Google Cloud

## サーバーレスを活用した、 立ち上げ期から、拡大期にわたる 「アーキテクチャの変遷」

株式会社 カウシェ Platform Team / Architect 伊藤 雄貴



日用品

## #シェア買いで

# みんなでお得に ショッピング

Application Architecture	01
Access Control / Network	02
Reliability / Observability	03
AI/ML	04

## **Application Architecture**

01



#### **Cloud Run**

Cloud Run is a managed compute platform that enables you to run containers that are invocable via requests or events.

Cloud Run is serverless: it abstracts away all infrastructure management...



- Everything runs on Cloud Run
- Everything runs as an API





e.g.) VS. Cloud Functions Trigger





## Everything is Managed as API Definitions

## Reuse same implementation logic as APIs

## **V** Use same Monitoring environments

## Everything is Managed as <u>API Definitions</u>

## **W** Reuse same implementation logic as APIs

## **V** Use same Monitoring environments

gRPC is a modern open source high performance Remote Procedure Call (RPC) framework that can run in any environment.





#### Architecture: 2020 ~



#### Architecture: 2021 ~



#### Architecture: 2022 ~



#### Architecture: 2022 ~



#### Offloading Cross-Cutting Concerns to the API Gateway

## Authentication / Authorization

## **V** Transcoding

## TLS / Domain / CDN / IP ...

#### Architecture: 2022 ~





Envoy is an L7 proxy and communication bus designed for large modern service oriented architectures. The project was born out of the belief that:

The network should be transparent to applications. When network and application problems do occur it should be easy to determine the source of the problem.





Extensibility with WebAssembly

V Dynamic Configuration

Widely used in the Cloud Native World



Access Control / Network

O2

#### Single Service









## **V** Restricting Ingress





## **V** Restricting Ingress









#### **Restricting Ingress**







### internal-and-cloud-load-balancing



Google Cloud






#### **Ingress Settings**

```
apiVersion: serving.knative.dev/v1
kind: Service
metadata:
  annotations:
    run.googleapis.com/ingress: internal
spec:
```

#### **Ingress Settings**

```
apiVersion: serving.knative.dev/v1
kind: Service
metadata:
  annotations:
    run.googleapis.com/ingress: internal
spec:
```

#### **Access Control**

Ingress settings and IAM authentication methods are two ways of managing access to a service. They are independent of each other. For a layered approach to managing access, use both.



#### Network

- Serverless VPC Access Connector
- Shared VPC Network

#### Network

### Serverless VPC Access Connector

Shared VPC Network



#### What does "internal" mean ...?

Setting	Description
Internal	Most restrictive. Allows requests from the following sources:
	<ul> <li>Internal Application Load Balancer, including requests from Shared VPC networks when routed through the internal Application Load Balancer</li> </ul>
	<ul> <li>Resources allowed by any VPC Service Controls perimeter that contains your Cloud Run service</li> </ul>
	<ul> <li>VPC networks that are in the same project or VPC Service Controls perimeter as your Cloud Run service</li> </ul>
	<ul> <li>Shared VPC ingress (Preview): The Shared VPC network that your revision is configured to send traffic to. For information about when Shared VPC traffic is recognized as "internal", see <u>Special considerations for Shared VPC</u>.</li> </ul>
	<ul> <li>The following Google Cloud products, if they are in the same project or VPC Service Controls perimeter as your Cloud Run service:</li> </ul>
	Cloud Scheduler (Preview)
	Cloud Tasks
	Eventarc
	• Pub/Sub
	Workflows
	• BigQuery
	Requests from these sources stay within the Google network, even if they access your service at the <b>run.app</b> URL. Requests from other sources, including the internet, cannot reach your service at the <b>run.app</b> URL or custom domains.
	There is no support for multi-tenancy, that is, multiple trust domains within the same project.

#### https://cloud.google.com/run/docs/securing/ingress#settings

#### By default, requests from other Cloud Run Services are not treated as "internal"



## For requests from other Cloud Run services in the same project, connect the service ... to a VPC network and route all egress through the connector



Serverless VPC Access makes it possible for you to connect directly to your Virtual Private Cloud (VPC) network from serverless environments such as Cloud Run ...





#### Network

### Serverless VPC Access Connector

Shared VPC Network

#### Network

- Serverless VPC Access Connector
- Shared VPC Network

#### Single VPC Network





Shared VPC allows an organization to connect resources from multiple projects to a common Virtual Private Cloud (VPC) network, so that they can communicate with each other securely and efficiently using internal IPs from that network

#### Delegating network responsibilities to administrators

#### Centralized control over network resources



# 03

## **Reliability / Observability**



#### Observability

...

our definition of "observability" for software systems is a measure of how well you can understand and explain any state your system can get into, no matter how novel or bizarre.

If you can understand any bizarre or novel state without needing to ship new code, you have observability.

# O'REILLY' Observability Engineering Achieving Production Excellence

**Reliability / Observability** 



#### to understand system states on Day 1.

#### Approaches

- Logs
- Traces
- Metrics
- •SLI / SLO

Approaches

- Logs
- Traces
- Metrics
- •SLI / SLO

## Request Logs

## Container Logs

#### **Request Logs**

```
24 ms 🚍
                                                    504 B
\sim 0
       2021-02-21 14:39:05.120 JST
                                     GET
                                           200
                                     Kauche X.X.X https://...
 ▼ {
          Hide log summary
                               Copy to clipboard
    insertId: "xxx"
  requestMethod: "GET"
     requestUrl: "https://..."
     requestSize: "1435"
     status: 200
     responseSize: "504"
     userAgent: "Kauche/X.X.X"
     remoteIp: "xxx.xxx.xxx.xxx"
     serverIp: "xxx.xxx.xxx.xxx"
     latency: "0.024070354s"
     protocol: "HTTP/1.1"
```

#### **Container (Application) Logs**

<b>&gt; *</b>	2021-02-21 14:39:05.098 JST 🔁 "g	Jrpc request"					
	Hide log summary = Expand nested fields	Copy to clipboard					
▼ {							
i	insertId: "xxx"						
🔻 j	sonPayload: {						
	logger: "grpc.request_logger"						
	<pre>method: "/customer.v1.CustomerService/GetXXX"</pre>						
	message: "grpc request"						
	level: "info"						
	timestamp: 1613885945098.689						
}	•						
▼ r	esource: {						
	type: "cloud_run_revision"	Structured Log					
•	labels: {5}						
}	•						

## **Request Logs**

#### +

#### **Container Logs**

**Google** Cloud

~	2021-02-21 14:39:07.846	5 JST GET	200 709 B 27	′ms \Xi Kauc	he X.X.X https://	Д
▼ { ir	nsertId: "xxx"	🖃 Hide log summ	mary = Expand neste	ed fields	Copy to clipboard	G⇒ Copy link
► ht	ttpRequest: {10}					
▶ re	<pre>resource: {2}</pre>					
ti	timestamp: "2021-02-21T05:39:07.846088Z"					
se	severity: "INFO"					
1a	<pre>labels: {1}</pre>					
10	logName: ""					
tı	trace: "projects//traces/xxx"					
re	receiveTimestamp: "2021-02-21T05:39:08.016290137Z" Request Logs					iest Logs
}						_
>  🏵	2021-02-21 14:39:07.826	JST 🔁	"http request"			
>  🛞	2021-02-21 14:39:07.82	1 JST \Xi	"grpc request"			
>  🛞	2021-02-21 14:39:07.82	1 JST 🗧	"authenticated, uid: x	(XX"	Contr	sinor Logo
>  🚯	2021-02-21 14:39:07.844	4 JST \Xi	"finished"		Conta	amer Logs
	SHOW MORE SHOW LESS SHOW ALL					

~	2021-02-21 14:39:07.846 JST	GET 200 709 B 27 ms 🚖 Kauche X.X.X https://				
▼ { i	msertId: "xxx"	e log summary = Expand nested fields □ Copy to clipboard = Copy link				
▶ h <sup>-</sup>	ttpRequest: {10}					
r	esource: {2}					
t	imestamp: "2021-02-21T05:39:07.	846088Z"				
S	everity: "INFO"					
1:	abels: {1}					
1	ogName: ""	Trace ID				
<pre>trace: "projects//traces/xxx"</pre>						
receiveTimestamp: "2021-02-21T05:39:08.016290137Z"						
}						
>  🏵	2021-02-21 14:39:07.820 JST	<pre>"http request"</pre>				
>  🏵	2021-02-21 14:39:07.821 JST	"grpc request"				
>  🏶	2021-02-21 14:39:07.821 JST	<pre>"authenticated, uid: xxx"</pre>				
>  🏵	2021-02-21 14:39:07.844 JST	= "finished"				
SHOW	MORE SHOW LESS SHO	V ALL				



#### Approaches

- Logs
- Traces
- Metrics
- •SLI / SLO




# OpenTelemetry

OpenTelemetry is a collection of tools, APIs, and SDKs. Use it to instrument, generate, collect, and export telemetry data (metrics, logs, and traces) to help you analyze your software's performance and behavior.







#### Approaches

- Logs
- Traces
- Metrics
- •SLI / SLO

#### Log Based Metrics



# **Log Based Metrics**

Counter

user/CloudRunCustomerAPIErrorLogs



72 B

344 B httpRequest.requestUrl=~"https://..." httpRequest.status!="200" httpRequest.status!="302" httpRequest.status!="404" resource.type="cloud\_run\_revision" resource.labels.project\_id="xxx"

## OpenTelemetry



#### Approaches

- Logs
- Traces
- Metrics
- SLI / SLO

# **Cloud Monitoring**



Google Cloud

# PromQL



04 AI/ML

#### **Architecture - Recommendation**



#### **Architecture - Recommendation**







```
resource "google_cloud_scheduler_job" "trigger-ai-pipeline" {
   schedule = "12 0 * * *"
   time_zone = "Asia/Tokyo"
   http_target {
     http_method = "POST"
     uri = "https://.../execute_pipeline"
     body = base64encode("{\"pipeline_spec_path\":\"<gcs bucket name>/path/to/pipeline_spec.json\"}")
     oidc_token {
        service_account_email = "..."
        audience = "..."
     }
   }
}
```





110

111	<pre>// Creates a PipelineJob. A PipelineJob will run immediately when created.</pre>
112 🗸	<pre>rpc CreatePipelineJob(CreatePipelineJobRequest) returns (PipelineJob) {</pre>
113	<pre>option (google.api.http) = {</pre>
114	<pre>post: "/v1/{parent=projects/*/locations/*}/pipelineJobs"</pre>
115	<pre>body: "pipeline_job"</pre>
116	};
117	<pre>option (google.api.method_signature) =</pre>
118	<pre>"parent,pipeline_job,pipeline_job_id";</pre>
119	}
120	

https://github.com/googleapis/googleapis/blob/ef2e2ea532248d6dc40a56bc6c95cea858ba31b6/google/cloud/aiplatform/v1/pipeline\_service.proto#L111C1-L119



Application Architecture	
Access Control / Network	02
Reliability / Observability	03
AI/ML	04





# Thank you.

Google Cloud