# OOI CI Data Team Update September 2017

M. Vardaro, M. Crowley, L. Belabbassi, L. Garzio, J. Kerfoot, F. Knuth, S. Lichtenwalner, M. Smith **Rutgers University, New Brunswick, NJ** 







### Overview

- 1. Introduction to OOI
- 2. Data Flow & Products
- 3. Data Review Procedures
- 4. Periodic Reviews & Documentation
- 5. Next Steps
- 6. Conclusions
  - a) A large amount of high value data is being collected
  - b) Data ingestion & review is our primary focus
  - c) Accelerating data review via development of specialized tools
  - d) Short-, medium-, and long-term goals to improve data quality
  - e) OOI is providing a curated, consistent data system that is delivering data and metadata to the community







### OOI By the Numbers









## **Rutgers CI Team**

Leadership & Oversight: 0 FTE

OOI Management: 1.0 FTE (ECR 1300-00525)

OOI Operations: 12.5 FTEs







### **Timeline Context**









### Data Team Primary Goals:

- 1. Monitor the operational status of data flowing through the OOI system
- 2. Ensure the availability of OOI datasets in the system (raw, processed, derived, and cruise)
- 3. Ensure that data delivered by the system meets quality guidelines
- 4. Identify availability and quality issues and ensure they are resolved
- 5. Communicate known data issues with end users
- 6. Report operational statistics on data availability, data quality, and issue resolution





# Day in the life of an OOI data evaluator



- Evaluator is assigned a specific array and works closely with the MIOs
- Range of expertise (biology, physics, geophysics)
- Data community including collocated non-OOI data experts
- Development of open-access tools to visualize and synthesize the data
- Developing inputs for automated QC
- Quick looks and deep data dives, updating OOINet asset sheets
- Interactions within team, with MIOs, with users & students, and specific SMEs
- Full-time effort is required





### Overview

- 1. Introduction to OOI
- 2. Data Flow & Products
- 3. Data Review Procedures
- 4. Periodic Reviews & Documentation
- 5. Next Steps
- 6. Conclusions







### Data Flow Example: Pioneer Profiler







### **Current Data Processing Flow**







### Data Types

- Telemetered Data
  - Data received through a transmission medium over distance (e.g. surface buoy to satellite, glider to satellite, acoustic modem); may be decimated
- Recovered Data
  - $\circ$  Data downloaded directly from a recovered instrument or data logger after the instrument has been recovered.
- Streamed Data
  - Data received via transmission over electro-optical cable. Streaming data are provided at full temporal resolution and near-real time.
- Shipboard Data
  - Shipboard data and water samples collected during OOI expeditions.
- Metadata
  - Info about the data record (e.g., time & location of collection, unique source & record description identifier, instrument serial #, etc.). OOI metadata follows the CF1.6 standard, with additional types and fields specific to OOI as necessary.





### **OOI** Data Product Levels

- Raw data: The datasets as they are received from the instrument
  - o May contain multiple L0, L1, or L2 parameters, data for multiple sensors, and be in native sensor units
  - Always persisted and archived by the OOI
  - Example: format 0 binary file from an SBE-37IM on a Global Flanking Mooring.
- Level 0 (L0): Unprocessed, parsed data parameter that is in instrument/sensor units and resolution
  - Sensor by sensor (unpacked and/or de-interleaved) and available in OOI supported formats (e.g., NetCDF)
  - $\circ~$  Always persisted and archived by the OOI
  - <u>Example</u>: SBE-37IM Temperature portion of the hex string
- Level 1 (L1): Data parameter that has been calibrated and is in scientific units
  - QC may be applied at this level, utilizing simple automated techniques or human inspection
  - Actions to transform Level 0 to Level 1 data are captured and presented in the metadata of the Level 1 data
  - Example: SBE-37IM Temperature converted from hex to binary and scaled to produce degrees C
- Level 2 (L2): Derived data parameter created via an algorithm that draws on multiple L1 data products
  - Products may come from the same or from separate instruments
  - $\circ~$  Data from all relevant instruments will be provided during download
  - <u>Example</u>: SBE-37IM Density and Salinity



### **OOI: Web Portals**







0	)	ER	DDAP	ncient fica	m							Broughtine you by NEAA MILES SATSO 1920
RD ck a	D/	<u>AP</u> atas	> Lis et	t of <i>I</i>	All Datasets							On the a Full Test Search for Desame: On Search for Search by Catagory, of all the search by Catagory, of all the search by Catagory, searched by Search for Catagory and Search Search Search On Search for Catagory with <u>Search Search Search</u> ®
H Su Su	et 1	DAP Date	Make W A N Graph S	Seama Data Filos	Title	Sum- mary	PODC, 150, Metadata	flack ground http	R55	1 ma	Institution	Dataset ID
21	8 :	2012	9/223		*The Ust of All Active Datasets in this ERODAP*		2	beconund			Ratgers Univers	alColasets
		data	9(303		05140 GPG-PC016-4A-CTDHTA108-streamed-ctdpf_optode_sample	ø	111	tectorium @	0.000	22	Ocean Observato 6	CEE405P6-P0318-44-CTCPFA188-streamed-cttp1_spixde_sample
		data	8923		CER4CSPS-FCI16-4C-PC02WA105-streamed-pcc2w_a_sem_data_record		ELB.	lazaran #	8162	B	Occan Ozacivala 6	CEL405PS-PC018-4C-PCC2014135-all carries-pco2e_a_sam_data_record
		data	92323		CE04CSP9-SF018-2A-CTDPFA107-streamed-ctdpf_ster8]_sample		ELH	tacionund @	C.L.	22	Ocean Observato 6	CEE406P6-5F019-2A-CTEPFA107-streamed-cttp1_ste42_sample
		data	97923		CENCOPS. SFDID-25. PRODUINIO-streamed-pheen_date_record	69	E1B	becomund @	2000	25	Ocean Observato 6	CEECOSPS. 57010-20-PHSCNA133-atreamed.phsen_data_record
		data	0/453		60940695 SEDIO 3A FLOREDIOI steased first_6_data_second	69	111	backymand #	12000	63	Ocean Observals	CEEKOSPS ST010 3A FLERTD134 streamed fort_6_082_record
		data	02823		CE04CSPS-SP010-4A-NUTNRA102-streamed-ratin_a_semple		KTR.	background #	1.69	8	Ocean Observato 6	C214CSPS-SP315-44-NUTNRA112-atrianed-cutre_x_sample
		data	82.823		CE04CSPS-SP018-4P-PC02044102-alreamed-pos2w_ac_eart_chala_record		2.18	background @	0.000	8	Ocean Observale 6	CEL4OSPS-SH18-4H-PC02XA302-almamed-pcs2w_a_aam_tala_record
		100	2223		CERECENT-RESIDENCE/FRAME-RESIDENCE-SEC_RE_RE_INFORMATION		ELB	torrand@	C.C.C.	63	Ocean Opacrysta 6	CELEOSSI-ROSS-04-VELPTABLO-LIKerkerkerkerkerkerkerkerkerkerkerkerkerke
		2012	21223		05140/0514-RD26-06-PH/SEND000-televenered-phaser_abodef_pol_instrument	0	ELB	toportund @	C.L.C.	62	Ocean Observato 6	CEE4055W-RC06-06-PHSEND003-bioinetored-phsen_abodef_00_instrument
		2052	9(323		050406504-R026-07-X4TNRB300-telenetered-nutric_b_dd_cono_instrument		ETR.	teconuel #	Catalog	22	Ocean Observato 6	# CECKOSSIX-RC06-07-NUTNRECCO-telenetered-nutry_b_dcl_core_instrument
		2404	0.423		010400538 RE27 02 FL 0RT0000 relevened flot_d_dcl instrument		111	tabigmint6	10000	2	Ocean Observato	CERIOSSIE ROOT 02 FLORTDOOD selemetered field (d) instrument
		2022	2223		CED4C10384-8022-403-C RDB+C000-ceenecored-ctdbp_cdef_dcl_instrument		218	terrorand @	1000	63	Cosan Observato 6	CERRORAMMONATION-CONTRAGENCE-CROST-Contragence - Crost-Crost-Section - Section - Contract - Cont
		data	97.923		000400014 RD27-04-DC0TA/D000 relevanced-dosta_shocks_dd_instrument	69	E18	becoround @	C.L.	R	Ocean Observato 6	# CEE4050/#.RI027-04.D00TX/D003 Internetined-doels_sbodys_ds(_ketument
		5854	97.833		CONCOM SECTION OF THE DANAGED IN A REAL PLACE AND A REAL PARTY AND A REAL PLACEMENT		118	Nonpoint #	1000	2	COMM COMMAND	ECOLOGICAL SCIENT IS: INTOINECCO INNERNEED By E. E. O.S. INVERSENT
		2853	02303		CONCOMPORTATION OF VELYTY/COD HIM MAINED WIRE JAD_DD_IMDFURWED		1.18	tacopued 6	1000	250	COLLET OSCENSIS	CEDACODAR-SOBITI-ON-VELIPINEED-GRAMMARKO-URDE_AD_GCI_KARSUNATE
		data	82.823		CE04CSSN-SE011-06-NE18C4000-skenetered-metol_a_dd_mahamani		2.18	Inclanated 6	U.L.C	83	Ocean Opervala 6	# CEL4OSSII-SED11-96-WETBKA000-Informational-mailak_ut_dol_mainsment
		2022	2223		CEDECESS-SED12-25-BY/CONCORPORTED-Spic_SCAL_matrixed		218	torund #	L.C.S.	02	Occan Opportato 6	CELEOSON-DEDT245-HYDORODON-Deventored-454_s_ds_struct
		2222	2223		CODECCENTER CODECCENTRY CODECCENTER CODECC		118	toroiuni @	L.C.	0.0	Cosan Cosanvalo 6	<ul> <li>CELEVICER CODARD CODE CODE CODE CODE CODE CODE CODE COD</li></ul>
		2222	9.223		CONVERSE OF A REPORT OF A DESCRIPTION OF		110	101001010	CHERRY	52	Cost Cost a	CONCEPTED A STATE AND A STATE
		2804	07833		Distances (e.a) of Prevenues Sentences parts in giver restorer		1.1.1	tao gauna	1000	1	DOMET DISCENSES	CERTIFICATION IS 311 D1 HORIZONDO SARTHARING DAMA IN DIDAL INCIDINATI
		-	02.833		Contraction of the Party and the second state whether a state being state			tablem of C	College of the local division of the local d	100	Other Observate	CERTIFICATION OF A DESCRIPTION AND AND AND AND AND AND AND AND AND AN
		444	0000		CONTRACTOR OF THE CONTRACTOR O		111	hadron of d		10	Ocean Observato	CERTING & MILLS CONTRACTOR STREET ALL CONTRACTORS
			9.434		Contractor and an exception of the second seco			A second	Carrow	60	Oran Observats	Contraction of the Contraction of the State of the State of the
			and a		Part of the second		111	an capitol a	Correst of the local division of the local d	100	Oran Character 6	CALCULATION OF THE ACCOUNTS AND ADDRESS OF THE ADDR
		200	12.A13		CEV/SYSW-Prose-cev/recorded and an and an and and and and and and			CALCULUM &	Caller Street	100	Control Control of	CEUSISSI-IPUS-IS-U2/IBW-INTERIOSPEZ/CEUCID
			2223		CONTRACTOR OF A CONTRACTOR OF		1.1	Manager and a	Contrast.	22	Contra Contrado	Construction of a provide the second second second second
		2002	9.213		CERTIFICATION AFFORTANCE CONTRACTOR AND A STORE		111	10202030	Canada	0.4	Cossi Osservasi	CEEPOID AND COURSE COURSES AND
		222	9.223		Contrance and contract contracted biser_abolat_co_netwinet		111	through the	1000	10	Cost Cost Nation	<ul> <li>Conversion and conversion and the state of t</li></ul>
		data	07303		07075HSM WF007 03 C700P0000 televisiwed citiba_cdef_dd_instrument		C.1.H	background #	2,572	$\simeq$	Ocean Observatio 6	CEETSUSE MITOR AD CTERPEORD relevated cidite_colef_doj_instrument



iompany Home > COI > Cabled Array	/ ➤ Cruise Data				
Cruise Data This view alows you to browse to Shipboard Data from Cabled Arr	he items in this space. ay Cruites		<b>3</b> (0)	Add Content Creater More Acti	ons 🐑 😑 Details View 💌
Browse Spaces					Items Per Page 200
Name A	Description			Created  Mo	dified a Actions
Cabled 1_TN-221_2008-7-22				28 October 2015 14:15 28 0	Xtober 2015 14:15 🗋 📄 💽
Cabled-2_TN-252_2010-7-26				28 October 2015 14:15 28 0	October 2015 14:15 🗋 📄 🐨
🕼 Cabled-3_TN-268_2011-8-11 🛈				28 October 2015 14:37 28 0	October 2015 14:37 🗋 📄 🐨
🐊 Cabled 4_TN-299_2013-06-30 🤇	2013 OCI Cabled Array dep	loyment cruise, R/V 1	Thompson (TN299), June 30-Aug	ast 23, 2013 24 October 2015 16:00 16 1	tay 2016 15:40 👔 🖻 💌
Cabled-5_TN-313_2014-7-13				28 October 2015 15:00 28 0	xtober 2015 15:00 🗋 📄 💌
Cabled-6_TN-326_2015-7-04				28 October 2015 15:16 28 0	October 2015 15:16 🗋 📄 🐨
		Page 1	of 1 14 4 1 1 14		
Content Items					Items Per Page 200
Name 🔺	Description	Size 🗰	Created @	Modified a	Actions
		4 KB	16 May 2016 10:13	16 May 2016 10:13	
DS_Store					







### Overview

- 1. Introduction to OOI
- 2. Data Flow & Products
- 3. Data Review Procedures
- 4. Periodic Reviews & Documentation
- 5. Next Steps
- 6. Conclusions







### First in Class Reviews: Jan-Aug 2016

- One example of each data stream (ingestion completed by Systems team)
- Review of 1207 (467 science) streams completed in August 2016
- Tested parsers, algorithms, ingestion, asset management and data product creation

	WBS	Task Name	%	Duration	Start	Finish
	*	•	.omplet(+	•	•	*
299	1.5.1	Data Ingestion	62%	130 days	Wed 1/20/16	Tue 7/19/16
300	1.5.1.1	First In Class for Cassandra Team	88%	62 days	Mon 1/25/16	Tue 4/19/16
301	1.5.1.1.1	Pioneer Coastal Glider, CP05MOAS-GL388	100%	8 days	Mon 1/25/16	Wed 2/3/16
308	1.5.1.1.2	Pioneer Central Inshore Profiler Mooring, CP02PMCI	100%	8 days	Thu 1/28/16	Mon 2/8/16
315	1.5.1.1.3	Endurance OR Offshore Surface Mooring - CE09OSSM	100%	30 days	Fri 1/29/16	Thu 3/10/16
322	1.5.1.1.4		84%	46 days	Thu 1/28/16	Thu 3/31/16
329	1.5.1.1.5	Cabled Slope Base Shallow Profiler Mooring - RS01SBPS	100%	43 days	Thu 1/28/16	Mon 3/28/16
336	1.5.1.1.6	Cabled Slope Base Deep Profiler Mooring - SRS01SBPD	100%	31 days	Fri 1/29/16	Fri 3/11/16
343	1.5.1.1.7	Irminger Sea Apex Surface Mooring, GI01SUMO	100%	36.95 days	Tue 2/2/16	Wed 3/23/16
350	1.5.1.1.8	Irminger Sea Apex Profiler Mooring (GI02HYPM)	100%	33.5 days	Tue 2/2/16	Fri 3/18/16
357	1.5.1.1.9	Imminger Sea Flanking Subsurface Mooring A (GI03FLMA)	63%	33 days	Wed 2/3/16	Fri 3/18/16
364	1.5.1.1.10	֎ Irminger Global Open Ocean Glider (GIO5MOAS-GL)	100%	19.33 days	Wed 2/3/16	Tue 3/1/16
371	1.5.1.1.11	Irminger Global Profiling Gliders (GI05MOAS-PG)	100%	31.5 days	Thu 2/4/16	Fri 3/18/16
378	1.5.1.1.12	Coastal Endurance OR Inshore Surface Piercing Profiler Mooring (CE01ISSP)	100%	30.5 days	Fri 2/5/16	Fri 3/18/16
385	1.5.1.1.13		31%	30 days	Wed 2/10/16	Tue 3/22/16
392	1.5.1.1.14	Cabled Seafloor Instruments	0%	12 days	Mon 4/4/16	Tue 4/19/16

	WBS	Task Name	% `amplat.=	Duration	Start	Finish
	•	•	.ompieu+	·	·	•
418	1.5.2	Data Verification & Validation	28%	255 days	Wed 2/10/16	Tue 1/31/17
419	1.5.2.1	First in Class	39%	123 days	Tue 3/1/16	Thu 8/18/16
420	1.5.2.1.1	Pioneer Coastal Glider, CP05MOAS-GL388	95%	34 days	Tue 3/1/16	Fri 4/15/16
421	1.5.2.1.2	Endurance OR Offshore Surface Mooring - CE09OSSM	70%	30 days	Fri 3/4/16	Thu 4/14/16
422	1.5.2.1.3	Pioneer Upstream Inshore Profiler Mooring, CP02PN	42%	31.8 days	Fri 3/18/16	Mon 5/2/16
423	1.5.2.1.4	Cabled Slope Base Deep Profiler Mooring - RS01SBPI	0%	20 days	Tue 6/7/16	Mon 7/4/16
424	1.5.2.1.5	Cabled Slope Base Low Power Jbox - RS01SLBS-LJ01A	0%	27 days	Fri 4/29/16	Mon 6/6/16
425	1.5.2.1.6	Cabled Slope Base Shallow Profiler Mooring - RS01SI	75%	33 days	Tue 3/15/16	Thu 4/28/16
426	1.5.2.1.7	Irminger Sea Apex Profiler Mooring (GI02HYPM)	100%	20 days	Fri 3/25/16	Thu 4/21/16
427	1.5.2.1.8	Irminger Global Open Ocean Glider (GIO5MOAS-GL)	0%	3 days	Fri 4/22/16	Tue 4/26/16
428	1.5.2.1.9	Irminger Sea Flanking Subsurface Mooring A (GI03FL	0%	18 days	Wed 4/27/16	Fri 5/20/16
429	1.5.2.1.10	Irminger Global Profiling Gliders (GI05MOAS-PG)	0%	5 days	Mon 5/23/16	Fri 5/27/16
430	1.5.2.1.11	Irminger (Or other global) Sea Apex Surface Mooring	40%	74 days	Mon 3/14/16	Thu 6/23/16
431	1.5.2.1.12	Coastal Endurance OR Inshore Surface Piercing Profi	0%	8 days	Fri 6/24/16	Tue 7/5/16
432	1.5.2.1.13	Coastal Endurance OR offshore BEP - CE04OSBP	0%	22 days	Fri 7/8/16	Mon 8/8/16
433	1.5.2.1.14	Cabled Axial Seamount Central Caldera Med Power J	0%	10 days	Fri 7/8/16	Thu 7/21/16
434	1.5.2.1.15	Cabled Seafloor Instruments	0%	20 days	Fri 7/22/16	Thu 8/18/16
435	1.5.2.1.16	AUVs	0%	6 days	Tue 8/9/16	Tue 8/16/16
436	1.5.2.2	HAGU Oceans Data Prep (THREDDS & GUI) - Reasonability	100%	30 days	Wed 2/10/16	Tue 3/22/16







### **Data Annotation**

- Annotations are the primary means of communication between data team and users
- Annotations can be directly entered via the GUI for specified data streams
- Annotation text appears in a tab on the data catalog/plotting page
- Annotation time ranges can be shown on plots (via "Options" interface)
- Annotations also included in downloaded data

Data Navigation	Plotting		Events		An	notations
ow Additional Parameters						
led Axial Seamount Axial Base Seafloor - Low-Power JBox (LJ03	A) - CTD strea	amed-ctdpf-optode-sample	X Y	🛗 August 22	2, 2014 - November 30, 2016	•
Ne, UTC		•				
inity Corrected Dissolved Oxygen Concentration, Umol/Kg				Plot Type:	X-Y	•
			)	Plot Style:	Scatter	•
ick To Add Another Parameter Input				Orientation	Horizontal	•
				Options:	Show Annotations	3 <b>v</b>
				QAQC:	Global Range	 • ]
	FICE			$\mathbf{\mathbf{N}}$		
Cal	oled Axial Seamount Axial Bas	GRAPH se Seafloor – Low–Power JI CTD	Box (LJO3	8A) – CTD		
750 —	oled Axial Seamount Axial Bas	GRAPH se Seafloor – Low–Power JI CTD	Box (LJO3	A) – CTD	Annota	tion ID 0
Cal	oled Axial Seamount Axial Bas	GRAPH Se Seafloor – Low-Power JI CTD	Box (LJO3	BA) – CTD	Annota	tion ID 0
Cal	oled Axial Seamount Axial Bas	GRAPH se Seafloor – Low–Power JI CTD	Box (LJO3	(A) – CTD	Annota	tion ID 0
Cal	oled Axial Seamount Axial Bas	GRAPH Se Seafloor – Low-Power JI CTD	Box (LJO3	ia) - ctd	Annota	tion ID 0

Annotation ID	Annotation	Reference Designator	Stream Name	Start Date	End Date	Exclude Data?
0	These data are suspect, possibly due to incorrect vendor calibration values. Raw phase data should be correct, but the derived O2 products should not be used from 7/12/16 onwards.	RS03AXBS-LJ03A- 12-CTDPFB301	streamed_ctdpf- optode-sample	Tue, 12 Jul 2016 00:00:00 GMT	Thu, 01 Dec 2016 23:41:00 GMT	false







## Current "Rest-in-Class" Reviews

#### Process:

- Check all deployments for presence & absence of all parameters
- Check science parameters for reasonableness
- Problem? Deep dive, report in Redmine, track, give feedback, check fixes, create annotations in QC Database

#### Challenge:

- Automated tools, Redmine questions, Cal sheets, raw data repository, modify ingest CSVs, testing UI fixes
- Upload and ingestion of data
- Delivery and archiving of Cruise Data
- Quality Assurance vs. Quality Control

#### **Expediting the Solution:**

 Populate QC database to automatically check for presence/ absence, gaps > 1 day, NaNs, negative values

- 1. Asset Management (MIOs & Data Team)
- Complete?
- Correct?
- 2. Data Delivery & Ingestion (MIOs, Systems, Data Team)
- Includes Cruise Data
- 3. Data Review
- Availability
- Quality

4. Investigate Gaps and QC failures

5. Communicate Issues (Annotation)





### **Rest in Class Data Review Workflow**







### **OOI** Automated QC Procedures

- 6 automated QC algorithms can produce 7 flags (including logical "or" which combines flags) which are plottable and are included in downloaded files
- Coded based on specifications written by OOI Project Scientists, derived from QARTOD manuals and other observatory experiences
- Algorithms refer to "lookup tables" assembled by OOI Project Scientists with input from subject matter experts: <u>https://github.com/ooi-integration/qc-lookup</u>
- 1. Global Range Test
- 2. Local Range Test
- 3. Spike Test
- 4. Stuck Value Test
- 5. Trend Test
- 6. Temporal Gradient Test
- 7. Spatial Gradient Test (Profile)







### **Rest in Class Data Status Categories**

Status	Description	QARTOD Code	QARTOD Description	Color
NOT_OPERATIONAL	Instrument not functional (no data expected)		Not operational	
NOT_AVAILABLE	Instrument functional, data lost in transmission	9	Missing data	
PENDING_INGEST	Instrument functional, data exists, Awaiting ingest			
NOT_EVALUATED	Instrument functional, data exists, Awaiting evaluation	2	Not evaluated, not available or unknown	
SUSPECT	Instrument functional, data exists and either failed a QC test or does not reflect environmental conditions	3	Questionable/suspect	
FAIL	Instrument functional, data exists but is known to be bad due to known instrument or calibration error	4	Bad	
PASS	Instrument functional, data exists, passed QC tests, is complete and looks reasonable	1	Good	
GOOD	Instrument functional, data exists, passed QC tests, is complete and has undergone validation with shipboard datasets and reached the highest level of QC that the OOI can provide			





# Data Availability

- 87.5% of all actively deployed platforms are providing all available data from the most recent deployments (1 Endurance glider and 3 Global platforms need updating)
- 75.3 Gb of processed data delivered in the last year
- 41.9 Tb of raw data delivered in the last year
- Recovered data and backlogged telemetered data are currently being ingested into the system sequentially
- Aiming to present full timeline "heat map" display by end of November early December 2017





# Ingestion

• Testing purge/ingest for complete deployments to generate estimated ingest schedule for all data sets

#### Dependencies

- Delivery of deployed and recovered data to the raw data server at Rutgers: mostly up to date (latest check is ongoing).
- Ingest currently requires a follow-up query to determine that ingestion has completed
- Playback functionality required for purge/reingestion of cabled data (fill gaps and improve provenance): PRIORITY
- Purge/reingestion of uncabled profilers can take 24 hours to complete. Surface moorings may take 1.5 days, gliders 2 days, depending on number of deployments and amount of recovered data = 197 days for 113 uncabled platforms.
- Need configuration changes to allow multiple simultaneous ingestions to run more quickly. This will allow the 113 ingestions to occur in a compressed time frame by multiple team members.
- Need modification to the ingestion delay that eliminated the race condition issue, without which all telemetered data will be 24-48 hours delayed (not a blocker, but it does prevent real-time delivery and data status alerts)
- Purge functionality is still limited to purging an entire instrument (by reference designator) for all deployments and all delivery methods. Improving purge to allow purge by stream and time range will speed up the process.



# Review

### **Current Status:**

• Data review interrupted while the team focuses on ingest. Previously run reviews resulted in annotations that largely deal with unexplained data gaps, which we hope will be filled in by the purge and reingestion process.

### Dependencies

- Completion of ingestion of recovered data sets that the data team is currently performing. Without ingestion, the sources of data gaps will not be able to be verified, requiring redundant checks.
- Data review can take 2-3 days to run per instrument (possibly less, depending on data volume and optimization of routines), and results in a report that can be used to enter annotations.
- Need to collate operational & cruise information from MIOs in order to compare to data availability/quality assessments and annotate properly





# Annotation

#### **Current Status:**

 Annotations from previous data reviews logged within Sage Database (ooi.visualocean.net/) and uFrame (only the most urgent notifications were entered manually). The Sage Database is the more comprehensive, and ingest of this group of annotations into uFrame is the focus of this topic. Any additional annotations will require data quality review to proceed.

#### Dependencies

- Data review is required to create a full list of annotations for each data set, using automated routines, and visual inspection of plots.
- An API to push annotations at the platform, node, instrument, and stream level exists, but the manual entry process via the UI only currently allows stream-level annotations.
- Once review is complete for an instrument, annotations can be pushed to the uFrame system using the API.
- MIO operators can also enter annotations via the UI, if they are granted permissions and training, but there may be conflict between annotations entered via API and manual entry.







# ERDDAP

#### **Current Status:**

 There are 846 streams available on the production server at <u>https://erddap-uncabled.oceanobservatories.org/uncabled/erddap/info/index.html?</u> <u>page=1&itemsPerPage=1000</u>

#### **Dependencies:**

- Login, edit and write access to production system:
  - Systems team must push the software updates as we do not have access to production.
  - $_{\odot}\,$  Timely responses from systems team on requests for info through Redmine
  - Developer requires direct access to production to edit existing ERDDAP XML Files
- Some active tickets need resolution to remove caveats and inconsistencies in the data sets, including #12544 (Variable naming conventions for external L1/L2 parameters) and #10745 (NetCDF format: use of coordinate attribute and associated issues)
- As data streams get split and parameters merged (DOSTA/CTD, ADCP, etc.) ERDDAP developer needs to be informed, so that templates and datasets can be updated.





# Demos

- 1. Main OOI website: oceanobservatories.org
- 2. Data Portal: ooinet.oceanobservatories.org
  - a) Navigation
  - b) Data catalog
  - c) Quick Plotting options
  - d) Data download
- 3. OOI ERDDAP server
- 4. M2M API and real-time plotting







### Overview

- 1. Introduction to OOI
- 2. Data Flow & Products
- 3. Data Review Procedures
- 4. Periodic Reviews & Documentation
- 5. Next Steps
- 6. Conclusions







### **Reviews and Reporting**

#### **Quality Timeline**



#### **Annotation Text**

Level	Deployment	StartTime	EndTime	Annotation	Status	Redmine#
ctdbp_no_sample	D00001	2014-08-15T00:12:00Z	2014-08-25T18:50:41Z		NOT_AVAILABLE	
ctdbp_no_sample	D00001	2014-08-31T19:13:27Z	2014-09-22T22:42:44Z		NOT_AVAILABLE	
ctdbp_no_sample	D00001	2014-11-04T16:05:51Z	2014-11-05T18:56:20Z		NOT_AVAILABLE	
ctdbp_no_sample	D00001	2014-11-14T18:36:38Z	2014-11-17T18:03:22Z		NOT_AVAILABLE	
CE04OSBP		2014-12-07T19:45:00Z	2014-12-16T00:00:00Z	PFE down. HVPS1 MOV explosion, 800A breaker tripped, investigation and restoration		12264
ctdbp_no_sample	D00001	2014-12-07T20:59:40Z	2014-12-16T22:29:37Z		NOT_AVAILABLE	
ctdbp_no_sample	D00001	2014-12-16T23:44:08Z	2014-12-29T20:30:01Z		NOT_AVAILABLE	
CE04OSBP		2015-01-07T07:32:00Z	2015-01-07T08:06:00Z	PNWGP Portland <-> Seattle outage		12264
CE04OSBP		2015-01-31T00:00:00Z	2015-02-04T00:00:00Z	Intermittent partial data loss due to storage drive problems at OTB		12264
ctdbp_no_sample	D00001	2015-01-31T23:59:59Z	2015-02-03T09:56:15Z		NOT_AVAILABLE	
ctdbp_no_sample	D00001	2015-03-03T02:16:22Z	2015-03-06T19:57:13Z		NOT_AVAILABLE	
ctdbp_no_sample	D00001	2015-03-06T19:58:48Z	2015-08-02T00:00:00Z		NOT_AVAILABLE	
CE04OSBP		2015-03-21T14:10:00Z	2015-03-22T04:20:00Z	PNWGP outage due to City of Seattle fiber cable work		12264
CE04OSBP		2015-06-13T00:00:00Z	2015-06-15T16:30:00Z	Network issues due to fire that damaged fibers between Portland and Seattle		12264
practical_salinity	D00002	2015-08-12T01:00:00Z	2015-08-14T02:00:00Z	Unusual drop in conductivity values.	SUSPECT	
density	D00002	2015-08-12T01:00:00Z	2015-08-14T02:00:00Z	Unusual drop in conductivity values.	SUSPECT	
ctdbp_no_seawater_conductivity	D00002	2015-08-12T01:00:00Z	2015-08-14T02:00:00Z	Unusual drop in conductivity values.	SUSPECT	
dissolved_oxygen	D00002	2015-08-12T01:00:00Z	2015-08-14T02:00:00Z	Unusual drop in conductivity values.	SUSPECT	
conductivity	D00002	2015-08-12T01:00:00Z	2015-08-14T02:00:00Z	Unusual drop in conductivity values.	SUSPECT	
CE04OSBP		2015-08-29T00:00:00Z	2015-08-29T00:30:00Z	Outage during major utility power failure in Seattle		12264
CE04OSBP		2016-01-07T06:10:00Z	2016-01-07T06:52:00Z	Four 1-minute outages between Portland and Seattle due to maintenance		12264
CE04OSBP		2016-03-10T23:06:00Z	2016-03-11T09:30:00Z	Fiber break between Portland and Seattle		12264
CE04OSBP		2016-05-20T16:33:00Z	2016-05-20T18:04:00Z	Fiber break between Portland and Pacific City		12264
CE04OSBP		2016-07-12T02:53:00Z	2016-07-12T03:51:00Z	Unexplained loss of power at Pittock Building in Portland		12264
ctdbp_no_sample	D00002	2016-07-18T00:42:58Z	2016-07-19T21:06:56Z		NOT_AVAILABLE	
ctdbp_no_sample	D00003	2016-07-22T22:50:00Z	2016-07-25T19:51:39Z		NOT_AVAILABLE	
CE04OSBP		2016-12-17T18:00:00Z	2016-12-17T19:00:00Z	Corvalis data center lost power		12264
CE04OSBP		2016-12-22T01:50:00Z	2016-12-23T12:44:00Z	Fiber break in Portland due to train crash		12264
CE04OSBP		2017-01-08T19:58:00Z	2017-01-08T21:41:00Z	Network outage during major Seattle utility power failure		12264
ctdbp_no_sample	D00003	2017-01-09T18:30:53Z	2017-01-11T01:16:53Z		NOT_AVAILABLE	
CE04OSBP		2017-01-09T18:32:00Z	2017-01-09T21:30:00Z	Lightning strike in Pacific City led to data interruption through isolation of both cable lines from shore station equipment.	NOT_OPERATIONAL	11776
CE04OSBP		2017-02-07T13:00:00Z	2017-02-07T15:00:00Z	Outage during PNWGP 1-hour router-maintenance		12264
ctdbp_no_sample	D00003	2017-02-15T14:43:05Z	2017-02-16T22:27:12Z		NOT_AVAILABLE	
CE04OSBP		2017-02-15T14:43:06Z	2017-02-15T17:16:00Z	On Wednesday, February 15, power to the North and South cable	NOT_OPERATIONAL	11998





## **QC** Database Tool

- Used for reference & statistics
- Includes status information, as well as a cruise data checklist
- Includes testing/review capability
- Annotation options
- http://ooi.visualocean.net

	Data Tean	n				ŝ.	Arrays	Instrume	nts Cl	asses	Cruises	Reference	e <del>-</del> Sign
Arrays / / 0	E02SHBP-L	J01D-06-CTE	BPN106 / stream	ed / ctdbp_no_	_sample								
Data Strea	m Repo	ort									00I S	ite Page 🕑	Data Portal 🖸
Instrument Reference Desi M	a Name CT gnator CE Method str Stream cto	D 202SHBP-LJ0 eamed lbp_no_samp	1D-06-CTDBPN10	8	Ufra Instrur	me Route Driver Parser nent Type	seabir Scien	d.sbe16plus ce	_v2.ctdbp	_no.driver			
Data Availab from 9/10/2014 3	<b>ility Plot</b> :43:00 PM to 5	/11/2017 5:15::	30 PM										
(	October	2015 A	pril July	October	2016 A	pril J	uly	October	2017	April			
Deployments						_							
Cassandra													
Annotations 30 Annotation: Metadata	Param S Start Date	eters End Date	Comment										
CE02SHBP-	9/10/14,	9/25/14,	Exact cause of da	ita gap at begir	nning of deployr	nent curren	tly unkn	own. Most lik	ely due to	testing as	s the instru	ment was co	ming online.
CTDBPN106	3.43 PW	0.17 PW	Todo: check: diffe By friedrich, on 3/29/	erence betweer	n data begin an	d deployme	nt begin	date is: 15 d	lays 02:34	:58			
Status: NOT_AVAILABLE Deployment: 1 Method: streamed Stream: ctdbp_no_sample				17									







### Deliverables

- Data Availability Reports
  - $_{\odot}$  (% completeness, streams/parameters reported, particles in the system)
- Data Quality Reports
- Redmine reporting

   Issues found, investigations, and Help Desk open/closed
- Deep dive investigation reports
- Annotations (to users)
- Download statistics
- Forum statistics (TBD)







### **Data Evaluation Daily Activities**



- Review the end-to-end operational status of online instruments and investigate any outages (e.g. instrument, telemetry, parsing, or ingestion failures).
- Review the operational status of other data archives (raw, cruise, ERDDAP)
- Look into and resolve new system alerts
- Follow up on any issue requests from users (via Redmine)
- Add annotations to notify users of operational status changes





### **Daily Review Workflow**







### **Periodic Data Team Activities**

- Meet with MIOs to discuss operational issues and data quality
- Instrument, stream, parameter and deployment completeness
- Conduct deep dives on datasets to review availability and quality
- Review & annotate full deployment data to assess data quality
- Develop new scripts, plotting tools, and quality checks
- Produce reports on the availability and quality of datasets
- Review appropriateness of QC flags
- Ensure asset, deployment, calibration, and ingestion configurations have been updated, and reports posted following every cruise
- Prototype and test new user interface and visualization features







### Data assurance/Data quality: *Pre and Post comparisons*







### Overview

- 1. Introduction to OOI
- 2. Data Flow & Products
- 3. Data Review Procedures
- 4. Periodic Reviews & Documentation
- 5. Next Steps
- 6. Conclusions







### Sept. 2016 Workshop Feedback: Short Term SW fixes

- i. No file aggregation prior to delivery. Fixed
- ii. Improved bathymetry. Fixed
- iii. Data Team **annotations** from MIO information about HW issues **Fixed:** Data team now able to enter annotations, view them in GUI (more work to go)
- iv. Large data download time outs, request lost, email response confusing. Fixed
- v. Depth, Lat/long and pressure for all deployed instruments. Fixed
- **vi. Status timeline** with metadata; overview from first deployment to present. Operator can perform manual updates. **Fixed**
- vii. Incorporate all the naming and labeling options (**vocabulary**) that the data team added to preload, and improve filtering in the GUI. Improved, but ongoing
- viii. Missing data & instruments in catalog (eg. MASSP, ZPLSC). Improved but ongoing: routes users to raw data archive, etc.; some analytical data still incoming





### Workshop Feedback: Long-Term SW fixes

- ix. Plotting from **multiple instruments** on 1 plot. Fix: Re-enabled, but needs GUI fix
- **x. Simplify Data Catalog** to start with empty "cart" so users are not overwhelmed with options. Future fix.
- **xi. Plotting clarifications**: users should be informed if data cannot be plotted in 2D and best way to plot. Future fix.
- **xii.Improved links and access** to raw data archive, documentation, and metadata info Fix in progress
- **xiii.Improve Overall Data Quality**. Long-term fixes: data team deep dives, secondary post-recovery calibration, external review by SMEs, Data Assembly Center (would require reallocation of funds)





### Data Availability and Completeness

- Some derived data products are still being added to the system or require additional processing
- HPIES, ZPLSC/G, ADCP, MASSP, BOTPT, HYDBB, VADCP, FDCHP, CTDBP/FLORT, CAMDS, CAMHD, OSMOI, PPSDN, RASFL

Seafloor Uplift and Deflation





1.0

-1.0

5-min Rate of Depth Change (BOTSFLU L2)



### **QC** Challenges & Solutions

- Local range values need statistical analysis of environmental data for each platform
  - Need to work with SOC to analyze and apply ranges and test algorithm
- Trend test may not work as designed, because it requires the system to compare data prior to the user request date – *analysis ongoing*
- Gradient test is complicated to apply, requires 2D dataset *analysis ongoing*
- Spike test is currently very simple needs tweaking to avoid false positives/negatives (especially in biological data) and to work with certain data types
- Not all QC algorithms apply to all data products *ongoing review with SOC*
- The QC algorithms do NOT trigger alerts in the system *Alerts/alarms only trigger when new data is telemetered/streamed* 
  - Can set alerts on L1/L2 data streams based on Global/Local range values





### **Options for Data Review Acceleration**

Option	Positive	Negative
MIO Operations Log at Rutgers	<ul> <li>Centralized log reduces time spent investigating issues</li> <li>All issues entered consistently</li> </ul>	<ul><li>Takes time to maintain</li><li>Some development time</li></ul>
Speed up ingestion	Fewer gaps to investigate	<ul><li>Currently requires FTE</li><li>Automated process not yet delivered</li></ul>
Data Team works only on RIC	More data reviewed faster	<ul> <li>No new data in system</li> <li>No bug investigation</li> <li>No QA testing</li> </ul>
Limit reviewed time period or stream type	<ul> <li>Data reviewed slightly faster, at high level</li> </ul>	<ul> <li>Review enhanced by looking at multiple deployments and trends</li> <li>Slows down future reviews</li> </ul>
Limit thoroughness of reviews	<ul> <li>Data reviewed faster, at high level</li> </ul>	<ul> <li>Unclear why gaps exist</li> <li>Quality issues not fully annotated</li> <li>Slows down future reviews</li> <li>Limits crowdsourcing options</li> </ul>
Crowdsourcing (enlist volunteer SMEs)	<ul> <li>Removes subset of datasets from review queue</li> <li>Assistance with complex data that requires expertise</li> </ul>	<ul> <li>Focus on specific interest, not whole of OOI</li> <li>Steep learning curve for advanced use of system (and knowledge of known issues)</li> <li>Pathway to triage and incorporate feedback</li> </ul>
Add employees or Data Assembly Center (DAC)	<ul><li>Data reviewed faster, in depth</li><li>Support for expert analysis</li></ul>	<ul><li>Requires additional funding</li><li>Setup and maintenance time</li></ul>



### Adding capability to OOINet experience

#### Ocean Observatories » Web Interface

evention security required inter bout outer and and becamento which they be	Overview	Activity	Roadmap	Issues	New issue	Gantt	Calendar	News	Documents	Wiki	Files	Set
---	----------	----------	---------	--------	-----------	-------	----------	------	-----------	------	-------	-----

#### Enhancement #10913 Color bar needs to scale to majority of data for binned pseudocolor plots Added by Daniel Maher 4 months ago. Updated 16 days ago. Status: Start date: Ready for Work **Priority:** Due date: Urgent Assignee: Daniel Maher % Done: 20% **Category:** Spent time: Target version: **Target Release: Array Affected:** Category 1: **Instrument Affected:** Severity: **CI Software Affected: Issue Closed:** Work Breakdown Structure (WBS): Northward UI Plot



016-10-05

016-10-06

Python Plots

ings

RS01SBPS-PC01A Stream: adcp\_velocity\_beam 2016.09.30T16.00.00 - 2016.10.07T16.00.00 Northward Sea Water Velocity, m s-1









#### Science Evaluation: Are ocean features encountered real? Outside local range

Taking advantage of all assets, even non-NSF, to assess data quality





76°W

74°W

72°W

70°W

**OOIFB Meeting Fall 2017** 

42

68°W



### Vicarious Calibration - Comparisons enabled by ERDDAP





Co-located and concurrent Temperature and Salinity Profiles Blue – Glider 387 Red – Coastal Pioneer Profiler Mooring





# **Outreach and Community efforts**

- Community tools section of website
- Hackathon to develop new visualization and processing tools
- Education meetings at Rutgers to get data into the classroom
- Data team presence at meetings and workshops (MTS, Ocean Sciences, Irminger)
- Engagement with SMEs







### Overview

- 1. Introduction to OOI
- 2. Data Flow & Products
- 3. Data Review Procedures
- 4. Periodic Reviews & Documentation
- 5. Next Steps
- 6. Conclusions







# Conclusions

- 1. A large amount of high quality data has been and is being collected, with high science value
- 2. Data review is finally our primary focus, given maturation of the system
- 3. Data team accelerating RIC review via development of specialized tools
- 4. Short-term, medium-term, and long-term goals for improving data quality and delivery
- 5. OOI is providing a curated, consistent data system that is delivering data and metadata to the community





### **Questions?**

- OOI Main Web site: <a href="http://oceanobservatories.org">http://oceanobservatories.org</a>
- Data Portal: <a href="http://ooinet.oceanobservatories.org">http://ooinet.oceanobservatories.org</a>

Mike Vardaro, Data Manager, OOI CI Data Team vardaro@marine.rutgers.edu

Mike Crowley, Program Manager, OOI CI Data Team Crowley@marine.rutgers.edu

**Acknowledgements**: NSF, COL, Rutgers University, University of Washington, WHOI, Oregon State University, RPS-ASA, Raytheon, UCSD/SIO





