CLOUDFLARE

# AI Inference at the Edge

**Logan Grasby**
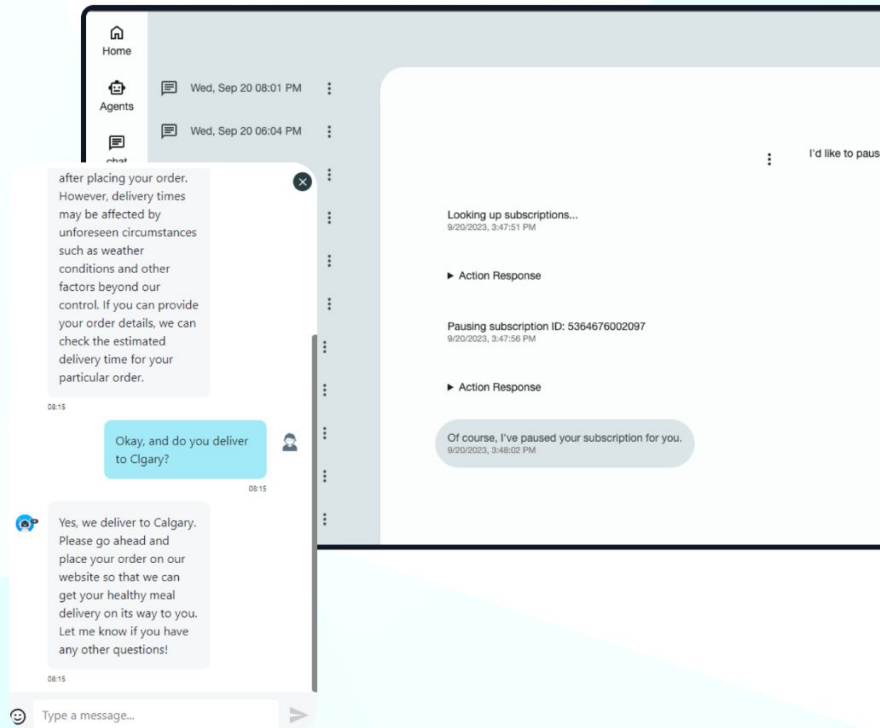**Developer Educator, AI at Cloudflare**

**Nvidia GTC 2024**

# Meet your newest employee

Azule delivers AI agents that your business can train to interact with your customers on your behalf.

Sign up >    Contact Us

**Tell me about the different types of technical indicators**

There are numerous technical indicators that can be used to analyze stock trends and help make decisions on when to buy or sell. Some popular types of technical indicators include moving averages, relative strength index (RSI), stochastic oscillator, Bollinger Bands, and MACD (Moving Average Convergence Divergence). Moving averages show the average price over a certain period of time while RSI measures the strength of a stock's price action. Stochastic oscillator is used to find oversold or overbought conditions, Bollinger Bands indicate volatility, and MACD shows trend momentum. It's important to note that no one indicator is perfect and they should be used in combination with other analysis tools to make informed decisions.

**Technical Indicators**



Type a message...

**Today Cloudflare serves on average over 55 million HTTP requests per second**

- **20% of internet traffic**
- **3 Terabytes of DNS queries per day**
- **8 Billion web pages crawled every 2 weeks**

# A history of AI/ML at Cloudflare

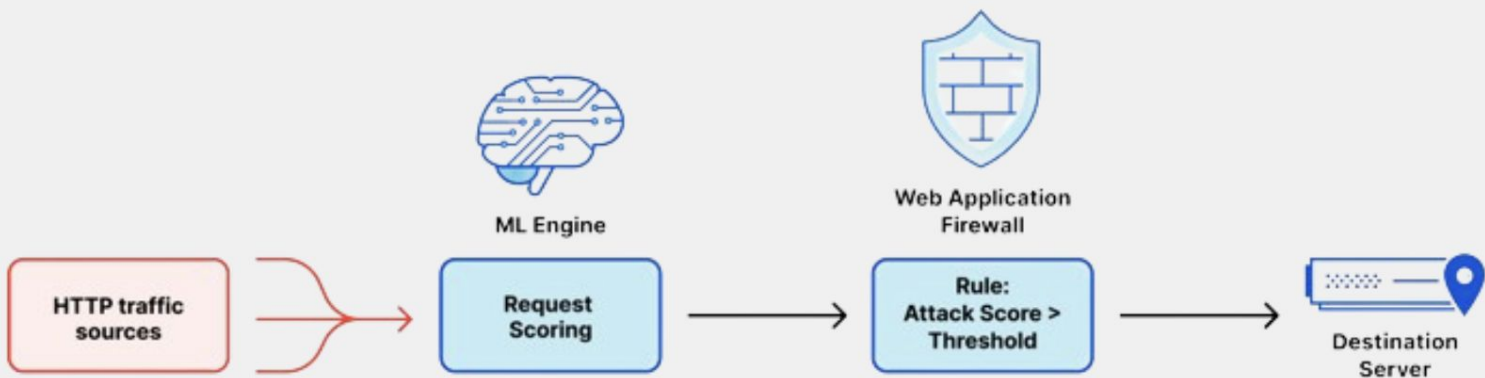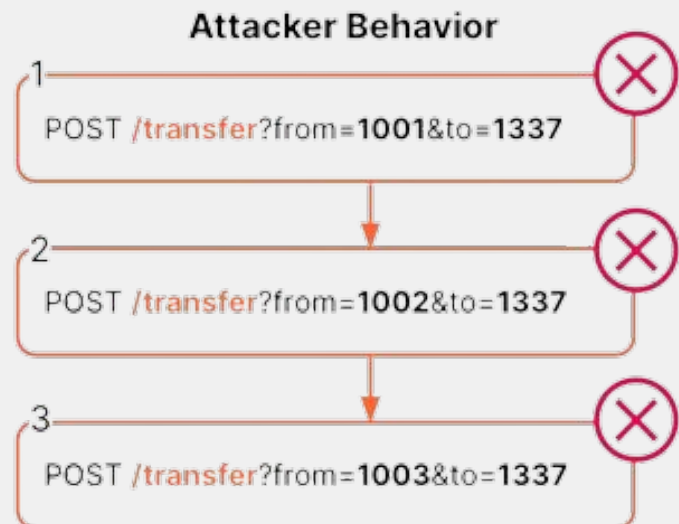## CloudFlare Uses Intelligent Caching to Avoid the Bot Performance Tax
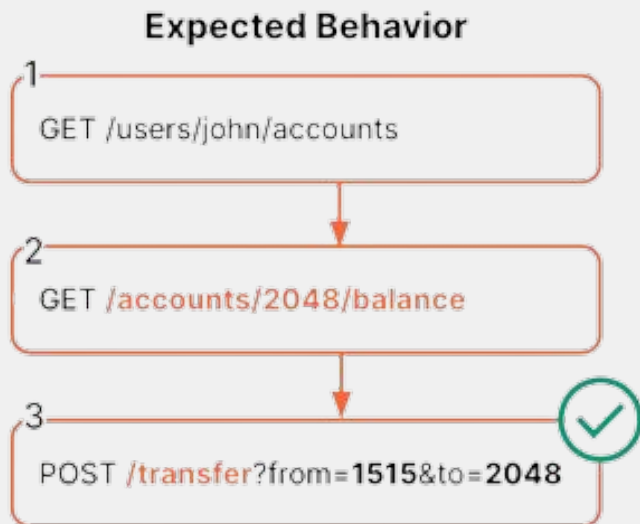
12/16/2011

Matthew Prince

# A history of AI/ML at Cloudflare

**AI/ML models power products like WAF and Bot Management at the edge**

# AI/ML at Cloudflare today

CLOUDFLARE

# Deploying GPUs around the world

- GPUs in 125+ locations globally
- Covering Cloudflare's entire network
- Global AI inference at the edge
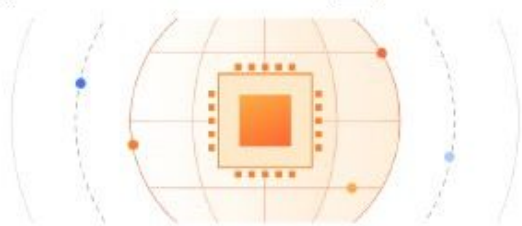
# Deploying GPUs around the world



Cloudflare Data Center City
Cloudflare GPU City
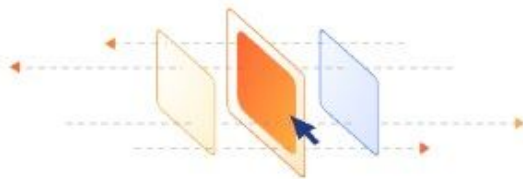
End of 2024

# Run inference on region: Earth

Build and deploy ambitious **AI applications** to Cloudflare's global network
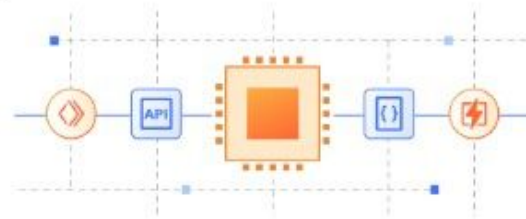
CLOUDFLARE

# AI accessible to everyone

### Serverless AI on GPUs

Run generative AI tasks on our global network of NVIDIA GPUs with no extra setup.

### Models Included

Choose from a variety of popular models in our catalog including Llama-2, Whisper, and ResNet50.
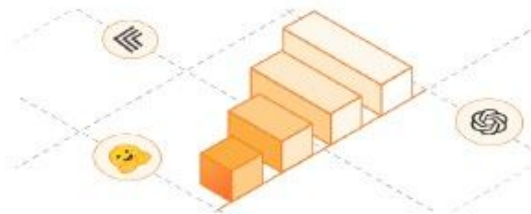
### Available everywhere

Run AI models from Workers, Pages, or anywhere via our REST API

CLOUDFLARE

# AI accessible to everyone



### Supercharge with Vectorize

Generate and store embeddings in a globally distributed vector database.



### AI Gateway

Improve reliability and scalability with caching, rate limiting, and analytics.



### Train with R2

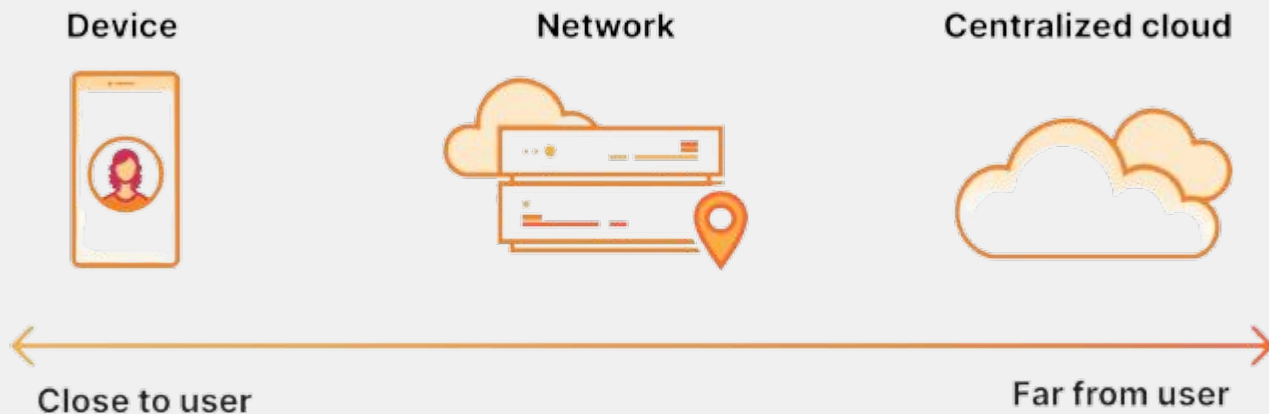Build multi-cloud training architectures with free egress.

# Announced this week:

Gemma 2b and 7b

# Why does AI inference belong on the edge?



Device     Network     Centralized cloud

Close to user     Far from user
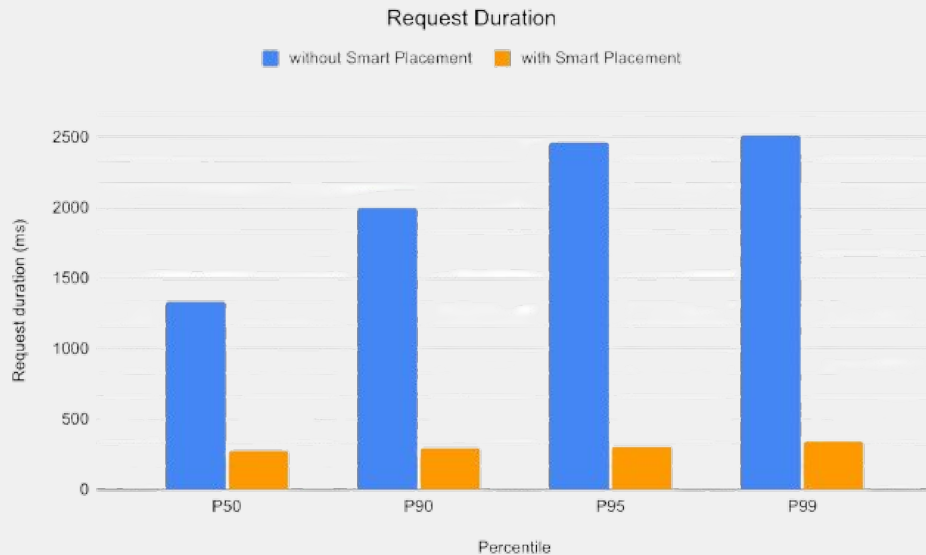
# Lowering inference latency

## Smart placement and GPU inference

- Avoid costly round trips to a centralized service

# Lowering inference latency

### Smart placement and GPU inference

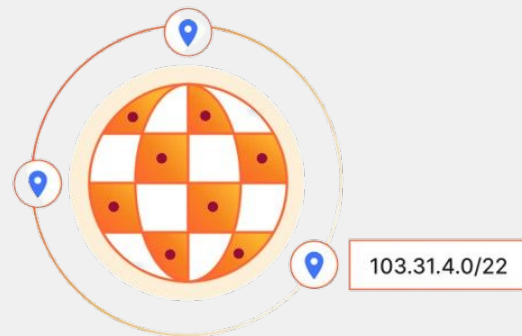- Avoid costly round trips to a centralized service

# Lowering inference latency

**Smart placement and GPU inference**



Request Duration

■ without Smart Placement   ■ with Smart Placement

# Helping products stay compliant

- GDPR
- EU AI Act
- Fine tuning with customer data



103.31.4.0/22

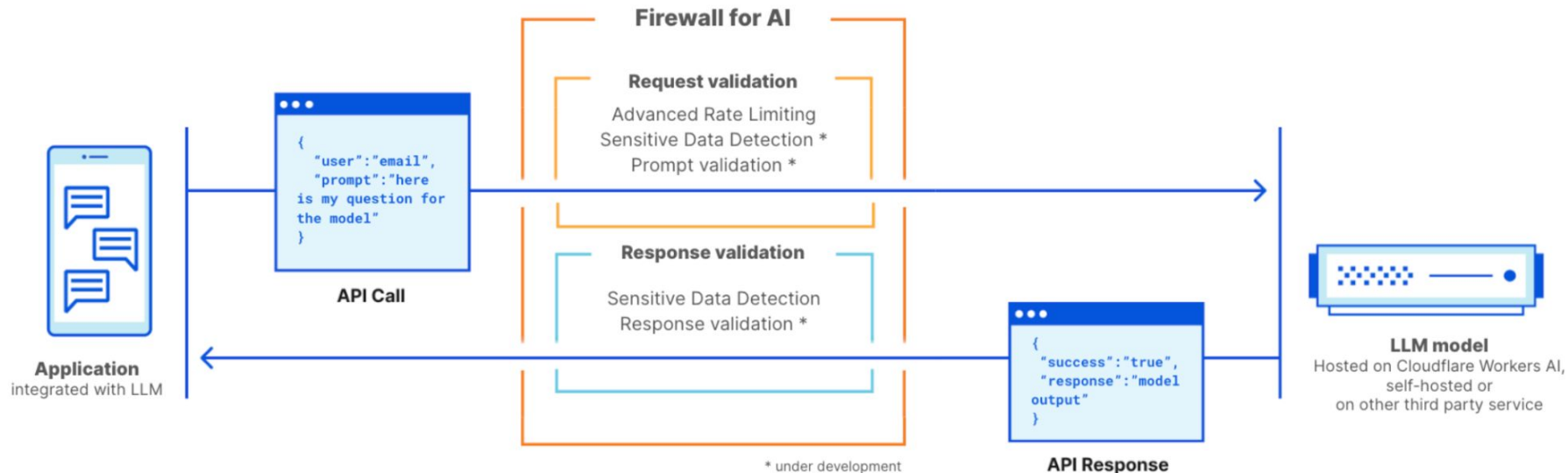CLOUDFLARE

# Enhanced protection with AI at the edge

**Firewall for AI**
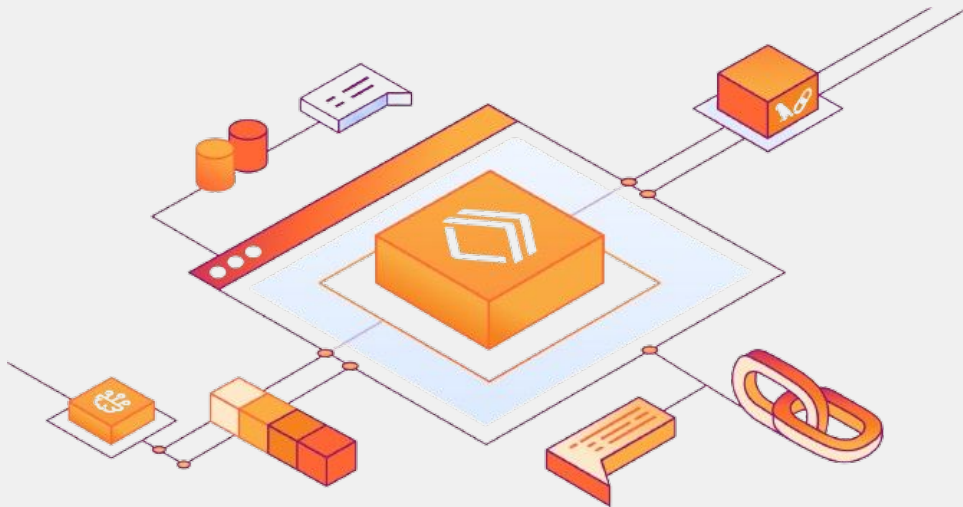
## LLM Security Concerns

- Prompt Injection
- Insecure Plugin Design
- Sensitive Information Disclosure
- Excessive Agency

# Enhanced protection with AI at the edge

# Not just AI, a whole developer platform

Build full stack, instantly scalable applications

**Storage Solutions**
- **Durable Objects**
- **D1 (SQLite)**
- **KV**
- **R2 (Zero egress object storage)**

**Bring your frontend framework**
- **Cloudflare Pages**

**Secure your applications**
- **API Gateway**
- **Bot Management**

**CLOUDFLARE**

# Connect with us!

**Logan Grasby**

**Developer Educator, AI at Cloudflare**

**X** **@LoganGrasby**

**in** **linkedin.com/in/logangrasby**

## At NVIDIA GTC

Booth 1634

## Online

meet-us.pages.dev