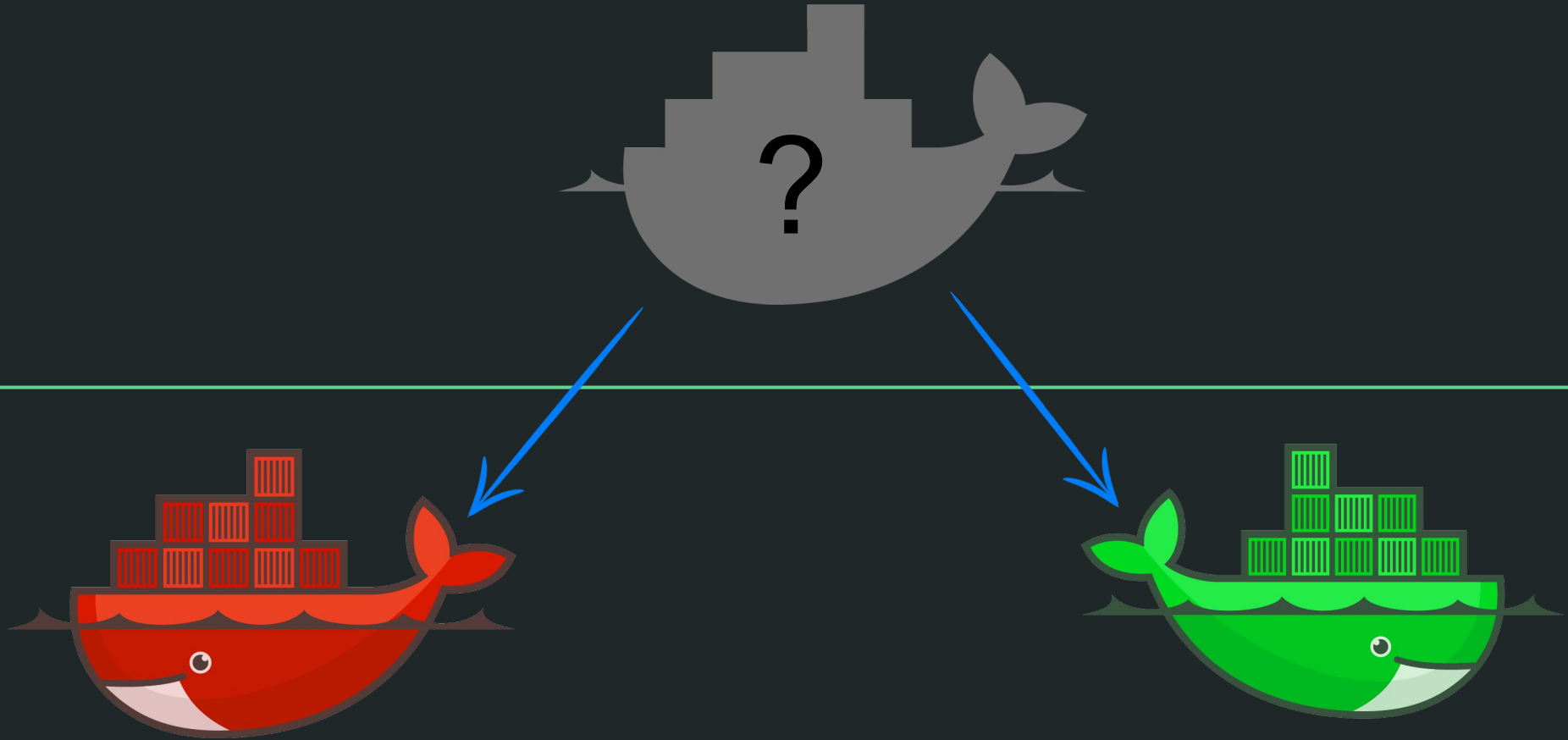


Container Code Classification

Gao, Jeremy, Kostas, Ozan, Rahul

Mentor: Sastry S Duri
(IBM Research)

The Goal: to Classify Safety



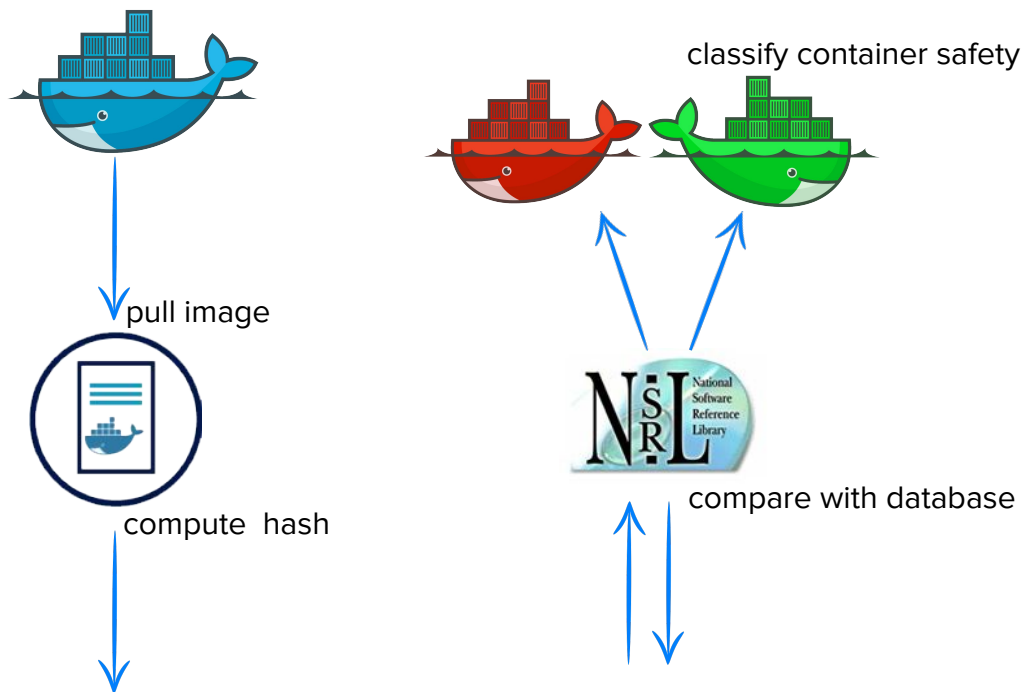
How can we detect suspicious files in a docker container?

Research Project Goal:

- Find a way to reduce the dataset to a number of suspicious files in an image
- We are testing different models and comparing them to find the most accurate prediction of the safety of a container

Procedure:

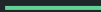
1. Retrieve
2. Hash
3. Compare
4. Classify

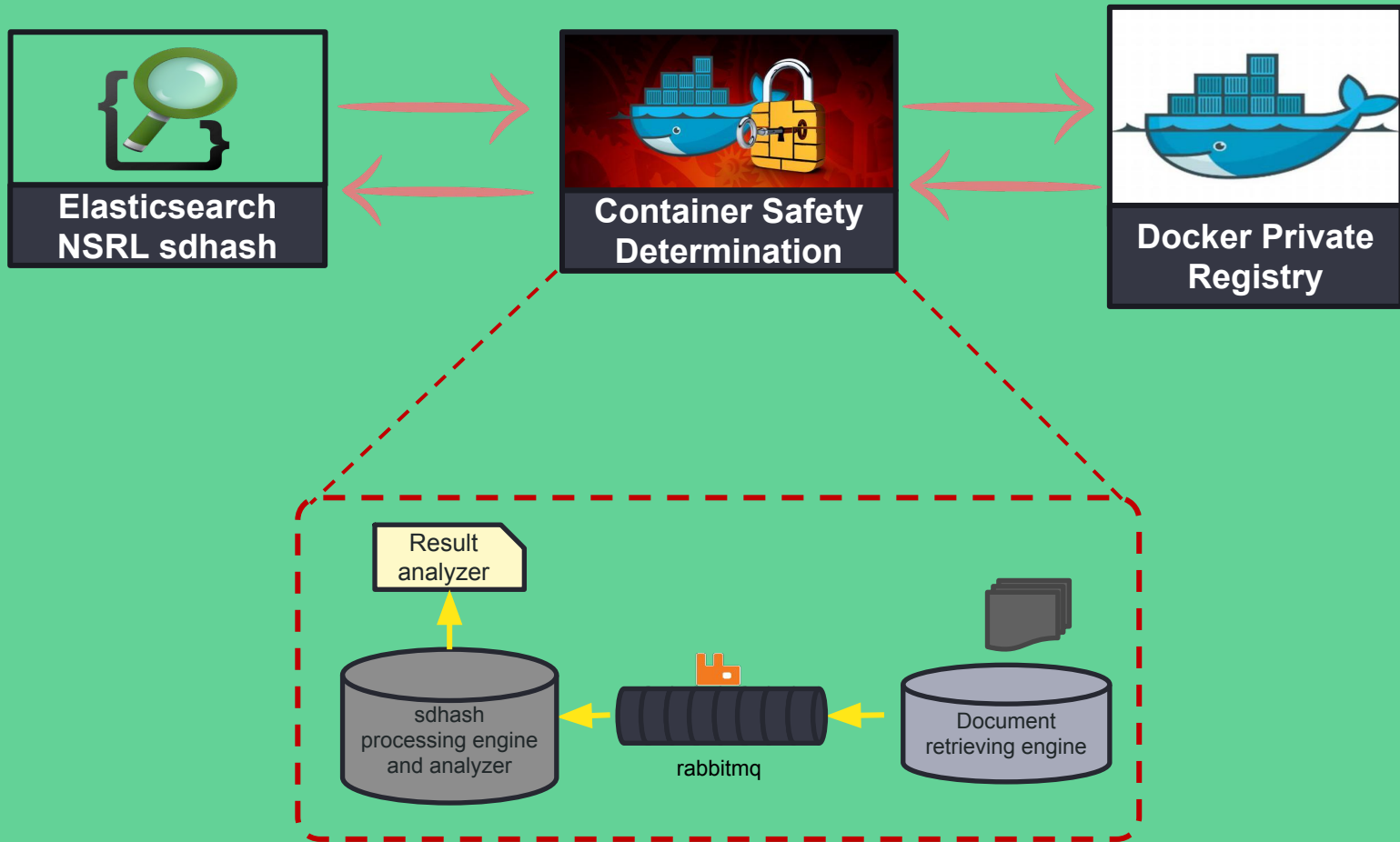


```
FAAAAAICAAARUAAAAAAAAQAAAAACEAAQQCAAQAQAAKAAEAQAAJAAAgEAAgIIABIQA  
AAAAoEBAgAAABABAAAAAAAAEAgAAgAAACAIAAAEAAAAAAAAABAAACAAAAAAgDAA  
AAAAAAAAAgAAABAAgAAAACQAAIAAAAAAAAAAAAAACAAIEAAAAAgwAAAAQAAAAAQAA  
QBAEAAKAACAAAAAAgDAAAAAAAAA
```

Technologies

- SDHash
- Docker
- ElasticSearch
- Rabbitmq





Why sdhash?

crypto-hashes not suited for similarity comparison:

- small change in a file will generate completely different crypto-hash.

sdhash: Similarity digest hash

Reduces the file size to 2-3% of the original

sdhash tool compares 2 files and generates a similarity index 0-100

Example: sdhashes for two text files (~14KB) with mostly same content except a few sentences deleted from file2.

file1 hash

```
1 sdbf:03:9:file1.txt:14030:shal:256:5:7ff:160:2:77:CIEVAjQGU
2 EAHAAStSBZowLB3B4FAZDQAADKGwFQNH11AalkDqmAA8AQOMJjIHQGSTAnA
3 cmBctDQDtJIwCoGInUkSDybgowgdEwQHf11BiIBYMyjKFYspoDzsAlxrFEL
4 iiLogDUQAQOEskIQAIBAAyCgQGaDBKAI0+LgAYE5DA5mQmmZgwAgIBAEmA9
5 AAMiLBAU0g0QoAiGrIg7oSJtG0boBIQEEJAawCQ4112ENeAUgna0yU4ICzq
6 ECVhowoxVUwJ2CAAXBkdqNyiiAoQA/E8FlrYQF4QRAAQqi21qDEtwSLKEAZ
7 ITuzD6gooJIr+KD2QSQIBzSFkasBAqQACz4AERqggAAAAEIALAUNoGCIAACT
8 AAAICCGWHIICwAQkgAkAEBBAQBAFEAgQagAAAAABABEGAMgBYiAAYAGAQAk
9 IAKFAQCBCwIgaUCBLDICAGCgIAURAAEAkgAIAEBAAMVxMALBQBAAQQCIigF
10 CACAAICAAgQEHFREAgggIAAgwEghgABEAgQihSAAFAoLATUAAADAASEAgABQ
11 BIUiABAiQAAEIAIDAAAAIggCJCARABgAJSApCaAwIAIBggwAAAUAUAINAI
12 ABKEBkAIAMAEcQEMikNGpAoCFAkBCQAAAFBgAgAAwAB2BAAUEgA6AADIVA
13 BEwMAIyEBQQAoAAAYAAggaEBGiA=
```

file2 hash

```
1 sdbf:03:9:file2.txt:13607:shal:256:5:7ff:160:2:71:CIEVAjQGF
2 OEAGAAStSBZowLB2BYFAZDQAADKGwFQJH11AalkDqmAA9AQOMJjIHwGSTAn
3 gcqBctDQDtJIwCoGInUkSDzaioxgdAwQHf11BiIBYMyjKFYspoDzsAlxrFE
4 KiiLogDUQAQOECKiQAIBAAyCgQGaDBKAIK+LgCYG5DA5mQimZgwBggBAEuA
5 9AAMiDBAU0g0QoAiGrIg7pSBtG0aoBIUEEpAAxCA4112ENeAEgna0zU4IKz
6 qECVhowoxVUwJOCEAXBkdqNyiiAoQA/E8Flj4QF4QRBBQQg21iDEtwSLKEE
7 ZATuzD6gooJIp+KD2QSSIBzCFOasBAqSACz8AERoggAAAAEIALAUNoGCIAAC
8 TAAAIcCGHHIICwAQkgAkAEBBAQBAFEAgQagAAAAABABEGAEGBYiAAYAGAQA
9 KIAKFAQCBCgAAECBLDICAGCgIAQRAAEAkGAIABEBAAMVRMALBQBAAQQCIig
10 FCACAAICAAgQEHFREAgggIAAgwEAhgAAEAgQihSAAFAoJATEAADAASEAgAB
11 QBAUiABAiAAAEIAIDAAAAAggCJCARABgAJSApCaAwIAIBAgwAAAUAUAINAI
12 EAAKEBkAIAMAEcQEMikNGpAgCFAkBCAAAAFBgAgAAwAB2AAAUEgAYAADIVA
13 IBewMAIQEBQQAoAAAYAAggaAGiA=
```

sdhash similarity index: 96

Why NSRL?

- Provides SDHash
- > 40,000,000 hashes
- Quarterly Release

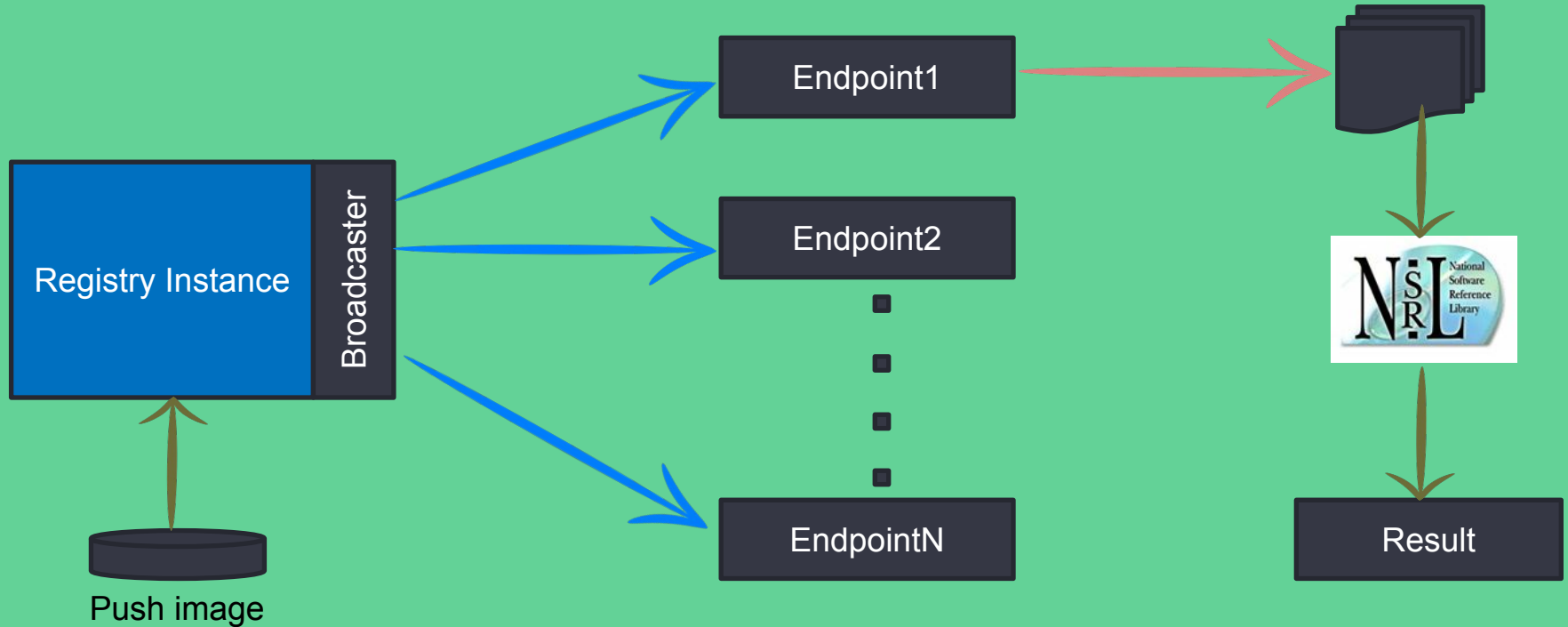
Goal: utilize the precomputed hashes provided by NSRL

Docker Registry

- Manage docker images
- Multiple storage backends (local, S3, Ceph, ...)
- Supports notifications to update its status to endpoints



Registry Notifications and sdhash



Sprint 1 tasks

- Private Docker registry
- Script to pull and extract files from Docker registry
- Setup private environment
- Sdhash compilation and installation
- Script to calculate sdhash of files extracted

Tasks done and in progress

In Progress Sprint 1

Understanding about elasticsearch
(basic setup and working)

✓ 2/3



RG

Understanding of rabbitmq



JM



Understanding the NSRL database



✓ 3/5



RG

Add a card...

Sprint 1 Done - Demo Feb 12

(4)Completing video tutorials on
Docker[4]

✓ 1

JM



RG

(6)Parse docker image to retrieve
files[6]

✓ 1/1

JM

(2)Calculate sdhash for each file in
docker image[2]

✓ 1/1

JM

(10)Setting up private docker
registry[10]

✓ 1

✓ 3/3



(4)Working with sdhash[4]

✓ 3/3

JM



RG

Add a card...

Sprint 2 [NoBurn]

Setup environment on OpenStack

✓ 0/4

JM



Write python Script

✓ 0/3

JM



RG

Tutorials of RabbitMQ

✓ 1

JM



RG

Add a card...

In Progress Sprint 2 [NoBurn]

(10)Configuring docker-registry for
notifications and setting up REST API
server to receive the notifications



Add a card...

Next sprint

- First run of data with elasticsearch
- Come up with the sdhash threshold value ranges
- Implement REST APIs for endpoints
- Integrate image push with hash-calculation

Thank you
&
Q/A