# Towards Consistent and Citable Data Quality Descriptive Information for End-Users

Ge Peng[1,2], Nancy Ritchey[2], Anna Milan[2], Sonny Zinn[3], Kenneth S. Casey[2], David Neufeld[4], Paul Lemieux[3], Raisa Ionin[3], Robert Partee[3], Donald Collins[2], Jason Shapiro[3], Aaron Rosenberg[4], Thomas Jaensch[3], and Philip Jones[3]

[1] Cooperative Institute for Climate and Satellites – North Carolina (CICS–NC), NC State University, Asheville, NC, USA. Ge.Peng@noaa.gov; [2] NOAA's National Centers for Environmental Information (NCEI); [3] ERT, Inc./NCEI; [4] Cooperative Institute for Research in Environmental Sciences (CIRES)/NCEI
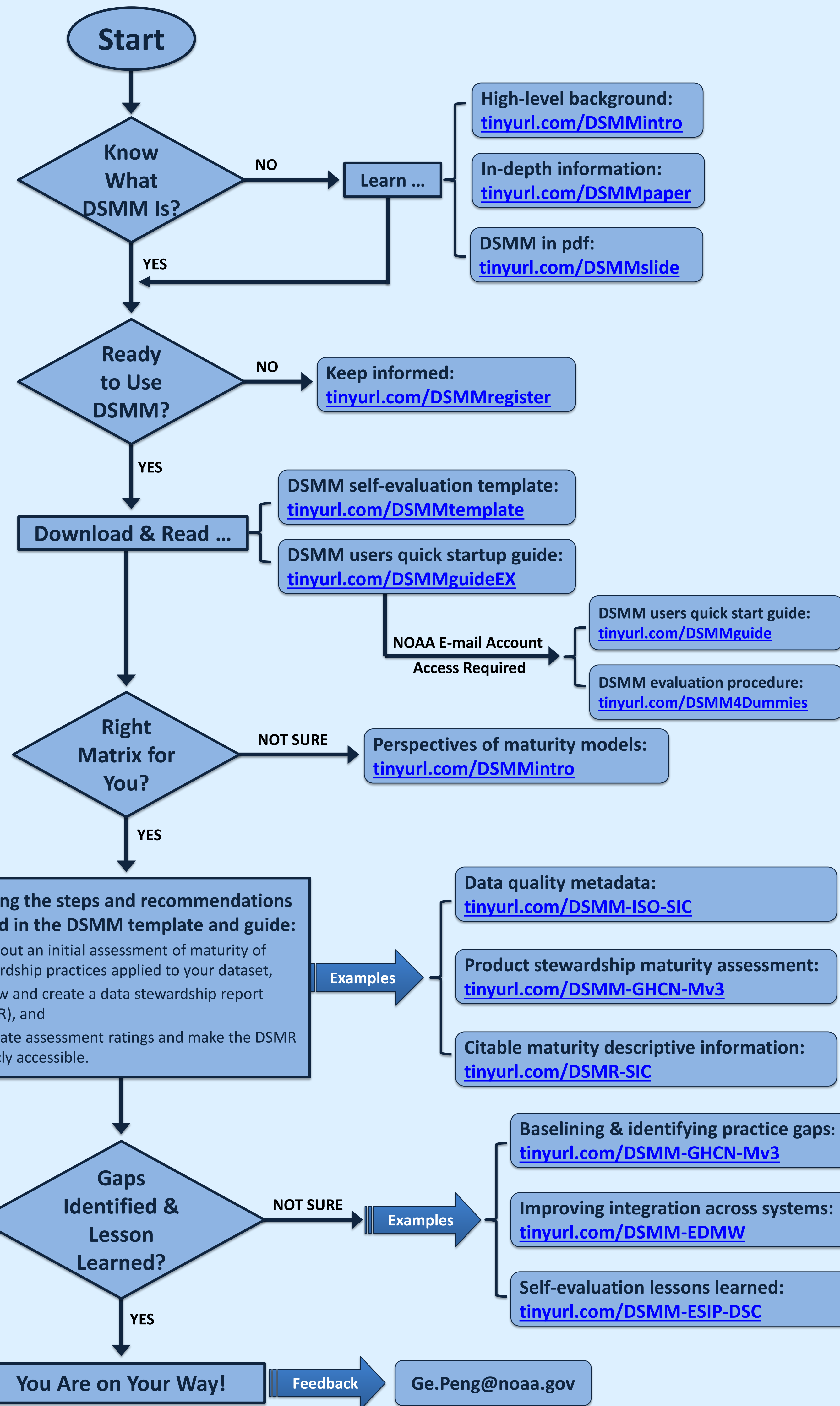
## Introduction

Curating quality descriptive information and metadata for datasets is a necessary step toward meeting the transparency requirement and helping establish the credibility and trustworthiness of individual data products. This, however, has been a difficult challenge for the data management community due to the lack of a consistent assessment framework, process, and workflow. Furthermore, developing and implementing these require multi-domain knowledge and close cross-disciplinary collaboration.

This presentation will first introduce a data stewardship maturity matrix (DSMM) as a reference framework for assessing stewardship maturity of individual digital datasets. Using the DSMM as an example, this presentation will then demonstrate that it is possible to consistently and systematically curate and publish data quality descriptive information both as citable documents and within ISO metadata records for human and machine end-users. These consistent and citable documents and metadata records can be readily integrated into or linked by other systems and tools to be used, for example, for enhanced data discoverability and usability.

This presentation will also outline the progress made under the auspice of the NOAA OneStop project in the area of consistently curate, publish, integrate, and display data quality information.

## DSMM: What, Scope, Who, How, …

### What Is the NCEI/CICS-NC Scientific Data Stewardship Maturity Matrix (DSMM)?

*A Unified Framework for Measuring Stewardship Practices Applied to Individual Digital Environmental Data Products That Are Publicly Available Online*

Leveraging Institutional Knowledge and Community Best Practices and Standards

### The Scope of Stewardship Practices

Measurable practices associated with the functional entities of the Open Archival Information System (OAIS) (within the shaded box)

(CCSDS Version: 650x0m2-2012; Image Source: Lotar International)

### DSMM Defines Measureable, Five-Level Progressive Practices

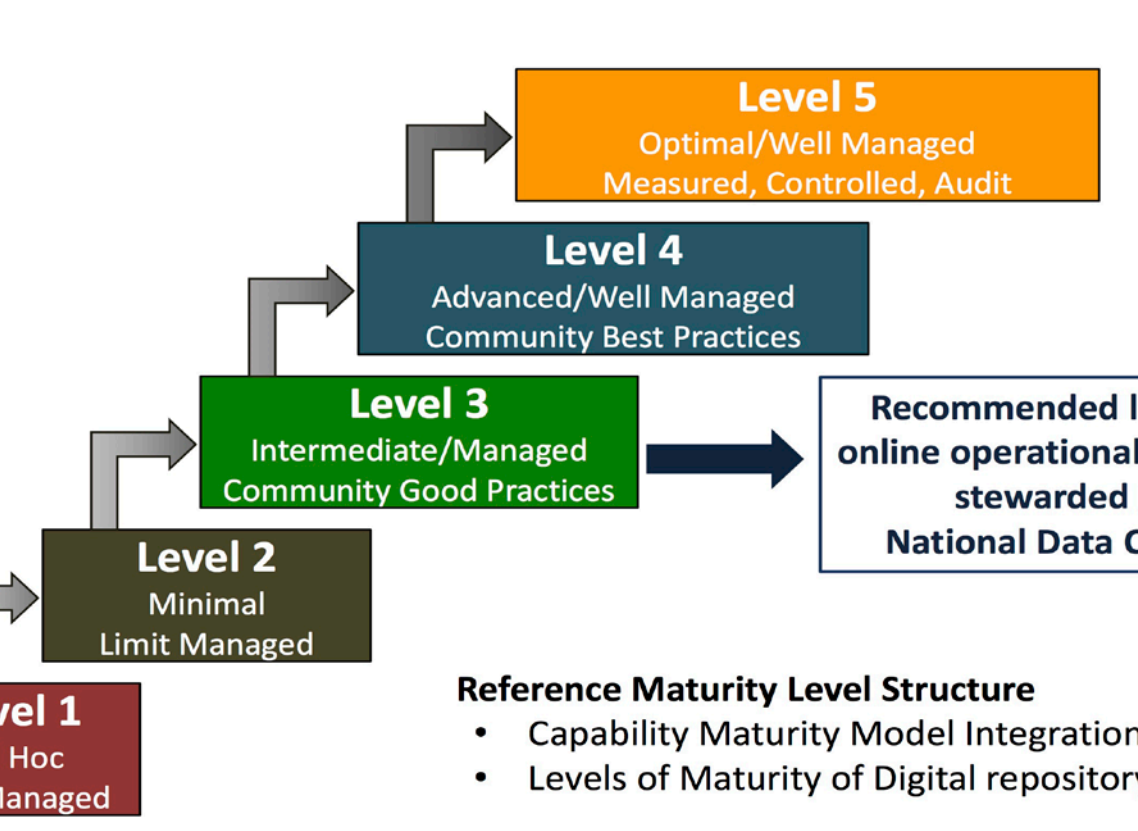| Maturity Scale / Key Component | Level 1 - Ad Hoc Not Managed | Level 2 - Minimal Limited | Level 3 - Intermediate Managed Defined, Partially Implemented | Level 4 - Advanced Well-Defined, Fully Implemented | Level 5 - Optimal Measured, Controlled, Audit |
|---|---|---|---|---|---|
| Preservability | The state of dataset being preservable | | | | |
| Accessibility | The state of dataset being publicly searchable and accessible | | | | |
| Usability | The state of data product being easy to understand and use | | | | |
| Production Sustainability | The state of data production being stable and extendable | | | | |
| Data Quality Assurance | The state of data product quality being assured/screened | | | | |
| Data Quality Control /Monitoring | The state of data product quality being controlled and monitored | | | | |
| Data Quality Assessment | The state of data product quality being assessed | | | | |
| Transparency /Traceability | The state of data product being transparent, trackable, and traceable | | | | |
| Data Integrity | The state of data integrity being verifiable | | | | |

(Data system integrity is also very important but not in the matrix due to potential security risks to the system.)

### DSMM Follows CMMI level Structure

- **Level 5** Optimal/Well Managed Measured, Controlled, Audit
- **Level 4** Advanced/Well Managed Community Best Practices
- **Level 3** Intermediate/Managed Community Good Practices
- **Level 2** Minimal Limit Managed
- **Level 1** Ad Hoc Not Managed

Recommended level for online operational products stewarded by National Data Centers

**Reference Maturity Level Structure**
- Capability Maturity Model Integration (CMMI)
- Levels of Maturity of Digital repository

### Who Could Use The Matrix?

- **Data providers and scientific stewards**
  - to evaluate and improve the quality and usability of their products against community best practices
- **Modelers, decision-support system users, and scientists**
  - to improve their products and uncertainty estimates
  - to make investment and use decision
- **Data managers/stewards of data centers and repositories**
  - to validate their compliance or lack of to community accepted stewardship practice or standards
  - to assess the current state
  - to create a roadmap forward to improve or enhance its stewardship maturity of practices applied to a certain product or all its holdings
- **General data users**
  - to make an educated choice on selecting or utilizing a dataset

### Ways to Utilize DSMM & Assessment Results

- To know the current state of your dataset(s) – maturity assessment (stewardship maturity scoreboard)
- To know where you want or need to be – stewardship requirements
- To know how to get there – roadmap forward (informed, actionable steps)
- A reference model for stewardship planning and resource allocation – informed decision-making support
- A consolidate source and transparency for information about stewardship practices – assessment with detailed justifications
- Content-rich quality metadata – enhanced discoverability and usability

Stewardship Maturity Scoreboard and Roadmap Forward

## Getting to Know and to Use the Data Stewardship Maturity Matrix (DSMM)

**Start**

**Know What DSMM Is?** — NO → **Learn …**
- High-level background: tinyurl.com/DSMMintro
- In-depth information: tinyurl.com/DSMMpaper
- DSMM in pdf: tinyurl.com/DSMMslide

YES ↓

**Ready to Use DSMM?** — NO → Keep informed: tinyurl.com/DSMMregister

YES ↓

**Download & Read …**
- DSMM self-evaluation template: tinyurl.com/DSMMtemplate
- DSMM users quick startup guide: tinyurl.com/DSMMguideEX
- **NOAA E-mail Account Access Required**
  - DSMM users quick start guide: tinyurl.com/DSMMguide
  - DSMM evaluation procedure: tinyurl.com/DSMM4Dummies

↓

**Right Matrix for You?** — NOT SURE → Perspectives of maturity models: tinyurl.com/DSMMintro

YES ↓

**Following the steps and recommendations outlined in the DSMM template and guide:**
- Carry out an initial assessment of maturity of stewardship practices applied to your dataset,
- Review and create a data stewardship report (DSMR), and
- Integrate assessment ratings and make the DSMR publicly accessible.

**Examples** →
- Data quality metadata: tinyurl.com/DSMM-ISO-SIC
- Product stewardship maturity assessment: tinyurl.com/DSMM-GHCN-Mv3
- Citable maturity descriptive information: tinyurl.com/DSMR-SIC

↓

**Gaps Identified & Lesson Learned?** — NOT SURE → **Examples**
- Baselining & identifying practice gaps: tinyurl.com/DSMM-GHCN-Mv3
- Improving integration across systems: tinyurl.com/DSMM-EDMW
- Self-evaluation lessons learned: tinyurl.com/DSMM-ESIP-DSC

YES ↓

**You Are on Your Way!** — **Feedback** → Ge.Peng@noaa.gov

## Pathway to Application of DSMM

### Communication and Use Case Studies

#### Selected NCEI Core Datasets

| Data Type | Dataset | Status |
|---|---|---|
| Satellite – polar ocean | NOAA/NSIDC Sea Ice Concentration CDR | Baselined |
| GIS - regional | NCEI-CO Digital Elevation Models (DEM) | Revised assessment draft review |
| Station - in situ - land | GHCN-M | Baselined |
| Station - gridded - land | National Climate Division (nCliDiv) | Not yet started |
| Satellite - global ocean | Optimum Interpolation Sea Surface Temperature (OISST) CDR | Baselined |
| Physical Records - In Situ Monthly Summaries | Local Climatological Data | Initial assessment draft review |
| Paleo – global land | NOAA/WDS International Tree-Ring Data Bank (ITRDB) | Baselined |

#### Selected ESIP Datasets

| Data Type | Dataset | Status |
|---|---|---|
| Model Reanalysis | NCAR Global Climate Four-Dimensional Data Assimilation Hourly 40km Reanalysis | Baselined * |
| Ecological Data | DataOne Member Node SBC LTER (Long Term Ecological Research) Network | Revised assessment draft review |
| Long-tail Data | NSF ACADIS (Advanced Cooperative Arctic Data and Information Service) | Initial assessment draft |
| Socioeconomic Data | NASA Socioeconomic Data | Initial assessment draft |
| Paleo Data | Australia Borehole Data | Not yet assessed |

### Dataset Stewardship Maturity Evaluation: Guidance, Template, Training, and Tools

**DSMM Applied:** # of NCEI Datasets by Data Groups

700+ datasets

## Consistent and Citable Data Quality Descriptive Information

Developing and Improving Evaluation, Curation, and Integration Workflows

### DSMM Integration Workflow

### Consistent Rating Displacement

### Data Stewardship Maturity Reports (DSMR)
- Data product quality descriptive information documents,
- Consistent document layout,
- Automated generation workflow,
- Unique persistent document identifier,
- To be published and archived by the NOAA Central Library.

### Data Stewardship Maturity Matrix (DSMM) Use Case Submission

### DSMM Data Flow Chart

DSMM Results → Google Form (To be replaced by a web application) → DSMM Automation Tools → DSMM Reports / DSMM diagrams / DSMM diagram templates / ISO metadata templates → DSMM Reports / ISO Metadata Records

### Integration with CEdit

### Consistent DSMM ISO 19115 Metadata Implementation

## Acknowledgment