

Deduplicating Cloud Functions - Demo 3

Members:

Beliz Kaleli

Vikash Sahu

Paritosh Shirodkar

Asutosh Patra

Mentor:

Shripad Nadgowda

Previously in De-duplicating...

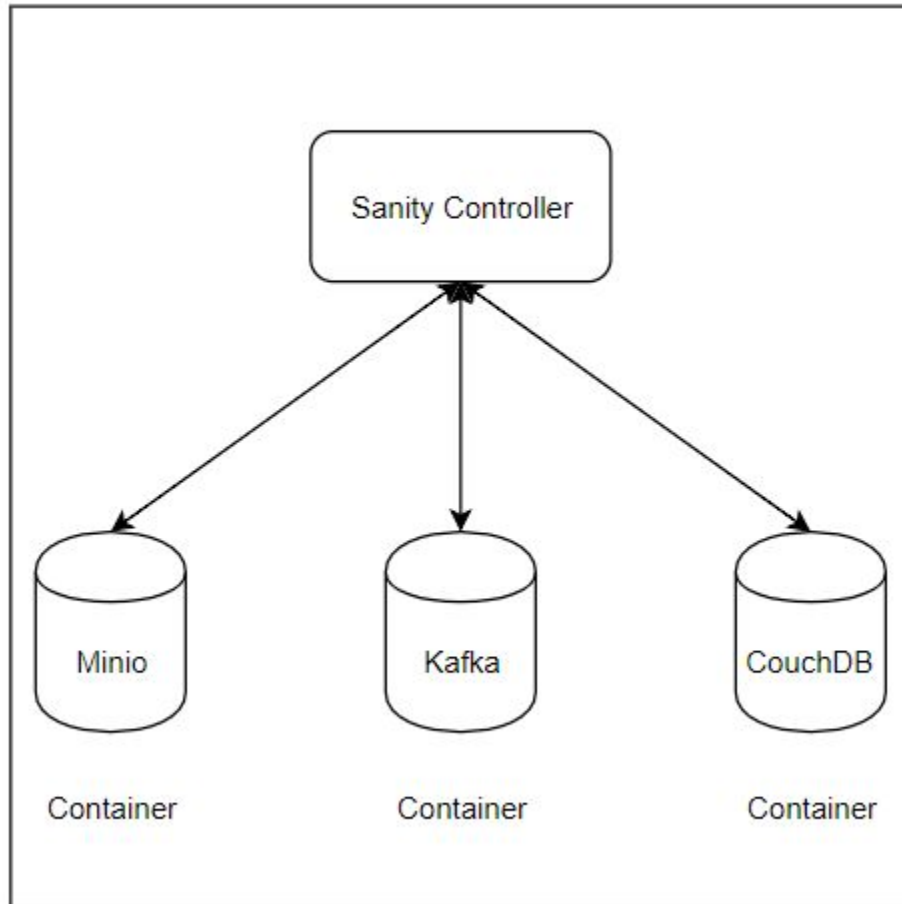
- Earlier, all the components (Minio, Kafka, CouchDB & OpenWhisk) were executing independently in local machines
- Interaction among the components was not done
- The feedback from the last sprint was to deploy the whole framework into cloud and execute

What did we achieve in this sprint?

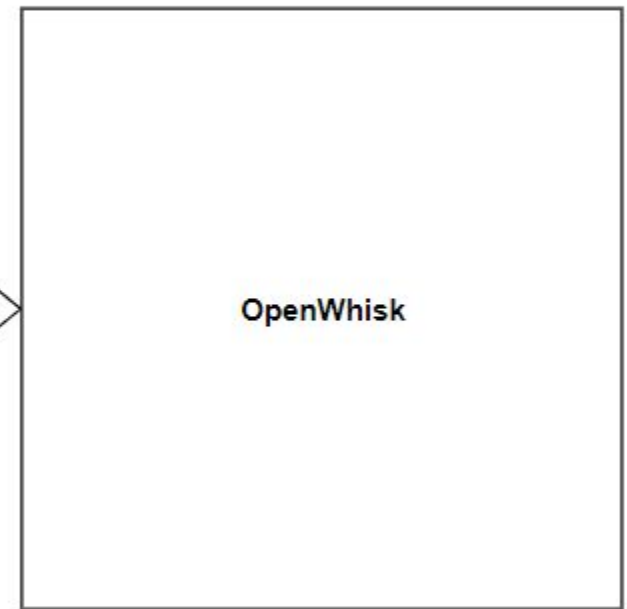
- Deployed all the individual components in the cloud
- Built the sanity framework to avoid deduplication of cloud function
- Successfully implemented **deduplication** of cloud function for **Image Thumbnail Use Case**

Now, Simplified working model in IBM Cloud

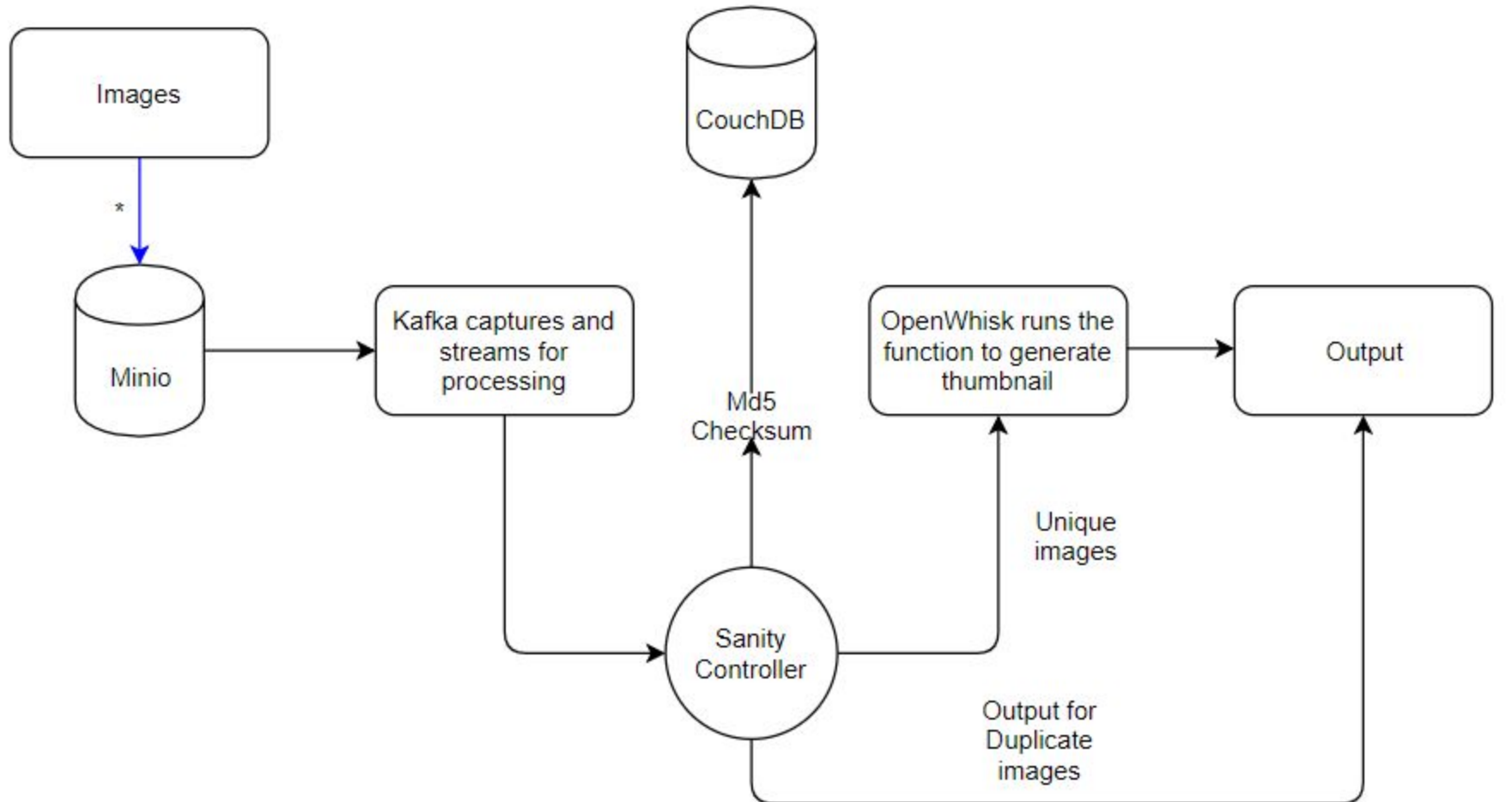
Server 1

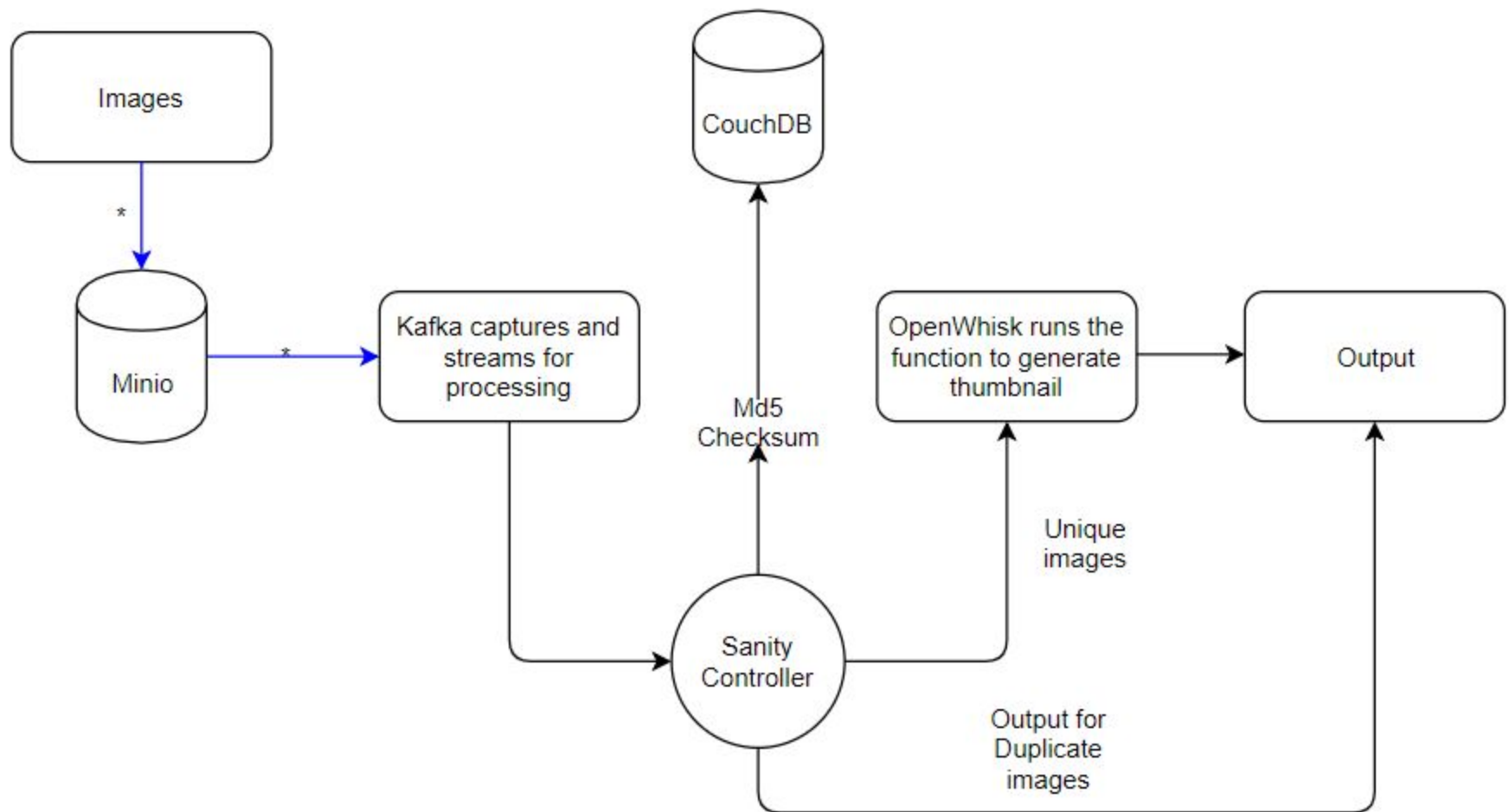


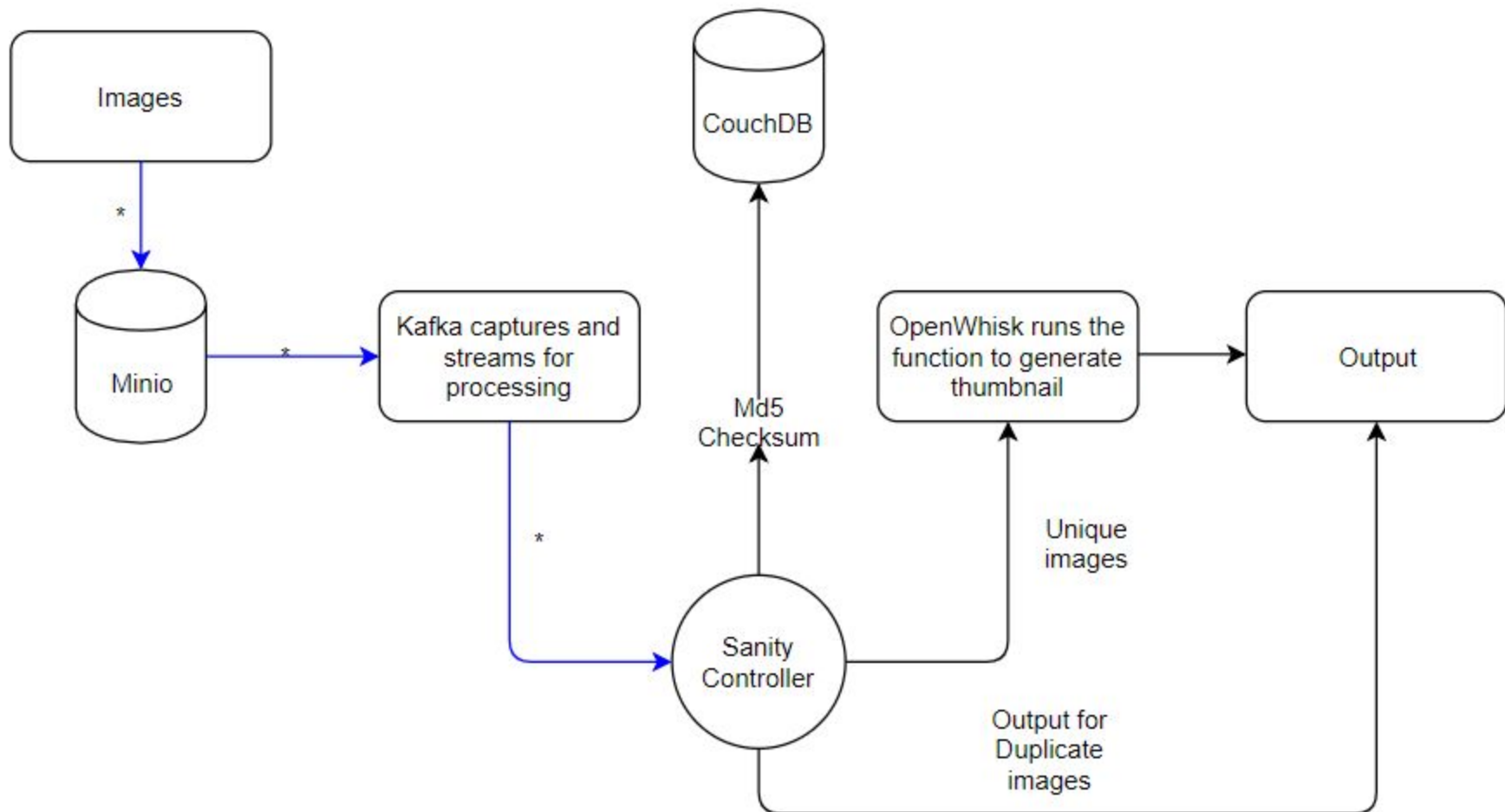
Server 2

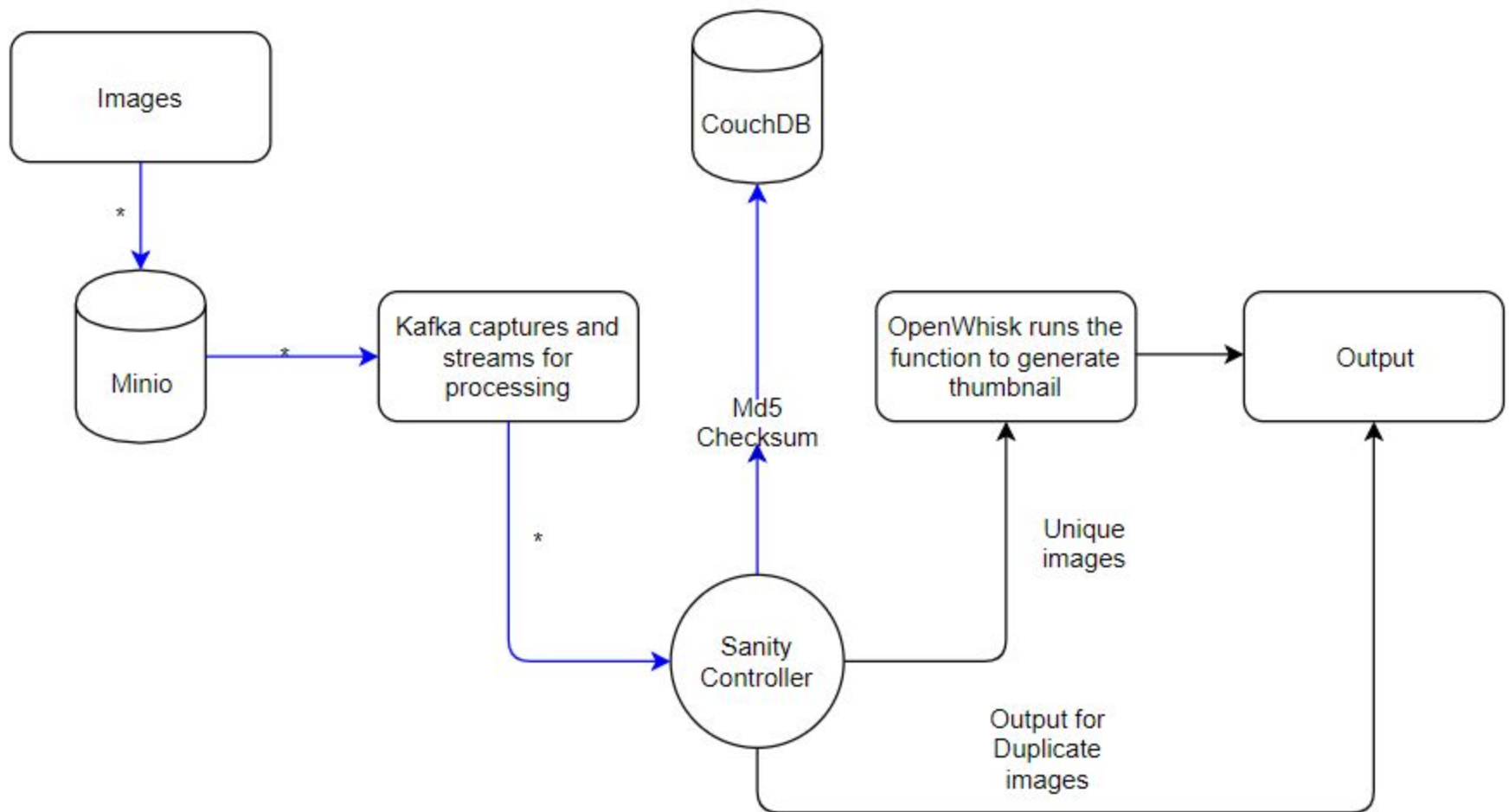


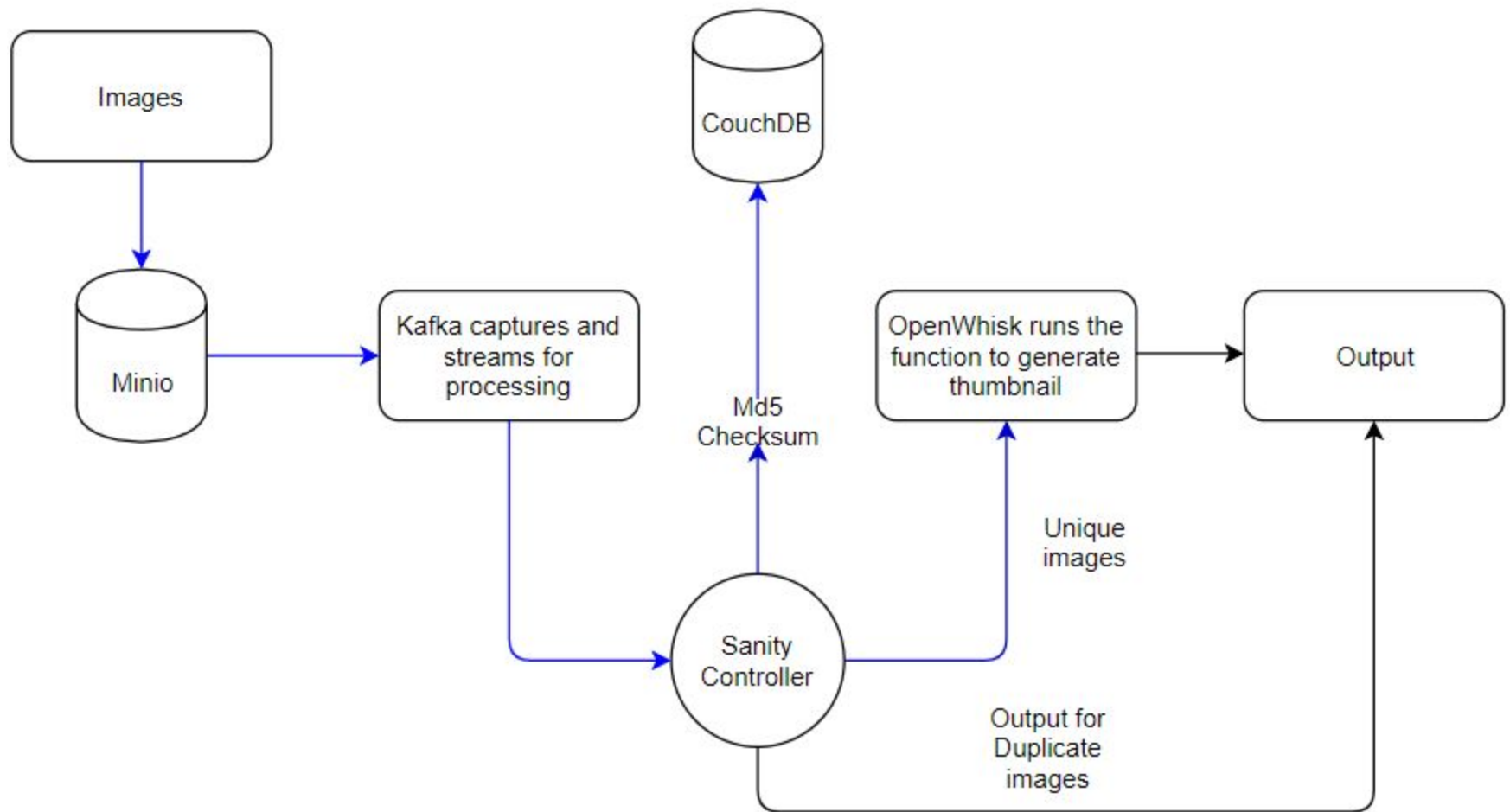
Architecture Diagram

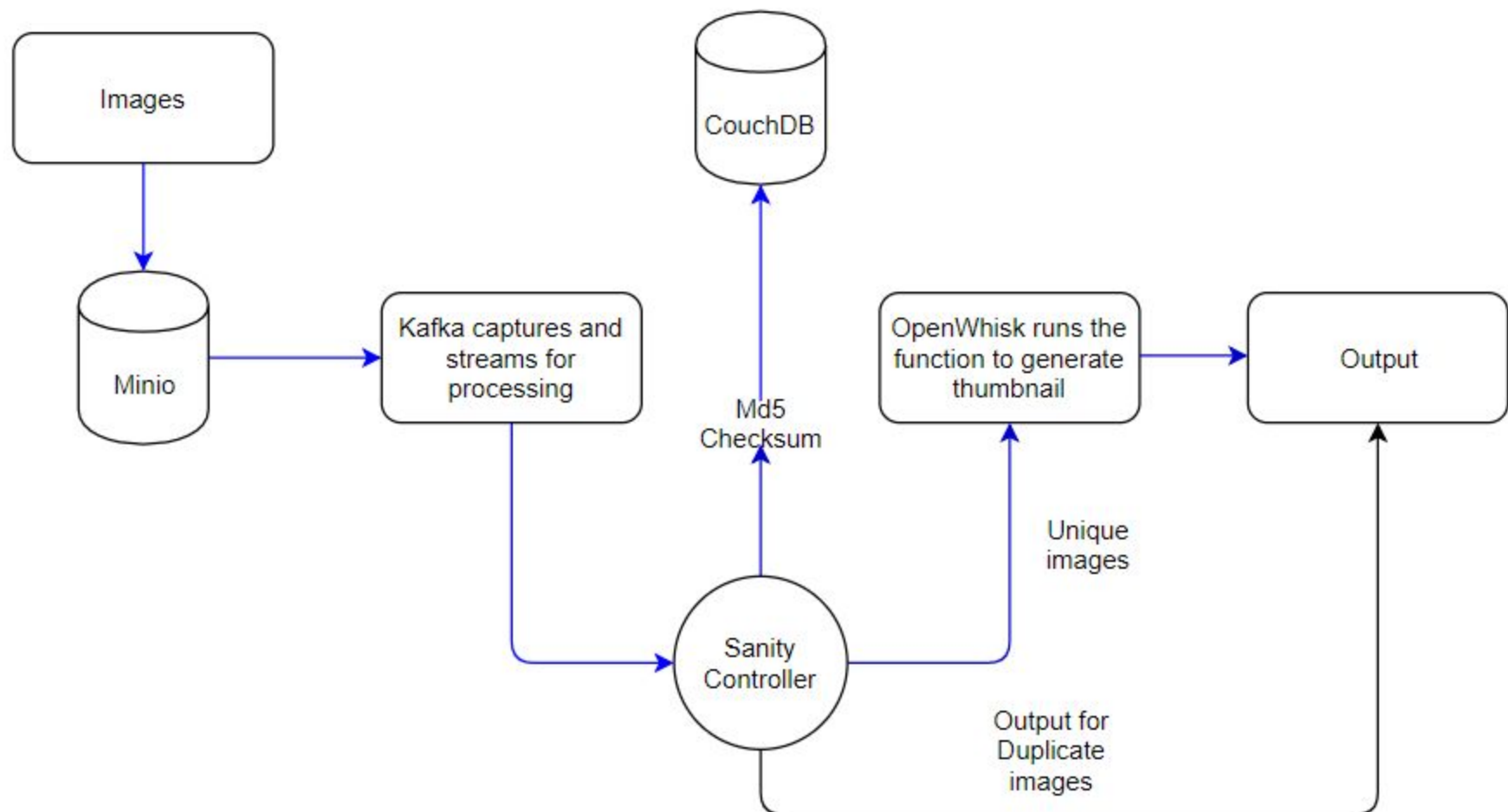


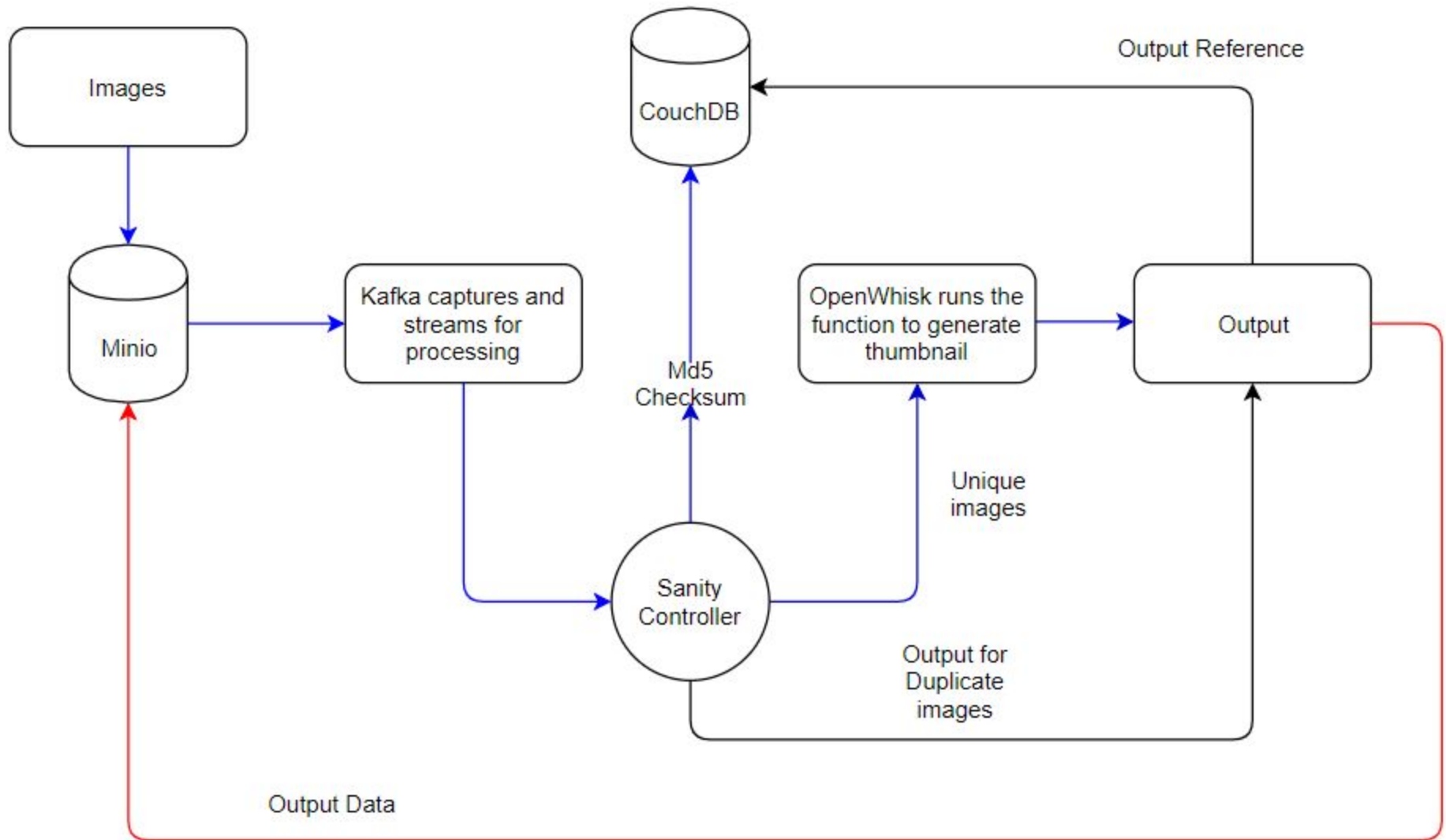


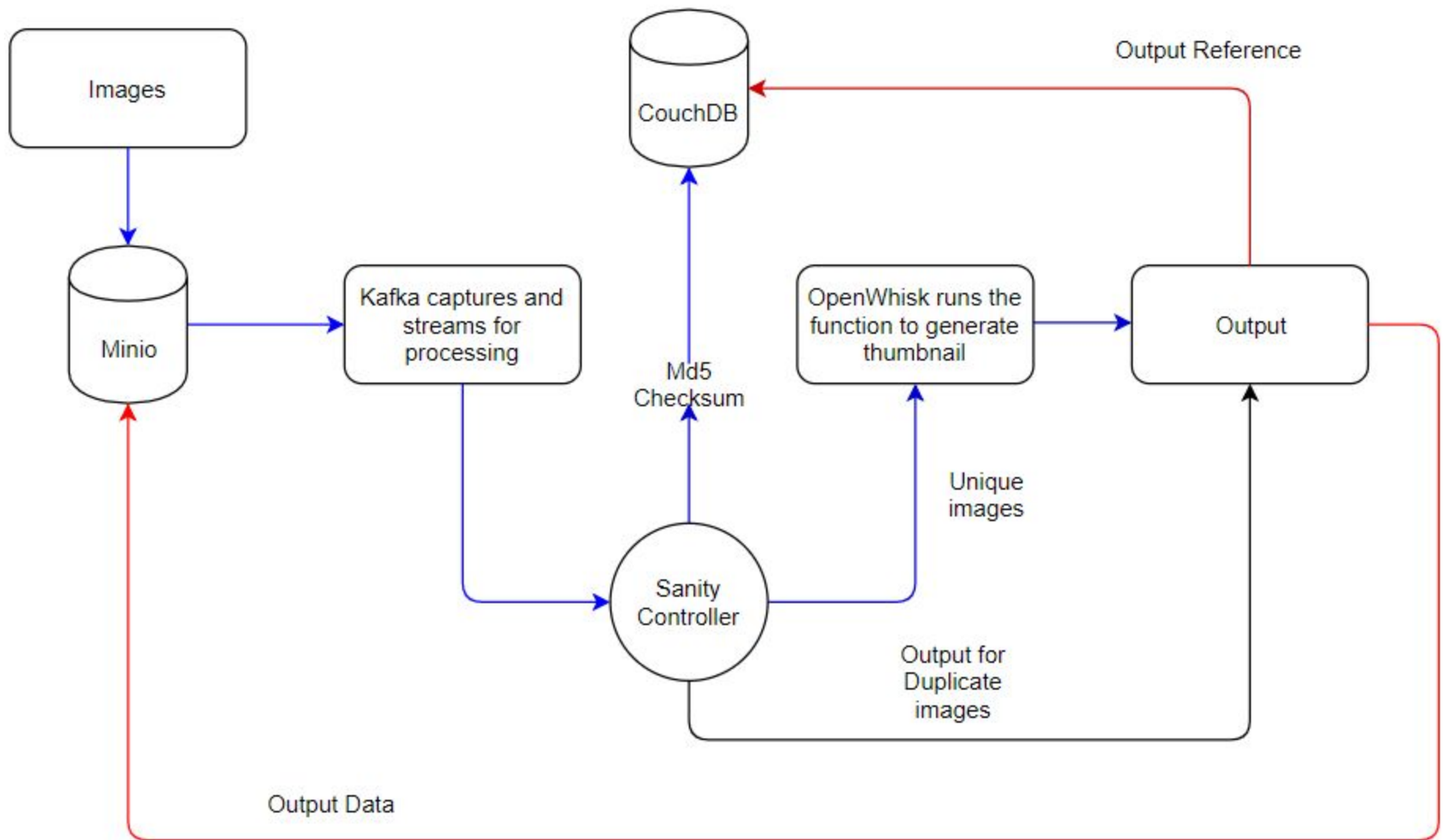


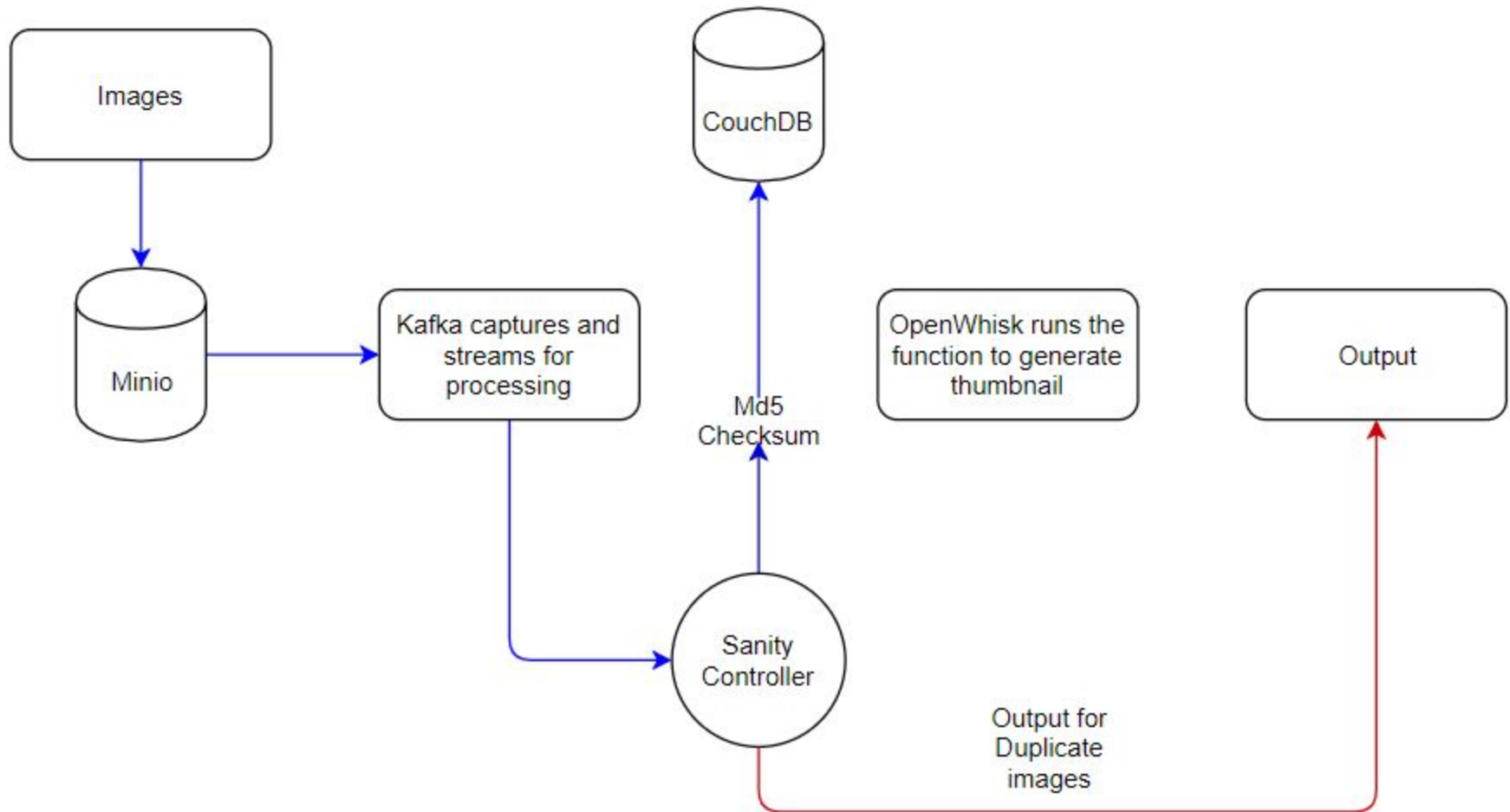












DEMO

How is deduplication achieved via Sanity Framework

- Upload a file to a Minio bucket (ex: inbucket, a.jpg)
> `mc cp a.jpg myminio/inbucket`
- Create a function (ex: myfunction) with wsk client
> `wsk -i action create myfunction.py myfunction`
- Invoke the function with the uploaded file
> `wsk -i action invoke myfunction inbucket/a.jpg outbucket`
- Result is automatically stored in Minio bucket(ex: outbucket)
- For duplicate images, we get the reference from Sanity and for unique images, openwhisk executes action and stores the thumbnail, hence de duplication of functions is achieved.

Couch DB

```
1 {  
2   "_id": "dae770ad388c898fa85dc140a0014a24",  
3   "_rev": "5-2b4a80a884faf4eab778ab5902f5e3ae",  
4   "function1hash": {  
5     input1forfunction1hash: "outputbucketfor_input1forfunction1hash/outputfilenamefor_input1forfunction1hash",  
6     input2forfunction1hash: "outputbucketfor_input2forfunction1hash/outputfilenamefor_input2forfunction1hash"  
7   },  
8   "function2hash": {  
9     input1forfunction2hash: "outputbucketfor_input1forfunction2hash/outputfilenamefor_input1forfunction2hash",  
10    input2forfunction2hash: "outputbucketfor_input2forfunction2hash/outputfilenamefor_input2forfunction2hash"  
11  }  
12 }
```

input_checksum

output_bucket/output_file

Deduplication In Progress..

test1 /

Used: 7.97 GB

Free: 85.84 GB

Name

Size



1.jpg

75.29 KB



2.jpg

75.29 KB



a.jpg

2.26 MB



duplication.jpg

16.24 KB



test.jpg

16.24 KB

Same file, different name!

After Running: Only one output for both files!

test2 /

Used: 7.97 GB

Free: 85.84 GB

Name

Size



1-thumbnail.jpg

1.42 KB



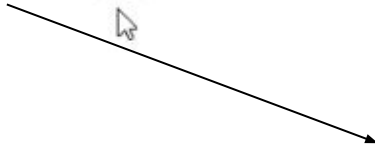
a-thumbnail.jpg

2.16 KB



test-thumbnail.jpg

3.36 KB



only output for
test.jpg and
duplication.jpg

Challenges in Current Sprint

- Configuring docker components in cloud and executing functions in them
- How can we invoke custom docker scripts in openwhisk?
- Creating network connection between two servers in IBM cloud
- Simulating the entire pipeline with all the components in place in cloud

Next Steps (Sprint - 4)

- Implementing the Command Line Interface
- Benchmarking performance savings for entire process
- Brainstorm with our mentor to improve our current framework

Burndown Chart

DEDUPLICATING-CLOUD-FU... SPRINT 3 28 FEB 2019-21 MAR 2019



100%

108 total points

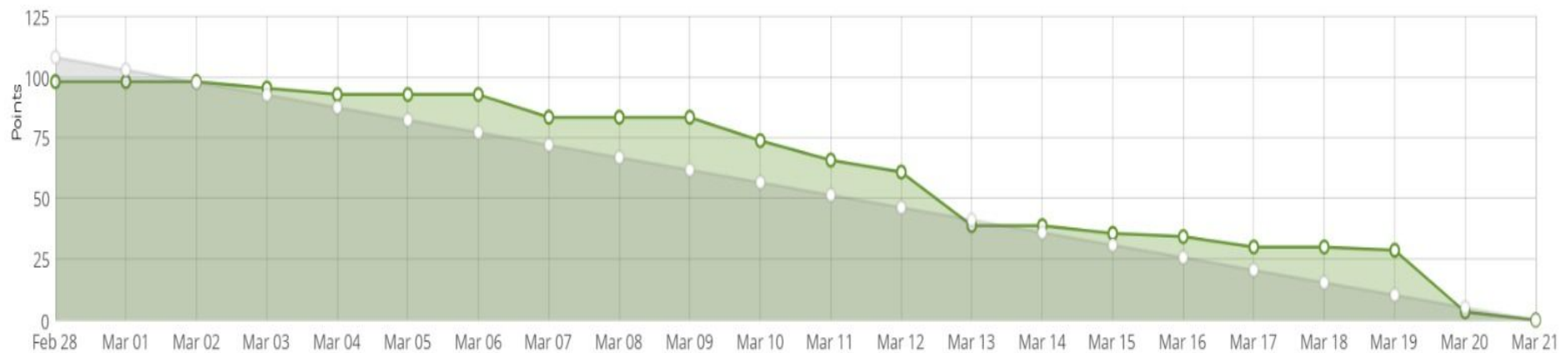
108 completed points

0 open tasks

31 closed tasks



0 cocaine doses



THANK YOU