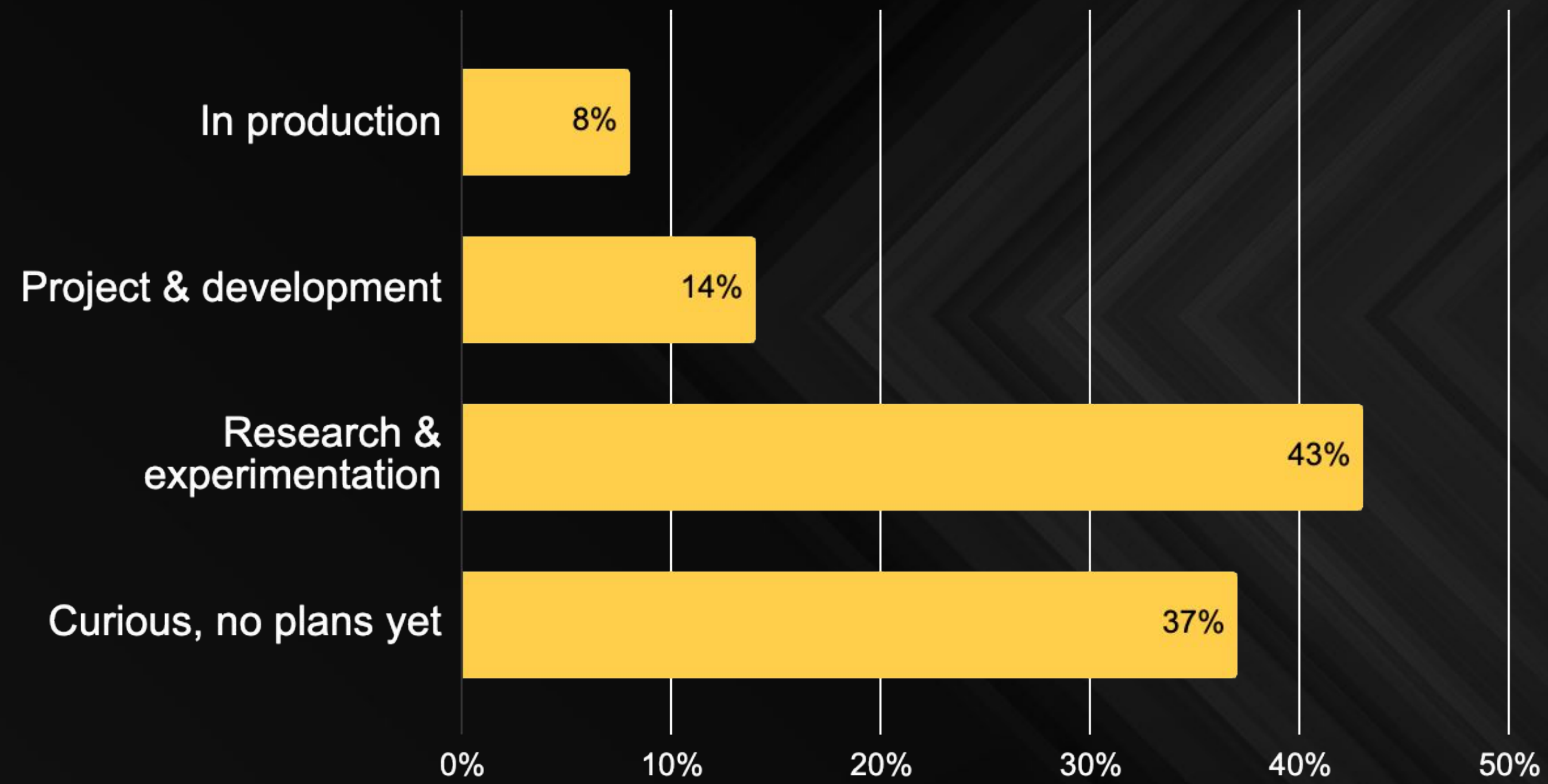AKKA

Webinar

# A Blueprint for Agentic AI Services

Best practices for designing and operating
agentic-scale services

Accelerate delivery.  Stay safe and be efficient.

# Poll question #1

## Where are you in your agentic journey?



n = 157 webinar poll respondents

In production — 8%
Project & development — 14%
Research & experimentation — 43%
Curious, no plans yet — 37%

# Today's discussion

**AKKA**

**01** Welcome
Darin Bartik, CMO, Akka

**02** The agentic opportunity and the move to a-tier architecture
Richard Li, AI Expert and Entrepreneur

**03** A blueprint for agentic services
Tyler Jewell, CEO, Akka

**04** Agentic stories and AI in practice
Real-time video augmentation, model-driven personalization, Google Earth AI inference

**05** Live Q&A
Q&A transcript and slides to be shared asap

# AI is **transforming** our lives

AKKA

## AI Assistant

A *user* app that understands natural language commands and uses a conversational AI interface to complete tasks on-demand.

ChatGPT

einstein
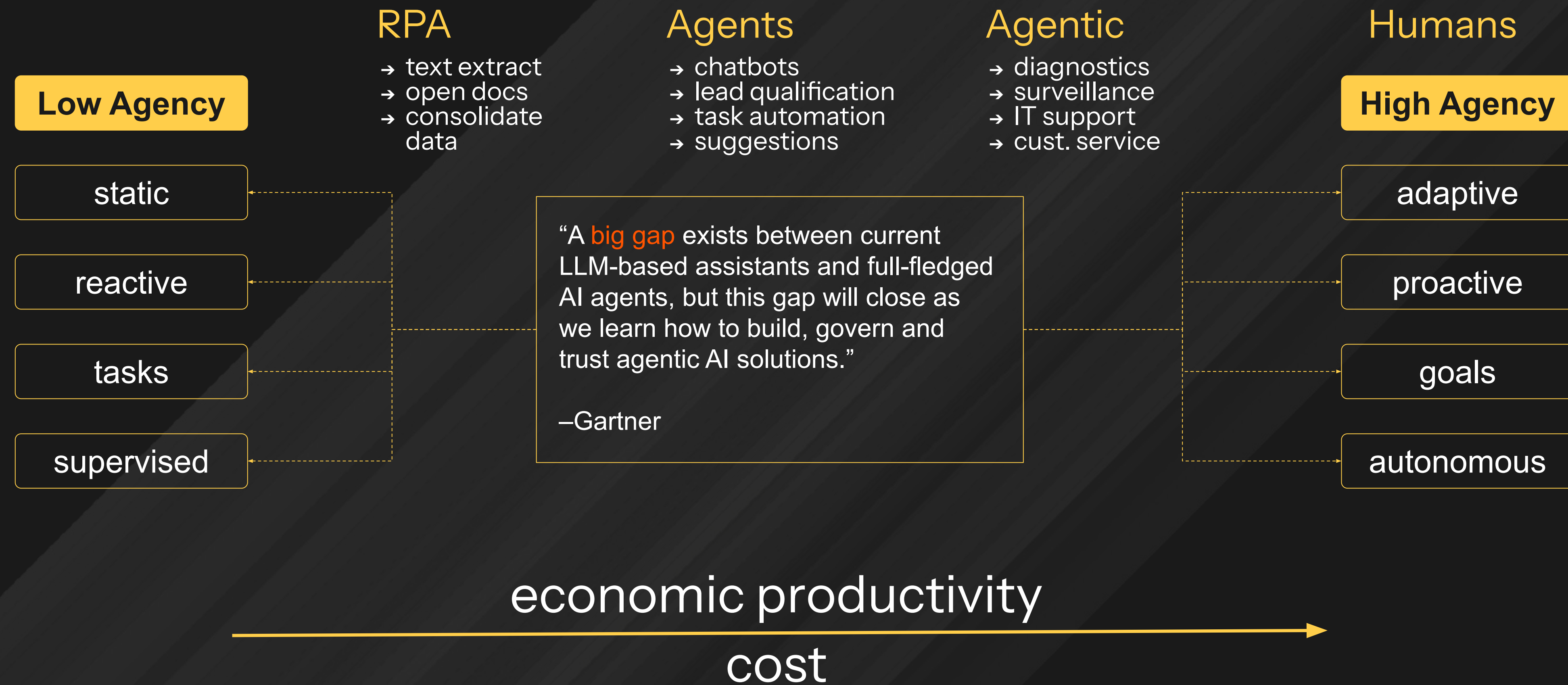
perplexity

Siri

## AI Agent

A *system* that can autonomously fulfill goals by interacting with other systems and agents.
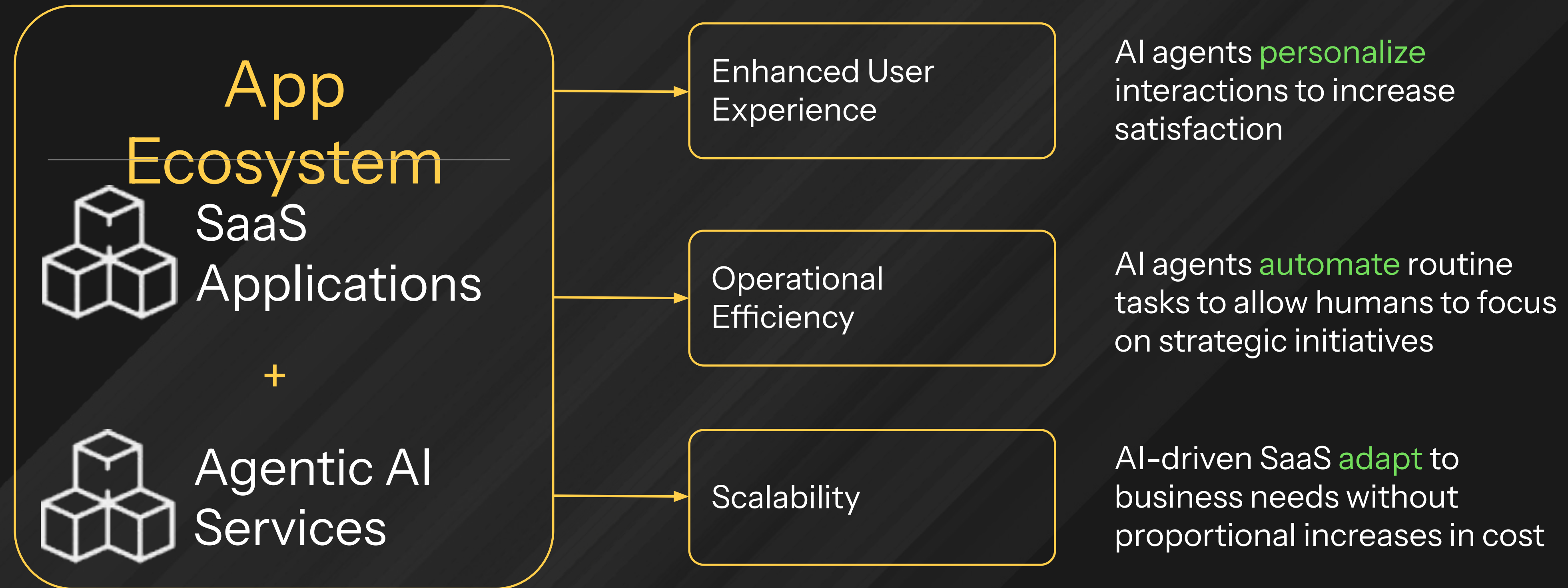
Agentforce

AI at ServiceNow

# AI **Agency**

Capacity to make meaning from your environment

AKKA

| | **RPA** | **Agents** | **Agentic** | **Humans** |
|---|---|---|---|---|

**Low Agency**

**RPA**
→ text extract
→ open docs
→ consolidate data

**Agents**
→ chatbots
→ lead qualification
→ task automation
→ suggestions

**Agentic**
→ diagnostics
→ surveillance
→ IT support
→ cust. service

**High Agency**

static

reactive

tasks

supervised

"A big gap exists between current LLM-based assistants and full-fledged AI agents, but this gap will close as we learn how to build, govern and trust agentic AI solutions."

–Gartner

adaptive

proactive

goals

autonomous

economic productivity
—————————————————→
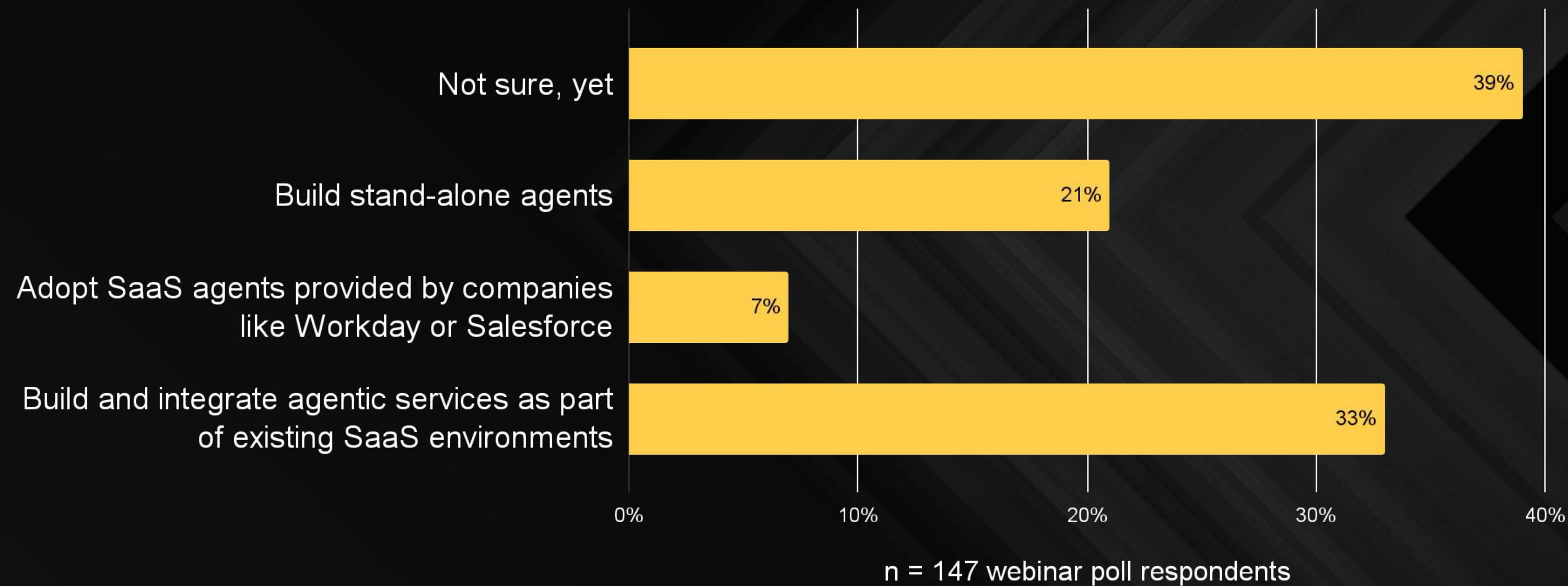cost

# A paradigm shift to AI-fueled app ecosystems

## AI agents and apps become part of a symbiotic existence

By 2028, 33% of enterprise software applications will include agentic AI, up from less than 1% in 2024.
Gartner, *TSP 2025 Trends: Agentic AI — The Evolution of Experience, 24 February 2025*

## App Ecosystem

SaaS Applications

**+**

Agentic AI Services

**Enhanced User Experience**

AI agents personalize interactions to increase satisfaction

**Operational Efficiency**

AI agents automate routine tasks to allow humans to focus on strategic initiatives

**Scalability**

AI-driven SaaS adapt to business needs without proportional increases in cost

# Poll question **#2**

## Which approach(es) is your organization considering about agentic apps?



n = 147 webinar poll respondents

# Agentic is the 5th wave of compute

Every human and device with dozens of sleepless assistants

|  | Mainframe | Web | Cloud | Mobile | Agentic |
|---|---|---|---|---|---|
| Users | thousands | millions | 10 millions | billions | trillions |
| TPS | 100 | 500 | 2,500 | 10,000 | 1,000,000 |
|  |  | 5x | 5x | 4x | 100x |

# Agents are orchestrated services

## Workflows: traceable, auditable, debuggable, with point-in-time recovery

AKKA

### Agents are workflows

reliable execution of AI tasks with visibility into request /
response data, built-in retries, and error compensation

### Task chaining

AI agents break complex workflows into smaller,
composable steps

# Agent types orchestrate **levels of agency**
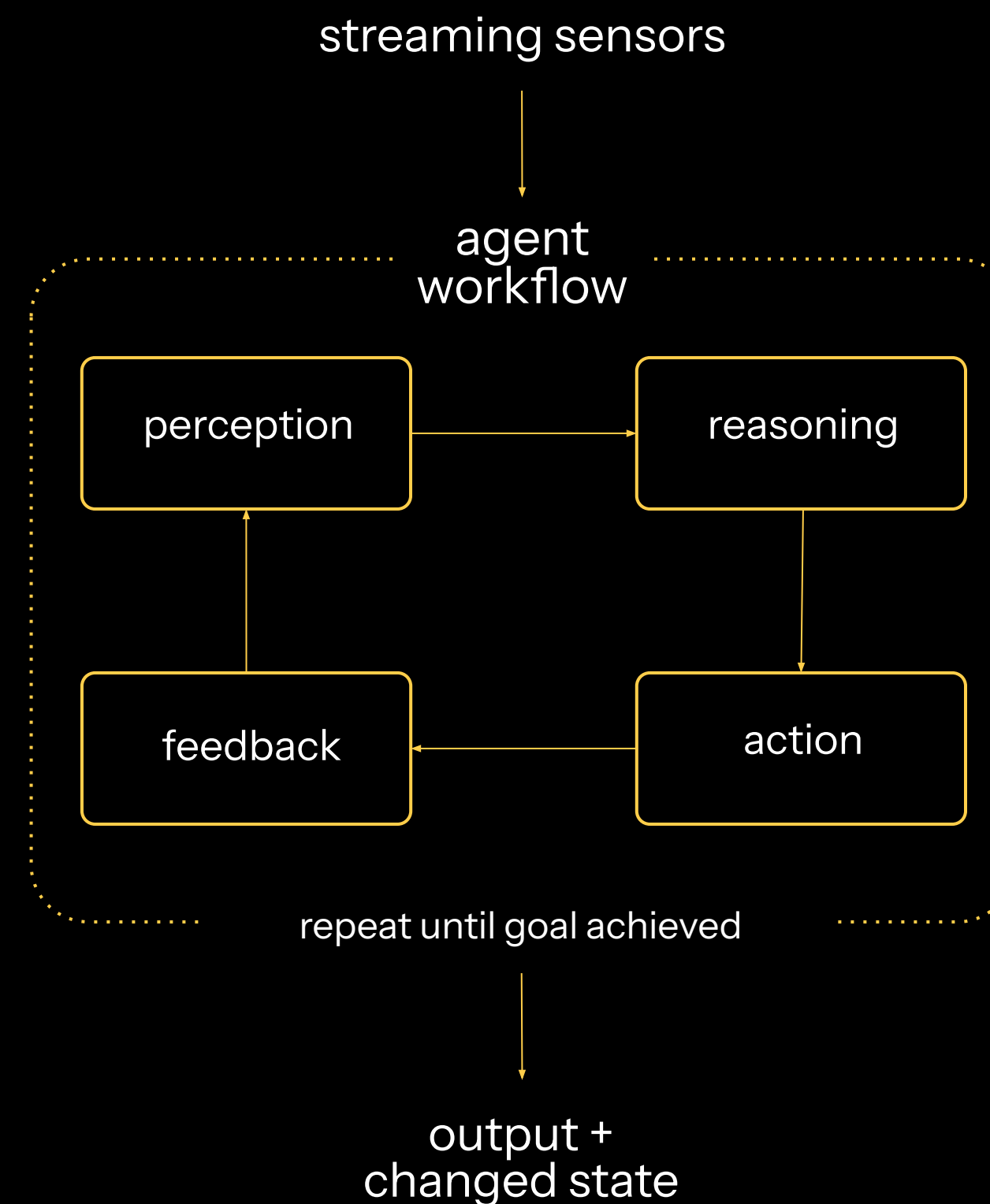
De-coupled, event-driven patterns and control loops

AKKA

## Retrieve - augment
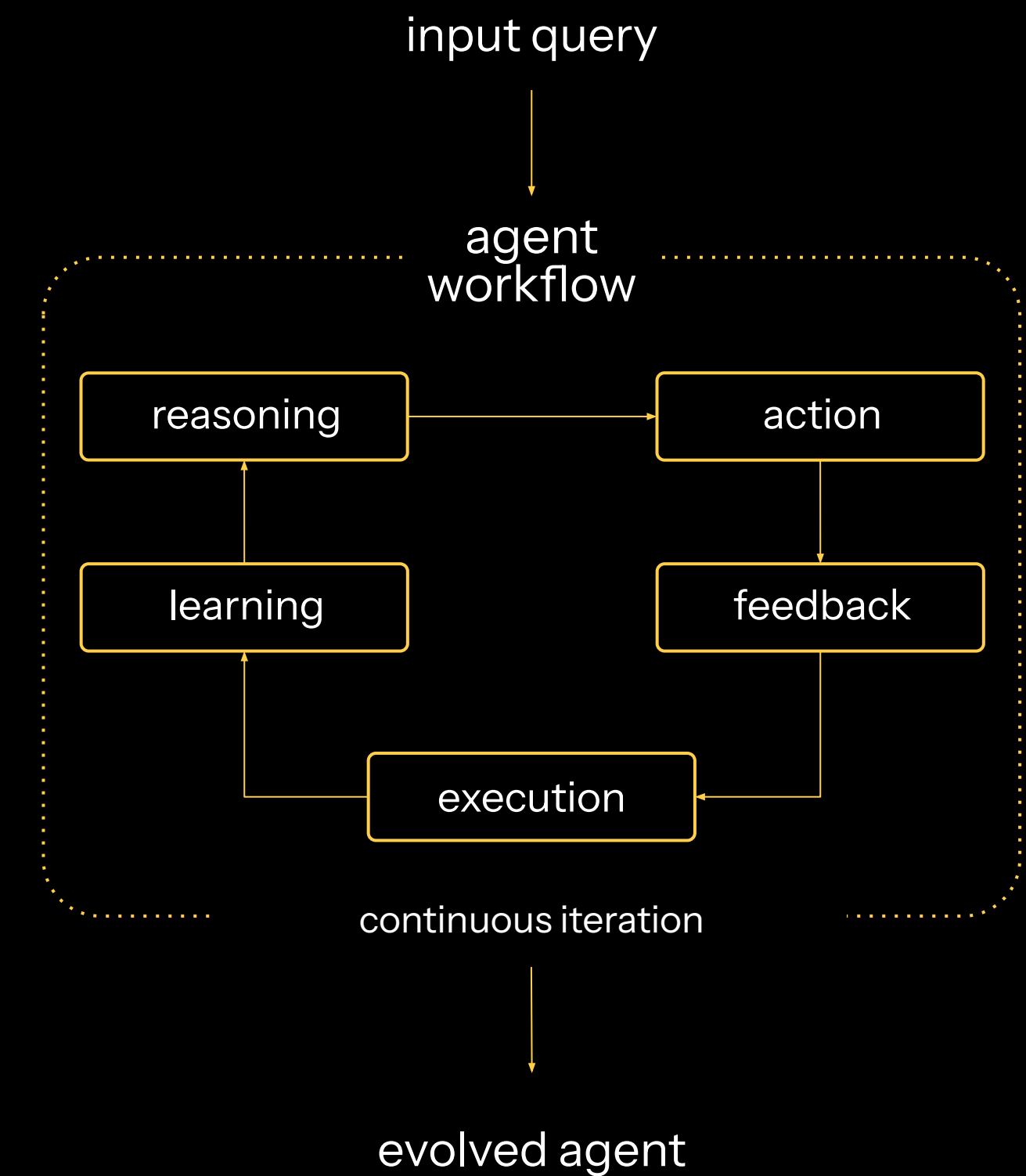agents that combine external knowledge with reasoning and action

trigger

agent workflow

knowledge retrieval

crawling fan-out

reasoning

augment prompt

repeat until desired outcome

output + action

## Environment controllers
control environments in real-time for robotics, edge, and automation

streaming sensors

agent workflow

perception → reasoning

feedback ← action

repeat until goal achieved

output + changed state

## Self learning
agents that improve themselves over time through self-reflection and environment adaptation

input query

agent workflow

reasoning → action

learning

feedback

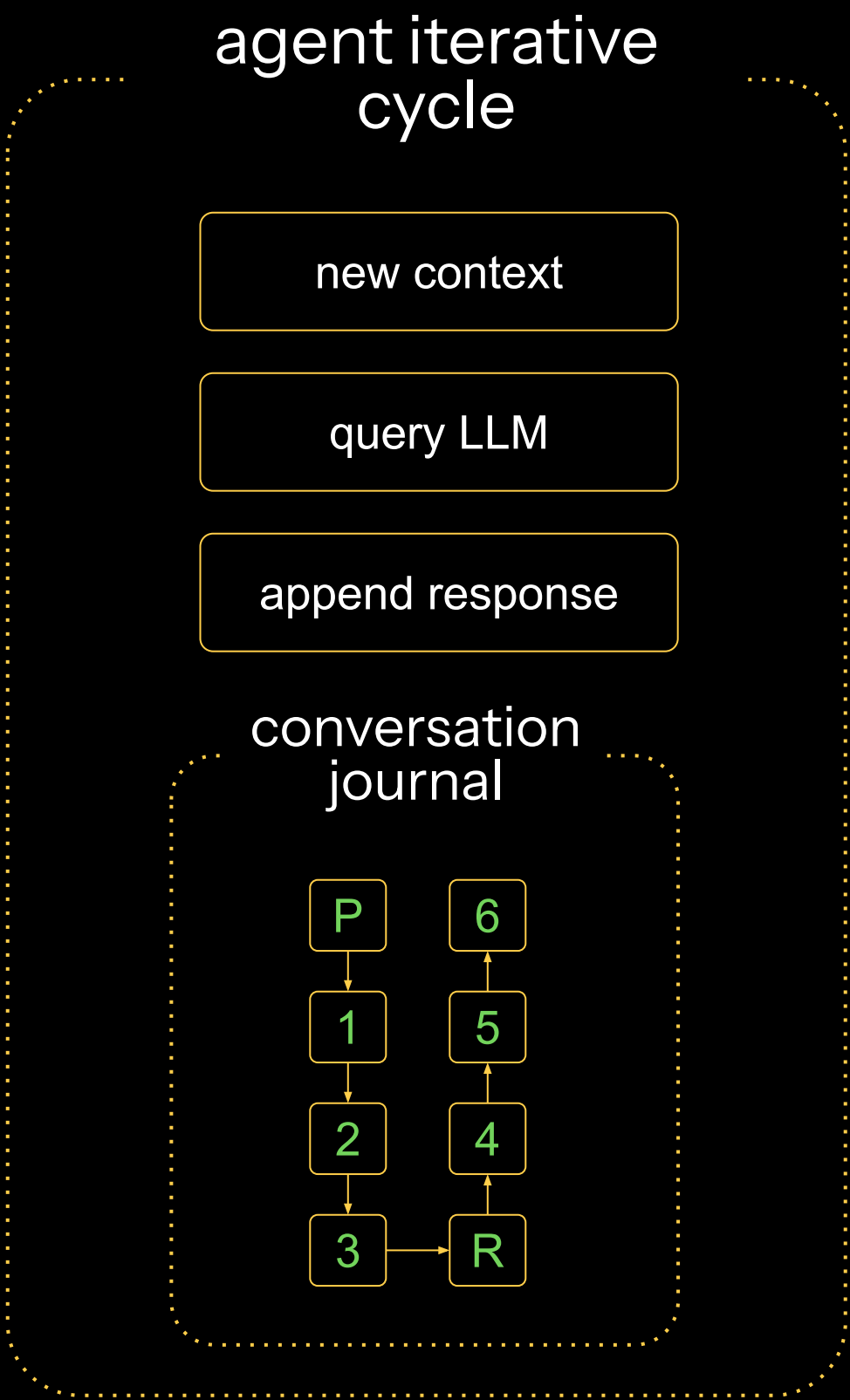execution

continuous iteration

evolved agent

# Agentic AI augmentation cycle

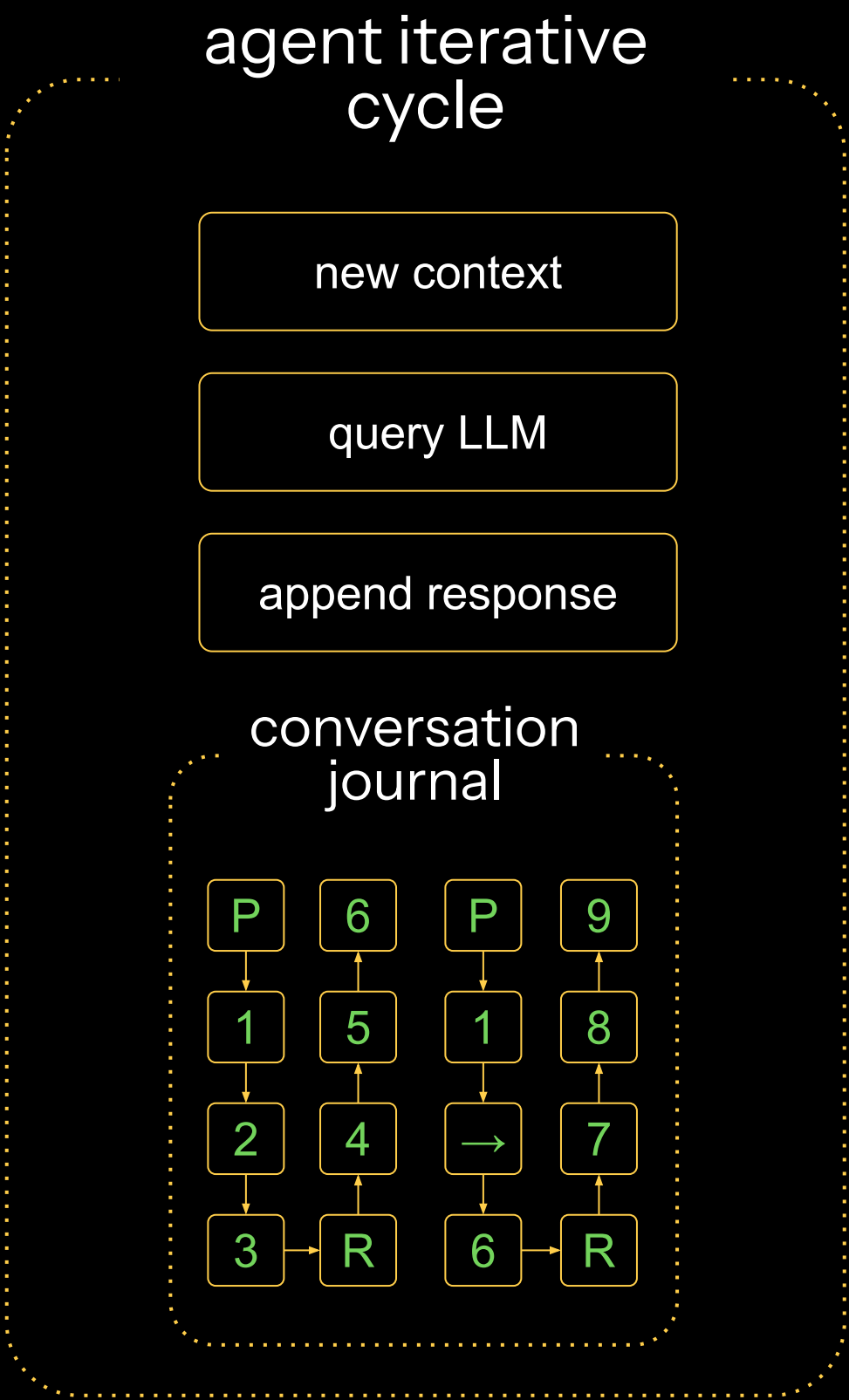Agents slow down on each iteration as context grows

## Agents start fast
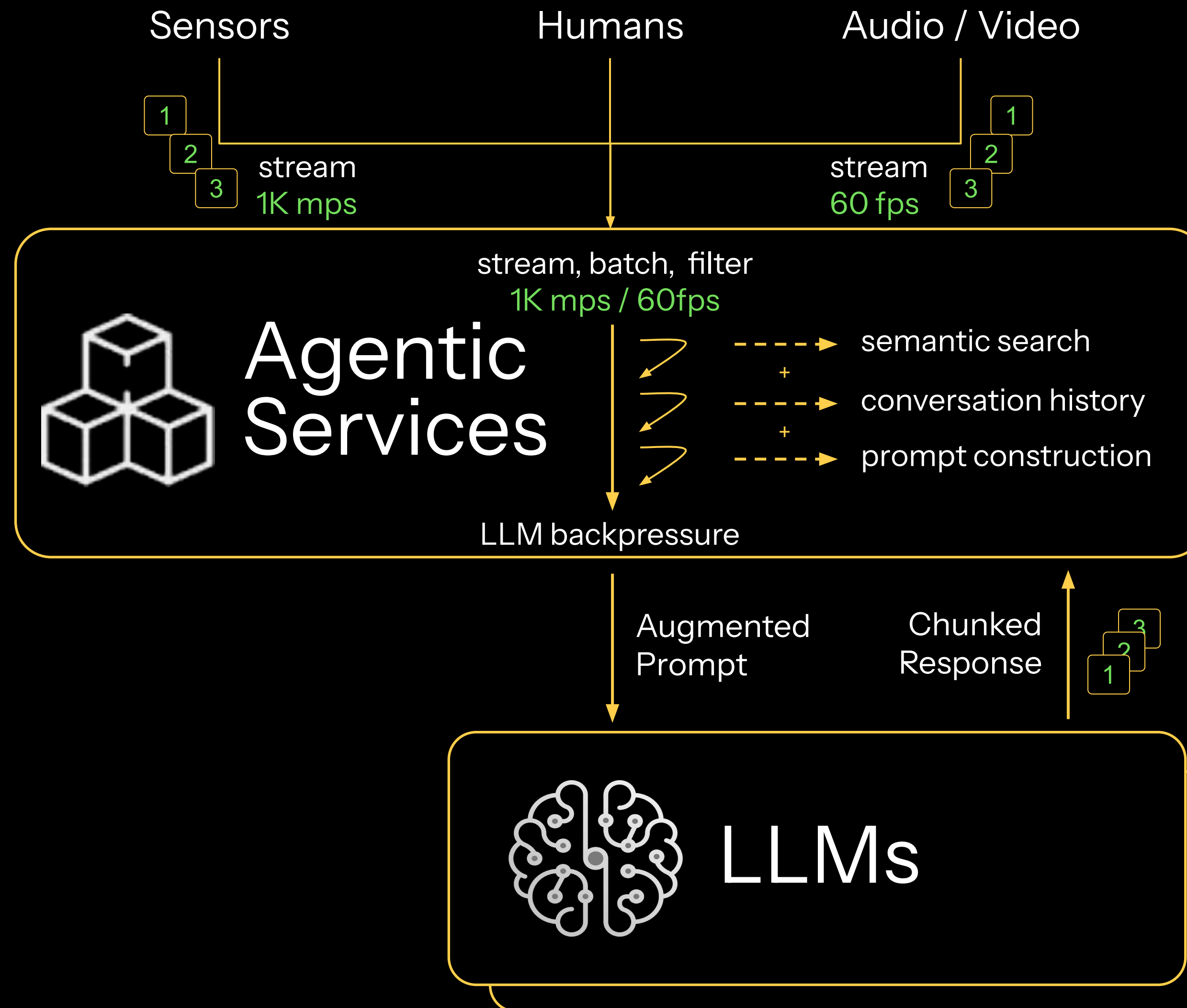small prompts, small conversations generate quicker responses

## Agent iterations grow slower
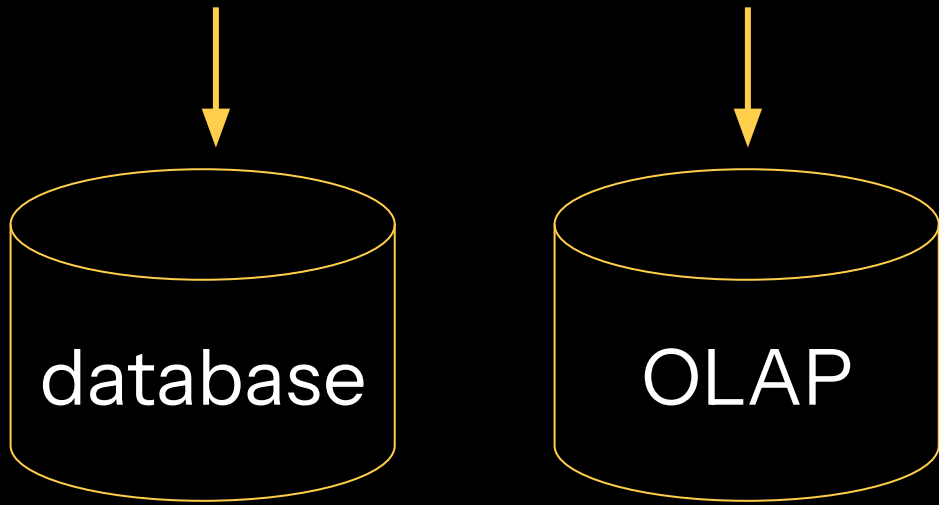Conversations and prompts grow, eventually hitting LLM token cap

# From n-tier to **a-tier** architecture

Humans + devices augmented with dozens of agent assistants that never sleep

AKKA

**Web Tier**

**API Tier**

database

OLAP

| | | |
|---|---|---|
| TPS | thousands | hundreds |
| p(99) latency | 10-50ms | 50-300 ms |

**Agentic Tier**

Agent lifecycle
AI orchestration
Context augmentation

Conversational
Streaming response
Poor latency
Not concurrent

Memory-hungry
Slow encoding algos

LLMs

vector db

event db

| | | | |
|---|---|---|---|
| TPS | 100x | 5x | 100x |
| p(99) latency | 150 - 3000ms | 50-200ms | 5-150ms |

# Agentic scale **requires efficiency**

More txs: each slower, less predictable and more costly

AKKA

|  | **SaaS** | **Agentic** |
|---|---|---|
| Users | billions | 20x |
| TPS | 10,000 | 100x |
| p(99) Latency | 10-80ms | 15-400x |
| Cost / LLM tx | cheap | 10-10,000x |

Mar 25: the best performing LLM @ 86% MMLU accuracy costs $98 / 1M tokens, or ~850,000x more expensive than the average database transaction. The worst performing LLM @ 36% MMLU accuracy costs $.01 / 1M tokens, or 7x more expensive.

# Bumpy path from POC to production

AKKA

**52%**

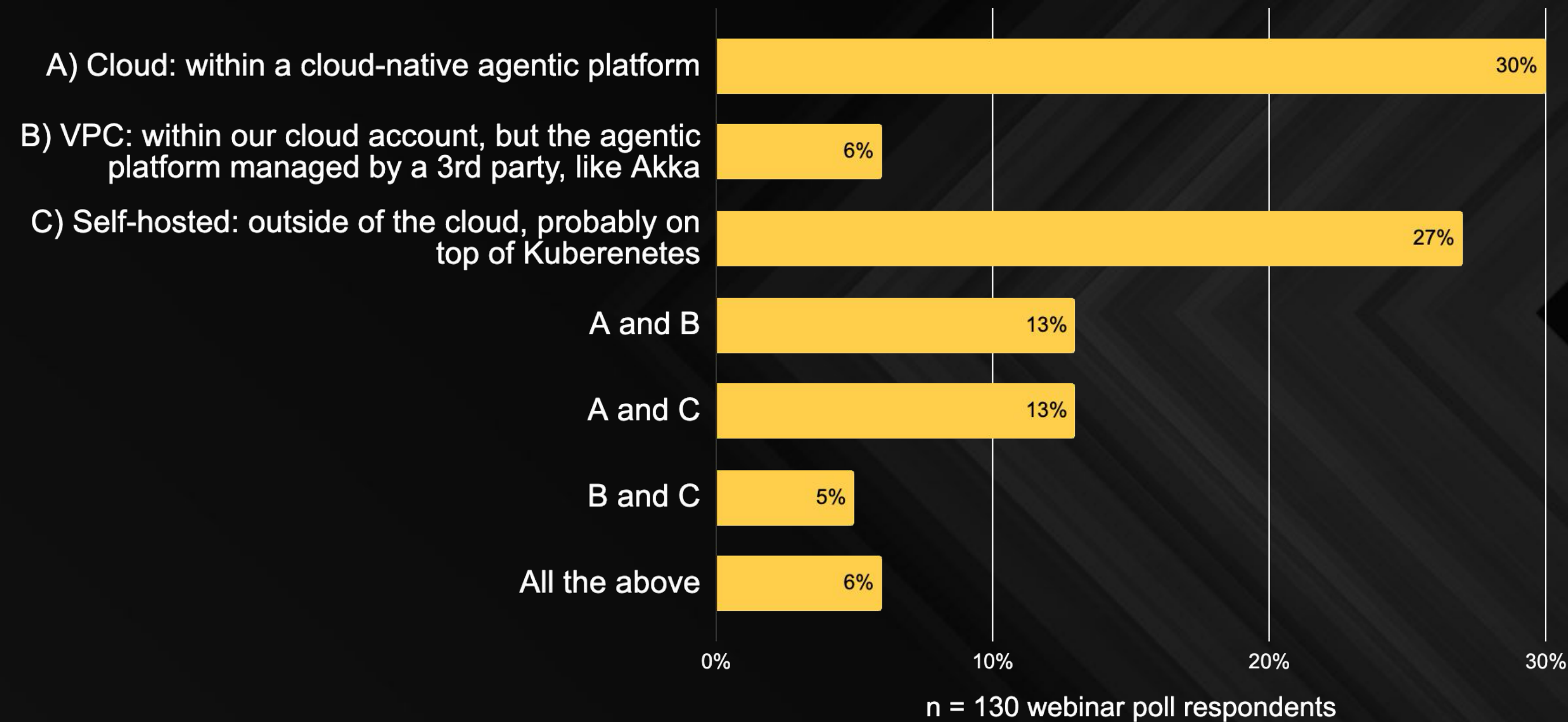fail to reach production

**8+ months**

POC to production

"Leaders reported that only 48% of AI POCs (Proof Of Concept) make it into production, and they take an average of 8.2 months to go from POC to production."
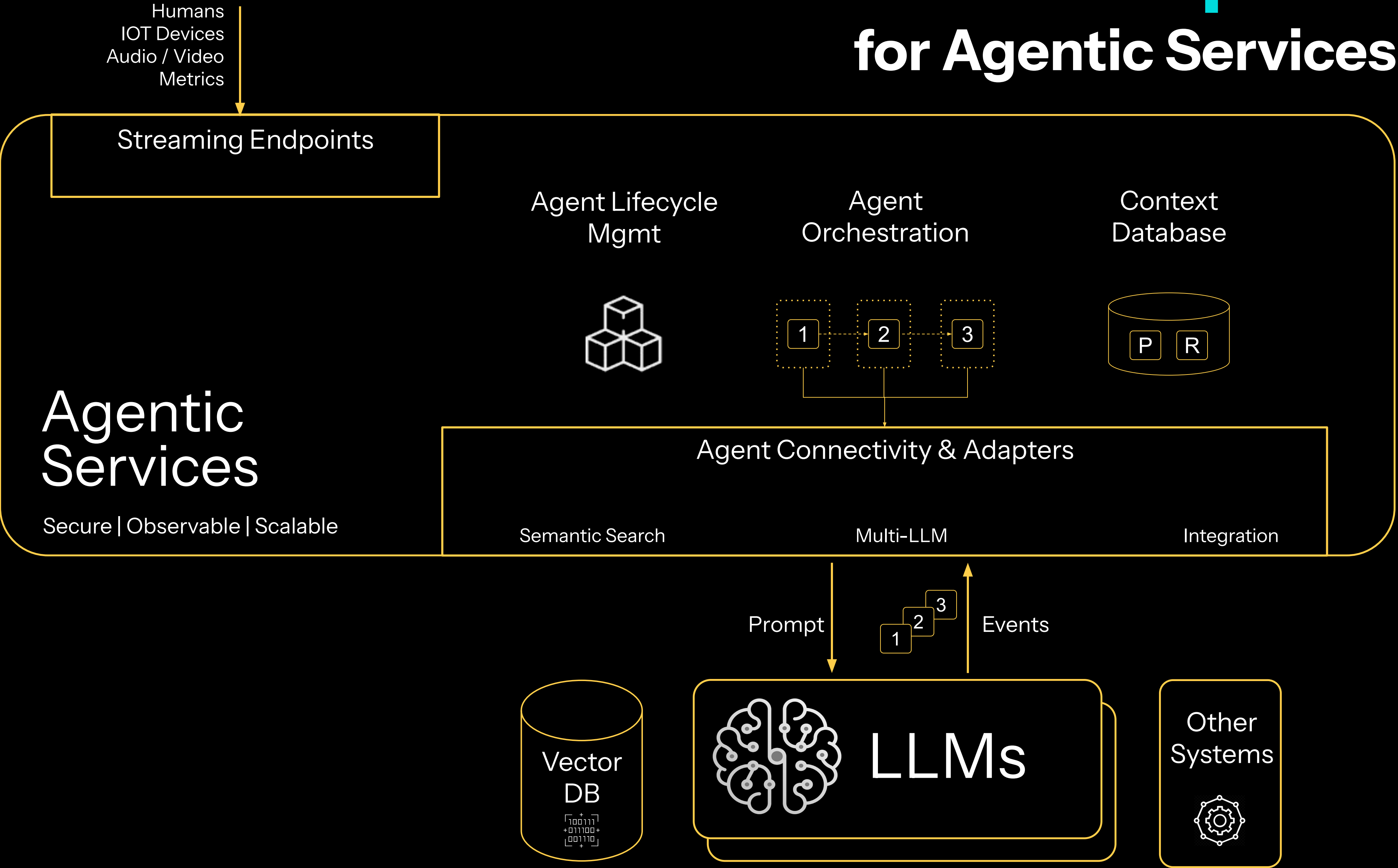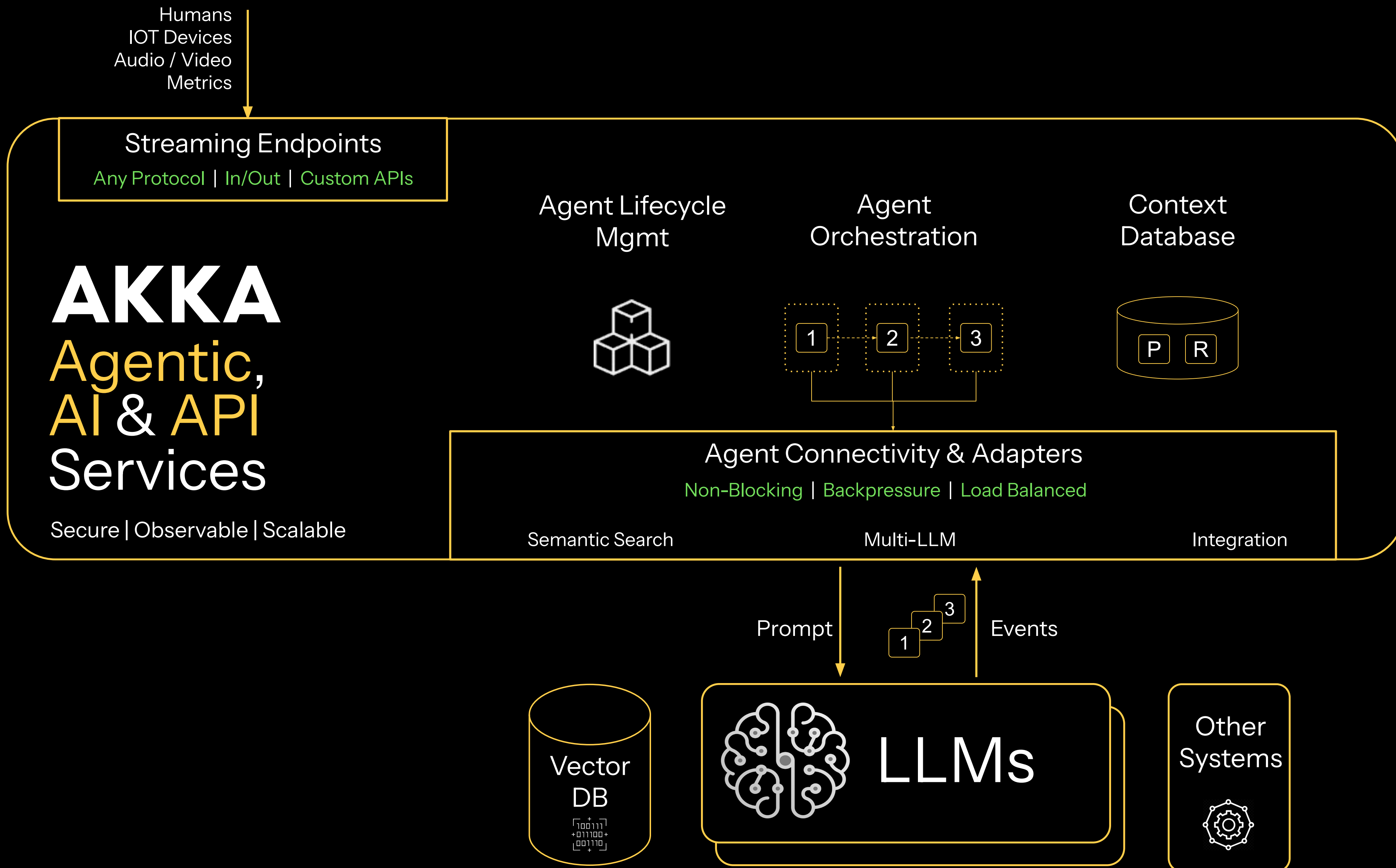
Gartner

# Akka **accelerates delivery** of agentic AI apps

## AKKA
### Shared SDK. Shared compute.

SaaS Applications

+

Agentic AI Services

**Enhanced User Experience**

AI agents personalize interactions to increase satisfaction

**Operational Efficiency**

AI agents automate routine tasks to allow humans to focus on strategic initiatives

**Scalability**

AI-driven SaaS adapt to business needs without proportional increases in cost

avoid the workflow island - orchestration without streaming, context database, or custom API endpoints

avoid the framework trap - dev tools with locking, concurrency, & memory not suited for 24/7 ops

crewai    Temporal    inngest    LangChain

# The Akka agentic advantage

✓ Agentic, AI, apps & data

✓ Hardened runtime

✓ Simple, expressive SDK

✓ Multi-region

✓ Automated ops

## Streaming endpoints

➔ Shared compute: agentic co-execution with API services
➔ HTTP and gRPC custom API endpoints
➔ Custom protocols, media types, and edge deployments
➔ Real-time streaming ingest, benchmarked to over 1TB

## Context database

➔ Agentic sessions with infinite context
➔ Context snapshot pruning to avoid LLM token caps
➔ In-memory context sharding, load balancing, and traffic routing
➔ Multi-region context replication
➔ Replication filters for region-pinning user context data
➔ Embedded context persistence with Postgres event store

## Agent connectivity & adapters

➔ Non-blocking, streaming LLM inference adapters with back pressure
➔ Multi-LLM selection
➔ LLM adapters & 100s of ML algos
➔ Agent-to-agent brokerless messaging
➔ 100s of 3rd party integrations

## Agent orchestration

➔ Event-driven runtime benchmarked to 10M TPS
➔ SDK with AI workflow component
➔ Serial, parallel, state machine, & human-in-the-loop flows
➔ Sub-tasking agents and multi-agent coordination

## Agent lifecycle management

➔ Agent versioning
➔ Agent replay
➔ Event, workflow, and agent debugger
➔ No downtime agent upgrades

# 2B people experience Akka daily

### SMILE

A fast ML engine with 100s of ML & LLM inference, powering Google Earth
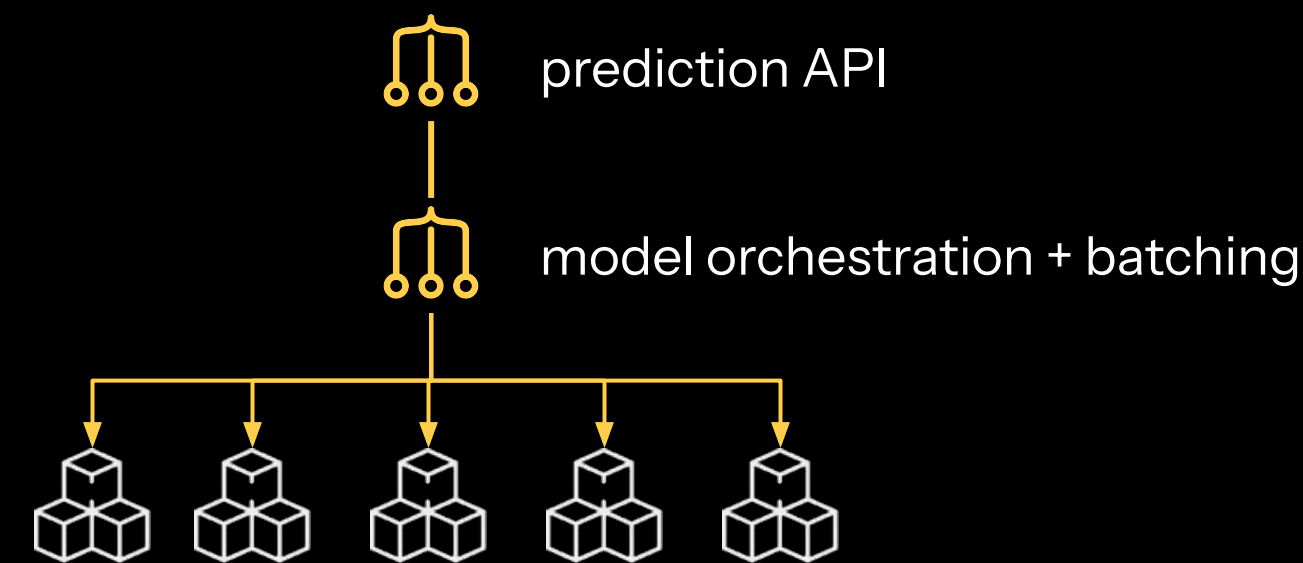
400K downloads / mo
6K GitHub stars

"Akka is used for streaming and back pressure - critical for hosted AI API inference. Akka enables event-driven inference exposed as HTTP efficiently, with low latency."

Haifeng Li – maintainer

### Swiggy

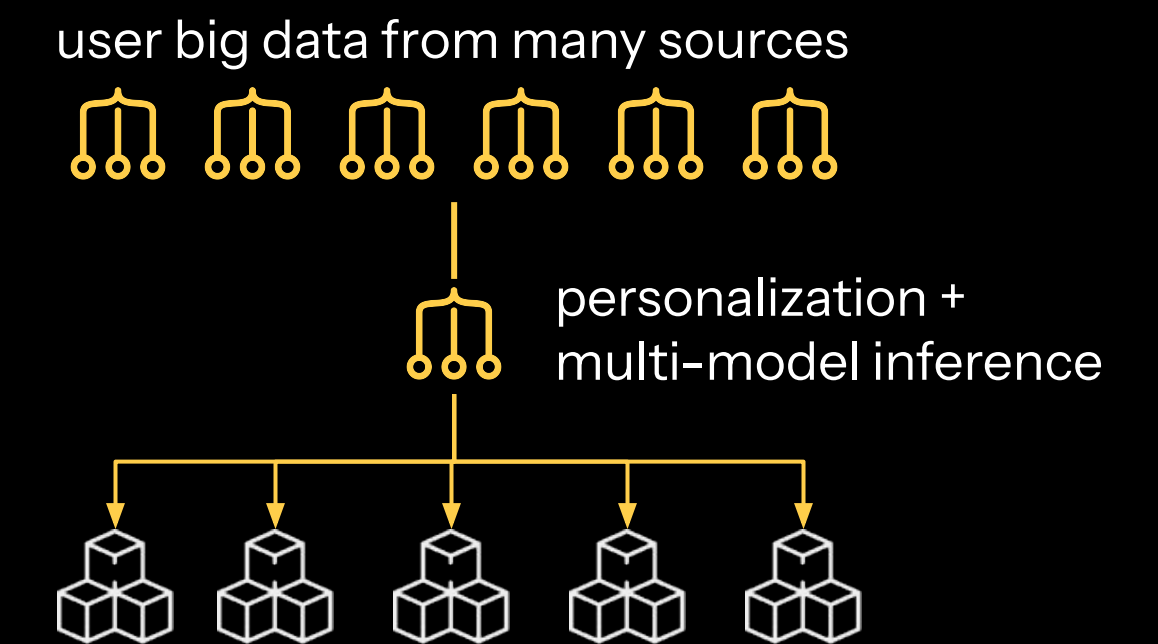API-driven predictions with multi-model fan-out and ultra-low latency

3+ million TPS
71ms p(99) latency

prediction API

model orchestration + batching

### Tubi [Fox]

Tubi applies ML models to real-time streams of data with in-memory, durable journals

2 month time to delivery

user big data from many sources

personalization + multi-model inference

### Horn

"Zero problems" augmenting high-performance audio and video streams on demand
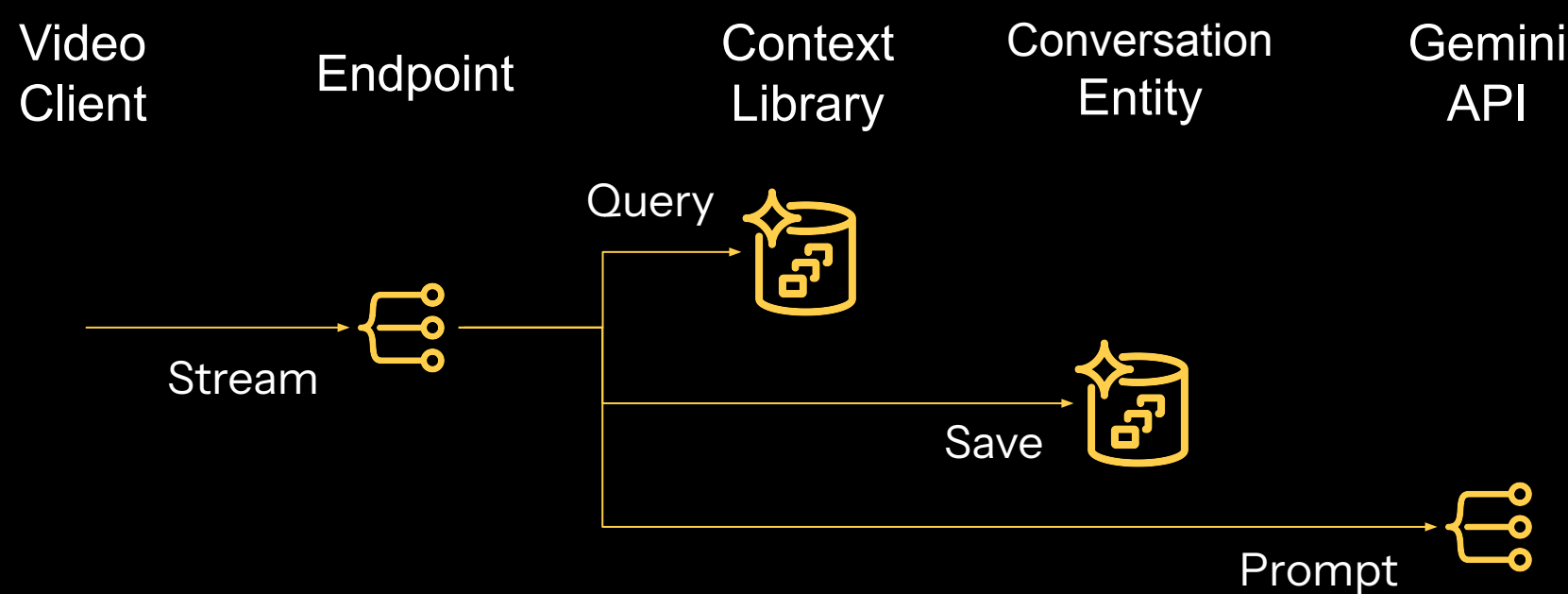Tomasz Wujec - Lead Developer

### Coho AI

"With Akka, we got to market 75% faster compared to other agentic solutions we had considered."
Michael Ehrlich – CTO

# Streaming Video Demo

| Video Client | Endpoint | Context Library | Conversation Entity | Gemini API |
|---|---|---|---|---|

Query · Stream · Save · Prompt

| Endpoint | gRPC API for receiving video |
|---|---|
| Context Library | Database of augmentation contexts saved as an Akka entity |
| Conversation Entity | Database of responses from Gemini |
| Gemini API | Google's API entry point |

## Velocity

### 145 LOC
2 components
gemini protocol client
4 integration tests
private GCP region
1 developer

### 2.5 days
concept to delivery

## Resilience

### 2x data redundancy
1 multi-cloud region; local data backup; add 2nd region for replication
99.9999% SLA for your apps → apps running across regions are nine-nines ready

### n/a ms
failover RTO

### 0ms
data RPO
Akka resilience guarantee – indemnities for reliability failures

## Cost Efficiency

### $150 / day
commodity cloud cost
10K TPS potential w/ this config

### 1M TPS
potential Akka write throughput
21M tx per $1 of cloud cost

# Concept-to-production in 8 weeks

| | |
|---|---|
| 1. Choose your agentic architecture | ➔ RAG, cooperative multi-agent, environment controller, or self-learning |
| 2. Select the right AI model | ➔ Prompt-based agents: GPT-4, Claude, Gemini, Mistra, Llama 2<br>➔ Embedding-based search agents: OpenAI Ada, Cohere, Google Vertex AI<br>➔ Fine-tuned industry models: Falcon, Mixtral |
| 3. Stand up agentic platform regions | InfoSec Review - Akka meets 19 levels of compliance including SOC 2 type 2<br><br>➔ Cloud: Akka Serverless<br>➔ Edge: Akka Edge<br>➔ Private: Akka BYOC |
| 4. Stand up AI inferencing | ➔ Cloud AI: OpenAI API, AWS Bedrock, Azure AI, Google Vertex AI<br>➔ Self-Hosted AI: Ollama, vLLM, TGI<br>➔ On-device AI: GTP4AII, LM Studio |
| 5. Build, test, debug and optimize | ➔ Build agents and agentic services offline with Akka's SDK<br>➔ Add human-in-the-loop features for oversight<br>➔ Run real-world performance, functional, and penetration simulations |
| 6. Deploy and observe | ➔ Setup API rate and cost limits to prevent abuse<br>➔ Record, track, and export performance or traces<br>➔ Monitor AI behavior for hallucinations or errors |

Live Q&A