

OCI Image Format

Glyn Normington

May 28, 2019

This document provides a formal model of some aspects of the *OCI image format*. It describes the objects which make up an OCI repository and how content addressing is used to refer to these objects.

Contents

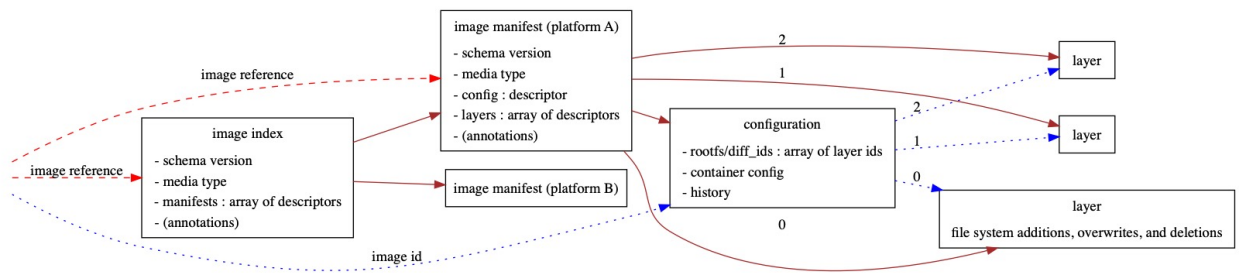
1	Introduction	1
2	Content Addressing	2
3	Change Sets (“Layers”)	4
3.1	Image Repositories (part 1 of 4)	5
4	Image Configurations	6
4.1	Image Repositories (part 2 of 4)	6
5	Image Manifests	8
5.1	Image Repositories (part 3 of 4)	8
6	Image Indices	10
6.1	Image Repositories (part 4 of 4)	10
A	References	12
B	Z Notation	13

1 Introduction

Docker Inc. ([1]) introduced *container images* and *registries* to hold them and these were later standardised as part of the Open Container Initiative ([2]). The reader is assumed to have a basic understanding of how images are used.

Registries are named collections of *repositories*. The relationship is described formally in *Image Registries* ([3]).

This document models how image are represented and stored in repositories. This is based on the *OCI Image Format Specification* ([4]). It does not cover all details, but aims to provide an overview of how images are structured as layers and gathered into a repository. The following diagram shows the main objects that are modelled.



The Z specification language is used to describe the model, but sufficient English text is also provided that readers who do not know Z should be able to get a reasonable idea of what's going on. The appendix contains a summary of the Z notation. For more information about Z, please consult the Z Manual ([5]). The model was type checked using **fuzz** ([6]).

2 Content Addressing

Before we describe the objects which make up a repository, we model how objects refers to each other using *content addressing*. In content addressing, a collision resistant function¹ is applied to the content of an object to produce a content address which can be used to refer to the object.

OCI uses a form of *content address*, known as a *digest*, which is a cryptographic hash of the *compressed* content of an object.

[*Digest*]

OCI also uses another form of content address, known as an *id*, which is a cryptographic hash of the *uncompressed* content of an object.

[*Id*]

Content addressable storage (CAS) in a repository has a collection of objects uniquely identified by digest and, sometimes, by id. Each object has a “raw” (possibly compressed) size in bytes.

$CAS[Object]$	_____
$deref : Digest \rightsquigarrow Object$	
$identify : Id \rightsquigarrow Object$	
$size : Digest \rightarrow \mathbb{N}$	
$\text{ran } identify \subseteq \text{ran } deref$	
$\text{dom } deref = \text{dom } size$	

As well as addressing objects by digest, some intra-repository references use *descriptors*. A descriptor augments a digest with extra descriptive information: media type, URLs, and strings (used in annotations), which are not modelled further.

[*MediaType, URL, String*]

A media type is used to distinguish between certain types of objects in a repository. URLs indicate where certain objects are hosted in a network. Annotations provide arbitrary metadata for certain objects.

¹collision resistant functions are designed so that it is extremely unlikely that distinct inputs produce the same output

A descriptor is a way of referring to an object in the same repository using a digest, a media type, a “raw” size, and zero or more URLs and annotations.

<i>Descriptor</i>
<i>digest</i> : <i>Digest</i>
<i>mediaType</i> : <i>MediaType</i>
<i>size</i> : \mathbb{N}
<i>urls</i> : seq <i>URL</i>
<i>annotations</i> : <i>String</i> \rightarrow <i>String</i>

We will model a repository as a collection of content addressable objects, which we will introduce in turn, along with their media types. These objects are change sets, image configurations, image manifests, and image indices.

3 Change Sets (“Layers”)

A changeset, or “layer”, is a means of creating a root file system for a container based on an image.

File system paths and files (including directories) are modelled in terms of their usage, but not their implementation.

$$[FSPath, File]$$

We introduce a simple model of a file system which will be sufficient to show how the file system of an image is built from layers.

A simplified file system is a collection of files indexed by file system path.

$$\boxed{\begin{array}{l} FS \\ fs : FSPath \rightarrow File \end{array}}$$

No mention is made of the hierarchical nature of file system paths. File modes and attributes and symbolic and hard links are not modelled.

An empty file system has no files.

$$EmptyFS \triangleq [FS \mid fs = \emptyset]$$

An *image layer file system changeset*, informally known as a “layer”, describes a particular type of modification to a file system where some files are first of all removed and then files are added or replaced.

$$\boxed{\begin{array}{l} ChangeSet \\ remove : \mathbb{P} FSPath \\ add : FSPath \rightarrow File \end{array}}$$

There is a media type for changesets which will be used in descriptors that reference changesets.

$$\mid LAYER : MediaType$$

A changeset is applied to a file system by removing files according to the changeset and then adding (or replacing) file according to the changeset.

<i>ChangeSetApply</i>	_____
ΔFS	
$c? : ChangeSet$	
$fs' = (c?.remove \triangleleft fs) \oplus c?.add$	

Applying a changeset to an empty file system results in a file system containing of just the additions from the changeset. So a sequence of changesets can be applied in turn starting with an empty file system to produce the “root” file system for a particular *container*.

The remainder of this document does not concern itself with combining changesets, except that it will show how an image refers to a sequence of changesets which can, as we have just seen, be used to form the root file system for each container constructed from the image.

3.1 Image Repositories (part 1 of 4)

We can now start to describe image repositories in terms of change sets. Other objects in an image repository will be described later.

An image repository has a content addressable collection of change sets.

<i>RepositoryChangeSets</i>	_____
$changesets : CAS[ChangeSet]$	

This will make more sense once we start to add other objects to the repository.

4 Image Configurations

Container configurations and (image history) logs are used by container configurations, but are not modelled further.

$[ContainerCfg, Log]$

There is a media type for image configurations which will be used in descriptors that reference image configurations.

| $CONFIG : MediaType$

An image configuration refers to the changesets (or “layers”, of which there must be at least one) of an image and records the container configuration and history of the image.

Config _____

$changesets : seq_1 Id$

$ccfg : ContainerCfg$

$history : seq Log$

4.1 Image Repositories (part 2 of 4)

We now add configurations to our description of an image repository and see how configurations relate to change sets. Other objects in an image repository will be described later.

An image repository has a content addressable collection of change sets and configurations addressed by id.

RepositoryConfigsAndLayers _____

RepositoryChangeSets

$configs : CAS[Config]$

$ids : \mathbb{P} Id$

$digests : \mathbb{P} Digest$

$\forall cfg : ran\ configs.deref \bullet ran\ cfg.changesets \subseteq dom\ changesets.identify$
 $\langle dom\ changesets.identify, dom\ configs.identify \rangle \text{ partition } ids$

$\exists other : \mathbb{P} Digest \bullet$

$\langle dom\ changesets.deref, dom\ configs.deref, other \rangle \text{ partition } digests$

All configurations refer to valid changeset ids.

Ids and digests are partitioned between configurations and changesets (and, in the case of digests, other objects still to be introduced). That is, each id and each digest refers to at most one object in content addressable storage.

The next step is to introduce the central object in a repository: the *image manifest*.

5 Image Manifests

The central object in an image repository is the *image manifest* which describes a single image.

Image manifests have schema versions, which are not modelled further.

[*SchemaVersion*]

There is a media type for image manifests which is used both in image manifests and in descriptors that reference them.

| *IMAGE_MANIFEST* : *MediaType*

An image manifest refers to an image configuration and at least one changeset (or “layer”).

<i>Manifest</i> <i>schemaVersion</i> : <i>SchemaVersion</i> <i>mediaType</i> : <i>MediaType</i> <i>config</i> : <i>Descriptor</i> <i>changesets</i> : seq ₁ <i>Descriptor</i> <i>annotations</i> : <i>String</i> \leftrightarrow <i>String</i> <hr/> <i>mediaType</i> = <i>IMAGE_MANIFEST</i>
--

5.1 Image Repositories (part 3 of 4)

We extend our model of an image repository to include image manifests.

$ \begin{array}{l} \textit{RepositoryManifestsConfigsAndLayers} \\ \textit{RepositoryConfigsAndLayers} \\ \textit{manifests} : \textit{CAS}[\textit{Manifest}] \\ \hline \textit{manifests.identify} = \emptyset \\ \exists \textit{other} : \mathbb{P} \textit{Digest} \bullet \\ \quad \langle \textit{dom changesets.deref}, \textit{dom configs.deref}, \\ \quad \quad \textit{dom manifests.deref}, \textit{other} \rangle \textbf{partition digests} \\ \forall \textit{man} : \textit{ran manifests.deref} \bullet \\ \quad \textit{man.mediaType} = \textit{IMAGE_MANIFEST} \wedge \\ \quad \textit{man.config.mediaType} = \textit{CONFIG} \wedge \\ \quad \textit{man.config.digest} \in \textit{dom configs.deref} \wedge \\ \quad \textit{man.config.size} = \textit{configs.size man.config.digest} \wedge \\ \quad (\forall \textit{cs} : \textit{ran man.changesets} \bullet \\ \quad \quad \textit{cs.mediaType} = \textit{LAYER} \wedge \\ \quad \quad \textit{cs.digest} \in \textit{dom changesets.deref} \wedge \\ \quad \quad \textit{cs.size} = \textit{changesets.size cs.digest}) \end{array} $
--

Manifests do not have ids.

Also, content addresses in the form of digests are partitioned between changesets, configs, manifests, and other objects (soon to be introduced). That is, each digest refers to at most one object in content addressable storage.

Also, all manifests have an appropriate media type and valid configuration and layer descriptors.

6 Image Indices

To complete our description of image repositories, we add another object, the *image index*. The main purpose of an image index is to enable a repository to contain multiple versions of an image for multiple platforms.

Image indices use extended descriptors which include platform information, which is not modelled further.

[*Platform*]

There is a media type for image indices which is used both in image indices and in descriptors that reference them.

| *IMAGE_INDEX* : *MediaType*

A manifest descriptor is a descriptor with an additional property denoting a platform.

<i>ManifestDescriptor</i> <i>Descriptor</i> <i>platform</i> : <i>Platform</i> <i>mediaType</i> ∈ { <i>IMAGE_MANIFEST</i> , <i>IMAGE_INDEX</i> }
--

A manifest descriptor may refer to image manifests or image indices.

An image index refers to zero or more image manifests and/or image indices.

<i>Index</i> <i>schemaVersion</i> : <i>SchemaVersion</i> <i>mediaType</i> : <i>MediaType</i> <i>manifests</i> : seq <i>ManifestDescriptor</i> <i>annotations</i> : <i>String</i> → <i>String</i> <i>mediaType</i> = <i>IMAGE_INDEX</i>

6.1 Image Repositories (part 4 of 4)

We can now complete our model of image repositories.

An image repository is a content addressable collection of change sets, configurations, image manifests, and image indices.

<i>Repository</i>
<i>RepositoryManifestsConfigsAndLayers</i> <i>indices</i> : CAS[Index]
<i>indices.identify</i> = \emptyset $\langle \text{dom } \textit{changesets.deref}, \text{dom } \textit{configs.deref}, \text{dom } \textit{manifests.deref}, \text{dom } \textit{indices.deref} \rangle \text{ partition } \textit{digests}$ $\forall \textit{idx} : \text{ran } \textit{indices.deref} \bullet$ $\textit{idx.mediaType} = \textit{IMAGE_INDEX} \wedge$ $(\forall \textit{man} : \text{ran } \textit{idx.manifests} \bullet$ $(\textit{man.digest} \in \text{dom } \textit{manifests.deref} \wedge$ $\textit{man.size} = \textit{manifests.size } \textit{man.digest}) \vee$ $(\textit{man.digest} \in \text{dom } \textit{indices.deref} \wedge$ $\textit{man.size} = \textit{indices.size } \textit{man.digest}))$ $\exists \textit{next} : \textit{Index} \leftrightarrow \textit{Index} \mid \forall i, j : \textit{Index} \mid i \mapsto j \in \textit{next} \bullet$ $(\exists m : \text{ran } \textit{i.manifests} \bullet \textit{indices.deref } m.\textit{digest} = j) \bullet$ $\textit{next}^+ \cap \text{id } \textit{Index} = \emptyset$

Image indices do not have ids.

Also, content addresses in the form of digests are partitioned between change-sets, configs, manifests, and indices. That is, each digest refers to at most one object in content addressable storage.

Also, all image indices have an appropriate media type and valid manifest descriptors (referring to image manifests and/or other image indices).

Also, there are no cycles of image indices.

A References

- [1] Various authors, *Docker Inc.*, <https://www.docker.com/>.
- [2] Various authors, *Open Container Initiative*,
<https://www.opencontainers.org/>.
- [3] Glyn Normington, *Image Registries*,
<https://github.com/glyn/riff-specs/raw/master/distribution/distribution.pdf>.
- [4] Various authors, *OCI Image Format Specification*,
<https://github.com/opencontainers/image-spec/blob/master/spec.md>.
- [5] Mike Spivey, *The Z Manual*, <https://www.cs.umd.edu/mvz/handouts/z-manual.pdf>.
- [6] Mike Spivey, *The fuzz type-checker for Z*,
<https://bitbucket.org/Spivey/fuzz>.

B Z Notation

Numbers:

\mathbb{N} Natural numbers $\{0,1,\dots\}$

Propositional logic and the schema calculus:

$\dots \wedge \dots$	And	$\langle\langle \dots \rangle\rangle$	Free type injection
$\dots \vee \dots$	Or	$[\dots]$	Given sets
$\dots \Rightarrow \dots$	Implies	$\prime, ?, !, 0 \dots 9$	Schema decorations
$\forall \dots \mid \dots \bullet \dots$	For all	$\dots \vdash \dots$	theorem
$\exists \dots \mid \dots \bullet \dots$	There exists	$\theta \dots$	Binding formation
$\dots \setminus \dots$	Hiding	$\lambda \dots$	Function definition
$\dots \hat{=} \dots$	Schema definition	$\mu \dots$	Mu-expression
$\dots == \dots$	Abbreviation	$\Delta \dots$	State change
$\dots ::= \dots \mid \dots$	Free type definition	$\Xi \dots$	Invariant state change

Sets and sequences:

$\{\dots\}$	Set	$\dots \setminus \dots$	Set difference
$\{\dots \mid \dots \bullet \dots\}$	Set comprehension	$\bigcup \dots$	Distributed union
$\mathbb{P} \dots$	Set of subsets of	$\# \dots$	Cardinality
\emptyset	Empty set	$\dots \subseteq \dots$	Subset
$\dots \times \dots$	Cartesian product	$\dots \subset \dots$	Proper subset
$\dots \in \dots$	Set membership	$\dots \text{ partition } \dots$	Set partition
$\dots \notin \dots$	Set non-membership	seq	Sequences
$\dots \cup \dots$	Union	$\langle \dots \rangle$	Sequence
$\dots \cap \dots$	Intersection	$\text{disjoint } \dots$	Disjoint sequence of sets

Functions and relations:

$\dots \Leftrightarrow \dots$	Relation	\dots^*	Reflexive-transitive closure
$\dots \rightarrow \dots$	Partial function	$\dots \langle \dots \rangle$	Relational image
$\dots \rightarrow \dots$	Total function	$\dots \oplus \dots$	Functional overriding
$\dots \mapsto \dots$	Partial injection	$\dots \triangleleft \dots$	Domain restriction
$\dots \mapsto \dots$	Injection	$\dots \triangleright \dots$	Range restriction
$\text{dom } \dots$	Domain	$\dots \triangleleft \dots$	Domain subtraction
$\text{ran } \dots$	Range	$\dots \triangleright \dots$	Range subtraction
$\dots \mapsto \dots$	maplet		
$\dots \sim \dots$	Relational inverse		

Axiomatic descriptions:

<i>Declarations</i>
<i>Predicates</i>

Schema definitions:

<i>SchemaName</i>
<i>Declaration</i>
<i>Predicates</i>

Decorations:

<i>(undecorated)</i>	Input state	?	Input parameter
\prime	Output state	!	Output parameter