# NUCAPS and ACSPO Reprocessing at CIMSS/SSEC/UW

Liam Gumley, Steve Dutcher, Bruce Flynn, Jim Davies
JPSS STAR Science Team Meeting, 8/25/2015

# Overview

Question: What would it take to reprocess the entire SNPP mission record to generate a consistent set of NUCAPS and ACSPO products?

- NUCAPS and ACSPO are the NOAA enterprise algorithms for Suomi NPP atmospheric profiles and sea surface temperature.

- STAR would like to have consistent calibration and retrieval algorithms, LUTs, and products for the entire mission. Therefore, start with RDRs.

# CIMSS Key Ingredients

- Complete archive of VIIRS, CrIS, and ATMS RDRs at CIMSS since start of SNPP mission.

- CSPP SDR processing software (to convert RDR to SDR) based on Mx 8.4, with LUTs for entire mission.

- Excellent collaboration with STAR NUCAPS and ACSPO enterprise algorithm development teams for preparing CSPP release packages.

- Cluster compute, storage, and expertise at CIMSS.

# CIMSS Cluster Overview

- 82 compute nodes (servers) with 8 to 16 cores per node. Total core count = 1184. 64-bit Linux CentOS.

- Approximately 4 GB to 8 GB of memory per core.

- Each compute node has between 0.3 - 1.0 TB of workspace.

- Total of 2 Petabytes of cluster storage.

- Network delivers data from storage to compute nodes speeds of up to 30 gigabits/second (aggregated).

# Job Management

- We use a PostgreSQL database to keep track of more than 50 million files from multiple satellites (Suomi-NPP, Aqua, Terra, Metop-A/B, Caliop,…)

- We use Condor as our workload management system which provides dynamic scheduling based on job requirements (e.g., memory, disk space).

- We use custom workflow manager (Flo) that scans the database for files to process and then uses Condor to submit the jobs.

# Data Volumes

**Input RDR volumes:**

   VIIRS RDR = 74 TB; CrIS & ATMS RDR = 18 TB

**Intermediate SDR volumes:**

   VIIRS SDR (GMTCO + 6 M-bands required by ACSPO) = 70 TB

   CrIS & ATMS SDR = 54 TB

**Level 2 product volumes:**

   NUCAPS L2 = 11 TB

   ACSPO L2, L2/L3 SQUAM, MICROS, L3U = 145 TB

# NUCAPS Overview

**Inputs (per 8 minutes)**

1 x CrIS RDR, 8m aggregated

3 x ATMS RDR, 8m aggregated

2 x GFS AVN

**Outputs** (32 second granules)

15 x NUCAPS EDR

15 x NUCAPS CCR

15 x statistic file for analysis

**Science Software**

CSPP SDR v2.0.1

CSPP NUCAPS v1.0.2

**Time Period Processed**

May 1, 2012 through Jul 31, 2015

Some data missed due to:

- inability to produce ATMS before May 1, 2012

- missing/failed CrIS/ATMS processing

- missing GFS AVN ancillary data

**Job Granularity**

Single 8 minute aggregated CrIS granule per job

Extra downloads of ATMS data due to processing of single granule, however …

Smaller job granularity parallelizes better on cluster increasing overall job throughput

1184 Days * 180 granules/day = ~66K jobs

# NUCAPS Processing

**ATMS SDR Processing**

~3 hours/year

**CrIS SDR Processing**

~9 hours/year

**NUCAPS Processing**

1184 days * 180 granules/day = 66K jobs

45 outputs per job = 3M outputs

32 minutes/job * 1184 cores = 30 hours/year

**Total Processing**

42 hours/year for ATMS/CrIS SDR + NUCAPS (Speedup: 208X)

~97% product yield

1 week to reprocess entire mission

# NUCAPS Lessons

**Yield can be improved**

- Need to investigate CrIS/ATMS SDR processing failure cases to improve yield (this includes ATMS SDR pre May 1, 2012).

- It appears NUCAPS does not handle CrIS granules with less than 4 scan lines (known issue).

**Helpful improvements for cluster environment**

- Specify all required inputs on command line.

- Assume software tree is read-only.

- Assume the network is not available, i.e. can't download ancillary data (the workflow manager will provide it).

# ACSPO Overview

**Inputs**

VIIRS RDR (86-second)

Ancillary: CMC, OSTIA, Daily Reynolds, iQUAM

**Outputs**

png/json files for web

product files

**Science Software**

CSPP VIIRS SDR v2.0.1 with VCST LUTS

ACSPO v2.40

**Time Period Processed**

May 1, 2012 through Jul 31, 2015.

Some data missed due to inability to produce VIIRS before March 1, 2012.

**Job Highlights**

1210 Days =1.4 million jobs.

Saving intermediate products from various stages produced 500 TB of data.

Network sustained 30 gigabits/second when delivering large dataset to compute nodes.

# ACSPO Processing

**VIIRS Processing (granularity = 86 seconds)**

~52 hours/year

**Aggregate/Destripe (granularity = 10 minutes)**

~6 hours/year

**ACSPO Processing (granularity = 1 day)**

ACSPO:  ~35 hours/year

ACSPO MICROS:  ~30 hours/year

ACSPO L2 SQUAM:  ~33 hours/year

ACSPO L3 SQUAM:  ~10 hours/year

**Total Processing**

165 hours to reprocess 1 year of through all ACSPO steps (Speedup: 53X)

3.5 weeks to reprocess the entire record

# ACSPO Lessons

**Yield can be improved**

Need to investigate VIIRS SDR processing failures for March 2012 and get a working set of LUTS allowing us to got back to Jan 2012.

**Helpful improvements for cluster environment**

Our biggest challenge was IDL integration. We needed to find a way to run IDL jobs in IDLRT cluster mode. ACSPO team and CIMSS worked together to solve this problem.

Numerous jobs had long run times (~10 hours), which makes the feedback loop for debugging purposes difficult.  Jobs that have smaller runtimes are highly preferred not only for debugging purposes but also it will parallelize better over the cluster.

# Reprocessing Summary

Numerous software packages need to be adapted to cluster environments,for example:
  - removing hard coded paths
  - dynamic ancillary selected prior to job submission
  - keeping the granularity small (~1 hour) helps in both debugging and optimal use of the cluster hardware.

Understanding the software such that jobs can be scheduled based on resources needed (cpu, memory, disk) is critical to reprocessing as fast as possible.

There is much attention to detail needed to achieve 99% data coverage over the life of the mission.

# Conclusion

Thanks to NOAA/NESDIS/STAR NUCAPS and ACSPO for their close collaboration.

CIMSS team (Steve Dutcher, Bruce Flynn, Jim Davies) put in many hours on this demonstration project.

We will work with STAR to resolve outstanding issues in the NUCAPS and ACSPO reprocessed datasets.

We look forward to future collaboration with STAR on reprocessing tasks.