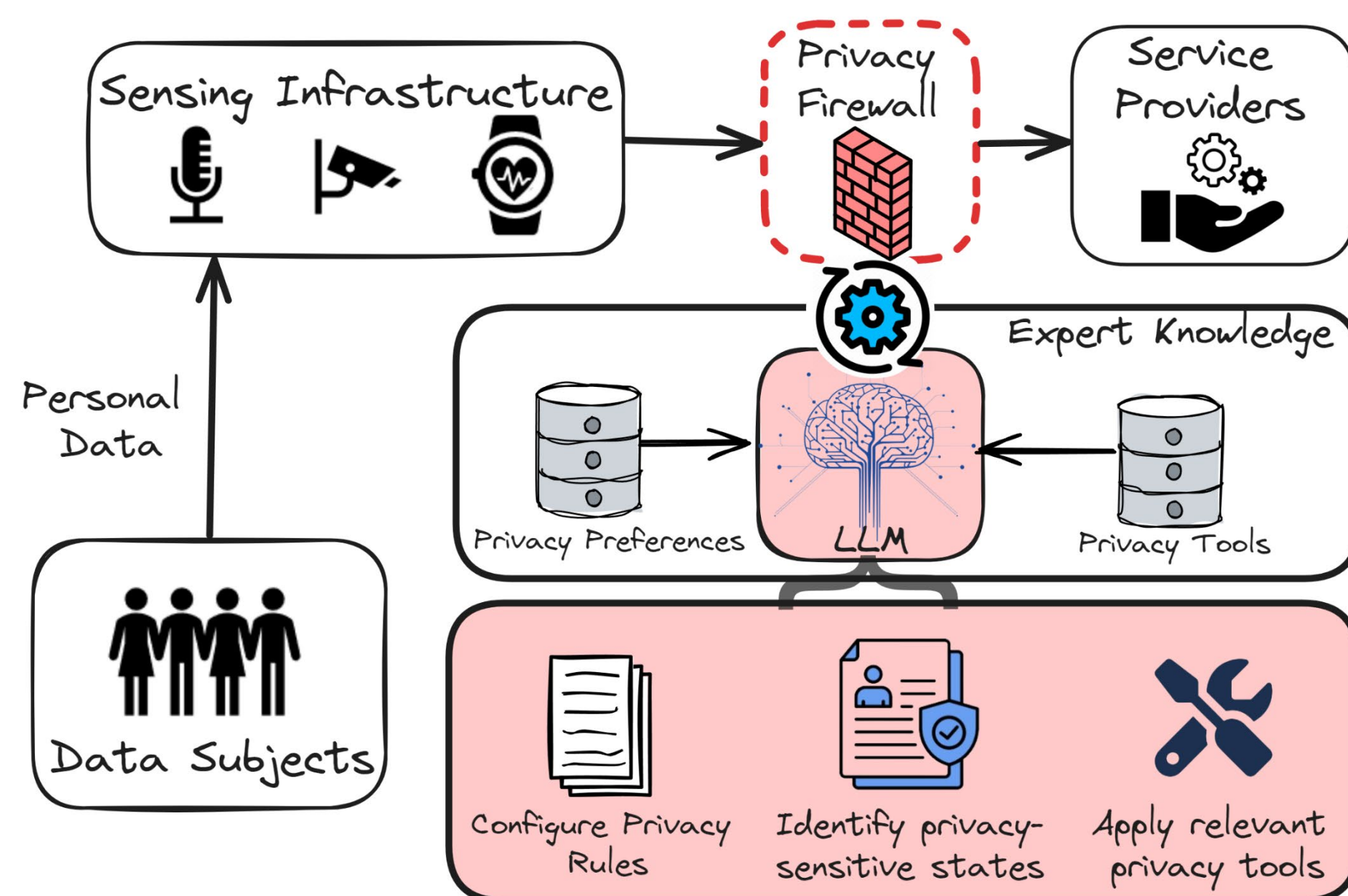
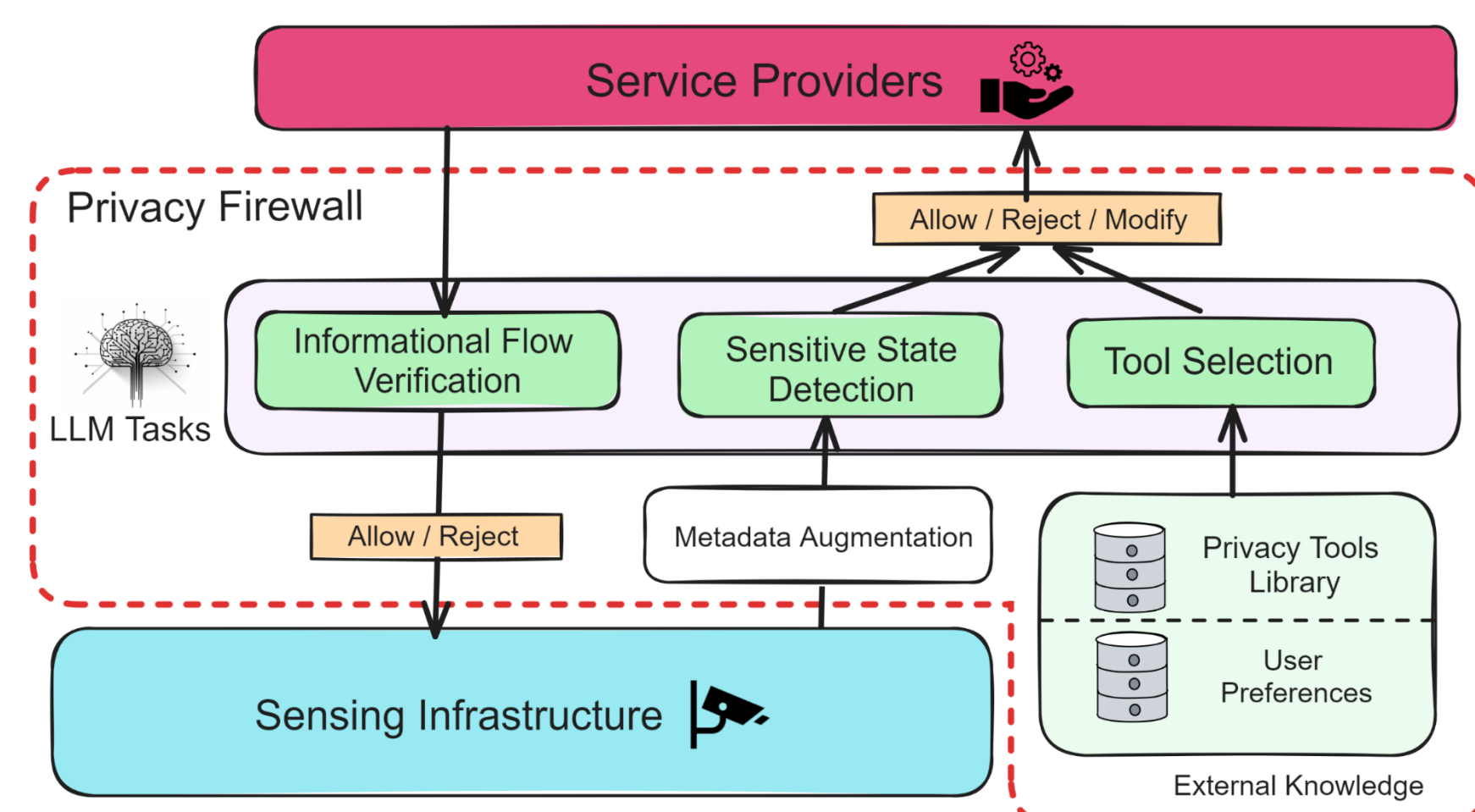


## Objectives

We propose a system for automatically managing privacy decisions of sensory data on behalf of data subjects in smart built environments,



# LLM-based Privacy Firewalls

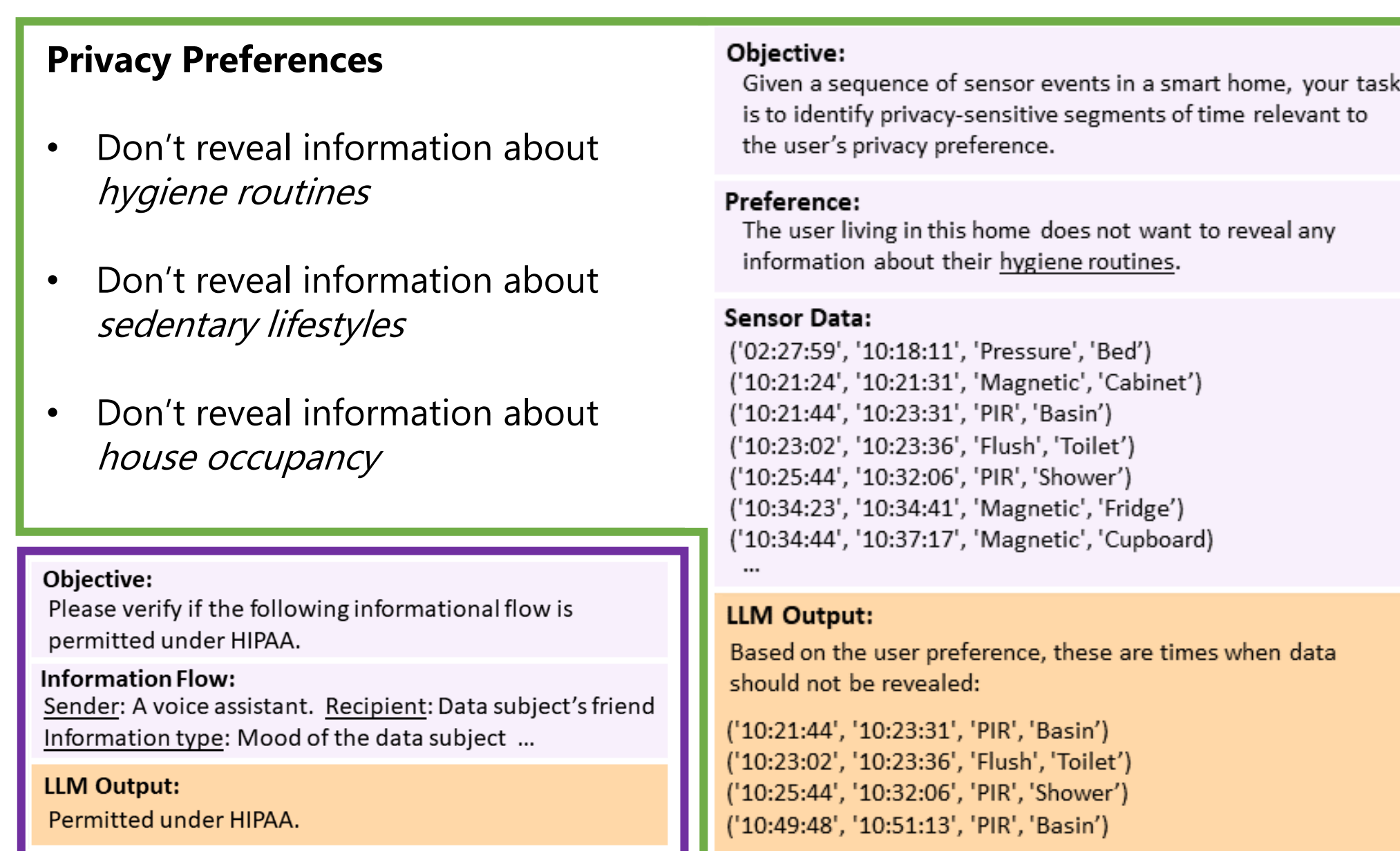


Our system, **PrivacyOracle**, accomplishes 3 tasks:

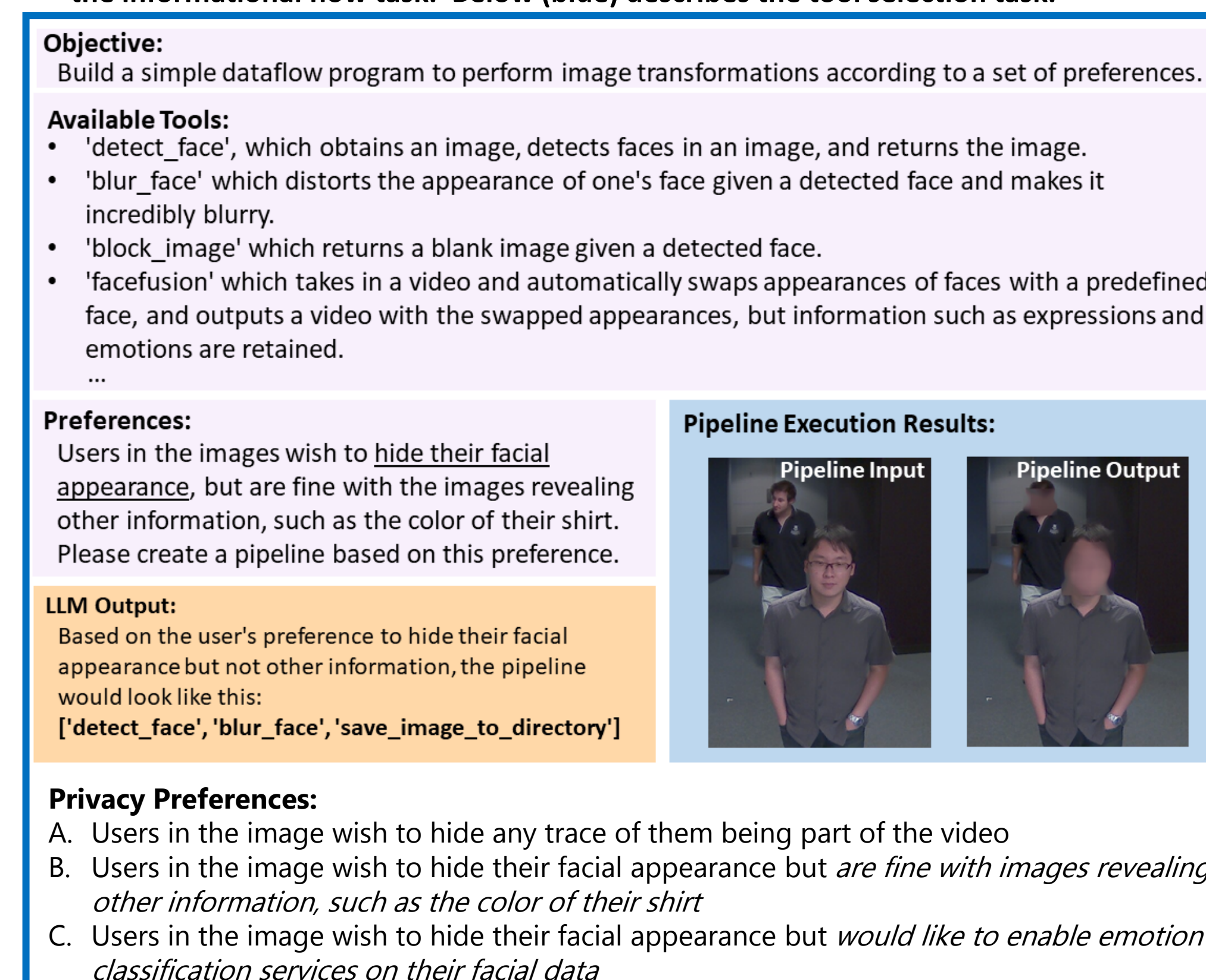
- *Informational flow verification* uses world knowledge of legal rules to identify violations of privacy
- *Sensitive state detection* uses lower level sensor information and preferences to identify privacy-sensitive states
- *Tool selection* identifies relevant privacy tools for a privacy preference and builds appropriate data processing pipelines

## Methods

We utilize GPT-3.5 and GPT-4 to generate a variety of privacy decisions, with examples shown below.



Top right (green) describes the sensitive state detection task. Left (purple) describes the informational flow task. Below (blue) describes the tool selection task.



## Results

Privacy-Sensitive State	IoU	F1
Hygienic Activities	0.684	0.831
Sedentary Lifestyle	0.844	0.983
House Occupancy	0.961	0.701

### Measuring agreement between LLM-inferred privacy sensitive states and ground truth sensitive states

For *informational flow verification*, we use a manually curated dataset of 16 HIPAA scenarios with various contextual integrity parameters, and achieve a false positive rate and false negative rate of **6.25%**.

Requirement	Age F1	Gender F1	Race F1	Emotion F
A	0.0	0.0	0.0	0.0
B	0.208	0.373	0.151	0.081
C	0.187	0.371	0.288	0.407

### Privacy (age, gender, race recognition) vs. Utility (emotion) on LLM-generated processing pipelines for each privacy preference

In the *Tool Selection* task, each tool has different privacy/utility costs and the LLM selects the appropriate tool for each preference.

## Future Work

- Examine enforcement of privacy decisions in non-cooperative sensing environments
- Evaluation of LLM results on different prompting strategies and validation mechanisms
- Managing conflicting privacy requires among users (democratization/negotiation)

## Acknowledgements

This research was sponsored in part by the National Science Foundation (NSF) under awards # 1705135, 2124252 and 2211301, the NIH mDOT Center under award #1P41EB028242, the DEVCOM ARL under Cooperative Agreement #W911NF-17-2-0196, and the DARPA/ANSR Program under Contract #FA875023C0519.