Developing a Process in Architecting Microservice Infrastructure with Docker, Kubernetes, and Istio

Yujing Wang

Department of Mechanical and Mechatronics Engineering
University of Waterloo
Waterloo, ON, Canada
yj9wang@edu.uwaterloo.ca

Darrel Ma

Cheriton School of Computer Science
University of Waterloo

Waterloo, ON, Canada
darrel.ma@edu.uwaterloo.ca

Abstract—As an application usage grows, its owner scales up vertically by replacing old machines with more powerful ones. This methodology is expensive and leads to resource waste. In response to the business needs, internet giants have developed the microservice architecture, which lets developers divide up their application into smaller units that can be hosted on multiple machines, thus enabling horizontal scale up. We propose a triphasic incremental process to transform a traditional application into a microservice application that guarantees stability during the operation. Then we demonstrated such methodology in a prototype microservice application based on an existing monolithic application. First, the developer splits a monolithic application into atomic services and aggregated services. Second, these services are packaged, containerized, and then deployed on Kubernetes. During this stage, Istio is deployed on the Kubernetes cluster to establish pod level communications, delegate traffic flows and filter requests, and enable the autoscaler. Other external add-ons, such as database connections, are defined in service entry. In the last stage, we developed an algorithm guideline to minimize inter-service calls by compiling all needed calls into a list and perform one finalized call. Although it increases memory usage, it avoided the wait time incurred during interservice calls. We then investigated managing configurations using config maps, recommended a pipeline being developed to perform automatic

Index Terms—Microservice, Kubernetes, Docker, Istio, Algorithm

I. INTRODUCTION

When an application usage grows, its owner scales it up to handle the increased traffic. Traditionally, companies scale up vertically by replacing current servers with more powerful ones [1]. This practice requires looking for more voluminous hardware resources at the time of needs, doesn't account for sudden traffic, and requires a major upfront capital investment. When the application does not use machines at their full capacity, resources are wasted. To compensate for the rigidity of vertical scaling, internet giants are promoting the microservice architecture that sees applications decoupled into logical units and then sliced into microservices. These services are packaged in Docker images and deployed on Kubernetes, which handles the hosting, scaling, and monitoring. Then aggregated services are used to facilitate inter-service communications [2].

Kubernetes is an open source orchestration system offered by Google for managing containerized services and facilitating declarative configurations and automation [2]. Although Pivotal offers cloud solution for Java Spring applications [?], adapting the microservice to a codebase specific platform would create a dependency on Java making it hard to move away should the team rewrites the program in other languages. In addition, adapting to codebase specific platform would be a bad example for developers that seeks a consistent guideline across all codebase. There is an urgent need for a replicable, scalable, and easy to follow process to transform monolithic applications to microservices. This process must be codebase agnostic host agnostic. This paper focuses on developing a methodology to transform traditional Java Spring monolithic backend applications to a network of microservice applications containerized with docker and hosted on Kubernetes. The report will dive deep into the findings during our research progress and discuss concern arose in the process of development.

A. Scaling

To elaborate on the previously introduced technical terms, vertical scaling and horizontal scaling are not exclusive ideas. Indeed, they are expansions in two different dimensions [1]. Horizontal scaling requires adding more hardware resources; while, vertical scaling requires finding hardware resources more powerful than current ones. Expanding on both directions yields maximum capacity. However, to attain the most cost-effective way requires finding a balance in the two directions. As Fig. 1 shows, the cost of using less than 15 processors in the server is cheaper than purchasing multiple servers. Having a single server would benefit from using less housing space. Horizontal scale-up adds cost to the storage space and machine maintenance; meanwhile, the cost of adding processors to a single machine grows exponentially.

B. Hosting Strategy

The aforementioned "On-premise" is one of the hosting strategies that clients host applications on their own infrastructure, on-premise [5]. It gives clients total control of their data but requires them to maintain network and security. "On cloud" is where the clients host their application on one of the providers' cloud networks such as Google Kubernetes Engine. Some cloud provider provides private cloud (IBM Bluemix),

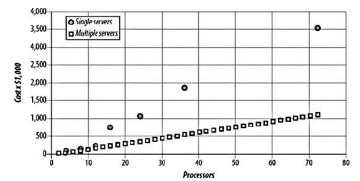


Fig. 1: Cost for Vertical vs Horizontal Scaling [4]

where the client's data is being put separately on machines on the cloud.

No upfront cost and flexible scaling make "On cloud" the most appealing option among startups that seek rapid growth and are willing to outsource their network security management to big firms. Compared to on-premise, cloud reduces the initial set up time and removes the need to organize a team of network engineers and security engineers. However, as usage grows, the cost of the cloud also increases exponentially. Snapchat, an instant picture chat startup that uses Google's Cloud Engine, pays Google 2 billion and Amazon 1 billion as of 2018 in cloud contracts [6].

II. PROBLEM DEFINITION

While we see a surge in demand for a microservices strategy, there lacks a concrete replicable yet scalable methodology to transform monolithic Java application to microservices application. The existing methodology does not offer a clear definition of minimum effort service discovery. Ones offered by Google's guidelines are vague for adapting Java Spring application. We seek to apply the emerging microservice architecture into production for efficient resource management and satisfy increasing client demands for a microservice strategy. We have decided to put the emphasis on platform agnosticism and codebase versatility. There is yet to be developed a portable, scalable, and continually deployable methodology to efficiently transform existing monolithic application to a network of microservices. Specifically, three questions need to be addressed:

- 1) How to perform services discovery and inter-service communication?
- 2) What data strategy would guarantee an efficient read and write capability? How to make such a data strategy compatible with existing data infrastructure?
- 3) How do we ensure any problems that arose in the operation of the microservice application are picked up in a speedy manner and sent quickly to the right team for a bug fix?

The transformed microservices must be able to perform functions that the original service does. The methodology should make the application more portable, more continually deployable, and more scalable than its original service. By looking at 12 factors app [7] as criteria, we will place emphasis on

- 1) service discovery
- 2) data strategy
- 3) canary development

III. Breaking up the Monolith

The first step is to break up an existing application into microservices. In an ideal microservices ecosystem, each service represents what a business unit does [8]. It exposes an API that the developer uses to communicate with other services. Each microservice should be stateless, meaning an enclosed lifecycle independent from the state of other services. Finally, each microservice enforces a team to program and maintain such a service independently. For demonstration purposes, a mock social media application "Userapp" is developed. It lets the user set up an account, add friends, and make posts. The structure of the java app is shown in Snippet 1.

Snippet 1: Monolithic Java Application Structure

```
userapp
 java
    src
      Application.java
      controller
         Controller java
         RestController.java
      service
         Userapp.java
         AccountService.java
         FriendService.iava
      Domain
         Account.java
         Friend.java
         FriendType.java
         Message.java
         Post.java
      repositories
         AccountRepositories.java
         FriendRepositories.java
         ControllerTest.java
```

In which, the service controls the business logic, the user-app.java is where interactions with the submodules are defined as shown in the 3 submodules shown in Table I: When the

TABLE I: Microservices Broken down from Monolith App

l	Userapp.Java	AccountService.java	FriendService.java	
	Userapp:	Userapp:	Userapp:	
	getAccount()	getAccount()	getAccount()	
	updateAccount()	updateAccount()	updateAccount()	
	post()	post()	post()	
	makeFriend()	makeFriend()	makeFriend()	

user makes a post, the post function in 'UserApp' will be called to generate the post, which then calls 'AccountService' that register the post onto one's own wall. Subsequently, it calls getAccount() to get a list of the users friends from FriendService to notify them that a post has been made by the said user. The repository module handle creates, read, update and delete (CRUD) operations.

Three classes in service categories are each responsible for one domain of functions and can be easily be separated into three services. Once the application has been broken down to this level, it cannot be further broken down because functions in each class are tightly coupled and inter-reliant. FriendService and AccountService cannot be further broken down because their functions belong to one set of logic tied to one database. Hence this kind of service is called atomic services. Meanwhile, Userapp does not have any resources attached but works by accepting and sending requests from other services. This kind of service is called an aggregated service, as Fig. 2 shows. In this specific case, userapp service interacts directly with the front end through API and hence is also an example of Backend for Frontend Service (BFFs).

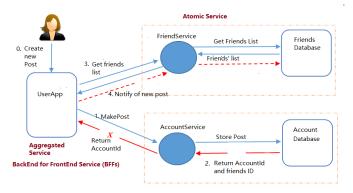


Fig. 2: Architecture View of "Userapp"

IV. A TRIAPHASIC INCREMENTAL APPROACH

Service built from the early days of Object-Oriented Programming (OOD) may not have such a clear cut of separation of concerns unlike 'Userapp'. For example, older applications may have placed friends and accounts in one database. As a result, we realized that the clear-cut idealistic approach is rather impractical, decides to take a progressive transformation in Fig. 3

Separation

- 1. Domain Separation
- Logical Separation
- Application Separation

Transition

- Package into Image
- Package into Deployment Create new Schema & data strategy
- Partial roll over
- Canary Testing

Completion

- Complete Roll over
- Complete transition microservice
- Remove resource dedicated to old application

Fig. 3: Triphasic Approach to Transform app. to Microservice

A. Seperation

Using this approach, the development team will first look at the 'domain' category in Table I. Then identify each domain as primary or helper object, which is shown in Snippet 2. A primary object has a set of functions to perform business logic. A helper object does not have its own functions. For example,

'Post' is a helper object that has no function of its own but is used by other functions to fulfill their business logic.

Snippet 2: Helper Object Used by Both AccountService and FriendService

account.makePost(post); friends.notify(posts);

Next, A hard logical separation must be placed between these primary objects. In a monolithic application, the account may directly call a friend's database. The developer must remove these cross-domain calls first. Ideally, services on the same hierarchy will never call each other. To bridge the communication, an aggregated service, in this case, Userapp, is used to bridge communication between the services. Because of this separation, some calls that could have been made directly now require exiting the originating service, finding the destination, and then entering the new services. As a result, communication now requires more calls and longer trips. This is a tradeoff that programmers need to make. Generally, developers shall aim for complete separation to avoid having to go back and reconsider logical separation. Once these logical units are packaged into microservices, the perks of microservices such as autonomous scaling, rapid development, and efficient service to service communication can compensate some efficiency lost in inter-service communication. Most importantly, separating the logical units also allows functionspecific development teams to develop independently without reliance on other services. Bugs in the service can be quickly allocated and fixed.

B. Transition

At phase two, "transition", the development team needs to work with the operation team to package the deployment, create new schemas for new atomic services. Each service corresponds to one schema, one database. The database is usually hosted separately from service because services are stateless, but databases are stateful. Next, the developer points the new microservices to the database and labels it as a test environment, running alongside the production environment. As the system matures, the group can decide whether to move a certain percentage of traffic into the new environment. In a progressive approach, this step will progress slowly to ensure minimum disturbance by having the existing monolithic application running alongside the new microservice app. Ultimately, this is a tradeoff between stability and resource. In the last phrase, the team will move everything to the new application and remove reference to the old. By now, the microservice application is mature, and the operation team will take over the app.

Breaking up the Monolith is the first step in transforming a monolith application into a microservice. A progressive approach prevents having to go back to re-engineer the service. Small changes are made incrementally towards building a big system. Developers should keep in mind the tradeoff they make when separating logical units and have reasonable logic before making each decision. The lower level the blocks split, the higher the cost to each call across the individual unit, but also, the easier it is to convert into microservice and scale-up. For large scale applications, it is always better to perform complete separation and have the infrastructure to handle the scaling for better performance. Relate back to the criteria and constraints, properly laying out the foundations for microservice allows for a portable and scalable and continually deployable methodology.

V. BUILDING DEPLOYMENTS

Packaging the sliced application from phrase one into deployments is a necessary step in phrase two. At this stage, both Docker and Kubernetes will be used to create deployment [2]. The deployer will only need the packaged image to complete this task. Ad-hoc modifications such as changing the port name or hard coding a destination in source code to aid deployment are strictly prohibited. Otherwise, changes made by one developer would result in the methodology being inadaptable to other applications.

A. Docker Deployment

Docker lets the developer run applications in any environment that has a Docker engine without having to worry about dependencies and the OS environment [9]. The user writes a Dockerfile, as shown in Snippet 4, to build the Docker image. The developer specifies the source code and operating system in the "FROM" line separated by ":". Then copies the original executable file using the "COPY" command into the image. Other relevant commands can be executed with "CMD". Lastly, expose a port for API communication with {application_port}:{Image_port} command found at the last line of Snippet 3.

Snippet 3: Dockerfile to build a Java Spring app. Image

```
FROM openjdk:8-jre-alpine
COPY userapp-0.0.1-SNAPSHOT.jar / app.jar
# run application with this command line
CMD ["/usr/bin/java", "-jar", "-Dspring.profiles.active=container", "/app.war"]
EXPOSE 8080:8080
```

Then run the Dockerfile with Snippet 4 in the terminal to build the docker image. In this statement, the author supplies an image tag (userapp) and version (latest). The tag identifies an image registry where the image will be stored. If no version has been supplied, the "latest" version will be the default.

Snippet 4: Command to Build a Docker Image with a Tag

```
$ docker build -t userapp:latest .
```

B. Kubernetes Deployment

Kubernetes creates deployments from Docker images, establish service communication with envoy gateways that looks up services' IP address [10]. When a service's usage increases to a preset threshold, Kubernetes Autoscaler replicates the service to fulfill the increased load. Kubernetes follows a master-workers architecture Fig. 4. The master node accepts commands, controls cluster, schedules pods, and store configurations using kube-apiserver, kube-controller, kube-scheduler,

and etcd. Then it uses Kubelet to perform the action on a node level. When the user executes Snippet 5 using kubectl.

Snippet 5: Command to Create Kubernetes Deployment

```
$ kubectl create —f deployment —s 8080
```

Kubernetes' API server receives the request from keyword 'kubectl' and delegates the controller manager to create a new deployment from the 'deployment' file. After creation, the scheduler takes over the newly created pod on an available thread and expose it to port '8080' [2]. The most fundamental building block of Kubernetes is a pod. Similar to the idea of the container from docker, pods are containerized units in Kubernetes with a specific IP address assigned by Kubernetes [2]. A pod can be created, cloned, terminated, and destroyed by the node. As shown in Figure 6, a pod may containerize an app, a volume, or both.

Kubernetes Architecture

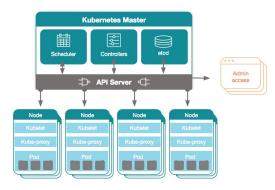


Fig. 4: Arcitecture of Kubernetes

To customize a kubernetes deployment, the developer writes an ymal configuration according to Snippet 6 [2]:

Snippet 6: Kubernetes Deployment

```
controllers/userapp.yaml
apiVersion: apps/v1
kind: Deployment
metadata:
  name: userapp-deployment
  labels:
   app: userapp
spec:
  replicas: 3
  selector:
    matchLabels:
      app: userapp
  template:
    metadata:
      labels:
        app: userapp
    spec:
      containers:
         name: nginx
         image: userapp:latest
         ports:
          - containerPort: 8080
```

The first line defines the Kubernetes API version. 'Kind' in the second line specifies the type of service this configuration contains. The "kind" field expects the type of Kubernetes object, the scheduler will expect. The most used "kinds" are pod, deployment, service, and PV. As mentioned before, 'pod' is a group of containers, including storage units that made up the most fundamental building block of Kubernetes [2]. When a pod is deployed, it only exists as a single instance and can be communicated through the kubectl API. It is usually used to host testing tools inside the Kubernetes container. For example, when the developer needs to test autoscaling, it becomes unrealistic to send 5000 requests in one second manually. A pod with 'busybox' image is created. Busy boy would continuously spam requests to the service at a specified interval with specified volume.

Very similar to the pod, "deployment" is another Kubernetes' kind' [2]. When creating a deployment, Kubernetes instantiates one pod with one replica set. If more traffic is being routed to this service, deployment can replicate its pods according to its replica set. Deployment lets the user control all instances created from one deployment at once. That way, an ordinary Kubernetes developer would not need to manually go through countless pods to apply the same change to one application. Detailed workflow of deployment-> replicaset-> pods are shown in Fig. 5.

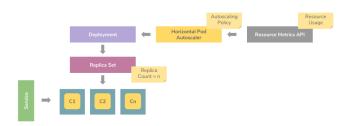


Fig. 5: Kubernetes Deployment Flow

Kubernetes offers more fine-grained resource control capability on deployment level than on the pod level. Take checking the status of pods as an example. Simply run "get pod" command in Snippet 7, the developer will find three pods belong to the same deployment, yet each pod has a different id.

Snippet 7: Get Pods with Label

```
$ kubectl get pods — show—labels
NAME READY STATUS RESTARTS AGE
nginx—deployment—75675f5897—7ci7o 1/1 Running 0 18s
nginx—deployment—75675f5897—kzszj 1/1 Running 0 18s
nginx—deployment—75675f5897—qqcnn 1/1 Running 0 18s
```

To update the deployment, the developer executes Snippet 8 as many times as there are pods.

Snippet 8: Manually Updating Pods One by One

```
kubectl delete pod nginx—deployment—{deployment_id}—{pod_id}
kubectl create pod {name_pod} —f {name_newapp}
```

However, the developer can work directly with deployment. First, execute Snippet 9 to check deployments and find that there are three pods associated with the nginx-deployment.

Snippet 9: Get Deployments

```
$ kubectl get deployments
NAME DESIRED CURRENT UP—TO—DATE AVAILABLE AGE
nginx—deployment 3 3 3 3 18s
```

The developer runs Snippet 10 to perform an update:

Snippet 10: Edit Deployments with New Configurations

```
$ kubectl edit deployment.v1.apps/nginx—deployment deployment.apps/nginx—deployment edited $
```

In order to see what has happened, the user runs the describe command Snippet 11.

Snippet 11: Describe Deployments to see Events and Logs

```
$ kubectl describe deployments
Name: nginx-deployment
Namespace: default
CreationTimestamp: Thu, 30 Nov 2017 10:56:25 +0000
Labels: app=nginx
Annotations: deployment.kubernetes.io/revision=2
Selector: app=nginx
Replicas: 3 desired | 3 updated | 3 total | 3 available | 0 unavailable
StrategyType: RollingUpdate
MinReadySeconds: 0
RollingUpdateStrategy: 25% max unavailable, 25% max surge
Pod Template:
  Labels: app=nginx
 Containers:
    Image: nginx:1.9.1
    Port: 80/TCP
    Environment: <none>
    Mounts: <none>
  Volumes: <none>
Conditions:
  Type Status Reason
  Available True MinimumReplicasAvailable
 Progressing True NewReplicaSetAvailable
OldReplicaSets: < none >
NewReplicaSet: nginx-deployment-1564180365 (3/3 replicas created)
 Type Reason Age From Message
 Normal ScalingReplicaSet 2m deployment-controller Scaled up replica set nginx-
            deployment-2035384211 to 3
  Normal ScalingReplicaSet 24s deployment-controller Scaled up replica set nginx-
           deployment - 1564180365 to 1
  Normal ScalingReplicaSet 22s deployment-controller Scaled down replica set nginx
             -deployment-2035384211 to 2
  Normal ScalingReplicaSet 22s deployment-controller Scaled up replica set nginx-
          → deployment - 1564180365 to 2
  Normal ScalingReplicaSet 19s deployment-controller Scaled down replica set nginx
            -deployment-2035384211 to 1
  Normal ScalingReplicaSet 19s deployment-controller Scaled up replica set nginx-
          → deployment – 1564180365 to 3
  Normal ScalingReplicaSet 14s deployment-controller Scaled down replica set nginx
         \hookrightarrow -deployment-2035384211 to 0
```

The developer immediately noticed the replica set section and found out that three pods are running, and all of them have been successfully updated. The exact step of an update is found in the 'Event' section. When the deployment was created initially, three pods have been created. Since the three pods were created from nothing, this action was classified as scale up. Then the user runs the update command; the second event fired up to create one copy of the new deployment, another scale-up. Kubernetes detected that four pods are now running and are therefore above the maximum threshold, it fires up the 3 rd event: a scale down from 3 to 2 for the original pod. This cycle runs until all the original pod is shut down and three new pods up and running.

This cycle that Kubernetes utilize is called "rollover" [2]. A strategy commonly practiced by network platform engineers to ensure stability during the update, meanwhile minimizing memory resource in an update. If there are N pods and all need to be updated, the maximum memory occupied is N+1 pod. And since the original pods already occupy N spaces. (N+1) N = 1 space is used, hence O(1), constant memory, which is the lowest memory usage achievable. This method also guarantees that at any time, there is at least N pod running. If three people are using the application during the update, one person would soon notice the change, then the second, then the third. None of them would be forced offline.

After creating the deployment, the developer creates service to expose deployments to the communication channel by executing Snippet 12.

Snippet 12: Creating Service on Kubernetes

```
apiVersion: v1
kind: Service
metadata:
  name: nginx
  labels:
   run: nginx
 type: NodePort
 ports:
    port: 8080
    targetPort: 80
    protocol: TCP
    name: http
   - port: 443
    protocol: TCP
    name: https
  selector:
    run: nginx
```

For each port definition, the user will provide an originating 'port' that listen for request and a 'targetPort' on deployments that the service forward to, and a connection protocol. Similar to how people make calls to www.google.com instead of https://www.google.com:80, in Kubernetes, the default protocol TCP and default port 80 can be left out. Note at spec.type, 'nodeport' is specified. This specifies the type of service. Kubernetes' default service type is Cluster-IP, which is only available within the cluster. Node Port exposes the service' node and port, as seen in Fig. 6.

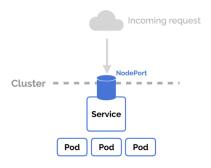


Fig. 6: Nodeport

The developer may directly access the service from their computers without ssh into the Kubernetes cluster. Nodeport exposes a vulnerability in the system that is prone to a Distributed Denial of Service (DDoS) attack, where the perpetrator uses a distributed system to spam requests to the node port system result in legitimate service lost in a flood of bot requests. Although some may argue that one can use nodeport during development and, subsequently, remove such vulnerability in production. This still cultivates a reliance on Nodeport for developers, which leads to a bad habit hard to change. This clear-cut solution removes human developers from developing such a tendency. The team has decided to ban the use of node port such that all traffic will be forced to go through a gateway coupled with a security layer, as Fig. 7 shows. Another explicit way to expose service is 'Load

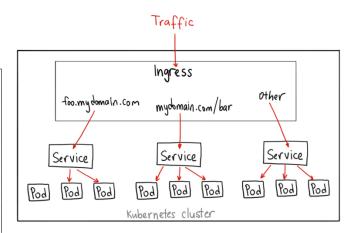


Fig. 7: An Ingress Gateway

Balancer'; it exposes the port to an external load balancer. Typically, this is used by services hosted on a cloud platform, which has provided a load balancing function, and the client need not perform load balancing task. Fig. 8 shows the architecture of one that hosts Kubernetes on Amazon Web Service. The web app supports entry from website1.com, website2.com, and website3.com through Amazon's Elastic Load Balancing (ELB) into service defined as "load balancer" service [11]. Since the Kubernetes is hosted on the cloud, traffic coming from ELB has already passed through Amazon's security layer. Clearly, this situation does not apply to the

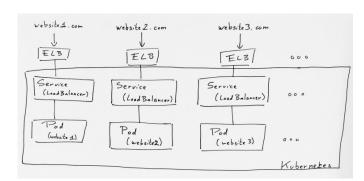


Fig. 8: Elastic Load Balancer

team that is developing a solution for hosting on-premise. A special type of service is called "virtual service" provided by Istio, a service mash manager developed by Google [12]. Virtual service, coupled with destination rules, defines a set of traffic routing rules to apply when a host is addressed. Each routing rule defines matching criteria for the traffic of a specific protocol. If the traffic is matched, then it is sent to a named destination service defined in the registry.

The last "kind" is Persistent Volume(PV). PV can be set up within a pod under the "spec" section or as its own pod [2]. Kubernetes uses volume to store data. To create volume, initiate a Persistent Volume Claim (PVC), a request for a PV. PVC defines the read and write authority of incoming traffic, and manage resources consumed by PV. Then the developer creates a PV that spawns up a stateful pod within the cluster where data is stored. A PV must exist corresponding to a PVC if the developer attempts to access the data. Otherwise, the PVC will dynamically spawn up PV pods to store the data. This method is used when Kubernetes wants to store some temporary data and only need the system to access it. However, the team discovered that storing stateful data in stateless pods and by labeling them stateful is rather paradoxical. The system would still treat these labeled stateful pods as stateless and shut down pods from time to time, resulting in data loss.

Moreover, spawning a persistent storage every time when the system initiates takes time. In Figure 13, only one pod containing a 'busybox' has been initiated and running, and it costs 400 MiB and at the peak of CPU usage 0.045 cores, as shown in Fig. 9. Thus, PV has not been investigated in detail for this project.



Fig. 9: Dashboard for Kubernetes

VI. SERVICE MESH

A service mesh is a dedicated infrastructure layer to run a fast, reliable, and secure network of microservices, a container orchestration system to provide a high level of deployment infrastructure [12]. Although the Kubernetes system provides the bare minimum backbone for service to service communication, it is far from easy to work with in terms of supporting a robust planetary-scale application. A couple of issues needs to be addressed:

- 1) How to communicate between services effectively
- 2) How to handle load balancing in the mesh without external Load Balancer

3) How to monitor each

In order to handle these issues, the team lead has investigated several solutions and decided to use Istio for its high degree of customization. The team is tasked with looking into extracting the useful aspect of Istio and applying them to the microservice prototype. Fig. 10 shows how Istio works by injecting sidecar into pods to perform mesh commands from Istio's master node. To install Istio, run script Snippet 13:

Before and After Istio

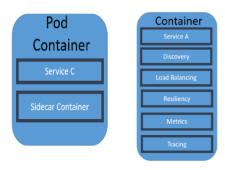


Fig. 10: Istio Sidecar Injection

Snippet 13: Istio Installation

```
#!/bin/bash
curl -L https://git.io/getLatestIstio | (sh -)
export PATH=$PWD/bin:$PATH
cd istio-1.0.0
kubectl apply -f install/kubernetes/helm/istio/templates/crds.yaml
kubectl apply -f install/kubernetes/istio-demo.yaml
kubectl get svc -n istio-system
kubectl get pods -n istio-system
```

The developer then runs Snippet 11 to check if service has been installed and Snippet 12 to check deployment.

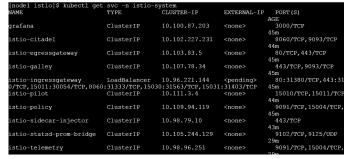


Fig. 11: Checking Istio Installation Services

If all pods are "ready" and "running", Istio has been deployed successfully. After deployment, Istio, together with Kubernetes, will consume approximately 1 GiB of memory as Snippet 14 shows. The memory consumption will increase logarithmically as more pods being instantiated. Istio works by injecting a sidecar into the pods, as shown previously in Fig. 10. The developer just needs to label the namespace with sidecar injection enabled by executing Snippet 14, Enabling Sidecar Injection on 'default' namespace.

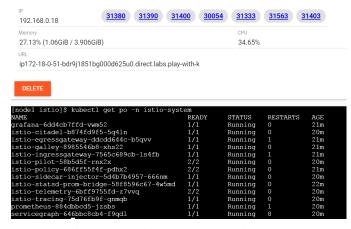


Fig. 12: Checking Istio Installation Pods

Snippet 14: Enable Istio Injection

\$ kubectl label namespace default istio-injection=enabled

The developer now can use Istio and all its functionality. Istio is made of two components (Fig. 13): a control pane and a data pane [12]. Much like Kubernetes master workers architecture. The control plane is the engine that offers the user an entry point; meanwhile, the data plane is the components injected into each pod to facilitate Istio functions. The architect of Istio does not want the developers to think of microservice in terms of master and workers. Instead, they what to enforce a parallel development mindset onto developers when developing the microservice architect. The developer may program with deployments and services simultaneously without having first to define all the services and bind them. When the developer spawns up deployment, a service is created by Istio and can be discovered by Envoy, thus removing the need to create service to expose deployments. Different from Kubernetes'

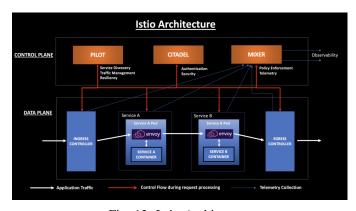


Fig. 13: Istio Architecture

service, at the core of Istio service discovery, is a technology called "Pilot" Pilot manages and configures all the envoy proxy deployed in the sidecar that the developer has previously embedded in each pod. Pilot lets the developer specify rules regarding traffic routing between envoy proxies and define a clear method to handle failure recovery in events such as time

out and circuit breaking. Note that the Pilot has embedded load balancing when the traffic travels through the Envoy. Figure 16 shows the working principle of Pilot [12]. SvcA does not need to have the knowledge of the deployment to access pods in deployment. Previously the service must have defined the access to SvcB Pod1 or svcB Pod2...svcB Pod4 to access the deployment.

VII. TRAFFIC ACCESS

A few common practices when it comes to defining traffic control are round-robin and weight. Round robin is the default mode, where each request will go to the subsequent instances, hence evenly distribute the traffic. In the weighted traffic access, the user assigns a weight to different deployments, as shown in Snippet 15, so the traffic is split according to user-defined weight.

Snippet 15: Defining Virtual Service to Configure Traffic Split Weight

```
apiVersion: networking.istio.io/v1alpha3
kind: VirtualService
metadata:
  name: reviews
spec:
  hosts:

    reviews

  http:
   - route:
      - destination:
         host: reviews
         subset: v1
      weight: 75
       destination:
         host: reviews
         subset: v2
      weight: 25
```

The user may also define a timeout and number of retries (Error! Reference source not found.) when the request initially fails to travel through. By industry standard, five retries and 200ms are designated. The time out is calculated based on relaying timeout. That is a gateway to service timeout $\frac{1}{2}$ 2000ms. Now, let L be a number of layers, and t be timeout defined in milliseconds, and user side time out $\frac{1}{2}$ 10s, and number of ties (R) $\frac{1}{6}$ = 1, one for the initial request, and a safety factor of 3, let C be constant time from user machine to Gateway, we uses a system of equation,

$$\begin{cases} tL < 2000ms \\ 3(tL) + C < 2000ms \end{cases} , \tag{1}$$

with limits obtained by,

$$Limits = 3(X < 2000) + C < 10000$$

$$= (3X < 6000) + C < 10000$$

$$= C < 10000 - (3X < 6000) .$$

$$= C < 10000 - 6000$$

$$= C < 4000$$
(2)

Hence assume $X = X_{max}$, C=4,000. The machine needs to connect in 4 seconds. However, C can be much lower if the host is a cloud provider. For example, Ping amazon at Toronto gives 6ms, so C= 4 is much smaller than 4000, and X can be much greater. However, for On-premise, we assume

the worst-case scenario for C that is the smallest delay at L. Since

$$c = 4000, 3tl = 6000, tl = 2000, t > 0, L > 0, L \in \mathbb{Z}$$
 (3)

To obtain L, count the number of layers including platform itself:

```
Kubernetes: L0, userapp: L1, account = L2, (4)

friends = L2: total = 3, 2000 = t(3), sot = 667 (5)
```

Apply the answer to write the virtual service in Snippet 16.

Snippet 16: Virtual Service with Calculated Timeout Limit

```
apiVersion: networking.istio.io/v1alpha3
kind: VirtualService
metadata:
 name: userapp
spec:
 hosts:

    userapp

  http:
   - route:
       destination:
         subset: v1
    retries:
      attempts: 3
      perTryTimeout: 667ms
       destination:
         host: account
         subset: v1
    retries:
      attempts: 3
      perTryTimeout: 667ms
       destination:
        host: friends
         subset: v1
      attempts: 3
      perTryTimeout: 667ms
```

VIII. DATA STRATEGY

Istio also offers three ways to handle data strategy [12]. Mesh expansion is the ideal way. It elevates the items from outside mesh to enjoy the same privilege as if it is inside the mesh. As it shows in Fig. 14-left, services in mesh(green) cannot access the database outside the Kubernetes. It can access services in Kubernetes but not in Istio through a gateway(black). After a service expansion (Fig. 14-right) on database and external service, the services can directly access each other. Meanwhile, the service inside Istio can access the database through a service entry(white). However, as of Istio 1.3.0, Mesh Expansion is not ready according to Google, it only works 50% of the time. The production environment cannot tolerate such instability. Another choice is to host data storage directly using the built-in Kubernetes capability. That is to utilize PV and PVC, as mentioned in 6.3. Kubernetes lets developers make claims to storage space, then spawn up pods to realize the storage facility. Such an option was rejected as storing data in Kubernetes' pod is basically storing stateful content in stateless infrastructure. Data can be lost when the system decides to shut down pods. As Figure 17 shows our last option: the service can connect to the outside database through an open gate (White) called "service entry". Service entry opens an "entry" into the service and lets only data to

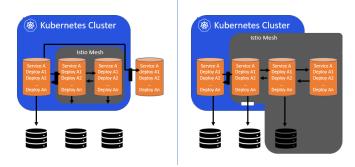


Fig. 14: Mesh Expansion

be transmitted. Action commands cannot be communicated through service entry; they must go through a gateway. The team decided to use service entry to connect the microservice to the existing database for the time being and leave the schema creation and database break down until Google fixed mesh expansion.

IX. CANARAY DEVELOPMENT

In the "updating pod" example, we demonstrated Kubernetes' capability of rolling out. Kubernetes recursively spawns a new pod then takes down an old one until all old pods are removed [12]. While Kubernetes' performance is stellar on a deployment level, Istio is needed to handle more robust yet fine-grained rollout on service level. Take the example of updating 1% of 1000 instances of account service: if a command to update 10 of the pods is ran, the developer does not know which ten pods are updated. To address this problem, the developer can choose to create a new deployment under a different label. However, under a new label, all configurations from the deployment are not inherited. The autoscaler from the old deployment does not control the new deployment. Istio offers control over the rollout process by limiting traffic going into new pods upon initialization, then gradually increases the traffic into the new pods as they mature. If errors happen, Istio rolls back to the previous version. More sophistically, Istio offers control over the region, user, or other properties at the developer's request.

Table ?? shows how Istio differs from Kubernetes. In plain Kubernetes, the developer needs to create two deployments with different labels, and each needs the specific number of pods specified at "replica". In Istio, the developer creates a new "kind" called virtual service, in which she specified the weight of traffic going into each respective version. Now a simple test can be performed to see the effect. First, run Snippet 17 to enable autoscaling:

Snippet 17: Enable Autoscaler on Kubernetes Deployment

```
$ kubectl autoscale deployment helloworld-v1 --min=1 --max=10
$ kubectl autoscale deployment helloworld-v2 --min=1 --max=10
```

Then after spawning a few minutes of request, one will notice that the first deployment scales up much faster than the second

plain	Istio		
kind: Deployment	kind: VirtualService		
metadata:	metadata:		
name: helloworld-v1	name: helloworld		
spec:	spec:		
replicas: 1	hosts:		
template:	 helloworld 		
metadata:	http:		
labels:	- route:		
app: helloworld	— destination:		
version: v1	host: helloworld		
spec:	subset: v1		
containers:	weight: 90		
image: helloworld-v1	— destination:		
	host: helloworld		
	subset: v2		
apiVersion: extensions/v1beta1	weight: 10		
kind: Deployment			
metadata:	apiVersion: networking.istio.io/v1alpha3		
name: helloworld-v2	kind: DestinationRule		
spec:	metadata:		
replicas: 1	name: helloworld		
template:	spec:		
metadata:	host: helloworld		
labels:	subsets:		
app: helloworld	- name: v1		
version: v2	labels:		
spec:	version: v1		
containers:	- name: v2		
image: helloworld-v2	labels:		
	version: v2		

TABLE II: Plain vs Istio Inject Kubernetes Config

one, corresponding to the 9:1 ratio specified in the virtual service (Snippet 18):

Snippet 18: Check Autoscaler Effect

```
$ kubectl get pods | grep helloworld helloworld-v1-3523621687-3q5wh 0/2 Pending 0 15m helloworld-v1-3523621687-73642 2/2 Running 0 11m helloworld-v1-3523621687-7hs31 2/2 Running 0 19m helloworld-v1-3523621687-dt7n7 2/2 Running 0 50m helloworld-v1-3523621687-gdhq9 2/2 Running 0 11m helloworld-v1-3523621687-gras4t 0/2 Pending 0 15m helloworld-v1-3523621687-jrs4t 0/2 Pending 0 19m helloworld-v1-3523621687-wwddw 2/2 Running 0 15m helloworld-v1-3523621687-wwddw 2/2 Running 0 15m helloworld-v1-3523621687-wrddw 2/2 Running 0 19m helloworld-v1-3523621687-xlt26 0/2 Pending 0 19m helloworld-v2-4095161145-963wt 2/2 Running 0 50m
```

To filter request by a specific case, the developer can use the match modifier in Snippet 19

Snippet 19: Uses Match Function

```
http:
- match:
- headers:
cookie:
regex: ""(.*?;)?(email=[^;]*@company.com)(;.*)?\$"
```

Snippet 21 shows the user filters incoming requests containing email ending with domain "@company.com", then delegate all traffic matching this condition to v1 or v2 as shown in Table 3. If the developer wants to delegate all traffic from the United States to V2 and Canada to V1, the developer can have programmed a match for "location", and in the app, but the location in the header whenever the user sends the request. The match will then intercept the location in each request and delegate the traffic to the corresponding destination.

A. Developing an Efficient Aggregated Service Algorithm

During the development of the BFF service, a couple of function challenged the developer's ability to program with microservice. One is notably getting total post count. Traditionally, a SQL script is run in the database that contains both the account and the friends. However, according to the aforementioned logical separation guideline, one should never have more than one domain in one database. The SQL solution is quickly eliminated. The first algorithm developed was Snippet 20. The developer does a parallel stream on posts service to get a list of posts the originating user's account. Whilst getting the posts, get posts' like from this persons' friend using the Post ID. The quantity with price and sum all up likes to return total likes.

Snippet 20: Algorithm to Get all Posts Likes from Friends

```
\begin{aligned} Double\ likeCount &= posts.parallelstream().filter(account\_id).collect((post) \ -> \ friends. \\ &\hookrightarrow getPosts(post).getLikes()).reduce() \end{aligned}
```

While there is only one call to post at ".parallelstream()" there are N calls, N = number of likes, to friends. However, since position and assets are two different services, N number of inter-service calls are made. Each Inter-service calls travel through 2 layers, as Fig. 15 shows. This then becomes N(2L) total runtime. To handle It is best if the number of calls can

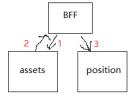


Fig. 15: Sigal Path

be inter-service call can be simplified from N to 1. Hence a new solution is developed in Snippet 21:

Snippet 21: Optimized Algorithm

```
ConcurrentHashSet<Post> set = posts.parallelstream()
    .filter(account_id).collect((posts))

Double likeCount = friends.getAllPosts(set).getLikes().parrallelstream().reduce()
```

Here although 2 streams are used the application run time increased from N to 2N, but since there is only one inter service call. The total runtime becomes 2N + L. Later algorithm is used to facilitate all aggregated services. This ensure a maximum of 2N + 2L = 2(N+L) run time anywhere in the network, where N = number microservices of concern L = layers between these microservices.

B. Memory Optimization

Upon finishing the code, memory optimization becomes the next issue. While duplicating the application's java dependencies in each microservice' container makes the application only slightly smaller than the original application. For example, if the original service costs 800 MiB, the four atomic services each costs 400 MiB, and 1 BFF cost 600. Summing the size of all microservices: 4*400 + 600 = 3200mb, new microservices use four times as many memory spaces as

the original application. When monitoring the size of the microservice application, we can use the guideline (6)

$$(\sum_{i=1}^{n} x) > X, x_i < X,$$
 (6)

, where x is the size of each microservice, n is the number of services, and X is the size of the original service. In short, this formula requires each microservice to be smaller the size of the monolithic application, but the total size will always be greater than the original application. The repetitive load balancing is unavoidable in containers. Since load balancing has been handled by Kubernetes' load balancer, the developers are safe to remove all load balancing mechanisms within the java app that handles loading as Kubernetes will handle the loading. First, remove the JPA library used to connect the database from the BFF because it does not connect to the database. Very quickly, the size decreased by 200 MiB. Next, the team looks at minimizing the memory size; two equations are used to model total memory use (7).

$$memory = heap + non - heap$$

 $non - heap = threads \times stack + classes \times 7 \div 1000$

By referring to the Spring website, the team finds Table III. The developer chooses to restrict the threads to 3 and let

APPLICATION	HEAP	NON HEAP	THREADS	CLASSES
Vanilla	22	50	25	6200
Plain Java	6	14	11	1500

TABLE III: Spring and Java Memory Usage

Kubernetes expands when it needs more pods. Next, since JPA has been taken out 5000 classes are removed. As the results in Fig. 16 show, eventually the memory drops to 375mbs with total heap usage cycles from 50 to 300.

C. Deployment with Kuberenetes

The last step is to package the deployment into Kubernetes files and figure out a way to manage the configuration. Spring offers a git config server that supports remote configuration. Put Snippet 22 in the application.yml file [12]:

Snippet 22: Application Level Environment Variables defined in "application.yml"

```
spring:
  cloud:
    config:
      server:
          uri: https://github.com/spring-cloud-samples/config-repo
           repos:
             development:
               pattern:
                    '*/development
                  - '*/staging'
               uri: https://github.com/development/config-repo
             staging:
               pattern:
                  – '∗/qa'
                    '*/production'
               uri: https://github.com/staging/config-repo
```

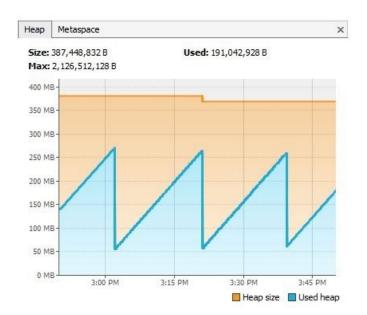


Fig. 16: Heap Usage After Optimization

Now the spring app will read configuration stored in a remote git repository https://github.com/development/config-repo. Git config allows the deployment team to make changes outside the deployment server and not have to work with the application code. In Kubernetes, the developer will use the config map to achieve the same goal. A config map is a user-defined key-value pair store in the clusters etcd. Config map is visible to all pods within the cluster. For example, the user can run Snippet to deploy a config map that contains key and password to a database:

Snippet 23: Kubernetes Pod level Environmental Variables Defined in Config Maps

```
kind: ConfigMap
name: db_info
metadata:
database: jdcb://db.cca.com:3306
name: db
password: -jyc9ep2
```

Then in the deployment file, add a reference to this config map as Snippet 24 does.

Snippet 24: Adding Reference to Config Map on Deployment

```
apiVersion: v1
kind: deployment
metadata:
name: account
spec:
containers:
- name: test
image: k8s.gcr.io/busybox
env:
- name: SPECIAL_LEVEL_KEY
valueFrom:
configMapKeyRef: db_info
restartPolicy: Never
```

Now in the application's "application.yaml", one can directly use values from the config map as a system environment variable. If multiple config maps are used with the same values, simply add the config map name prepended to the variable name in Snippet 25.

Snippet 25: Dereferencing Environmental Variables

\${name} \${db_info1.password} \${db_info2.password}

Configuration can be easily managed outside of the deployment. Similar to the idea of the rollout, every time the configuration is edited, the deployer will run a restart command to Kubernetes, which will bring up a new pod with the new configuration meanwhile bring one old pod down until all old pods are replaced by the new ones. Towards the conclusion of the prototype construction, an approach was discovered to build a pipeline (Fig. 17) to connect the remote config repository with the Kubernetes cluster such that whenever the remote config repository is changed, no human involvement is needed to trigger the rollout. The pipeline will recognize a new change and trigger a pod restart command to Kubernetes' Istio framework. One such plugin to perform the task is Weave Flux [13]. The developers may choose to build their own pipeline from scratch.

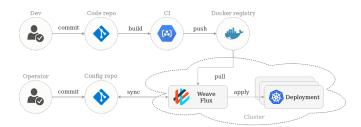


Fig. 17: Continous Integration Pipeline with Kubernetes and GitOps

X. CONCLUSION

In this paper, we introduced a process to transform an existing application to a microservice system that is scalable, portable, and continually deployable. We also discussed service to service communication, data strategy, and Canary development in the new set up. Eventually, we constructed a prototype in Java using a microservice strategy. It involves adding Istio, a mesh service platform for Kubernetes. In Istio, service communicates by calling each's IP by name, when a pod exceeds its predetermined threshold, Istio automatically scale up according to the rule specified by the user. This rule lets developers filter incoming requests and set the percentage of traffic going into different versions of pods, a useful feature in the A/B split test. Although the team wishes to use Mesh expansion to incorporate new data schema into the service in the future, the team decided to settle with using service entry to open a channel for data operation. During the application, a couple of side problems was addressed. First, a methodology of minimizing inter-service calls by combining all information needed into one call, enabling a 2(N + L) runtime, where N = number of microservices, and L = number of layers between services. A quick tweak was performed to remove the unneeded library in spring to minimize jar memory usage during deployment. Lastly, the team decided to use Config map a Kubernetes object to supply needed configuration data, thus allowing the deployment team to edit configuration from a remote repository quickly.

A. Recommendation

As mentioned in 7.3, we have not had a chance to build the pipeline to enable automatic configuration rollover. Having such a feature would benefit the fluidity amongst the deployment team greatly. When two developers working on the same microservice changed a configuration at the same time, some data from the microservice may be lost due to concurrency. If A pipeline is implemented, it will establish a queue to roll over one set of changes before rolling over the next set and keep a log of all changes. When the developer team wants to roll back, a quick reversal action can be executed from the pipeline to check out the git history of the old configurations manually.

REFERENCES

- [1] "Horizontal vs Vertical Scaling," [Online]. Available: http://pudgylogic.blogspot.com/2016/01/horizontal-vs-vertical-scaling.html. [Accessed 1 January 2019].
- [2] G. Inc., "Kubernetes Concept," Google Inc., [Online]. Available: https://kubernetes.io/docs/concepts/services-networking/ingress/. [Accessed 2 January 2019].
- [3] pivotal.co, "Spring.io," Pivotal, [Online]. Available: https://spring.io/. [Accessed 2 January 2019].
- [4] S. Stein, "Scaling out with Microsoft Azure SQL Database," Microsoft Inc., [Online]. Available: https://docs.microsoft.com/en-us/azure/sql-database/sql-database-elastic-scale-introduction. [Accessed 2 January 2019].
- [5] Z. Hale, "Cloud ERP vs. On-Premise ERP," Software Advice Inc., [On-line]. Available: https://www.softwareadvice.com/resources/cloud-erp-vs-on-premise/. [Accessed 2 January 2019].
- [6] T. Townsend. "This what Google is Snap is paying billion for," Vox Media, [Online]. Available: https://www.recode.net/2017/3/1/14661126/snap-snapchat-ipospending-2-billion-google-cloud. [Accessed 2 January 2019].
- [7] A. Wiggins, "The Twelve Factor App," Heroku inc., [Online]. Available: https://12factor.net/. [Accessed 2 January 2019].
- [8] Z. Dehghani, "How to break a Monolith into Microservices," Throught Works, [Online]. Available: https://martinfowler.com/articles/breakmonolith-into-microservices.html. [Accessed 2 January 2019].
- [9] Docker.com, "Docker Get Started," Docker Inc., [Online]. Available: https://www.docker.com/get-started. [Accessed 2 January 2019].
- [10] I. Gunaratne, "Beginner's Guide to Kubernetes," Google Inc., [Online]. Available: https://medium.com/containermind/a-beginners-guide-to-kubernetes-7e8ca56420b6. [Accessed 2 January 2019].
- [11] "Elastic Load Balancing," Amazong Web Services, Inc., [Online]. Available: https://aws.amazon.com/elasticloadbalancing/. [Accessed 2 January 2019].
- [12] I. Authors, "Istio Manual," Google, Inc., [Online]. Available: https://istio.io/. [Accessed 2 January 2019].
- [13] "What Flux Does," Weaver, [Online]. Available: https://github.com/weaveworks/flux. [Accessed 2 January 2019].