



# MobileRAG: The Future of On-Device Retrieval-Augmented Generation

A Fast, Memory-Efficient, and Energy-Efficient Pipeline for Privacy-First AI

Presented by: AI Research Team

Based on: *MobileRAG: A Fast, Memory-Efficient, and Energy-Efficient Method for On-Device RAG* (Park et al., KAIST, 2025)

# The Mobile Challenge: Constraints for On-Device RAG

Running Retrieval-Augmented Generation (RAG) directly on mobile devices is crucial for safeguarding user privacy, especially with personal data like photos and messages. However, this approach faces significant hurdles due to the inherent hardware limitations of mobile platforms compared to powerful server-based implementations.



## Memory Wall

Mobile devices typically offer 4–12GB of RAM, much of which is consumed by the operating system. Storing massive vector indices, such as those used by HNSW (Hierarchical Navigable Small Worlds) graphs, in this limited RAM is simply impractical.



## Power Drain

The continuous, CPU-intensive computations required for distance calculations during retrieval and the subsequent large language model (LLM) generation rapidly deplete battery life. This also leads to thermal throttling, severely impacting performance and user experience.



## Latency Cycle

"Time to First Token" (TTFT) is a critical metric for user experience. Standard RAG pipelines often feed overly large inputs (e.g., 2,000 tokens) to the language model, resulting in unacceptably slow inference times and a degraded user experience.

# Why Current RAG Solutions Fall Short on Mobile

Existing RAG approaches, while effective in server environments, present critical flaws when adapted to the constrained mobile landscape. They fail to address the unique challenges of limited memory, battery, and processing power.

Naive RAG	Feeds top-k documents directly to LLM.	<b>Slow Generation:</b> Uses full 2K-token inputs.	High inference latency, poor UX.
Advanced RAG	Uses a secondary Re-Ranker model for precision.	<b>Heavy Computation:</b> Re-Ranker adds significant overhead.	Increased memory and processing demands.
EdgeRAG	Optimizes indexing (IVF-DISK) to save RAM.	<b>Inefficient Generation:</b> Still feeds full documents to LLM.	Drains power during the critical generation phase.

Key Insight: We need a solution that optimizes *both* the retrieval (Index) and the generation (Input Size) phases for true mobile efficiency.

# Introducing MobileRAG: A Unified On-Device Solution

MobileRAG is a novel, fully on-device pipeline specifically engineered to overcome the limitations of mobile AI. It integrates two innovative components that work in tandem to optimize both retrieval and generation stages, delivering privacy-first AI directly to your device.



## Two Core Innovations:

### EcoVector (Retrieval Engine)

A hybrid RAM-Disk vector search architecture designed to minimize memory footprint without compromising search speed. It intelligently manages index storage and retrieval for optimal performance.

### SCR (Generation Optimizer)

Selective Content Reduction, a sophisticated method to filter and compress retrieved text content *before* it reaches the Small Language Model (sLM). This significantly reduces the LLM's input size.

## MobileRAG's Impact:

- Achieves 100% offline operation, eliminating network dependency.
- Ensures robust user privacy by keeping all data on-device.
- Drastically reduces battery consumption, enabling practical daily use.

# EcoVector Architecture: Divide and Conquer

EcoVector redefines on-device vector search through its intelligent hybrid RAM-Disk strategy. It partitions large vector spaces into manageable clusters, optimizing for both speed and minimal memory footprint.

## Cluster Partitioning

Vectors are grouped into distinct clusters using a technique like k-means. Each cluster is represented by a central point, or "Centroid," simplifying the overall search space.

## Dual-Graph Structure

**Centroids Graph (RAM):** A small, efficient HNSW graph comprising only the cluster centroids. This graph is stored entirely in RAM, allowing for lightning-fast identification of the most relevant clusters for a given query.

## Inverted Lists Graph (Disk)

The actual data vectors within each cluster are organized into independent HNSW graphs. These larger, detailed graphs are stored on disk. Only the necessary cluster's graph is loaded when required, significantly saving RAM.

# EcoVector: The Strategic RAM-Disk Approach

EcoVector's innovative RAM-Disk strategy directly addresses the memory constraints of mobile devices by intelligently managing where and when index data is accessed. This approach minimizes the memory footprint while delivering comparable or superior performance.



## ☐ Partial Loading for Efficiency

Unlike traditional methods that load entire, often billion-scale, indices into RAM, EcoVector only loads the specific cluster graph pertinent to the current query. This drastically reduces the active memory usage.

## ☐ Leveraging Modern Mobile Storage

While disk I/O is inherently slower than RAM, modern mobile flash storage (e.g., UFS 4.0) offers high speeds. EcoVector capitalizes on this, making disk access a viable component of the search pipeline.

## ☐ Net Energy Savings

Crucially, CPU computations consume significantly more power than disk I/O on mobile devices. By relying on smaller, clustered graphs, EcoVector reduces intensive CPU distance calculations, leading to a net reduction in energy consumption even with disk access.

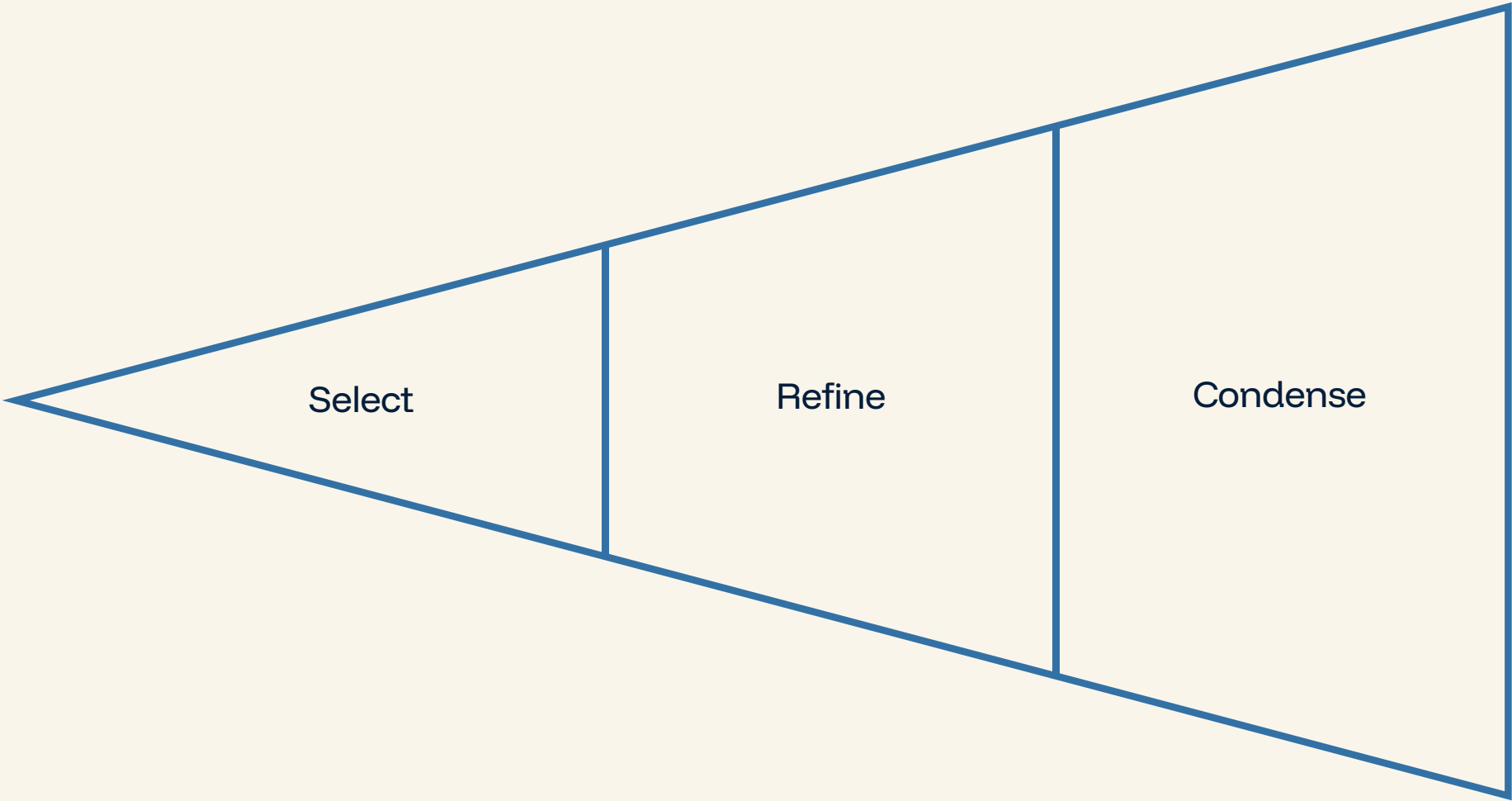
## ☐ Measurable Memory Reduction

This strategy keeps RAM usage exceptionally manageable, often reducing memory consumption from hundreds of megabytes to less than 10MB for the indices. This frees up vital resources for other applications and the LLM itself.



# Selective Content Reduction (SCR): Optimizing Generation

Retrieving the most relevant documents is only half the battle. The costly process of reading and processing these documents by the LLM is where SCR delivers significant optimization, tackling the latency and power consumption head-on.



## 1. Similarity Computation

The retrieved document is first segmented into smaller, sliding windows, typically individual sentences. Each of these segments then undergoes a re-calculation of its similarity score against the original user query.



## 2. Selecting & Merging

Only those sentence windows that achieve a high similarity score are retained. To maintain coherence and context, "context extensions"—sentences immediately preceding or following the selected windows—are strategically added to ensure logical flow.



## 3. Reordering

The now-condensed chunks of text are re-ranked based on their new, granular similarity scores. This ensures that the most relevant and coherent information is presented to the LLM in an optimized order, minimizing extraneous data.

# Why SCR Outperforms Traditional Chunking

While chunking documents into smaller pieces is a common strategy, it often comes with significant drawbacks. MobileRAG's Selective Content Reduction (SCR) offers a superior approach, preserving context and accuracy without heavy computational overhead.

## Limitations of Traditional Approaches:

### → Context Loss with Naive Chunking

Simply breaking documents into small, fixed-size chunks frequently severs vital contextual connections. This often leads to fragmented information and ultimately results in lower quality or inaccurate answers from the LLM.

### → Inefficiency of Model-Based Compressors

Utilizing heavy model-based summarizers, such as BERTSUM, for text compression introduces significant computational burden. These models are often too large and resource-intensive for mobile devices, and their use can lead to substantial degradation in accuracy.

## The MobileRAG Advantage:

### → Full Context Retrieval

MobileRAG initially retrieves documents with their full context intact, ensuring high recall and reducing the risk of missing critical information during the initial search phase.

### → Lightweight Local Refinement

Instead of heavy model inference for compression, MobileRAG employs lightweight vector mathematics for local refinement and content reduction. This approach is highly efficient for on-device processing.

### → Significant Input Size Reduction

The SCR process effectively reduces the LLM input size by **42%** for datasets like SQuAD and **31%** for TriviaQA, all while maintaining **zero accuracy loss**.



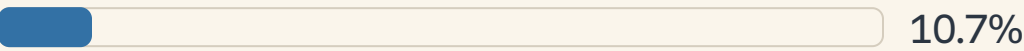
# Empirical Results: Speed and Memory Efficiency

Rigorous benchmarking on a Samsung Galaxy S24 (Exynos 2400, 8GB RAM) demonstrates MobileRAG's significant advancements in search latency and memory footprint, proving its efficacy for on-device RAG.



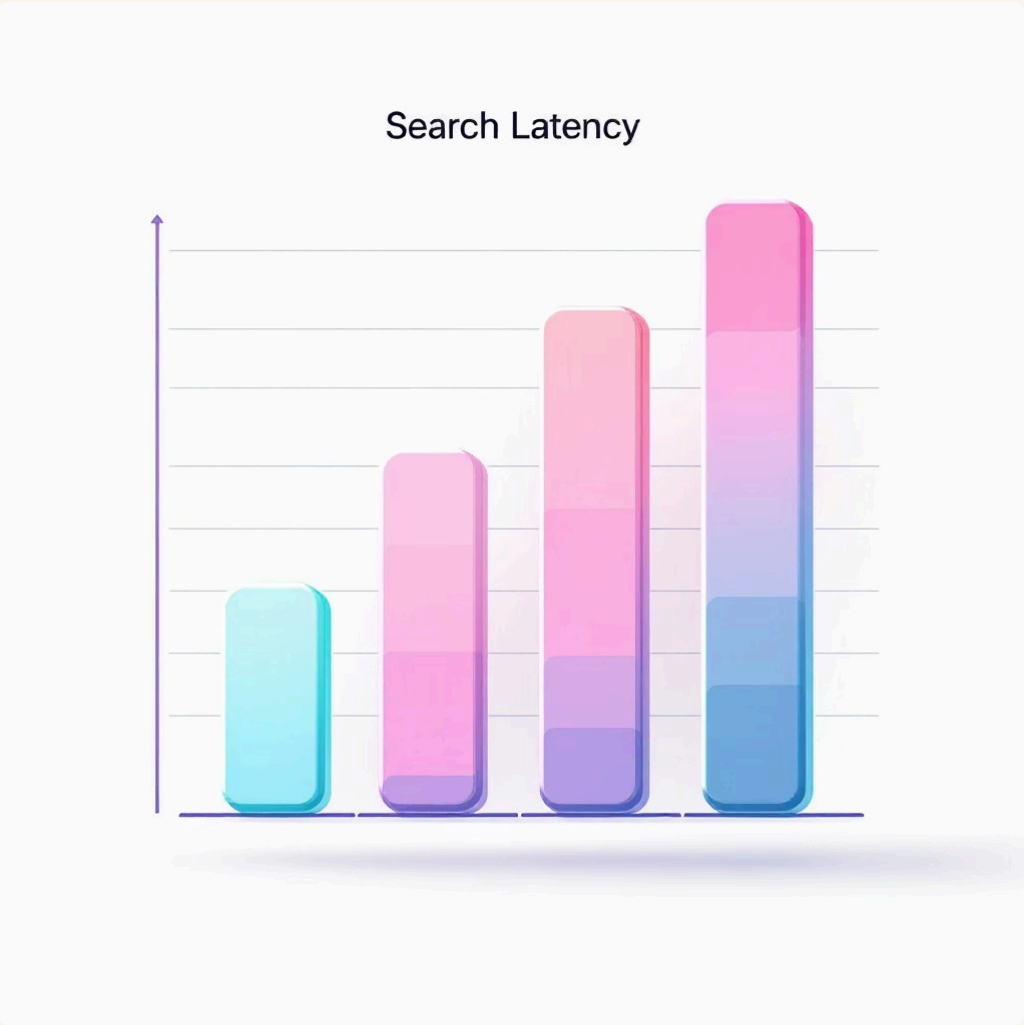
## Faster Search

EcoVector is **1.72x - 8.89x faster** than standard IVF/HNSW baselines, leading to quicker retrieval times and a more responsive user experience.



## Memory Reduction

MobileRAG reduces memory consumption by an impressive **10.7% - 54.5%**, freeing up critical RAM for other system operations and larger sLLMs.



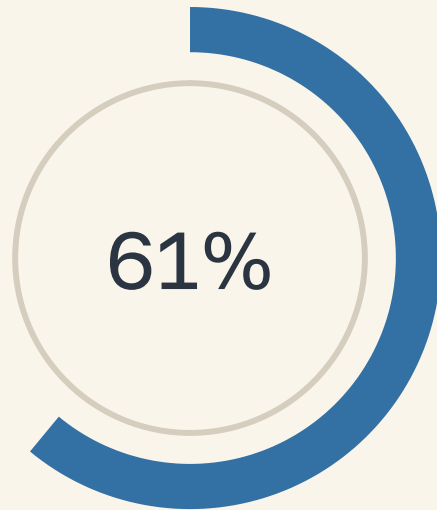
## Validation of the RAM-Disk Strategy

Even with the inclusion of Disk I/O, the strategic reduction of the search space through small, clustered graphs makes EcoVector **faster than brute-forcing** through a large, entirely RAM-resident index. This confirms that intelligently managed disk access can be more efficient than overwhelming limited RAM.

- These results highlight MobileRAG's ability to provide server-grade retrieval speed and efficiency within the constraints of a mobile device.

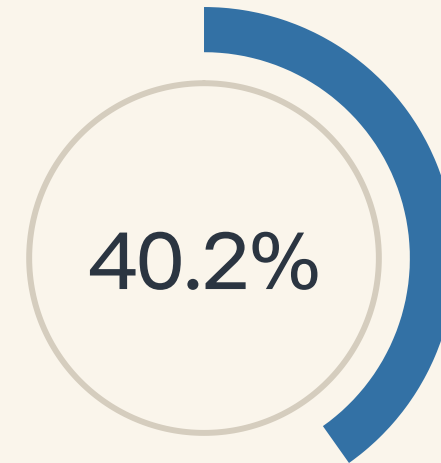
# Empirical Results: Unparalleled Power Efficiency

Power consumption is the most critical metric for mobile applications. MobileRAG delivers substantial reductions in energy usage across the entire RAG pipeline, making privacy-first AI truly sustainable for everyday mobile use.



Retrieval Power Saved

MobileRAG slashes retrieval power consumption by an average of **61-75%** when compared to EdgeRAG (IVF-DISK) implementations, optimizing the most intensive phase of the process.



Total Pipeline Power Saved

Across the full RAG process, including both retrieval and generation, MobileRAG reduces overall power usage by up to **40.2%** against advanced RAG solutions. This comprehensive efficiency is a game-changer.

## The Core Reason: Smaller Inputs, Faster Generation

The primary driver behind MobileRAG's superior power efficiency is the significant reduction in the Large Language Model's input size. Smaller inputs translate directly to faster generation times, meaning the CPU spends less time running at maximum frequency, thus conserving precious battery life.

- Less computational load on the CPU during generation.
- Reduced thermal output, preventing throttling and maintaining performance.
- Extended battery life for users, making on-device RAG a practical reality.