# USGIN DRAFT SPECIFICATIONS FOR GEOSCIENCE RESOURCE METADATA AND IDENTIFIERS

by Stephen Richard, Wolfgang Grunberg, and Ryan Clark, US Geoscience Information Network, Arizona Geological Survey, Tucson, Arizona
Geological Society of America, National Meeting, Nov. 1, 2010, Denver, Colorado

The catalog: a key system component in which data providers register and describe for discovery and access by users.

The metadata offered by the catalog must
-- guide search engines to locate requested resources
-- describe the resource content, provenance, and quality so users can determine if the resource will serve for intended usage
-- enable the user to obtain or access the resource
-- enable software to automatically access the resource

The USGIN project has been working on specifications to reduce the heterogeneity of 'standard' metadata and service implementations such that a single client can search against different catalogs, and successfully use the results. See **http://lab.usgin.org/profiles/usgin-iso-metadata-profile** for more information.

Discovery of resources across a wide spectrum of related geoscience disciplines requires a metadata content model with concepts that span those disciplines.

Specification-- instructions on how to do something
Standard-- a specification that alot of people use

**Summary of USGIN metadata content recommendations**.
See **http://lab.usgin.org/profiles/doc/metadata-content-recommendations** for complete discussion.

## Catalog system



**Producer cloud** — Heterogeneous
**Homogeneous content** — **Homogeneous format**
**User cloud**

producers announce what they've got

Cross-domain generalized metadata and specialized discipline specific descriptions can use same pattern with different content models

Users search using many different applications

## Resource Identifiers

**Stable, dereferenceable identifiers are an essential component of a linked data system**

Stable: mapping between identifier string and identified thing is immutable.

Dereferenceable: a well know system exists to present user with an explanation or representation of what is identified.

| Resource type | notes | examples |
|---|---|---|
| concrete information resource | a specific representation of a resource. Equivalence is determined by bitwise matching of the representation bitstream | monaLisa.tif, USGIN_URI_scheme_1.0.1.doc |
| abstract information resource | an information resource defined by a particular content scope, includes most kinds of documents. Equivalence is determined by matching of content meaning or intention (depending on how the resource is defined) | Image of Mona Lisa, current version of USGIN URI scheme, the Declaration of Independence |
| concrete non-information resouce | An identifiable physical object or event. | The original Mona Lisa painting, The original Declaration of Independence parchment, a particular rock sample, a particular car, a particular meeting, a particular project. |
| abstract non-information resource | A concept defined by some human intention, may categorize other resources, or represent some abstract idea. | painting, love, the US Congress |

**Kinds of resources that may be identifier by a URI**

**Dereferencing URI for information resource:**
--- Present the normative (canonical form) actual resource.

**Dereferencing URI for a non-information resource:**
--- Present a representation of the resource. There may be many valid representations of non-information resources, useful in different contexts. Each of these is an identifiable information resource and the canonical or default representation may be context dependent.

## URI Scheme (see http://lab.usgin.org/profiles/doc/usgin-uri-recommendations-v10)

"http:" "//" uriHost "/uri_gin/" nameAuthority "/" resourcePath ["/" resourceSpecificString ]["/"/"/" representationPart ]]

**uriHost** -- host system on the internet that will de-reference the following identifier. Not considered part of the USGIN resource identifier for string comparisons of identifier equivalence, but usage of a known host conveys trust. Syntax follows standard web conventions, e.g. "azgs.az.gov" .
*The sequence of characters between "uri_gin/" and the last "/" in the URI constitute the identifier for a resource.*

**nameAuthority** Name of the naming authority for the identified resource, taken from the register with USGIN URI. A name authority defines the collection of resource categories for which it maintains identifiers.

**resourcePart** -- A path that identifies a resource category and specific instance of a resource of that type. Includes:

**resourcePath** Specifies the resource category, taken from the register at http://resources.usgin.org/uri_gin/usgin/register/resourceType/. Each resource type definition should specify what the default and canonical representations are (they may be the same or different), the known valid representations, and the dereferencing behavior for the resource type. The resourcePath is hierarchical to allow representation of resource hierarchy when that is useful.

**resourceSpecificString** A string that identifies a particular instance of the resource type in the resource path, unique in the scope of the name authority. This string may have syntax specially scoped for a particular resource type. Definition of resourceSpecificString syntax should be in the resource type definition.

**representationPart** is intended to identify a representation of the preceding resource part. Dereferencing behavior for such an identifier is defined in the associated resource type definition. Note that an abstract document is an information resource, so in this case the representationPart is semantically equivalent to a resourceSpecificString.

## Metadata model



See: http://lab.usgin.org/groups/metadata-interest-group/conceptual-model-general-purpose-resource-metadata

Based on ISO 19115 and FGDC content.
Modified to be explicit while avoiding implementation-specific artifacts where possible
Focus is on attributes that will apply to broad spectrum of resources
Discipline-specific details are implemented by extensions from DescriptionDetails abstract class
Various details of associated with the various base attribute classes are not shown in this diagram. See http://lab.usgin.org/groups/ metadata-interest-group for the full model UML.
Implementation of this model as an ontology and mapping of various schemes in use to a single model would enable better metadata interoperability.

## Applications



Online tool for generating metadata in USGIN profile; will export FGDC xml or ISO19139 xml. Access at http://mw.usgin.org



Light weight catalog client will search a variety of metadata catalogs. Code project started by IDEC (Zaragoza, Spain) with contributions from USGIN. http://catalog.usgin.org/catalogConnector
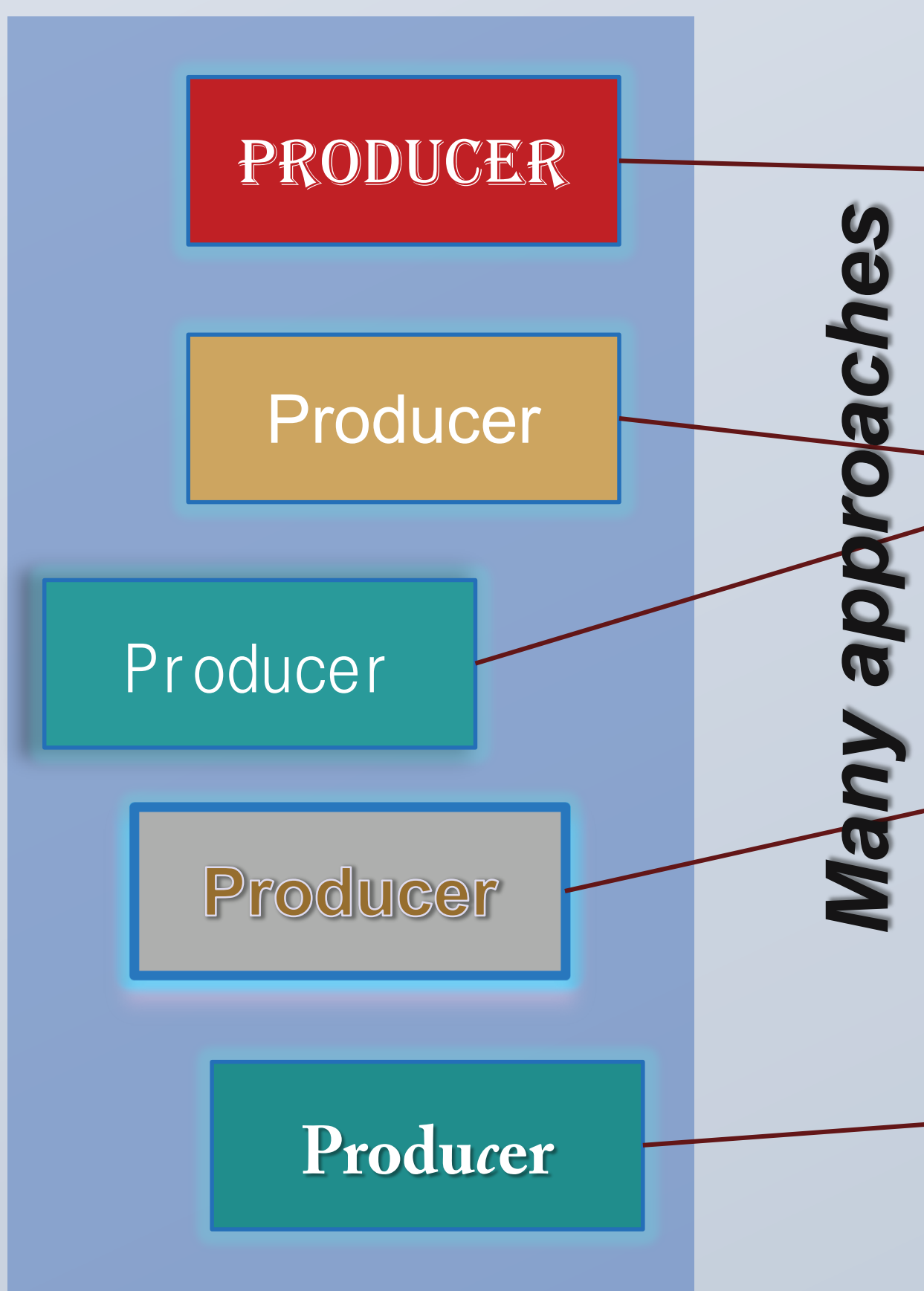
### Resource
- **Title** (1 entry): Succinct (preferably <250 characters) name of the resource.
- **Description** (1 entry): Information about the resource's content as well as its context.
- **Originators** (1 to many entries): Authors, editors, or corporate authors/curators of the resource.
- **Publication Date** (1 entry): Publication, origination, or update date (not temporal extent) for the resource. Use ISO 8601 date and time format. If no publication date is known, estimate the publication date range, enter the oldest year as the publication date, and include the estimated date range in the Description field.
- **Geographic Extent - Horizontal** (1 entry, point or minimum bounding rectangle): North and South Bounding Latitude or Point Latitude, East and West Bounding Longitude or Point Longitude. Values given in decimal degrees using the WGS 84 datum. Some resources may not be usefully described by an extent; if no extent is specified the default is Earth. If a particular encoding scheme requires a bounding box, a minimum bounding rectangle will be created if only a point coordinates are given.
- **Author or Intellectual Originator** (0 to 1 entry): The primary party responsible for creating the resource. Organization Name, Person Name, Position Name; contact point (Telephone, E-mail, postal).
- **Bibliographic Citation** (0 to 1 entry): Full bibliographic citation if the resource has been published.
- **Subject Keywords** (0 to many entries): Thematic, spatial and temporal free-form subject descriptors for the resource. If possible, submit keywords in separate Thematic, Spatial, and Temporal keyword categories.
- **Resource Language** (0 to 1 entry): Language of resource content. Use three letter ISO 639-2 language code (defaults to "eng" for English).
- **Resource ID** (0 to many entries): Resource identifier(s) following any public or institutional standard. Identifier consists of an identifier string and if applicable a Resource ID Protocol identifier string that specifies the protocol for the resource ID standard. For example: unidefined, ISBN-10, ISBN-13, ISSN, URN, URI, IRI, DOI, HTTP, SSN, etc. Examples: doi:10.1000/182; isbn:0-671-62964-6; issn:1935-6862; azgs:OFR-10-02. Many protocols build the identifier for the protocol into the identifier string.
- **Geographic Extent – Vertical** (0 to 1 entry): Compound content includes Datum Elevation, Datum Type, Maximum Elevation, Minimum Elevation. Maximum and Minimum Elevations are relative to the reported datum elevation, which will typically be the Earth surface at the location of the resource or sea level. Datum Elevation must be reported relative to mean sea level (MSL) in meters using EPSG::5714 geodetic parameters (WGS 84). Datum type must be a controlled vocabulary (Earth surface, MSL, Kelly bushing, etc.). The maximum is always numerically greater than the minimum elevation. For boreholes with datum at the earth surface, depth below surface is reported as a negative number. Example: core from borehole at depths between 100 and 470 feet, borehole collar at 4787 feet above sea level. Vertical extent could be reported in either of the following ways: {0, "MSL", 1420, 1308} or {1450.6, "Earth surface", -30.3, -142.4}.

- **Temporal Extent** – Temporal range over which the resource was collected or is valid. If the resource pertains to specific named geologic time periods, those terms should be entered as keywords (preferable as Temporal Keywords). Start Date (0 to 1 entry), End Date (0 to 1 entry; required if start date exists; use ISO 8601 date and time format.
- **Quality Statement** (0 to 1 entry): Text specification of the quality of the resource.
- **Lineage Statement** (0 to 1 entry): Text description of the resource's provenance. Complex discipline-specific lineage (e.g. analytical procedures) are described in free text, and may be supplemented by discipline-specific content models.

**Access**
- **Access Statement** (1 entry): Text instructions for how to access the resource.
- **Distribution Contact** (0 to 1 entry): Party to contact about accessing the resource. Organization Name, Person Name, Position Name; contact point either telephone, e-mail, or postal address. In general, a contact for distribution should be required for physical resources.
- **Link to resource** (0 to many entries): A URL pointing to a resource or resource webpage. Element content includes URL, Link Function, and Representation Format. URL is minimum content required if a link is included. Link Function term from controlled vocabulary specifies what a HTTP GET using the URL will invoke. The link might return an html page, electronic document in some other format, an end point for a service, an online application that requires user interaction, etc. Representation Format is a controlled vocabulary term specifying the format (MIME media types) of a file-based response if applicable.
- **Constraints Statement** (0 to 1 entry): describe the resource's legal and usage constraints in free text.
- **Distribution Keywords** (0 to many entries): keywords describing the physical form of the resource (core, rock sample, digital file, book, journal article), formatting of resource content (file format, e.g. tiff, xls, MIME type), or physical distribution media (film, floppy disk, online service, hard copy). Use of these keywords allows users to search for particular kinds of artifacts.

**Metadata**
- **Metadata Date** (1 entry): Last metadata update/creation date-time stamp in ISO 8601 date and time format.
- **Metadata Contact** (1 entry): The party to contact with questions about the metadata itself. Organization Name, Person Name, Position Name; contact point either telephone, e-mail, or postal address.
- **Metadata Specification** (0 to 1 entry): Identifier string for the metadata specification used to create the metadata record encoding this content. Should indicate the base standard and version, as well as any profile that applies to the content or encoding. Ideally the identifier could be dereferenced to obtain information about the applicable specification. Identifiers for metadata encoding specifications will have to be formally defined and registered for such identifiers to be broadly useful.
- **Metadata UUID** (0 to 1 entry): A Universally Unique Identifier (UUID) for the metadata record; required to avoid duplicate entries across multiple metadata catalogs. The UUID format provides unique identificatio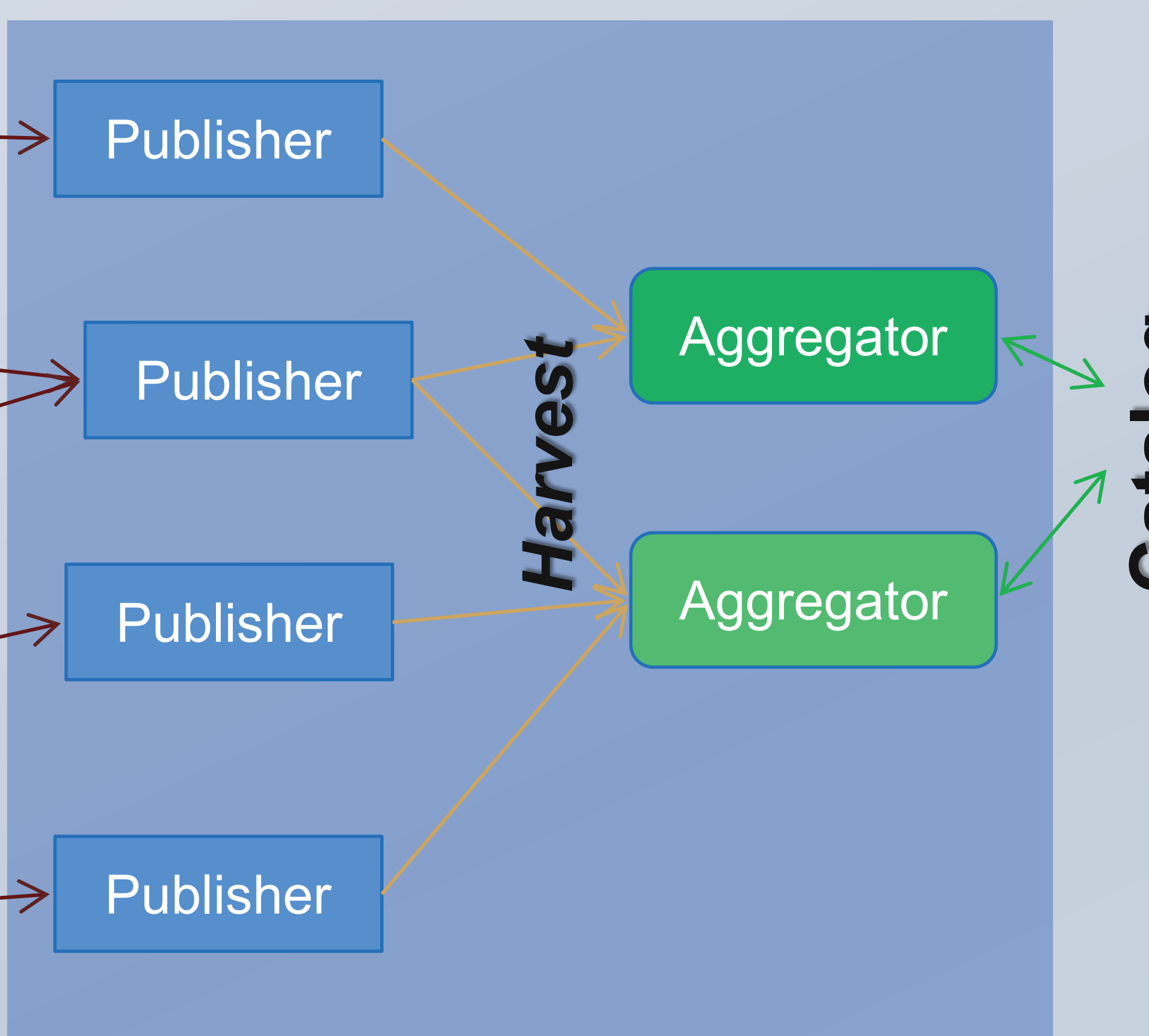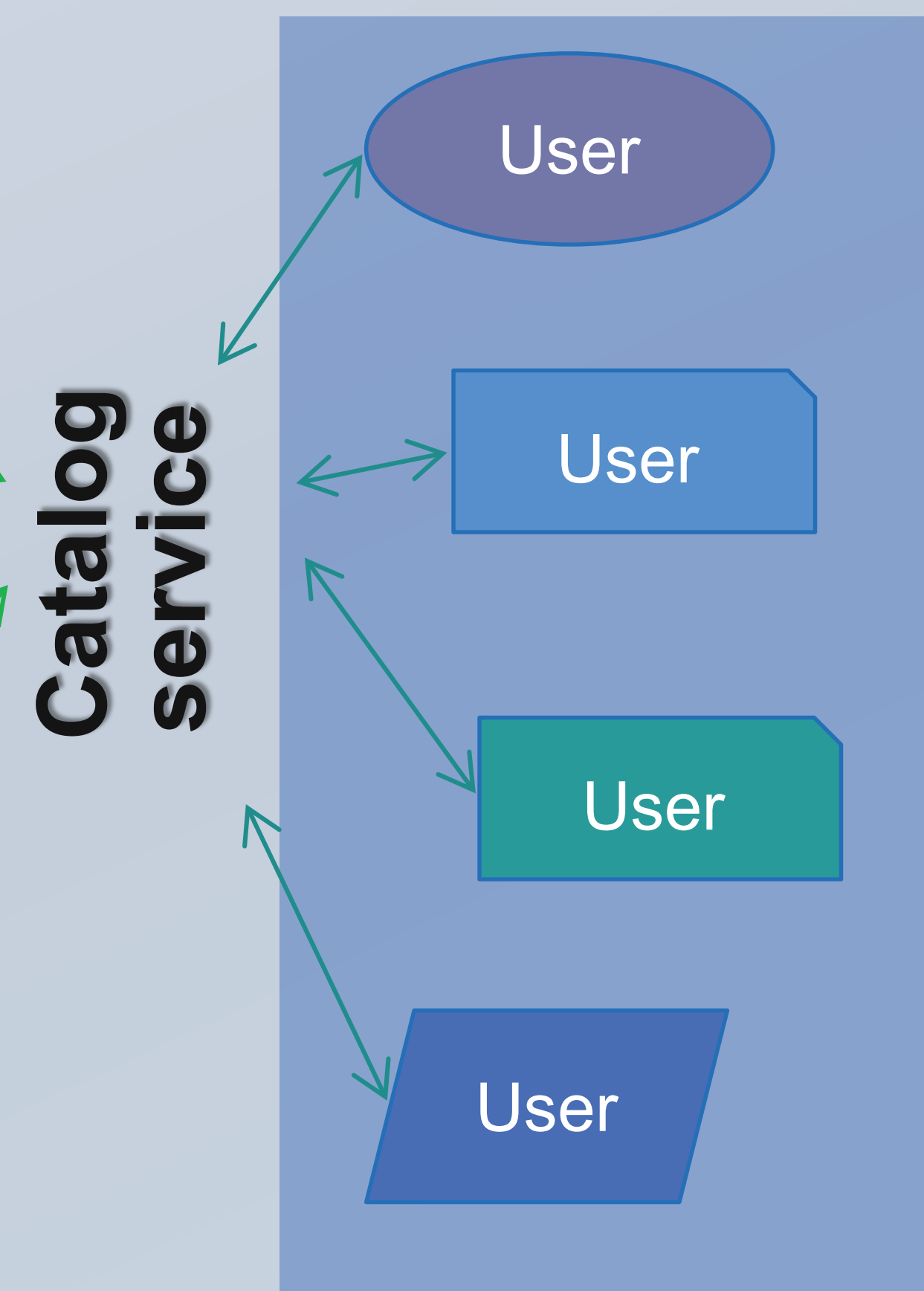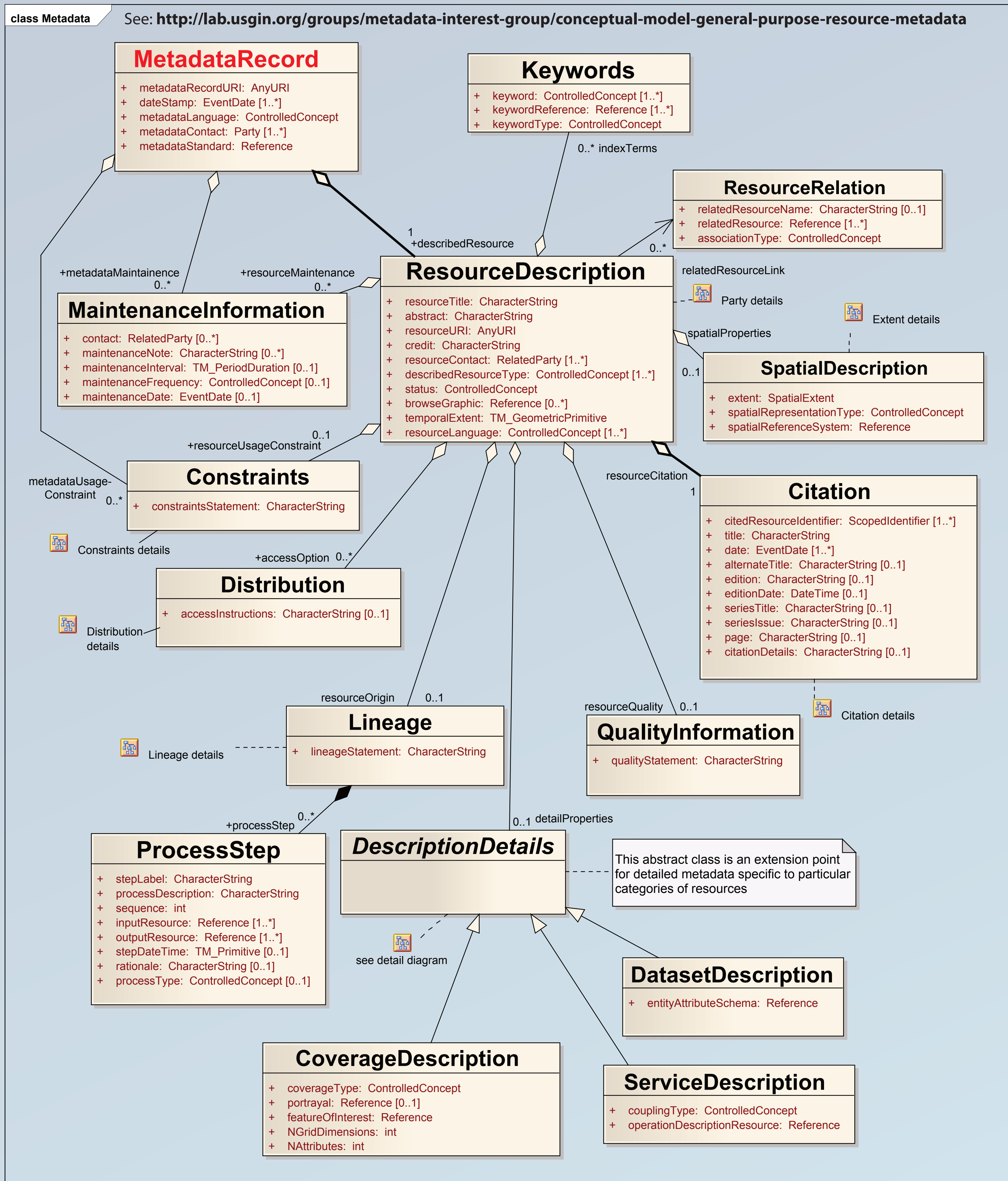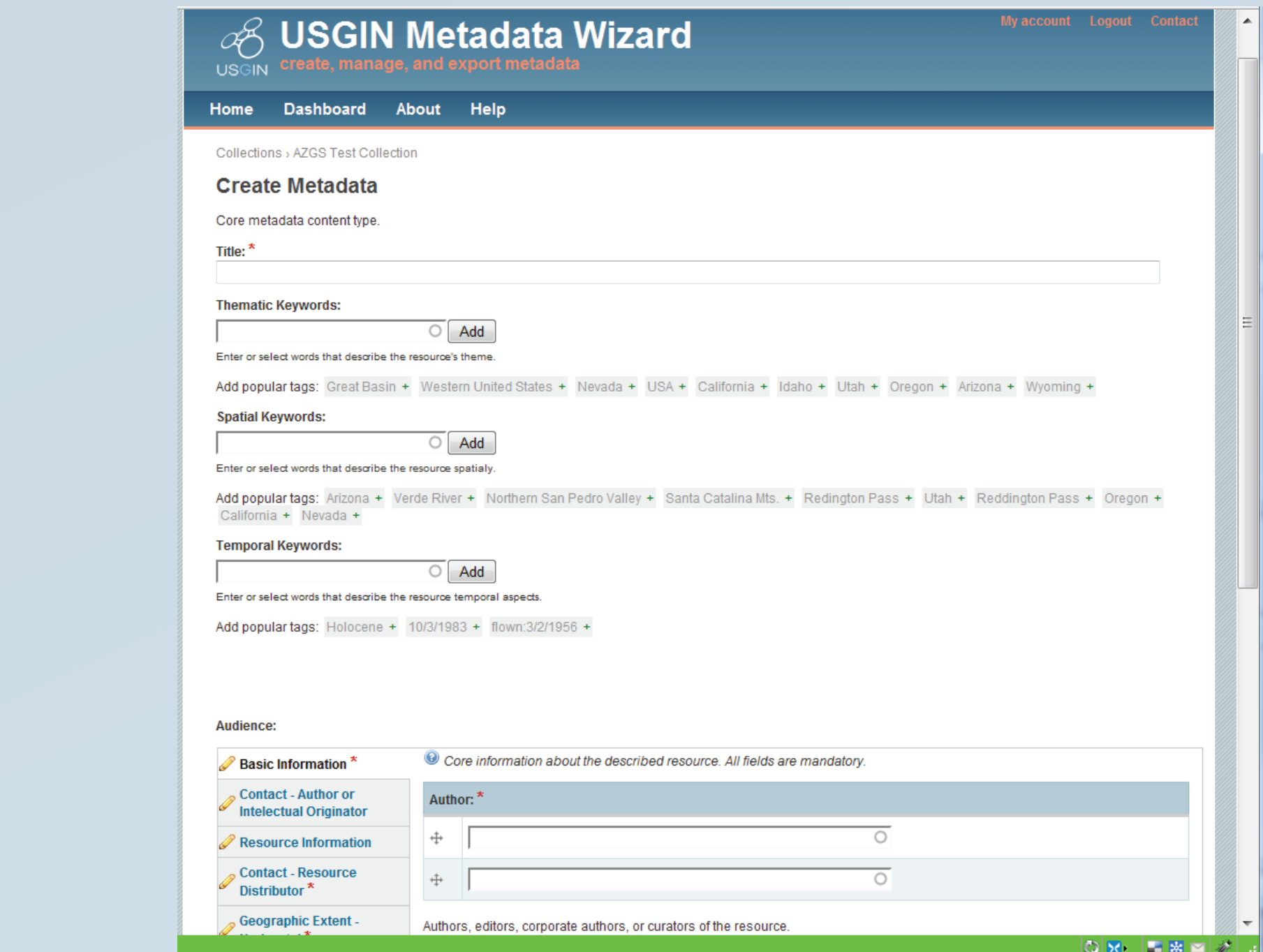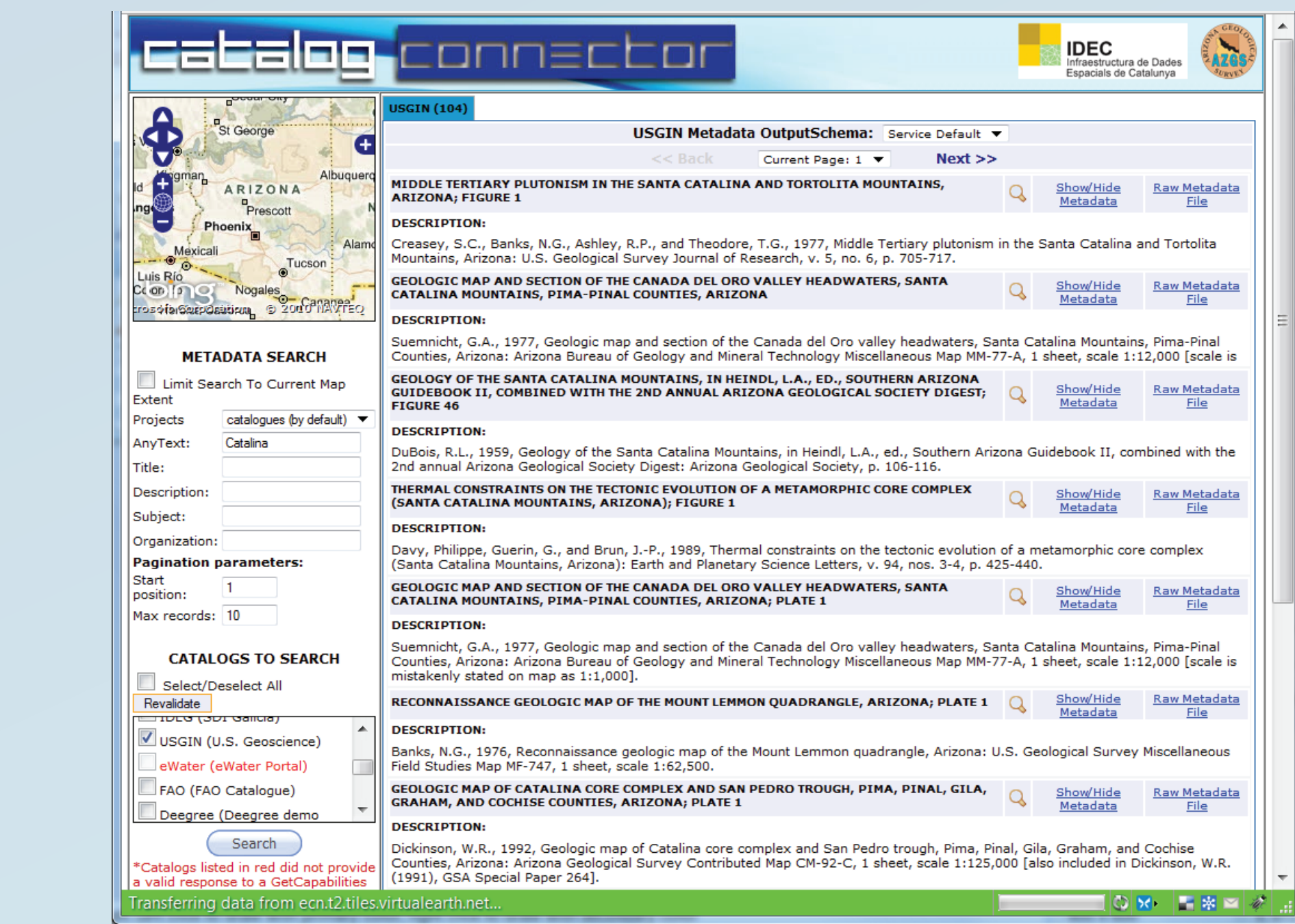n without centralized coordination.