



Scan Here!

Project Page

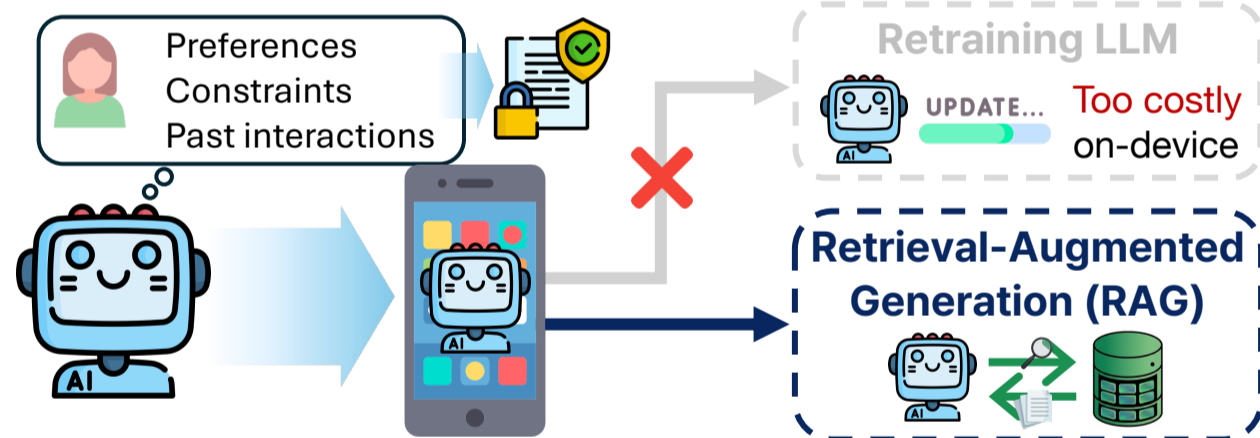
Key Takeaway

For on-device personalized RAG, the key is not to store more, but to store what matters.

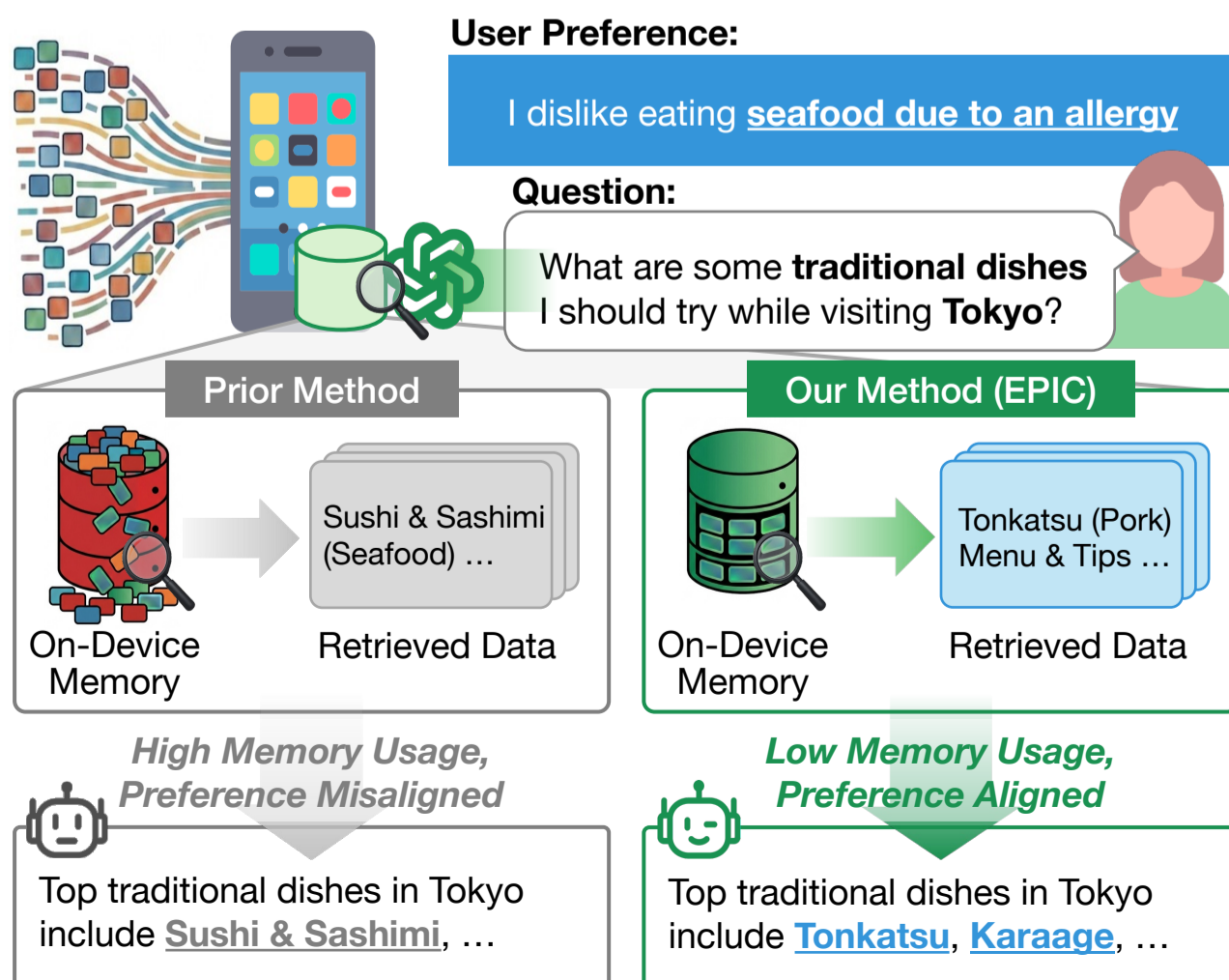
How can we build practical personal agents that run entirely on-device?

Motivation

Personalization needs private, on-device memory

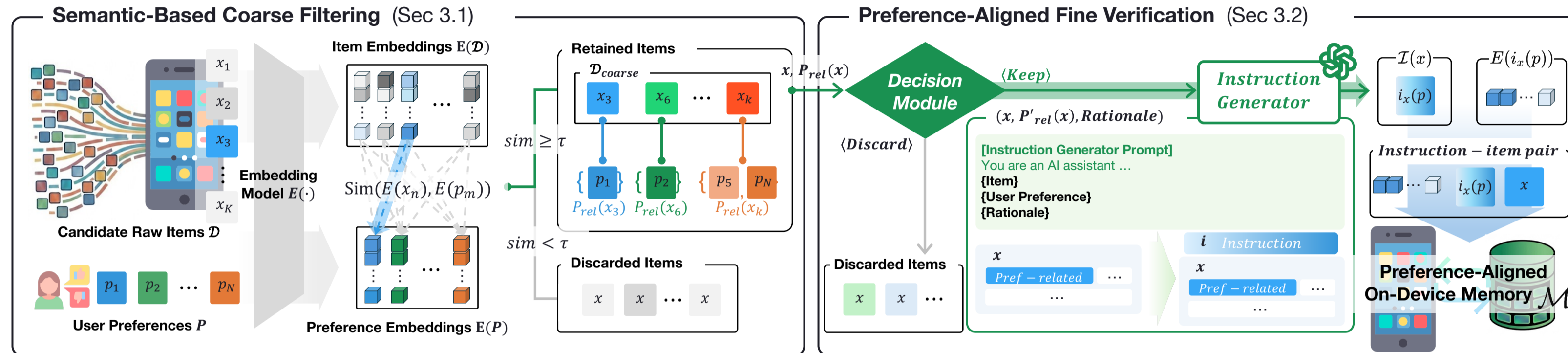


From "how to use" to "what to store"



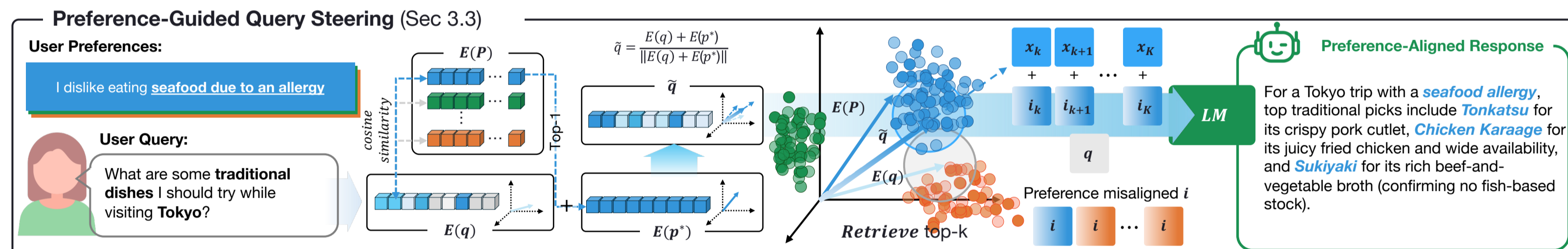
Our Method: EPIC: Efficient Preference-aligned Index Construction

EPIC works in three steps



Retain Preference-related candidates

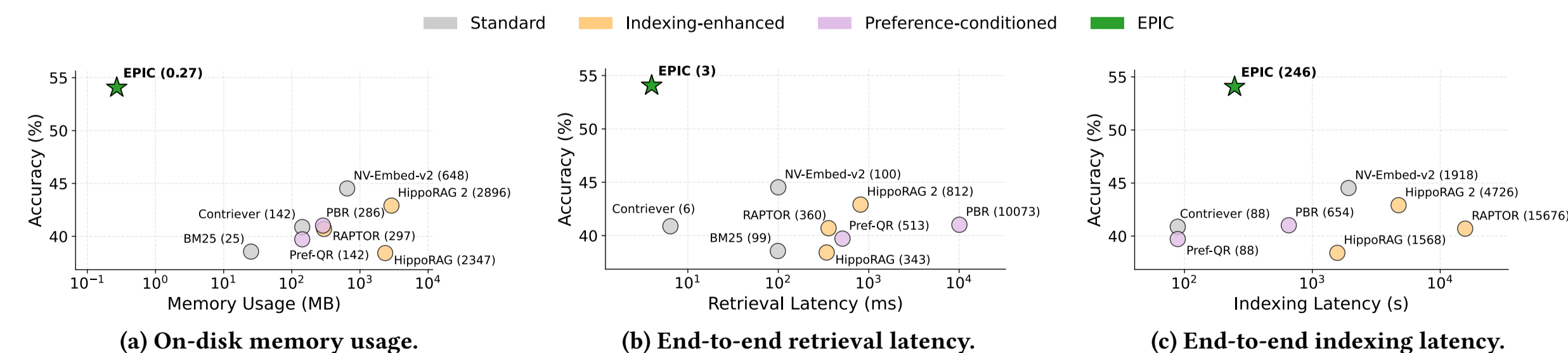
Verify and generate preference-guided instructions



Steer retrieval toward the relevant preference

Experiment Results

Server-Side Results



2,404x
Less Memory

33.33x
Faster Retrieval

7.80x
Faster Indexing

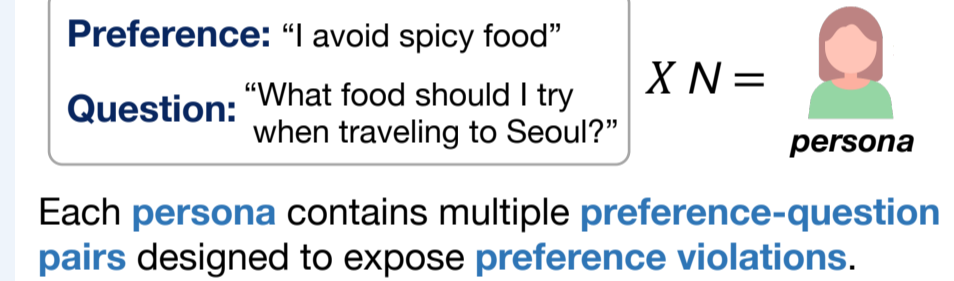
On-Device Results

Component	Jetson Orin Nano	MacBook Pro M4	Galaxy Z Flip 6
<i>Retrieval Latency (ms)</i>			
Embedding query	29.03	3.36	3.85
Preference-Guided Query Steering (Sec. 3.3)	0.18	0.03	0.55
FAISS retrieval	0.14	1.84	0.81
Total Retrieval Latency	29.35	5.23	5.21
<i>Indexing Latency (ms)</i>			
Embedding items	29.99	4.31	4.01
Semantic-Based Coarse Filtering (Sec. 3.1)	0.02	0.01	0.10
Preference-Aligned Fine Verification (Sec. 3.2)	70.73	47.09	245.65
Embedding instruction	0.02	0.01	0.00
Build FAISS	0.00	0.00	0.00
Total Indexing Latency	102.67	51.42	249.76

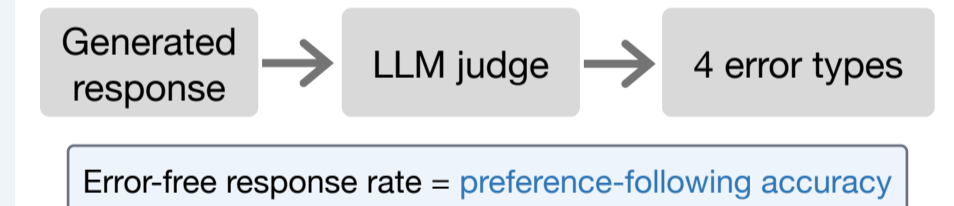
Benchmarks

Diverse Corpora	ours	existing
1. Static Knowledge	PrefWiki (Recommendation)	Wikipedia
2. Web Footprint	PrefELI5 (Explanation)	Common Crawl
3. Conversation Logs	PrefEval ^[1] (conversation)	LMSYS-Chat

Per-Persona Indexing



Evaluation



Reference [1] "Do LLMs Recognize Your Preferences? Evaluating Personalized Preference Following in LLMs". ICLR 2025.