# Advancing the Search for Dark Energy with Parsl and HPC

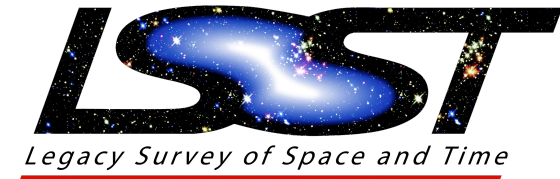Tom Glanzman - SLAC National Accelerator Laboratory
glanzman@stanford.edu

➡ in close collaboration with Ben Clifford who adapted an existing DESC workflow to Parsl and continues to partner in this endeavor

# The Rubin and DESC Projects

- Vera C. Rubin Observatory [formerly LSST] (DOE+NSF)
  - Sited on a mountain top (Cerro Pachon) in Chile
  - 8.3 meter diameter primary mirror
  - WIDE field of view (10x10 degrees)
  - Worlds largest digital camera (3.2 Gpixels)
  - Begin operation ~2022-3 with 10-year whole-sky survey program

- What is Dark Energy?
  - "Dark energy is the name given to the mysterious force that's causing the rate of expansion of our universe to accelerate over time, rather than to slow down." [ref]

- Dark Energy Science Collaboration (DOE)
  - >1000 scientist collaboration started in 2012
  - Exploit Rubin data to study clues to dark energy

Mountain top observatory (Chile)



Telescope mount (Spain)



Grinding the 8.3m lens
(Steward Observatory, Tucson, AZ)

Photos courtesy of Rubin Observatory, LSST Project/NSF/AURA

6 Oct 2020

Final raft of sensors being
installed: January 2020

- Camera focal plane
- 189 science sensors (4k x 4k pixels)
- 12 special purpose sensors (focus,pointing)

Photos courtesy of Rubin Camera Team

T.Glanzman                                          Parslfest                                          6 Oct 2020

# The DESC Data Challenges

- No data yet! (Not until ~2022-3)
- Must hit the ground running.  Therefore,
    - simulate (part of) the sky,
    - exercise the LSST project (DM) software to convert raw images into catalogs,
    - develop and test DESC-specific algorithms on the result.
- Data Challenge 2 (DC2)
    - ~300 sq. degrees of the sky (about 0.7% of entire sky)
    - 5 years of observation (one-half the Rubin survey program)
- Computational steps involved (simplified):

Subject of this talk

| Sim catalog generation | ⇨ | Image simulation | ⇨ | Image processing | ⇨ | Observation catalog generation |
|---|---|---|---|---|---|---|

- Natural parallelization: images(exposures),sensors,patches of sky, etc.
- DC2 generates >1PB of data and consumes 10's of millions of CPU hours
- DOE has provided cycles at **NERSC** and ALCF to support this work
- Image simulation step managed by Parsl at NERSC & ALCF (and presented at last year's Parslfest)

# Parsl @NERSC

- Cori-KNL (primary HPC machine at NERSC)
  - 9,688 nodes each with 68 cores x 4 hardware threads
  - Modest clock speed 1.4 GHz
  - 96 GB memory per node
- Storage = GPFS ($HOME) + Lustre ($SCRATCH)
- Batch access via SLURM
- Challenges:
  - Relatively little memory/core (or hyperthread), ~1.5 GB/core
  - Disk I/O can be problematic, slow, erratic
  - SLURM queue often experiences **very large dispatch latencies (hours to days)**, even for small jobs which can be a problem for development and production throughput
  - Rubin/LSST codes are single-threaded

# Data Release Pipeline (DRP) à la Parsl

| Task name (parsl app) | Executor | Instances (est.) |
|---|---|---|
| make_tract_list | batch-2 | 1 |
| make_patch_list_for_tract | batch-2 | 173 |
| visits_for_tract_patch_filter | batch-2 | 43,506 |
| coadd_parsl_driver | local | 43,506 |
| make_coadd_temp_exp | batch-3 | 360,595 |
| assemble_coadd | batch-4 | 43,008 |
| detect_coadd_sources | batch-4 | 43,008 |
| multiband_parsl_driver | local | 50,568 |
| merge_coadd_detections | batch-4 | 8,428 |
| deblend_coadd_sources | batch-4 | 50,568 |
| measure_coadd_sources | batch-5 | 50,568 |
| merge_coadd_measurements | batch-4 | 8,428 |
| forced_phot_coadd | batch-5 | 50,568 |

DRP consists of Rubin DM project algorithms (python/C++) representing all of the processing from raw camera images to catalogs of sky objects.

The primary Parsl apps used in this workflow.

Using *multiple HTEX executors* to **match tasks to needed resources**.

| Executor | # nodes/block | # workers/node | Clock limit |
|---|---|---|---|
| batch-2 | 1 | 200 | 9:00:00 |
| batch-3 | 400 | 22 | 10:00:00 |
| batch-4 | 50 | 20 | 10:00:00 |
| batch-5 | 100 | 50 | 24:00:00 |

# wstat - <u>w</u>orkflow <u>stat</u>us reporting tool

- A python script to read and interpret Parsl's **monitoring.db**
- Produce various (text-based) reports and plots.
- General tool -- *not tied to any specific workflow*
- Tabular reports including all runs, all tasks, full task history, etc.
- Full references to log files
- Example reports in the *Backup Slides*
- *Very much a work in progress* - if there is interest, contact me for github info

Example execution timing histogram

Report Header



```
Workflow summary at 2020-10-02 07:59:01.909485
================================================
+--------------------------+-----------------------------------------+
| workflow name            | DRPtest                                 |
| run                      | 001      <<-most current run->>         |
| run start                | 2020-09-27 10:22:05                     |
| run end                  | *pending*                               |
| run duration             | *pending*                               |
| tasks completed          | 3                                       |
| tasks completed: success | 3                                       |
| tasks completed: failed  | 0                                       |
| ----------               | ----------                              |
| workflow user            | descdm@cori20                           |
| workflow rundir          | /global/cscratch1/sd/descdm/ParslRun/dr2 |
| MonitorDB                | ./monitoring.db                         |
+--------------------------+-----------------------------------------+
```

# Parsl Wish List

- Ability to "roll back" selected task(s) within workflow
  - To expedite development of both the workflow and its component tasks
  - In production to surgically redo selected task - and it's downstream dependencies
- Improved executor with better control over task assignment to batch nodes
  - Do not start task requiring 3 hours on a batch node with only 1 hour left
  - Do not start task requiring 4 GB of memory on node with only 1 GB remaining
  - Flexibility to request #nodes/job according to task backlog
  - (User must specify these limits!)
- Monitoring
  - Extend "monitoring" to all executors
    - Very difficult to collect performance statistics without monitoring data
  - Make monitoring data *reliable*
    - Data are lost!  For example, runtime (task_time_running)
  - Better accounting of batch jobs
    - For calculating efficiency, need record of #idle workers vs time, data for tasks that fail due to batch job running out of time (ref github issue #1658)
  - Record task failure codes (return codes or time-out or crash or …) (ref issue #1453)
- Command/Control communication with running workflow
  - E.g., refresh executor parameters or other config without usual {**^c**, edit, restart} cycle
- Support for application-level checkpointing, e.g., dmtcp
  - Long-running, or batch time-outs can be restarted for better efficiency

# Backup Slides
## (intended to be viewed full screen)

# wstat - <u>W</u>orkflow <u>STAT</u>us

wstat is a very basic text-oriented report generator using data from the Parsl monitoring.db.  These reports are intended to provide a quick overall status of a running (or completed) workflow.

Obviously, monitoring must be enabled for this to work.  Currently, only the HTEX (high-throughput executor) supports monitoring 😖.

You are welcome to take wstat out for a spin.

Github repo:  https://github.com/TomGlanzman/Perp

Caveats:

- Work in progress - some features may not quite work right: consider this a *prototype tool*
- Plot function is at a very early stage of development
- You may encounter extraneous debug statements
- Monitoring.db schema can change - and foul up the SQL in wstat

# wstat "help" - listing reports and options

```
(Sat 13:35) descdm@cori20 $ python wstat -h
3.7.5 (default, Oct 25 2019, 15:51:11)
[GCC 7.3.0]
usage: wstat [-h] [-f FILE] [-r RUNNUM] [-s] [-t TASKID] [-S TASKSTATUS]
             [-l TASKLIMIT] [-d DEBUG] [-v]
             [reportType]

A simple Parsl status reporter. Available reports include:['shortSummary',
'taskSummary', 'taskHistory', 'runNums', 'runHistory', 'plot']

positional arguments:
  reportType            Type of report to display (default=shortSummary)

optional arguments:
  -h, --help            show this help message and exit
  -f FILE, --file FILE  name of Parsl monitoring database file
                        (default=./monitoring.db)
  -r RUNNUM, --runnum RUNNUM
                        Specific run number of interest (default = latest)
  -s, --schemas         only print out monitoring db schema for all tables
  -t TASKID, --taskID TASKID
                        specify task_id (taskHistory only)
  -S TASKSTATUS, --taskStatus TASKSTATUS
                        specify task_status_name
  -l TASKLIMIT, --taskLimit TASKLIMIT
                        limit output to N tasks (default is no limit)
  -d DEBUG, --debug DEBUG
                        Set debug level (default = 0)
  -v, --version         show program's version number and exit
```

Report types

Various options
(mostly to limit
output)

# wstat -- "shortSummary" example

```
3.7.5 (default, Oct 25 2019, 15:51:11)
[GCC 7.3.0]
wstat - Parsl workflow status (version  1.0.0 , written for Parsl version 1.0.0:lsst-dm-202005)

Workflow summary at 2020-09-26 13:19:54.488504
==============================================
+------------------------------+---------------------------------------------+
| workflow name                | DRPtest                                     |
| run                          | 000      <<-most current run->>             |
| run start                    | 2020-09-26 10:40:24                          |
| run end                      | *pending*                                   |
| run duration                 | *pending*                                   |
| tasks completed              | 33705                                       |
| tasks completed: success     | 33623                                       |
| tasks completed: failed      | 82                                          |
| ----------                   | ----------                                  |
| workflow user                | descdm@cori20                               |
| workflow rundir              | /global/cscratch1/sd/descdm/ParslRun/dr2    |
| MonitorDB                    | ./monitoring.db                             |
+------------------------------+---------------------------------------------+

Node usage summary:
+----------+------------+
| Node     | #running   |
|----------+------------|
| nid02525 |         48 |
| nid02526 |         50 |
| nid03040 |        241 |
| nid04629 |         22 |
```

*(Header)*

*(List of nodes currently running - and the number of tasks running on each)*

```
| nid09640 |          8 |
| nid09641 |         11 |
| nid09642 |         23 |
| nid09643 |         11 |
+----------+------------+
   Number of active nodes =  110
   Number of running tasks =  2799
```

*(Statistics for all tasks for each Parsl "state")*

```
Task status matrix:
+-----------------------------+---------+----------+---------+---------+---------+---------+----------+-----------+--------+----------+---------------+--------+
|                             | pending | launched | joining | running | unsched | unknown | exec_done | memo_done | failed | dep_fail | fail_retryable | TOTAL  |
+-----------------------------+---------+----------+---------+---------+---------+---------+----------+-----------+--------+----------+---------------+--------+
| make_tract_list             |       0 |        0 |       0 |       0 |       0 |       0 |        1 |         0 |      0 |        0 |             0 |      1 |
| make_patch_list_for_tract   |       0 |        0 |       0 |       0 |       0 |       0 |      173 |         0 |      0 |        0 |             0 |    173 |
| process_patches             |       0 |        0 |     173 |       0 |       0 |       0 |      173 |         0 |      0 |        0 |             0 |    173 |
| visits_for_tract_patch_filter |     0 |    26929 |       0 |     241 |       0 |       0 |    16336 |         0 |      0 |        0 |             0 |  43506 |
| coadd_parsl_driver          |   27170 |       79 |   13976 |       0 |       0 |       0 |     2281 |         0 |      0 |        0 |             0 |  43506 |
| multiband_parsl_driver      |    6952 |        0 |     282 |       0 |       0 |       0 |        0 |         0 |     17 |        0 |             0 |   7251 |
| combine                     |     455 |        0 |       0 |       0 |       0 |       0 |        0 |         0 |      0 |       17 |             0 |    472 |
| make_coadd_temp_exp         |       1 |   110882 |       0 |     824 |       0 |       0 |     9180 |         0 |      0 |        0 |             0 | 120887 |
| assemble_coadd              |    8199 |     5404 |       0 |     355 |       0 |       0 |     2111 |         0 |      0 |        0 |             0 |  16069 |
| detect_coadd_sources        |   13958 |        1 |       0 |      17 |       0 |       0 |     2093 |         0 |      0 |        0 |             0 |  16069 |
| merge_coadd_detections      |       0 |        1 |       0 |      38 |       0 |       0 |      260 |         0 |      0 |        0 |             0 |    299 |
| deblend_coadd_sources       |     234 |       39 |       0 |     535 |       0 |       0 |      986 |         0 |      0 |        0 |             0 |   1794 |
| measure_coadd_sources       |     801 |       72 |       0 |     776 |       0 |       0 |      145 |         0 |      0 |        0 |             0 |   1794 |
| merge_coadd_measurements    |     278 |        0 |       0 |       2 |       0 |       0 |       19 |         0 |      0 |        0 |             0 |    299 |
| forced_phot_coadd           |    1680 |        0 |       0 |      11 |       0 |       0 |       38 |         0 |     65 |        0 |             0 |   1794 |
| TOTAL                       |   59728 |   143407 |   14431 |    2799 |       0 |       0 |    33623 |         0 |     82 |       17 |             0 | 254087 |
+-----------------------------+---------+----------+---------+---------+---------+---------+----------+-----------+--------+----------+---------------+--------+
wstat elapsed time =  0:00:10.051733
```

# wstat - "taskSummary" example



```
(Sat 21:25) descdm@cori20 $ python wstat taskSummary -S failed
3.7.5 (default, Oct 25 2019, 15:51:11)
[GCC 7.3.0]
wstat - Parsl workflow status (version  1.0.0 , written for Parsl version 1.0.0:lsst-dm-202005)

Workflow summary at 2020-09-26 21:26:05.022466
==============================================
+--------------------------+------------------------------------------------+
| workflow name            | DRPtest                                        |
| run                      | 000     <<-most current run->>                 |
| run start                | 2020-09-26 10:40:24                            |
| run end                  | *pending*                                      |
| run duration             | *pending*                                      |
| tasks completed          | 114240                                         |
| tasks completed: success | 113963                                         |
| tasks completed: failed  | 277                                            |
| ----------               | ----------                                     |
| workflow user            | descdm@cori20                                  |
| workflow rundir          | /global/cscratch1/sd/descdm/ParslRun/dr2       |
| MonitorDB                | ./monitoring.db                                |
+--------------------------+------------------------------------------------+
```

Report consisting of every failed task

Header

```
+--------+---------------------+---------+--------+----------+-----+-------+---------------------+---------------------+---------------------+----------------+-----------------------------------------+
| task_id| task_name           | run_num | status | hostname | try | #fails| submitTime          | startTime           | endTime             | runTime        | stdout                                  |
+--------+---------------------+---------+--------+----------+-----+-------+---------------------+---------------------+---------------------+----------------+-----------------------------------------+
|   3042 | multiband_parsl_driver |    000 | failed |          |  0  |   1   | 2020-09-26 19:50:28 |                     | 2020-09-26 19:50:28 |                | None                                    |
|   3068 | multiband_parsl_driver |    000 | failed |          |  0  |   1   | 2020-09-26 19:29:50 |                     | 2020-09-26 19:29:50 |                | None                                    |
|   3095 | multiband_parsl_driver |    000 | failed |          |  0  |   1   | 2020-09-26 19:59:57 |                     | 2020-09-26 19:59:57 |                | None                                    |
|   3115 | multiband_parsl_driver |    000 | failed |          |  0  |   1   | 2020-09-26 20:07:57 |                     | 2020-09-26 20:07:57 |                | None                                    |
| 107934 | forced_phot_coadd   |     000 | failed | nid08920 |  2  |   3   | 2020-09-26 11:19:01 | 2020-09-26 11:19:02 | 2020-09-26 11:22:24 | 0:03:22.335702 | /global/cscratch1/sd/descdm/ParslRun/dr2/ |
| 107935 | forced_phot_coadd   |     000 | failed | nid08920 |  2  |   3   | 2020-09-26 11:19:01 | 2020-09-26 11:19:02 | 2020-09-26 11:22:24 | 0:03:22.335742 | /global/cscratch1/sd/descdm/ParslRun/dr2/ |
| 107937 | forced_phot_coadd   |     000 | failed | nid08920 |  2  |   3   | 2020-09-26 11:19:01 | 2020-09-26 11:19:02 | 2020-09-26 11:22:24 | 0:03:22.330392 | /global/cscratch1/sd/descdm/ParslRun/dr2/ |
| 107939 | forced_phot_coadd   |     000 | failed | nid08920 |  2  |   3   | 2020-09-26 11:19:01 | 2020-09-26 11:19:02 | 2020-09-26 11:22:24 | 0:03:22.340348 | /global/cscratch1/sd/descdm/ParslRun/dr2/ |
| 107955 | forced_phot_coadd   |     000 | failed | nid08920 |  2  |   3   | 2020-09-26 11:19:01 | 2020-09-26 11:19:02 | 2020-09-26 11:22:24 | 0:03:22.347536 | /global/cscratch1/sd/descdm/ParslRun/dr2/ |
| 107957 | forced_phot_coadd   |     000 | failed | nid08920 |  2  |   3   | 2020-09-26 11:19:01 | 2020-09-26 11:19:02 | 2020-09-26 11:22:24 | 0:03:22.274380 | /global/cscratch1/sd/descdm/ParslRun/dr2/ |
| 107959 | forced_phot_coadd   |     000 | failed | nid08920 |  2  |   3   | 2020-09-26 11:19:01 | 2020-09-26 11:19:02 | 2020-09-26 11:22:24 | 0:03:22.346121 | /global/cscratch1/sd/descdm/ParslRun/dr2/ |
| 107962 | forced_phot_coadd   |     000 | failed | nid08920 |  2  |   3   | 2020-09-26 11:19:01 | 2020-09-26 11:19:02 | 2020-09-26 11:22:24 | 0:03:22.348449 | /global/cscratch1/sd/descdm/ParslRun/dr2/ |
| 107976 | forced_phot_coadd   |     000 | failed | nid08935 |  2  |   3   | 2020-09-26 11:19:44 | 2020-09-26 11:19:45 | 2020-09-26 11:23:07 | 0:03:22.321341 | /global/cscratch1/sd/descdm/ParslRun/dr2/ |
```

Useful information for investigating problems

Including location of logs

# wstat - "tasksHistory" example

```
(Fri 11:38) descdm@cori20 $ python wstat taskHistory -t 7604
```

Report consisting of full history for a specific task

```
Workflow summary at 2020-10-02 11:38:25.691188
===============================================
+------------------------+------------------------------------------+
| workflow name          | DRPtest                                  |
| run                    | 000    <<-most current run->>            |
| run start              | 2020-08-05 15:18:12                      |
| run end                | 2020-08-06 11:19:25                      |
| run duration           | 20:01:13.490372                          |
| tasks completed        | 33084                                    |
| tasks completed: success | 33070                                  |
| tasks completed: failed | 14                                      |
| ----------             | ----------                               |
| workflow user          | descdm@cori20                            |
| workflow rundir        | /global/u1/d/descdm/tomTest/DRPtest/workDir |
| MonitorDB              | archive/50.gen/monitoring.db             |
+------------------------+------------------------------------------+
```

Header

| task_id | task_name | run_num | status | hostname | try | #fails | submitTime | startTime | endTime | runTime | stdout |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 7604 | make_coadd_temp_exp | 000 | pending | nid07828 | 0 | 3 | 2020-08-05 15:27:56 | 2020-08-05 18:44:03 | 2020-08-05 18:47:25 | 0:03:21.707510 | /global/u1/d/descdm/tomTest/DRPtest/wo |
| 7604 | make_coadd_temp_exp | 000 | launched | nid07828 | 0 | 3 | 2020-08-05 15:27:56 | 2020-08-05 18:44:03 | 2020-08-05 18:47:25 | 0:03:21.707510 | /global/u1/d/descdm/tomTest/DRPtest/wo |
| 7604 | make_coadd_temp_exp | 000 | running | nid07828 | 0 | 3 | 2020-08-05 15:27:56 | 2020-08-05 18:44:03 | 2020-08-05 18:47:25 | 0:03:21.707510 | /global/u1/d/descdm/tomTest/DRPtest/wo |
| 7604 | make_coadd_temp_exp | 000 | fail_retryable | nid07828 | 0 | 3 | 2020-08-05 15:27:56 | 2020-08-05 18:44:03 | 2020-08-05 18:47:25 | 0:03:21.707510 | /global/u1/d/descdm/tomTest/DRPtest/wo |
| 7604 | make_coadd_temp_exp | 000 | pending | nid07829 | 1 | 3 | 2020-08-05 18:47:25 | 2020-08-05 18:47:26 | 2020-08-05 18:50:48 | 0:03:21.658725 | /global/u1/d/descdm/tomTest/DRPtest/wo |
| 7604 | make_coadd_temp_exp | 000 | launched | nid07829 | 1 | 3 | 2020-08-05 18:47:25 | 2020-08-05 18:47:26 | 2020-08-05 18:50:48 | 0:03:21.658725 | /global/u1/d/descdm/tomTest/DRPtest/wo |
| 7604 | make_coadd_temp_exp | 000 | running | nid07829 | 1 | 3 | 2020-08-05 18:47:25 | 2020-08-05 18:47:26 | 2020-08-05 18:50:48 | 0:03:21.658725 | /global/u1/d/descdm/tomTest/DRPtest/wo |
| 7604 | make_coadd_temp_exp | 000 | fail_retryable | nid07829 | 1 | 3 | 2020-08-05 18:47:25 | 2020-08-05 18:47:26 | 2020-08-05 18:50:48 | 0:03:21.658725 | /global/u1/d/descdm/tomTest/DRPtest/wo |
| 7604 | make_coadd_temp_exp | 000 | pending | nid07832 | 2 | 3 | 2020-08-05 18:50:48 | 2020-08-05 18:50:48 | 2020-08-05 18:54:10 | 0:03:21.688318 | /global/u1/d/descdm/tomTest/DRPtest/wo |
| 7604 | make_coadd_temp_exp | 000 | launched | nid07832 | 2 | 3 | 2020-08-05 18:50:48 | 2020-08-05 18:50:48 | 2020-08-05 18:54:10 | 0:03:21.688318 | /global/u1/d/descdm/tomTest/DRPtest/wo |
| 7604 | make_coadd_temp_exp | 000 | running | nid07832 | 2 | 3 | 2020-08-05 18:50:48 | 2020-08-05 18:50:48 | 2020-08-05 18:54:10 | 0:03:21.688318 | /global/u1/d/descdm/tomTest/DRPtest/wo |
| 7604 | make_coadd_temp_exp | 000 | failed | nid07832 | 2 | 3 | 2020-08-05 18:50:48 | 2020-08-05 18:50:48 | 2020-08-05 18:54:10 | 0:03:21.688318 | /global/u1/d/descdm/tomTest/DRPtest/wo |

Full history of this task's attempt to run through Parsl

# wstat - "runHistory" example

```
(Sun 11:26) descdm@cori20 $ python wstat runHistory
3.7.5 (default, Oct 25 2019, 15:51:11)
[GCC 7.3.0]
wstat - Parsl workflow status (version  1.0.0 , written for Parsl version 1.0.0:lsst-dm-202005)

+----------+----------------+---------+--------+---------------------+---------------------+-------------+-------------+-------------+----------------------------------------------------------+
| RunNum | workflow_name  | user   | host   | time_began          | time_completed      | RunDuration |   #tasks_good |   #tasks_bad | rundir                                                   |
|----------+----------------+---------+--------+---------------------+---------------------+-------------+-------------+-------------+----------------------------------------------------------|
|   000 | DRPtest        | descdm | cori20 | 2020-09-26 10:40:24 | 2020-09-27 09:58:16 | 23:17:52    |      181874 |         386 | /global/cscratch1/sd/descdm/ParslRun/dr2/runinfo/000 |
|   001 | DRPtest        | descdm | cori20 | 2020-09-27 10:22:05 | -> incomplete <-    |             |           3 |           0 | /global/cscratch1/sd/descdm/ParslRun/dr2/runinfo/001 |
+----------+----------------+---------+--------+---------------------+---------------------+-------------+-------------+-------------+----------------------------------------------------------+
wstat elapsed time =  0:00:00.725282
```