Google Cloud

Language and image recognition capabilities of AI systems have improved rapidly

# Timeline of images generated by artificial intelligence

**Our World in Data**

These people don't exist. All images were generated by artificial intelligence.

**2014**

Goodfellow et al. (2014) – Generative Adversarial Networks

**2015**

Radford, Metz, and Chintala (2015) – Unsupervised Representation Learning with Deep Convolutional GANs

**2016**

Liu and Tuzel (2016) – Coupled GANs

**2017**

Karras et al. (2017) – Progressive Growing of GANs for Improved Quality, Stability, and Variation
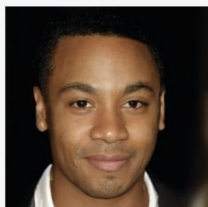
**2018**

Karras, Laine, and Aila (2018) – A Style-Based Generator Architecture for Generative Adversarial Networks

**2019**

Karras et al. (2019) – Analyzing and Improving the Image Quality of StyleGAN

**2020**

Ho, Jain, & Abbeel (2020) – Denoising Diffusion Probabilistic Models

**2021** Image generated with the prompt: *"a couple of people are sitting on a wood bench"*

Ramesh et al. (2021) – Zero-Shot Text-to-Image Generation (OpenAI's DALL-E 1)

**2022** Image generated with the prompt: *"A Pomeranian is sitting on the King's throne wearing a crown. Two tiger soldiers are standing next to the throne."*

Saharia et al. (2022) – Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding (Google's Imagen)

OurWorldinData.org – Research and data to make progress against the world's largest problems.    Licensed under CC-BY by the authors Charlie Giattino and Max Roser

## Nov 2023

Image that I generated using Stable Diffusion with the prompt:

*"A photo of Albert Einstein in a space suit watching a solar eclipse."*

https://stablediffusionweb.com/

Citation: The brief history of artificial intelligence: The world has changed fast – what might be next?

# The future of GenAI is even more **incredible**



Prompt: **A teddy bear running in New York City.**

This is an example from imagen.research.google/video



Prompt: **A glass bead falling into water with a huge splash. Sunset in the background.**

This is an example from imagen.research.google/video

Citation: https://imagen.research.google/video/

# Generative AI will Transform Live Service Games into Living Games

**Global Generative AI in Gaming Market**

Size, by Function, 2022-2032 (USD Million)

- Non-player characters
- Level generation
- Image enhancement
- Scenarios and stories
- Balancing in-game complexity

922, 1,137, 1,367, 1,770, 2,236, 2,690, 3,075, 3,792, 4,561, 5,624, 7,105

2022 2023 2024 2025 2026 2027 2028 2029 2030 2031 2032

The Market will Grow At the CAGR of: **23.3%** The forecasted market size for 2032 in USD: **$7,105M** market.us ONE STOP SHOP FOR THE REPORTS

GenAI is being used in **game development** use-cases first, and ultimately will be eclipsed by **new game experiences**

# **Classifications** of GenAI use cases in Games

*We will focus on inference during gameplay*

## Improving Productivity during Game Development

Use Generative AI to accelerate time-to-launch by creating content and simplifying development

- 2D & 3D assets (characters, props)
- Audio & video assets generation
- Code generation
- AI-based game testing

**Turnkey (VertexAI, Sagemaker) | Kubernetes (GKE)**

## Improving Player Experience during GamePlay

Use AI/ML & GenAI to adapt the gameplay and empower players to generate game content in realtime.

- Smart NPCs (bots)
- Dynamic in-game content
- Customized player experiences
- User-generated content
- Endless worlds

**Kubernetes (GKE)**

# User pain-points for GenAI in Games

## Platform

### Cost
At-scale cost efficiency to ensure financial feasibility for AAA games

### Latency
Low latency to ensure smooth gameplay & user experience.

### Platform Selection
Platform(s) with performance, & access to models without lock-in.

## AI Maturity

### LLM Unpredictability
Need a coherent, relevant, and contextually appropriate inference

### Avoiding Bias
Training & fine-tuning should not propagate biases & stereotypes

### Content Filtering
Content moderation is needed to ensure safe & inclusive gameplay

## Gameplay

### Creativity vs Boundaries
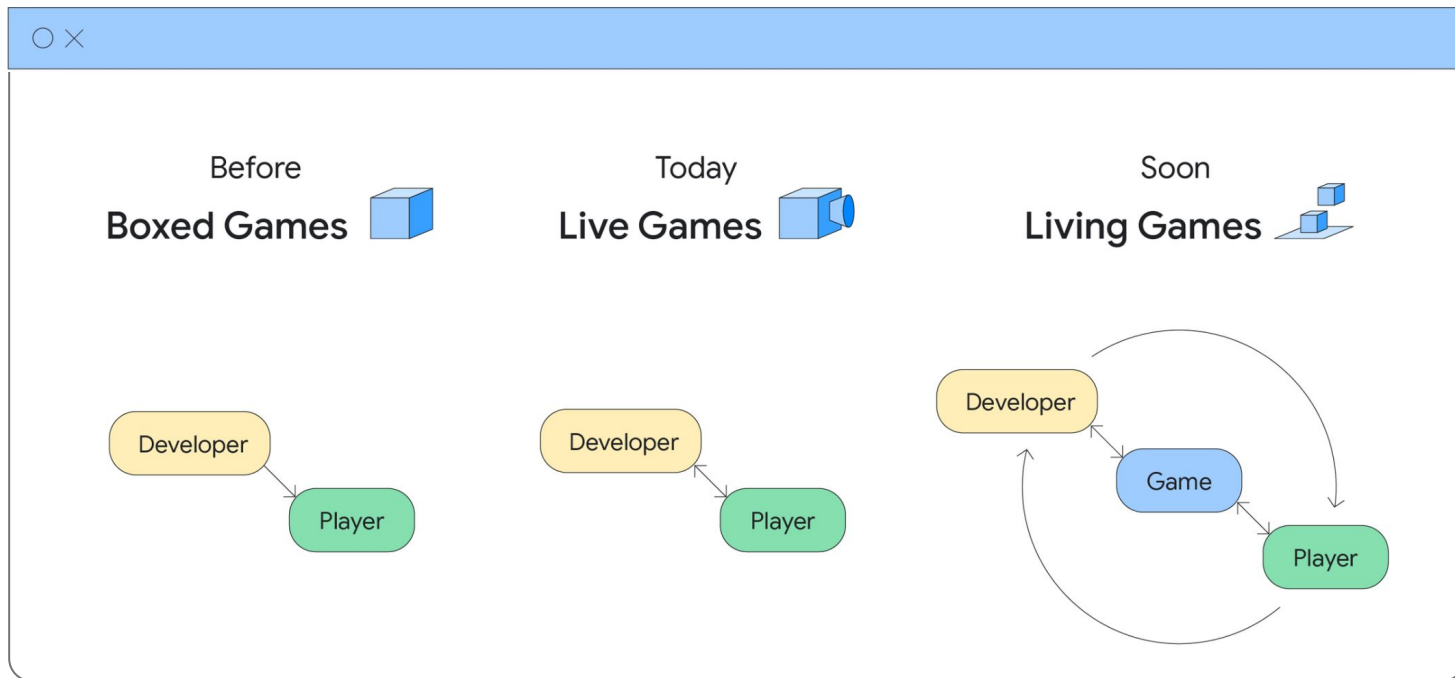Need to balance user generated content with game lore & structure

### GenAI model constraints
Some games need content for gameplay which LLMs filter out

### Procedural generation
Procedural generation with GenAI still requires human supervision

# Generative AI is evolving the games industry

# Kubernetes is a great computing solution for games

Solves **majority** of the IT operations problems:

- Scheduling, health-checking, deployment methods, autoscaling and rollbacks

- Centralized logging & monitoring

- Extended by a massive ecosystem of tooling

- Declarative paradigm - (say what you want instead of manipulating what you have)

- Primitives for isolation

**Challenge:** Kubernetes on it's own, does not understand how game servers work.

For game servers, we need:

- to maintain in-memory state

- to start and shut down game servers on demand

- to protect running servers from shutting down (even for upgrades!)

- to scale based on demand - location, # of players - rather than CPU

# Agones

**Agones** enables hosting, running, and scaling dedicated game servers **on Kubernetes**

**Agones** is an open source, batteries-included, multiplayer dedicated game server scaling and orchestration platform that can run anywhere Kubernetes can run.

Get all the benefits of the Kubernetes operations, but now for game servers as well:

- **Termination:** understand game matches
- **Scaling:** understand player load
- **Networking:** Multiple UDP/TCP ports per node
- **Hot-spares:** Tunable warm-up parameters
- **Open source:** No vendor lock-in

```
$
```

*Learn more at* **agones.dev**

High-level architecture of a live service game

# Integrating GenAI inference with gameservers

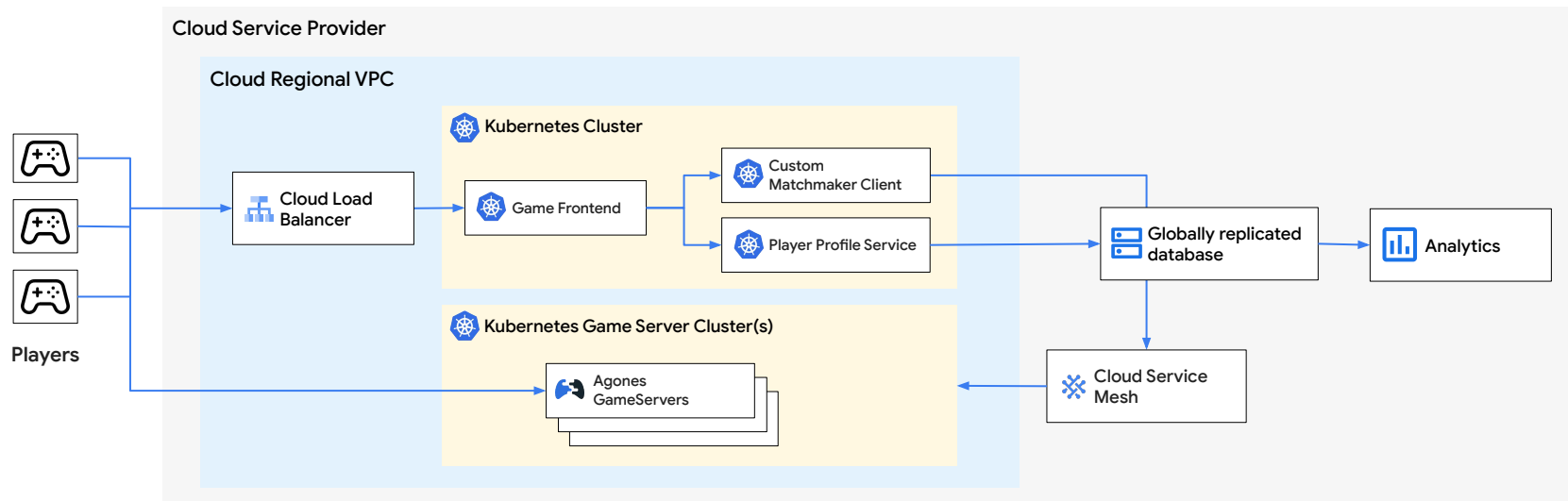**Turnkey Solution**
(E.g. VertexAI, Sagemaker, etc.)

Kubernetes Game Server Cluster(s)

Agones GameServers

Agones GameServers

GenAI Turnkey Solution API

**DIY Solution with Kubernetes**

## Dedicated GenAI Inference Kubernetes Nodes

Kubernetes Game Server Cluster(s)

Agones GameServers

Agones GameServers

Agones GameServers

Kubernetes GenAI Inference Cluster

GenAI Inference Servers

## GenAI inference Sidecar

Kubernetes Game Server Cluster(s)

Agones GameServers → GenAI Inference Server

Agones GameServers → GenAI Inference Server

# Discussion on integrating GenAI inference with gameservers

Advantages of using a Turnkey solution:
- Game development use-cases
- Improving time-to-value: POCs for realtime use-cases during gameplay
- Specific models only available through Turnkey APIs

DIY solution with Kubernetes for GenAI in games:
- Increasing number of openly available models that can run in containers
- Cost optimization at-scale: k8s can be more cost-effective than pay-per-use APIs for high usage scenarios (game launches)
- Dedicated inference k8s nodes are easy to set up with k8s features such as HPA, scheduling with taints/ tolerations, etc.
- GenAI sidecars may have a slight advantage in latency, but are costly (1:1)

| GenAI Inference Deployment Method | Latency | |
| --- | --- | --- |
| | Image Generation (Stable Diffusion) | Text Generation (Bloom) |
| Dedicated Inference Kubernetes Nodes | ~1s-1.3s | ~146-147 ms |
| Inference Sidecar | ~1s-1.3s | ~144-145 ms |

Today, inference latency overpowers any difference between different Kubernetes deployment methods. Dedicated inference k8s nodes provide the most versatility, ease of use, and flexibility.

# Using Kubernetes for GenAI in games

## Portability

**Write Once.
Run Everywhere.**

- Train and serve the same model(s) across clouds and on-premises
- Open standards prevent vendor lock-in



### Kubernetes

Industry standard compute orchestration platform available anywhere you need it

## Flexibility

**Choose the right framework(s) for the job**

- Meet the needs of multiple teams with their framework of choice
- Customize the platform to meet your structure and requirements

Spark

beam

DASK

RAY

RAPIDS

XGBoost

Vibrant ecosystem of frameworks from which to choose

## Scalability & Performance

**Fine tune performance and scale the platform**

- Hyper-optimize the architecture for peak performance
- Scale the platform to meet the needs of all of your ML workloads

## Cost & Efficiency

**Pay for what you need when you need it**

- Higher utilization of compute resources (CPUs, GPUs, TPUs) and cost savings with Spot
- Reduced operational costs for unified platform

Agones

**Run alongside game servers on Kubernetes.**

- Improve latency & performance by running AI inference alongside game servers using Agones on Kubernetes.
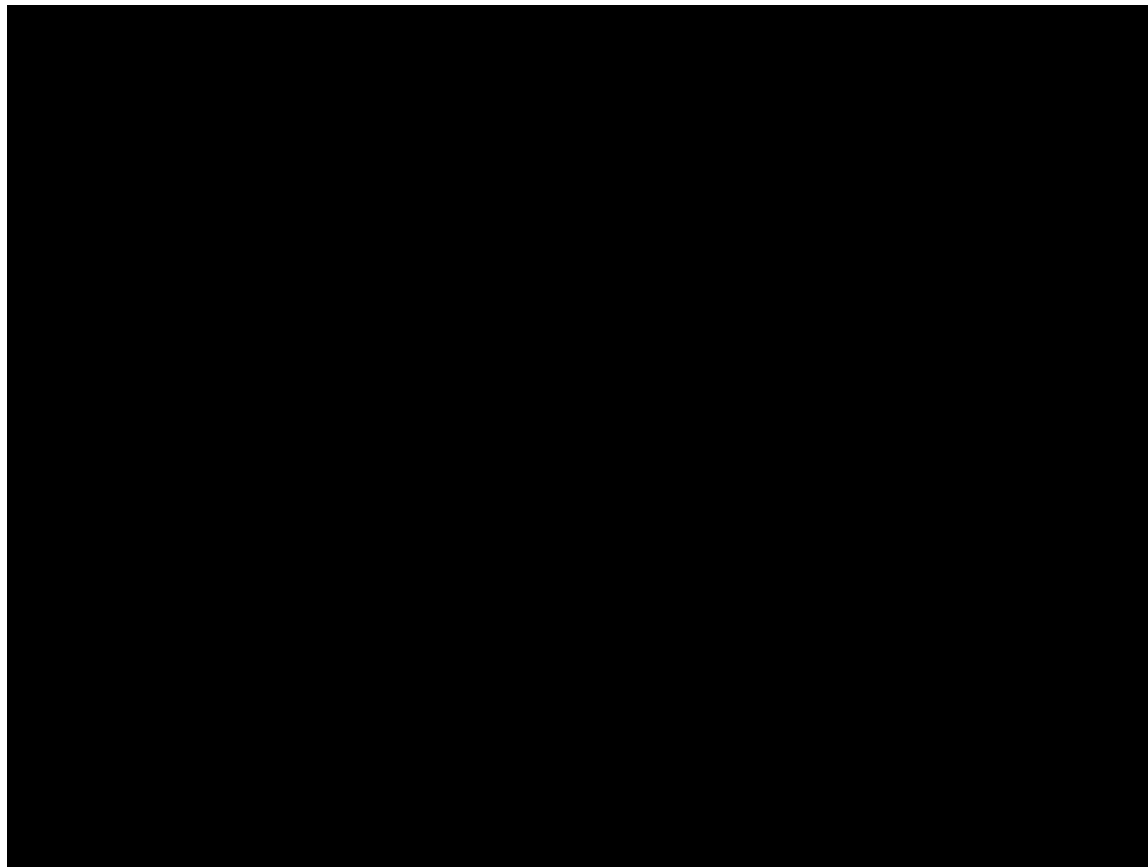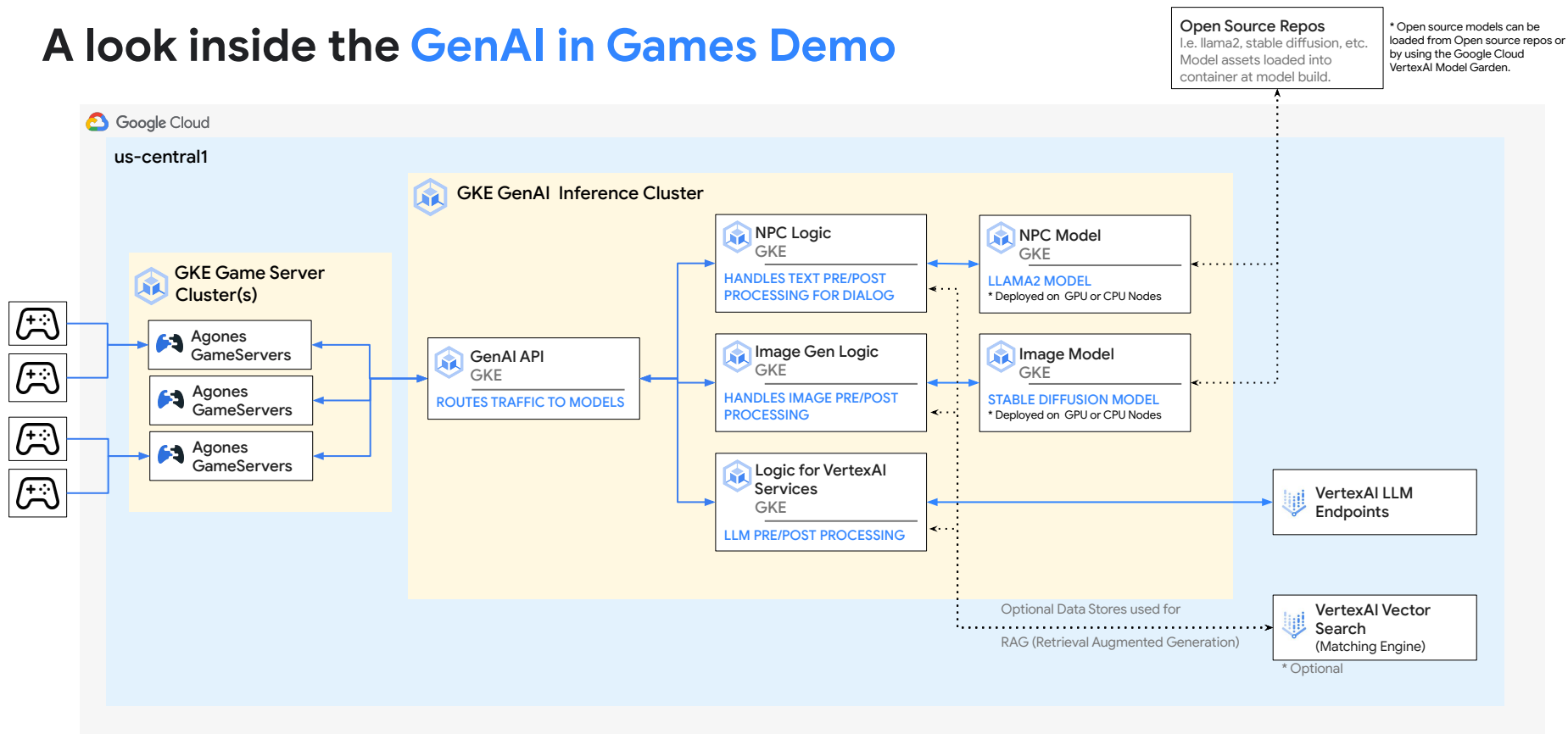- Reduce management overhead by using fully managed K8s (GKE Autopilot) for game servers

Demo

# A look inside the GenAI in Games Demo

**KubeCon | CloudNativeCon**
North America 2023

**Open Source Repos**
I.e. llama2, stable diffusion, etc.
Model assets loaded into
container at model build.

\* Open source models can be
loaded from Open source repos or
by using the Google Cloud
VertexAI Model Garden.

## Google Cloud

### us-central1

#### GKE GenAI Inference Cluster

**GKE Game Server Cluster(s)**

Agones GameServers

Agones GameServers

Agones GameServers

**GenAI API**
GKE
ROUTES TRAFFIC TO MODELS

**NPC Logic**
GKE
HANDLES TEXT PRE/POST
PROCESSING FOR DIALOG

**NPC Model**
GKE
LLAMA2 MODEL
\* Deployed on GPU or CPU Nodes

**Image Gen Logic**
GKE
HANDLES IMAGE PRE/POST
PROCESSING

**Image Model**
GKE
STABLE DIFFUSION MODEL
\* Deployed on GPU or CPU Nodes

**Logic for VertexAI Services**
GKE
LLM PRE/POST PROCESSING

**VertexAI LLM Endpoints**

Optional Data Stores used for

RAG (Retrieval Augmented Generation)

**VertexAI Vector Search**
(Matching Engine)
\* Optional

# Acknowledgements

| Google Kubernetes Engine (GKE) Team ||
|---|---|
| Technical leadership | Robert Bailey |
| GKE Games Engineering | Zach Loafman<br>Max Gong<br>Ivy Gooch |
| GKE AI Inference Team (benchmarking) | Kellen Swain<br>Abdullah Gharaibeh |
| GKE User Research | Anna Poznyakov |
| Leadership | Alex Bulankou<br>Iftach Ragoler<br>Maulin Patel<br>Jack Buser |

| Generative AI in Games Demo ||
|---|---|
| **Google Cloud team** | **Globant team** |
| Patrick Smith<br>Dan Zaratsian<br>Michael Bychkowski<br>Giovane Moura Jr<br>Sebastian Weigand<br>Mark Mandel | Laura Gonzalez<br>Cristian Giovagnoli<br>Diego Salatino<br>Rafael Martins<br>Arturo Salinas<br>Fede Vezzoso<br>Jesus Gonzalez<br>Matias Rodriguez<br>Maggie Avallone |

As you explore integrating Generative AI in your games, consider deploying your services on Kubernetes - matchmaking, game servers, and generative AI inference servers.

We would love to connect with you:

- **g2x**@google.com
- linkedin.com/in/sharmai
- twitter.com/IshantheSharma

Key links:

- Google Cloud for Games: goo.gle/cloudforgames
- GKE: cloud.google.com/kubernetes-engine
- Agones: agones.dev



**Please scan the QR Code above to leave feedback on this session**

# Thank you

Google Cloud