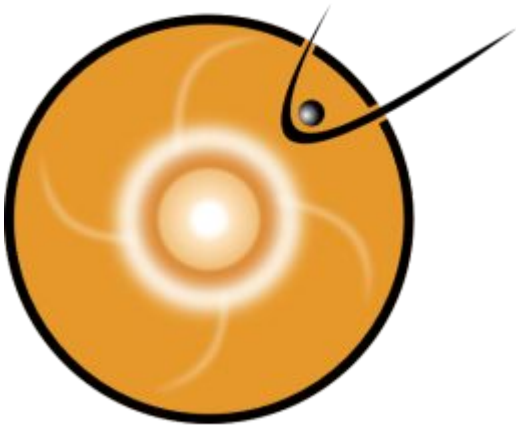# Streamlining Data Management and Operations with Airflow for Heliophysics Research and Space Weather Forecasting
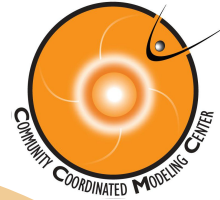
Presented by Damian Barrous-Dume

CCMC/Navteca LLC

Multi-agency strategic investment in US space weather program

## CCMC Goals

*Facilitate space weather & space science* **research & model development**

*Support transition of advances in research to* **space weather operations**

**Established in 2000**

Original team: Michael Hesse (founding director), Masha Kuznetsova (deputy), Lutz Rastaetter

First equipment: 3 Sun Workstations
First model: SWMF (U. Michigan)

# Models at CCMC

**Corona**
SWMF.SC+EEGGL+CME
AWSoM    EEGGL    SRPM
PFSS.Petrie   ANMHD
PFSS.Macneice   FLAMPA
PFSS.Luhman
MagPy
MAG4   UMASEP   SPRINTS-SEP
ASAP   ASSA   AMOS
WSA   NLFFF   GSU All Clear
MAGIC   SNB3GEO   FISM2
GCR   BON   NOVICE
NAIRAS   CARI-7

**Heliosphere**
WSA-ENLIL
WSA-ENLIL+Cone
WSA-ENLIL+EPREM
WSA-ENLIL+SEPMOD
HESPERIA REleASE
PREDICCS   EMMREM
SEPSTER
iPATH   ZEUS+IPATH
SAWS-ASPECS
CORHEL
CORHEL- CME
Heltomo IPS
GAMERA/Helio
DBM   SEPSTER2D
SWMF.SH
DIPS

**Magnetosphere**
MAGE/GAMERA+REMIX+RCM
LFM-MIX   GIC
OpenGGCM+CTIM
SWMF+RCM+deltaB
SWMF+RCM
SWMF+RCM+RBE
SWMF+RCM+CRCM
LFM-MIX-TIEGCM
WINDMI   LANLstar
IGRF   Tsyganenko
PS VP   Weigel-deltaB
AACGM   Apex
AMPS   GUMICS

**Local Physics**
VPIC
PAMHD
PIC-Hesse

**Inner Magnetosphere**
RCM
VERB
AMPS
Fok.CIMI
Li's Rad Belt
PINE   BSPM
UPOS RB
AE-8/AP-8
AE-9/AP-9
IMPTAM
RAM-SCB
SHELLS
ORIENT

**Ionosphere/ Thermosphere**
WACCM-X   WAM-IPE
SAMI3/WACCM-X
NCAR DART
GMAT   CTIPe
DTM2020   IDA4D
TIE-GCM   USU-GAIM
SAM   SWACI-TEC
ABBYNormal
NRLMSISE
SAMI-3   GITM
PBMOD
WBMOD
Weimer IE
Weimer-deltaB
IRI   JB2008
IMPACT
COSGROVE-PF
Ovation Prime
TRIPL-DA

# Services at the CCMC

We offer a range of run services that contribute to building a network of simulation services and run result archives. Our services support forecasting, model validation, education, and space science research.

We provide multiple different ways to perform model runs:

### Runs-on-Request

Long running simulations by request. Submit a request and get notified when the results are ready.

### Continuous Run

Continuously running models. View the latest results of several models without a separate request.

### Instant Run

Instantly executing models. Select a model and see the results instantly in your browser.

# CCMC's Hybrid Environment

We utilize a Hybrid architecture to help manage all of the different services we help support.

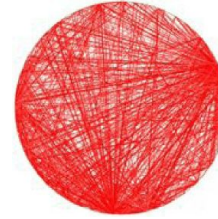But how do we manage all of the data?


Taylor-SQL
@tayloramurphy

you need DataOps so you don't make a DataOops

12:46 PM · Dec 19, 2021


HECC/Pleiades    On Premise

aws

CCMC Cloud

# Data and Operational Challenges



Big Ball of Mud of Pipelines
Look Familiar?

**Current Challenges**

- Manage multiple file transfer pipelines and syncs across environments.
- Improve visibility for scattered cron jobs to ensure proper execution.
- Increase overall reliability for users and data archives.

**What's the Impact?**

- Inefficient operations and delays in delivering critical data to the services that might need them.

# Future Challenges: The Reality of Big Data

## Big Data



**Volume**
The large amount of data generated and collected.

**Velocity**
The speed at which data is generated, collected, and processed.

**Variety**
The different types of data (structured, unstructured, semi-structured).

## More Data, More Costs



Storage Costs Smoothed by Year

- Trendline for Total Approx Storage Cost
- Hypothetical Linear Trendline

Storage Available by Year (Culmative)

- Total Available Capacity (TB)
- Trendline for Total Available Capacity (TB)

## Data Tiering



**Hot Data**
Frequently Accessed

**Warm Data**
Infrequently Accessed

**Cold Data**
Sporadically Accessed

Performance

Data Volume

# Role of Apache Airflow

**Key Benefits**

- Distributed workflow orchestration tool used by many organizations like AirBnb, Uber, Department of Energy Labs.
- Agents can run within servers and containers.
- Centralizes and schedules jobs across environments, which can typically be Python but supports other types depending on the Operator used.
- Efficient detection and alerting for pipeline failures.
- Allows non-airflow experts to manage, start, and restart the pipelines they own.

*Different Roles for Different Folks*

# CCMC Implementation of Airflow Demo

# Improves our interconnected networks of simulation services and run result archives

- Enhanced **monitoring and error detection** for improved reliability.
- **Standardized metadata** capture to build a metadata database for tiering and improved DataOps.
- Support for complex workflows and data transformations.
- **Improved data reliability and accessibility** for researchers and forecasters for all of the services we offer.
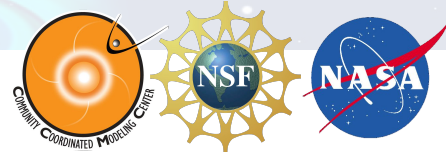
# Leading Open Science at NASA:
# CCMC's Commitment

- Pioneering the use of AWS since 2019 to build collaborative, distributed systems.
- Open-sourcing tools like **Kamodo** via NASA's open science initiative.
- Actively contributing to global open-source projects such as **OpenSpace**.

With this in mind, we're looking at how to share our distributed Airflow setup to help others facing similar challenges.
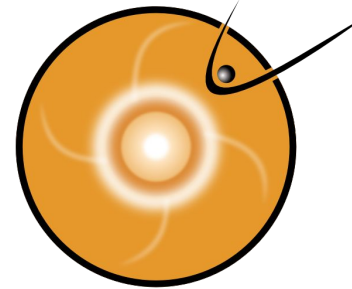
Reach out to us if interested!

# Future Outlook and Conclusion

- Continue migrating more jobs and file transfers to Airflow; collaborate with researchers to help them create and manage their own workflows.
- Develop a metadata database to enhance data captured through Airflow.
- Explore additional ways to improve the setup; Airflow remains a flexible tool within a centralized distributed workflow management system.
- Overall, Airflow has met our needs at CCMC and will be a key component in developing interconnected networks of simulation services and run result archives.

# Thank you
# for your time!
## Any Questions?

**Community Coordinated Modeling Center**

**For more information on CCMC, please visit our website:**

Contact Email:
**dbarrous@navteca.com**

NASA

Goddard
Space Flight Center