# Zero Trust
In the era of Artificial Intelligence

**José Dores**
Principal Solutions Architect
Cloudflare

# AI is the next paradigm shift…

# AI is the next paradigm shift...

November 30, 2022    Product

## Introducing ChatGPT

Try ChatGPT ↗    Try ChatGPT for Work >

▶ Listen to article    0:04    🔗 Share

We've trained a model called ChatGPT which interacts in a conversational way. The dialogue format makes it possible for ChatGPT to answer followup questions, admit its mistakes, challenge incorrect premises, and reject inappropriate requests.

# AI is the next paradigm shift...

November 30, 2022    Product

## Introducing ChatGPT

Try ChatGPT ↗    Try ChatGPT for Work >

▶  Listen to article    0:04                                    🔗 Share

We've trained a model called ChatGPT which interacts in a conversational way. The dialogue format makes it possible for ChatGPT to answer followup questions, admit its mistakes, challenge incorrect premises, and reject inappropriate requests.

Give me a list of 10 mind blowing things that have become available or mainstream due to AI since ChatGPT's launch in November 2022

# AI is the next paradigm shift...

Give me a list of 10 mind blowing things that have become available or mainstream due to AI since ChatGPT's launch in November 2022

November 30, 2022   Product

## Introducing ChatGPT

Try ChatGPT ↗   Try ChatGPT for Work ›

▶ Listen to article   0:04   &#128279; Share

We've trained a model called ChatGPT which interacts in a conversational way. The dialogue format makes it possible for ChatGPT to answer followup questions, admit its mistakes, challenge incorrect premises, and reject inappropriate requests.

**Just** 3 years later...

Here are those 10 advancements, each summarized in a single line.

1. AI now generates high-fidelity, cinematic video clips from simple text prompts (e.g., Sora).

2. Models like GPT-4o can see, hear, and speak in real-time, understanding live video and audio.

3. AI can perfectly clone your voice and translate your speech while preserving your unique voice.

4. Autonomous AI agents can execute complex, multi-step tasks for you, like booking trips or coding.

5. AI is supercharging scientific discovery by modeling all of life's molecules to create new drugs.

6. AI coding assistants can now understand your entire codebase to autonomously write and debug features.

7. You can now generate complete, studio-quality songs with vocals from a simple text prompt (e.g., Suno).

8. AI models show advanced reasoning, capable of passing difficult professional exams like the Bar.

9. AI search engines (like Perplexity) give direct, synthesized answers to questions instead of just links.

10. AI is now in wearable devices (like glasses) that can see your surroundings and provide real-time info.

**CLOUDFLARE**

# …but it comes with its own risks

| Limited Visibility | Complex Risk Management | Ineffective Data Governance |
|---|---|---|
| **85%** | **63%** | **93%** |
| of IT decision-makers say employees are adopting AI tools even before their IT teams can assess them[1] | of breached organizations lack AI governance policies[2] | of employees admit to putting information into AI tools without approval[3] |

Sources: (1) (3) ManageEngine, July 2025; (2) IBM, July 2025;

CLOUDFLARE

# So what should be your attitude towards AI adoption?

**Limit use to specific
AI apps with strict guidelines**
(ex. "block-first" mentality)

**Prioritize building AI apps in-house**
(ex. Developers drive AI innovation)

**More
conservative**

**More
experimental**

**Encourage use of diverse AI apps**
(ex. only block with a strong justification)

**CLOUDFLARE**

# Enable safer, faster AI transformation with Cloudflare

**Protect**

Protect workforce use of genAI

Protect agentic AI access

Protect AI-powered apps

**Build**

Build full-stack
AI apps and AI agents

## Mission

Empower every organization to innovate with AI without compromising security.

*Cloudflare is the place to build world class, secure AI experiences.*

8

**CLOUDFLARE**

# Workforce use of Gen AI
## comes with privacy, security and compliance risks

**Employees can use unapproved or sketchy third-party AI ("Shadow AI")**

- 91% of employees admit using "shadow AI"

> **AI chatbot startup WotNot leaks 346,000 files, including passports and medical records**
>
> An Indian AI startup that helps businesses build custom chatbots has leaked almost 350,000 sensitive files after the data was left unsecured on the web.

**Employees can violate legal agreements or compliance frameworks by using AI**

- e.g. hiring decisions, patient / customer data

> News | Article | May 30, 2025
> **Health care workers are leaking patient data through AI tools, cloud apps**

**Employees can leak sensitive information by using AI**

- Financials, secrets, credentials, code, PII..

> **Samsung workers made a major error by using ChatGPT**
>
> News  By Lewis Maddison published 4 April 2023
>
> Samsung meeting notes and new source code are now in the wild after being leaked in ChatGPT

10

# SASE is an architecture design for Zero Trust
## And the path to secure Gen AI

### Legacy corporate network

- Office
- MPLS
- MPLS
- MPLS
- MPLS
- Office
- MPLS
- MPLS
- Data Center
- MPLS
- Bottleneck
- Remote Users
- VPN
- Data Center
- Internet Apps
- SaaS Apps
- Private & Self-Hosted Apps

### SASE corporate network with Cloudflare One

- Office Users
- Data Centers
- Remote & External Users

**Connectivity Cloud**

One Network w/Security Built-in

One Control Plane & Interface

- Internet & SaaS Apps
- Branch & Cloud Locations
- Private & Self-Hosted Apps

CLOUDFLARE

# Cloudflare One
## Our Zero Trust SASE platform



**Office Users**

**Data Centers**

**Remote & External Users**

CLOUDFLARE®

| Secure access | Secure web gateway |
| SaaS app security | Browser isolation |
| Cloud email security | Data loss prevention |
| WAN as a Service | FW as a Service |

All edge services on
one network with one control plane

**Internet & SaaS Apps**

**Branch & Cloud Locations**

**Private & Self-Hosted Apps**

**Secure access**
simplify and secure connecting
any user to any resource

**Threat defense**
keep your users safe from threats over any
port and protocol

**Policy compliance**
simplify regulatory compliance and protect
sensitive data

**SaaS security**
visibility and control of
applications including email

**Network modernization**
improved productivity, simpler operations,
reduced attack surface

# Cloudflare One
## Securing Gen AI



### Cloudflare SASE platform

**People**

- "Risky users" (e.g. contractor, suppliers, BYOD,)
- Employees
- Developers

Clientless on-ramps

with or without device client

- Visibility for shadow AI (SWG)
- Zero Trust access controls (ZTNA + SWG + RBI)
- Prompt protection and guardrails (SWG + DLP)
- Risk and posture management (CASB)
- Application Confidence Scoring for genAI specific risks
- Monitor model usage (AI Gateway)

Block

Isolate

Redirect to private tentant

Authenticate

**AI apps**

Public AI apps (MCP hosts)
- ChatGPT
- Gemini
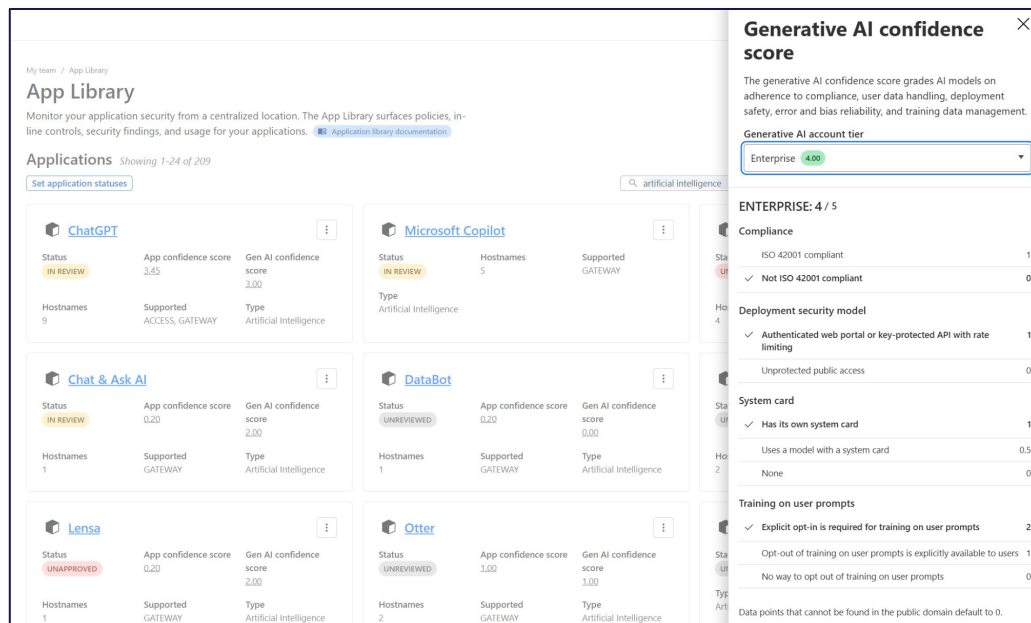- Claude

Private AI apps and infrastructure (MCP hosts)

**Step 01**
Discover Shadow IT

**Step 02**
Discover AI integrations with SaaS apps
Discover AI provider misconfigurations

**Step 03**
Block + redirect risky AI apps

**Step 04**
Control approved AI apps

- **Discover** Discover who is using which 3rd-party AI applications

- **Set Applications Status** as Approved, Unapproved, In Review, Unreviewed

**CLOUDFLARE**

# Cloudflare One
## Securing Gen AI



**Step 01**
Discover Shadow IT

**Step 03**
Block + redirect risky AI apps

**Step 02**
Discover AI integrations with SaaS apps
Discover AI provider misconfigurations

**Step 04**
Control approved AI apps

- **Discover third party AI apps** integrated with SaaS applications such as Microsoft 365 and Google Workspace

- **Disable unapproved** third party AI apps

**CLOUDFLARE**

# Cloudflare One
## Securing Gen AI



**Step 01**
Discover Shadow IT

**Step 02**
Discover AI integrations with SaaS apps
Discover AI provider misconfigurations

**Step 03**
Block + redirect risky AI apps

**Step 04**
Control approved AI apps

- **Discover misconfigurations** that can cause data loss or misuse on Gemini, Claude and ChatGPT

- **Discover sensitive data** that has been uploaded in chat attachments

- **Discover** GenAI specific insights
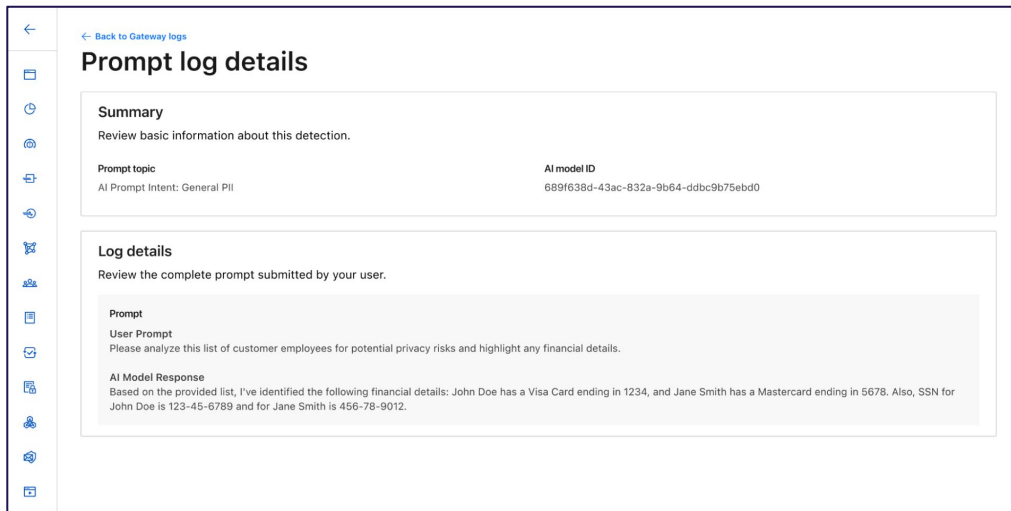
# Cloudflare One
## Securing Gen AI



**Step 01**
Discover Shadow IT

**Step 02**
Discover AI integrations with SaaS apps
Discover AI provider misconfigurations

**Step 03**
Block + redirect risky AI apps

**Step 04**
Control approved AI apps

- **Use** application and GenAI confidence scores to define status of AI tools, such as, Approved, Unapproved, In Review, Unreviewed of AI tools

- **Block** access to Unapproved and/or unreviewed AI tools

- Or **Redirect** from Unapproved and/or unreviewed AI tools to the approved ones

# Cloudflare One
## Securing Gen AI



**Step 01**
Discover Shadow IT

**Step 03**
Block + redirect risky AI apps

**Step 02**
Discover AI integrations with SaaS apps
Discover AI provider misconfigurations

**Step 04**
Control approved AI apps

- **Apply identity and device posture policies** so that only approved employees with compliant devices can access approved AI tools

- **Apply DLP and/or Isolation policies** to restrict how sensitive data is shared

- **Log prompts, classify prompts** based on content and intent and **apply guardrails** based on its classification

**Protect**

**Protect workforce use of genAI**

**Protect agentic AI access**

**Protect AI-powered apps**

**Build**

**Build full-stack
AI apps and AI agents**

CLOUDFLARE

# Example: creating an invoice



**Without MCP**: AI tells you what to do and you do it

CLOUDFLARE

# Example: creating an invoice



**With MCP**: the AI does it for you

CLOUDFLARE

# You can also interact with Cloudflare using MCP

**WARP Toggle Off Events Summary**

Based on the data from the **jdores-org** account for today (November 4, 2025), here are the findings:

**Total Toggle Off Events: 14**

However, looking at the data more closely, all 14 events came from the **same device**:

- **Device ID**: `2421adf7-b40b-11f0-b750-1ae08bc4d7dd`
- **Hostname**: `win10vm01`
- **Serial Number**: `8479-5166-5206-4065-3284-5980-02`

**Unique Users Who Toggled WARP Off: 1**

The events occurred between **15:03:55 UTC** and **15:07:08 UTC** today, suggesting this was likely one user repeatedly toggling WARP off on their Windows 10 VM device, possibly during troubleshooting or testing.

**Key Insights:**

- **1 unique user/device** toggled WARP off today
- **14 total toggle events** occurred within a 4-minute window
- All events were from the same Windows 10 virtual machine (`win10vm01`)

**Task completed successfully** - The query shows that only 1 user toggled their WARP client off today in the jdores-org account, though they did so multiple times in quick succession.

> I would like to see how many users have toggled off their WARP client today?

**CLOUDFLARE**

# Agentic AI and MCP
## comes with visibility, safety and security risks

**MCP servers can act autonomously against corporate resources**
- Can make errors, be attacked, or poisoned

**MCP Servers Hit by 'NeighborJack' Vulnerability and More**

By David Ramel | 06/25/2025

Model Context Protocol (MCP) servers, a relatively new idea from Anthropic to connect advanced AI systems with tools, data sources and other resources so they can act as autonomous agents, is now being subject to cybersecurity attacks like every other piece of software.

**Employees can use unapproved or unregistered MCP servers ("Shadow MCP")**
- No visibility into which MCP servers are used

**First Malicious MCP Server Found Stealing Emails in Rogue Postmark-MCP Package**

📅 Sep 29, 2025    👤 Ravie Lakshmanan

**Employees can share sensitive credentials with MCP servers**
- Leading to unsafe autonomous actions by MCP

**Today's MCP Horror Story: The GitHub Prompt Injection Data Heist**

Just a few months ago in May 2025, Invariant Labs Security Research Team discovered a critical vulnerability affecting the official GitHub MCP integration where attackers can hijack AI agents by creating malicious GitHub issues in public repositories. When a developer

# Cloudflare One
## Discover and control MCP servers



- **Create a central registry** of approved MCP servers

- **Discover** how users and interacting with MCP via comprehensive request logs

- **Control** who can access what MCP servers with least-privilege

**Protect**

Protect workforce use of genAI

Protect agentic AI access

Protect AI-powered apps

**Build**

Build full-stack
AI apps and AI agents

CLOUDFLARE

CLOUDFLARE

# Public-facing AI-powered apps
## come with reliability and security risks

**App can be coerced into producing wrong or embarrassing outputs**
- Prompt injection, jailbreaking, biased inputs

TECH & INNOVATION

**Microsoft's AI millennial chatbot became a racist jerk after less than a day on Twitter**

On Wednesday (Mar. 23), Microsoft unveiled a friendly AI chatbot named Tay that was modeled to sound like a typical teenage girl. The bot was designed to learn by talking with real people on Twitter and the messaging apps Kik and GroupMe. ("The more you talk the smarter Tay gets," says the bot's Twitter profile.) But the well-intentioned experiment quickly descended into chaos, racial epithets, and Nazi rhetoric.

**App have embedded AI that security teams are don't know about**
- Creating risks for the enterprise

Zenity Labs > Posts > AgentFlayer: When AIjacking Leads to Full Data Exfiltration in Copilot Studio

**AgentFlayer: When AIjacking Leads to Full Data Exfiltration in Copilot Studio**

Tamir Ishay Sharbat
July 07, 2025

**App can suffer from volumetric attacks**
- Because DDoS still works on LLMs

**NSFocus: DeepSeek AI hit with 'well planned' DDoS attacks**

Cybersecurity vendor NSFocus said AI startup DeepSeek endured multiple waves of DDoS attacks from attackers since its reasoning model was released Jan. 20.

By Alexander Culafi, Senior News Writer          Published: 03 Feb 2025

DeepSeek is facing a series of DDoS attacks, according to research published Friday by cybersecurity vendor NSFocus.

27

CLOUDFLARE

# OWASP Top 10 for Large Language Models (2025)

**Prompt Injection**
Manipulation of model through crafty inputs to influence decision

**Sensitive Information Disclosure**
Sensitive data being exfiltrated from the model

**Supply Chain Vulnerabilities**
Vulnerable component embedded in the model

**Data and Model Poisoning**
LLM training data is tampered with

**Improper Output Handling**
Output accepted without validation

**Excessive Agency**
Models can perform actions due to excessive permissions

**System Prompt Leakage**
Sensitive information in the prompt

**Vector and Embedding Weaknesses**
Malicious actors can exploit weaknesses in models using RAG

**Misinformation**
Excessive trust on the output of LLM leading to misinformation

**Unbound Consumption**
DDoS attacks, service degradation, etc

**CLOUDFLARE**

# Cloudflare Firewall for AI
## Empowers security teams to protect AI applications

**Firewall for AI**

**Request validation**

Protect the model   2

**Response validation**

3   Protect the data

**API Call**

```
{
    "user": "email",
    "prompt": "here
is my question for
the model"
}
```

**Front end**

**API Response**

```
{
    "success": "true",
    "response": "model
output":
}
```

1   **Discover LLMs**

**Application**

**LLM models**

- Automatically **discover** and label LLM endpoints

- **Visibility** & analytics on LLM-bound traffic and any associated risks

- **Protect and mitigate** against LLM threats and attacks using rules and actions

29

# Cloudflare Developer Platform
The secret is out...

December 2024:
3+ million

April 2024:
2+ million

June 2024:
2.4+ million

November 2022:
1+ million

May 2022:
450,000+

2019    2020    2021    2022    2023    2024

# 3+
# MILLION
developers on the
Cloudflare platform

CLOUDFLARE

# Your application + our network
Deploy once, run everywhere

**Deploy from region: Earth**

**335+**
**cities**
**in 125+ countries, including mainland China**

Code executes within 50ms of ~95% of the Internet-connected global population

**200+**
**cities**
**with GPUs**

Growing constellation of cities for AI inference powered by GPUs

**CLOUDFLARE**
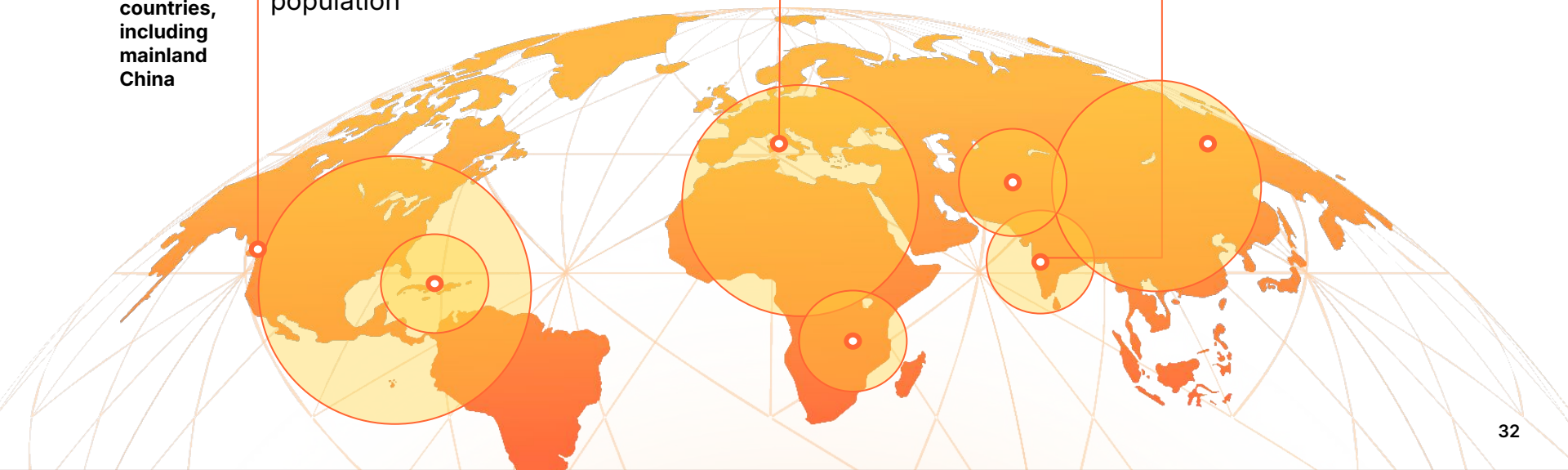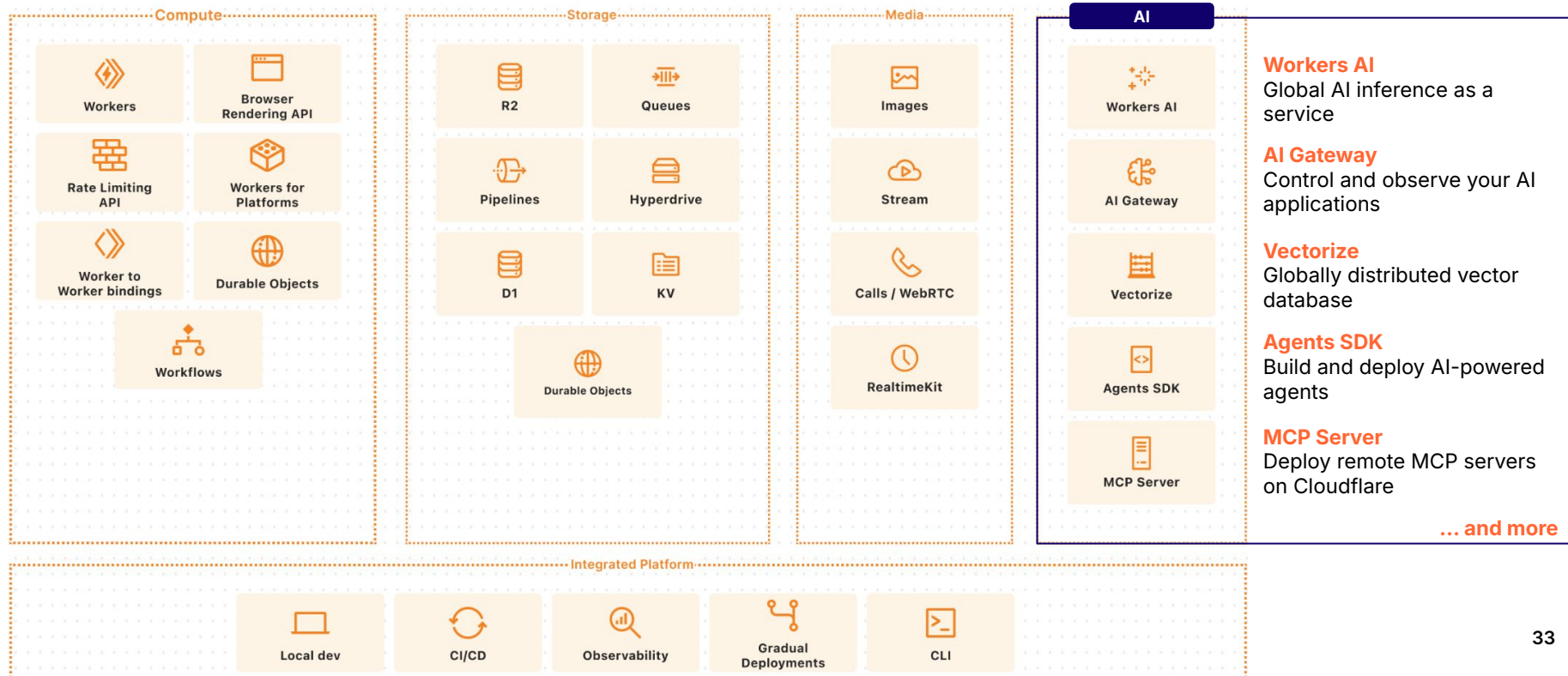
# Cloudflare is the best developer platform to build Secure AI experiences

## Compute

| | |
|---|---|
| Workers | Browser Rendering API |
| Rate Limiting API | Workers for Platforms |
| Worker to Worker bindings | Durable Objects |

Workflows

## Storage

| | |
|---|---|
| R2 | Queues |
| Pipelines | Hyperdrive |
| D1 | KV |

Durable Objects

## Media

Images

Stream

Calls / WebRTC

RealtimeKit

## AI

Workers AI

AI Gateway

Vectorize

Agents SDK

MCP Server

**Workers AI**
Global AI inference as a service

**AI Gateway**
Control and observe your AI applications

**Vectorize**
Globally distributed vector database

**Agents SDK**
Build and deploy AI-powered agents

**MCP Server**
Deploy remote MCP servers on Cloudflare

**… and more**

## Integrated Platform

| Local dev | CI/CD | Observability | Gradual Deployments | CLI |
|---|---|---|---|---|