

EXTENDS *FiniteSets*, *Sequences*, *Naturals*, *TLC*

```
* Licensed to the Apache Software Foundation (ASF) under one
* or more contributor license agreements. See the NOTICE file
* distributed with this work for additional information
* regarding copyright ownership. The ASF licenses this file
* to you under the Apache License, Version 2.0 (the
* "License"); you may not use this file except in compliance
* with the License. You may obtain a copy of the License at
*
*   http://www.apache.org/licenses/LICENSE-2.0
*
* Unless required by applicable law or agreed to in writing, software
* distributed under the License is distributed on an "AS IS" BASIS,
* WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
* See the License for the specific language governing permissions and
* limitations under the License.
```

This defines the YARN registry in terms of operations on sets of records.

Every registry entry is represented as a record containing both the path and the data.

It assumes that

1. operations on this set are immediate.
2. selection operations (such as \forall and \exists are atomic)
3. changes are immediately visible to all other users of the registry.
4. This clearly implies that changes are visible in the sequence in which they happen.

A multi-server Zookeeper-based registry may not meet all those assumptions

1. changes may take time to propagate across the *ZK* quorum, hence changes cannot be considered immediate from the perspective of other registry clients. (assumptions (1) and (3)).
2. Selection operations may not be atomic. (assumption (2)).

Operations will still happen in the order received by the elected *ZK* master

A stricter definition would try to state that all operations are eventually true excluding other changes happening during a sequence of action. This is left as an exercise for the reader.

The specification also omits all coverage of the permissions policy.

CONSTANTS

<i>PathChars</i> ,	the set of valid characters in a path
<i>Paths</i> ,	the set of all possible valid paths
<i>Data</i> ,	the set of all possible sequences of bytes
<i>Address</i> ,	the set of all possible address n-tuples
<i>Addresses</i> ,	the set of all possible address instances
<i>Endpoints</i> ,	the set of all possible endpoints

<i>PersistPolicies</i> ,	the set of persistence policies
<i>ServiceRecords</i> ,	all service records
<i>Registries</i> ,	the set of all possible registries
<i>BindActions</i> ,	all possible put actions
<i>DeleteActions</i> ,	all possible delete actions
<i>PurgeActions</i> ,	all possible purge actions
<i>MknodeActions</i>	all possible mkdir actions

the registry
VARIABLE *registry*

Sequence of actions to apply to the registry
VARIABLE *actions*

Tuple of all variables.

$vars \triangleq \langle registry, actions \rangle$

Persistence policy
 $PersistPolicySet \triangleq \{$
 "","
 "permanent",
 "application",
 "application-attempt",
 "container"
 $\}$

Undefined; field not present. PERMANENT is implied.
persists until explicitly removed
persists until the application finishes
persists until the application attempt finishes
persists until the container finishes

Type invariants.
 $TypeInvariant \triangleq$
 $\wedge \forall p \in PersistPolicies : p \in PersistPolicySet$

An Entry is defined as a path, and the actual data which it contains.

By including the path in an entry, we avoid having to define some function mapping $Path \rightarrow entry$. Instead a registry can be defined as a set of *RegistryEntries* matching the validity criteria.

$RegistryEntry \triangleq [$

The path to the entry

$path : Paths,$

the data in the entry

$data : Data$

]

An endpoint in a service record

$Endpoint \triangleq [$

API of the endpoint: some identifier

$api : STRING,$

A list of address n-tuples

$addresses : Addresses$

]

Attributes are the set of all string to string mappings

$Attributes \triangleq [$

$STRING \mapsto STRING$

]

A service record

$ServiceRecord \triangleq [$

ID – used when applying the persistence policy

$yarn_id : STRING,$

the persistence policy

$yarn_persistence : PersistPolicySet,$

A description

$description : STRING,$

A set of endpoints

$external : Endpoints,$

$Endpoints$ intended for use internally

$internal : Endpoints,$

Attributes are a function

$attributes : Attributes$

]

Action Records

```

putAction  $\triangleq$  [
  type : "put",
  record : ServiceRecord
]

deleteAction  $\triangleq$  [
  type : "delete",
  path : STRING,
  recursive : BOOLEAN
]

purgeAction  $\triangleq$  [
  type : "purge",
  path : STRING,
  persistence : PersistPolicySet
]

mkNodeAction  $\triangleq$  [
  type : "mknode",
  path : STRING,
  parents : BOOLEAN
]

```

Path operations

Parent is defined for non empty sequences

$\text{parent}(\text{path}) \triangleq \text{SubSeq}(\text{path}, 1, \text{Len}(\text{path}) - 1)$

$\text{isParent}(\text{path}, c) \triangleq \text{path} = \text{parent}(c)$

Registry Access Operations

Lookup all entries in a registry with a matching path

$\text{resolve}(\text{Registry}, \text{path}) \triangleq \forall \text{entry} \in \text{Registry} : \text{entry.path} = \text{path}$

A path exists in the registry iff there is an entry with that path

$\text{exists}(\text{Registry}, \text{path}) \triangleq \text{resolve}(\text{Registry}, \text{path}) \neq \{\}$

A parent entry, or an empty set if there is none

$\text{parentEntry}(\text{Registry}, \text{path}) \triangleq \text{resolve}(\text{Registry}, \text{parent}(\text{path}))$

A root path is the empty sequence

$$isRootPath(path) \triangleq path = \langle \rangle$$

The root entry is the entry whose path is the root path

$$isRootEntry(entry) \triangleq entry.path = \langle \rangle$$

A path p is an ancestor of another path d if they are different, and the path d starts with path p

$$\begin{aligned} isAncestorOf(path, d) &\triangleq \\ &\wedge path \neq d \\ &\wedge \exists k : SubSeq(d, 0, k) = path \end{aligned}$$

$$\begin{aligned} ancestorPathOf(path) &\triangleq \\ &\forall a \in Paths : isAncestorOf(a, path) \end{aligned}$$

The set of all children of a path in the registry

$$children(R, path) \triangleq \forall c \in R : isParent(path, c.path)$$

A path has children if the *children()* function does not return the empty set

$$hasChildren(R, path) \triangleq children(R, path) \neq \{\}$$

Descendant: a child of a path or a descendant of a child of a path

$$descendants(R, path) \triangleq \forall e \in R : isAncestorOf(path, e.path)$$

Ancestors: all entries in the registry whose path is an entry of the path argument

$$ancestors(R, path) \triangleq \forall e \in R : isAncestorOf(e.path, path)$$

The set of entries that are a path and its descendants

$$\begin{aligned} pathAndDescendants(R, path) &\triangleq \\ &\vee \forall e \in R : isAncestorOf(path, e.path) \\ &\vee resolve(R, path) \end{aligned}$$

For validity, all entries must match the following criteria

$$\begin{aligned} validRegistry(R) &\triangleq \\ &\text{there can be at most one entry for a path.} \\ &\wedge \forall e \in R : Cardinality(resolve(R, e.path)) = 1 \\ &\text{There's at least one root entry} \\ &\wedge \exists e \in R : isRootEntry(e) \\ &\text{an entry must be the root entry or have a parent entry} \\ &\wedge \forall e \in R : isRootEntry(e) \vee exists(R, parent(e.path)) \\ &\text{If the entry has data, it must be a service record} \\ &\wedge \forall e \in R : (e.data = \langle \rangle \vee e.data \in ServiceRecords) \end{aligned}$$

Registry Manipulation

An entry can be put into the registry iff its parent is present or it is the root entry

$$\begin{aligned} canBind(R, e) &\triangleq \\ &isRootEntry(e) \vee exists(R, parent(e.path)) \end{aligned}$$

'bind()' adds/replaces an entry if permitted

$$\begin{aligned} bind(R, e) &\triangleq \\ &\wedge canBind(R, e) \\ &\wedge R' = (R \setminus resolve(R, e.path)) \cup \{e\} \end{aligned}$$

mknode() adds a new empty entry where there was none before, iff

- the parent *exists*
- it meets the requirement for being “bindable”

$$\begin{aligned} mknodeSimple(R, path) &\triangleq \\ &LET \ record \triangleq [path \mapsto path, data \mapsto \langle \rangle] \\ &IN \quad \vee exists(R, path) \\ &\quad \vee (exists(R, parent(path)) \wedge canBind(R, record) \wedge (R' = R \cup \{record\})) \end{aligned}$$

For all parents, the *mknodeSimple*() criteria must apply. This could be defined recursively, though as TLA+ does not support recursion, an alternative is required

Because this specification is declaring the final state of a operation, not the implemental, all that is needed is to describe those parents.

It declares that the *mkdirSimple* state applies to the path and all its parents in the set *R'*

$$\begin{aligned} mknodeWithParents(R, path) &\triangleq \\ &\wedge \forall p2 \in ancestors(R, path) : mknodeSimple(R, p2) \\ &\wedge mknodeSimple(R, path) \end{aligned}$$

$$\begin{aligned} mknode(R, path, recursive) &\triangleq \\ &IF \ recursive \ THEN \ mknodeWithParents(R, path) \ ELSE \ mknodeSimple(R, path) \end{aligned}$$

Deletion is set difference on any existing entries

$$\begin{aligned} simpleDelete(R, path) &\triangleq \\ &\wedge \neg isRootPath(path) \\ &\wedge children(R, path) = \{\} \\ &\wedge R' = R \setminus resolve(R, path) \end{aligned}$$

Recursive delete: neither the path or its descendants exists in the new registry

$$\begin{aligned} recursiveDelete(R, path) &\triangleq \\ &\text{Root path: the new registry is the initial registry again} \\ &\wedge isRootPath(path) \Rightarrow R' = \{[path \mapsto \langle \rangle, data \mapsto \langle \rangle]\} \end{aligned}$$

Any other entry: the new registry is a set with any existing entry for that path is removed, and the new entry added

$$\wedge \neg isRootPath(path) \Rightarrow R' = R \setminus (resolve(R, path) \cup descendants(R, path))$$

Delete operation which chooses the recursiveness policy based on an argument

$$delete(R, path, recursive) \triangleq$$

IF *recursive* THEN *recursiveDelete*(*R*, *path*) ELSE *simpleDelete*(*R*, *path*)

Purge ensures that all entries under a path with the matching *ID* and policy are not there afterwards

$$purge(R, path, id, persistence) \triangleq$$

$$\wedge (persistence \in PersistPolicySet)$$

$$\wedge \forall p2 \in pathAndDescendants(R, path) :$$

$$(p2.attributes["yarn:id"] = id \wedge p2.attributes["yarn:persistence"] = persistence)$$

$$\Rightarrow recursiveDelete(R, p2.path)$$

resolveRecord() resolves the record at a path or fails.

It relies on the fact that if the cardinality of a set is 1, then the CHOOSE operator is guaranteed to return the single entry of that set, iff the choice predicate holds.

Using a predicate of TRUE, it always succeeds, so this function selects the sole entry of the resolve operation.

$$resolveRecord(R, path) \triangleq$$

$$LET \ l \triangleq resolve(R, path) IN$$

$$\wedge Cardinality(l) = 1$$

$$\wedge CHOOSE \ e \in l : TRUE$$

The specific action of putting an entry into a record includes validating the record

$$validRecordToBind(path, record) \triangleq$$

The root entry must have permanent persistence

$$isRootPath(path) \Rightarrow (record.attributes["yarn:persistence"] = "permanent"$$

$$\vee record.attributes["yarn:persistence"] = "")$$

Binding a service record involves validating it then putting it in the registry marshalled as the data in the entry

$$bindRecord(R, path, record) \triangleq$$

$$\wedge validRecordToBind(path, record)$$

$$\wedge bind(R, [path \mapsto path, data \mapsto record])$$

The action queue can only contain one of the sets of action types, and by giving each a unique name, those sets are guaranteed to be disjoint

$$\begin{aligned}
\text{QueueInvariant} &\triangleq \\
&\wedge \forall a \in \text{actions} : \\
&\quad \vee (a \in \text{BindActions} \wedge a.type = \text{"bind"}) \\
&\quad \vee (a \in \text{DeleteActions} \wedge a.type = \text{"delete"}) \\
&\quad \vee (a \in \text{PurgeActions} \wedge a.type = \text{"purge"}) \\
&\quad \vee (a \in \text{MknodeActions} \wedge a.type = \text{"mknode"})
\end{aligned}$$

Applying queued actions

$$\begin{aligned}
\text{applyAction}(R, a) &\triangleq \\
&\quad \vee (a \in \text{BindActions} \wedge \text{bindRecord}(R, a.path, a.record)) \\
&\quad \vee (a \in \text{MknodeActions} \wedge \text{mknode}(R, a.path, a.recursive)) \\
&\quad \vee (a \in \text{DeleteActions} \wedge \text{delete}(R, a.path, a.recursive)) \\
&\quad \vee (a \in \text{PurgeActions} \wedge \text{purge}(R, a.path, a.id, a.persistence))
\end{aligned}$$

Apply the first action in a list and then update the actions

$$\begin{aligned}
\text{applyFirstAction}(R, a) &\triangleq \\
&\quad \wedge \text{actions} \neq \langle \rangle \\
&\quad \wedge \text{applyAction}(R, \text{Head}(a)) \\
&\quad \wedge \text{actions}' = \text{Tail}(a)
\end{aligned}$$

$$\text{Next} \triangleq \text{applyFirstAction}(\text{registry}, \text{actions})$$

All submitted actions must eventually be applied.

$$\text{Liveness} \triangleq \Diamond(\text{actions} = \langle \rangle)$$

The initial state of a registry has the root entry.

$$\begin{aligned}
\text{InitialRegistry} &\triangleq \text{registry} = \{ \\
&\quad [\text{path} \mapsto \langle \rangle, \text{data} \mapsto \langle \rangle] \\
&\}
\end{aligned}$$

The valid state of the “registry” variable is defined as Via the *validRegistry* predicate

$$\text{ValidRegistryState} \triangleq \text{validRegistry}(\text{registry})$$

The initial state of the system

$$\begin{aligned}
\text{InitialState} &\triangleq \\
&\quad \wedge \text{InitialRegistry} \\
&\quad \wedge \text{ValidRegistryState} \\
&\quad \wedge \text{actions} = \langle \rangle
\end{aligned}$$

The registry has an initial state, the series of state changes driven by the actions, and the requirement that it does act on those actions.

$$\begin{aligned} RegistrySpec &\triangleq \\ &\wedge InitialState \\ &\wedge \Box [Next]_{vars} \\ &\wedge Liveness \end{aligned}$$

Theorem: For all operations from that initial state, the registry state is still valid

THEOREM $InitialState \Rightarrow \Box ValidRegistryState$

Theorem: for all operations from that initial state, the type invariants hold

THEOREM $InitialState \Rightarrow \Box TypeInvariant$

Theorem: the queue invariants hold

THEOREM $InitialState \Rightarrow \Box QueueInvariant$
