Overview

# Kong AI Gateway

Accelerating the AI transformation

# API adoption is driven by **digital demand**



SOA
2001+

MOBILE
2007

CLOUD
2007+

MICROSERVICES
2014

COVID
2020+

GenAI
2022+

Applications — Cloud LLM

AI Agents — Self Hosted

ChatBots — Fine tuned

Other — Vector DBs

Applications

SEMANTIC CACHING
TOKEN RATE LIMITING
PROMPT FIREWALL
SECURITY
FAIL OVER
AI TRANSFORMATION
MORE

Cloud LLM

AI Agents

MULTI-LLM
AUTHN/Z
PROMPT TEMPLATING
ANALYTICS
CREDENTIALS MGMT
SECURITY
MORE

Self Hosted

ChatBots

OFF TOPIC DETECTION
SEMANTIC CACHING
RAG
PROMPT SANITIZATION
PROFANITY DETECTION
MORE

Fine tuned

Other

ANALYTICS
SECURITY
AUTHN/Z
RETRIES
MORE

Vector DBs

Proprietary and Confidential, Kong Inc

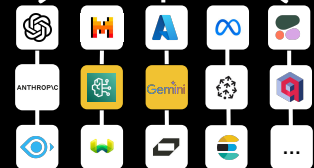API GATEWAY

INGRESS CONTROLLER

SERVICE MESH

KONG MESH

AI GATEWAY

# Kong Konnect

**Platform Applications**

- API Security (COMING SOON)
- API Analytics
- API Portal
- API Products (COMING SOON)
- API Logging (COMING SOON)
- Infra Diagnostics (COMING SOON)

**Core Platform**

- API Governance
- API Registry
- API Observability

**Plugins, Policies & AI**
- Caching
- AuthN/Z
- Traffic Control
- Observability
- Serverless
- AI Gateway
- Zero-Trust
- Routing
- EBPF
- + More

**Unified Control Plane**
- Gateway Manager
- Mesh Manager
- Insomnia Manager
- Billing
- Teams
- SSO
- RBAC
- Audit

**Service Hub**
- API Catalog

**API Infra, Automation & DevX**

- API Gateway
- AI Gateway
- Ingress Controller
- Service Mesh

**API Debug, Design & Testing**

Clients → API Gateway → Mesh

**API Automation**
- DecK
- Insomnia
- KonnectCTL

- Hybrid Self-Hosted
- Cloud Managed

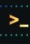Proprietary and Confidential, Kong Inc

**Applications**

**AI Gateway**

All Kong Plugins
AuthN/Z - Rate Limiting, Retries, Caching and more.

AI Governance | AI Observability
AI Credentials | AI Traffic Control

Unified API Interface | AI Proxy
AI Prompt Guard | AI Request Transformer
AI Prompt Template | AI Response Transformer
AI Prompt Decorator | Data Sanitization
AI Azure Content Safety | AI Rate Limiting Advanced

**AI Technologies**

Vector DBs

Pinecone | Qdrant | Milvus
Weaviate | Vespa | Elastic

LLMs

Cohere | Anthropic | LLaMA
OpenAI | Azure | Mistral
AWS Bedrock | Google Gemini | More...

One API to rule them all

Cohere
Anthropic
OpenAI
Azure
LLaMA
Mistral
AWS Bedrock
Google Gemini

AI Security &
AI Observability

AI Credentials

AI Traffic Control

AI Observability

And more

Cohere

Anthropic

OpenAI

Azure

One API to
rule them all

LLaMA

Mistral

AWS Bedrock

Google Gemini

AI Security &
AI Observability

AI Credentials

AI Traffic Control

AI Observability

And more

Cohere

Anthropic

OpenAI

Azure

One API to
rule them all

AI Prompt Guard

AI Prompt Decorator

AI Prompt Templator

AI Compliance &
AI Abuse prevention

LLaMA

Mistral

AWS Bedrock

Google Gemini

AI Security &
AI Observability

AI Credentials
AI Traffic Control
AI Observability
And more

One API to
rule them all

AI Request
Transformer

AI Response
Transformer

AI Prompt Guard
AI Prompt Decorator
AI Prompt Templator

AI Compliance &
AI Abuse prevention
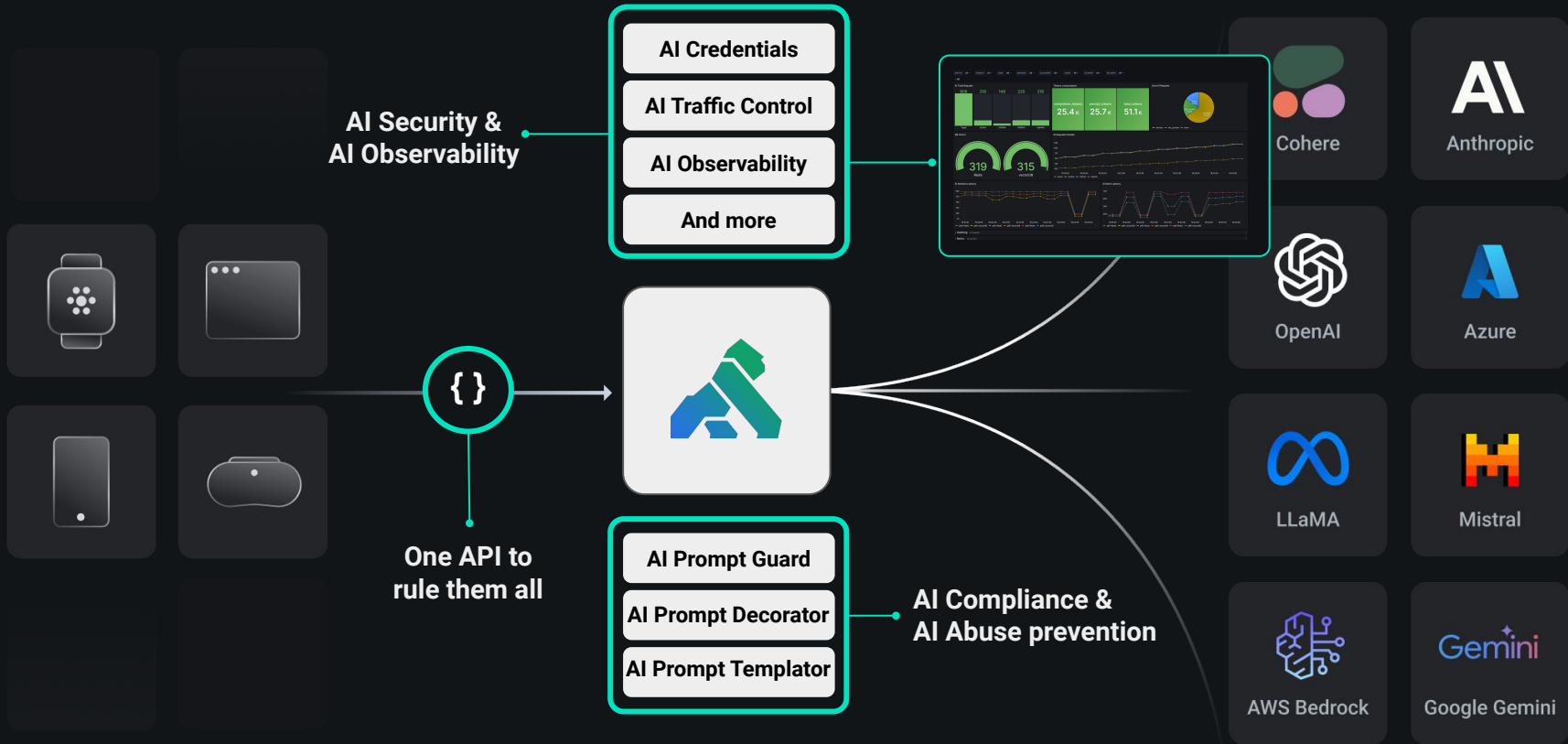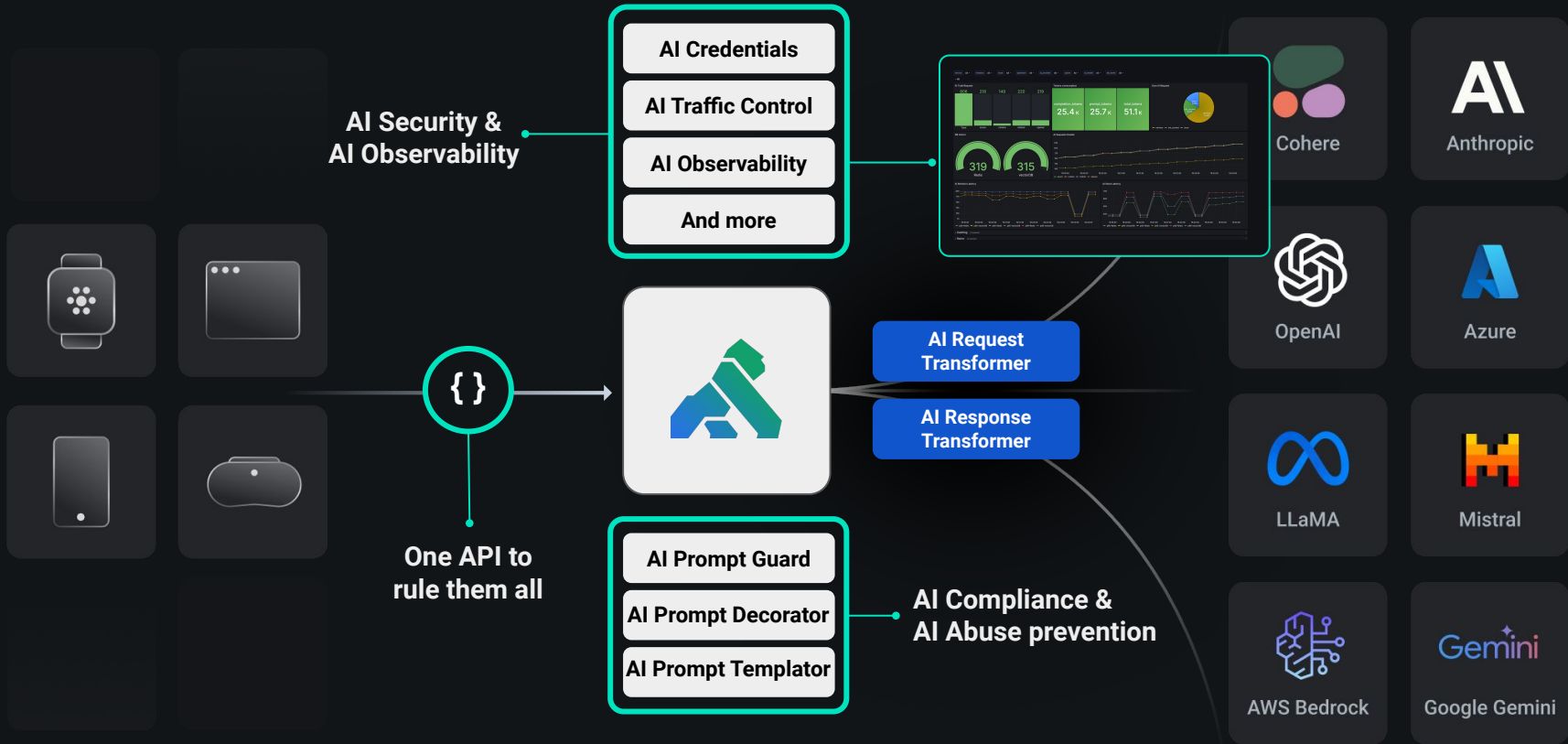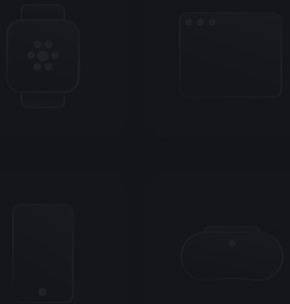
Cohere
Anthropic
OpenAI
Azure
LLaMA
Mistral
AWS Bedrock
Google Gemini

AI Security &
AI Observability

AI Credentials

AI Traffic Control

AI Logging

And more

One API to
rule them all
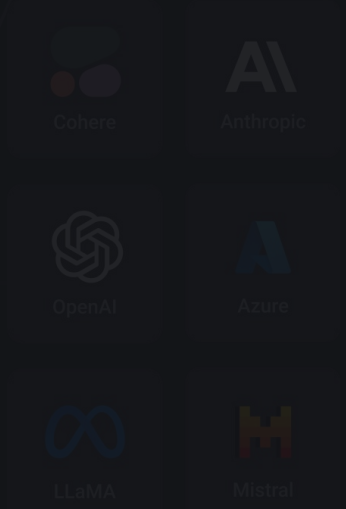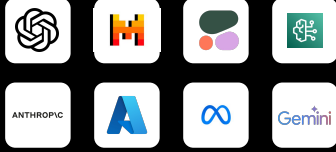
Prompt Guard

Prompt Decorator

Prompt Templator

AI Request
Transformer

AI Response
Transformer

Use AI without
coding

Cohere

Anthropic

OpenAI

Azure

LLaMA

Mistral

Client

ALL API + AI TRAFFIC

MODEL TRAINING

DATA ENRICHMENT

RAG

API / LLM

Applications

AI Agents

ChatBots

Copilot

HIGHER DEVELOPER PRODUCTIVITY

ENFORCING GOVERNANCE & COMPLIANCE

OBSERVABILITY AND OPTIMIZING AI SPEND

Cohere

Anthropic

OpenAI

Azure

LLaMA

Mistral

AWS Bedrock

Google Gemini

**May 2024 (3.7)**

- STREAMING SUPPORT
- AI TOKEN RATE LIMITING
- AZURE CONTENT SAFETY
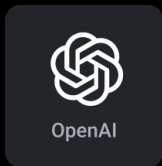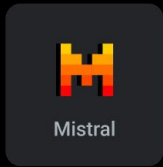- URL-SOURCED ROUTING
- CLAUDE 2.1 MESSAGES
- OPENAI SDK COMPATIBILITY

**June 2024 (3.7.1)**

- GCP VERTEX LLM SUPPORT
- AWS BEDROCK LLM SUPPORT
- ADVANCED LOAD BALANCING

**Q3-Q4 (3.8+)**

- SEMANTIC CACHING
- SEMANTIC PROMPT GUARD
- SEMANTIC ROUTING
- PROMPT JAILBREAKING SEC
- HALLUCINATION MITIG.
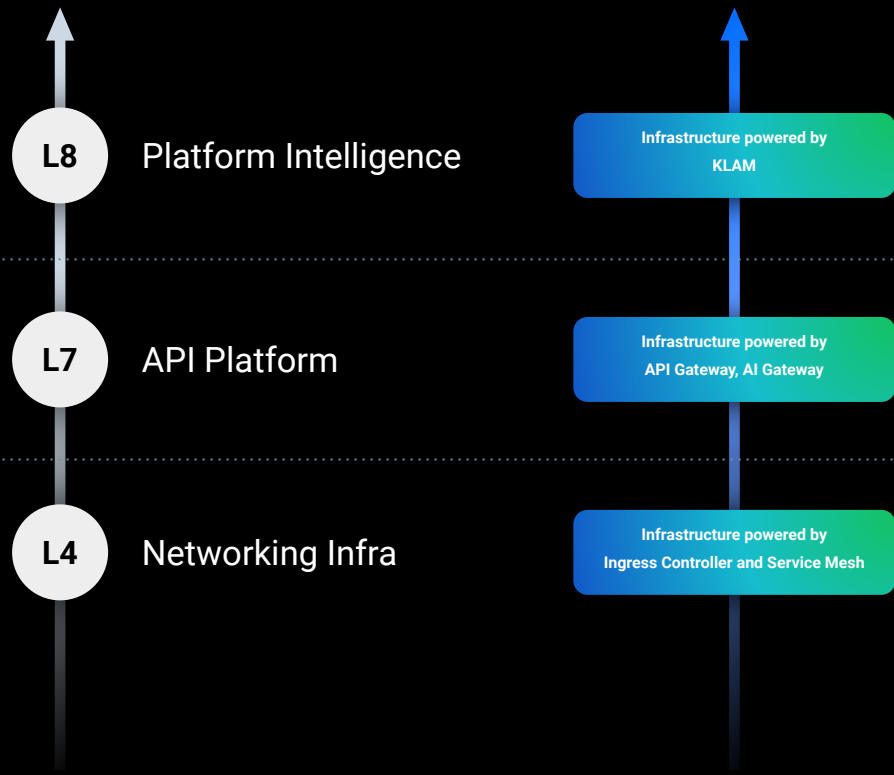- PROMPT INJECTION SEC
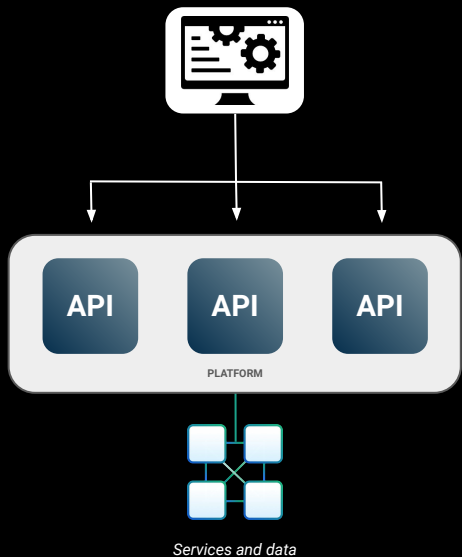- PROFANITY MITIG.
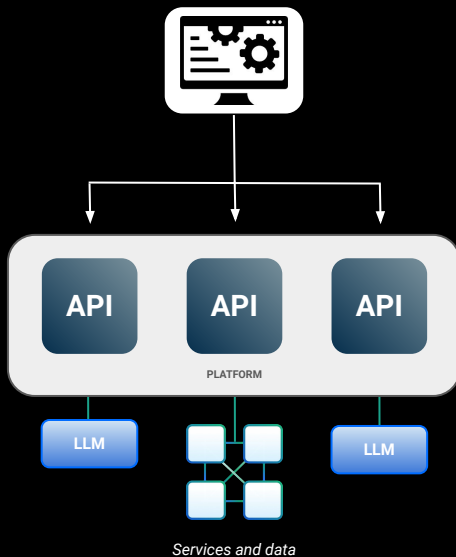- AI TRAIN PLUGIN
- SEMANTIC AI RAG

**Semantic AI Plugins (Aug 2024)**

**L8** Platform Intelligence
Infrastructure powered by KLAM

**L7** API Platform
Infrastructure powered by API Gateway, AI Gateway

**L4** Networking Infra
Infrastructure powered by Ingress Controller and Service Mesh

# Kong AI Gateway

konghq.com