



# CGIAR Data Management Task Force (DMTF) Annual Meeting Report

June 17, 2016; Montpellier, France

Author(s): CGIAR System Office & DMTF members



## CGIAR Data Management Task Force (DMTF) June 17, 2016 Annual Meeting Report

### Contents

Agenda .....	3
Session 1: Brainstorming .....	4
Session 2: Center/CRP updates .....	9
Session 3: Amazon Web Services.....	11
Session 4: CRP mapping project .....	11
Session 5: Next Steps & Meeting Closing.....	12

## Agenda

Friday, June 17, 2016	
9:00 - 9:30	<b>Introductions, overview, goals</b>
9:45 - 11:00	<b>Brainstorming</b> <ul style="list-style-type: none"> <li>In a perfect world, where do you see CGIAR and your Center in the next 3-5 years in terms of open data, and being able to leverage big data capabilities?</li> <li>What are your top 3-5 challenges in achieving this?</li> <li>What are some solutions to overcome those challenges – actions from Center, CO, others?</li> </ul>
11:00 - 11:15	Tea/coffee
11:15 - 12:30	<b>Center discussion: OA/OD &amp; policy compliance; data capture &amp; management tools (10 min max) – not formal, more learning/sharing opps; no need for polished slides</b> <ul style="list-style-type: none"> <li>Current status towards full compliance in 2018 (including CG core implementation/mapping etc.)</li> <li>What's working (including cross-Center, cross-CRP initiatives)</li> <li>What's not working so well/areas you could use help</li> <li><u>QUICK demos</u> of data capture/curation/quality tools and approaches</li> </ul>
12:30 - 13:30	Lunch
13:30 - 15:00	<b>Center discussion: OA/OD &amp; policy compliance; data capture &amp; management tools (10 min max) – not formal, more learning/sharing opps; no need for polished slides</b> <ul style="list-style-type: none"> <li>Current status towards full compliance in 2018 (including CG core implementation/mapping etc.)</li> <li>What's working (including cross-Center, cross-CRP initiatives)</li> <li>What's not working so well/areas you could use help</li> <li><u>QUICK demos</u> of data capture/curation/quality tools and approaches</li> </ul>
15:00 - 15:30	Tea/coffee
15:30 - 16:00	Amazon Web Services (Eddie Romanzo; Jawoo Koo)
16:00 - 16:30	CRP mapping project (Jawoo Koo)
16:30 - 17:15	Legacy data and data prioritization discussion (Medha Devare)
17:30 - 18:00	Next steps; meeting closing

## Participants

	Institution	Name
1	AfricaRice	Ibnou Dieng
2	Biodiversity International	Elizabeth Arnaud
3	CCAFS	David Abreu
4	CIAT	Leroy Mwanzia
5	CIFOR	Sufiet Erlita
6	CIFOR	Usman Muchlish
7	CIMMYT	Kate Dreher
8	CIP	Henry Juarez
9	ICRAF/DS CRP	Paul Baraka
10	ICRISAT	Srinivasarao Chukka
11	IITA	Martin Mueller
12	ILRI	Jane Poole
13	IWMI	Salman Siddiqui
14	ICARDA/DS CRP	Enrico Bonaiuti
15	CGIAR System Office	Medha Devare
16	CGIAR System Office	Michelle Fotsy

## Session 1: Brainstorming

Group work focusing on key questions in three areas: 1) aspirations, 2) challenges, and 3) solutions. Post-its with individual points (bullets in tables below) were clustered by the group into broad themes, which form the headings of columns in each table.

### 1) Aspirations

Where do you see CGIAR and your Center in the next 3-5 years in terms of open data, and also in terms of being able to leverage big data capabilities? What can we realistically accomplish in 3 years? The caveat is that the OADM policy goes into full effect in 2018; the intent of the policy is to try and move us in the same direction, recognizing that not everybody will be in the same place at the same time.

	Usability /access	Interoperability	Discoverability	Culture/Incentives
	<ul style="list-style-type: none"> <li>• Great quality of reporting from project leaders</li> <li>• Shortened data lifecycle (curation/cleaning/storage, etc.)</li> <li>• High quality public data, tagged with appropriate caveats or qualifiers</li> <li>• Big data: providing broader analytical capabilities</li> <li>• Ontologies that can be easily updated by the community (with moderation) and fed back into systems that use them</li> <li>• GPS tagged data collected in fields around the world (by the farmers, collaborators, etc.) -- which can be fed back into CG data systems and used to make</li> </ul>	<ul style="list-style-type: none"> <li>• Platform launched and used by centers and CRPs because consensus is reached to use similar standards and tools/toolboxes</li> <li>• Public datasets are tagged with ontology terms based on manual or automated annotation</li> <li>• Seamless interaction between the three platforms</li> <li>• Compliance with linked open data</li> <li>• Climate data collected at micro &amp; macro level that can be easily passed into breeding decision systems, modeling tools, etc.</li> </ul>	<ul style="list-style-type: none"> <li>• Central workplace</li> <li>• Effective data harvesting system between center repositories</li> <li>• One stop-shop access to data across CRPs/Centers</li> <li>• All publicly available CG datasets can be found through a CG &amp; other aggregating search portals based on metadata</li> <li>• A platform capable of harvesting OD from all centers (metadata only)</li> <li>• Data are captured in internal systems &amp; pushed automatically into public warehouses</li> </ul>	<ul style="list-style-type: none"> <li>• Having OA/OD culture first implemented in the teams → projects → programs → center → CGIAR → CG stakeholders</li> <li>• Centers giving clear guidance and applying robust ethical rules for sharing confidential data + clear link to IP &amp; research ethics committees</li> <li>• Scientists being recognized for publishing data &amp; tools as much as they are for other information products (like papers)</li> <li>• Data management &amp; analyses will be cost effective</li> </ul>

	Usability /access	Interoperability	Discoverability	Culture/Incentives
	<p>key breeding, targeting, forecasting decisions, etc.</p> <ul style="list-style-type: none"> <li>• Data collection using mobile technology becoming the norm (including direct beneficiary input)</li> <li>• Data stored cheaply and securely and connected to cheap high power computing facilities</li> <li>• Data can be accessed and understood by people in <u>all</u> languages</li> </ul>	<ul style="list-style-type: none"> <li>• Availability of data as web services or linked open data</li> <li>• All platforms interoperable at data level</li> <li>• Data is interoperable, harvestable (machine readable) according to I.O. standards</li> <li>• A strong community for documenting different types of datasets (controlled vocabularies, ontologies)</li> <li>• Data can pass easily between a suite of common systems &amp; tools used by the CG and pushed into a wide array of analytical pipelines</li> <li>• Big data leveraged to provide data integration possibilities</li> <li>• Substantial progress in semantics: Agrovoc; ontologies; subject terms; automatic tagging; profiling people and objects; linking (research) data models with data distribution /management</li> <li>• All (most with exceptions) data is open within policy timeframe, documented with metadata at all levels</li> </ul>	<ul style="list-style-type: none"> <li>• Sharing services (Amazon, DOIs, etc.) and harmonization (minimum standards)</li> <li>• Data playground to leverage big data</li> <li>• Ability to ask questions to the big data platform and find the right datasets to answer those</li> </ul>	

## 2) Challenges

What are your top 3-5 challenges in achieving this? What are some stumbling blocks?

	Infrastructure, standards, Data Quality & Prioritization	Cost/Resources	Culture and Capacity	Impact
	<ul style="list-style-type: none"> <li>• Prioritizing legacy data [for “open”]</li> <li>• Improving quality of reporting</li> <li>• Common standards -- as tools, people, donors change, and ‘international’ standards also change (e.g. country naming!)</li> <li>• Lack of ontologies in common</li> <li>• Data generated within center and by partners are not standardized and are of variable quality</li> <li>• Even within a center, lack of communication leads to some data siloing &amp; lack of interoperability</li> <li>• Different priorities challenge collaboration across centers (sometimes)</li> <li>• Some key primary databases are not ready – and the whole system of sharing data with the public will draw on these</li> </ul>	<ul style="list-style-type: none"> <li>• Convincing donors to support OD/OA</li> <li>• Cost of making data open: legacy data; resources for curation</li> <li>• Resource mobilization for sharing services</li> <li>• Time/people vs money</li> <li>• Not adequate staff time: training scientists; testing systems; migrating data; curating data; writing policies; scoping new tools, etc.</li> <li>• Resource cost sharing by CG for common activity; e.g. data manager salary can be shared by center (70%) as well as CG (30%) [i.e., need for sharing some key services]</li> </ul>	<ul style="list-style-type: none"> <li>• Lack of expertise</li> <li>• Big data/interoperability</li> <li>• Capacity and time from our researchers devoted to data management</li> <li>• Frequently changing/ advancing tech domain</li> <li>• Capacity building &amp; sustainability</li> <li>• Practices during the project (scientist behavior)</li> <li>• Motivation of researchers to share data</li> </ul>	<ul style="list-style-type: none"> <li>• Tracking OD usage &amp; impact →not just downloads (e.g. citations, new tool using data)</li> <li>• Tracking impact of the data &amp; the platform</li> <li>• Measuring the impact of sharing data</li> <li>•</li> </ul>

- Core data management systems in centers must have the data push feature and be compliant to CG metadata (e.g; phenotype data from BMS)
- Big Data computing infrastructure required
- Inflexibility of hosted open data platforms, platforms; e.g; hosted DataVerse does not support CGCore. No SLA agreements with hosted OD providers e.g. Harvard.

### 3) Solutions

What are some solutions to overcome those challenges – including possible actions from Center, SO, others?

	Infrastructure, standards, Data Quality & Prioritization	Cost/Resources	Culture and Capacity	Shared Services
	<ul style="list-style-type: none"> <li>• Publish only metadata of the legacy, non-curated, data → discoverable [then assess which data sets to make open based on usage/hits]</li> <li>• Legacy data:               <ul style="list-style-type: none"> <li>○ Priority setting protocols, incl. how far we “go back” -- should be led by key priority research question</li> <li>○ Have metadata on the status of the data</li> <li>○ Accessible upon request</li> <li>○ Crowd-source cleaning of metadata for legacy data</li> </ul> </li> <li>• Enabling feedback by data user (fitness &amp; quality)</li> </ul>	<ul style="list-style-type: none"> <li>• Centers write proposals for data curation and infrastructure</li> <li>• SO to help mobilize resources for data curation and infrastructure</li> <li>• More time + resources to:               <ul style="list-style-type: none"> <li>○ Change culture (budgeting for DM, cognition of data as product)</li> <li>○ Provide incentives</li> <li>○ Collaborate</li> <li>○ Provide support (tech + non-tech)</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>• Quality of reporting as part of project evaluation</li> <li>• SO to take lead in building capacity of centers to try to bring them at same level</li> <li>• SO to take lead on capacity building to target technical advances</li> <li>• Capacity building workshops on OA/OD, data management good practices</li> <li>• Training for data managers on platforms and tools</li> <li>• Motivation: adopting incentives based data sharing policy, e.g. performance evaluation</li> </ul>	<ul style="list-style-type: none"> <li>• Enabling centers have control of OD platforms, e.g.:               <ul style="list-style-type: none"> <li>○ CGIAR support financially center hosted platforms, or</li> <li>○ CG hosted solutions: CGKAN &amp; CGverse</li> </ul> </li> <li>• CG-wide cloud computing facility supporting processing and storage reducing transaction time &amp; cost</li> <li>• CG-wide (shared services) access to data tracking / usage; e.g. WOS (data citation index) to track the impact</li> <li>• Pool/shared consultancy available with different expertise areas</li> </ul>



## Session 2: Center/CRP updates

- CIMMYT: Hoping to comply with 2018 data release deadline; working on data prioritization to figure out which are the highest priority ones to make available; OA/OD impl plan being worked on.
- CCAFS: Resource-based management system that tracks all outputs from planning to reporting and beyond to monitor OA compliance; this generates a set of qualifiers; though there may be some flaws in reporting accuracy and quality, it is more like a subsite to interact with the community to get those numbers more accurate and go forward; OA/OD implementation plan has been revised a couple of times, and is now being done together with CIAT; comments being finalized before submitting to SO.
- ICARDA: Impl. plan submitted and revised based on SO comments; plan passed on to Executive Committee and needs to be endorsed by Board of Trustees. Being the last center/CRP to have a metadata schema, we adopted CG core. Pulled out of CGspace but agreed that if CG core is in place and a governance structure for CGspace and Dataverse, ICARDA would like to participate in some decisions and again share some of the budget. Having issues due to different installations of DataVerse in programs; planning to have only one installation for all. The recovery of lost data gave back 10000 datasets to register and upload; progress at approximately 100 items/month, and keep an inventory; given the lack of resources, it will likely take a couple of years.
- IITA: Efforts have been focused on infrastructure; in July CKAN will be in place. IITA is restructuring and the outcome is an opportunity to give more power to OA/OD steering; e.g. data processing, data managers, statisticians are put together and closer to ICT. IITA uses SharePoint as a collaboration platform, with project sites, repositories and will be putting workflows in place during the second half of the year. OA/OD impl plan was submitted in August 2015 and will be reviewed and updated.
- ILRI: Impl. plan was in final draft by August last year but still needs to be submitted. CKAN is up-to-date with CGCore; CGspace has had some technical issues. ILRI set up an institute committee (headed by Peter Ballantyne) established around compliance and monitoring of OA. Data portal has about 60 datasets; project datasets are in different places, and ILRI has plans to manage those.
- CIP: Impl. plan was submitted in September 2015; CIP is OA-compliant; has a data management plan working. CIP has a dedicated [OA site](#) on [cgiar.org](#) with CIP's OA & DM Policy, guidelines and procedures, as well as training modules; support is good at the DG level (e.g., yesterday their DDG research sent a memo to researchers to make open datasets a requirement). Dataverse newly implemented – not many records yet, but expect to have made much progress by 2017 end.
- IFPRI: Impl plan still in development; hoping to have ready by summer's end. IFPRI Dataverse has about 400 datasets; also have sub-Dataverses for CRPs. IFPRI is developing 3 ontologies: agriculture and nutrition, technology and value chain.

- CIFOR: Impl plan submitted June 2015. Have a Dataverse data system called 'PMC', and a Dspace publication system called 'My Cifor'; have been implementing CG Core in DSpace.
- ICRAF: Impl plan still under discussion and being written; currently with management. Have been focusing on defining the workflows. Project start up and project closure meetings are in place to ensure that all the project data is captured. HR clearance is also being implemented, with links to Dataverse for all datasets as a requirement. An ICRAF performance dashboard is also being developed as a data reporting management tool which captures all the different tools at ICRAF. Library catalog has changed over 3 years; now implementing KOHA system (same as IRRI); as Dataverse is using DOIs linking both systems is possible. Landscape portal stores spatial data. CG core has not yet been implemented; crosswalk done, and will see what other centers have done.
- Bioversity: Impl plan submitted to the SO. 2015 external audit conducted on research management and publishing processes, looking at workflows, bottlenecks, how scientists enter data, etc.; this enabled the creation of a resource data management plan; first version of guidelines being discussed with scientists. Bioversity going through portfolio restructuring, and internal reorganization with changes that will affect these processes. Champion users from different units and projects identified; datasets underlying published papers being uploaded to Dataverse and Dspace; champions' feedback on Dataverse is being compiled. Bioversity more or less in compliance with CG Core. Using the ESRI tool but lack a solid strategy for geospatial data. Re: ontologies, Bioversity works with breeders (mainly CRP breeding clusters) on the Crop Ontology and is working on an agronomy ontology; both ontologies have similar concepts, so can link them; have started collaboration with SDG ontology.
- A group discussion was engaged re: risks of Harvard Dataverse; in case of crash, even though there are back-ups, there is a risk of not having data readily available; additionally, there is no legal agreement between a Center and Harvard Dataverse (as opposed to service-level agreements with other providers like Amazon).
- IWMI: Impl plan done and submitted; now under review by management. Modified IWMI OA/OD policy based on CGIAR OADM Policy; have data management workflow in place. CG Core implemented; besides adopting CG Core, IWMI also follows ISO standards (e.g., [Water Data Portal](#) includes ISO and CG Core standards). A few OA seminars were held and plan to do more of those.
- ICRISAT: Currently making a shorter version of OA data management policy. Landscape analysis on current and legacy data is being done; scientists, project and team leads have been asked to give their top five datasets, as well as datasets linked to past research publications; these will be uploaded to data portal. Currently trying to convert metadata to CG Core standard. ICRISAT has not adopted a data implementation plan within projects yet, but this is in discussion with the leadership team. ICRISAT has had recent structural changes, and the data management unit is now with biometrics and bioinformatics enabling checking of products before making them public. Discussions with scientists have been initiated to try to understand more about datasets associated with published data and interlink the OA and OD repositories.

- CIAT: A common impl plan for CIAT and CCAFS submitted to SO; a final draft will be submitted to their team. Harvard Dataverse is data platform. Still working on CGspace and waiting to see what ILRI will be doing to conform to CG Core. There is now one team for data information and knowledge; library and data managers are working together to support OA/OD; also creating a CoP for data managers within projects to support their work; core team wants to provide support through a data quality management system; reaching out to institutions like Reading to see how this can be done; CIAT is open to suggestions and any other center interested in this initiative is welcome to join. Datasets related to all the publications are being requested. Working with the project management team to follow up on the project status; for flexibility reasons versions of the management plan is adapted to different project stakeholders; final product is uploaded to Dataverse as open. A new Strategic IT Committee is being created which will also look at OA implementation.
- AfricaRice: Impl plan was submitted to SO. Because many people are involved in managing information products (KM, Comms, DM, legal, editorial and publication committee), need to discuss among them as well as with DDG on activities related to OA/OD; periodically a memo is sent as reminder that all information products generated by AfricaRice staff is property of AfricaRice. Data management plan has not yet been established. A wiki is in place, with an external and an internal part for project management. Publication repository was developed internally, and discussion is ongoing whether to move to DSpace. For data, moving from Dataverse to CKAN; close to 500 datasets already in CKAN; have both AfricaRice metadata and the CG Core in place.

### Session 3: Amazon Web Services

[Amazon 17June2016 CGIAR Overview.pdf](#) Presentation by Eddie Romanzo, Nonprofit Account Manager & Dustin Sell, Senior Solution Architect

### Session 4: CRP mapping project

Discussion led by Jawoo Koo, around a global system like the WorldBank and USAID have, very transparently showing where people work, on what project, and what they generate. For CGIAR, there is no such global system yet; there were/are several efforts, e.g. ILRI had the regional mapping, CIMMYT has Maize Atlas & Wheat Atlas, the Gender Network has the gender map, etc. It would be good to have a better idea of those efforts and be able to search across them from a central place. There is the IATI (International Aid Transparency Initiative) with CGIAR as official member; some time ago, IFPRI was asked to develop a system to map all CRP projects. Now there is funding for this through PIM. The objective would be to have a dynamic map with indicators, of all CRPs' activities.

For centers to support the project, the needs, the audience must be defined. Jawoo Koo will draft something and share with the group.

## Session 5: Next Steps & Meeting Closing

### 1) Proposed institutional capacity building (face to face, online tutorials, webinars etc)

- Tools / specific platforms – e.g. DV, CKAN, etc. (e.g. Open Knowledge Foundation)
- Data management (best practices etc.)
- Technologies:
  - LOD & ontologies
  - New tech / agility – Big Data
  - Interoperability protocols at data level
- Development of a consultancy pool to help with these

### 2) Priorities / to-do

	Task	Who takes lead
1	Harvesting across repositories	Medha +CIAT, IFPRI, ICARDA
2	DV consulting – talk to CoP	Kate + Leroy + Medha
3	CGverse exploration	Kate + Leroy + Medha
4	Data-level interoperability; CG core mapping	David + Medha (with consultant)
5	DMP tool for CGIAR	Leroy + Jane (+ Medha in loop)
6	Impact of open data <ul style="list-style-type: none"> <li>○ Data citation index</li> <li>○ (Altmetrics)?</li> </ul>	Chukka (+ Medha in loop)
7	Ontology working group	Elizabeth + Kate + Martin + Enrico + Chukka + Henry + Leroy + Sufiet + Medha
8	DMTF meetings: <ul style="list-style-type: none"> <li>○ Virtual: 1 x 3 months (Sept, Dec, Mar)</li> <li>○ Face-to-face: 1 x 12 months –@ 3 days coupled with capacity building/professional development event (target: <u>RDA 9<sup>th</sup> plenary meeting</u>, April 5-7, 2017 in Barcelona)</li> </ul>	Michelle + Medha
9	Webinars	Rodrigo, Medha