

# **Project Group (SS26+WS26/27)**

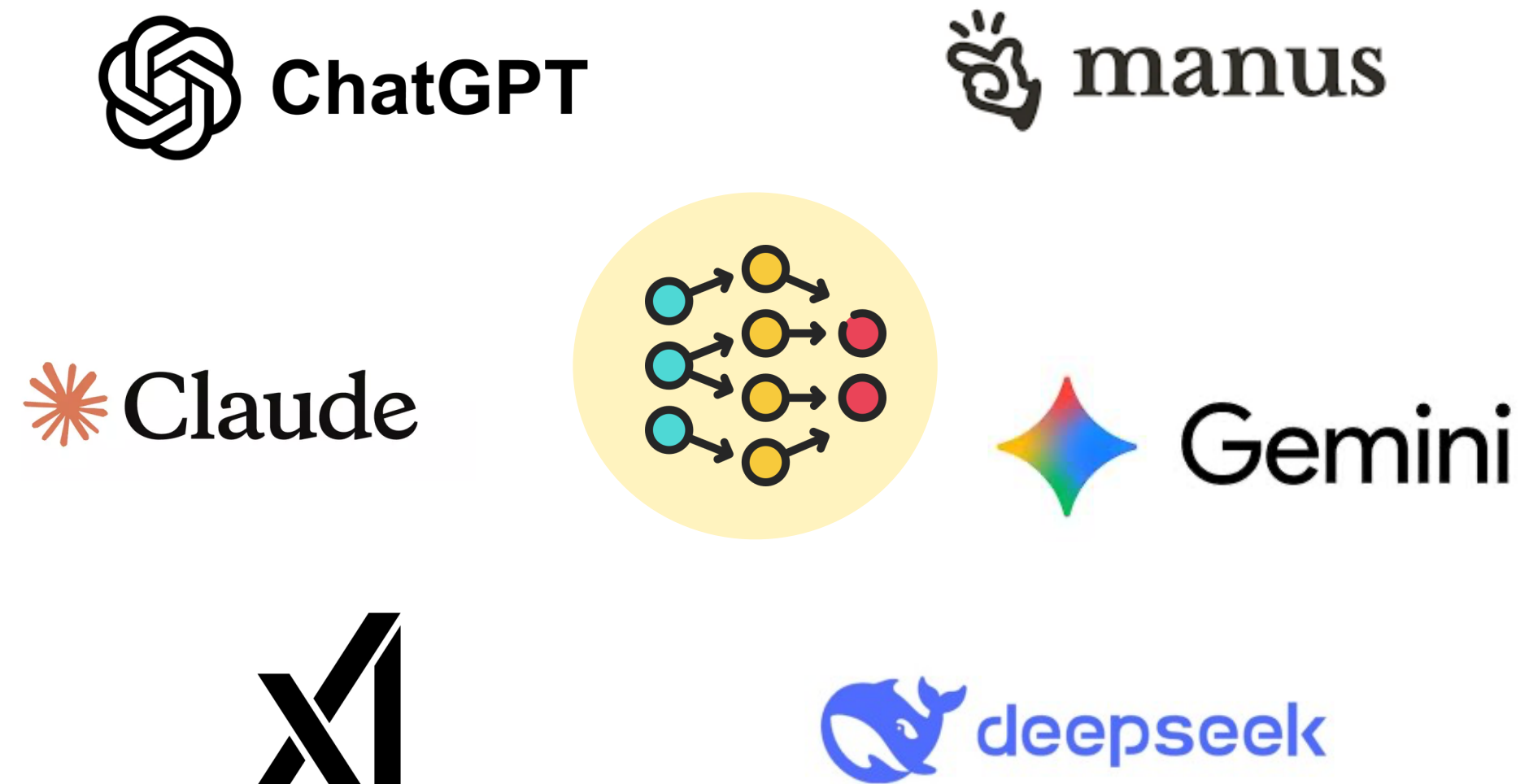
## **Benchmarking and Optimizing Mobile RAG Systems**

Prof. Dr. Lin Wang, Huzaifa Shaaban Kabakibo

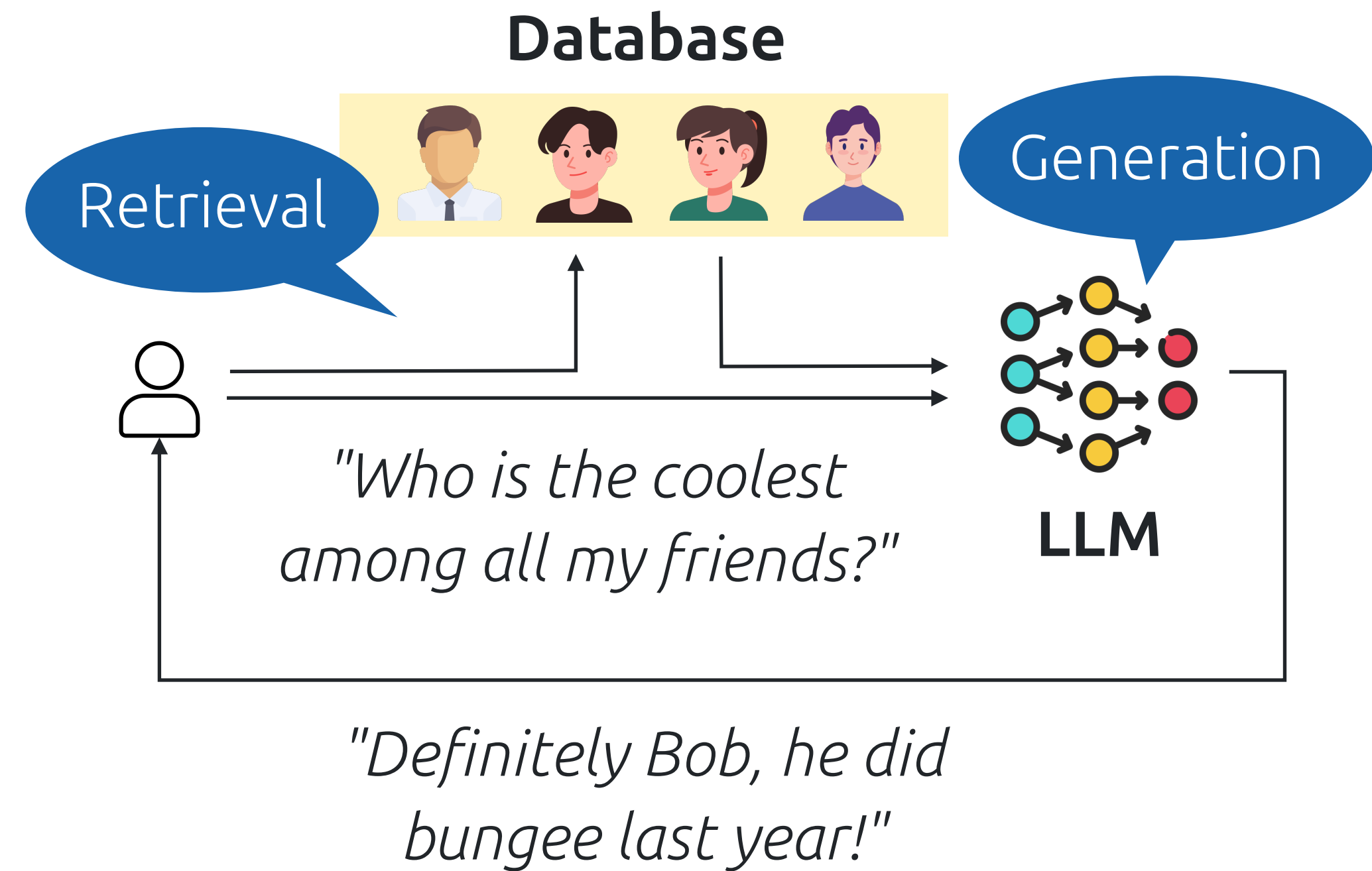
Computer Networks Group  
Paderborn University

# LLMs and RAG

Large language models (LLMs) have enabled many cool applications

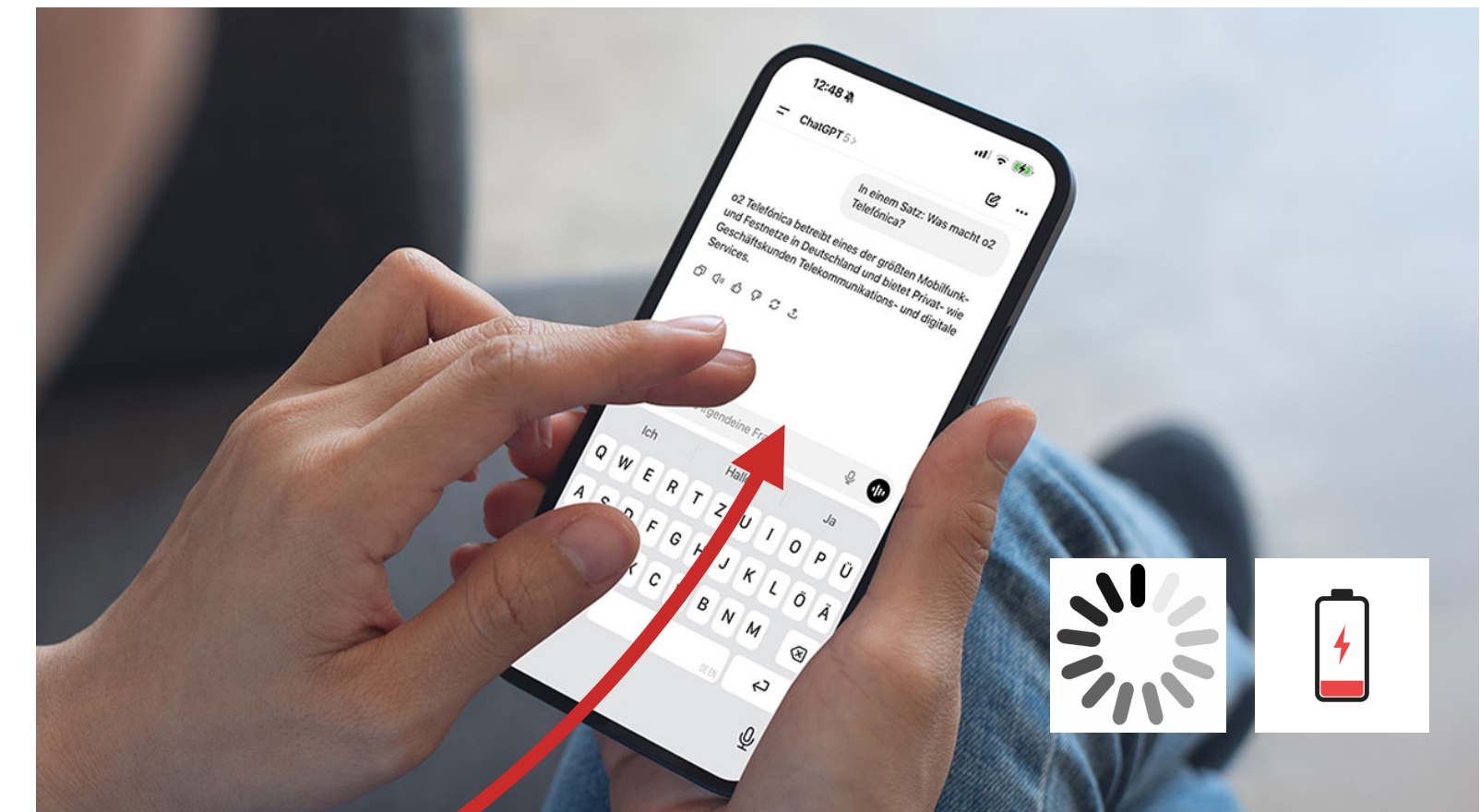
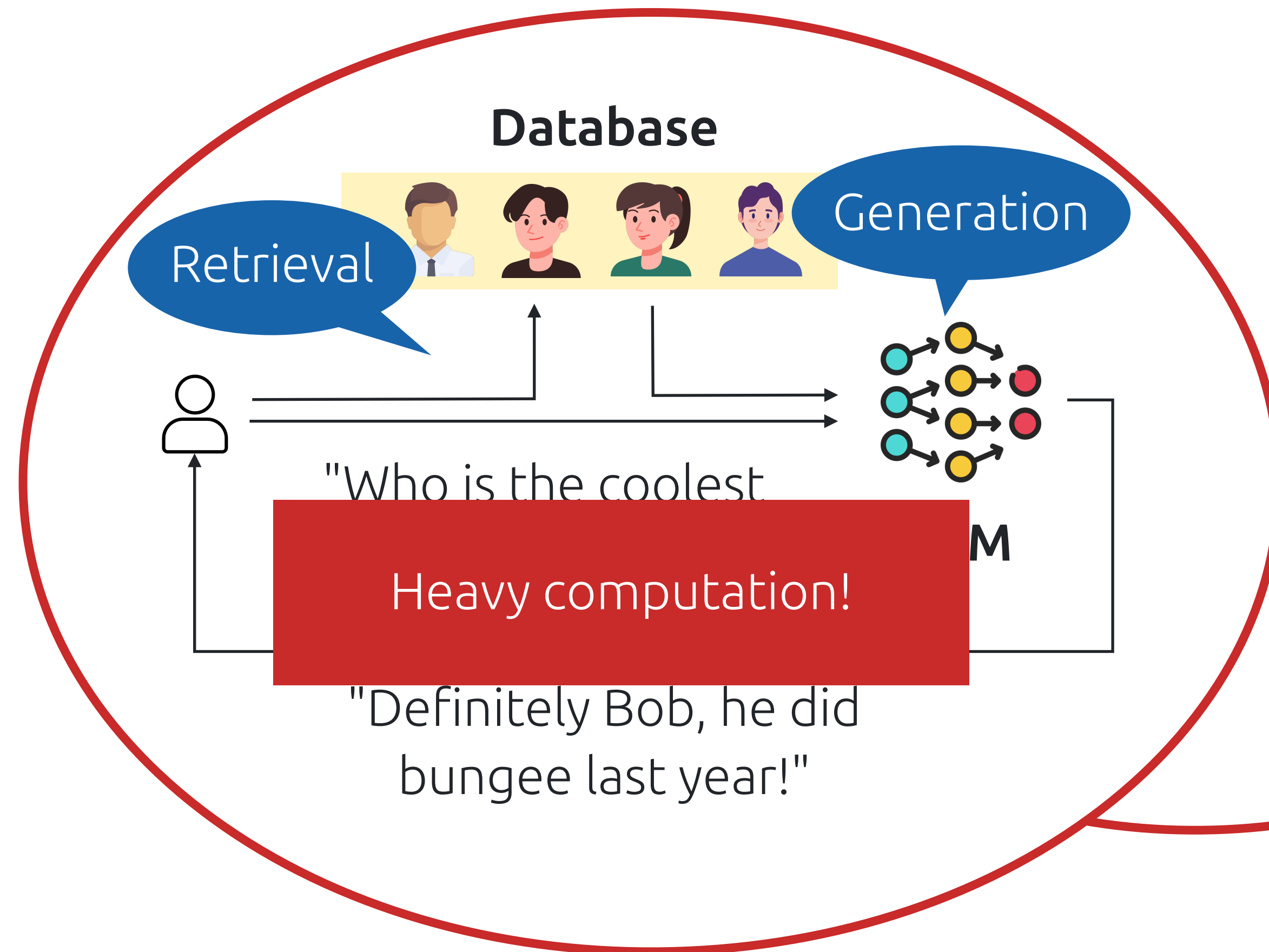


Retrieval augmented generation (RAG) makes LLMs even better!



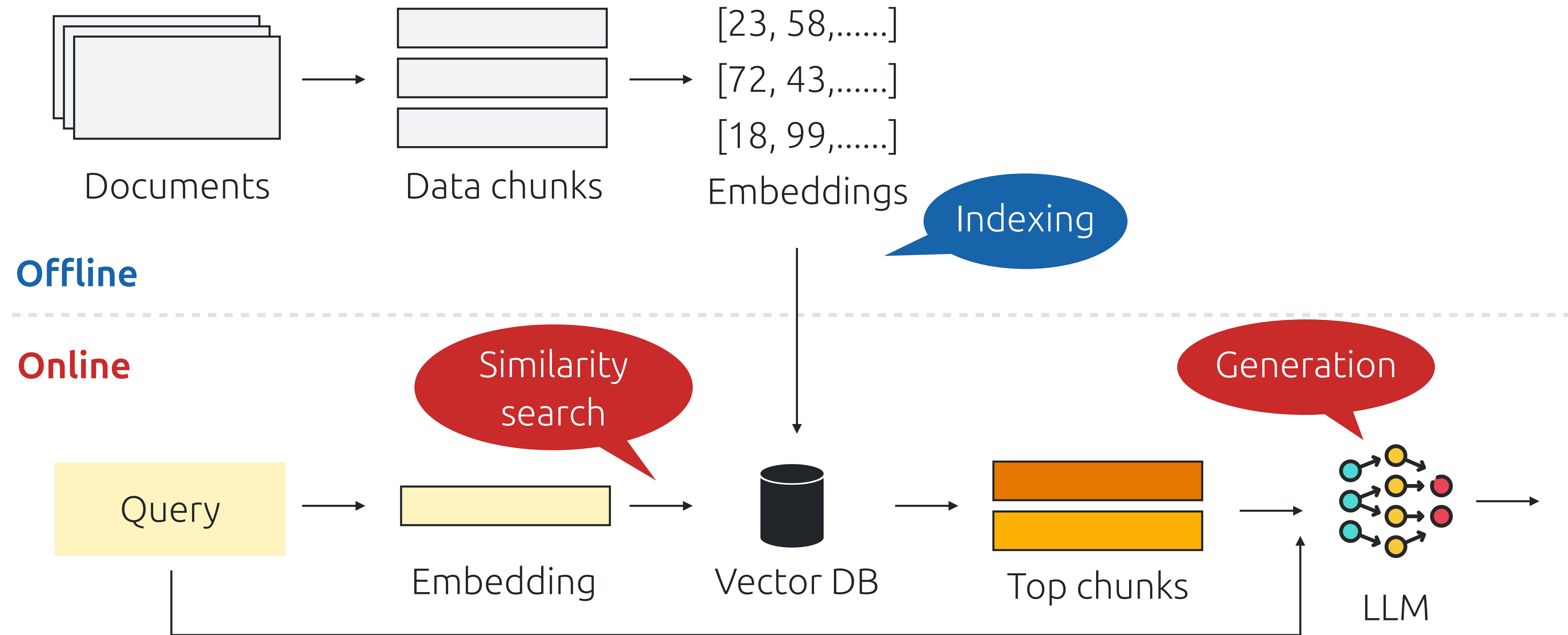
# MobileRAG

Everything on-device: **better privacy** and **more usable** under poor network

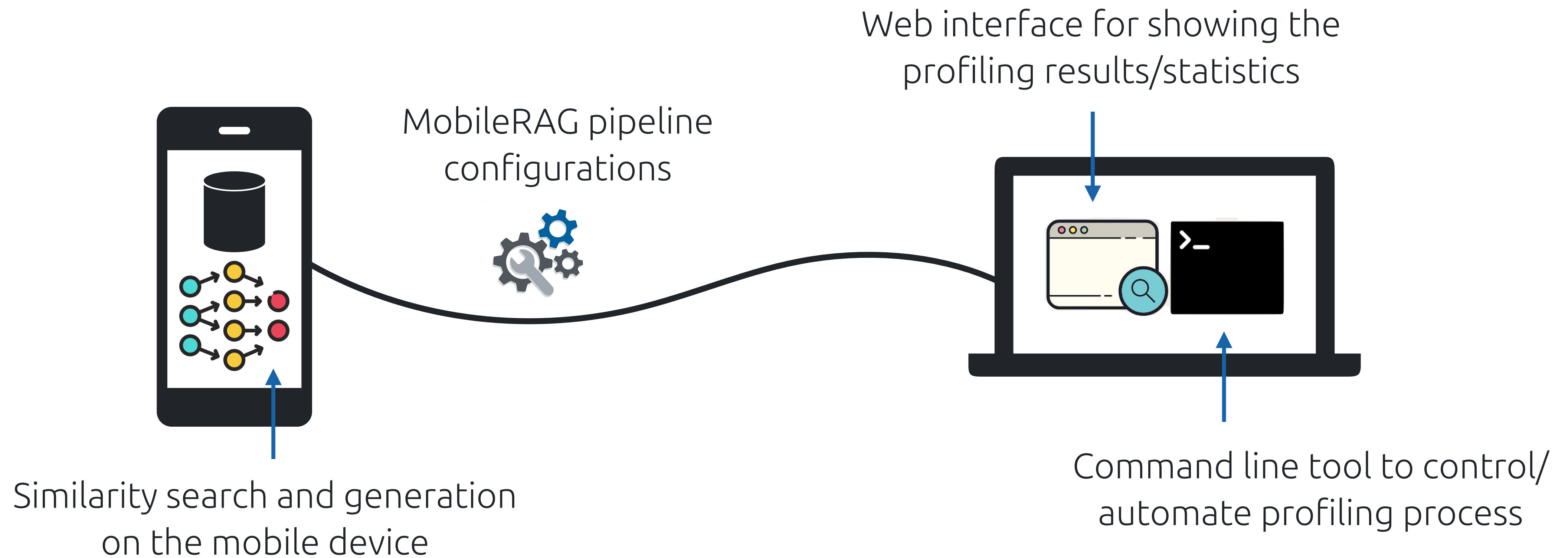


**Question: how to make  
MobileRAG fast and efficient?**

# MobileRAG pipeline



# Benchmarking framework



# Current status



## Configurable parameters

- Supports multiple datasets
- Supports different LLMs
- Supports 3 different indexing methods
- User-tunable parameters



## Automated benchmarking for different datasets

## Results generation and reporting for different metrics

- Per-component latency, throughput
- Accuracy
- Memory

```
1 {
2   "downstream_task": {
3     "name": "squad",
4     "sampling_method": "first_n",
5     "limit": 2000
6   },
7   "rag_pipeline": {
8     "embedding": {
9       "backend": "cpu",
10      "model_name": "all-minilm-16-v2",
11      "dtype": "int8",
12      "chunker": {
13        "method": "token",
14        "size": 256,
15        "overlap_enabled": true,
16        "overlap_size": 50
17      }
18    },
19    "faiss": {
20      "method": "ivf",
21      "backend": "cpu",
22      "metric": "IP",
23      "top_k": 3,
24      "config": {
25        "nprobe": 8,
26        "nlist": 16,
27        "num_training_vectors": 150
28      },
29      "use_cache": false
30    },
31    "llm": {
32      "aug_method": "concatenation",
33      "backend": "cpu",
34      "model_name": "qwen2.5-0.5B",
35      "use_sampling": false,
36      "dtype": "int8",
37      "kv_window": 4096,
38      "prefill_chunk_size": 1024,
39      "max_tokens": 96,
40      "ignore_eos": false,
41      "generate_until": ["\n", "\n\n"],
42      "system_prompt": "You are a helpful assistant."
43    }
44  }
45 }
```



# Ongoing and future directions

## Basic

- Add support for more embedding models and LLMs, more datasets, more indexing methods
- Add hyper-parameters for each component
- Add UI to the configuration and result reporting

## Intermediate

- Add augmentation methods (e.g., re-ranking, filtering)
- Add support for more metrics (CPU/GPU utilization, power consumption)
- Intensive benchmarking on multiple mobile phones

## Advanced

- Code migration to native languages (e.g., C)
- Add support for GPU/NPU acceleration for each component



