

# KumbhDoot: A Scale-Ready, LLM-Bounded Architecture for Mass-Gathering Public-Service Assistants

Saurabh Sakalkar, Cisco  
Abhishek Singh, MIT  
Ramesh Raskar, MIT

---

## Executive Summary

KumbhDoot is designed as a public-service assistant for a massive, multilingual event where infrastructure **cost** is a concern to support 80M users, **network connectivity** may be poor and some questions are safety-critical.

Instead of sending every query to an LLM, the system uses an escalation ladder:

1. Look for a semantically similar previously approved question and answer.
2. Classify intent by comparing the query with approximately 520 labelled examples.
3. Retrieve official information using lexical and embedding search.
4. Format a response directly from approved records.
5. Use an LLM only for low-confidence retrieval, multi-part synthesis, personalization, or live external information.

Its central object is the InfoBin: an approved information unit containing content, source, language, location, validity dates, freshness, safety rules, and an embedding.

The architecture has three broad workload bands:

- Band A: direct retrieval or cached answer; no generative model.
- Band B: retrieval plus coded orchestration across records; no generative model.
- Band C: bounded LLM synthesis over retrieved information.

It also proposes an offline SQLite bundle containing selected InfoBins, FTS5 search, and compact embeddings, with Redis and larger embeddings on the server. Overall, the system shows that **LLM use can be eliminated** in nearly all UAT (user acceptance testing) scenarios.

---

# Abstract

Mass religious gatherings such as the Kumbh Mela concentrate tens of millions of people into a single region over a few weeks, producing intense, repetitive, multilingual, and safety-critical demand for information. The default response, a conversational assistant that routes every query to a large language model (LLM), is poorly matched to this setting: it is costly at scale, slow on emergency paths, prone to hallucination on facts that can cause physical harm, and unusable when connectivity fails.

We describe KumbhDoot, an agentic pilgrim assistant for the Nashik Simhastha Kumbh Mela built on a different first principle.

- It operates on a foundational design principle that prioritizes semantic similarity over starting with an LLM.
- Generative models are invoked only in instances where similarity-based retrieval is insufficient to produce a correct answer.
- The system utilizes a "semantic cache"—an embedding-indexed store—as a single retrieval primitive, which handles intent routing, answer caching, offline lookups, and multi-agent retrieval.
- A custom three-tier agent architecture operates directly on this store, ensuring decision paths remain inspectable and avoiding the use of generic multi-agent frameworks that would trigger implicit per-step LLM calls.

We present the architecture, an analytical cost model for its per-query economics, and an honest account of where similarity is sufficient and where generative reasoning remains necessary. We argue that for bounded, high-stakes, low-connectivity public-service domains, a similarity-first and LLM-bounded design is not merely cheaper but architecturally more appropriate than an LLM-default one.

---

## 1. Introduction

A public-service assistant for a Kumbh-scale event operates under a set of constraints that rarely appear together. It must support millions of concurrent users on low-end devices, over intermittent connectivity, in several languages and in code-mixed speech, while some fraction of queries concern genuine emergencies: lost children, medical need, crowd danger, route closures. In this setting the cost of a wrong or slow answer is not a degraded user experience but a safety failure.

The prevailing pattern for conversational assistants, sending each user utterance to an LLM for understanding and generation, is attractive for a prototype and acceptable for a low-stakes demo. At event scale it introduces several coupled risks. Repeated questions, of which there are many in a pilgrimage context, incur repeated model and compute cost. Generative latency is unacceptable on emergency paths. Free generation over facts such as ghat access, medical contacts, and ritual schedules invites hallucination where hallucination does harm. And an LLM-default design can degrade in on-device settings due to hardware heterogeneity.

This paper argues that the appropriate response is to invert the usual order of operations. The foundational question of information retrieval, how to represent meaning as numbers and compare them, predates generative AI by decades and is computationally cheap, deterministic, and auditable. KumbhDoot makes that operation, cosine similarity over approved content, the default path for the large majority of queries, and treats the LLM as a bounded fallback reserved for genuinely generative or low-confidence tasks.

Our contribution is a systems contribution rather than an algorithmic one. The individual techniques we use, dense sentence embeddings, hybrid lexical-semantic retrieval, cross-encoder reranking, and embedding nearest-neighbor intent routing, are all established. The contribution is their integration into a single similarity-first retrieval substrate, the discipline of bounding LLM use behind explicit escalation gates, and a custom, inspectable agent orchestration layer, applied to a constrained public-service domain where reliability and offline resilience matter more than open-ended capability. We also contribute an explicit account of the boundary between what similarity can and cannot do, which we believe is more useful to practitioners than a uniformly favorable system description.

The remainder of the paper is organized as follows. Section 2 motivates the architectural inversion. Section 3 reviews the intellectual lineage and related industry patterns. Section 4 presents the system architecture, including the semantic cache, intent routing, the offline tier, and the three-tier agent model. Section 5 develops the analytical cost model. Section 6 examines where similarity matches generative utility and where it does not. Section 7 states limitations and the planned evaluation, and Section 8 concludes.

---

## 2. Why an LLM-Default Architecture Is Not Enough

A conversational front-end backed by an LLM treats every query as an act of generation. For a public-service system at Kumbh scale, six risks follow directly from that assumption.

**Cost.** A pilgrimage produces enormous query redundancy: the same questions about bathing dates, nearest ghats, transport, and helpline numbers recur across millions of users. Regenerating each answer with a model converts a fundamentally cacheable workload into a recurring per-query expense.

**Latency.** Emergency requests cannot wait for multi-step generation. A user reporting a lost child or seeking an ambulance must receive immediate, deterministic guidance, not a model completion.

**Hallucination.** When a model invents details about route access, medical facilities, or official schedules, the error is not cosmetic. Grounding answers in pre-approved content rather than free generation is therefore a safety control, not only a quality control.

**Connectivity.** Dense crowds degrade mobile networks. Offline behavior is not an edge case in this domain; it is a core requirement, and an LLM-default design cannot satisfy it.

**Privacy.** Location, health status, family composition, and journey stage are sensitive signals. A design that ships every utterance to a cloud model maximizes exposure of exactly the data that should be minimized.

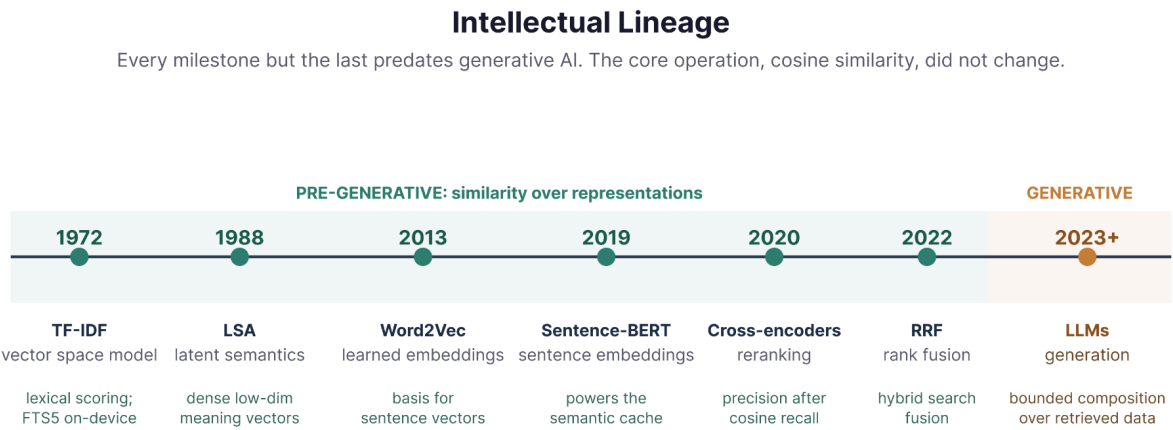
**Operations.** Administrators need aggregate operational intelligence, crowd concentration, emerging concerns, not raw individual conversations or default per-user tracking.

These risks motivate treating the assistant as trusted event infrastructure with a layered design, rather than as a conversational shell over a model. The central design principle is to use the cheapest, fastest, and safest path that can answer a query reliably, escalating to a model only when no cheaper path suffices.

---

# 3. Background and Related Patterns

## 3.1 Intellectual lineage



KumbhDoot does the heavy lifting with the green lineage and adds the LLM only as the final composition layer.

Figure 5. Six decades of representation learning; the operation that compares them stays constant.

KumbhDoot is built on a long line of work in representing and comparing meaning, most of which predates generative AI. The vector space model and TF-IDF (Salton, 1972) established lexical scoring. Latent Semantic Analysis (Deerwester et al., 1988) introduced dense low-dimensional representations. Word2Vec (Mikolov et al., 2013) learned dense word embeddings, and Sentence-BERT (Reimers and Gurevych, 2019) extended dense representation to full sentences, which is the level at which KumbhDoot operates. Cross-encoder reranking and Reciprocal Rank Fusion provide, respectively, precision after initial recall and a principled way to combine lexical and semantic rankings. Across this history the quality of the representation improved at each step while the core operation, cosine similarity, did not change. Large language models add a final capability, turning retrieved records into fluent natural language, but in our design they sit atop this lineage rather than replacing it.

## 3.2 Related industry patterns

KumbhDoot adapts several established patterns to the public-service, low-bandwidth, high-scale setting. Retrieval-Augmented Generation grounds model output in external content; KumbhDoot inverts the emphasis, answering from approved records first and invoking generation only when needed. Prompt and response caching reduce repeated model cost; KumbhDoot extends this to a paraphrase-aware semantic cache and, for common intents, avoids the model call entirely. Offline-first mobile data stores cache content and sync on reconnect; KumbhDoot extends offline behavior beyond data to include intent routing and retrieval. On-device inference runtimes make small-model execution feasible on the handset; KumbhDoot uses lightweight on-device components where practical without assuming high-end hardware. Finally, LLM observability platforms trace prompts, latency, and cost; KumbhDoot extends observability to non-LLM paths so that rule, cache, routing, and fallback decisions are equally visible.

The distinguishing position is that KumbhDoot is not attempting to be a general answer engine. It is an event-specific assistant with bounded knowledge, strict fallbacks, and offline resilience, and that scoping is what makes a similarity-first design viable.

## 4. System Architecture

### 4.1 The layered runtime

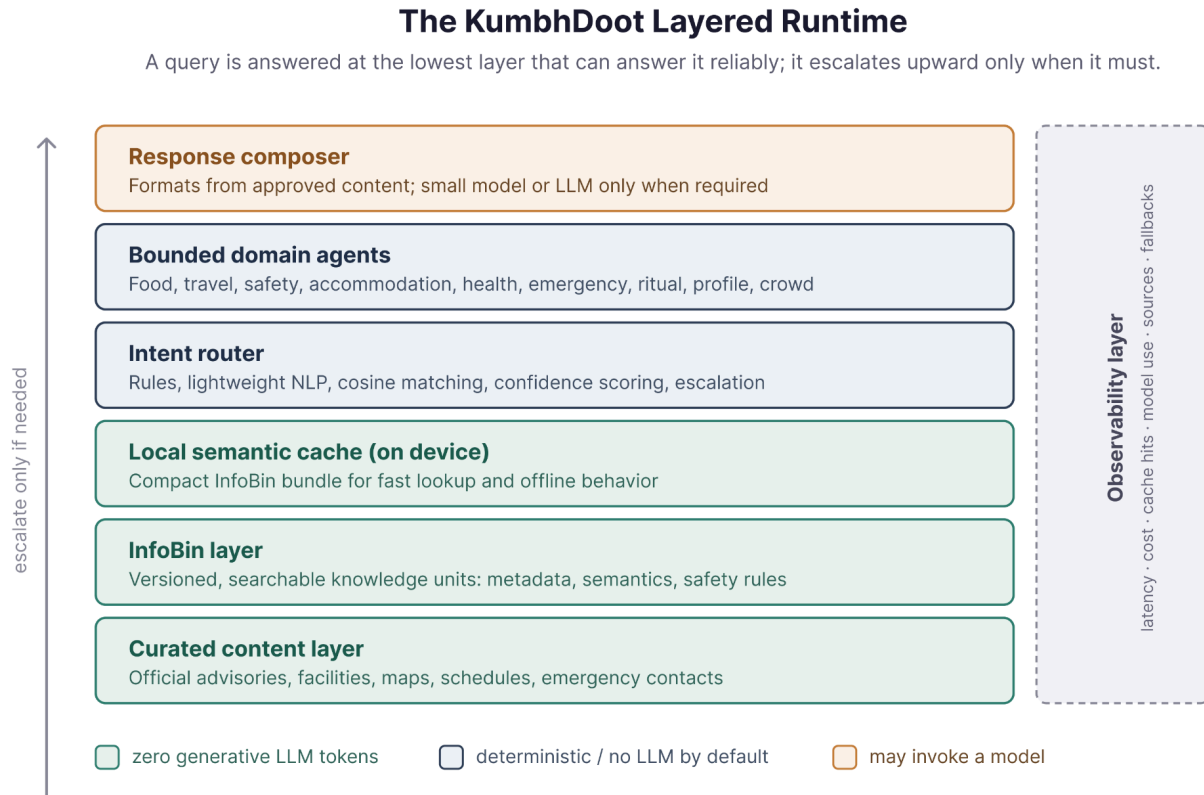


Figure 1. The layered runtime. Observability spans every layer, covering both LLM and non-LLM paths.

The KumbhDoot runtime is a layered stack in which each layer answers as many queries as it can before escalating. From the bottom up: a *curated content layer* of official advisories, facility details, maps, schedules, and emergency contacts; an *InfoBin layer* that structures this content into versioned, searchable knowledge units; a *local semantic cache* that ships a compact partition of those units to the device; an *intent router* that uses rules, lightweight NLP, and semantic matching with confidence scoring; a layer of *bounded domain agents* for food, travel, safety, accommodation, health, emergency, ritual, profile, and crowd concerns; a *response composer* that formats answers from approved content and may call a small model or LLM only when required; and an *observability layer* tracking latency, cost, cache hits, model use, sources, safety flags, and fallback paths. The design is deliberately narrower than a general-purpose agent platform: it uses agentic structure but keeps each agent bounded, auditable, and grounded in approved content.

### 4.2 From Info Cards to InfoBins

Curated official information typically appears in a basic app as Info Cards. For an agentic system, those cards must become structured, versioned, and locally usable objects, which we call *InfoBins*. An InfoBin bundles approved content with metadata (domain, language, source authority, zone, validity period, last-updated time), a semantic representation that allows natural-language matching, and a safety layer covering source labels, freshness, revocation, and escalation rules. The runtime's first question for any

query is which approved InfoBin already contains the answer; only when no local match is found does the system escalate.

### 4.3 The semantic cache as a unified retrieval primitive

#### One Retrieval Primitive, Five Functions

Each function is cosine similarity against a different partition of the same embedding-indexed store.

Conventional design: five separate systems

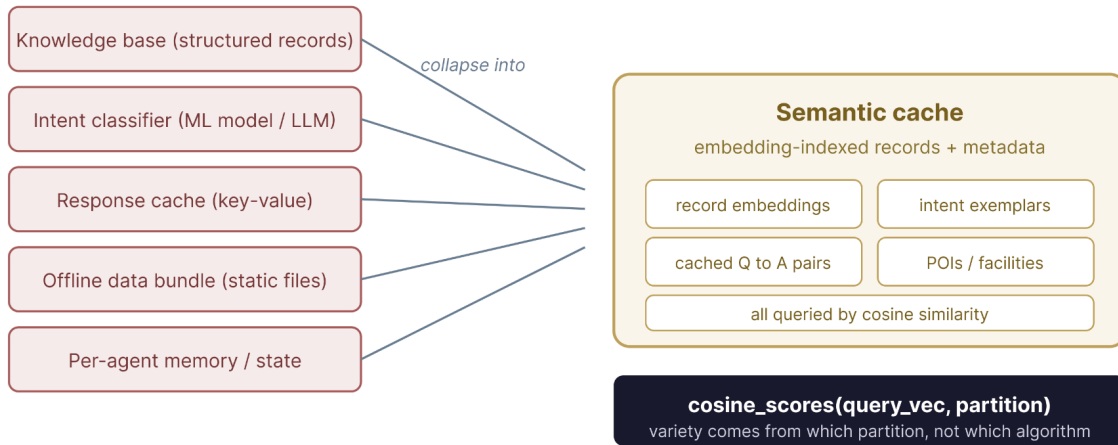


Figure 4. Five traditionally separate systems expressed as one retrieval primitive over partitions of a shared store.

The central technical idea is that a single embedding-indexed store, which we term the semantic cache, can serve five functions that conventional systems implement separately: a knowledge base of curated records, an intent classifier, a paraphrase-aware response cache, an offline data bundle, and per-agent memory. Each function is expressed as cosine similarity against a different partition of the same store. The architectural variety comes from which partition is compared against, not from different algorithms.

We are deliberate about what this claim does and does not assert. The semantic cache is an embedding-indexed retrieval layer, not a single physical database, and the system does still use conventional storage backends. In practice the cache is realized differently on each tier. On the server, domain partitions are held in memory and cache-backed stores, including a Redis store for paraphrase-aware query-to-answer caching, using higher-dimensional embeddings that are updated with live content. On the device, a curated offline partition ships as a SQLite bundle of roughly 700 KB that combines a full-text (FTS5) index with precomputed lower-dimensional embeddings stored as binary blobs, enabling hybrid retrieval with no network call. Both tiers expose the same abstraction, approved records queried by cosine similarity and hybrid lexical-semantic search, even though they use different embedding dimensions and backends. We note as a current limitation that the two tiers do not share a single physical store and that live synchronization from server indexes to the device bundle is not yet implemented; the offline bundle is materialized at build time.

### 4.4 Intent routing as nearest-neighbor classification

Most agentic systems classify intent by calling an LLM, which is expensive, slow, and opaque. KumbhDoot reframes intent classification as a similarity problem: whether a query is close to questions that humans

have labelled with a known intent. We curate approximately 520 labelled exemplar queries across 26 intent categories in English, Hindi, and Marathi, embed them once at startup, and at query time embed the user query and compare it against all exemplars in a single matrix multiplication. The nearest exemplar yields intent and domain labels together with a traceable numerical similarity score. This is k-nearest-neighbor classification over multilingual sentence embeddings. It consumes zero generative LLM tokens, though the embedding model itself still runs, and it produces an auditable score rather than an opaque model decision. Query complexity is determined not by the exemplar match but by a separate heuristic component that counts domain-keyword hits, conjunction words across languages, and query length. The router precedes this with an emergency check that detects urgent phrases by keyword and routes them onto a rule-based fast path, and with language normalization for transliteration and common speech-recognition errors.

### 4.5 Offline and low-bandwidth resilience

Because dense crowds make connectivity unreliable, the device tier is designed to remain useful offline. Emergency contacts, key facilities, saved itinerary, basic maps, FAQs, phrasebooks, and previously synced InfoBins are available without a network. When connectivity is absent the system degrades gracefully: voice falls back to text, agentic answers fall back to approved static guidance, and live data is labelled with its last-updated status. Non-critical requests are queued and resumed on reconnect, while emergency updates receive sync priority. Answers are explicitly labelled as available offline, last updated, or pending refresh, so users can judge freshness. On the device, retrieval is hybrid: a lexical FTS5 search is fused with semantic similarity over the embedding blobs using rank fusion. When a quantized on-device embedding model is present, the client embeds live queries locally; when those assets are absent, it falls back to word-overlap matching against pre-baked query embeddings, preserving basic offline function without live embedding.

### 4.6 The three-tier agent model

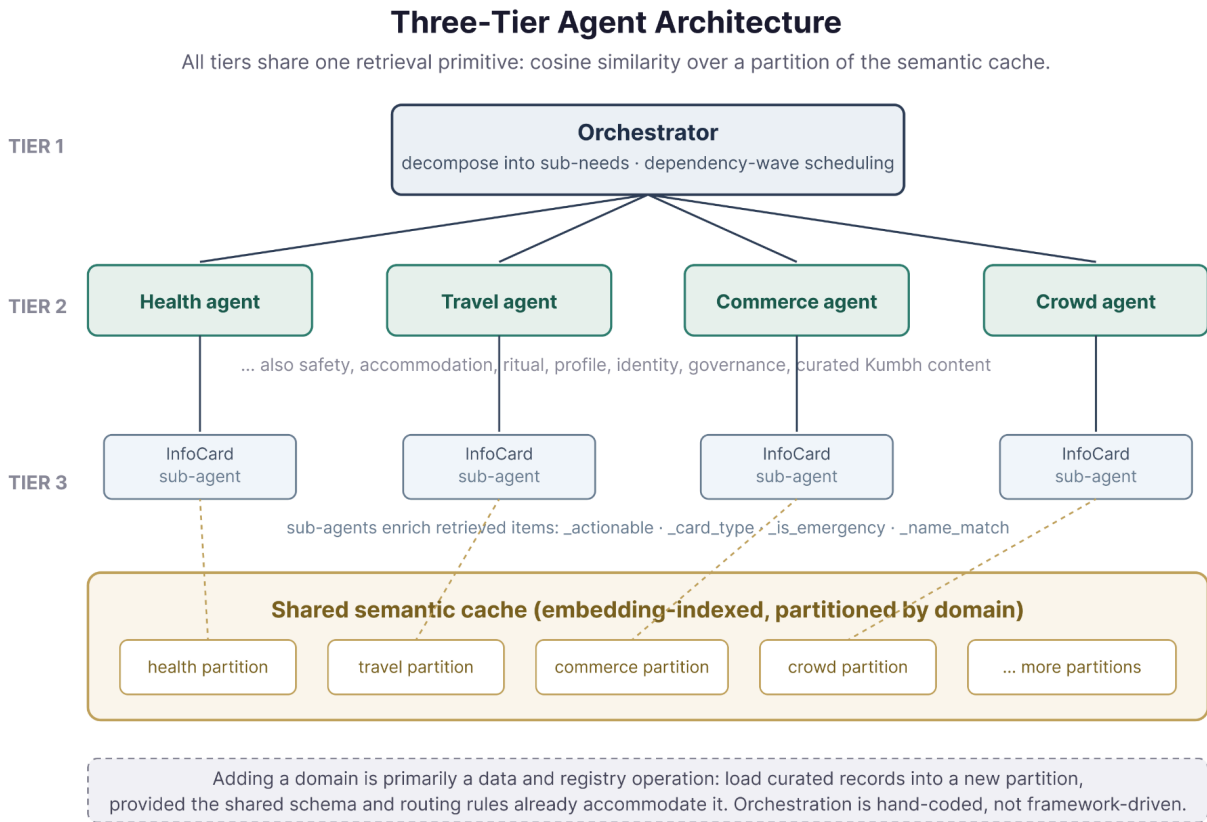


Figure 3. Tiers share a retrieval primitive, not one physical store; the LLM is reached only under explicit routing conditions.

For queries that require composition across domains, KumbhDoot uses a custom three-tier agent architecture in which all tiers share the same semantic retrieval layer. An *orchestrator* decomposes a complex request into sub-needs with dependency edges and schedules them into parallel execution waves, so that independent sub-needs run concurrently while dependent ones wait for their prerequisites and receive prior results as context. *Domain agents*, one per domain such as health, travel, commerce, crowd, identity, governance, and curated Kumbh content, each operate over a partition of the cache. *Record-level sub-agents* enrich retrieved items into structured answers with actionability signals, card type, emergency flags, and name matching.

Two design choices distinguish this architecture. First, the orchestration logic is hand-written and domain-specific rather than built on a general-purpose multi-agent framework. Because the agents are wired to the semantic cache first, an LLM is reached only under explicit routing conditions rather than as a default at every orchestration step, which keeps each agent's decision path inspectable. Second, adding a domain is primarily a data and registry operation, loading curated records into a new partition and extending a data-driven registry, provided the shared schema and routing rules already accommodate that domain; it is not, in the common case, a matter of writing new orchestration code. We state both claims carefully: the tiers share a retrieval *primitive*, not one physical store, and new domains are *primarily*, not entirely, a data operation.

## 4.7 The escalation ladder

The runtime tries methods in cost order and reaches the LLM only when the cheaper layers cannot answer. A query first hits the paraphrase-aware query-to-answer cache; a confident single-record hit is returned by a coded formatter with no LLM call. Failing that, cosine-based intent routing selects a domain at zero generative cost. For confident single-domain retrieval, multiple records are formatted into a grounded, cited answer, still without an LLM. The LLM is invoked only when an explicit escalation gate fires: a complex multi-part synthesis, retrieval confidence below threshold, or a need for live or external-source data. Even then, generation is constrained to grounded, retrieved records, and a grounding-validation step strips any citation that does not correspond to a record the retriever actually returned, removing invented identifiers before the answer is shown.

# The Escalation Ladder

Methods are tried in cost order. The LLM is reached only when an explicit gate fires.

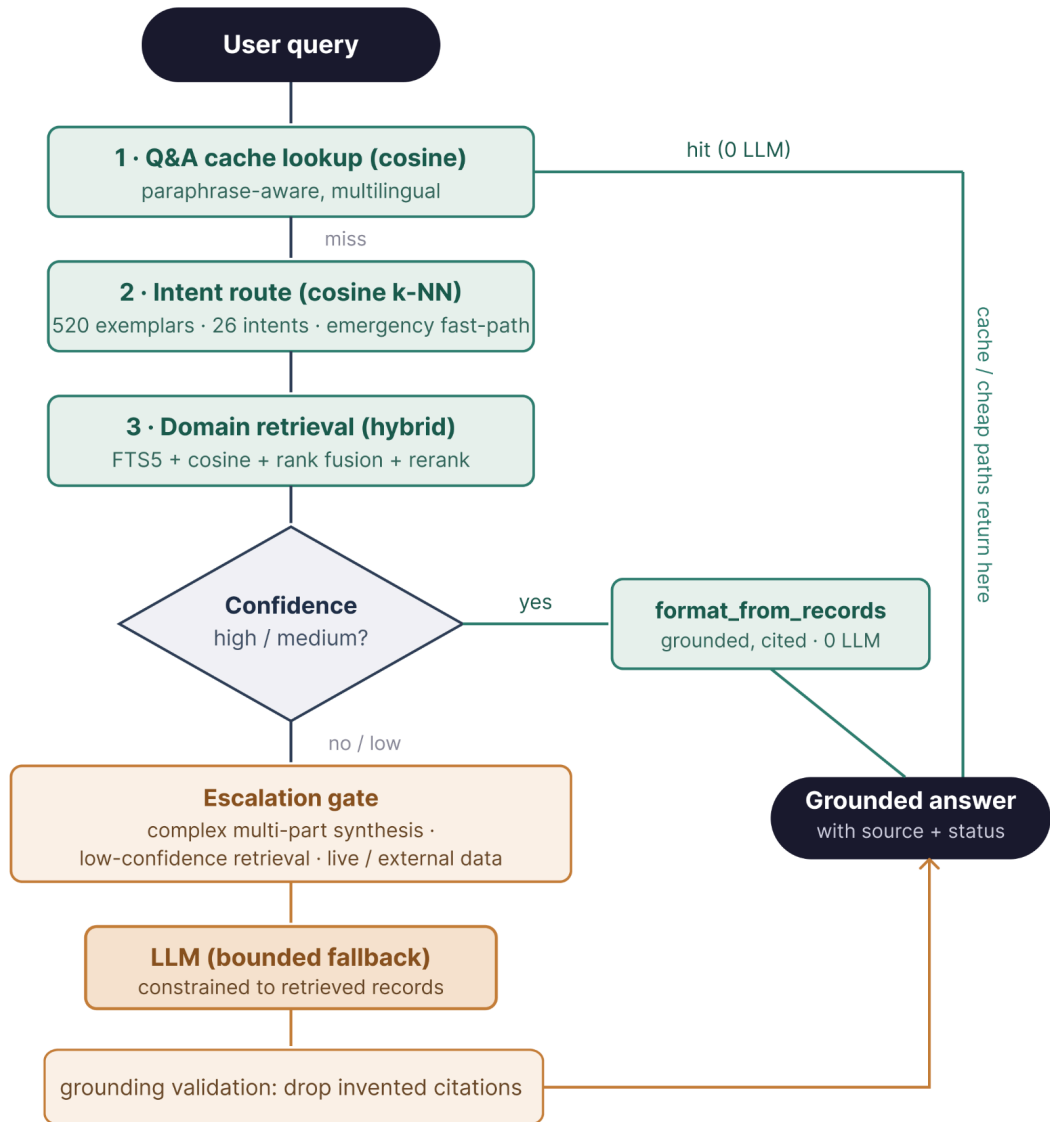


Figure 2. Green paths return without a generative model. The amber path is the only one that invokes the LLM.

## 5. Cost Model: Per-Query Economics

The similarity-first design produces a distinctive cost structure that can be reasoned about analytically, independent of any particular benchmark. In a pipeline where every processing step is an LLM call, cost accrues at each of intent classification, complexity assessment, decomposition, per-domain reasoning, and synthesis. In KumbhDoot, deterministic components handle each of these steps, and the model is invoked only as a bounded fallback.

It is useful to distinguish the cost of the operations themselves. Cosine similarity over a fixed exemplar set is a single matrix multiplication; paraphrase matching in the response cache is a dot product; hybrid retrieval combines vector math with a lexical index and a small local reranker. None of these consume generative tokens. The embedding step does consume compute, so we are careful to say that common paths cost zero *generative* LLM tokens rather than zero compute. On the device tier, where vectors are pre-baked and shipped in the bundle, even the embedding step can be avoided for cached content.

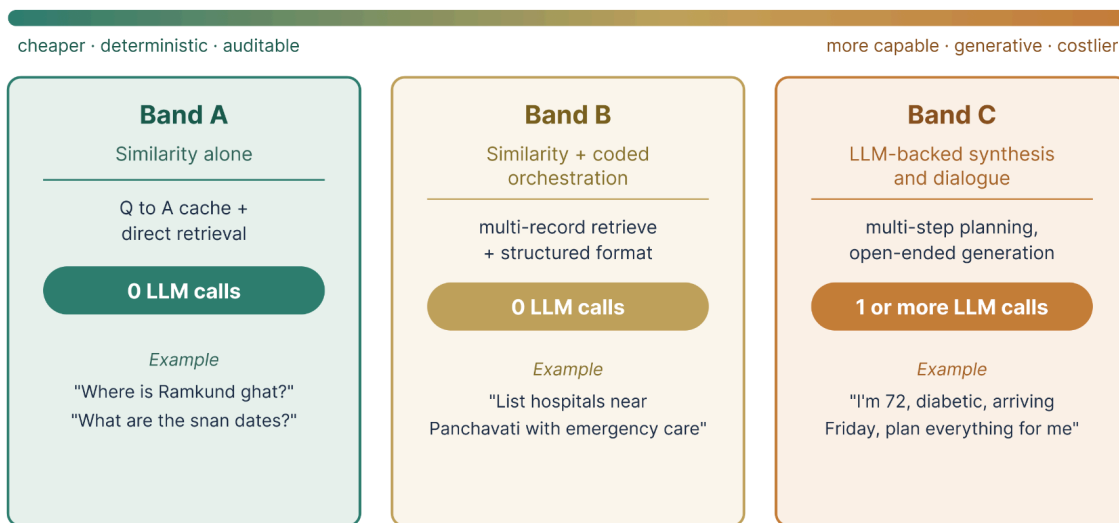
This structure yields a per-query cost profile that scales favorably with volume. The first time any question is asked, the cost is likely to be one embedding and a retrieval. Every subsequent semantically equivalent question, in any of the supported languages and in any phrasing, is likely to be served from the cache. Because pilgrimage query distributions are heavily concentrated on a relatively small set of recurring needs, the share of queries answerable without a model call grows as the cache warms, and the marginal cost of the common case approaches the cost of similarity computation alone. The LLM cost is incurred only on the long tail of genuinely novel, multi-step, or externally-sourced queries.

We present this as an analytical argument about cost structure rather than as a measured result. The empirical validation of the resulting cost and latency figures, and of the fraction of production traffic that the cheaper paths actually absorb, is the subject of the planned evaluation described in Section 7.

## 6. Where Similarity Suffices, and Where It Does Not

### Where Similarity Suffices, and Where It Does Not

Answer each query in the cheapest band that can answer it correctly. Escalation to Band C is deliberate and gated.



Hard limits on A and B: nothing-in-nothing-out, retrieval is not reasoning, no novel synthesis, actions need tools.

Figure 6. The three-band model. Hallucination is mitigated, not eliminated, since Band C still invokes a model.

A central claim of this paper is that the boundary between similarity and generation should be stated explicitly rather than blurred. In KumbhDoot's bounded and curated domain, semantic similarity combined with coded formatters is sufficient for a substantial class of needs: intent classification over a fixed label set, paraphrase-tolerant retrieval, FAQ-style answers, and structured factual responses drawn from approved records. For these, a generative model adds latency, cost, and hallucination risk without adding correctness.

Similarity is not, however, a substitute for reasoning or open-ended generation, and we identify several hard limits. Queries about records absent from the corpus yield low-confidence retrieval and must deflect or escalate rather than fabricate: nothing in, nothing out. Retrieval finds the closest text, which does not guarantee correct intent at decision boundaries. Multi-constraint planning across domains requires orchestration and often genuine prose generation. Coded formatters produce accurate structured output but not fluent narrative, which requires either stored text or a model. Transactional actions such as booking or identity verification need tool-using agents, not cosine similarity. And multilingual coverage, while it helps, does not eliminate degradation from code-mixing, dialect variation, and speech-recognition error.

These limits map onto a three-band view of the workload. In Band A, similarity alone answers the query through the response cache or direct retrieval, with no model call: for example, asking where a ghat is or what the bathing dates are. In Band B, similarity plus coded orchestration handles multi-record, structured requests, again without a model: for example, listing hospitals with emergency care near a locality. In Band C, the system escalates to LLM-backed synthesis and dialogue: for example, an open-ended personalized plan for an elderly traveller arriving on a given day. The design goal is not to push everything into Band A but to ensure that queries are answered in the cheapest band that can answer them correctly, and that escalation to Band C is a deliberate, gated decision rather than a default.

On hallucination we are correspondingly careful. By driving agents from curated, embedding-indexed records and using coded formatters for the majority of responses, KumbhDoot reduces opportunities for unconstrained generation and thereby *mitigates* hallucination relative to architectures that invoke a model at every step. It does not *eliminate* hallucination: the system still calls a model for some queries, and retrieval-based systems can themselves surface incorrect records. The grounding-validation step that removes unsupported citations is a mitigation, not a guarantee.

---

## 7. Limitations and Planned Evaluation

This paper presents an architecture and an analytical cost argument; it does not yet present a controlled empirical evaluation, and we are explicit about what remains to be measured.

**Evaluation of the core cost and latency claims.** The central quantitative claims, the fraction of production queries answerable without a generative model call, the realized per-query cost and latency by band, and the speedup relative to an LLM-default pipeline, require a controlled benchmark over a representative query distribution rather than a handful of illustrative queries. We plan a tiered benchmark spanning simple, multi-record, and complex multi-domain queries, each paired with semantically equivalent paraphrases to test cache behavior, with full per-step tracing of tokens, latency, and cost. The baseline must be characterized honestly: an "LLM-for-everything" pipeline that forces a model through every step is a worst-case rather than a representative comparison, and a fair evaluation should also include a conventional RAG baseline that uses a model only for synthesis.

**Emergency systems need to have more escape hatches.** A keyword fast path is useful, but it cannot be the principal emergency safety mechanism. There is a chance of missing indirect statements, transcription errors and code-mixed expressions, and it may create false positives. The production design should additionally include:

- a permanent one-tap emergency control,
- locally stored emergency numbers and instructions,
- explicit confirmation of location and emergency type,
- direct telephone/SMS paths,
- human escalation,

- and a safe fallback when classification confidence is uncertain.

Cell broadcast is especially relevant for official one-to-many emergency alerts because it is designed to operate with limited impact from network congestion.

**Routing accuracy across languages.** The accuracy of cosine-based intent routing across English, Hindi, Marathi, and code-mixed input is asserted but not yet measured. Embedding quality degrades under dialect variation and speech-recognition error, and the confidence thresholds that govern answer, clarify, deflect, and escalate decisions need to be tuned and reported against a labelled multilingual test set.

**Offline coverage.** The minimum InfoBin bundle that yields useful offline behavior, and the cache hit rate achievable on-device, are open empirical questions. The favorable scaling argument in Section 5 depends on a query distribution concentrated enough for the cache to warm quickly, which should be validated on real traffic rather than assumed.

**Device-server synchronization.** As noted in Section 4.3, live synchronization from server indexes to the device bundle is not implemented; the offline partition is materialized at build time. Closing this gap, including signed updates and revocation for stale or corrected guidance, is necessary before the offline tier can be trusted for fast-changing advisories.

**Freshness and safety operations.** Cached guidance can become stale during road closures, crowd changes, or emergencies. The system design includes last-updated labels, tiered refresh cadences, and revocation lists, but the operational discipline of keeping safety-critical InfoBins current, and the audit standards appropriate to a public-sector agentic application, require field validation.

**Privacy at scale.** The intended privacy posture, bucketed rather than exact attributes, local-first storage, coarse and consented location, aggregate-by-default dashboards, and audit logs, is a design commitment that needs to be verified in deployment, including the question of how much personalization is useful before its privacy cost outweighs its benefit.

The next phase of work converts these limitations into a research agenda: a controlled benchmark with honest baselines, a labelled multilingual routing evaluation, measurement of offline cache coverage and hit rate on representative traffic, implementation and evaluation of device-server sync, and a privacy and audit assessment suitable for public-sector deployment.

---

## 8. Conclusion

KumbhDoot addresses the specific challenges of mass-gathering environments, where large-scale, multilingual, safety-critical public services must operate despite limited connectivity and high cost sensitivity. By rejecting the standard LLM-default architectural pattern, the system prioritizes reliability and efficiency through an inverted operational order: semantic similarity over a curated knowledge base serves as the foundational principle, while large language models are restricted to a bounded, secondary role.

- **Environmental Constraints:** The system is engineered to handle mass-gathering scenarios characterized by millions of concurrent users, intermittent network connectivity, code-mixed and multilingual input, and safety-critical information needs where errors carry physical consequences.
- **Architectural Inversion:** Rather than generating every response via an LLM, the system centers on semantic similarity. It utilizes a single retrieval primitive—cosine similarity over an embedding-indexed semantic cache—to perform five key functions: intent routing, paraphrase caching, offline data lookup, and multi-agent retrieval.

- **Custom Agent Framework:** The implementation employs a hand-coded, three-tier agent architecture that operates directly on the semantic store. This bypasses the opacity of general-purpose multi-agent frameworks, ensuring that decision paths remain fully inspectable and that LLM calls are gated rather than automatic.
- **Rigorous Analytical Basis:** The authors provide an analytical cost argument to support the design, while remaining transparent about the current implementation's limitations and the necessity of future empirical work to establish quantitative claims and validate research gaps.
- **Architectural Thesis:** For bounded, high-stakes, and low-connectivity domains, a similarity-first and LLM-bounded approach is both economically superior and architecturally appropriate. The value of generative AI in such contexts is maximized by reserving it exclusively for tasks that cannot be solved through reliable, deterministic retrieval.

## References

1. Salton, G. (1972). The vector space model and TF-IDF for lexical scoring.
2. Manning, Christopher D. *Introduction to information retrieval*. Synpress Publishing,, 2008.
3. Deerwester, Scott, Susan T. Dumais, George W. Furnas, Thomas K. Landauer, and Richard Harshman. "Indexing by latent semantic analysis." *Journal of the American society for information science* 41, no. 6 (1990): 391-407.
4. Mikolov, Tomas, Kai Chen, Greg Corrado, and Jeffrey Dean. "Efficient estimation of word representations in vector space." *arXiv preprint arXiv:1301.3781* (2013).
5. Reimers, Nils, and Iryna Gurevych. "Sentence-bert: Sentence embeddings using siamese bert-networks." In *Proceedings of the 2019 conference on empirical methods in natural language processing and the 9th international joint conference on natural language processing (EMNLP-IJCNLP)*, pp. 3982-3992. 2019.
6. Cheng, Zhiyuan, Longying Lai, Yue Liu, Kai Cheng, and Xiaoxi Qi. "Enhancing financial report question-answering: A retrieval-augmented generation system with reranking analysis." *arXiv preprint arXiv:2603.16877* (2026).
7. Lewis, Patrick, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler et al. "Retrieval-augmented generation for knowledge-intensive nlp tasks." *Advances in neural information processing systems* 33 (2020): 9459-9474.
8. Yao, Shunyu, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. "React: Synergizing reasoning and acting in language models." *arXiv preprint arXiv:2210.03629* (2022).
9. Hinton, Geoffrey, Oriol Vinyals, and Jeff Dean. "Distilling the knowledge in a neural network." *arXiv preprint arXiv:1503.02531* (2015).
10. McMahan, H. Brendan, Eider Moore, Daniel Ramage, Seth Hampson, and B. Agüeray Arcas. "Communication-efficient learning of deep networks from decentralized data (2016)." *arXiv preprint arXiv:1602.05629* (2016).
11. Han, Song, Huizi Mao, and William J. Dally. "Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding." *arXiv preprint arXiv:1510.00149* (2015).
12. Brown, Tom, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D. Kaplan, Prafulla Dhariwal, Arvind Neelakantan et al. "Language models are few-shot learners." *Advances in neural information processing systems* 33 (2020): 1877-1901.
13. Guu, Kelvin, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. "Retrieval augmented language model pre-training." In *International conference on machine learning*, pp. 3929-3938. PMLR, 2020.
14. Karpukhin, Vladimir, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. "Dense passage retrieval for open-domain question answering." In *Proceedings of the 2020 conference on empirical methods in natural language processing (EMNLP)*, pp. 6769-6781. 2020.
15. Shuster, Kurt, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. "Retrieval augmentation reduces hallucination in conversation." In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pp. 3784-3803. 2021.
16. Gao, Tianyu, Xingcheng Yao, and Danqi Chen. "Simcse: Simple contrastive learning of sentence embeddings." In *Proceedings of the 2021 conference on empirical methods in natural language processing*, pp. 6894-6910. 2021.