

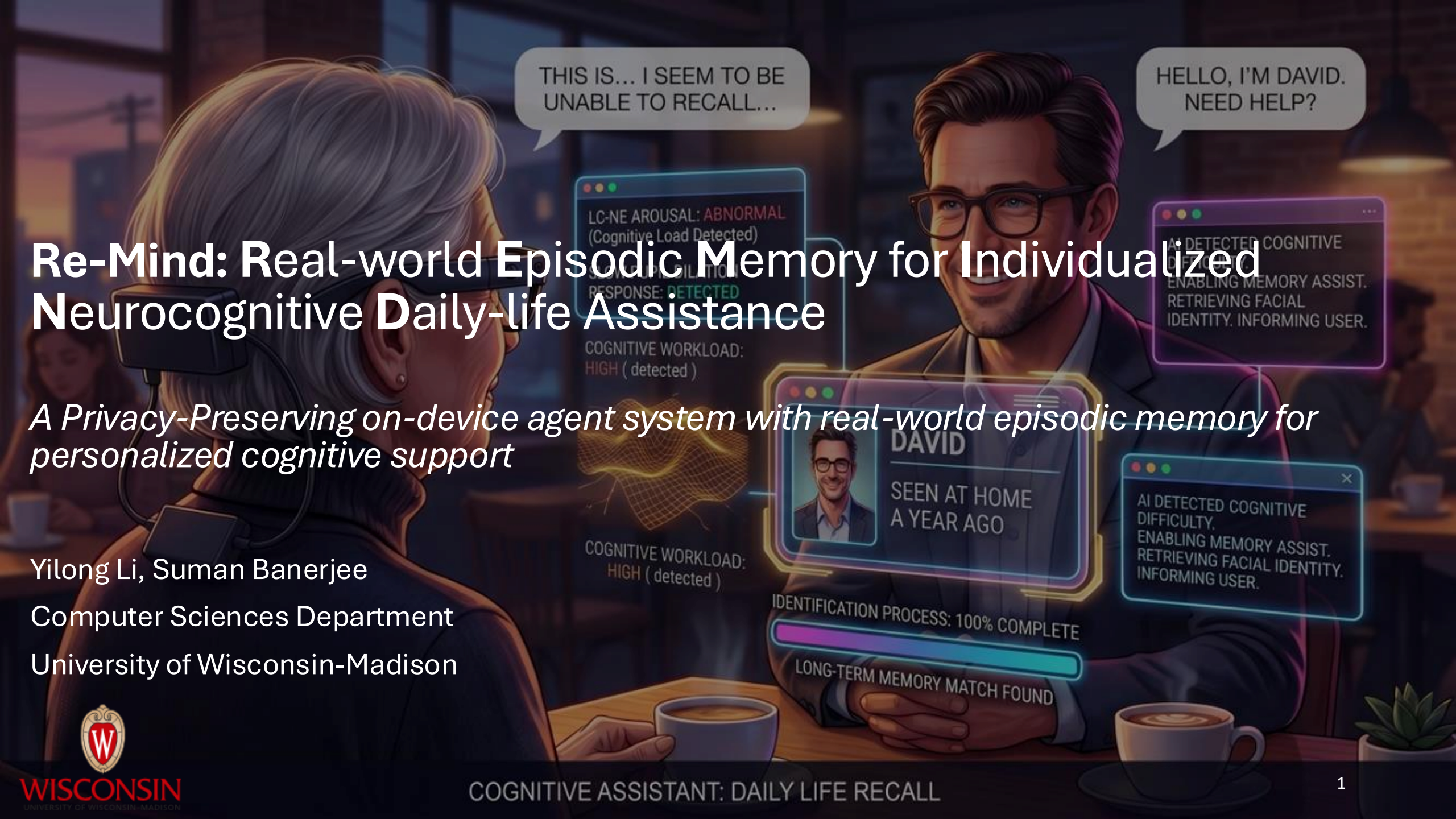
Re-Mind: Real-world Episodic Memory for Individualized Neurocognitive Daily-life Assistance

A Privacy-Preserving on-device agent system with real-world episodic memory for personalized cognitive support

Yilong Li, Suman Banerjee

Computer Sciences Department

University of Wisconsin-Madison



Motivation



AI Models and Agents are now powerful,
but still lacking connection to people



Not Privacy-Preserving

No Real-Life Perception Interface

No Personalized Long-Term Memory

“Powerful, but still behind the screen”

Motivation

Ubiquitous Surveillance

Continuous Collecting Data

Cloud Data Mining

Cloud Power

The Privacy vs. Capability Dilemma

On-device AI is better in privacy-preserving but often slow speed



Privacy Secured

Small Battery Life

Hardware Limited

Compute Constrained

High Latency

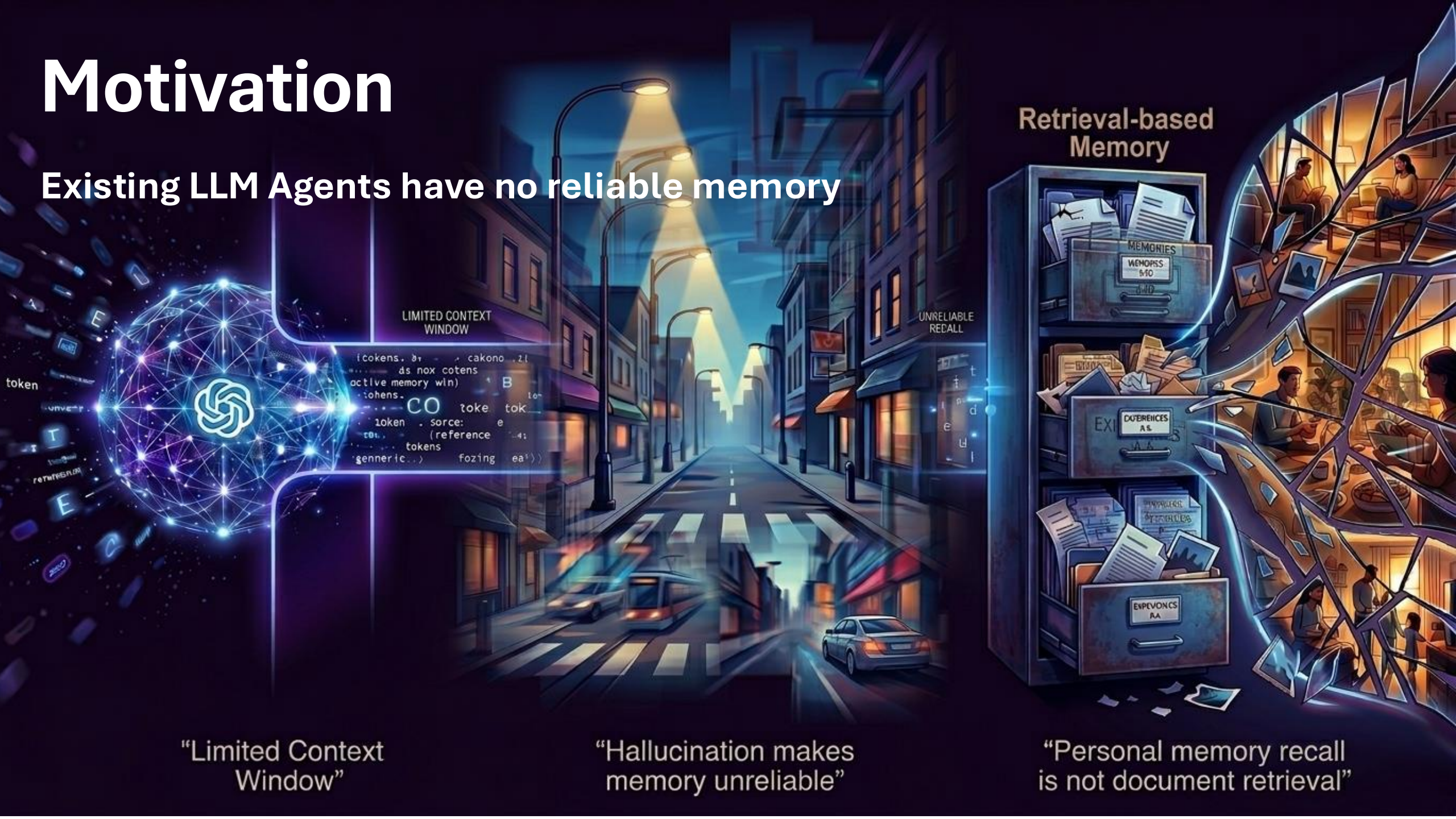
Slow Vision Inference

Frontier AI models are powerful but intrusive

“Powerful but intrusive cloud AI vs. Private and local AI on low-power devices”

Motivation

Existing LLM Agents have no reliable memory



“Limited Context Window”

“Hallucination makes memory unreliable”

“Personal memory recall is not document retrieval”

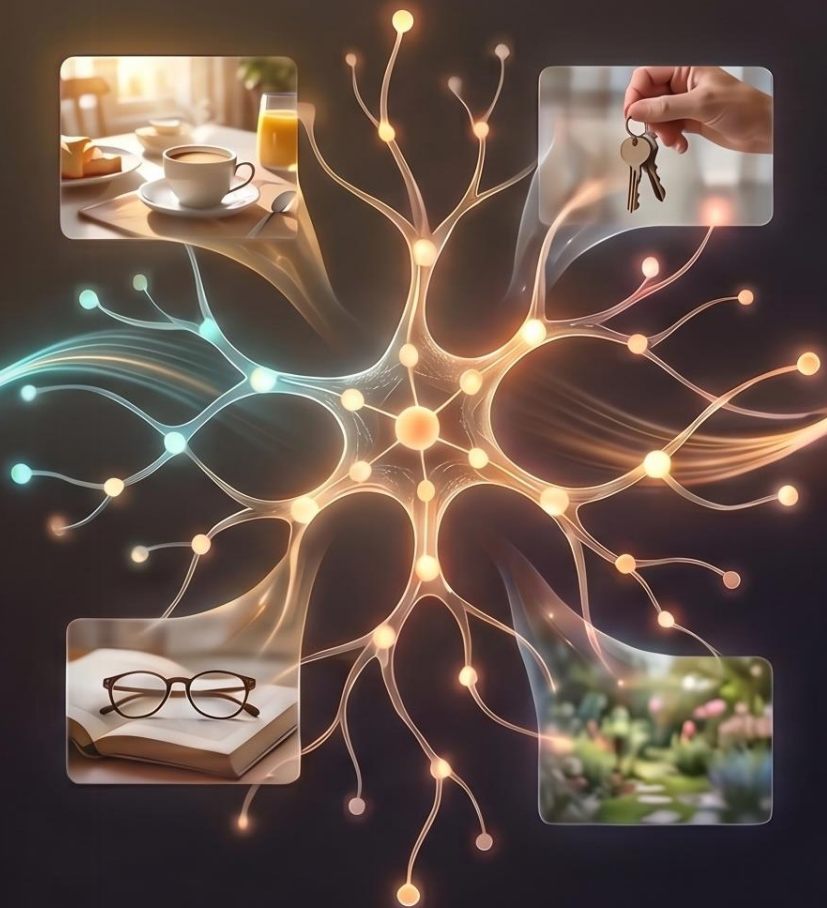
Solution



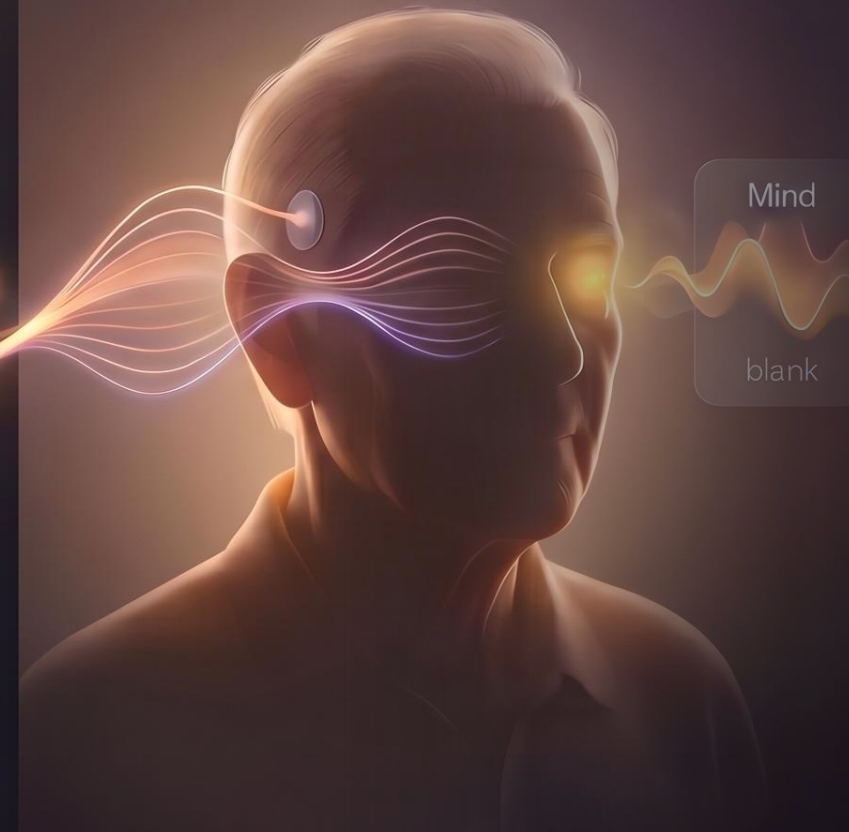
Efficient On-Device AI

What We Build

A Dedicated, Privacy-Preserving Cognitive Assistant Agent with Personalized Memory



Privacy-Preserving Personal Memory



Cognitive-State-Aware Assistance

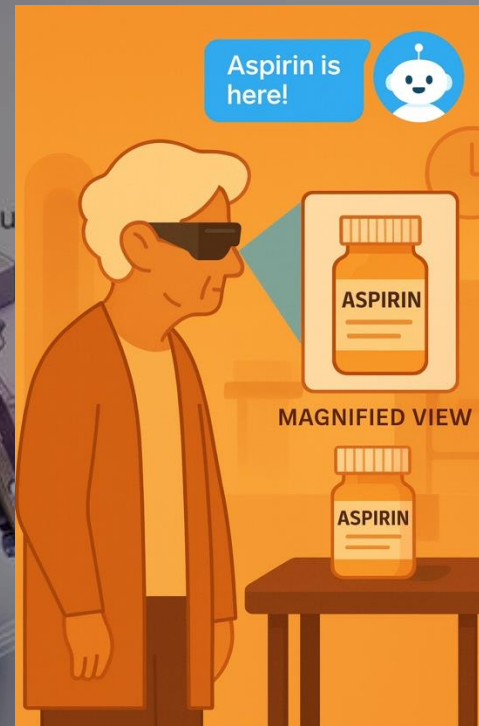
Solution

What We Want:

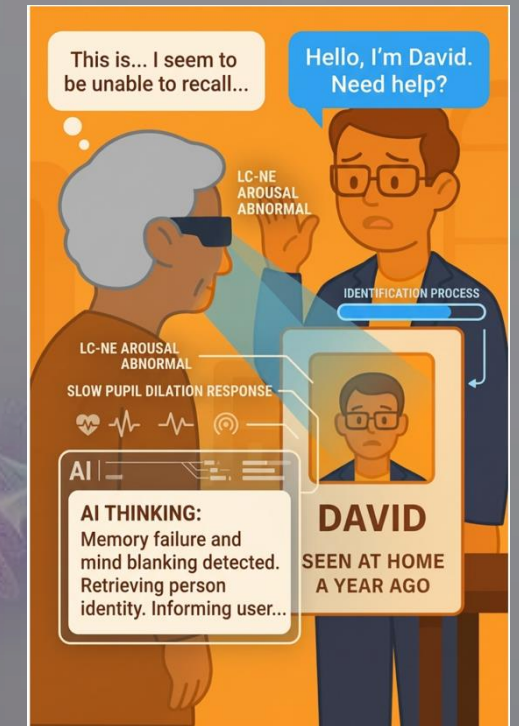
A Device for People's Cognitive Assistant and Accessibility



Visual Understanding



Find Item



Memory-failure Awareness

Not Another Platform. The Right Device.

Solution



Home Keys

MEETING DAVID: 1 YEAR AGO

Coffee Mug - 8 AM Routine

Solution

- Vision-Language Model
- Face Recognition
- Object Detection



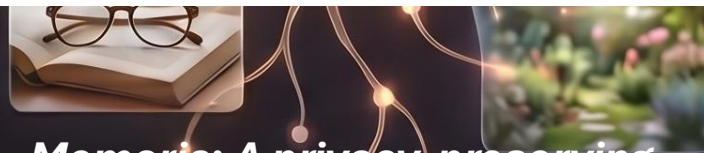
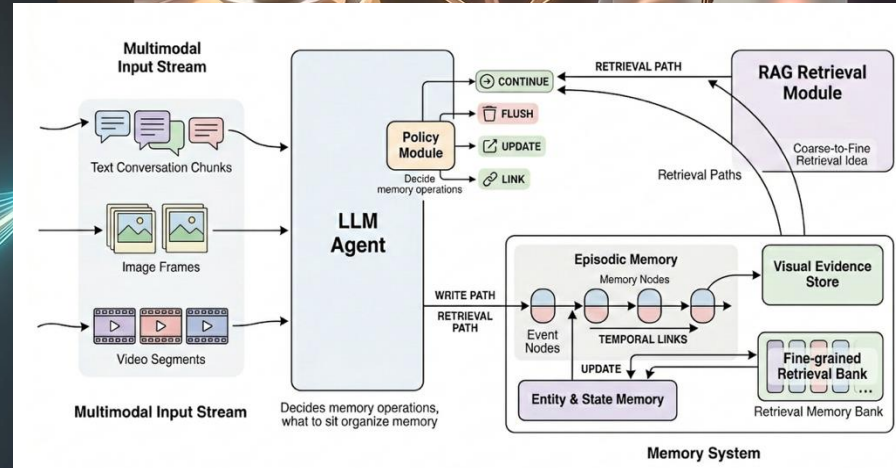
Platform of Efficient Vision Language Model Inference on a battery-powered device

Efficient On-Device AI

What We Have Now

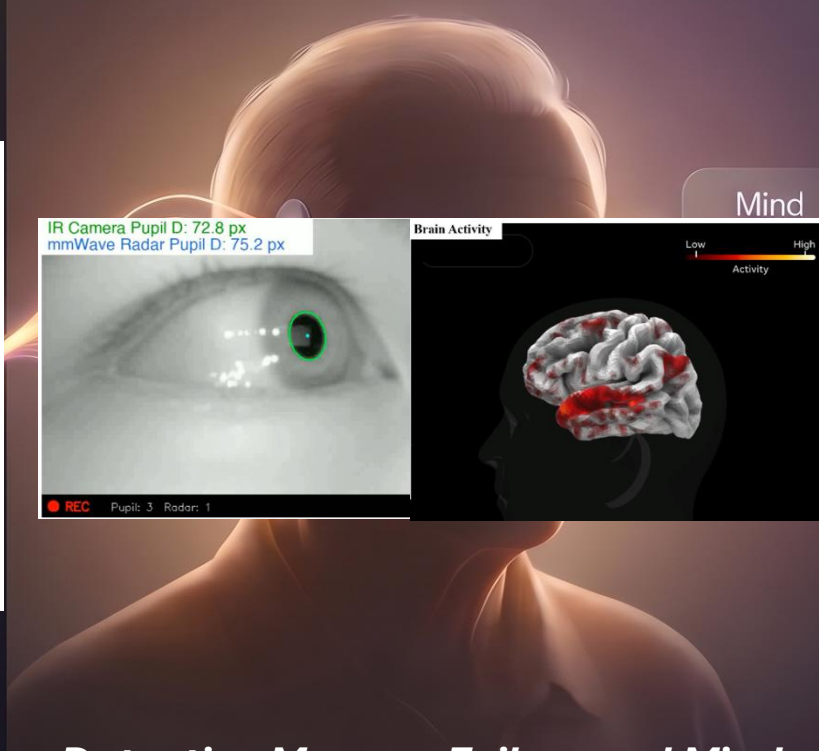
A Real-Time Demo of Visual Understanding and Face-and-Item Memory on Our Custom-Built Hardware

Developed Techniques for Cognitive-State Estimation and Memory-Failure Awareness



Memoria: A privacy-preserving, continuously evolving personal agentic memory

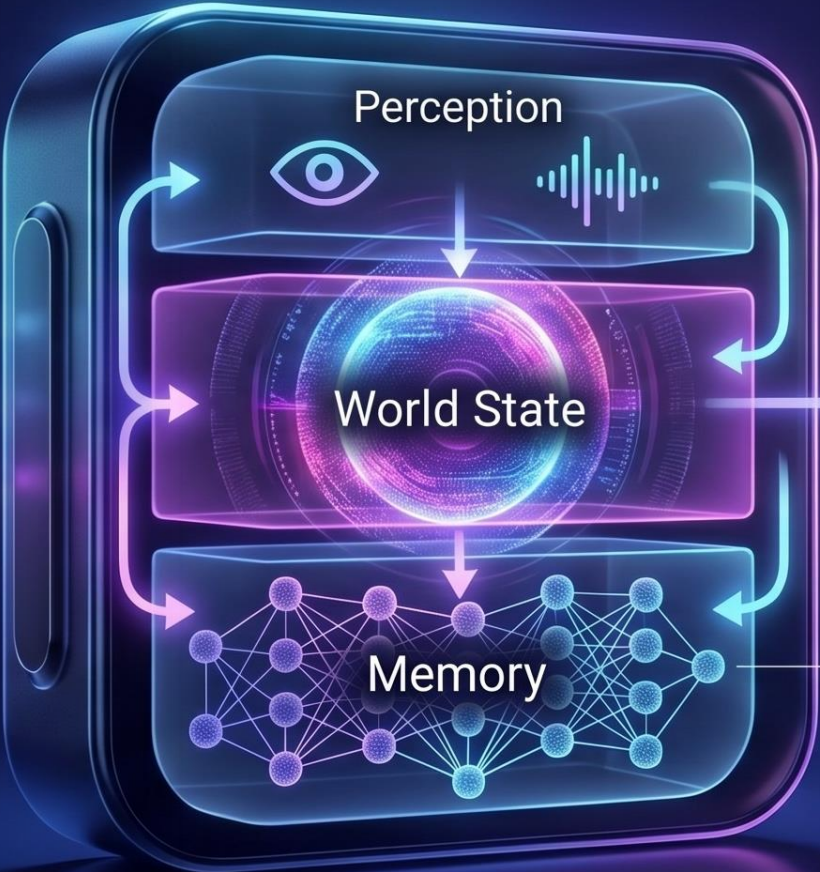
Privacy-Preserving Personal Memory



Detecting Memory Failure and Mind Blanking via LC-Linked Pupillary Responses

Cognitive-State-Aware Assistance

Privacy-Preserving Agent



Search



Identify



Focus



Face

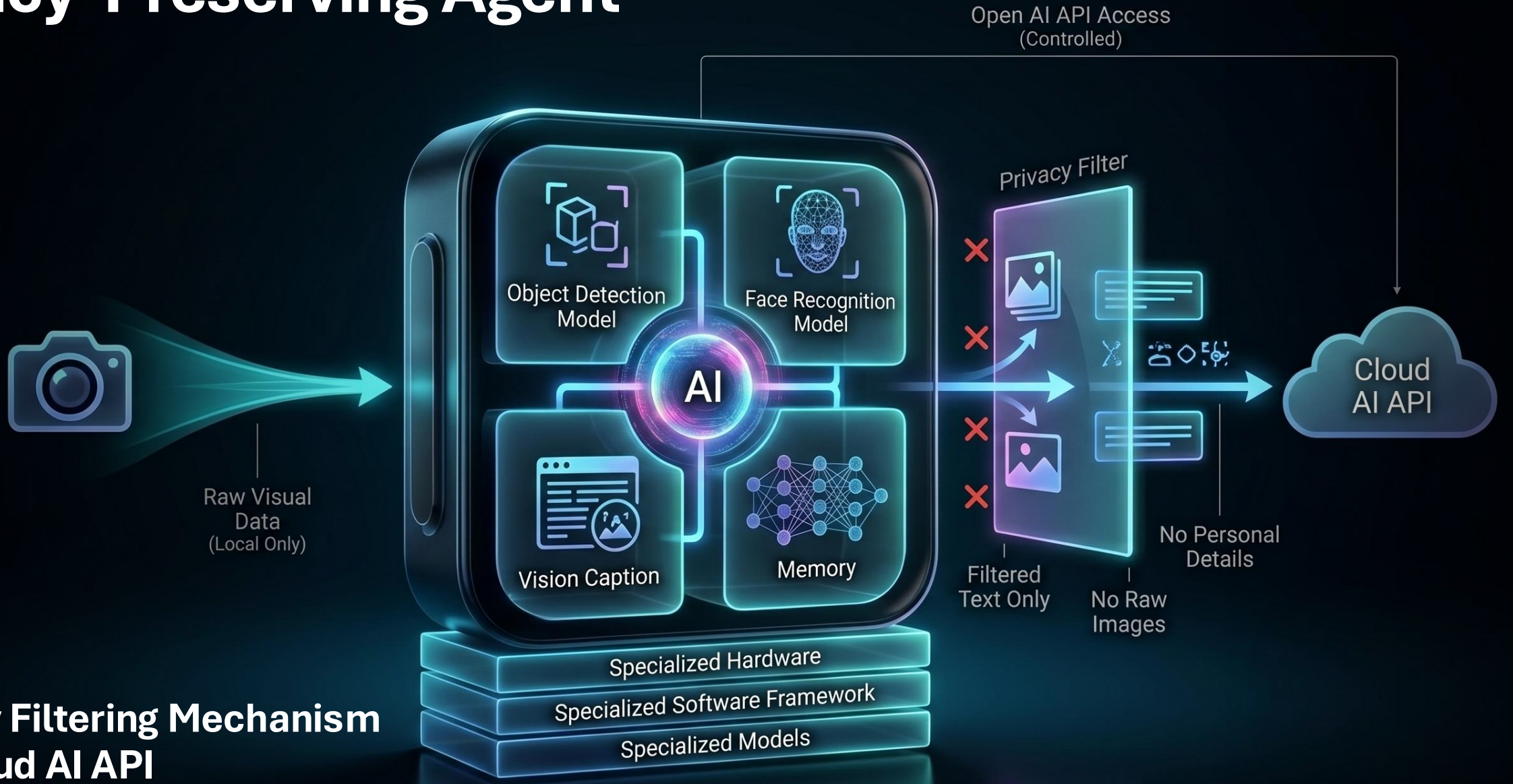


Object

Learned from daily life

Personal Agent

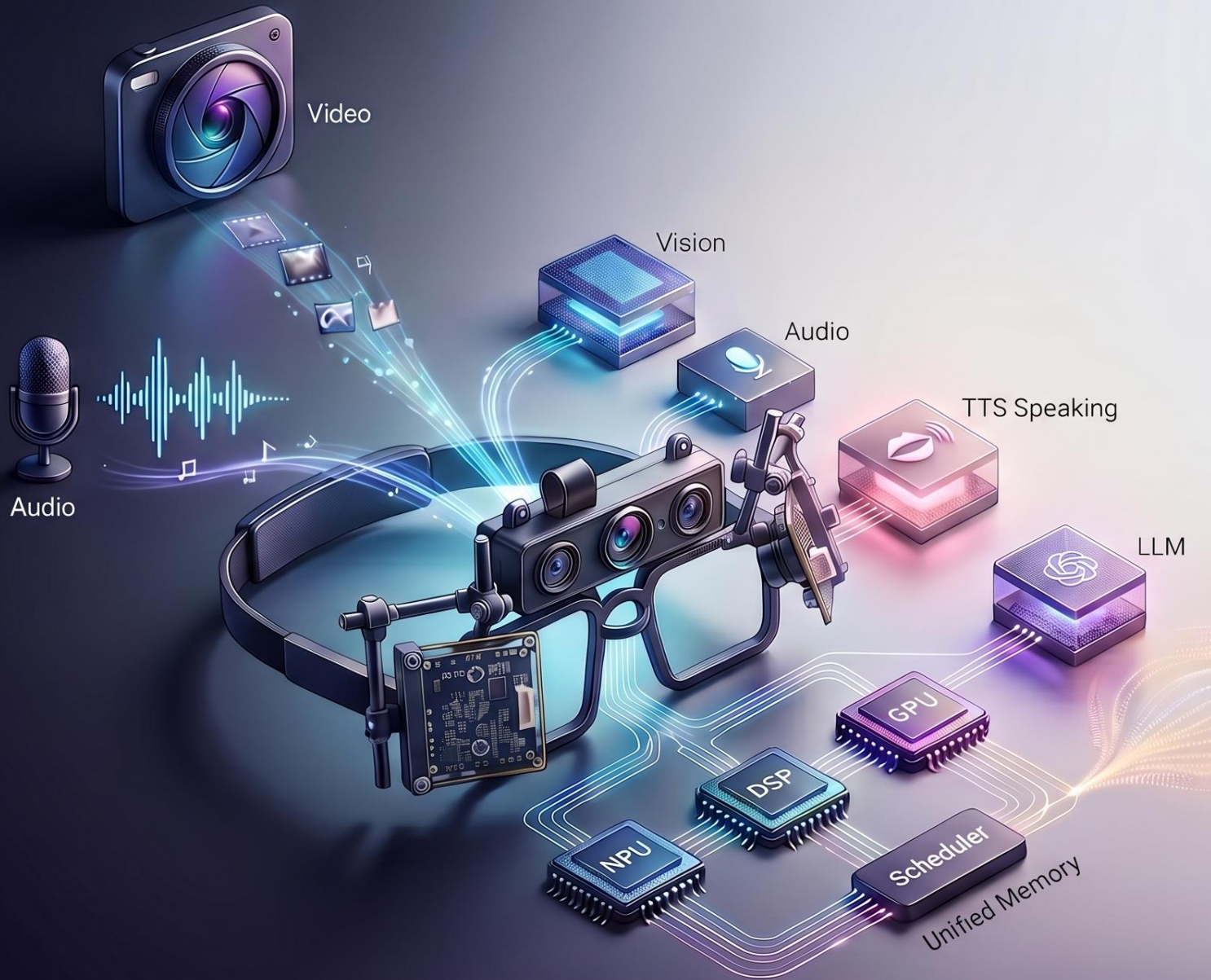
Privacy-Preserving Agent



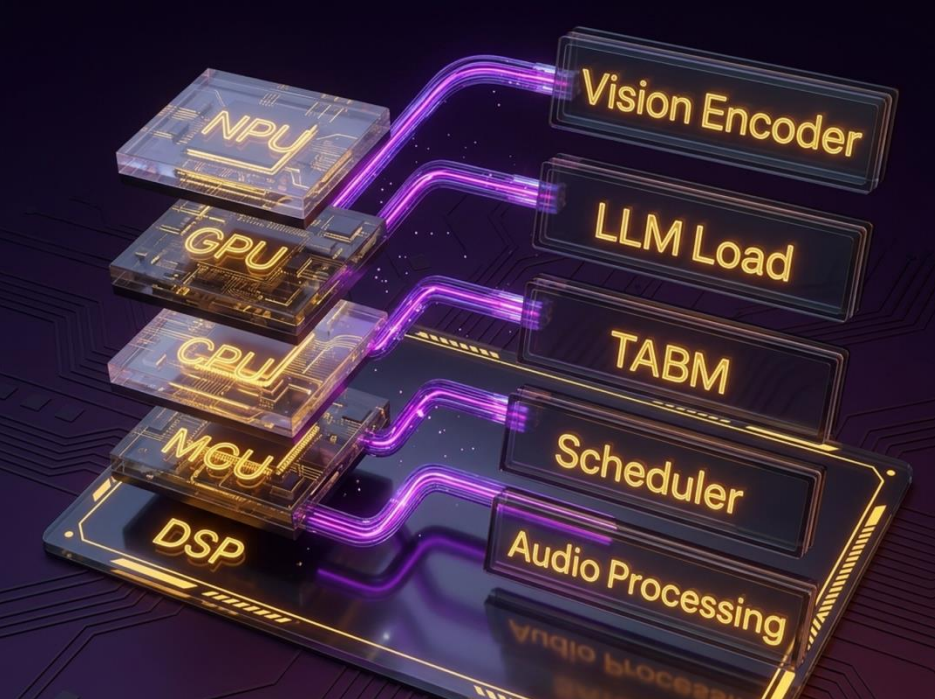
Privacy Filtering Mechanism for Cloud AI API

"Not Another Platform. The Right Device."

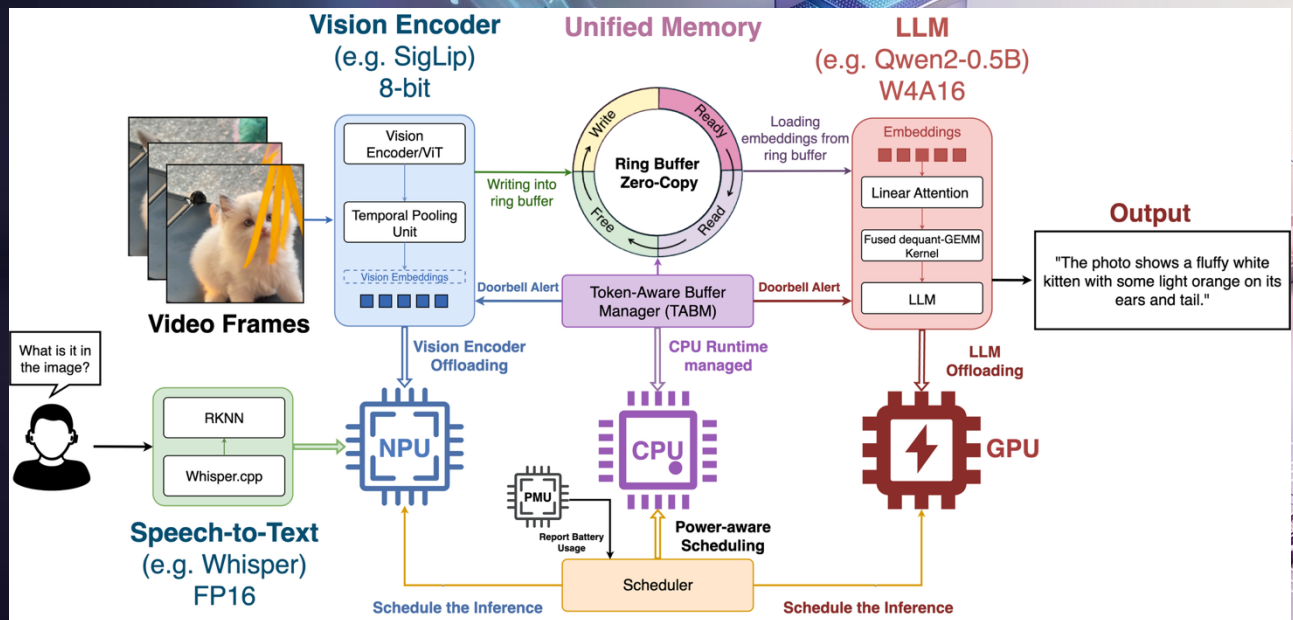
Efficient Cross-Accelerator Inference



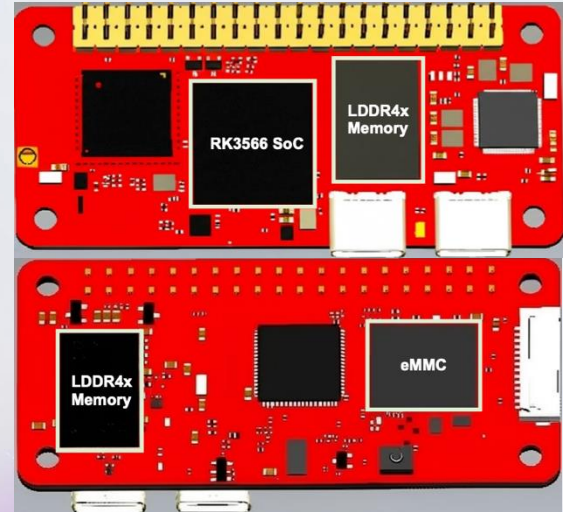
Many companies are building expensive custom AI chips for wearables. **In contrast, we rely only on low-cost COTS and readily available components.** By adopting a cross-accelerator design, we partition AI models across different cheap accelerators and connect them through our specialized memory architecture.



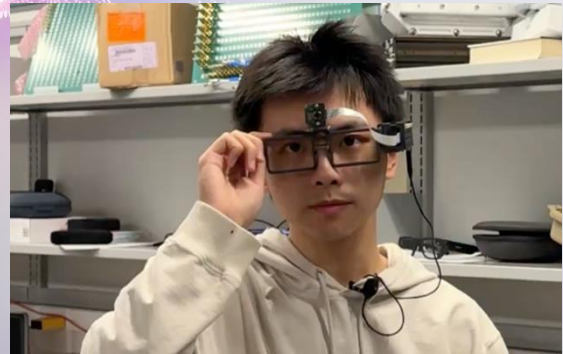
Efficient Cross-Accelerator Inference



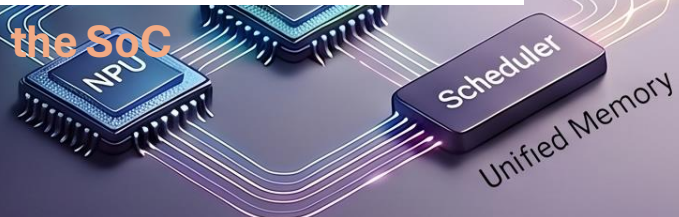
Use more accelerators on the SoC



Custom Low-Power Hardware Design

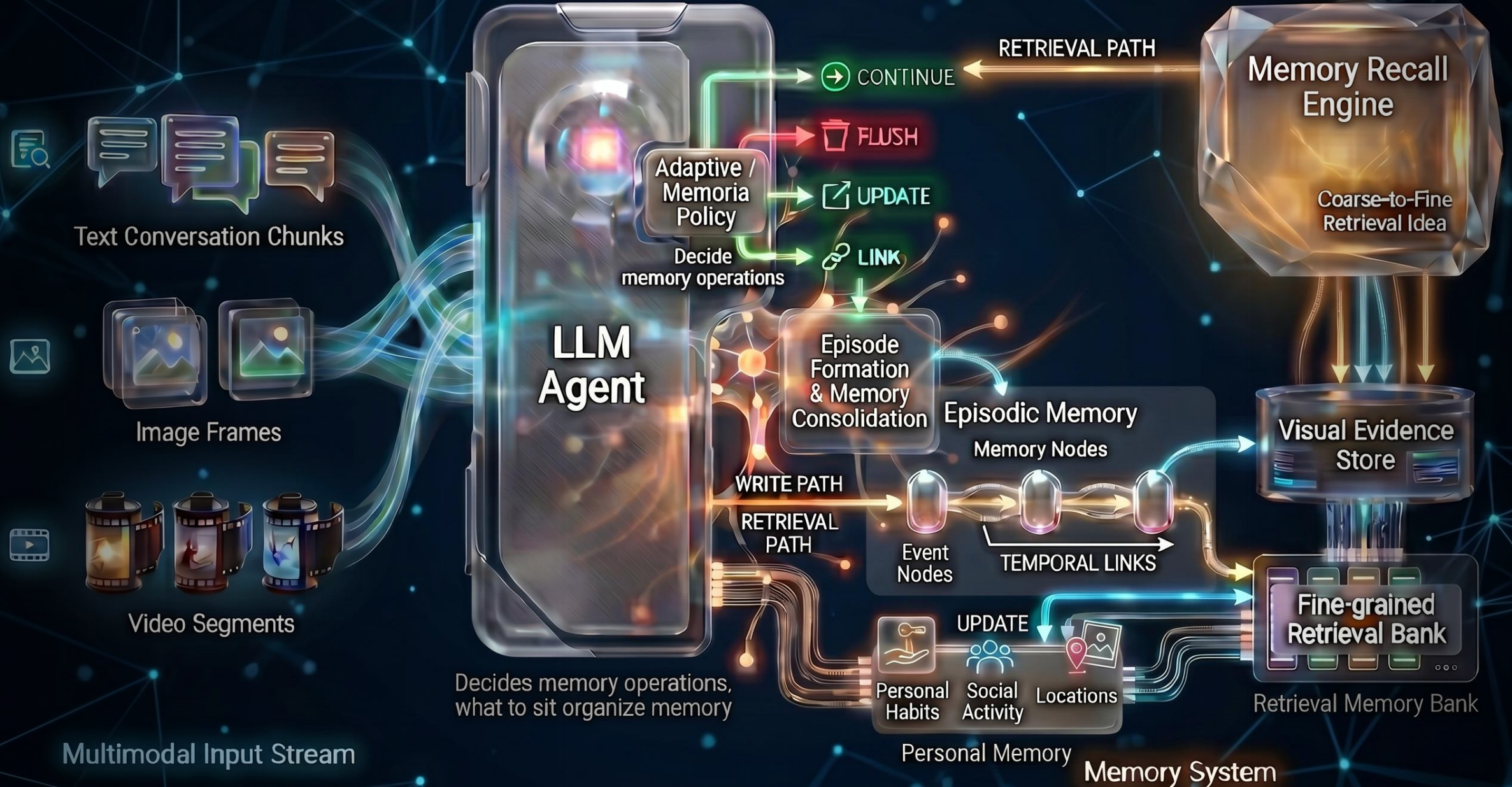


Prototype with Auto-focus Camera



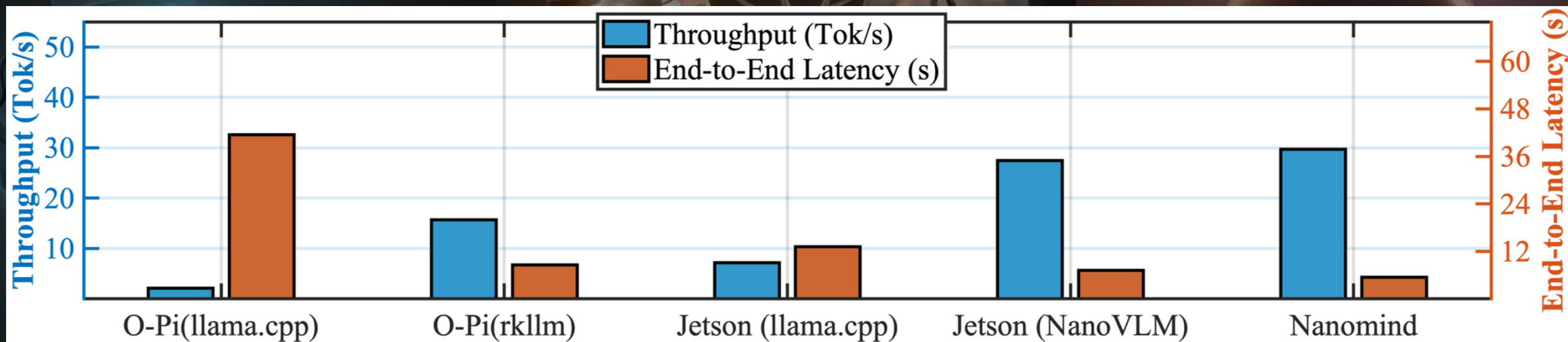
Private Memory

Memoria: Self-Evolving Personal Agentic Memory



Evaluation

Throughput (Tok/s)



Throughput (tokens/s) and end-to-end latency (s) for Qwen2-VL-2B-Instruct

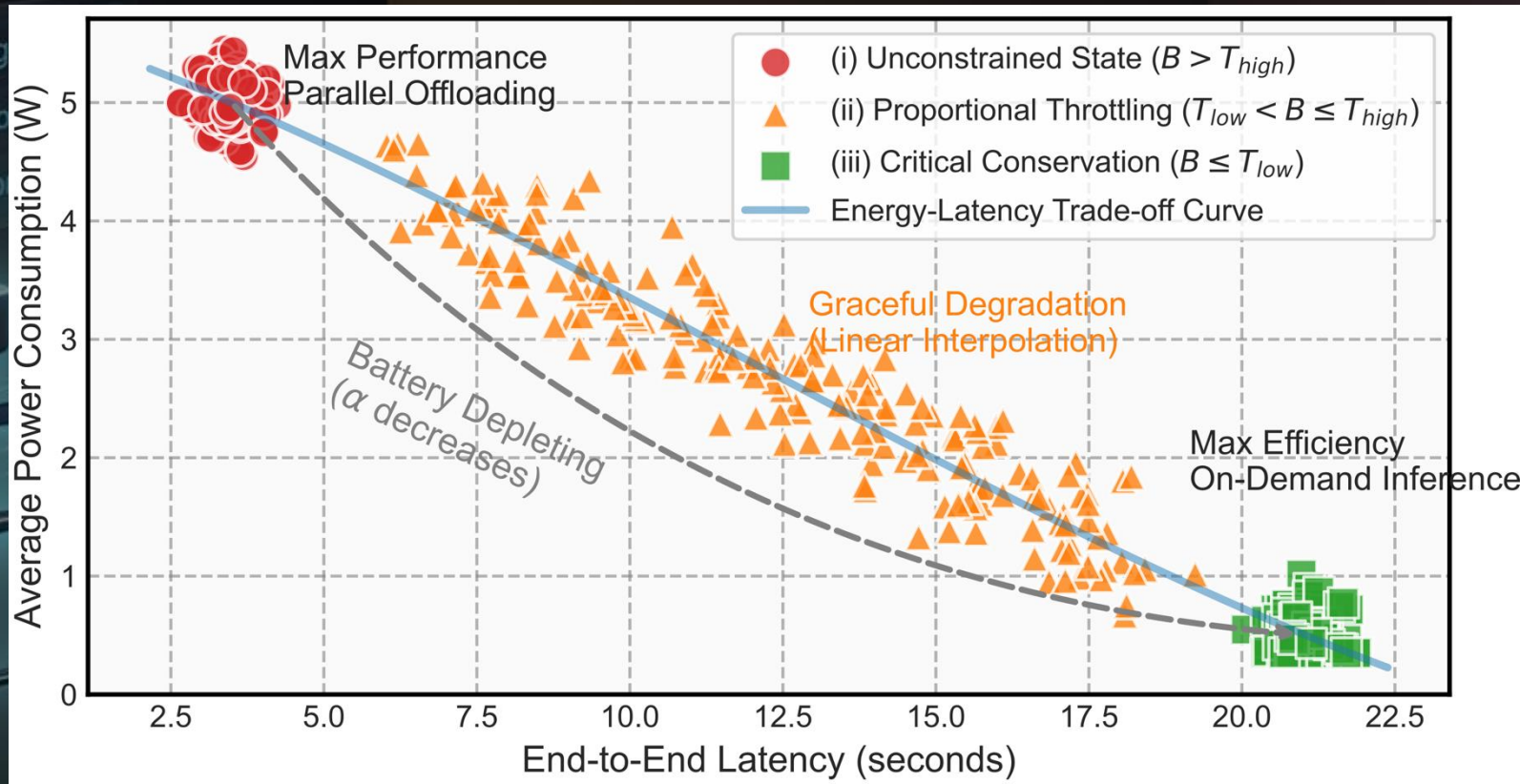
Efficient On-Device AI

Privacy-Preserving Personal Memory

Cognitive-State-Aware Assistance

Evaluation

Energy-Latency Trade-off Across Different Modes



Energy-Latency Trade-off Across Three Power Modes. The curve illustrates how the system adapts to the battery level (B). (1) In the Unconstrained State, parallel acceleration delivers low latency at higher power. (2) In the Proportional State, the system linearly throttles frame rate and memory bandwidth as B decreases, producing a continuous latency-power trade-off trajectory. (3) In the Critical State, the system transitions to the low-power On-Demand Cascade pipeline.

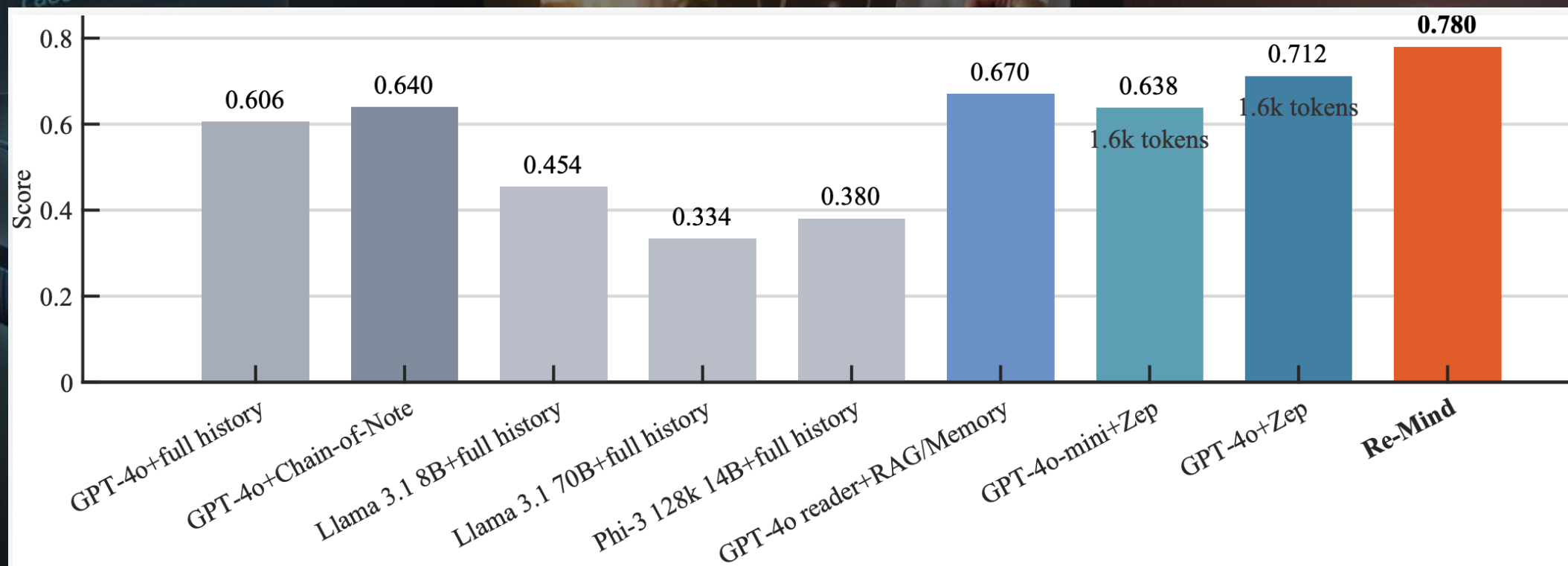
Efficient On-Device AI

Privacy-Preserving Personal Memory

Cognitive-State-Aware Assistance

Evaluation

Evaluation of Long-Term Memory



LongMemEval-S long-context Evaluation Scores

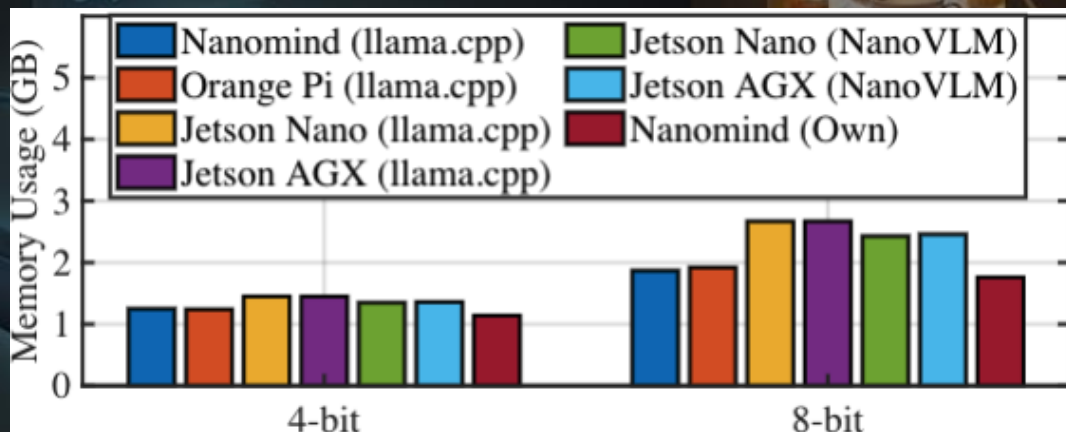
Efficient On-Device AI

Privacy-Preserving Personal Memory

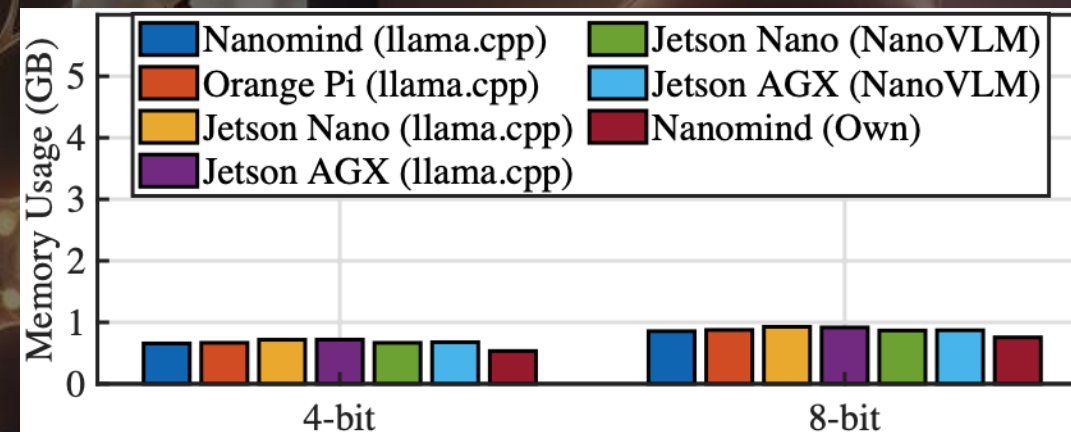
Cognitive-State-Aware Assistance

Evaluation

Memory Usage



Memory Usage (GB) while running Qwen2-VL-2B



Memory Usage (GB) while running Llava-onevision-0.5B