

A large school of fish, possibly jackfish, swimming in the ocean. A diver is visible in the lower left, holding a camera, suggesting a documentary or scientific context. The water is clear and blue.

Re-envisioning data repositories for the 21st century

Diggs, S. and Kinkade D.














EGU General Assembly 2020
2020-05-07

EGU2020-20826
<https://doi.org/10.5194/egusphere-egu2020-20826>

Credit: Octavio Aburto



FIFTY YEARS OF OCEAN PROFILE DATA

OCEAN SCI.							
							
	- GEOSECS	- SAVE (Atlantic) - CTD on Rosette - ADCP - TAO Array	- WOCE - LADCP - PALACE Floats	- CLIVAR - Argo Array - Pinniped CTD	- GO-SHIP	- Deep Argo - BIO-Argo	
TECH / ORG	1960	1970	1980	1990	2000	2010	2020
	- NODC est. - Data by Mail - Inquiries by Phone - Acoustic Modem	- NOAA est. - DARPA net - Service Argos - Minicomputers	- Personal Computers - Internet/Email - FTP - 9.6k Modem	- Civilian GPS - Email on Ships - WWW	- Internet on Ships - Iridium Satellites - Data Web Sites	- Personal Mobile Dev. - Data Web Sites - Cloud Storage - Social Media	
							
	1960	1970	1980	1990	2000	2010	2020
#of ocean profiles		3M	4.6M	6.4M	8.2M	12.14M	13.5M
disk cost/mb		\$260	\$192	\$5.28	\$0.01	---	
memory cost/mb		\$7.4K	\$6.4K	\$98	\$1	\$0.01	

Evolution Of A Data Center



1990



2000



2010



2020

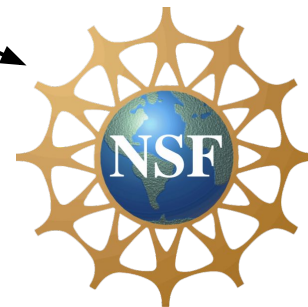
BCO-DMO

Biological & Chemical Oceanography Data Management Office

locatedAt



primarilyFundedBy



providesDataTo



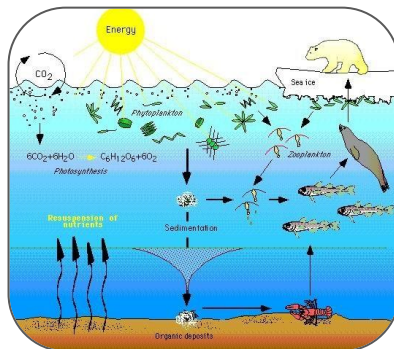
public

submitsDataTo



scientists

conductsResearch



ocean ecosystem dynamics

BCO-DMO Data Types

Science Perspective

- Biological
- Chemical
- Biogeochemical
- Physical
- Geophysical

- *In situ*
- Laboratory
- Remotely sensed
- Synthetic/derived

- Molecular to Megafaunal
- Local to Global
- Discrete to continuous /synoptic

fromDiscipline

usedMeasurementTechnique

Highly Heterogeneous!

digitalDocument

storesValuesFor

DATASETS = >9000

Curation Perspective

- ASCII Text (tabular)
- Binary (e.g., NetCDF)
- Images
- Acoustics
- Application (e.g., Matlab)
- Links to other data

- Variable organization
- Varying metadata
- Local parameter terms
- Emerging data types
- Distributed complementary info

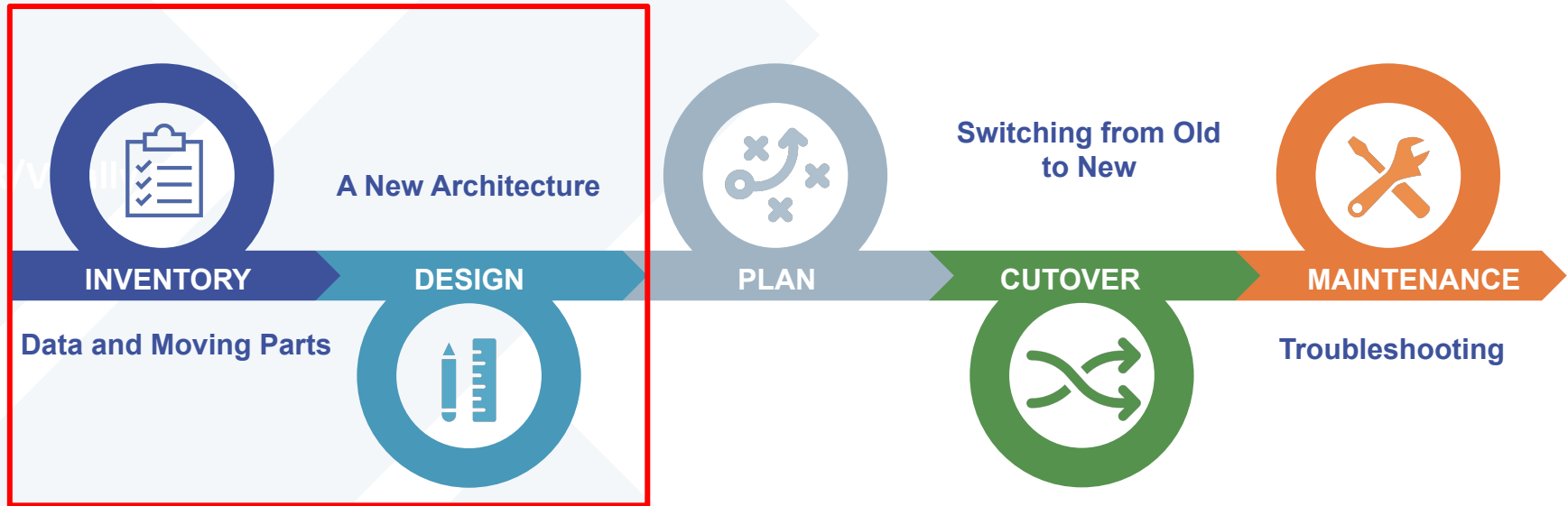


Many repository improvements are driven by the state of technology.

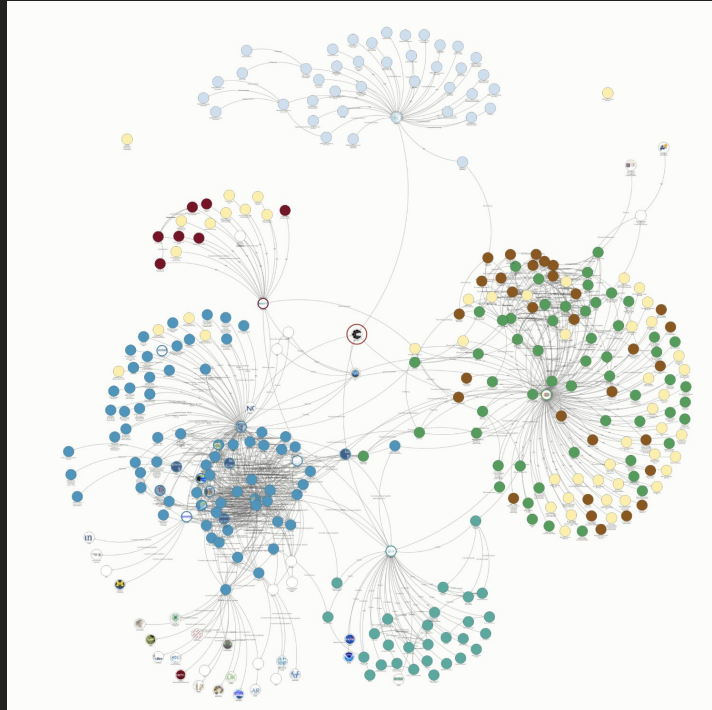
This is reactive.

How do we become proactive?

GENERALIZED REPOSITORY TRANSITION MODEL (IOOS)



Intentional evaluation of current technologies and best practices



BCO-DMO

Biological & Chemical Oceanography Data Management Office



Submission Software
Curation Software
Access Software
Custom APIs

currentlyDefiningModel

definesMetadataModel



DigitalObject

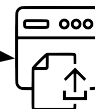


Provenance

DigitalObjectAccess



DigitalObject
Submission



reducesAssumptionsInCode

describeWhatWeDoAndKnow

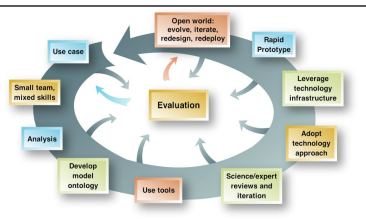
informedBy



[Best Practices for Publishing Linked Data](#)



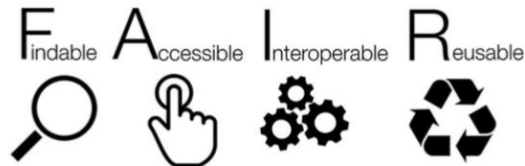
[Data on the Web Best Practices](#)



Schema.org Value

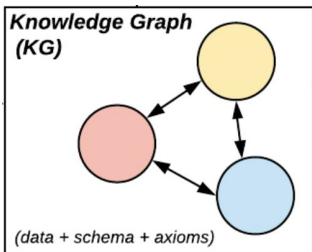
Why add additional information to websites?

- Improve *findability*
 - Google and other indexers
- Facilitate data *access*
 - Programmatic support for resource access
- Enhance *interoperability*
 - Well described resources
- Promote *reuse*
 - Simplified programmatic access



Result: new data model reflects the current state of original and adjacent domains.

lod.bco-dmo.org



constructedWith



usesModel

followsPattern

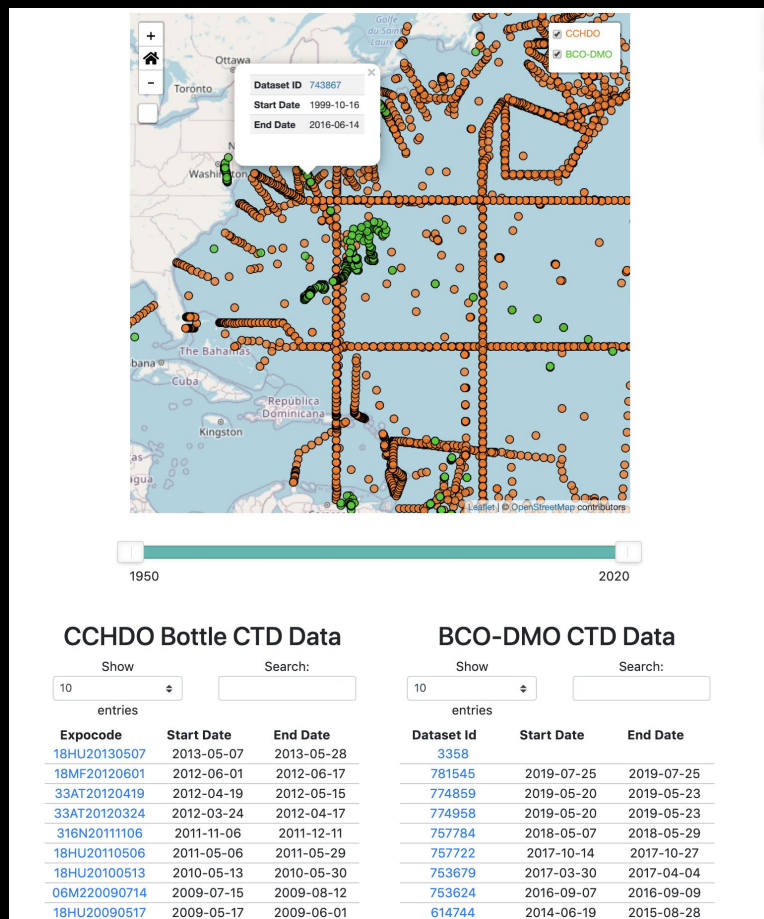
supportsQueriesFor

- Dublin Core
- Dataset Catalog (DCAT)
- PROV-O
- Schema.org
- Bibliographic Ontology
- GeoSPARQL
- NERC Parameter Usage Vocabulary
 - Environment Ontology
 - CHEBI
 - *and more...*
- *NanoPublications*
- *Ocean Data Ontology (BCO-DMO)*

- Dataset & Data Type
- Funding Source & Award
- Geospatial & Temporal Coverages
- Instrumentation
- Observed Properties
- People & Roles
- Persistent Identifiers (DOIs, ORCID, etc.)
- Related Publications (by PID too)
- Research Topic
- Supplemental Documents

- Science-on-schema.org (ESIP)
 - Dataset
 - Repository Services
- Data Type Registry (RDA)
 - Ocean Proteomics
 - *more coming*
- *PID Kernel Information Profile (RDA)*

Coordination of technologies and practices lay a solid foundation for desirable outcomes





The Journey
Continues...

@scdiggs
@danie_b_k