

Add Service Container type to NodeManager in YARN

1. Motivation

From our work to support running OpenMPI on YARN (<https://issues.apache.org/jira/browse/MAPREDUCE-2911>), we found that it's important to have framework specific daemon process manage the tasks on each node directly. The daemon process, most likely similar in other frameworks as well, provides critical services to tasks running on that node(for example “wireup”, spawn user process in large numbers at once etc) . In YARN, it's hard, if not possible, to have the those processes to be managed by YARN.

We propose to extend the container model on NodeManager side to support “Service Container” to run/manage such framework daemon/services process. We believe this is very useful to other application framework developers as well.

2. Idea

Currently in YARN during job execution time, as far as YARN is concerned, there is no distinction among all the processes/containers on each node, all are directly managed by each node manager, regardless the process roles each one plays in the application.

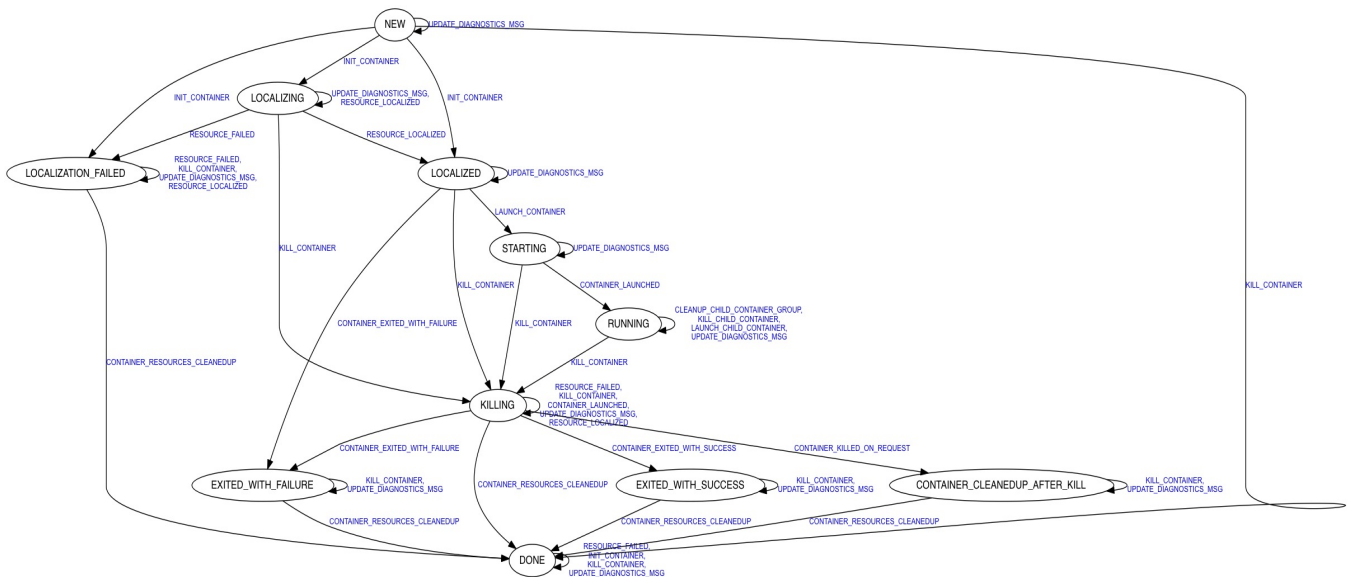
We propose to add a “new” container type modeled after daemon process/framework services on NodeManager side to provide the necessary constructs that framework developers can use. This new container:

1. Can be used to run/track framework services
2. It has capability to handle the lifecycle of other containers, in another word, NM can delegate container management to framework services containers.
3. Accordingly, container management protocol needs to be extended to allow start, stop, query the containers that are managed by the framework services containers.
4. Above are core requirements, should be efficient for process launching, killing etc.

To minimize the change needed and to be backwards compatible, we propose following implementation.

3. Implementation

3.1 Container State



1. At RUNNING state, service containers can accept additional events to support managing task containers(child containers in the diagram), those key events are: LAUNCH_CHILD_CONTAINER, KILL_CHILD_CONTAINER
2. Nice-to-have enhancements: a new STARTING state to allow start daemon process asynchronously.

1. is created by Service Container when user wants to start and manage the child containers in groups
2. by default, all child containers are in the same group managed by the service container

-
- ```
graph TD; LAUNCHING([CHILD_CONTAINER_GROUP_LAUNCHING]) -- "LAUNCH_CHILD_CONTAINER" --> RUNNING([CHILD_CONTAINER_GROUP_RUNNING]); RUNNING -- "UPDATE_DIAGNOSTICS_MSG" --> LAUNCHING; RUNNING -- "UPDATE_DIAGNOSTICS_MSG" --> KILLING([KILLING_CHILD_CONTAINER_GROUP]); RUNNING -- "KILL_CHILD_CONTAINER" --> KILLING; RUNNING -- "CHILD_CONTAINER_EXITED" --> EXITING([CHILD_CONTAINER_GROUP_EXITING]); EXITING -- "KILL_CHILD_CONTAINER" --> KILLING; EXITING -- "CHILD_CONTAINER_EXITED" --> DONE([CHILD_CONTAINER_GROUP_DONE]); KILLING -- "UPDATE_DIAGNOSTICS_MSG" --> LAUNCHING; KILLING -- "UPDATE_DIAGNOSTICS_MSG" --> KILLING; KILLING -- "CHILD_CONTAINER_EXITED_WITH_FAILURE, CHILD_CONTAINER_KILLED_ON_REQUEST, CHILD_CONTAINER_EXITED_WITH_SUCCESS" --> DONE;
```
- The diagram illustrates the lifecycle of a child container group with the following states and transitions:
- CHILD\_CONTAINER\_GROUP\_LAUNCHING** (yellow oval) transitions to **CHILD\_CONTAINER\_GROUP\_RUNNING** (green oval) via the event **LAUNCH\_CHILD\_CONTAINER**.
  - CHILD\_CONTAINER\_GROUP\_RUNNING** (green oval) transitions to **CHILD\_CONTAINER\_GROUP\_LAUNCHING** (yellow oval) via the event **UPDATE\_DIAGNOSTICS\_MSG**.
  - CHILD\_CONTAINER\_GROUP\_RUNNING** (green oval) transitions to **KILLING\_CHILD\_CONTAINER\_GROUP** (red oval) via the event **UPDATE\_DIAGNOSTICS\_MSG**.
  - CHILD\_CONTAINER\_GROUP\_RUNNING** (green oval) transitions to **KILLING\_CHILD\_CONTAINER\_GROUP** (red oval) via the event **KILL\_CHILD\_CONTAINER**.
  - CHILD\_CONTAINER\_GROUP\_RUNNING** (green oval) transitions to **CHILD\_CONTAINER\_GROUP\_EXITING** (blue oval) via the event **CHILD\_CONTAINER\_EXITED**.
  - CHILD\_CONTAINER\_GROUP\_EXITING** (blue oval) transitions to **KILLING\_CHILD\_CONTAINER\_GROUP** (red oval) via the event **KILL\_CHILD\_CONTAINER**.
  - CHILD\_CONTAINER\_GROUP\_EXITING** (blue oval) transitions to **CHILD\_CONTAINER\_GROUP\_DONE** (grey oval) via the event **CHILD\_CONTAINER\_EXITED**.
  - KILLING\_CHILD\_CONTAINER\_GROUP** (red oval) transitions to **CHILD\_CONTAINER\_GROUP\_LAUNCHING** (yellow oval) via the event **UPDATE\_DIAGNOSTICS\_MSG**.
  - KILLING\_CHILD\_CONTAINER\_GROUP** (red oval) transitions to **KILLING\_CHILD\_CONTAINER\_GROUP** (red oval) via the event **UPDATE\_DIAGNOSTICS\_MSG**.
  - KILLING\_CHILD\_CONTAINER\_GROUP** (red oval) transitions to **CHILD\_CONTAINER\_GROUP\_DONE** (grey oval) via the event **CHILD\_CONTAINER\_EXITED\_WITH\_FAILURE, CHILD\_CONTAINER\_KILLED\_ON\_REQUEST, CHILD\_CONTAINER\_EXITED\_WITH\_SUCCESS**.

We propose to reuse the current container manager protocol and only extends ContainerLaunchContext code as following:

- 3

