

Joseph C. Jacob, Brian D. Wilson and Huikyo Lee

Jet Propulsion Laboratory, California Institute of Technology, Pasadena, CA 91109-8099

Abstract

The need to better understand climate change has driven model simulations to greater fidelity with improved spatiotemporal resolution (e.g., < 10 km at sub-hourly cadence). The rise of these high-fidelity climate models coincides with the emergence of cloud computing as a viable platform for scientific analytics. However, it is not cost- or time-effective to move the High-Performance Computing (HPC)-based model computations and data to the cloud. Thus, there is a need for scalable model evaluation compatible with both the cloud and HPC platforms like NASA Center for Climate Simulation (NCCS). To fill this need we have extended the analytics component of the Apache Science Data Analytics Platform (SDAP) with a streamlined adaptation that specifically targets high-resolution science data products and climate model outputs on a regular coordinate grid. Gridded inputs (as opposed to other data structures like point clouds or swath-based measurements supported by SDAP), enable offsets to particular grid cells to be directly computed, allow for processing on the original NetCDF or HDF granules, do not require a second tiled copy of the data, and accommodate a simpler technology stack since no geospatial database is required for lookups or tile storage. Our core module, Parmap, abstracts the map-reduce model so that users can select from a variety of map computational modes, including Spark, Dask, serverless AWS Lambda, PySparkling, and Python multiprocessing.

Parmap Key Features

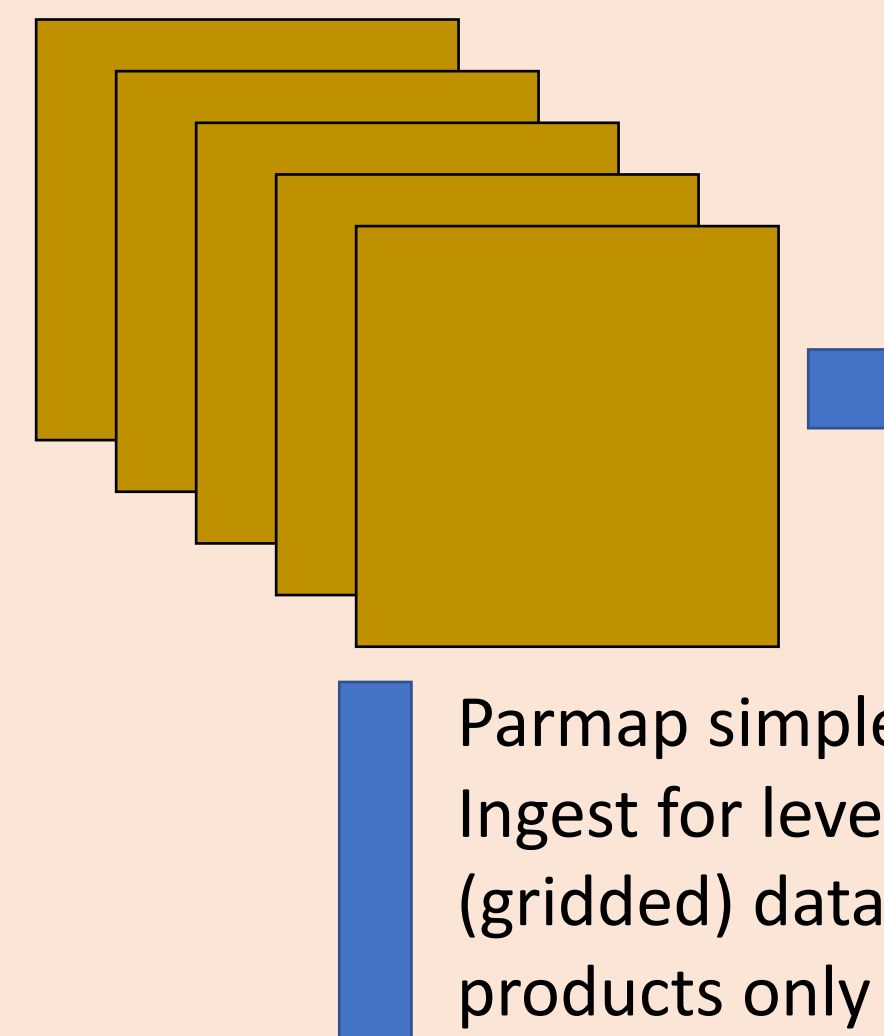
- Enables **simple python coding** for computing analytics in parallel
- Extends SDAP with a **streamlined workflow** optimized for data products on a regular coordinate grid.
- Operate on the original data granules so **no data duplication** is needed and no databases are required.
- Nearly **no data ingest** mechanism needed. Simply prepare a configuration file for a dataset to indicate how the date or timestamp can be extracted from the data granules or their filenames and the dataset is ready for Parmap analytics.
- No persistent services** (like databases) required means it is easy to deploy Parmap to shared HPC environments as well as public and private clouds.
 - ✓ Multicore (single node)
 - ✓ PySparkling (pure Python, single node)
 - ✓ Full Apache PySpark Cluster (multi-node)
 - ✓ Dask Cluster (multi-node)
 - ✓ AWS Lambda (serverless)
- Fully compatible with the existing SDAP architecture. Use full SDAP for selected datasets that have been ingested and Parmap as default for regularly gridded data.
- SDAP/Parmap architecture allows access via four interfaces:
 - ✓ Web service call
 - ✓ Python client library
 - ✓ Jupyter Notebook
 - ✓ Web GUI with map displays, built-in analytics

Acknowledgements

Parmap is being developed under NASA's Computational Modeling Algorithms and Cyberinfrastructure (CMAC) Program. SDAP was initially sponsored by NASA's Advanced Information Systems Technology (AIST) Program under the OceanWorks project (PI: Thomas Huang, JPL).

New Parmap workflow

Local file system or Object Store (e.g., AWS S3):
Granules (e.g., NetCDF, HDF)



Simply add a configuration file that describes how the date/time can be extracted from the filename (e.g., regular expression) or from the file contents (e.g., variable name).

Parmap Data Reader

Data array "slices" (subsets) – read only the data needed

Parmap Orchestrator

Parmap Mapper

AWS Lambda
Dask Cluster
Apache Spark Cluster
Python Multiprocessing
PySparkling

Flexible parallelism modes for single- and multi-node platforms

Job requests

Results from Analytics (e.g., JSON)

Job parameters

Spark job submission

Tile data retrieval

Geospatial search for tiles in bounding box

Spark In-Memory Parallel Computation

Web Server and Application Framework

Metadata

Data

Tiles

Solr

Cassandra

Standard SDAP/NEXUS workflow (previously reported)

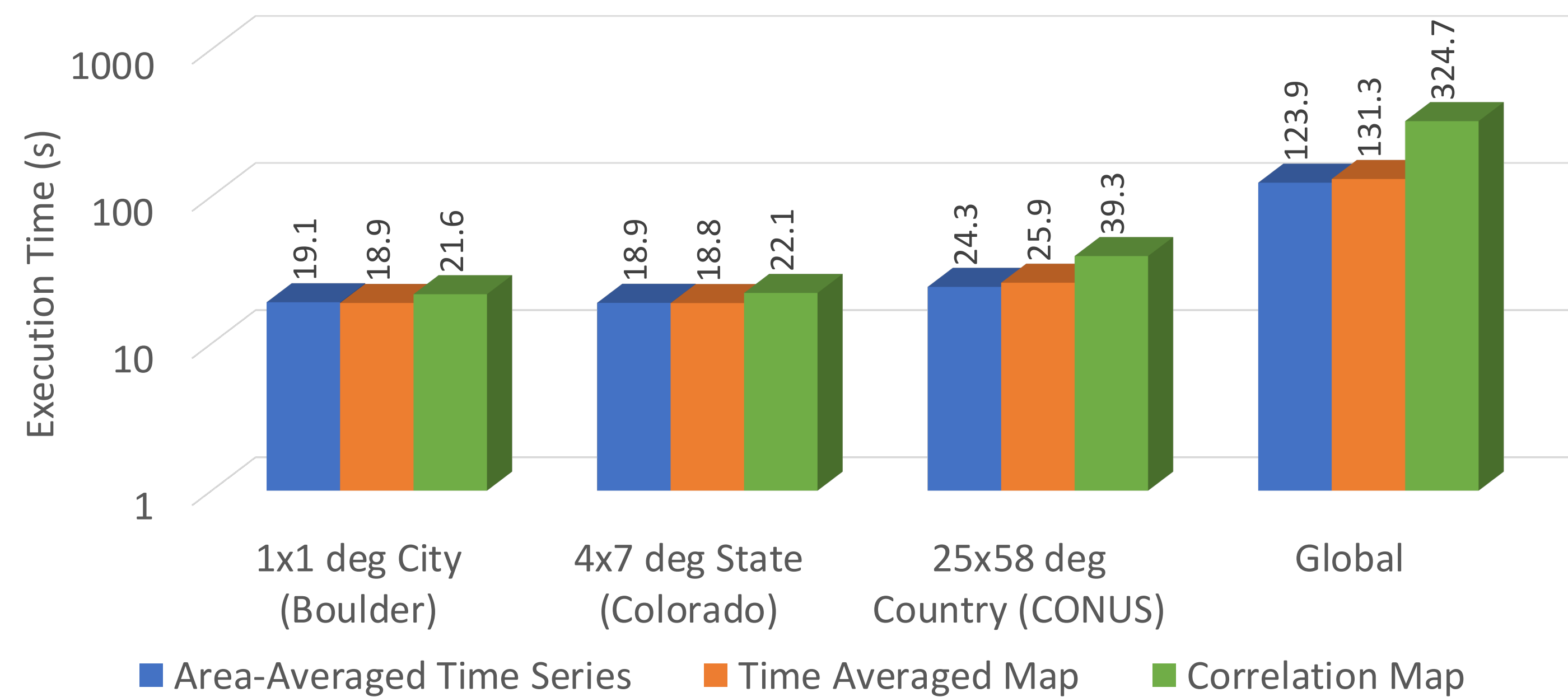
Parmap Analytics Benchmarks from GEOS5

Nature Run (G5NR):

- Platform: 8-node Spark cluster
- Dataset: GEOS5 Nature Run (G5NR)
 - June 2005 – June 2007
 - 7 km spatial resolution
 - 30 minute model outputs
 - ~4 PB (>40 variables)
- Input data used for these benchmarks:
 - 30 days of G5NR: 9/1/2006 – 9/30/2006
 - 1 variable for area-averaged time series and time averaged map: *total aerosol extinction AOT [550 nm]*
 - 2nd variable for correlation map: *dust surface mass concentration [pm 2.5]*
- Various spatial extents:
 - City – 1x1 deg Boulder, Colorado, USA
 - State – 4x7 deg Colorado, USA
 - Country – 25x58 deg Continental United States (CONUS)
 - Global
- Plan to be horizontally scalable for full 2-year simulation

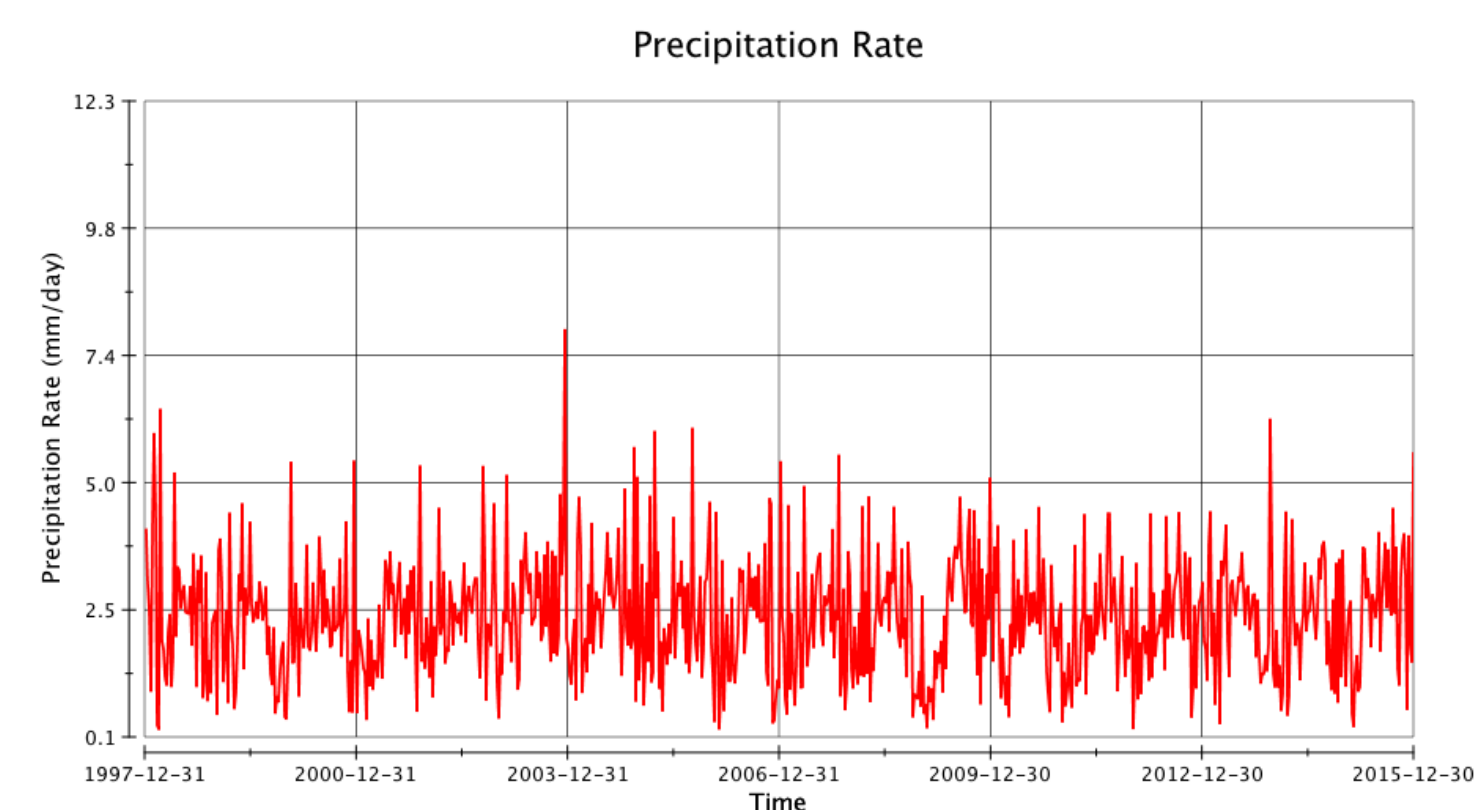
Parmap Analytics Benchmarks

Spark Cluster, 16 Processors

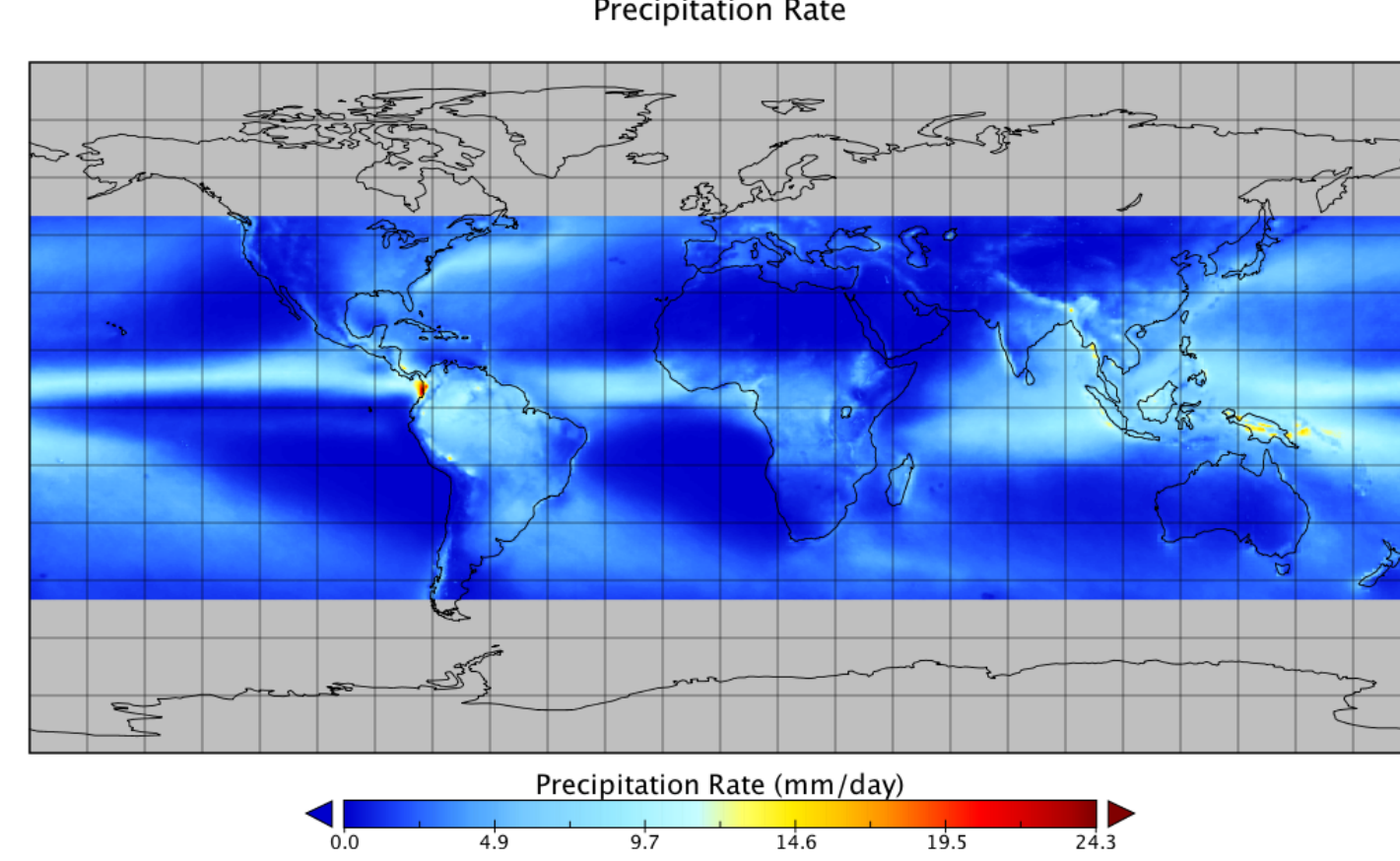


SDAP analytics currently available in Parmap:

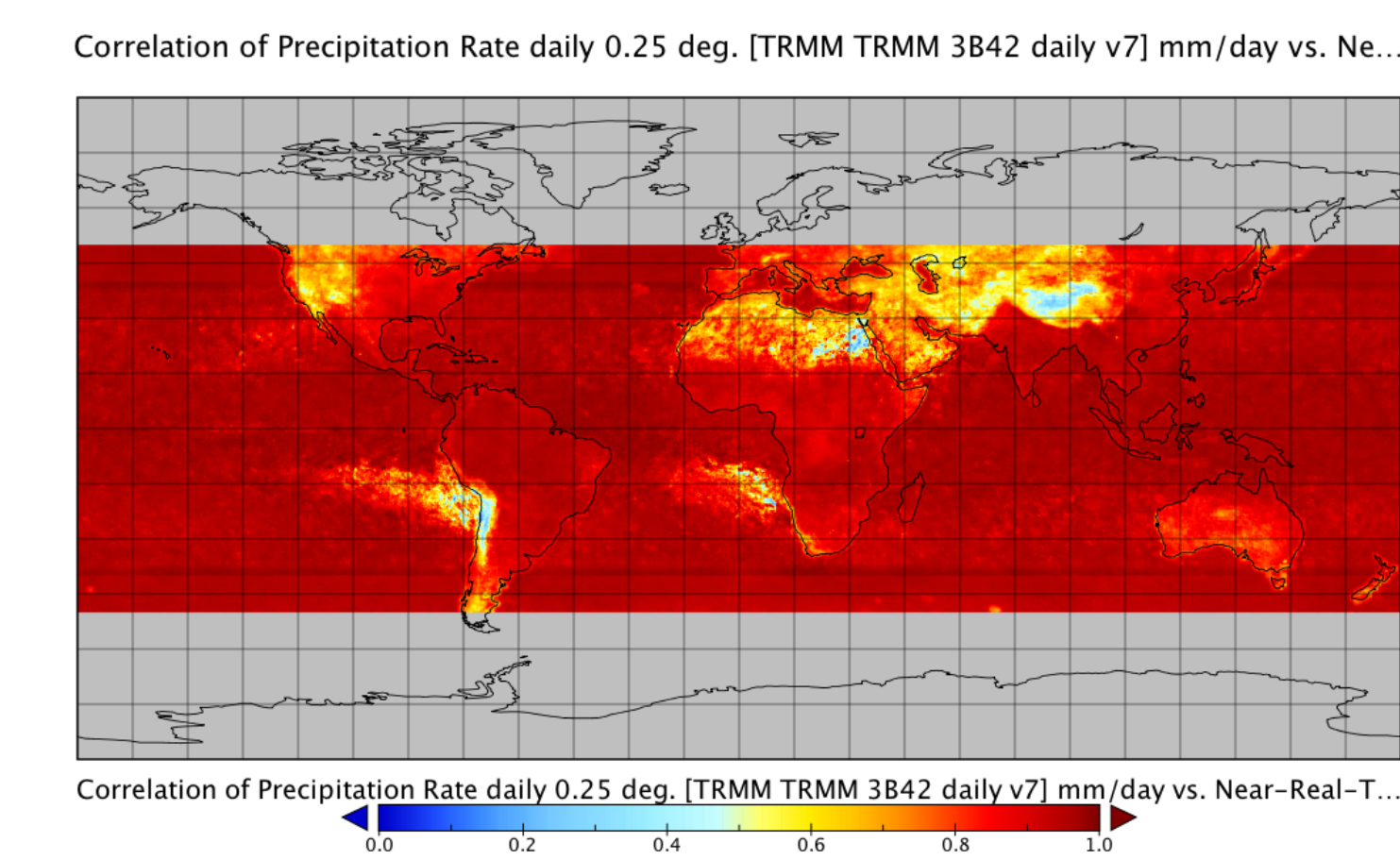
Area-Averaged Time Series



Time-Averaged Map



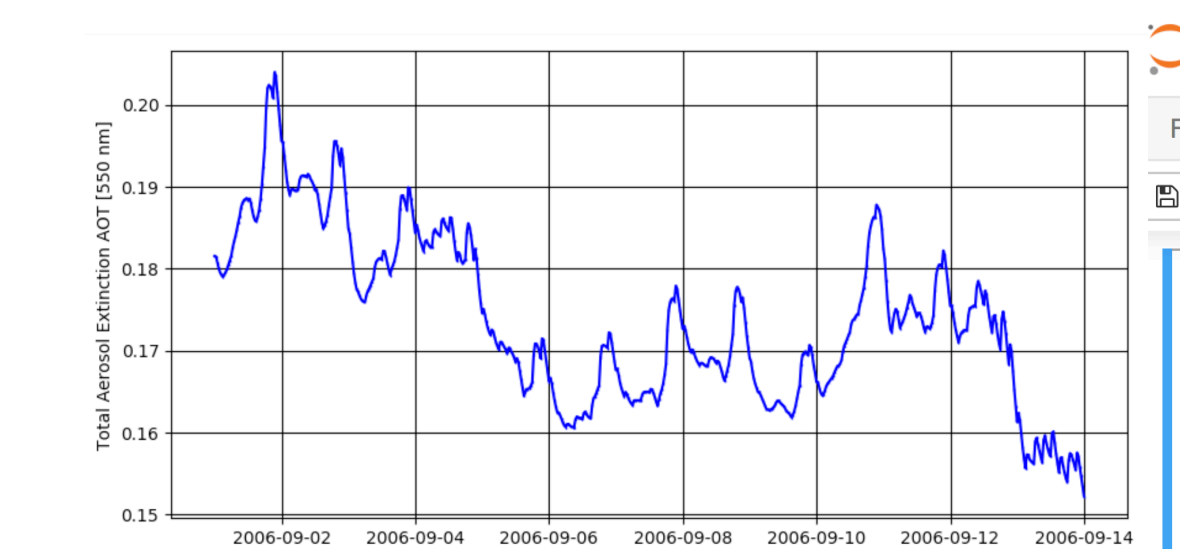
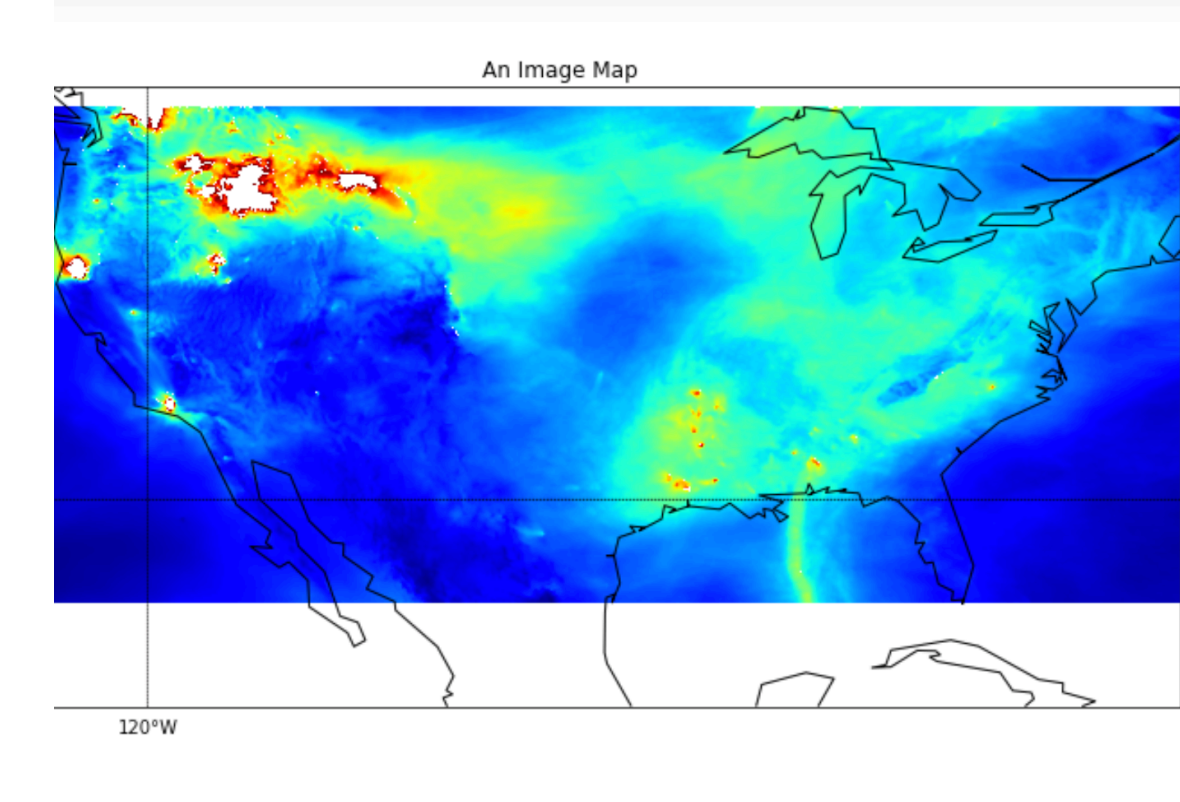
Correlation Map



More to come:

- Climatological Map
- Latitude-Time Hovmoller Map
- Longitude-Time Hovmoller Map
- Daily Difference Average

Notebook Analytics: CONUS Time Series and Time-Averaged Map



```
jupyter G5NR_Notebook (unsaved)
File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 O
In [11]: # 24 deg, X 55 deg, Continental USA subset
        bbox = box(-125, 25, -70, 49)

        # Compute time-averaged map using the SDAP/NEXUS web/HTTP interface
        start = time.perf_counter()
        dt_format = "%Y-%m-%dT%H:%M:%S"
        url = "{0}/timeAvgMapSpark?ds=G5NR_TOTEXTAU_500&mi
              nLon=-125.0&minLat=25.0&maxLon=-70.0&maxLat=49.0&startTime=2006-09-01T00:
              00:00&endTime=2006-09-14T00:00:00Z

        format(base_url, datasets[0], "bbox.bounds, start_time, strftime(dt_fo
              rm, print(url)
              print()
              tan = requests.get(url).json()[["data"]]
              print("Time averaged map took {} seconds".format(time.perf_counter() - start))

http://54.190.30.161:32139/nexus/timeAvgMapSpark?ds=G5NR_TOTEXTAU_500&mi
nLon=-125.0&minLat=25.0&maxLon=-70.0&maxLat=49.0&startTime=2006-09-01T00:
00:00&endTime=2006-09-14T00:00:00Z

jupyter G5NR_Notebook Last Checkpoint: Yesterday at 1:24 PM (unsaved changes)
File Edit View Insert Cell Kernel Widgets Help Trusted Python 3 O
In [19]: # 24 deg, X 55 deg, Continental USA subset
        bbox = box(-125, 25, -70, 49)

        # Compute area-averaged time series using the SDAP/NEXUS command line interface
        start = time.perf_counter()
        ts = nexutil.time_series(datasets, bbox, start_time, end_time, spark=True)[0]
        print("Area-averaged time series took {} seconds".format(time.perf_counter() - start))
```

Parmap Time Averaged Map Performance

G5NR (7 km 30 min) Global 30 Day Subset

