

Lilo Engine

Production-Grade Therapeutic AI Platform

Multi-Agent AI for Elderly Mental Health in Senior Living Facilities

[65K+ Lines of Code]

[15 Microservices]

[100% Crisis Recall]

[<1s Response]

Platform At-A-Glance — Key Metrics

CODEBASE

65,000+

lines

SERVICES

15

microservices

CRISIS RECALL

100%

zero false neg

RESPONSE TIME

<1 second

30x regulatory

INTENT ACCURACY

92-95%

214 prototypes

THERAPEUTIC AGENTS

7

evidence-based

CACHE HIT RATE

60-80%

<5ms on hit

GENERATION LATENCY

~7.6s

streaming

[Python 3.12]

[Go 1.25]

[PyTorch 2.8]

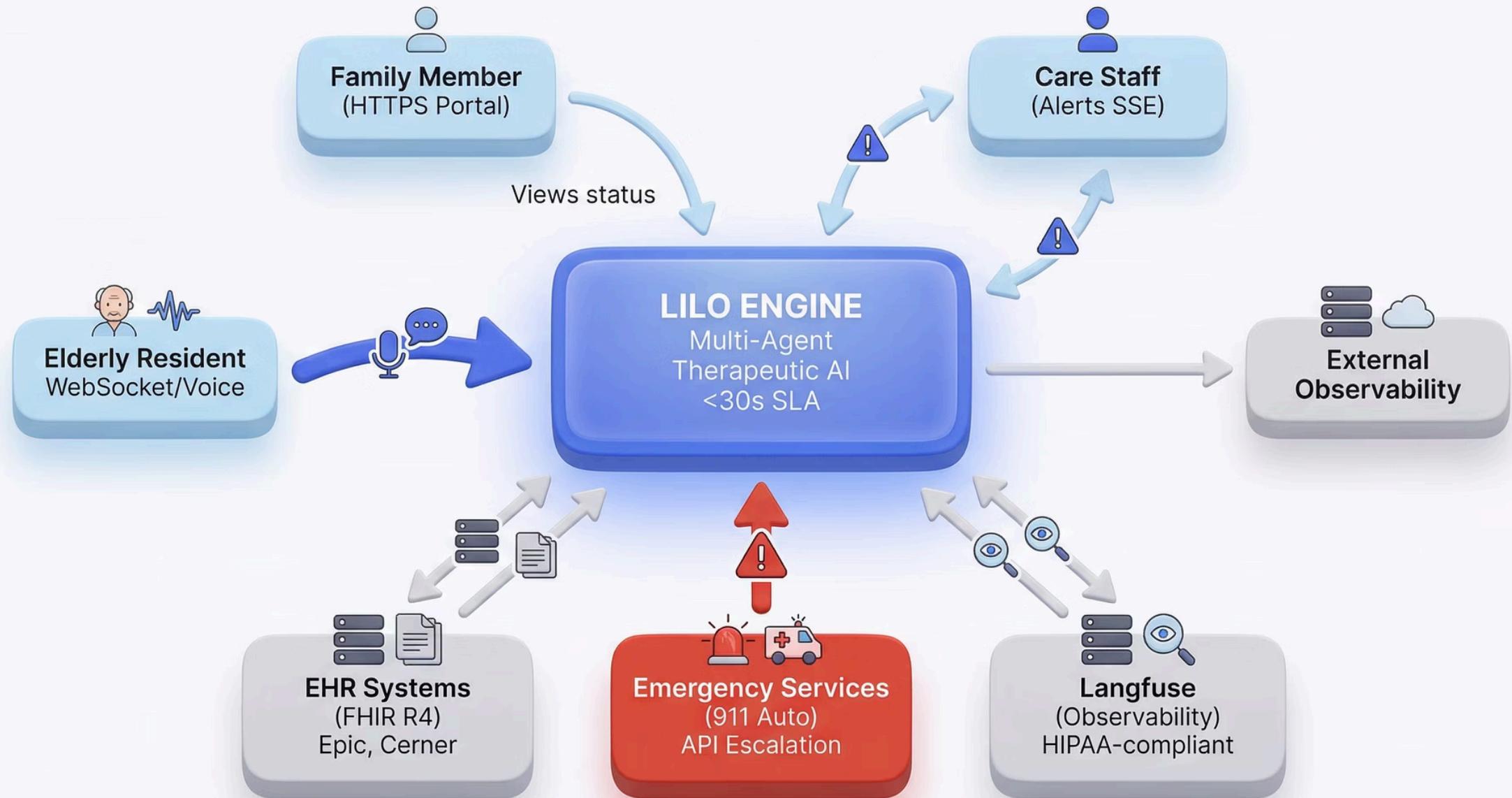
[PostgreSQL 16 + pgvector]

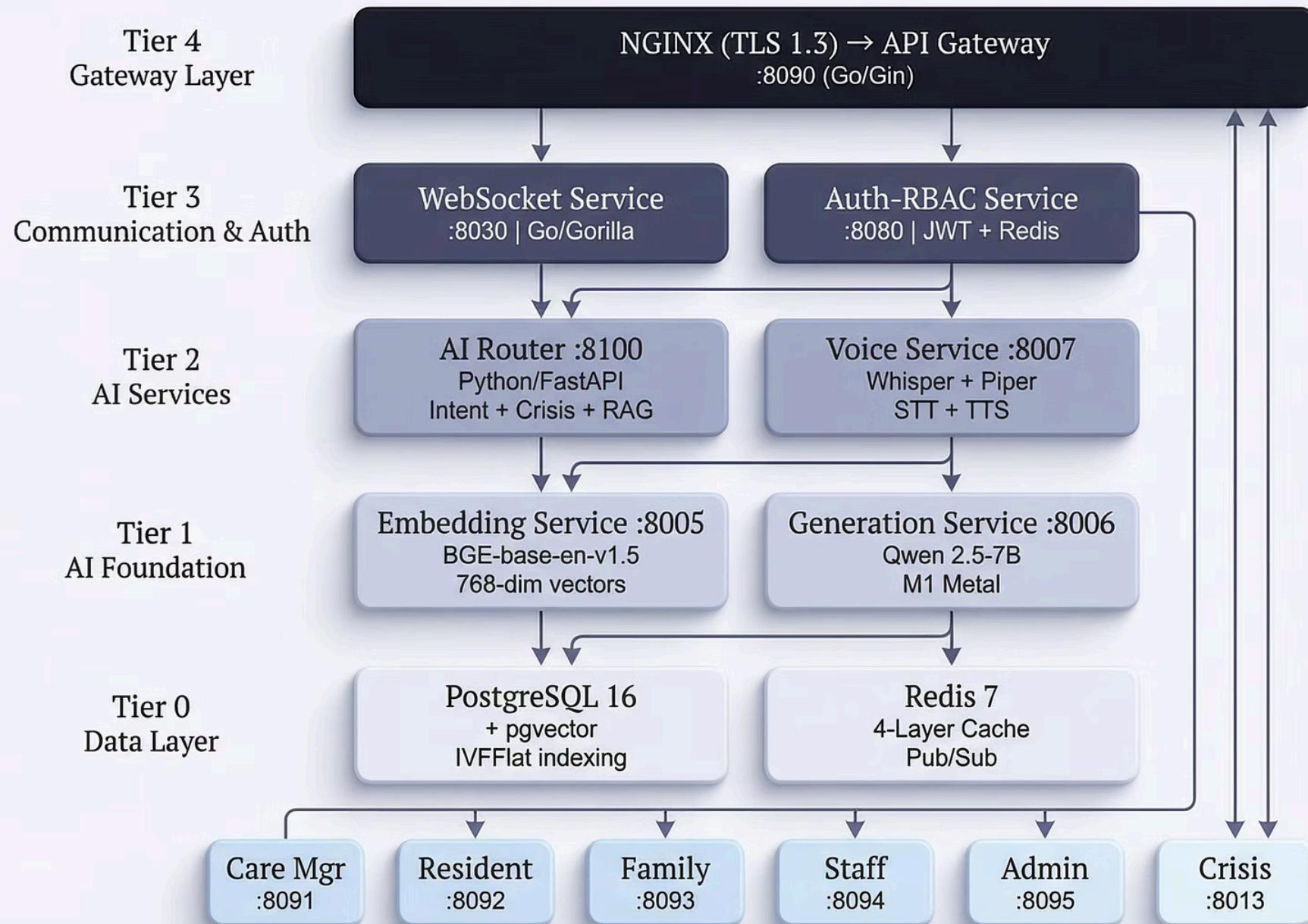
Redis 7]

Docker

Qwen 2.5-7B]

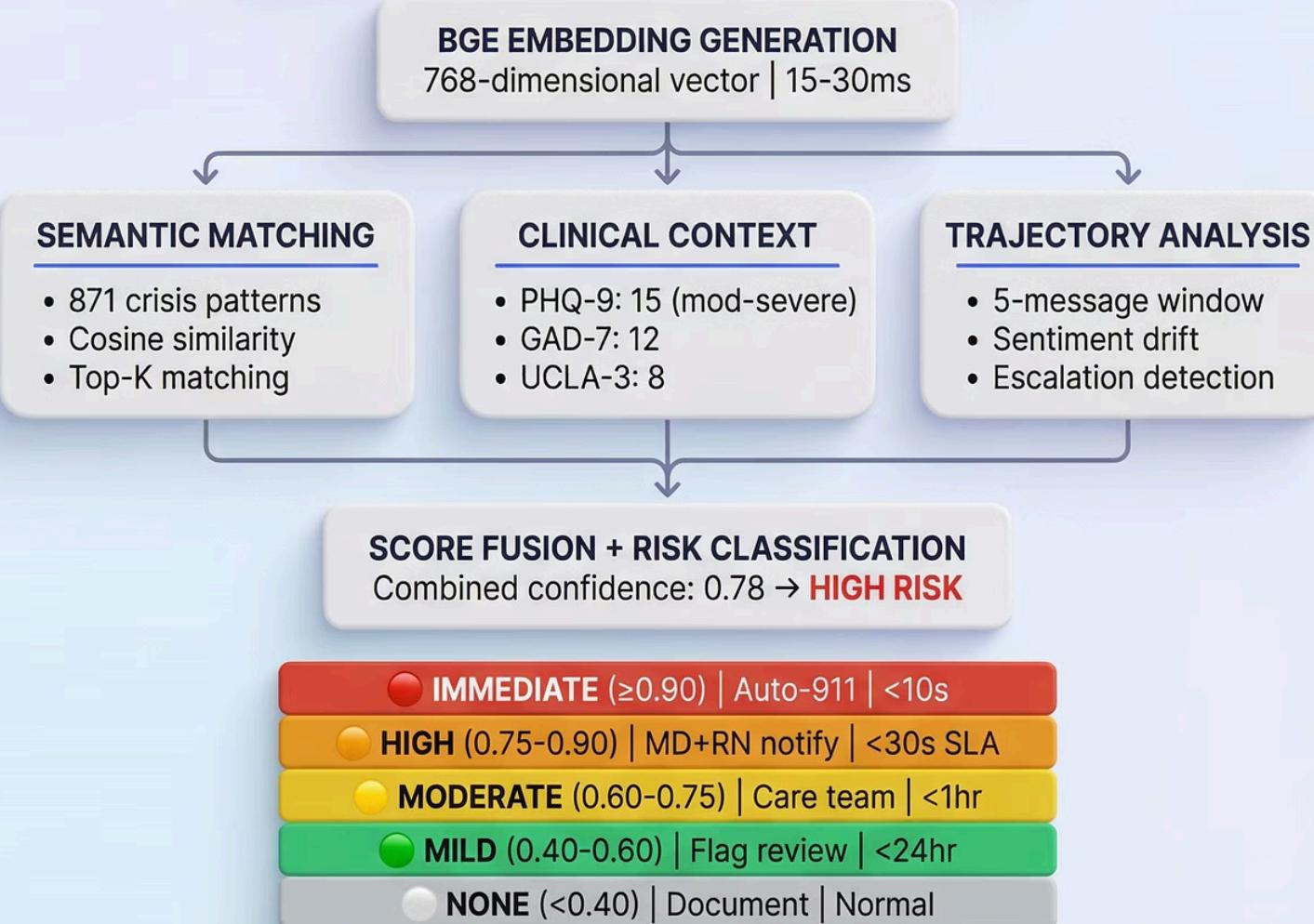
BGE Embeddings







INPUT: Sometimes I wonder if it's worth going on



Lilo Engine	Industry Standard
Recall: 100%	60-70%
Precision: 95%+	70-80%
Response Time: <1s	15-30min
False Positives: <5%	20-30%

INTENT CLASSIFICATION

Primary: 0.85
Secondary: 0.72

CRISIS CHECK

(Safety Override)

SAFETY ASSESSMENT

(C-SSRS Protocol)

Usage: 8%

CLINICAL_PRIORITY

NORMAL AGENT SELECTION

CONVERSATIONAL

(CBT foundations)

Usage: 45%

WEB SEARCH

(Real-time info)

Usage: 1%



REMINISCENCE

(Life review therapy)

Usage: 22%



-15% depression

BRIDGE

(Social connection)

Usage: 2%



-2 UCLA-3



GROUNDING

(40-60% anxiety reduction)
Usage: 10%



BEHAVIORAL ACTIVATION

(35% PHQ-9 reduction)

Usage: 12%

CLINICAL_PRIORITY:

Safety agent exclusively

SEQUENTIAL:

Primary → 2 secondary max

PARALLEL:

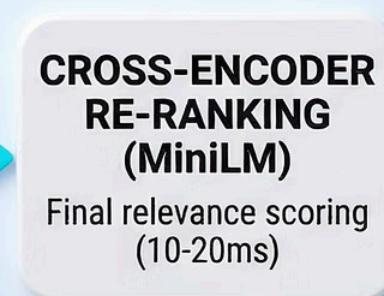
Multiple agents simultaneous

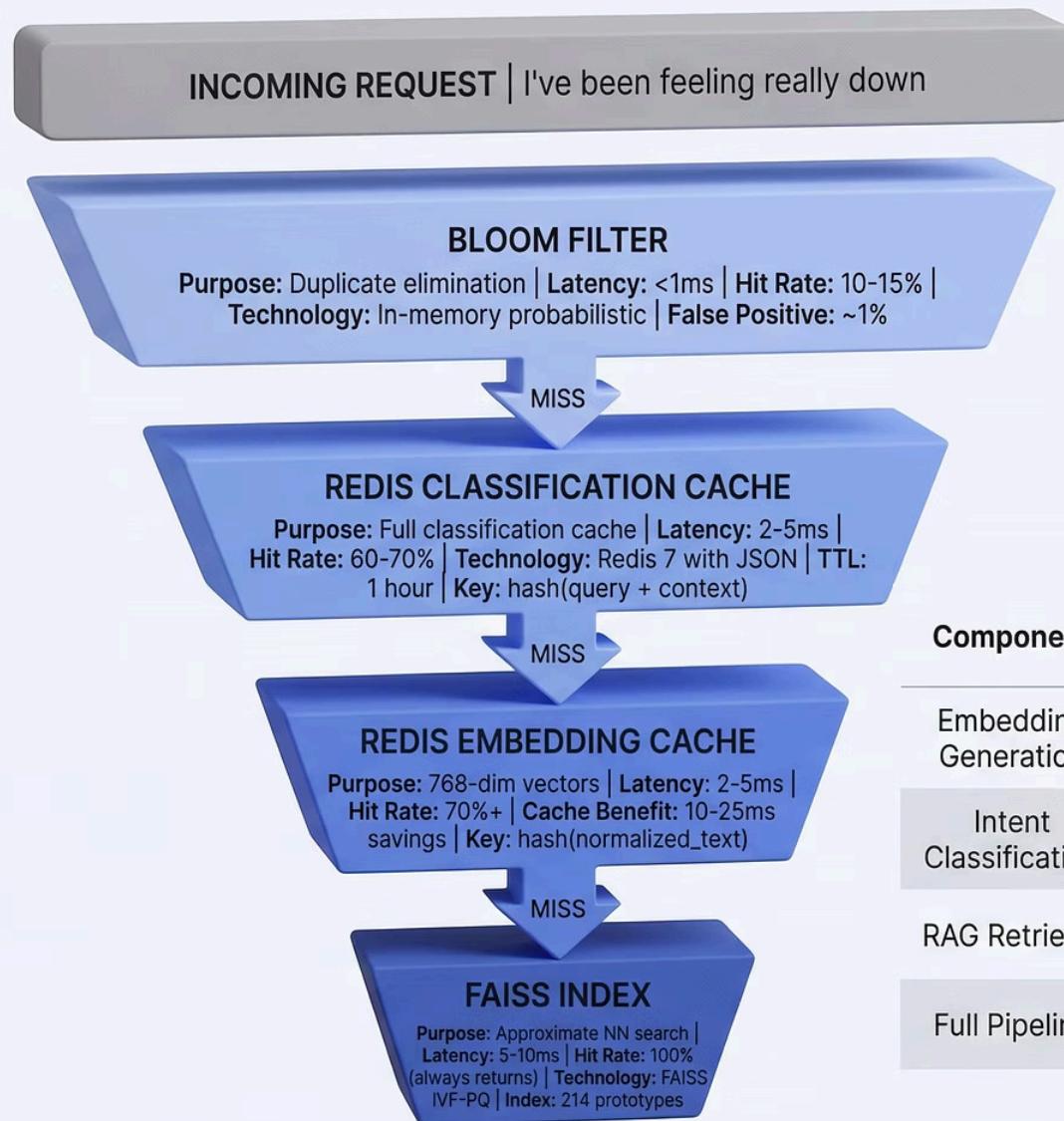
QUERY: I miss going to church with my husband



BGE EMBEDDING: 768-dimensional vector (15-30ms)

asyncio.gather() — PARALLEL EXECUTION | Total: ~45ms





Component	Without Cache	With Cache	Savings
Embedding Generation	15-30ms	2-5ms	10-25ms
Intent Classification	40-50ms	10-15ms	30-40ms
RAG Retrieval	60-70ms	30-45ms	15-25ms
Full Pipeline	150-200ms	80-120ms	~50%

Technical Differentiators — Why Lilo Engine Wins

DIFFERENTIATOR	LILO ENGINE	INDUSTRY STANDARD
Crisis Detection Recall	100% Zero false negatives	60-70% Missed crises common
Response Time	<1s 30x regulatory	15-30 minutes Manual review
False Positive Rate	<5% Low alert fatigue	20-30% Staff overwhelmed
HIPAA Compliance	Built-in Day 1	Often retrofitted
Edge Deployment	8GB 90/10 split	N/A Cloud-only
Therapeutic Agents	7 Evidence-based	1 generic chatbot
Clinical Integration	PHQ-9 GAD-7 C-SSRS	Limited
RAG Personalization	5 parallel streams	Generic only

