
Optimizing Feature Extraction for On-device Model Inference with User Behavior Sequences

Chen Gong, Zhenzhe Zheng, Yiliu Chen, Sheng Wang, Fan Wu, Guihai Chen

Shanghai Jiao Tong University & ByteDance Client AI

2026-05-12



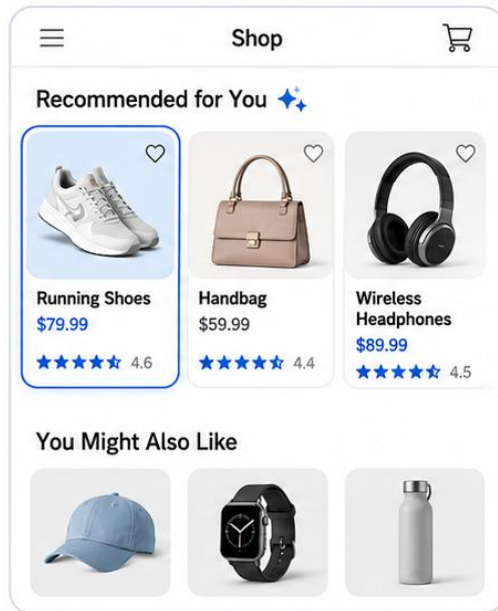
The Rise of Mobile ML



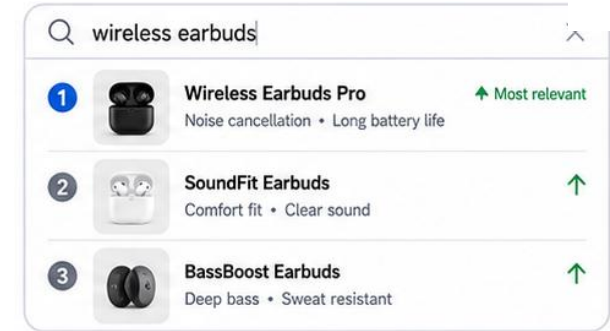
ML models are deeply integrated into modern mobile apps



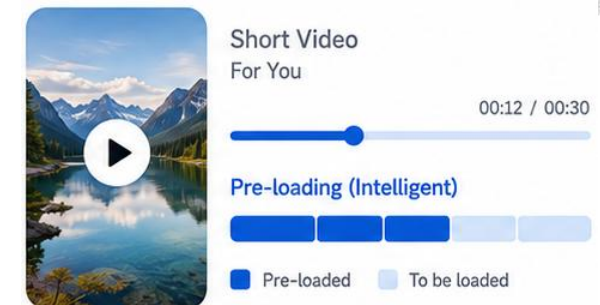
Product Recommendation



Search Ranking



Video Preloading



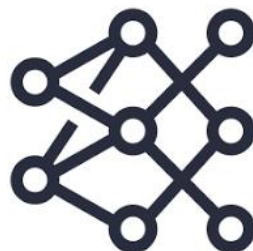
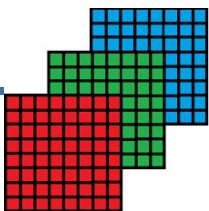
Dynamic vs. Static Features



Traditional CV/NLP models: static input features



Image Pixels



Token Embeddings

Token String	Token ID	Embedded Vector
"The"	101	[0.5, 0.8, ..., 0.65]
"cat"	2598	[0.1, 0.3, ..., -0.82]
"sat"	4450	[0.41, -0.53, ..., 0.8]
"."	119	[0.17, 0.75, ..., 0.2]



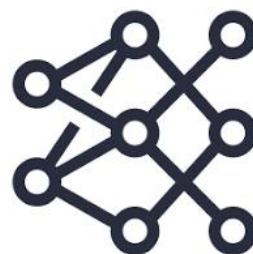
On-device ML models: dynamic user features

T-2: Search

T-1: Watch

T: Favorite

T+1: ?

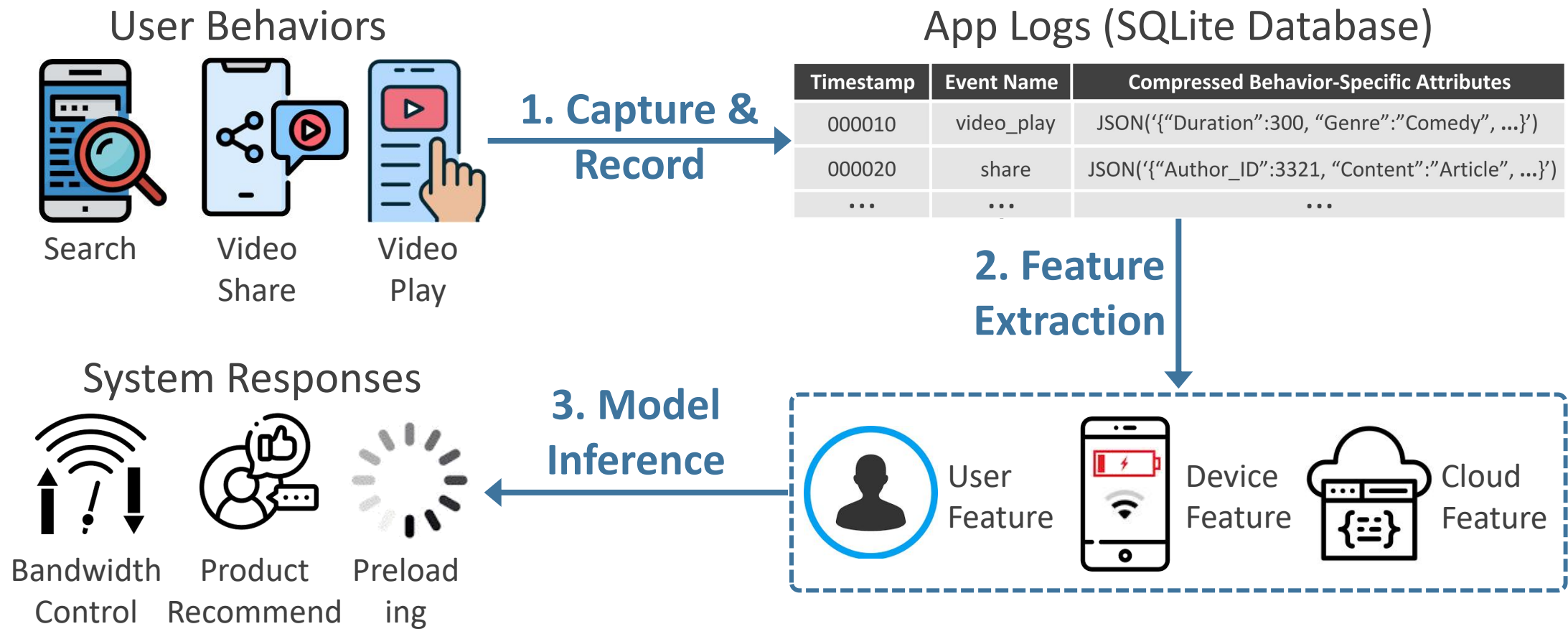


Evolving User Behaviors

On-Device Workflow



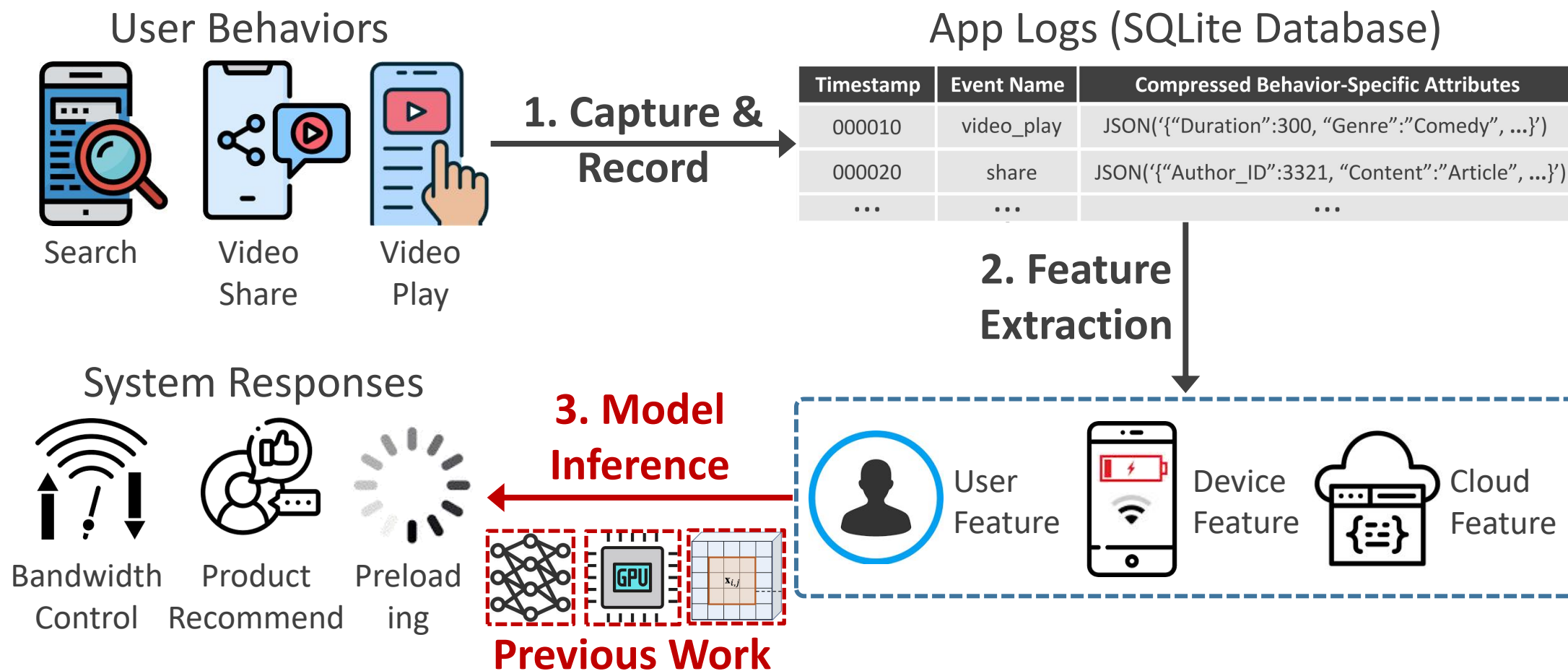
User Behavior → App Logs → Features → App Response



Focus of Prior work



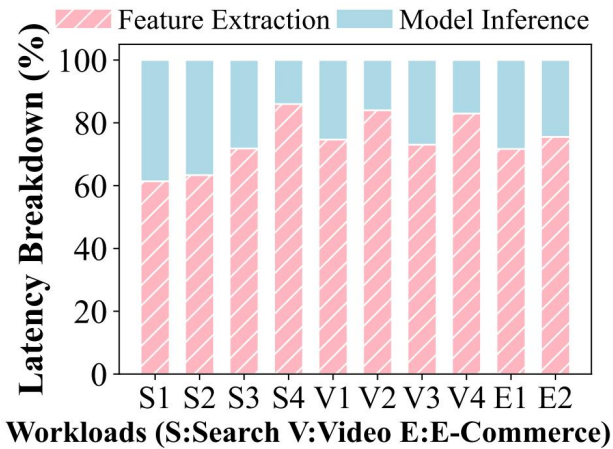
Optimize model inference stage.



Overlooked Bottleneck



Feature extraction accounts for 61% - 86% of execution time



↑ ≤ 39%:
Inference

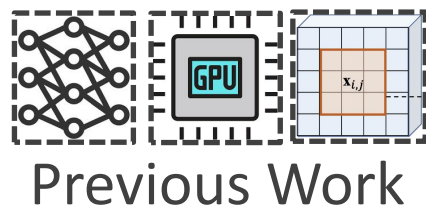
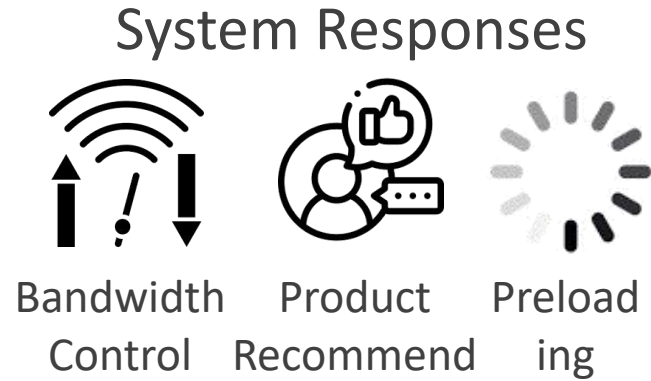
↓ ≥ 61%:
Feature Extraction

App Logs (SQLite Database)

Timestamp	Event Name	Compressed Behavior-Specific Attributes
000010	video_play	JSON({"Duration":300, "Genre":"Comedy", ...})
000020	share	JSON({"Author_ID":3321, "Content":"Article", ...})
...

2. Feature Extraction

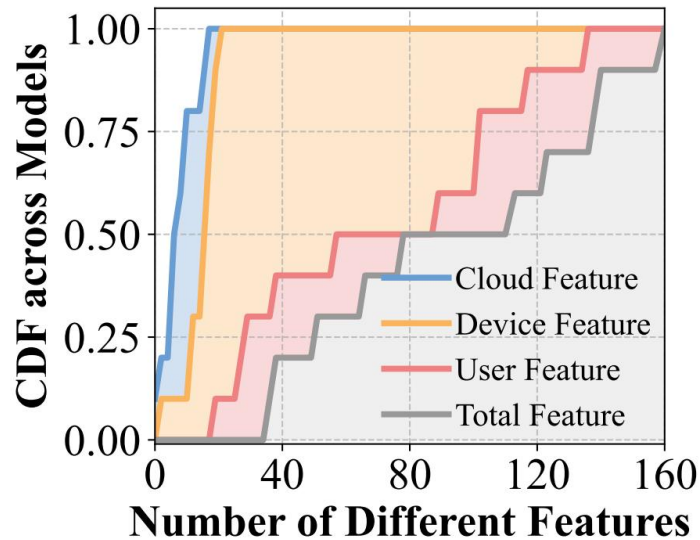
3. Model Inference



Why is extraction so slow?



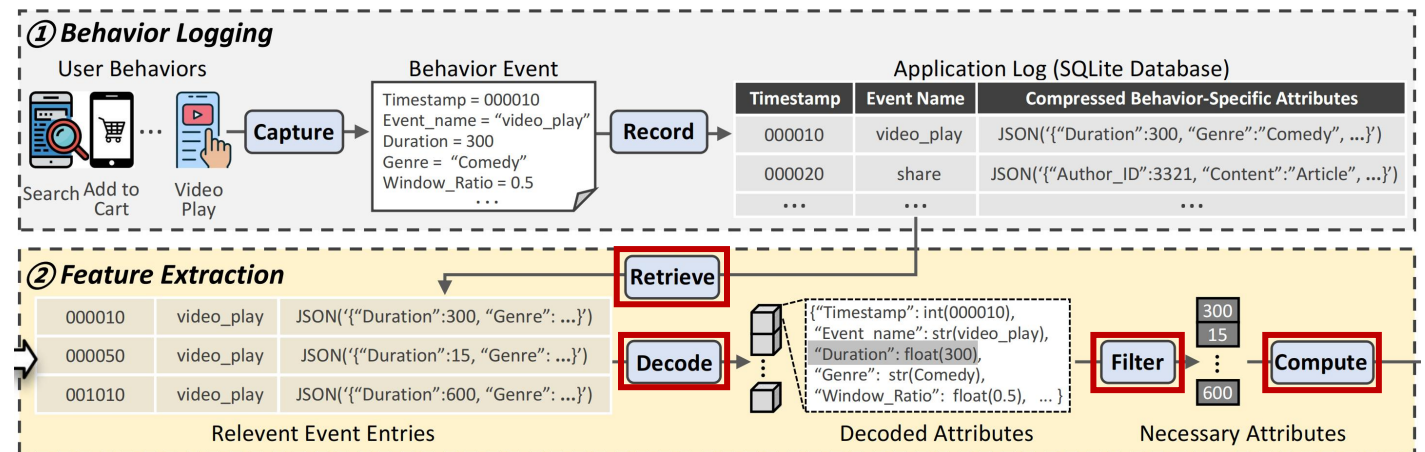
1. High proportion of user features



On average, **74%** of input features are user features.

2. Cumbersome feature extraction

- Retrieve, Decode, Filter, Compute



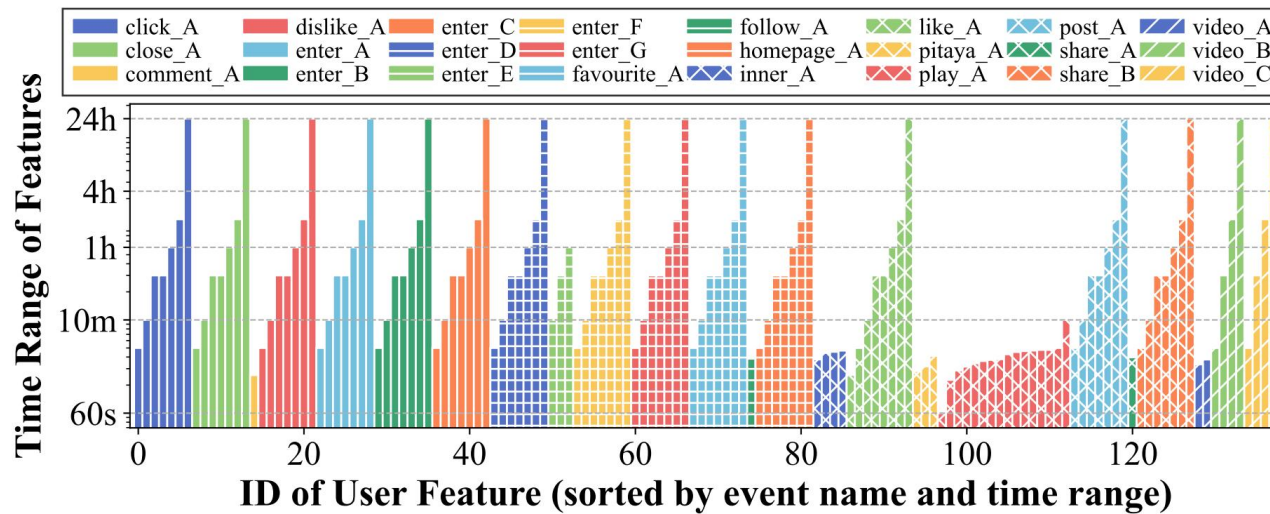
3. Fast inference: ms-level latency

- Model size: MB-level
- Mature optimization

Optimization Opportunity

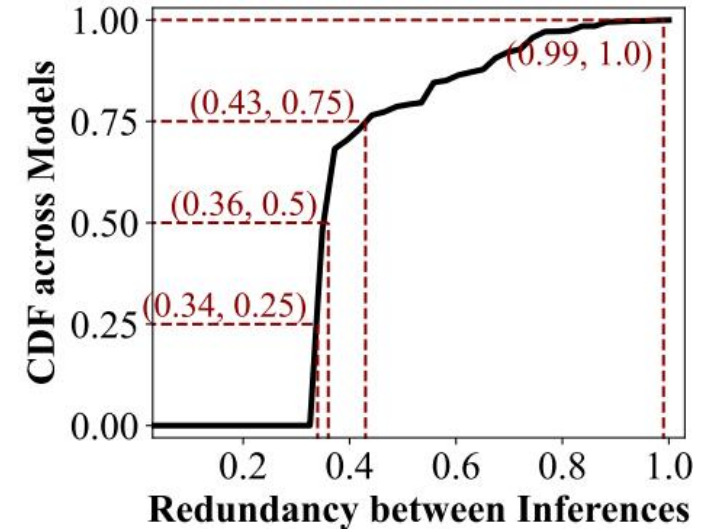


Key Observation: Many user features are extracted from overlapping user behavior data in app logs.



Inter-Feature Redundancy

- **134 features** extracted from **24** distinct behavior types with different time ranges



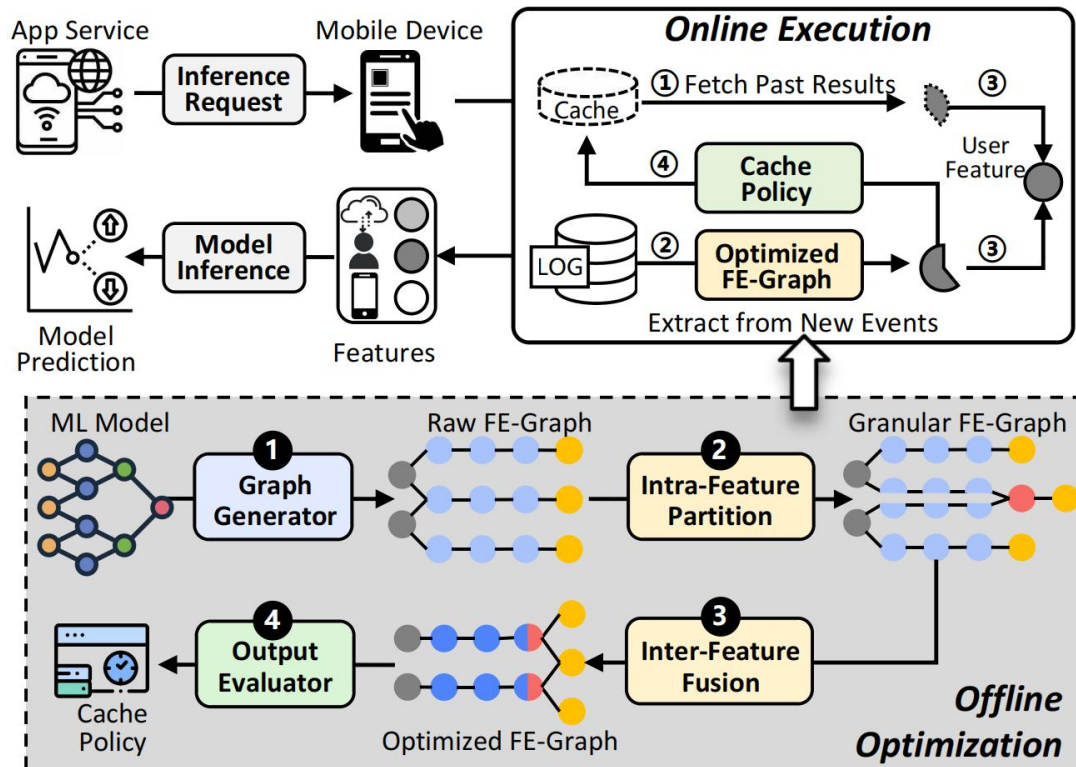
Cross-Inference Redundancy

- **75%** online model inferences have \geq **43%** overlapping data

AutoFeature Overview



Accelerate on-device feature extraction without compromising model inference accuracy.



Main Challenges:

#1: How to systematically **identify** redundancy?

2: How to effectively eliminate **inter-feature** redundancy?

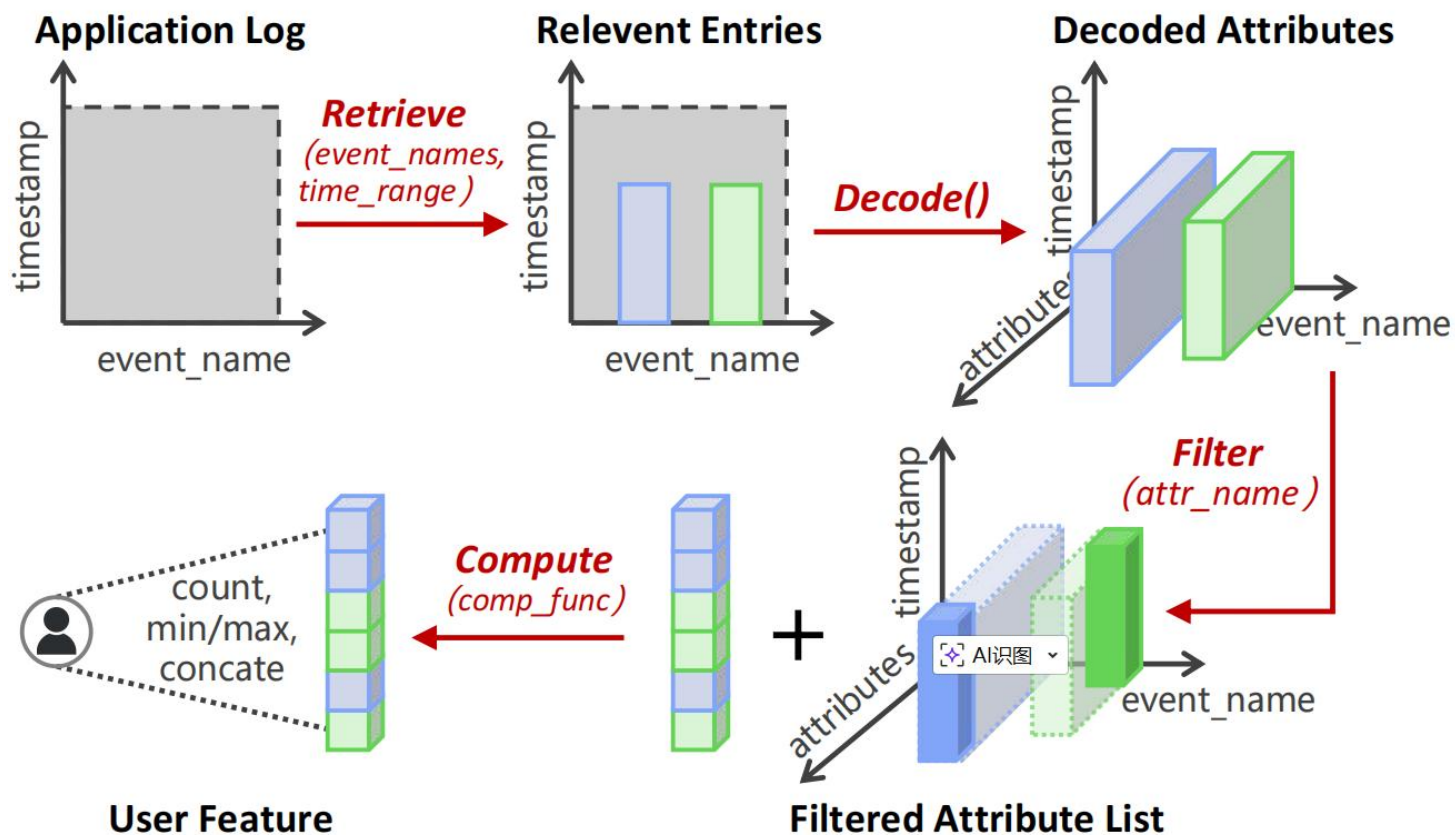
3: How to efficiently minimize **cross-inference** redundancy?

Solution 1 - Graph Generator



How to systematically identify redundancy?

- Model feature extraction workflow as a DAG



Source Node: App Log Data

Target Node: Each Feature

Connection: 4 atomic operation nodes

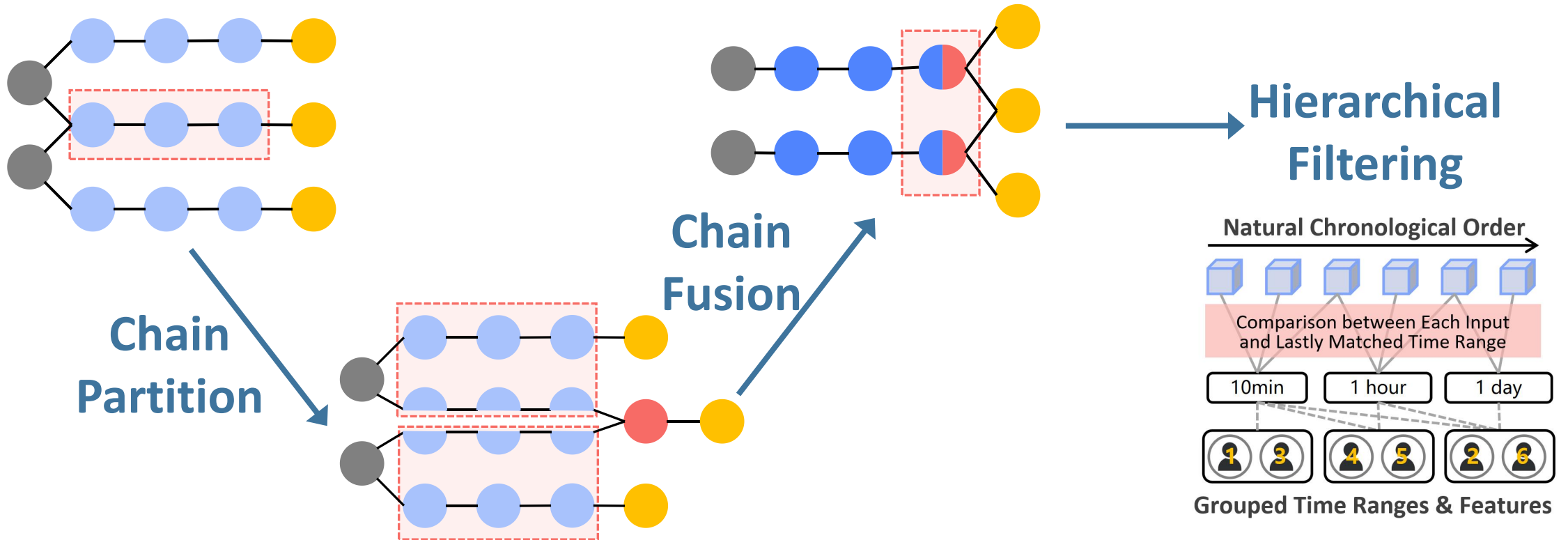
Quantify all redundancy via condition intersection !

Solution 2 - Graph Optimizer



How to effectively eliminate inter-feature redundancy?

- How to fuse? Chain Partition & Fusion
- How to separate? Hierarchical Filtering



Solution 3 - Efficient Caching



How to efficiently minimize inter-inference redundancy?

- Formulate as a knapsack packing problem

Target: Maximize savings

$$\mathcal{P}^* = \arg \max_{\mathcal{P}_i \in \{0,1\}} \sum_{i=1}^N [\mathcal{P}_i \times U(E_i)]$$

$$\text{s.t.}, \sum_{i=1}^N [\mathcal{P}_i \times C(E_i)] \leq M,$$

Constraint: device memory

Greedy Policy:

Based on saving-to-cost ratio

$$\begin{aligned} \frac{U(E_i)}{C(E_i)} &\stackrel{(a)}{=} \frac{\text{Time_Overlap}(E_i) \times \text{Freq}(E_i) \times \text{Cost_Opt}(E_i)}{\text{Time_Range}(E_i) \times \text{Freq}(E_i) \times \text{Size}(E_i)} \\ &= \underbrace{\frac{\text{Time_Overlap}(E_i)}{\text{Time_Range}(E_i)}}_{\text{Dynamic Term}} \times \underbrace{\frac{\text{Cost_Opt}(E_i)}{\text{Size}(E_i)}}_{\text{Fixed Term}}, \end{aligned}$$

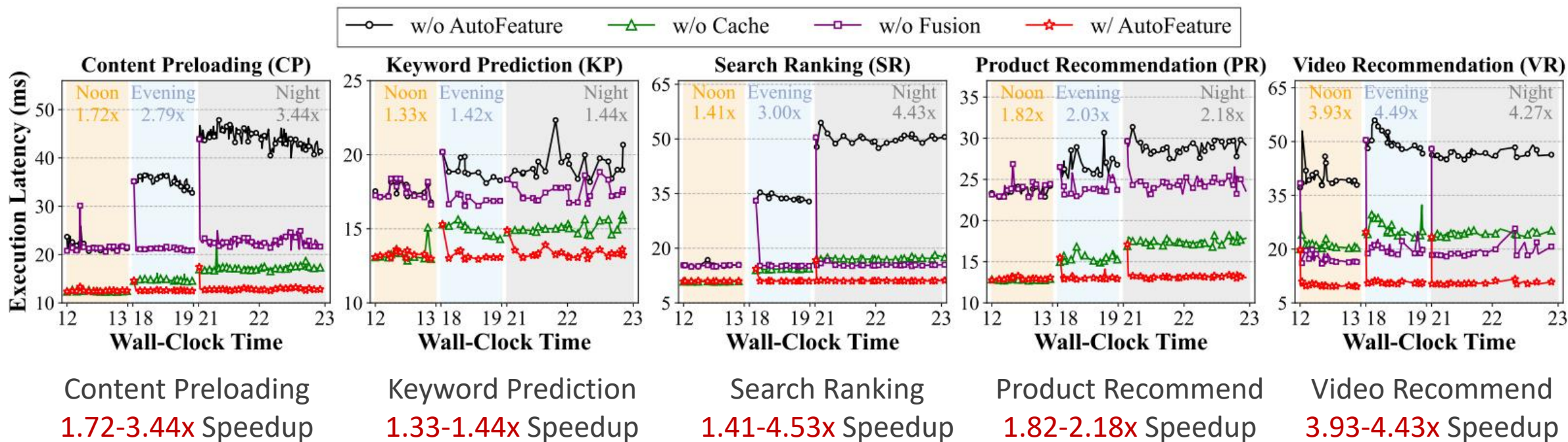
Dynamic Term: $O(1)$ Time
Fixed Term: Precompute

Evaluation & Results



5 Industrial Services & Real-world Users

- Consistently reduce model execution latency to ≤ 20 ms
- **Daytime: 1.33-3.93x** speedup, **Night: 1.44-4.43x** speedup



Conclusion



Problem:

Unveil feature extraction bottleneck in on-device model inference

Solution:

AutoFeature optimizes feature extraction without compromising model inference accuracy

Result:

Remarkable speedups in online evaluation of industrial services.

Thank You for Your Attention !

Chen Gong
gongchen@sjtu.edu.cn

