# Decimate: a portable and fault-tolerant scheduler.

# Application to big ensembles data assimilation and forecasting in the Red Sea Circulation

Samuel Kortas, Habib Toye, Peng Zhan, Ibrahim Hoteit

KAUST   IXUPG Middle East Conference 2018

# As core number grows
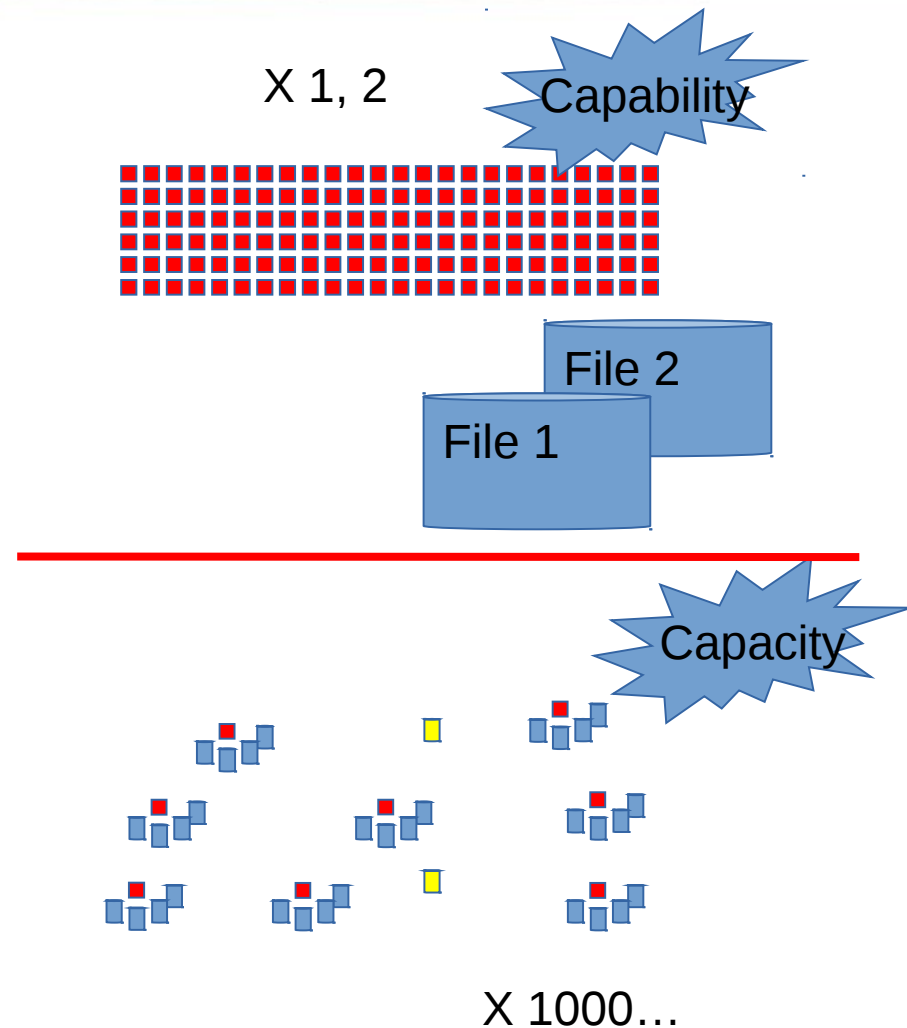## Launch of thousands of jobs became an option...

- Some of our users use Shaheen to run workflows composed of thousands of jobs saving thousands of temporary files

- Need a result in a guaranteed time

- Are not HPC experts, but are challenging problem in terms of scheduling and file system stress

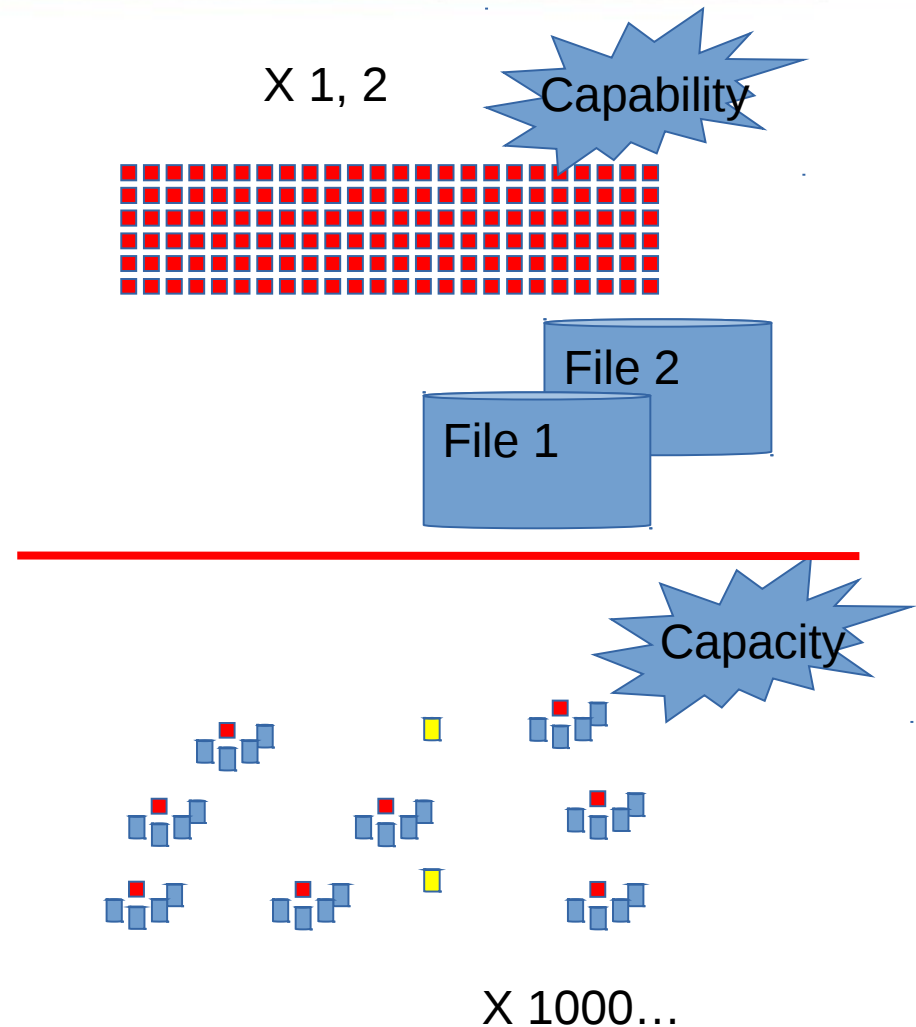- Resources are Shaped/Tuned for one of these two cases:
  - Capability (big job, big files)
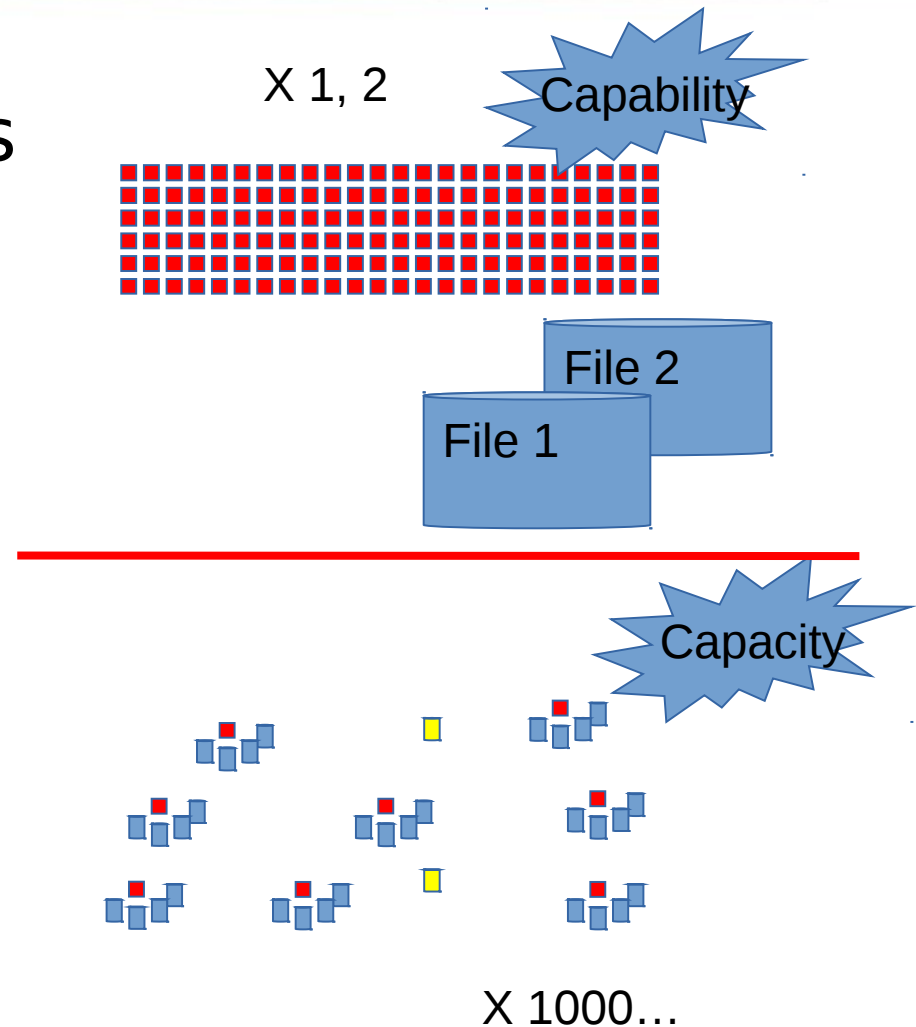  - Capacity (numerous jobs, small files)

X 1, 2

Capability

File 2

File 1

Capacity

X 1000…

# Our Strategy (1/2)

- Pack 'many jobs' together to make them appear as big ones

- Reduce stress on filesystem by using Ramdisk and messages.

X 1, 2

Capability

File 2

File 1

Capacity

X 1000…

- Extend the scheduler with tools transparent to users

- Allow them to tune easily part of workflow with parametric sweep on number of threads

X 1, 2

Capability

File 2

File 1

Capacity

X 1000…

# What is Decimate? (1/2)
## Swiss knife of many-jobs workflows

Makes It easier to

- Generate
- Submit
- balance
- Monitor
- Control
- Cure
- Check
- Terminate

- Many-jobs
- Dependent

Workflows

- Parametric
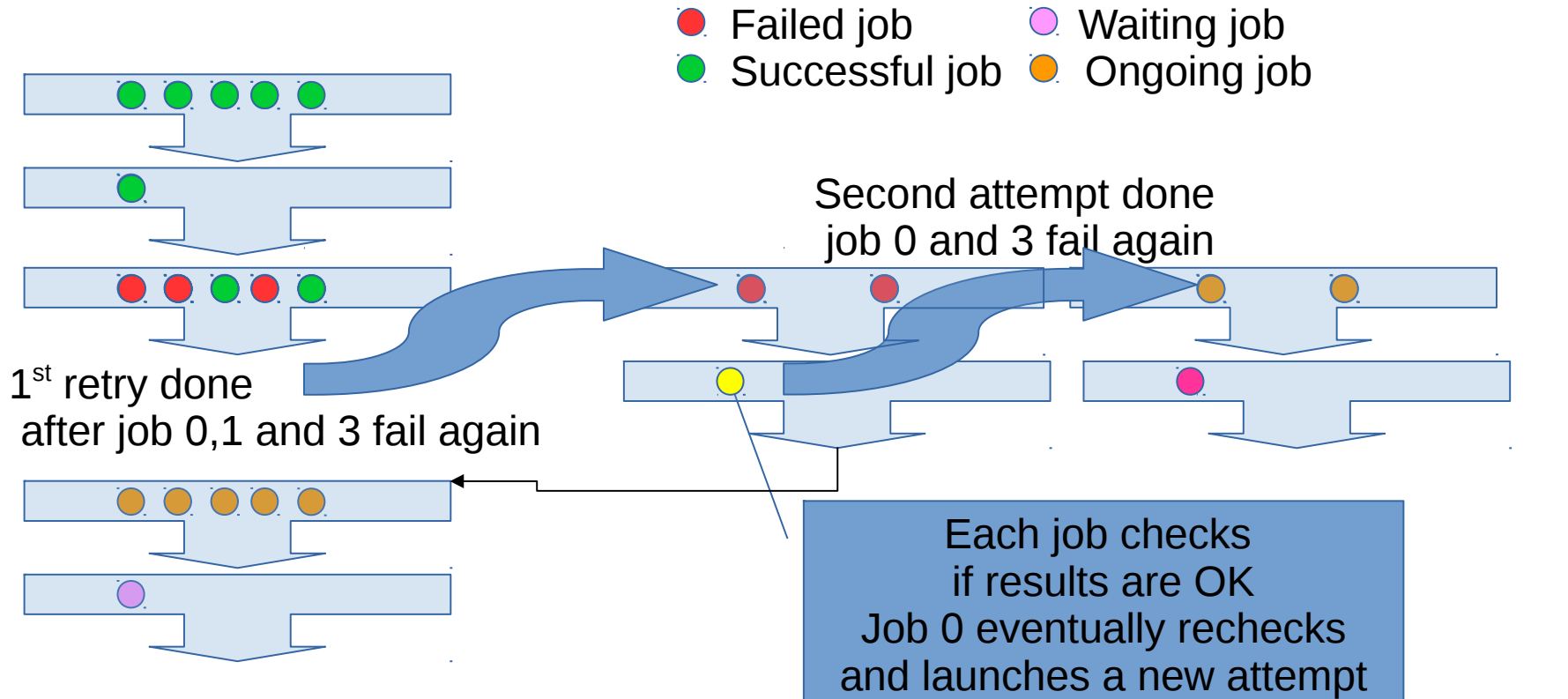- Unstable

# What is Decimate? (2/2)
## A Scheduler extension easy to install

- A scheduler extension (SLURM) written in python 2.7 fully documented at http://decimate.readthedocs.org

- Installable with no special priviledges:
  pip install decimate   or   conda install -c hpc4all decimate

- Available as an open source software at https://github.com/KAUST-KSL/  (FreeBSD License)

- Installed on Shaheen II, Ibex, portable on any machine with SLURM

- Under development but already used heavily but some of our users. 250 kjobs submitted in production via *Decimate*

- Maintained by KSL (samuel.kortas (at) kaust.edu.sa)

## Automated restart in case of failure



Failed job     Waiting job
Successful job     Ongoing job

Second attempt done
job 0 and 3 fail again

1<sup>st</sup> retry done
after job 0,1 and 3 fail again

Each job checks
if results are OK
Job 0 eventually rechecks
and launches a new attempt

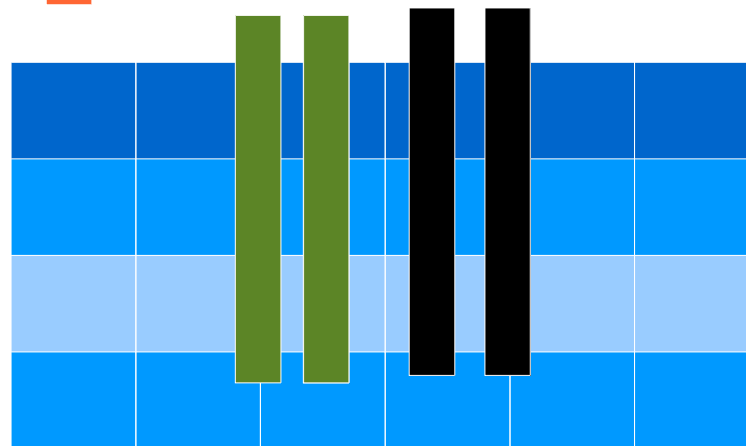dbatch   --array 1-5 --check check1.sh --max-retry 3  step1.sh

## Automated feeding of the job queue

200 jobs submitted

20 active jobs only in the queue
Others are submitted on the fly
as chunks as dependent jobs
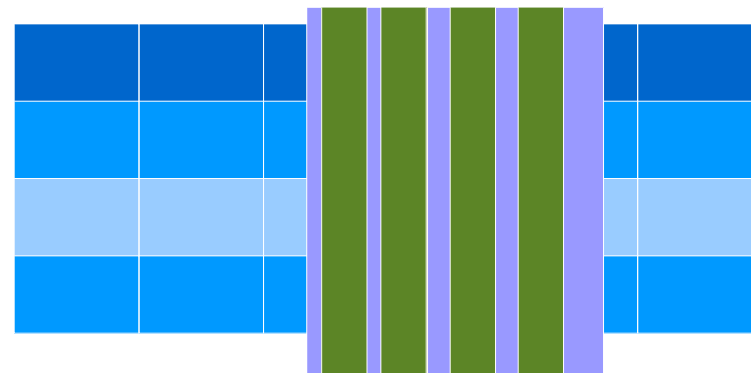
dbatch  –max-jobs 20  --array 1-200  step1.sh

## Execution on a pool of nodes

200 jobs submitted

20 active jobs in parallel on resources booked for a long duration

dbatch   –parallel-runs 20  --array 1-200 --nodes=3 step1.sh

# Decimate's Features (4/4)
## Extended support of parametrized jobs

params.txt

my_job.sh

```
nodes ntasks
  1      32
  1      16
  2      64
  4      128
#DECIM nthreads=nodes*32/ntasks
#DECIM COMBINE dim =[10,100,1000]
```

```
#SBATCH…

export OMP_NUM_THREADS=$nthreads
srun -N $nodes -n $ntasks \
     --cpu-per-tasks=$nthreads \
        my_app.exe $dim
```

dbatch  –param-file params.txt my_job.sh

Decimate computes every job possible, group them in blocks, submit them and manage them in a fault tolerant environment

**Fig. 1** *A schematic description of the MITgcm/DART assimilation system.*

Ocean model

Kalman Filter to reroute the simulation back on the right track fitting better the observations

**MITgcm**
Integrates ensemble in time

*Forecast ensemble*

*Analyzed ensemble*

**DART**
Updates ensemble with Data

**Observations**

**Fig. 1** *A schematic description of the MITgcm/DART assimilation system.*

Observations taken from satelites

**Fig. 1** *A schematic description of the MITgcm/DART assimilation system.*

**Ocean model**

**Kalman Filter to reroute the simulation back on the right track fitting better the observations**

**Observations taken from satelites**

**Fig. 1** *A schematic description of the MITgcm/DART assimilation system.*

Ocean model

Step 2: Filter & Correction

**MITgcm**

Integrates ensemble in time

*Forecast ensemble*

*Analyzed ensemble*

**DART**

Updates ensemble with Data

Kalman Filter to reroute the simulation back on the right track fitting better the observations

**Observations**

**Fig. 1** *A schematic description of the MITgcm/DART assimilation system.*

Observations taken from satelites

Ocean model

Forecast ensemble

**MITgcm**
Integrates ensemble in time

Analyzed ensemble

**DART**
Updates ensemble with Data

to reroute the simulation back on the right track fitting better

**Observations**

**Fig. 1** *A schematic description of the MITgcm/DART assimilation system.*

Observations taken from satelites

**Fig. 1** *A schematic description of the MITgcm/DART assimilation system.*

Ocean model

Step 1: Prediction

**MITgcm**
Integrates ensemble in time

*Forecast ensemble*

*Analyzed ensemble*

**DART**
Updates ensemble with Data

**Observations**

Kalman Filter to reroute the simulation back on the right track fitting better the observations

Observations taken from satelites

Ocean model

Kalman Filter to reroute the simulation back on the right track fitting better the observations

**MITgcm**

Integrates ensemble in time

*Forecast ensemble*

*Analyzed ensemble*

**DART**

Updates ensemble with Data

**Observations**

**Fig. 1** *A schematic description of the MITgcm/DART assimilation system.*

Observations taken from satelites

**Fig. 1** *A schematic description of the MITgcm/DART assimilation system.*

Ocean model

Step 1: Prediction

Kalman Filter to reroute the simulation back on the right track fitting better the observations

Observations taken from satelites

Ocean model

Step 1: Prediction
Step 1: Prediction
Step 1: Prediction
Step 1: Prediction
Step 1: Prediction

Forecast ensemble

Analyzed ensemble

**Observations**

DART
Updates ensemble with Data

Kalman Filter to reroute the simulation back on the right track fitting better the observations

**Fig. 1** *A schematic description of the MITgcm/DART assimilation system.*

Observations taken from satelites

# Typical MITgcm/DART assimilation workflow

First set of MITgcm        (1000 x 3-node runs)
–        barrier        –
Apply the filter (DART)   (1 x 16-node run)
–        barrier        –
second set of MITgcm (1000 x 3-node runs)
–        barrier        –
Apply the filter (DART)    (1 x 16-node run)
–        barrier        –
3[rd]  set of MITgcm        (1000 x 3-node runs)
–        barrier        –
Apply the filter (DART)    (1 x 16-node run)

- Some intermediate steps may break
    - → dependency will break
    - → the workflow will remain idle,

# Decimate is the solution! user-define result checking function

- the user can set his own rules thanks to a python function.

```python
################################################################
# checking job correct completion
################################################################

def check_job(self,what,task_id,running_dir,output_file,error_file,is_job_completed):

    with working_directory(running_dir):

        everything_ok = True

        for error in ['STOP ABNORMAL END','Problem opening','Problem while opening file','ERROR','
            is_error = self.greps(error,error_file,exclude_patterns=['[INFO','[DEBUG'])
            #print >> sys.stderr, 'error', is_error
            #self.log_info('user error detected --> '+pprint.pformat(is_error))
            if (is_error):
                s = '%s detected %d times in %s --> \n\t\t ERROR DETECTED: %s ' % \
                    (error,len(is_error),error_file,pprint.pformat(is_error[0]))

                # special treatment of EXTREME pot
                if is_error[0].find('S/R MON_SOLUTION, stops due to EXTREME Pot.Temp')>=0:
```

```
kortass@cdl3:/project/k1029/Sam/run_1k> d -sa

       #########################################
       #                                       #
       # Welcome to dart_mitgcm v  0.4.1!  #
       #    (using DECIMATE Framework 0.9)     #
       #                                       #
       #########################################

       running on cdl3 (shaheen)
              python /project/k1029/Sam/dart_mitgcm/0.4.1/sles11.3_gnu5.1.0/dart_mitgcm/dart_mitgcm.py -sa

[INFO ] !!!! WARNING WARNING mismatched tag single_restart_file_in values:
                              >.false.< in ensemble_manager_nml.
                              >.true.< in restart_file_tool_nml.
[INFO ] ------------> deleting non namespaced value........ "sin
[INFO ] !!!! WARNING WARNING mismatched tag restart_in_file_name
                              >perfect_ics< in perfect
                              >assim_model_state_ud< i
                              >smoother_ics< in smooth
[INFO ] ------------> deleting non namespaced value........ "res
[INFO ] input,nml.template did not change since last time, no need to regenerate
[INFO ] launch-0!0:2 unconsistent steps were found: [12-mitgcm-3,11-mitgcm-6]
 100%...[INFO ] launch-0!0:no active job in the queue, changing all WAITING in ABORTED???
[MSG ] step 2-filter-0:1-1                    SUCCESS    SUCCESS:  100%    FAILURE:   0% -> []
[MSG ] step 2-mitgcm-0:1-999                  MIXED      SUCCESS:   99%    FAILURE:   0% -> [614,617-620,622,624-625]
[MSG ] step 2-mitgcm-1:614,617-620,622,624-625 SUCCESS   SUCCESS:  100%         FAILURE:   0% -> []
[MSG ] step 3-filter-0:1-1                    SUCCESS    SUCCESS:  100%    FAILURE:   0% -> []
[MSG ] step 3-mitgcm-0:1-999                  SUCCESS    SUCCESS:  100%    FAILURE:   0% -> []
[MSG ] step 4-filter-0:1-1                    SUCCESS    SUCCESS:  100%    FAILURE:   0% -> []
[MSG ] step 4-mitgcm-0:1-999                  MIXED      SUCCESS:   74%    FAILURE:  25% -> [124,235-236,238,241,243-245,751-999]
[MSG ] step 4-mitgcm-1:124,235-236,238,241,243-245 SUCCESS   SUCCESS:  100%    FAILURE:   0% -> []
[MSG ] step 5-filter-0:1-1                    SUCCESS    SUCCESS:  100%    FAILURE:   0% -> []
[MSG ] step 5-mitgcm-0:1-999                  MIXED      SUCCESS:   99%    FAILURE:   0% -> [861,875]
[MSG ] step 5-mitgcm-1:861,875               SUCCESS     SUCCESS:  100%    FAILURE:   0% -> []
[MSG ] step 6-filter-0:1-1                    SUCCESS    SUCCESS:  100%    FAILURE:   0% -> []
[MSG ] step 6-mitgcm-0:1-999                  MIXED      SUCCESS:   99%    FAILURE:   0% -> [345,350,352]
[MSG ] step 6-mitgcm-1:345,350,352           SUCCESS     SUCCESS:  100%    FAILURE:   0% -> []
[MSG ] step 7-filter-0:1-1                    FAILURE    SUCCESS:   0%    FAILURE: 100%  -> []
```

Second step of mitgcm was restarted once after 1% of failure

With no manual intervention…
5 steps made in 3h30  on a crowded machine

# Example of Mails sent by Decimate

To: Samuel Kortas <samuel.kortas@kaust.edu.sa>;

Workflow has just been submitted

Job just started

To: Samuel Kortas <samuel.kortas@kaust.edu.sa>;

ok everything went fine for the step 4-filter (1)
Step 4-mitgcm (1) is starting...

Step is successful. Moving further!

To: Samuel Kortas <samuel.kortas@kaust.edu.sa>;

---------- problem in output file ------
String "filter finished" not found in s11/mitgcmstep11numens860.3338556.out.task_860-attempt_2 --> something went wrong
------------------------------- ------

Previous step failed and has restarted

User error detected!!! for step 11-mitgcm  task 860  attempt 2
Will restart the uncomplete step and fix the workflow

# Want to know more?



http://decimate.hpc4all.org
decimate@hpc4all.org

help@hpc.kaust.edu.sa