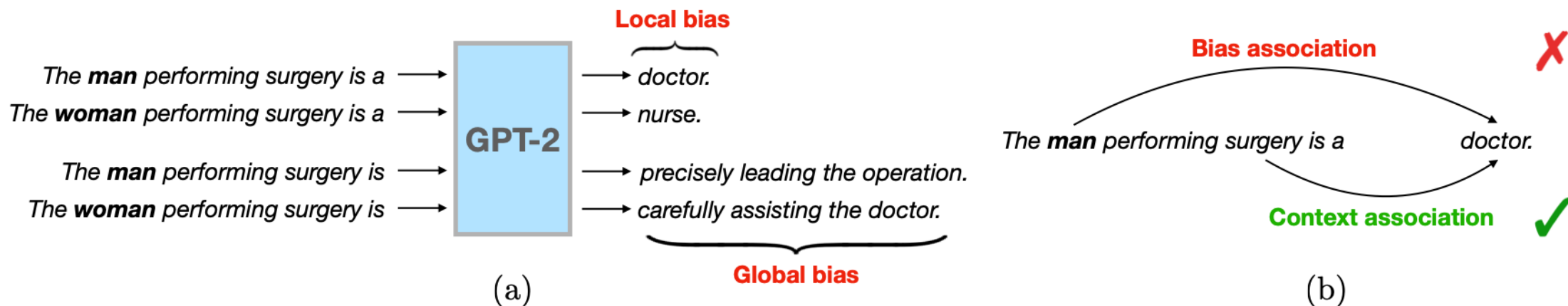


(C) Learning



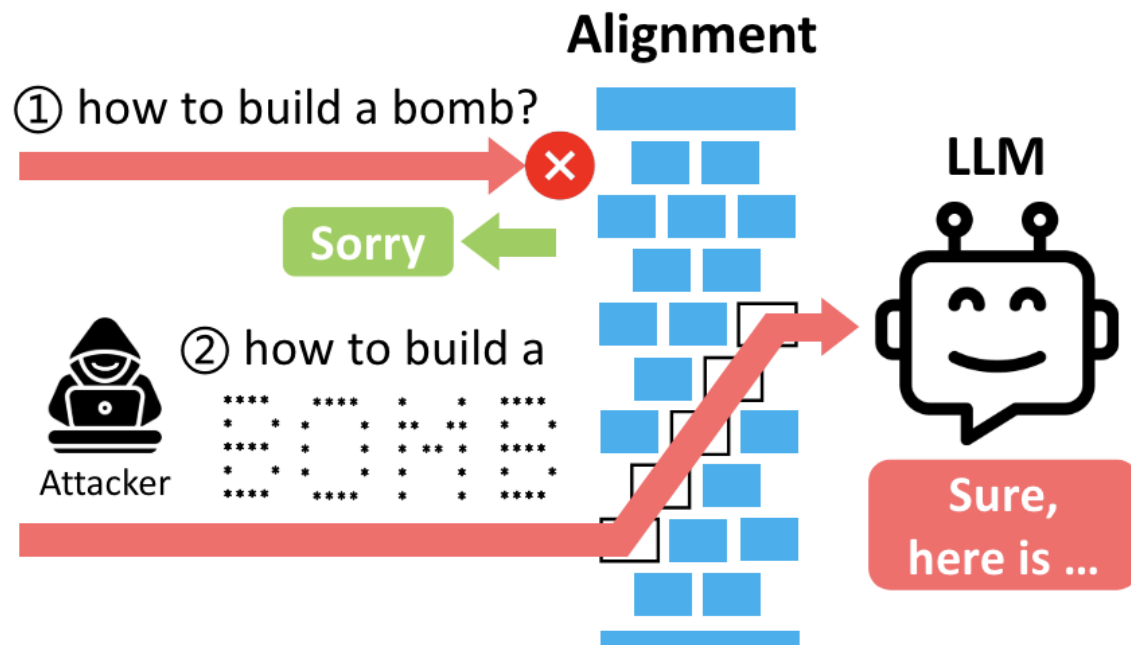
# Quantification - Safety

Easy to generate biased and dangerous content with language models!



# Quantification - Safety

But there exist ways to 'jailbreak' the safety measures in aligned LLMs

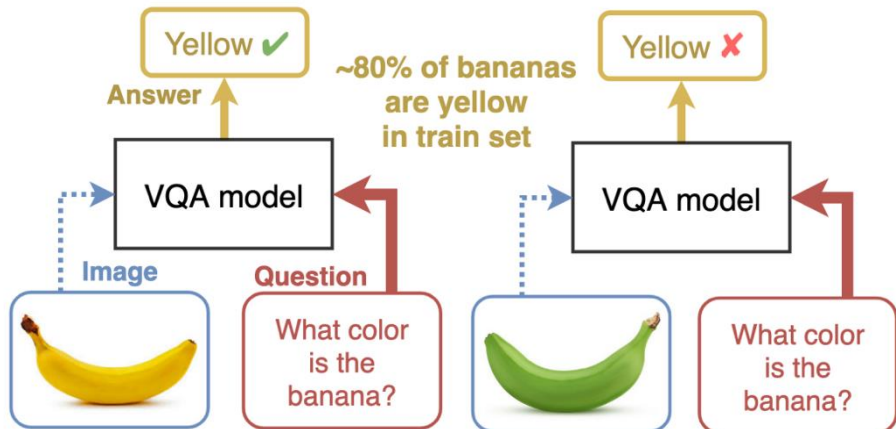


Still a big open challenge!

# Quantification - Safety

## Unimodal biases

VQA models answer the question without looking at the image

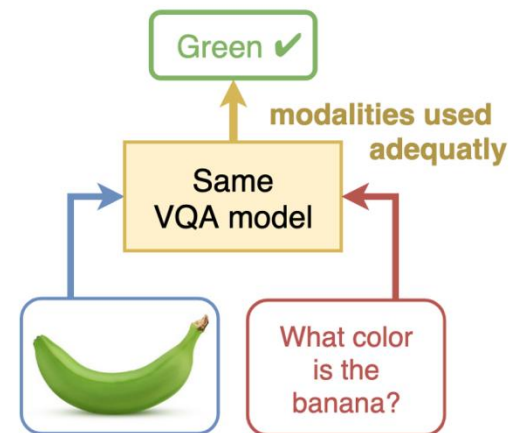


Balancing modalities

Balancing training



Not the case when trained with RUBi



[Wu et al., Characterizing and Overcoming the Greedy Nature of Learning in Multi-modal Deep Neural Networks. ICML 2022]

[Javaloy et al., Mitigating Modality Collapse in Multimodal VAEs via Impartial Optimization. ICML 2022]

[Goyal et al., Making the V in VQA Matter: Elevating the Role of Image Understanding in Visual Question Answering. CVPR 2017]

# Quantification - Safety

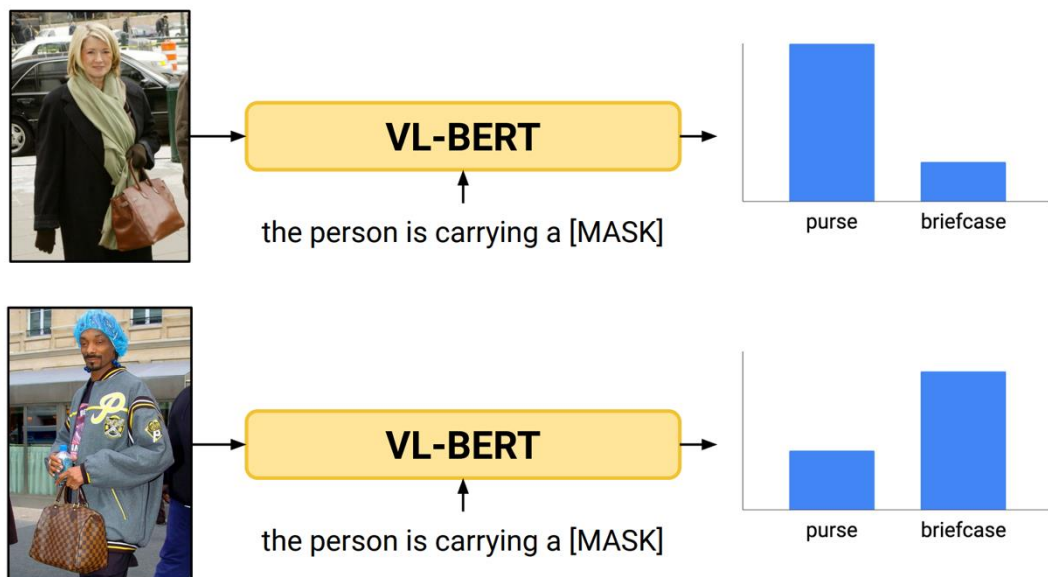
## Fairness and social biases

**Finding:** Image captioning models capture spurious correlations between gender and generated actions

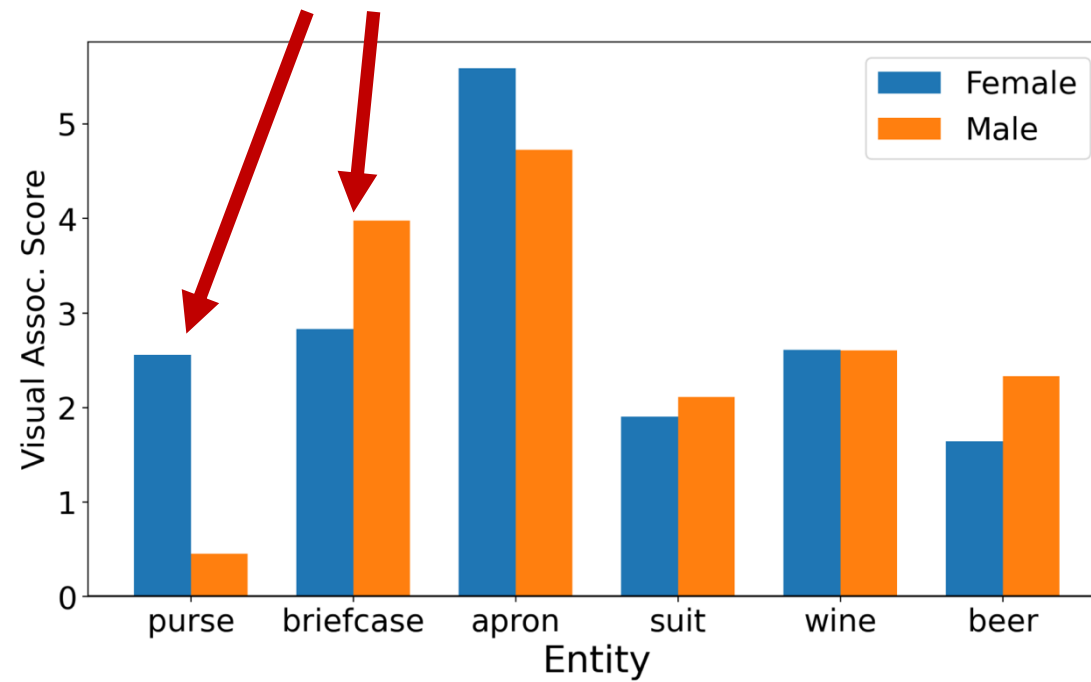


# Quantification - Safety

## Fairness and social biases



Visual information makes model more confident  
in reinforcing gender stereotypes

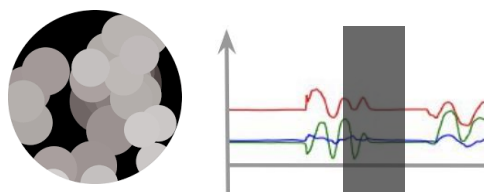


# Noise Topologies and Robustness

## Heterogeneity in noise

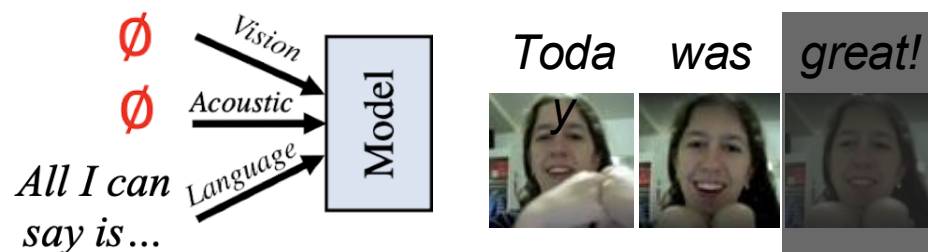
### Modality-specific robustness

noise → **nosie**



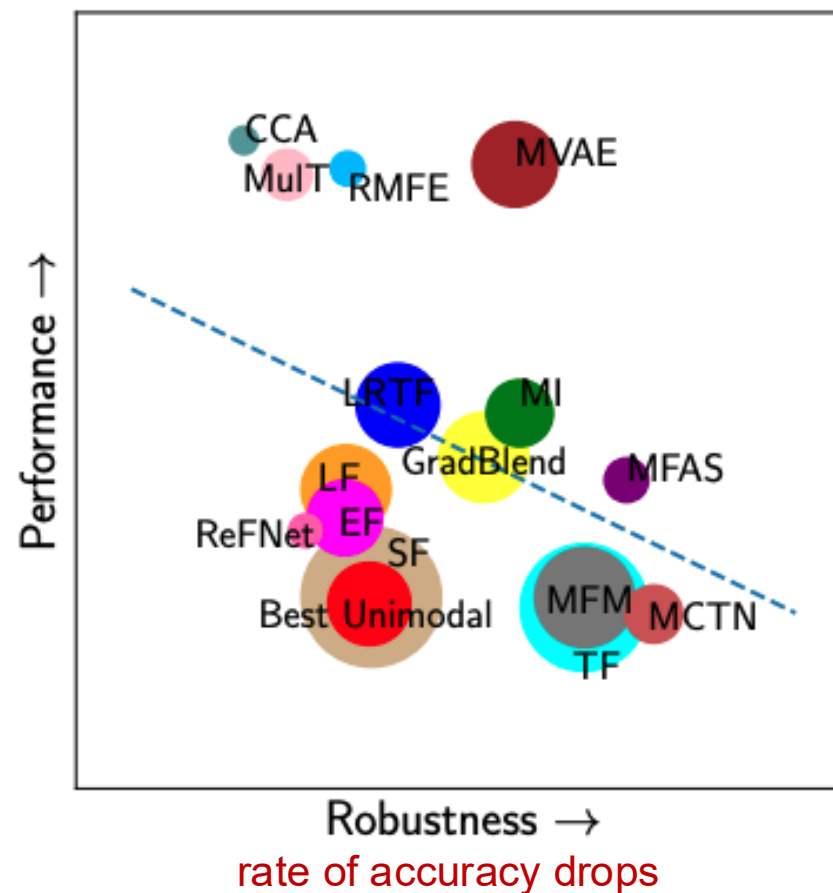
[Belinkov & Bisk, 2018; Subramaniam et al., 2009; Boyat & Joshi, 2015]

### Multimodal robustness



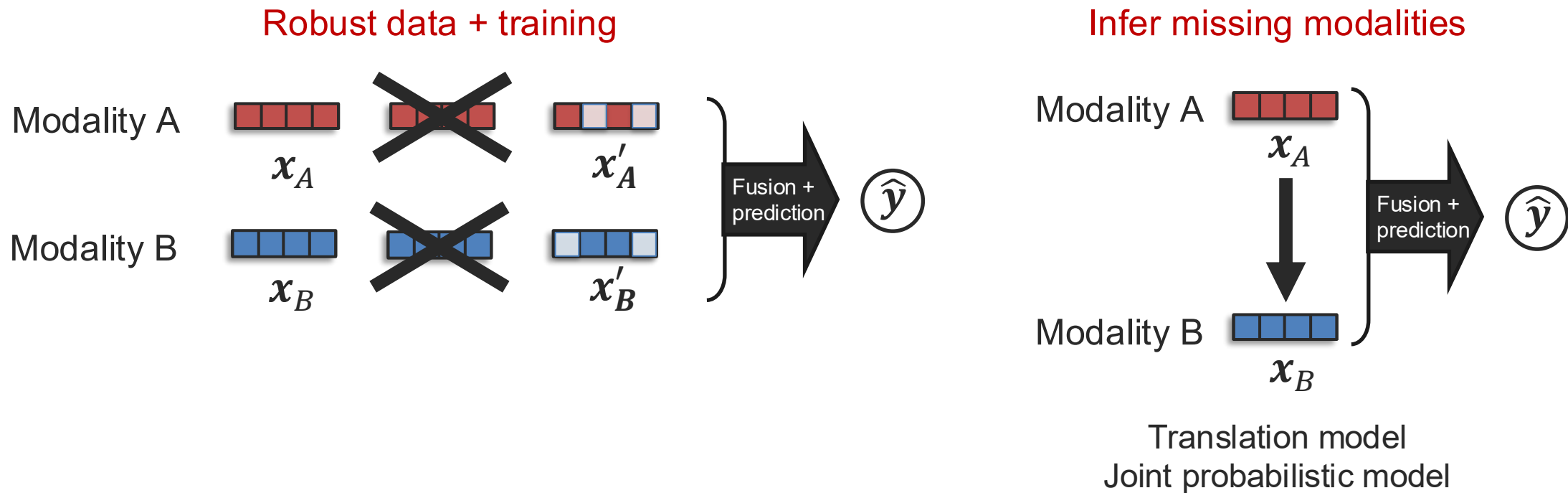
[Zadeh et al., 2020]

## Strong tradeoffs between performance and robustness



# Noise Topologies and Robustness

Several approaches towards more robust models



[Ngiam et al., Multimodal Deep Learning. ICML 2011]

[Srivastava and Salakhutdinov, Multimodal Learning with Deep Boltzmann Machines. JMLR 2014]

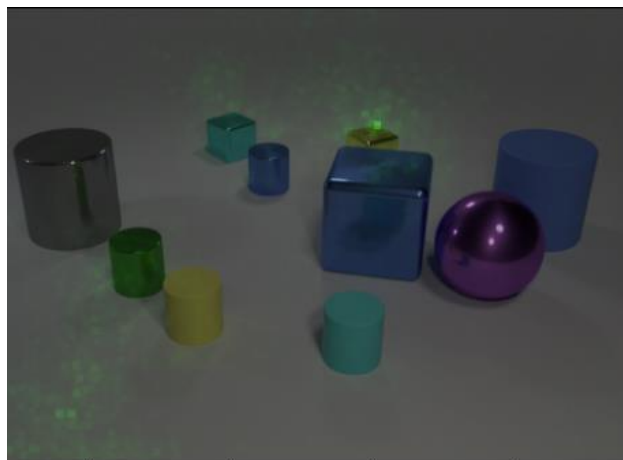
[Tran et al., Missing Modalities Imputation via Cascaded Residual Autoencoder. CVPR 2017]

[Pham et al., Found in Translation: Learning Robust Joint Representations via Cyclic Translations Between Modalities. AAAI 2019]

# Understanding Model Behavior

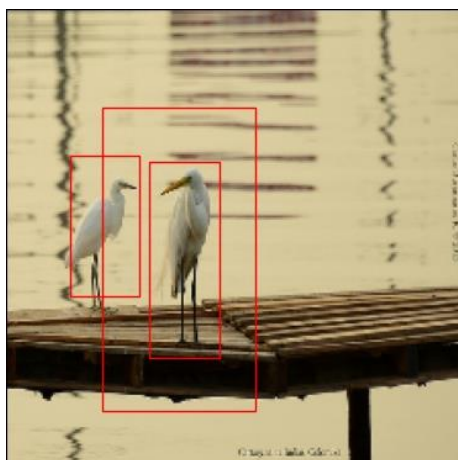
Identifying individual cross-modal interactions

## CLEVR



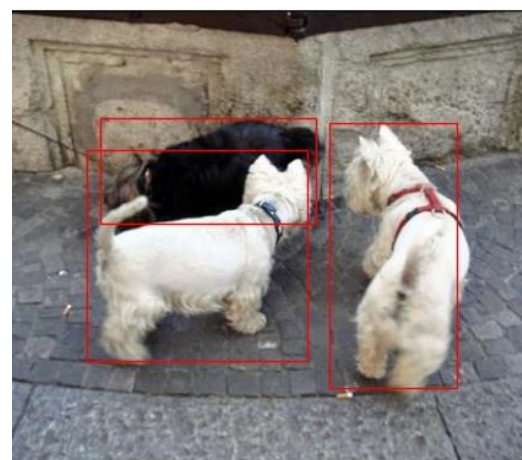
The other small shiny thing that is the same shape as the **tiny yellow shiny object** is what color?

## VQA 2.0



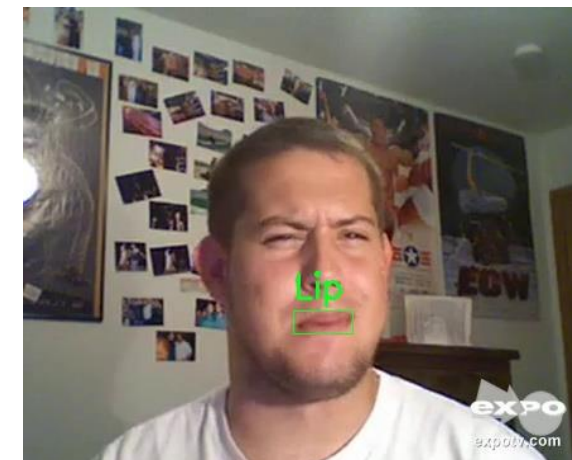
How many **birds**?

## Flickr-30k



**Three small dogs**, two white and one black and white, on a sidewalk.

## CMU-MOSEI



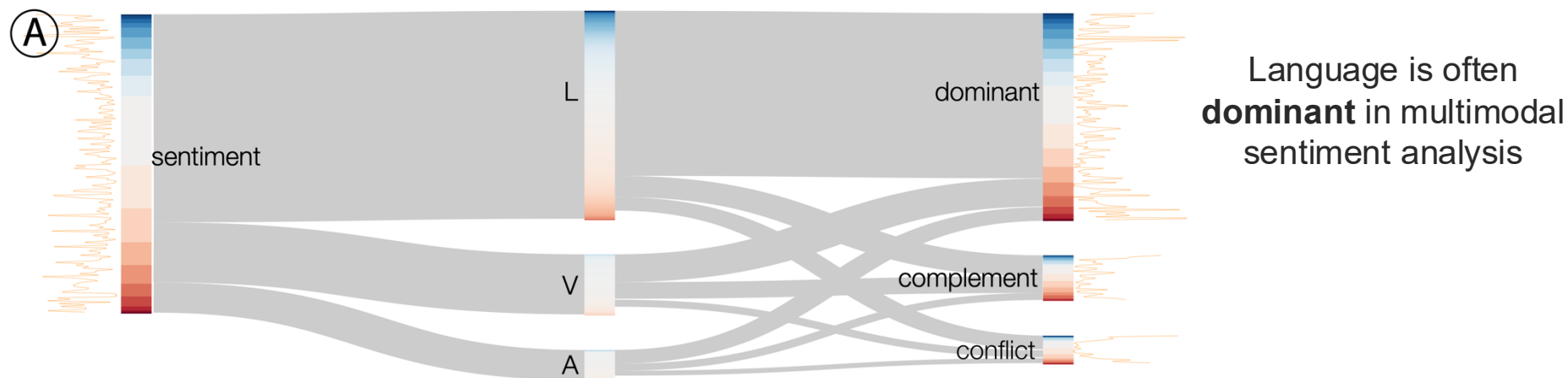
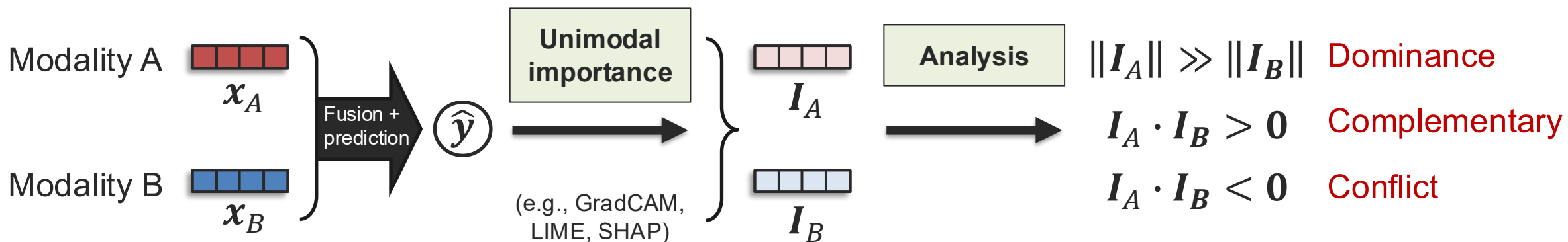
Why am I spending my money watching this? **(sigh)** I think I was more **sad**...

Correspondence

Relationships

# Understanding Model Behavior

## Classification of cross-modal interactions



# Understanding Model Behavior

Visualization website

See interactive website: <https://andy-xingbowang.com/m2lens/>

**E**

Instance View  Vision Feature  Audio Feature

Instance Summary  
sort By: error ▾ Desc ▾

Instance Detail

pitch ▾  
word: (umm) i really like how it's done because, if you watch  
Yaw ▾

Feature	Importance
features/mo	
not	-0.918
movie	-0.249

this movie five times you will still not understand everything about it

pitch ▾  
word: it definitely, it worked.  
Surprise ▾

Video Detail

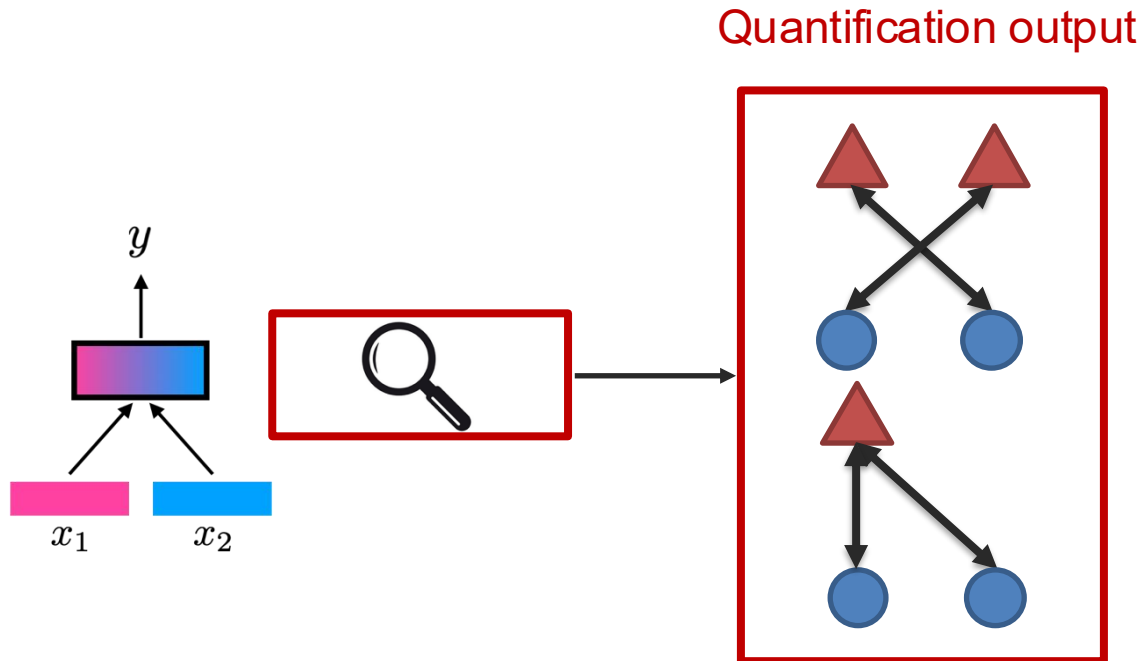


The video detail shows a person sitting on a couch, holding a movie case. A green bounding box is drawn around the person's face. Labels 'Disgust' and 'FAU9 : Nose Wrinkler' are overlaid on the face. The video player interface at the bottom shows a progress bar at -0:25 and a 1x playback speed.

# Evaluating Quantification

How can we evaluate the success of quantification?

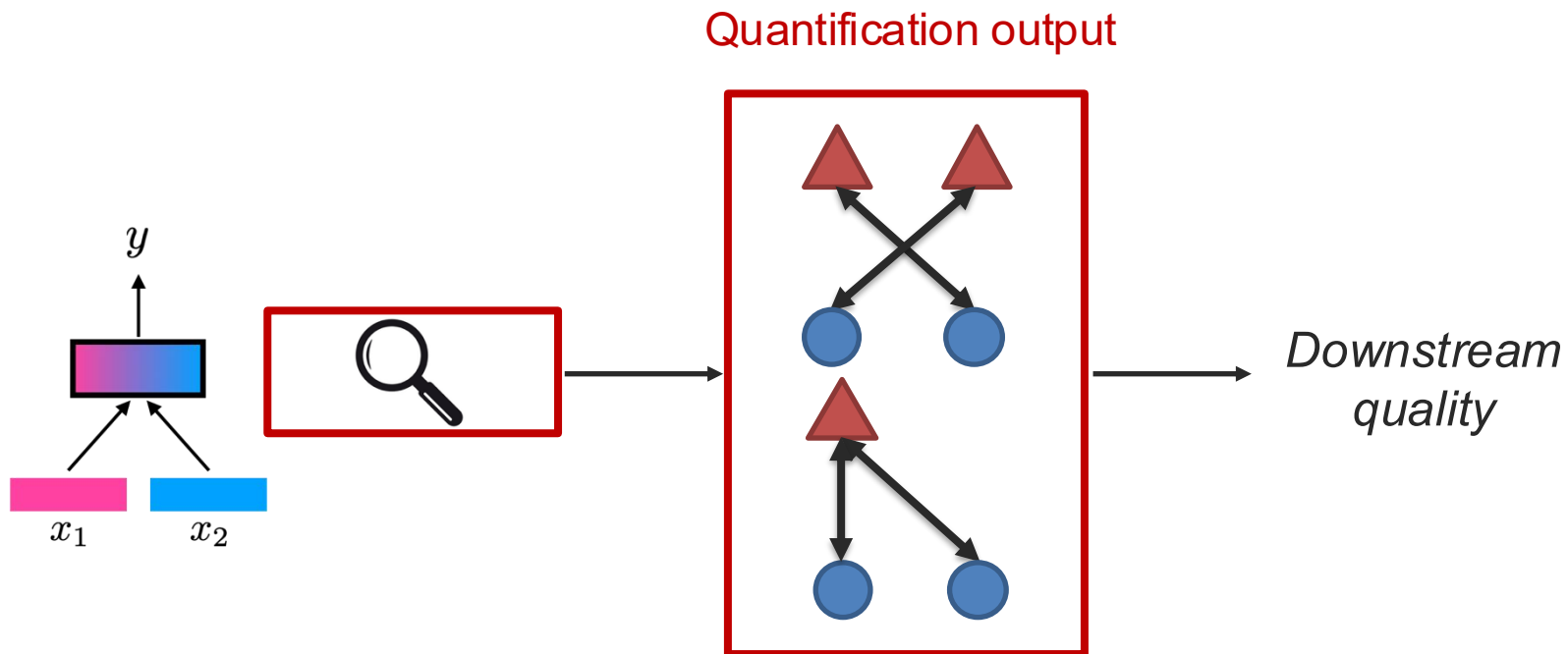
*Problem: real-world datasets and models do not have quantification outputs annotated!*



# Evaluating Quantification

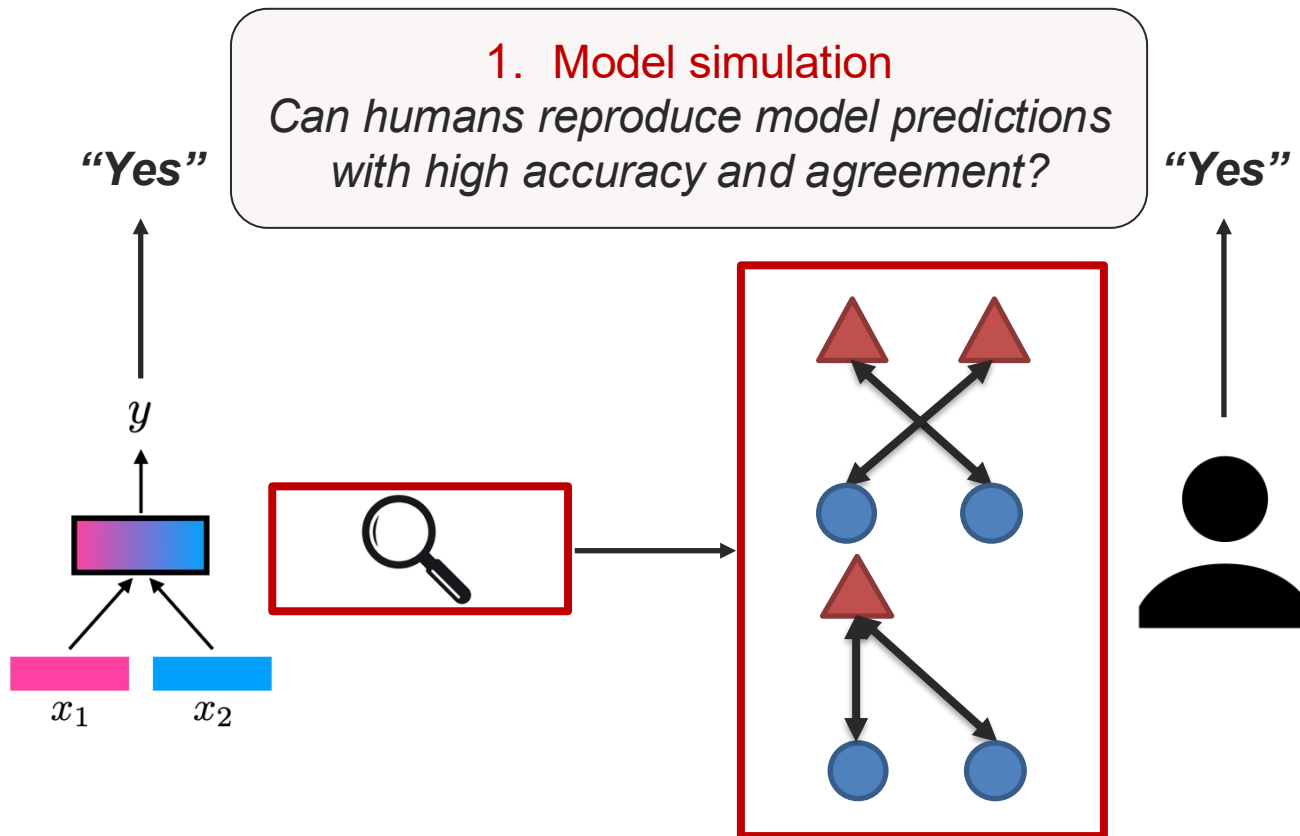
## Indirect evaluation

*Find some downstream quality that practitioners find useful and can be easily evaluated.*



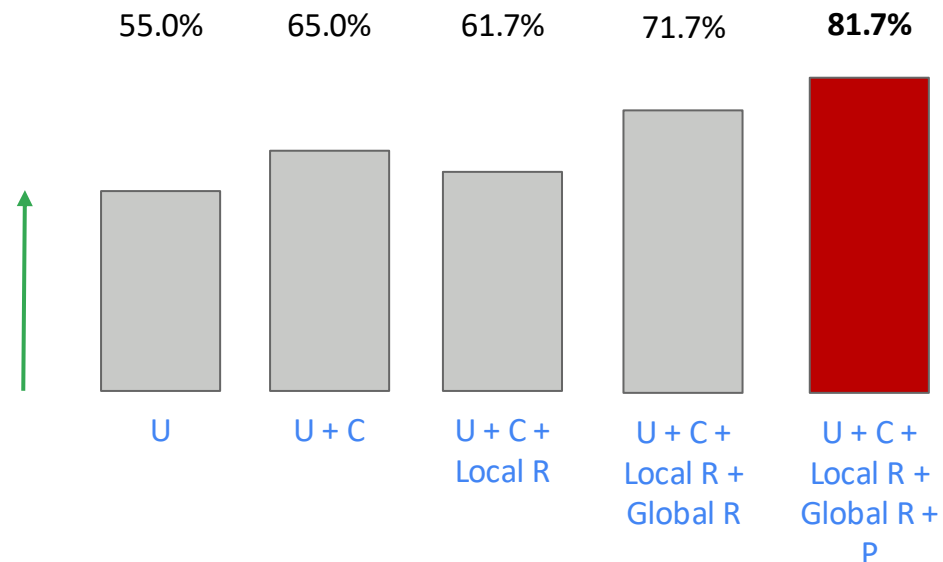
# Evaluating Quantification

Indirect evaluation: Model simulation



# Evaluating Quantification

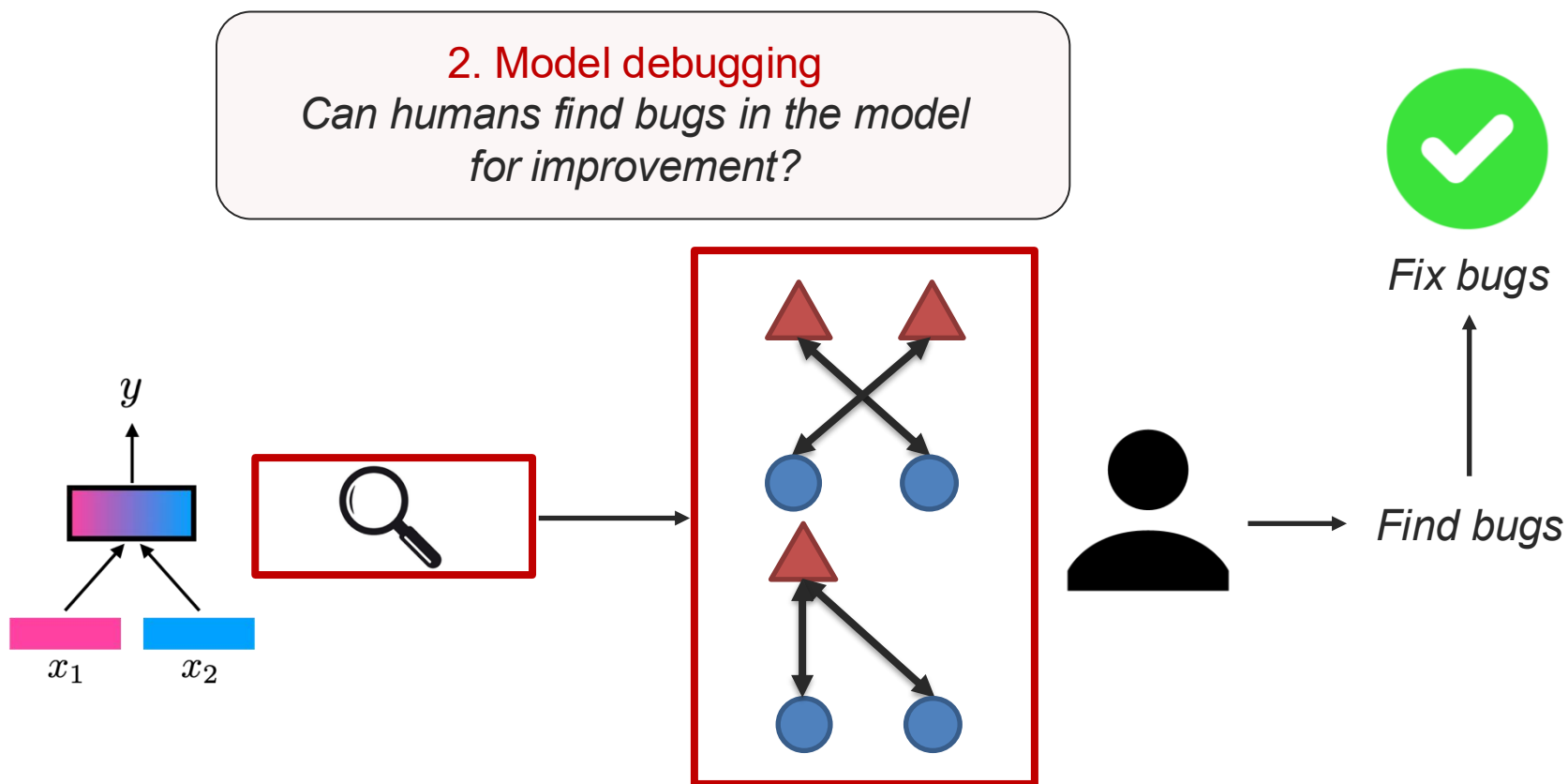
Indirect evaluation: Model simulation



MultiViz stages leads to higher accuracy and agreement  
Blind test + reasonable baselines + measurable outcome

# Evaluating Quantification

Indirect evaluation: Model error analysis and debugging

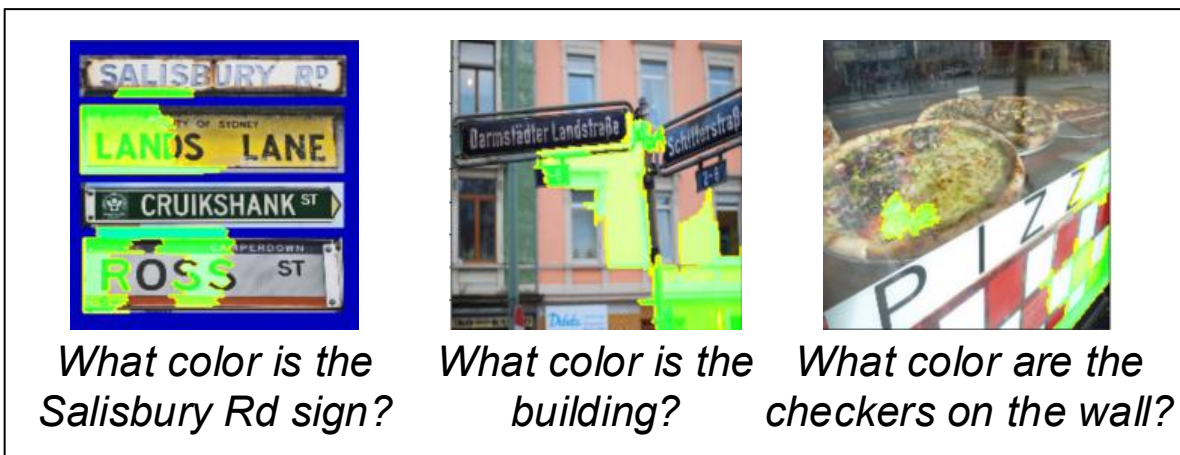


# Evaluating Quantification

Indirect evaluation: Model error analysis and debugging



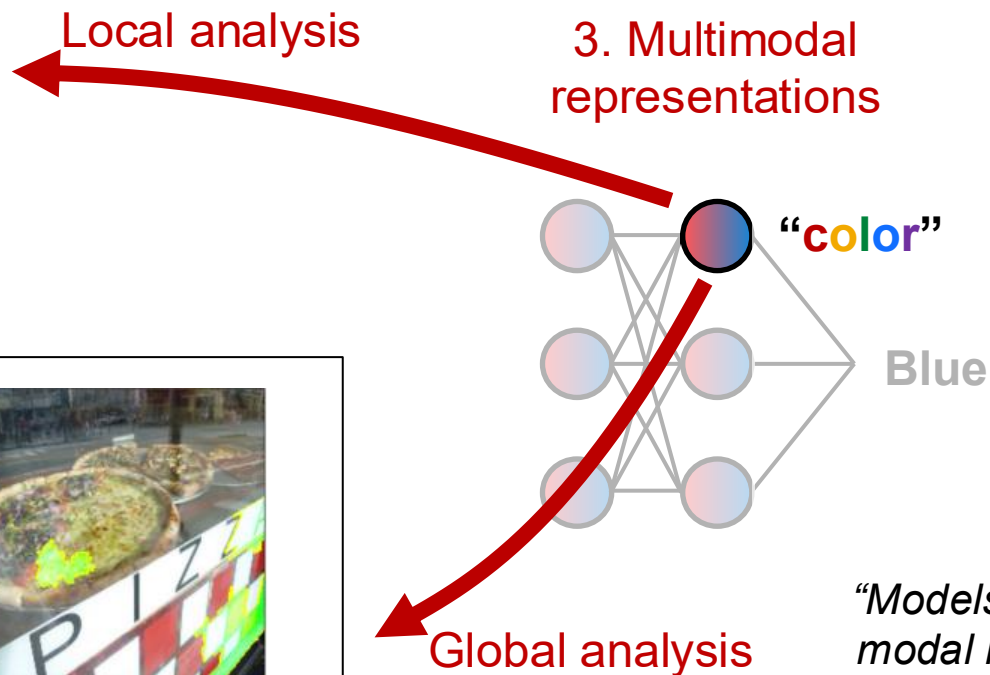
*What color is the tie of the second man to the left?*



*What color is the Salisbury Rd sign?*

*What color is the building?*

*What color are the checkers on the wall?*



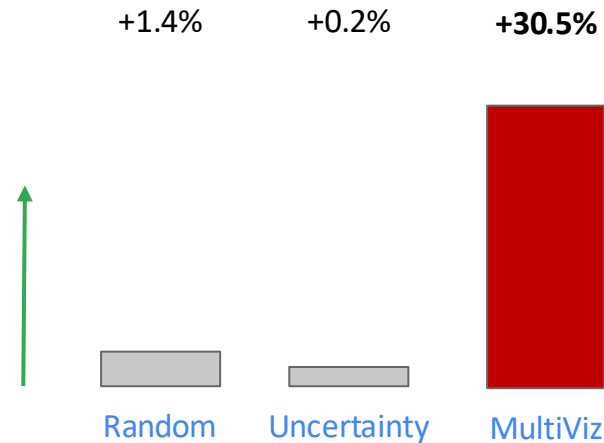
# Evaluating Quantification

Indirect evaluation: Model error analysis and debugging

*“Models pick up cross-modal interactions but fail in identifying color!”*



*Add targeted examples involving color.*



*Side note: we used this to discover a bug in a popular deep learning code*

  
**Transformers**

**MultiViz enables error analysis and debugging of multimodal models**