

Proudly Operated by Battelle Since 1965

ARM's Data Integration Framework

How ARM (Atmospheric Radiation Measurement) Focuses Its Resources on Science, Not Data Pre-processing

KRISTA GAUSTAD

GABE GIBLER

ARM Infrastructure

https://github.com/ARM-DOE/ADI

https://github.com/ARM-DOE/adi-macosx

Oct 2, 2023

ARM Data Challenge: Diverse Products Collected and Processed Continuously



- 3 fixed sites
- 3 mobile facilities deployed at over 36 locations
- Aerial deployments
- 460+ instruments of approx. 100 types
- 250+ higher level algorithms (VAPs)
- 200–2000 Mb data per day
- 20k–40k files processed per day
- Spanning 30+ years





ARM Data Challenge:

Processing Complex, Interdependent Data

- Ease process of integrating data from diverse datasets with different sampling intervals.
- Simplify development of scientific algorithms.
- Ensure adherence to data standards.





ARM Data Integrator (ADI)





PCM (Process Configuration Manager) User Interface







- Automate repetitive data tasks:
 - Reading input data
 - Transforming and consolidating data
 - Writing data to output
 - Applying or generating quality control flags
- Embed data standards into data design interface.
- Reduce cost, complexity of ingest and scientific algorithm development.
- Improve software robustness and reliability.
- Support multiple programming languages (Python, C).
- Support local installations and open-source code.

ADI Architecture





PCM Process: Overview



Proudly Operated by Battelle Since 1965



3: Run Data Consolidator

Output QC ∨

4: Create ADI Project

Revision History

PCM Process: Retrieval Rules



Proudly Operated by Battelle Since 1965

🕫 Process: doe_data_days

📲 | Mode: Standard 🗸 | 🖺 🗸 | 💰 |

1: Describe Process	0 Re	trieval Rule S	iets (4) 🗙	+ Add Rule	Set							
	-	» Name	» Rules	5					» Varia	ble Name in Input	Datastream	
2: Define Inputs & Outputs		coil						~ +				~ 1
Overview	U	Cell			1929						1223	
Inputs				Deineiter	If process runs	Deturen	And	Retrieve		Assisted Name	ceil.b1	ceil.b1
🔺 Variables 🗸 🛛 🔸				Phonty	ac	Between:	And:	aata from:	U	Assigned Name	(priority 1)	(priority 2)
Retrieval Rules				1	SGP E13	Any	Any	SGP C1 >		backscatter	backscatter	backscatter
Input QC						ume	ume	CEILDI		1		
Coordinate Systems +				2	Any	Any	Any	Same				
half_min_grid: time					sites/facilities	time	time	location				
half_min_grid: range								as				
metb1: time								line >				
Outputs								ceil.b1				
Output Datastreams 🗸 🛛 🕂				1		-	1					
Output QC 🗸			Offse	ts: 0 Star	t: 0 seconds 0 E	nd: 0 secon	ids					

Preferred and alternate data sources by:

- Datastream
- Location
- Time

User-defined variable names can be assigned to the inputs.

PCM Process: Data Consolidation



Proudly Operated by Battelle Since 1965

📲 | Mode: Standard 🗸 | 🖺 🗸 | 💰 🥰 Process: doe_data_days Dimension Shape A 1: Describe Process How should time coordinate data be determined? Regular Interval A 2: Define Inputs & Outputs time values will be defined using the following parameters: Overview Interval: 30 6 Start: 0 Inputs 0 End: 86370 ▲ Variables ∨ O Number of intervals: 2880 **Retrieval Rules** Input QC **Coordinate Systems** Transformation half_min_grid: time O What algorithm should be used to transform variable data dimensioned by time? half_min_grid: range Auto (average/interpolate) metb1: time O Bin Averaging O Bilinear Interpolation Outputs O Nearest Neighbor (i.e., subsample) Output Datastreams ∨ O Passthrough (i.e., no change in values) Output QC V Do you want to override this value by variable? O Yes No 3: Run Data Consolidator 0 Output data will be moved to center bin alignment. Do you want to change this setting? 4: Create ADI Project Bin Alignment: Start Revision History O Do you want to adjust the output bin width? O Yes I No O When reading the input data, how far should the interpolator look for the next valid data point in the time dimension (i.e., the range)? Default range for all applicable datastreams: 300 seconds Do you want to override this value by datastream? () Yes () No Do you want to override this value by variable? O Yes O No

PCM DOD: Data Output Design

Onventions

process_version



Proudly Operated by Battelle Since 1965

ARM DOD Manage	r ▼		🕇 New 🚽 🔔 Im	port -				() L	ogout (gaustad) 🖂	Contact Us 😯 Help	
DODs	🔹 🔺 kigt	tutor	ial30s.c1 v1.0 ×								
tutorial		ວ	С 🛯 - Т	Mode 👻	ଓ ≓ ।	🗹 - I 🖉 I 🔁				O Submit for Review	
▶ Filters	•	um	ensions (3)							+:	
Advanced Search		Dimension Length									
		time	9		Unlimi	ted +					
tutorial 😒		bound				Value - 2					
Showing 25 of 4955:		rang	ge		Value	- 801					
ectutorial30s.c1											
jwmtutorial.c1	• •	Vari	ables (13) F	ilter			Config	ure visible columns 👻		+:	
▼ 🦲 klgtutorial30s.c1	L 1	•	name	type	dimension(s)	long_name	units		standard_name	ancillary_variables	
DOD v1.0	ø	0	base_time	int		Base time in Epoch	seconds 0:00	since 1970-1-1 0:00:00	<att not<br="">Defined></att>	time_offset	
qitutorial.c1		1	time_offset	double	time	Time offset from base_time	<set at="" f<="" td=""><td>Runtime></td><td><att not<="" td=""><td>base_time</td></att></td></set>	Runtime>	<att not<="" td=""><td>base_time</td></att>	base_time	
qitutorial30s.c1		~							Defined>		
rkntutorial.c1	• *	2	time	double	time	Time offset from midnight	<set at="" h<="" th=""><th>Runtime></th><th>time</th><th><att defined="" not=""></att></th></set>	Runtime>	time	<att defined="" not=""></att>	
rkntutorial30s.c1		3	time_bounds	double	time, bound	Time cell bounds	<att not<="" th=""><th>Defined></th><th><att not<br="">Defined></att></th><th><att defined="" not=""></att></th></att>	Defined>	<att not<br="">Defined></att>	<att defined="" not=""></att>	
sltutorial.c1	1	4	range	float	range	Distance to the center of the	m		<att not<="" td=""><td><att defined="" not=""></att></td></att>	<att defined="" not=""></att>	
sltutorial30s.c1						corresponding range bin			Defined>		
tutorial.c1	1	5	range_bounds	float	range, bound	Coordinate_variable cell bounds	<att not<="" th=""><th>Defined></th><th><att not<br="">Defined></att></th><th><att defined="" not=""></att></th></att>	Defined>	<att not<br="">Defined></att>	<att defined="" not=""></att>	
tutorial30s.c1	A /	6	temperature	float	time	Air temperature	к		air_temperature	qc_temperature	
tutorial30smgd.c1	A /	7	gc_temperature	int	time	Quality check results on field: Air	unitless		<att not<="" td=""><td><att defined="" not=""></att></td></att>	<att defined="" not=""></att>	
tutorialcmh.c1						temperature			Defined>		
tutorialcmh30s.c1	A Ø	8	backscatter	float	time, range	Backscatter	1/(sr*km	*10000)	<att not<br="">Defined></att>	qc_backscatter	
tutorialidr.c1	A /	9	qc_backscatter	int	time, range	Quality check results on field:	unitless		<att not<="" th=""><th><att defined="" not=""></att></th></att>	<att defined="" not=""></att>	
tutorialmgd.c1		10	1-4	flt		BackScatter			Defined>		
tutorialst.c1		10	iat	float		North latitude	degree_l	-	latitude	<att defined="" not=""></att>	
tutorialtsq.c1		11	lon	float		East longitude	degree_l	-	longitude	<att defined="" not=""></att>	
tutorialzcg.c1	1	12	alt	float		Altitude above mean sea level	m		altitude	<att defined="" not=""></att>	
wigtutorial.c1	-	Glob	al Attributes (1	3)							
wigtutorial30s.c1		Attr	ibute	-,		Туре		Value			
xctutorial.c1	I '	con	nmand_line			char		<set at="" runtime=""></set>			

- xctutorial30s.c1
- zftutorial.c1
- zftutorial30s.c1

ARM-1.1

<Set at Runtime>

char

char

PCM DOD: Data Output Validation



	DOD Manage	er 🔻	🕂 New 👻 📩 Import 🗸								
DODs		• 🛦 kigt	Itorial30s.c1 v1.0 ×								
tutorial Validation Report for klgtutorial30s.c1: 1.0											
▶ Filters											
Advanced Search		3 errors, 9 w	varnings, 0 disallow (10 can be auto fixed; 0 are ignored))		_					
tutorial O		Туре	Message	Field							
Showing 25 of 4955:	0	Warning	The calendar variable attribute is optional if the calendar used is Gregorian.	time.calendar	1	1					
ectutorial30s.c1	A	Error	"cell_transform" variable attribute value should generally be null, because it is set automatically by ADI.	temperature.cell_transform backscatter.cell_transform	1	1					
 klgtutorial30s.c1 DOD v1.0 qitutorial30s.c1 	•	Warning	The variable attributes "valid_min", "valid_max" and "valid_range" are intended to describe physical or mathematical limits. For QC test limits, it is suggested to use different attributes, located on the QC variable instead of on the data variable. Suggested attribute names include: "fail_min", "fail_max", "fail_range", "warn_min", "warn_max", "warn_range".	temperature.valid_min temperature.valid_max		%					
rkntutorial.c1	0	Warning	Use of the value "unitless" is deprecated. Use "1" instead.	qc_temperature.units qc_backscatter.units	×	1					
 rkntutorial30s.c1 sltutorial.c1 	0	Warning	QC variable must have a long_name of "Quality check results on variable: <data variable's<br="">long_name attribute value>".</data>	qc_temperature.long_name qc_backscatter.long_name	1	ø					
situtorial30s.c1	0	Warning	It is recommended that all QC variables have a "standard_name" attribute with a value of "quality_flag".	qc_temperature.standard_name c_backscatter.standard_name							
 tutorial30s.c1 tutorial30smgd.c1 tutorialcmh.c1 	•	Error	QC variables using flag_method "bit" and described at the field level must have a description of "This variable contains bit-packed integer values, where each bit represents a QC test on the data. Non-zero bits indicate the QC condition given in the	qc_temperature.description qc_backscatter.description		Ø					
 tutorialcmh30s.c1 tutorialldr.c1 		CSV)		Autofix All Solution Ignore All Viol	ations	Close					
tutorialmod et		-	Bac	kscatter	111033						





Decreased development time and cost.

- CSV ingest decreased up to 80% (from weeks to days).
- VAP algorithm decreased up to 60% (from several to few months).
- Improved reliability through code-reuse.
- Software maintenance simplified.
- The processes that generate our data products are more accessible.
- Production processing and reprocessing streamlined.