

Le Système d'Information Theia/OZCAR :

Un portail d'accès à l'ensemble des données in-situ des surfaces continentales

Equipe projet

Coordination scientifique:

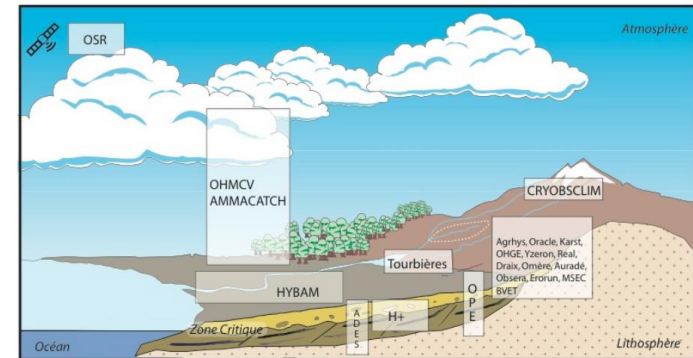
- Sylvie Galle (IRD/IGE, CMI INSU-SIC données in situ, membre BE Theia)
- Isabelle Braud (Irstea Lyon, co-animatrice de l'IR OZCAR)

Equipe Développement:

- Véronique Chaffard (IRD/IGE, responsable technique)
- Charly Coussot (OSUG, ingénieur développement, CDD depuis le 1/10/2017)

Le réseau d'observatoires de la Zone Critique

- **Infrastructure de Recherche (IR) française sur la zone critique** créée en 2016.
- OZCAR met en réseau les observatoires de long terme de la zone critique depuis la haute montagne jusqu'à la mer.
- **21 "observatoires élémentaires" labélisés** collectent **>300 variables différentes** (physiques et géochimiques) sur **~100 sites**



- Chaque observatoire a une histoire différente et a son propre système de gestion et diffusion des données
- OZCAR fait partie de l'**IR européenne European Long Term Ecological Research** (eLTER-RI 2018)

OZCAR : un réseau de 21 observatoires



- Documente ~ 60 sites
- En France et dans les pays du Sud



Une grande diversité de mesures in situ

Observations

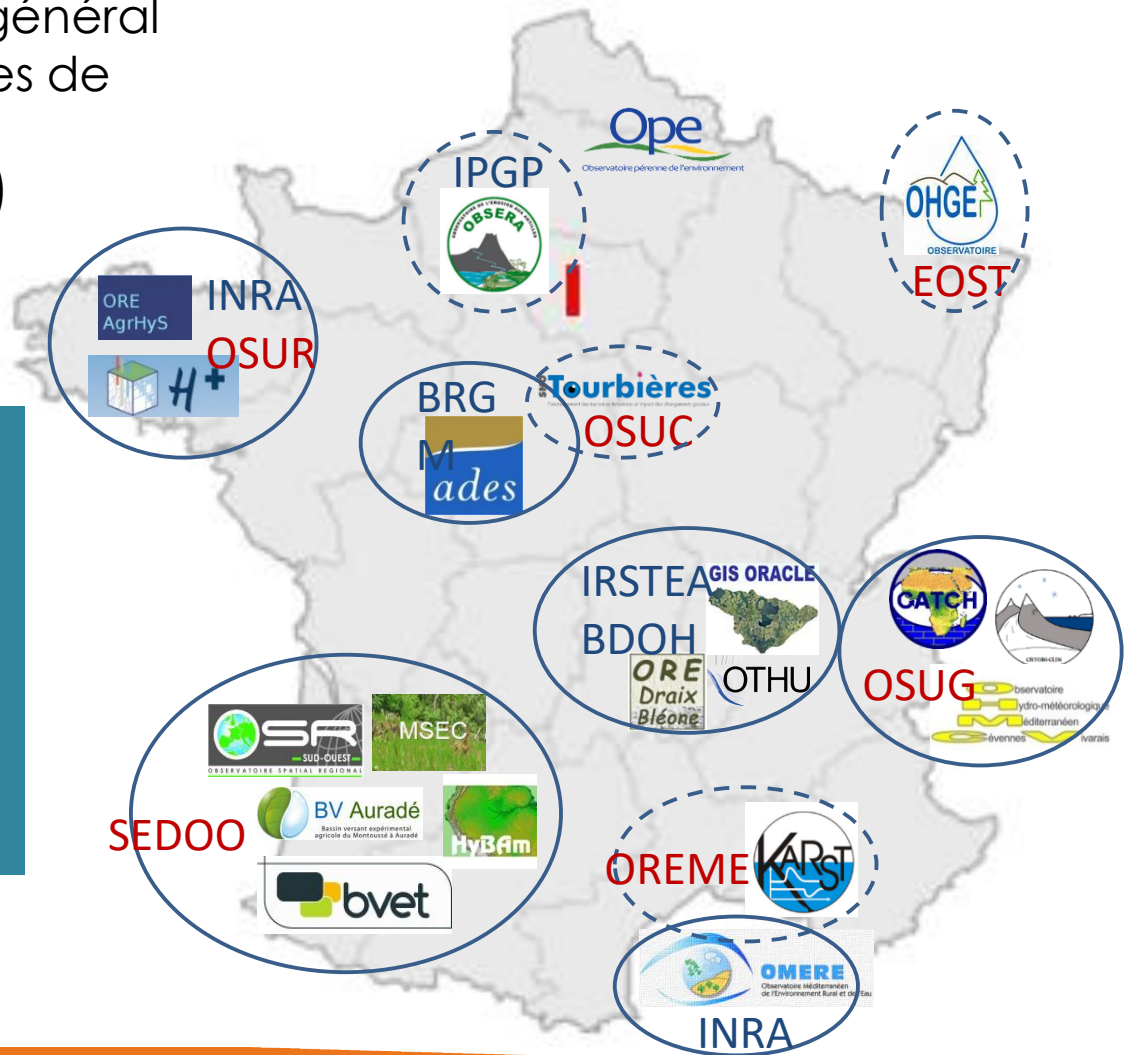
- Transport de l'eau, des sédiments et des solutés
- Bilans d'énergie de surface
- Exploration géophysique
- Occupation des sols

Objets d'intérêt

Bassins versants, aquifères, rivières, glaciers, permafrost...

Le panorama des SI des observatoires d'OZCAR

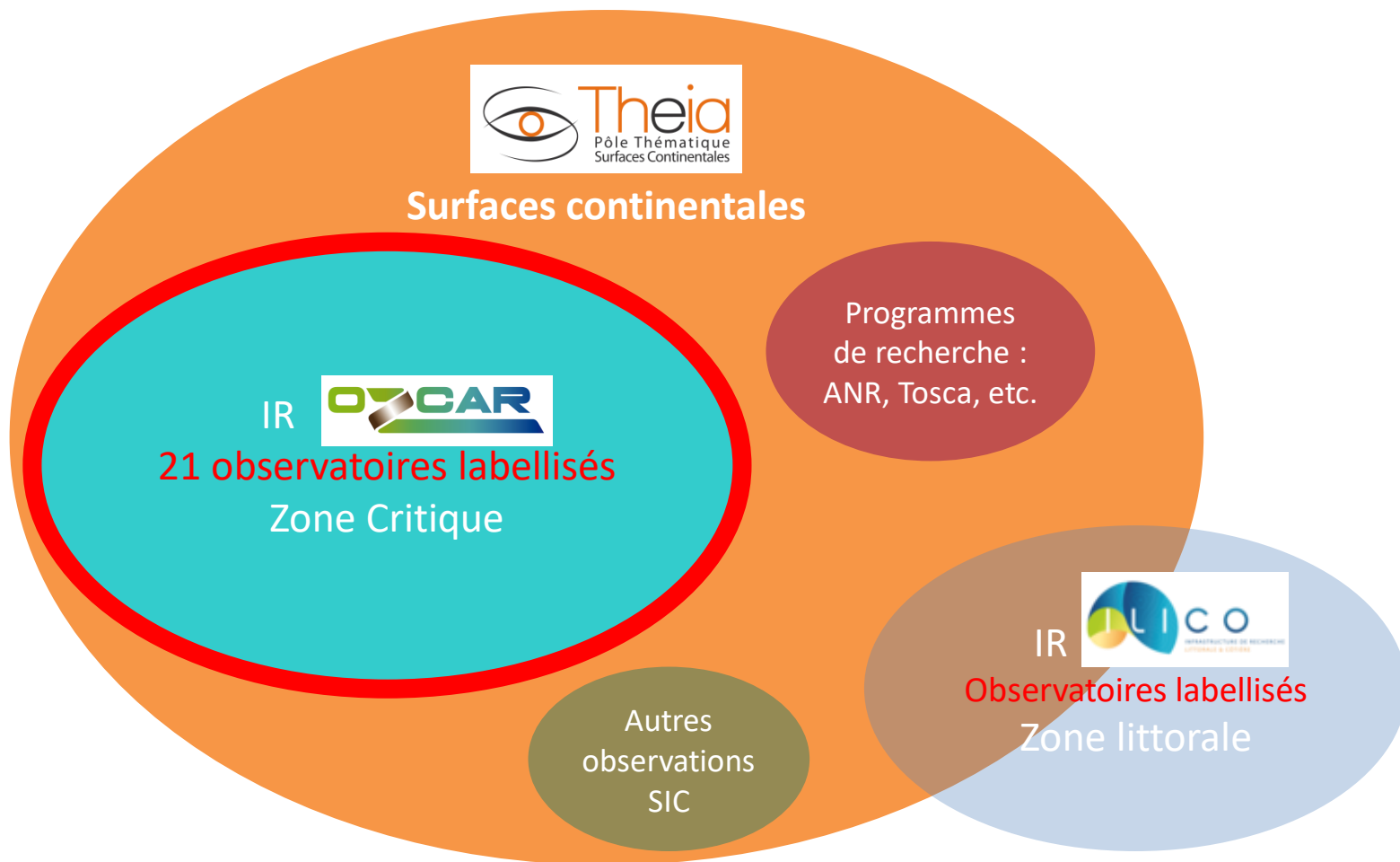
Des observatoires en général
rattachés à des centres de
données
(OSU ou institutionnels)



Une grande
hétérogénéité dans la
structuration des
données (granularité),
dans la maturité des SI,
dans les noms de
variables

OZCAR / Theia : quelle différence?

OZCAR et Theia n'ont pas le même périmètre

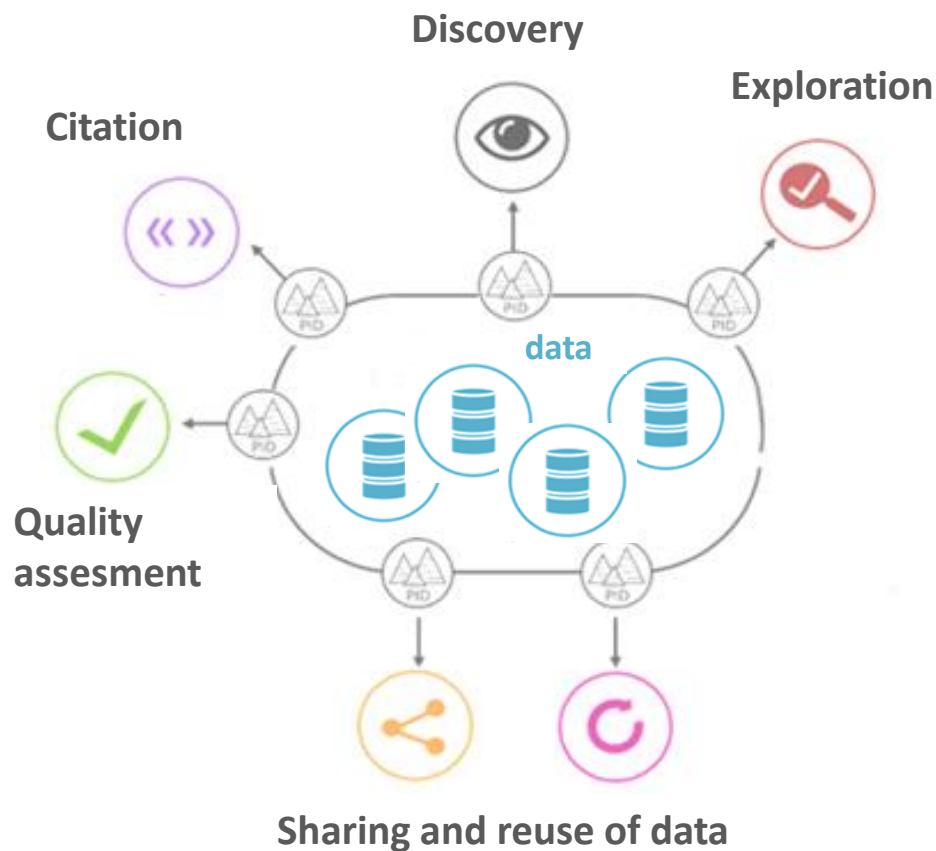


Objectifs du SI Theia/OZCAR

- **Un portail unique** des données d'observation in-situ des surfaces continentales qui permettra un accès **transparent** aux utilisateurs
- Un portail qui ne stocke pas les données qui sont déjà gérées par ailleurs mais qui met en place des flux d'informations qui les collecte
- Un système d'information qui respecte les **critères d'accessibilité et d'interopérabilité** de la Directive Inspire et de standards internationaux
- **Un système interopérable** avec les systèmes d'information français (pôle de données Theia, métacatalogue IR Data Terra) et européen en cours de construction (e-LTER european Long Term Ecological Research)
- Encourager la déclaration de **DOI de données**

Objectifs du SI Theia/OZCAR

Favoriser la **découverte** et l'**exploration** des données, leur **partage** et **réutilisation**, leur **citation**



Stratégie de construction:

- Dialogue avec les scientifiques et les informaticiens des observatoires
- Définition d'un **modèle de métadonnées commun (modèle pivot)** pour mettre en place en entrée du SI Theia/OZCAR des **flux de données** avec le SI des producteurs (éléments de métadonnées standards)
- Mettre en place en sortie du SI des **standards d'échange** (web service)
- Utilisation d'un **vocabulaire contrôlé commun pour les noms de variable** (à partir de référentiel existant)

Mise en oeuvre

1. Construction d'un **thésaurus** pour les noms et catégories de variable
2. Travail sur les **métadonnées**: analyse des **standards** et identification des **flux d'informations** à mettre en place
3. Définition d'un **format pivot (data model)** pour l'échange de données
4. Architecture du système et développement d'un prototype de portail

Construction d'un thésaurus pour les noms et catégories de variables

Objectif:

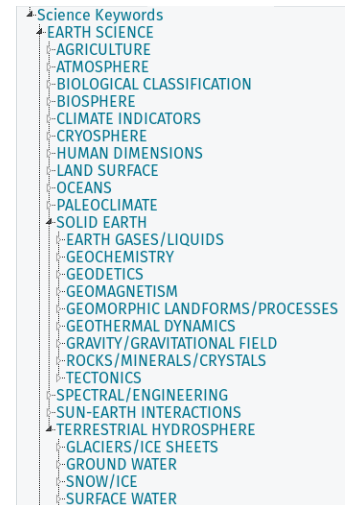
Permettre de rechercher la donnée par un filtrage sur les variables.

Les étapes de construction: démarche bottom-up

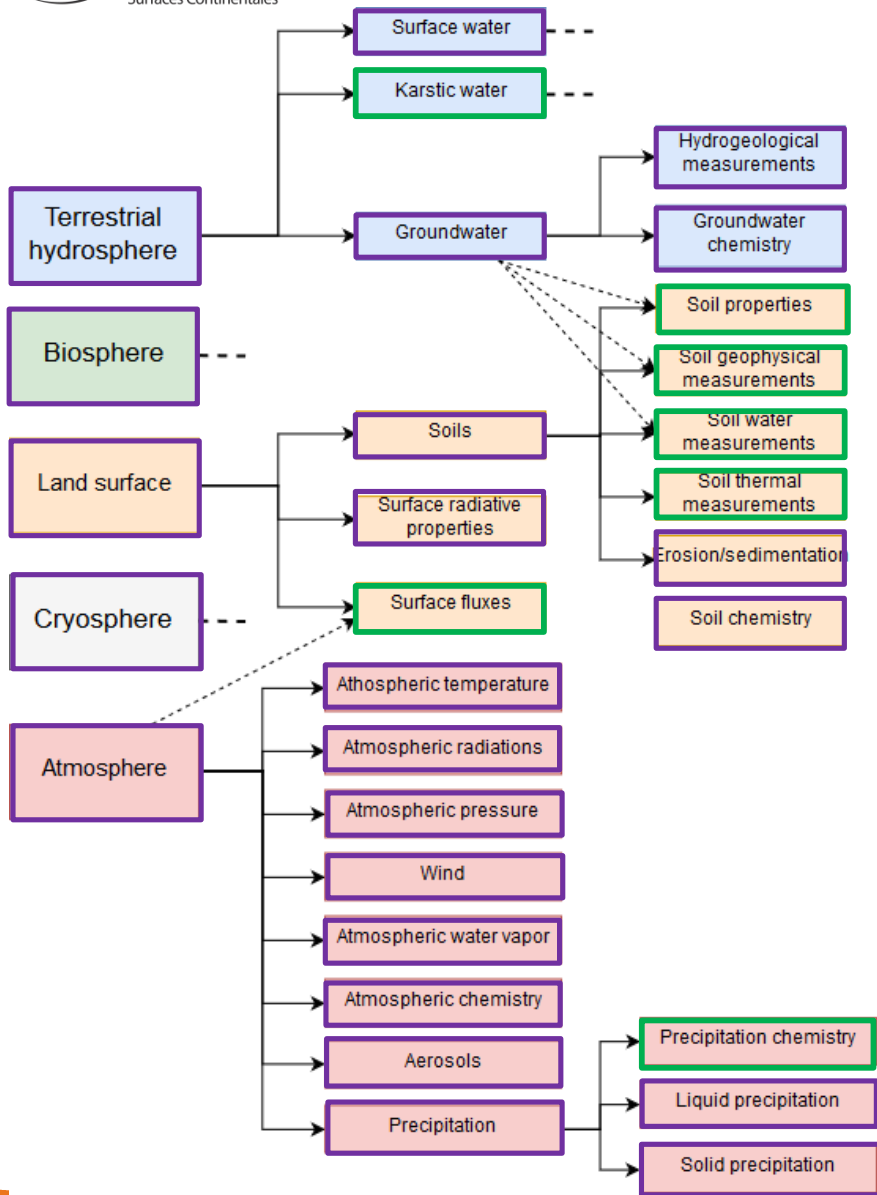
1. Collecte des noms variables producteurs

Catégorisation/classification de l'ensemble des noms de variables vis-à-vis du thésaurus de la [NASA GCMD Science keywords](#).

Structure l'information à l'aide de concepts **hiérarchisés**



2. Création d'un **vocabulaire de noms de variables Theia/OZCAR** en se basant sur le thésaurus GCMD et mapping avec nom producteur avec les scientifiques



3. Hiérarchisation par catégories de variables

La taxonomie GCMD pas toujours adaptée aux besoins Theia/OZCAR (pas assez précise)

-> **nécessité d'ajouter des catégories propres à Theia/OZCAR**

<https://wiki.earthdata.nasa.gov/display/CMR/GCMD+Keyword+Access>

GCMD categories

<https://earthdata.nasa.gov/about/gcmd/>

Added categories

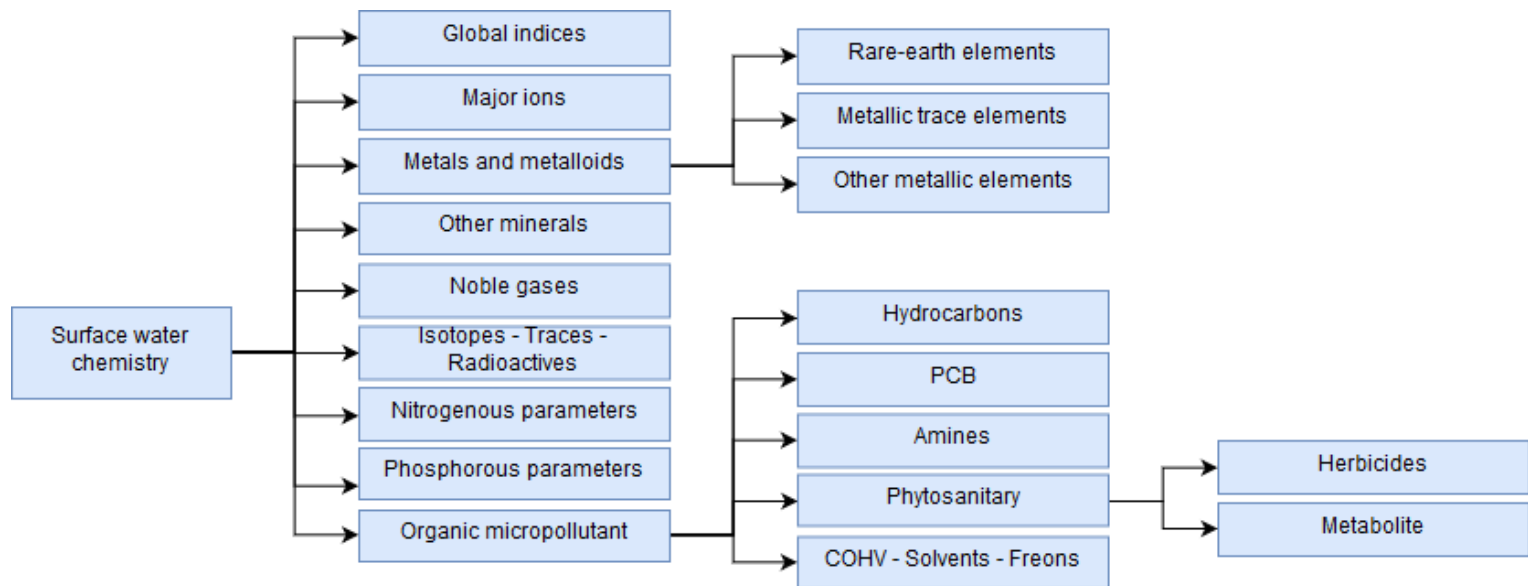
Relevant for OZCAR thematics

3. Hiérarchisation par catégories de variables

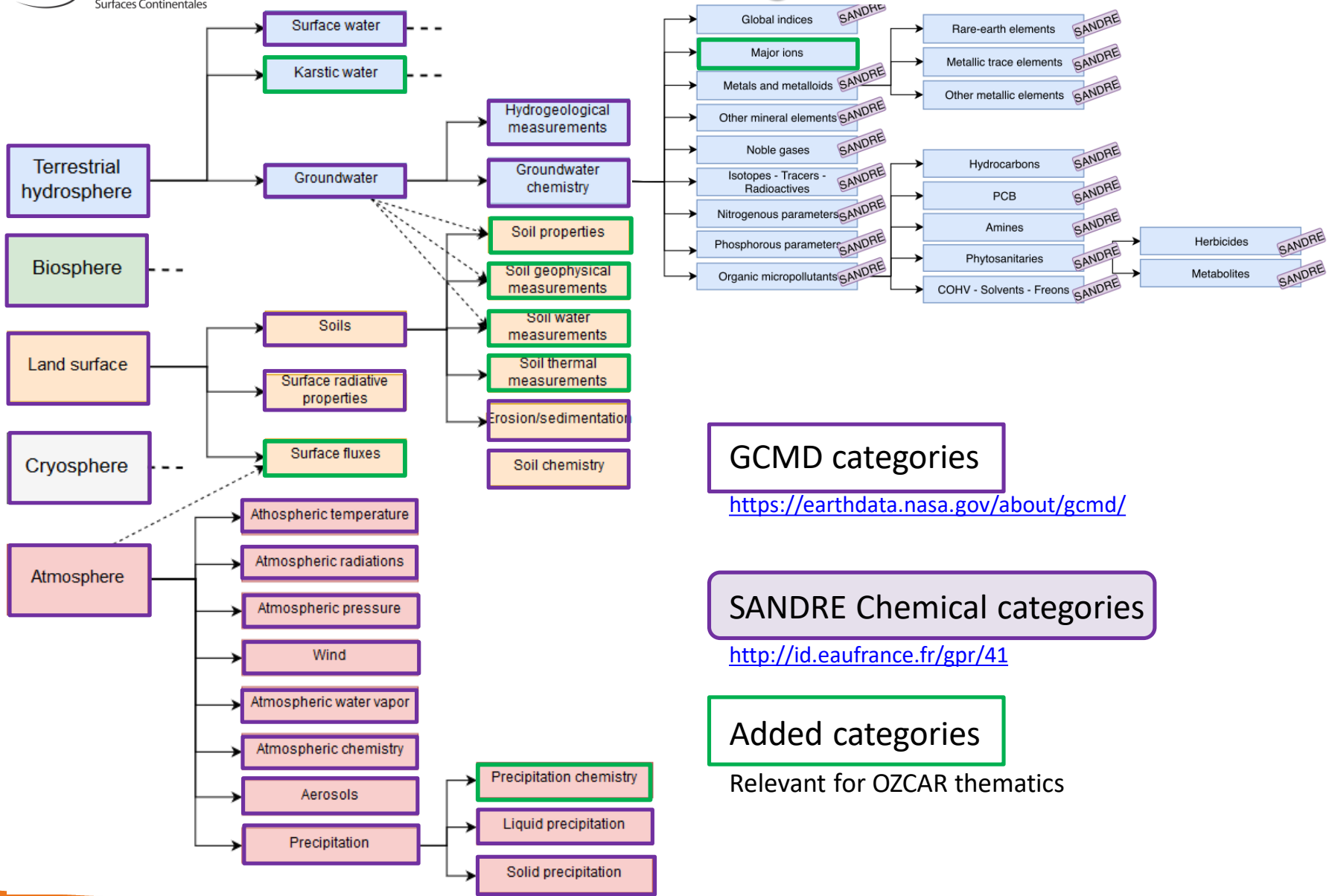
La taxonomie **GCMD** est très **inadaptée** pour **catégoriser** les **variables issues de mesures chimiques**.

Le [thesaurus du SANDRE](#) est repris pour organiser les variables issues de mesures chimiques.

SANDRE: Service National des Données et Référentiels sur l'Eau



Thésaurus final : catégories de variables



GCMD categories

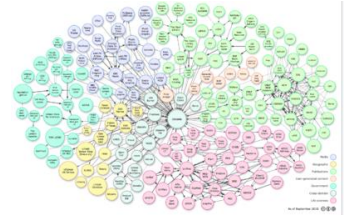
<https://earthdata.nasa.gov/about/gcmd/>

SANDRE Chemical categories

<http://id.eaufrance.fr/gpr/41>

Added categories

Relevant for OZCAR thematics



4. Publication du thésaurus en Linked Open Data

- **Alignement des concepts avec des thésaurus internationaux du domaine:** AGROVOC (FAO), GEMET (Agence européenne pour l'environnement, EnvThes (LTER), ANAEE,....).
Outil: OnaGUI
- Formalisation en **SKOS**:
 - définition de prefLabel, relations sémantiques: skos:Broader, propriétés d'alignement skos:ExactMatch, skos:closeMatch, skos:related, collection (catégories variables/noms de variables)
 - Outil: Excell + [Skos Play convertisseur RDF/SKOS](#)
- Publication sur le Web: outil [Skosmos](#) (s'appuie sur Apache Jena Fuseki SPARQL server)

http://in-situ.theia-land.fr/vocabularies/Skosmos/theia_ozcar_thesaurus/en

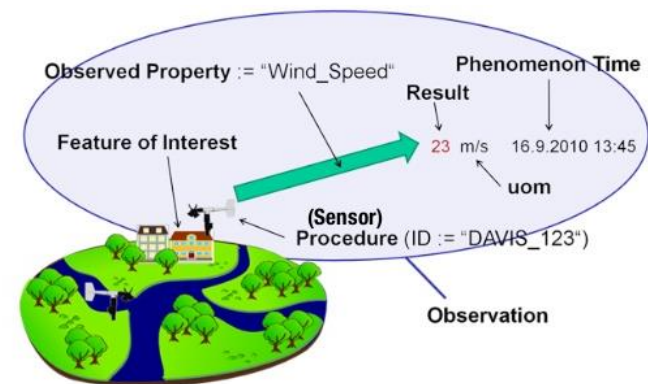
Travail sur les métadonnées: analyse des standards et identification des flux d'informations à mettre en place

Métadonnées nécessaires pour mettre en place des **services d'interopérabilité** et des **services de citation** ?

Standards analysés:

- Dublin Core
- ISO 19115/Inspire
- DataCite
- O&M et SensorML

- Data Catalog Vocabulary DCAT (schéma RDF)
- Schema.org/dataset (schéma créé Google, Yahoo, Bing, nécessaire pour être référencé dans Google Dataset Search)

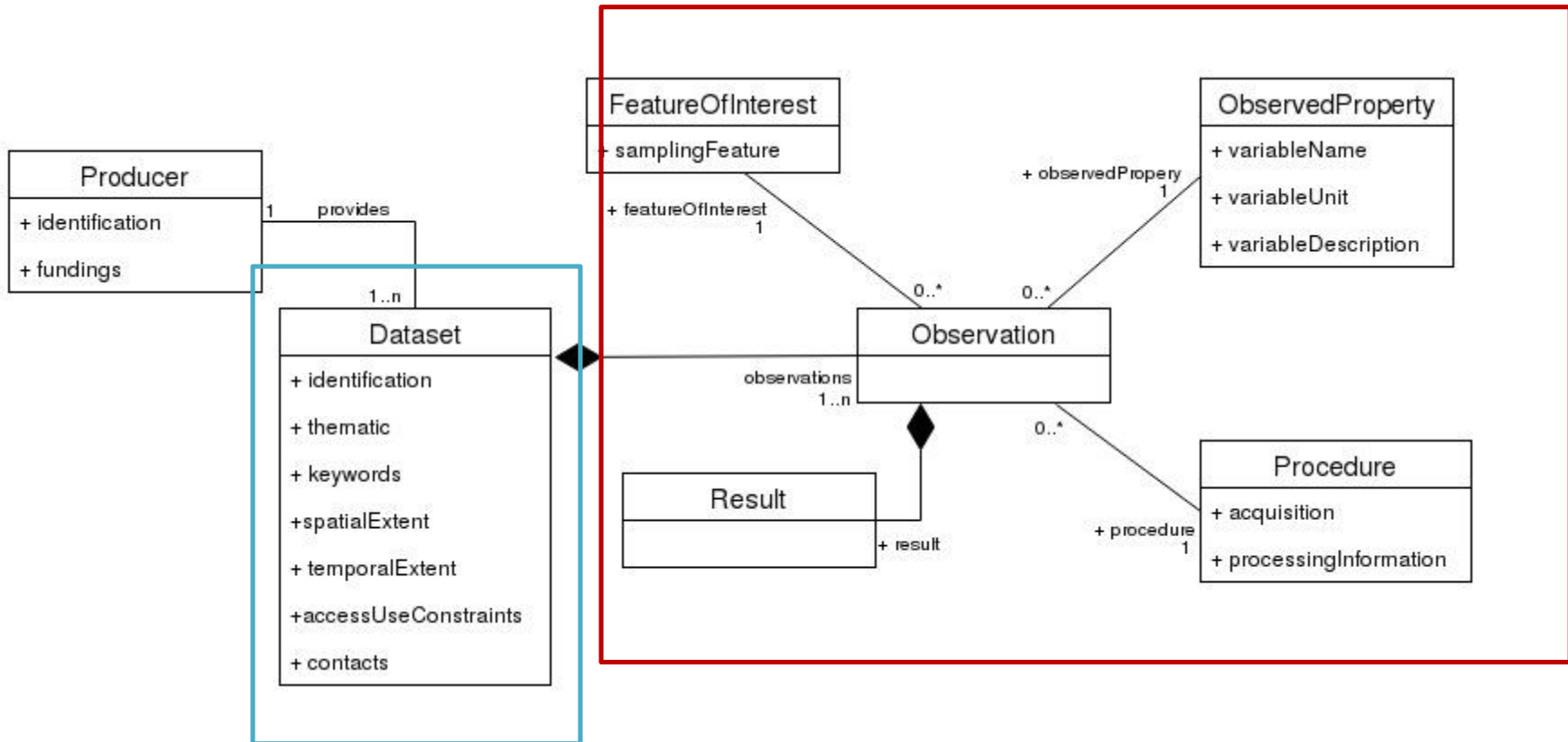


Travail sur les métadonnées: analyse des standards et identification des flux d'informations à mettre en place

- Mapping entre les **différentes normes** (pour générer les flux standards en sortie du SI (webservice OGC CSW &SOS, déclaration de doi)
- Identification d'un **set minimum** de **métadonnées** à échanger
- Choix d'un **format pivot maison** et non le format d'un standard en particulier :
 - Pour contenir toutes les informations nécessaires pour répondre à la fois aux différents standards d'interopérabilité et de citation et aux fonctionnalités du système
 - Pour offrir une liberté évolutive
 - Pour offrir un format simple (moins verbeux que XML)

Format pivot (data model)

O&M



ISO 19115 / Inspire

Modèle conceptuel

Format pivot (data model) 1/2

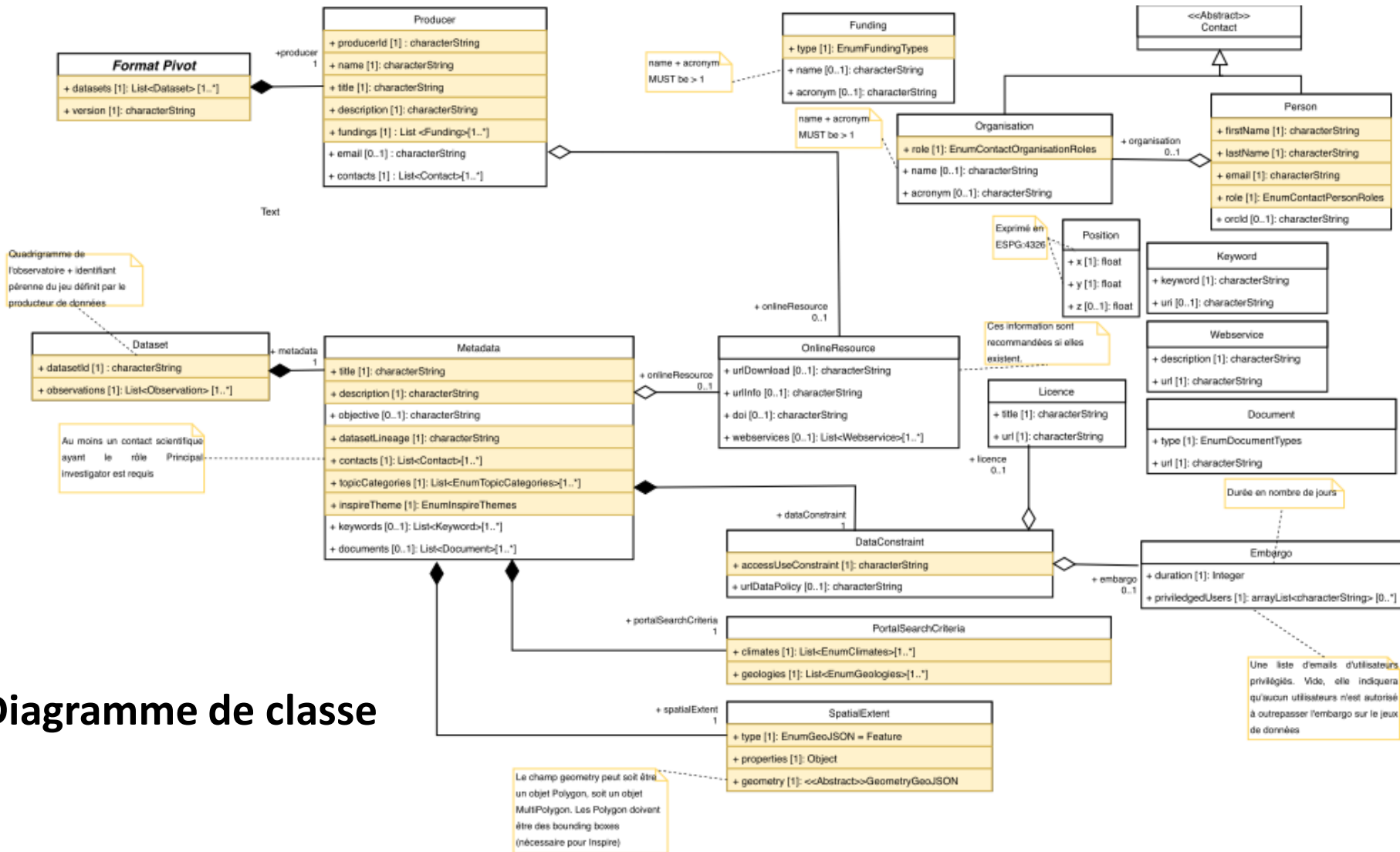
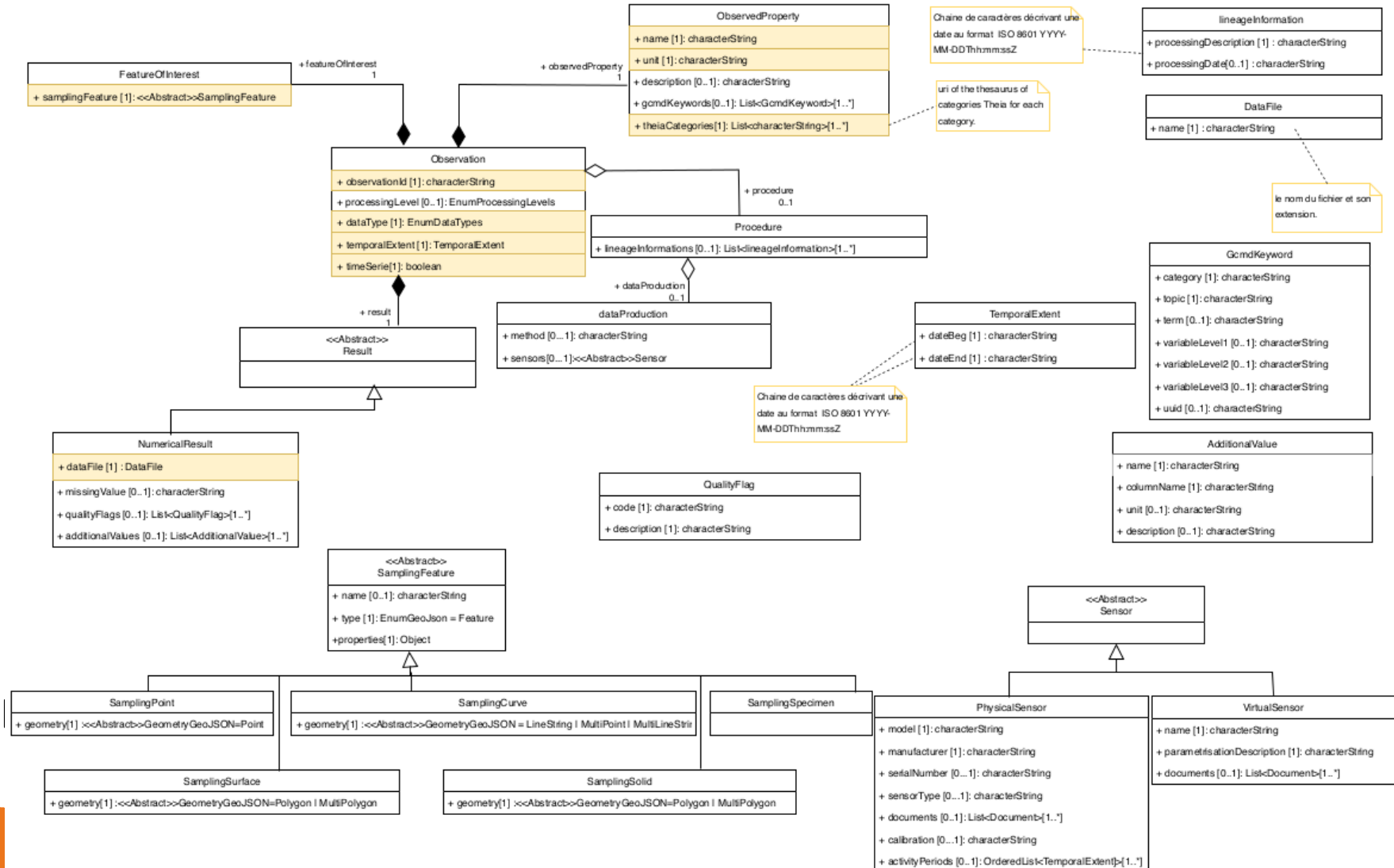


Diagramme de classe

Format pivot (data model) 2/2



Portail de données

The screenshot displays the OZCAR data portal interface. On the left, a sidebar contains several filter sections: 'Categories of variable' with expandable options for Atmosphere (1646), Biosphere (564), Cryosphere (455), Land surface (862), and Terrestrial hydrosphere (2593); 'Temporal extent' with 'From' and 'To' date pickers; 'Producer' with a list of OZCAR-RI projects and their counts, all checked; 'Full text search' with a search input field; and 'Geologies' with a list of rock types and their counts. A yellow arrow points to the 'Producer' section. The main area shows a world map with colored circles representing data points, each with a number (15, 5, 2, 38, 999, 470). A scale bar indicates 3000 km. At the top right, there is a 'Select a base layer' dropdown menu. At the bottom right, there is a 'show measurement list' button with an upward arrow icon.

Reset or Submit selection

Categories of variable

- Atmosphere (1646)
- Biosphere (564)
- Cryosphere (455)
- Land surface (862)
- Terrestrial hydrosphere (2593)

Temporal extent

From [] []

To [] []

Producer

- OZCAR-RI AMMA-CATCH (4189)
- OZCAR-RI CRYOBS-CLIM (642)
- OZCAR-RI ERORUN (25)
- OZCAR-RI MSEC (52)
- OZCAR-RI SNO KARST (267)
- OZCAR-RI SNO Tourbières (5)
- OZCAR-RI SO-HYBAM (58)

Full text search

Search...

Geologies

- Carbonate rocks (267)
- Metamorphic rocks (3856)
- Other sedimentary rocks (2769)
- Plutonic rocks (1659)

3000 km

Select a base layer

15 5 2 38 999 470

show measurement list

Recherche par facette:

- par zone géographique
- par période temporelle
- par variable
- par observatoire
- Par organisme financeur

Portail de données

Reset or Submit selection

- Categories of variable
- Atmosphere (1646)
 - Biosphere (564)
 - Cryosphere (455)
 - Land surface (862)
 - Terrestrial hydrosphere (2593)

Temporal extent

From

To

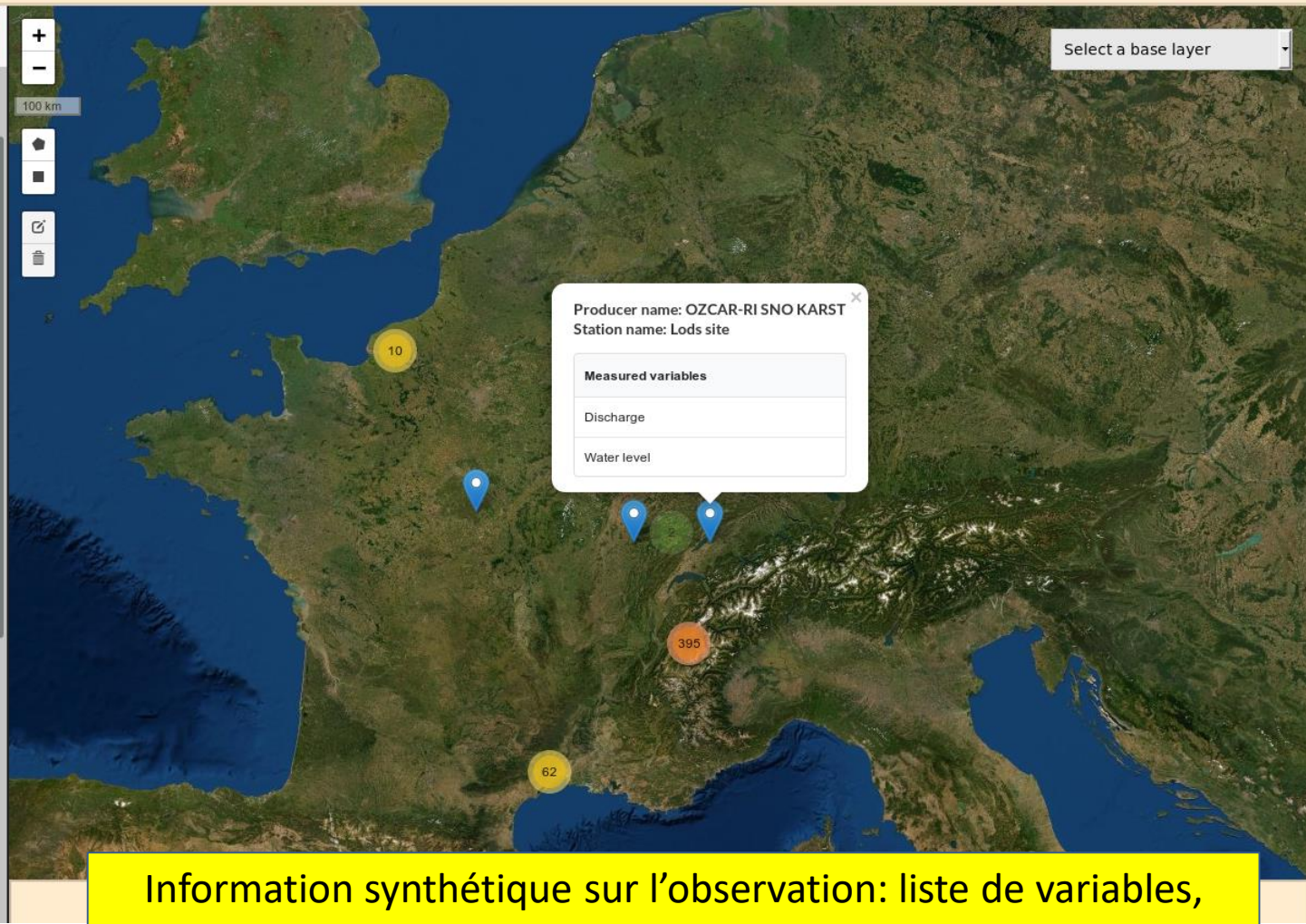
- +

- Producer
- OZCAR-RI AMMA-CATCH (4189)
 - OZCAR-RI CRYOBS-CLIM (642)
 - OZCAR-RI ERORUN (25)
 - OZCAR-RI MSEC (52)
 - OZCAR-RI SNO KARST (267)
 - OZCAR-RI SNO Tourbières (5)
 - OZCAR-RI SO-HYBAM (58)

Full text search

Search...

- Geologies
- Carbonate rocks (267)
 - Metamorphic rocks (3856)
 - Other sedimentary rocks (2769)



Information synthétique sur l'observation: liste de variables, noms de station, producteur

Portail de données

Measurement :

Variables:

Theia variable name:

Discharge

Theia categories:

Terrestrial hydrosphere > Karstic water > Karst hydrology

	Producer variable name	Unit	Temporal extent	Description	GCMD science keywords
1	Discharge	m3/s	2016-01-01T00:00:00Z 2016-12-31T23:30:00Z	Discharge Lods	EARTH SCIENCE > LAND SURFACE > GEOMORPHIC LANDFORMS/PROCESSES > KARST PROCESSES > KARST HYDROLOGY EARTH SCIENCE > TERRESTRIAL HYDROSPHERE > SURFACE WATER > SURFACE WATER PROCESSES/MEASUREMENTS > DISCHARGE/FLOW

Location:

Station name: Lods site

Longitude : 6.239739
Latitude : 47.04832
Altitude : 370

Sensors - variable 1

Model : Orpheus Mini
Manufacturer : OTT

Other variables at this location:

- Conductivity
- Turbidity
- Water temperature
- Nitrate
- Fluorescence
- Water level
- Organic carbon
- Total organic carbon (TOC)

Included in dataset :

Information détaillée sur l'observation: liste de variables, noms de station, producteur, jeux de données

Prochaines étapes

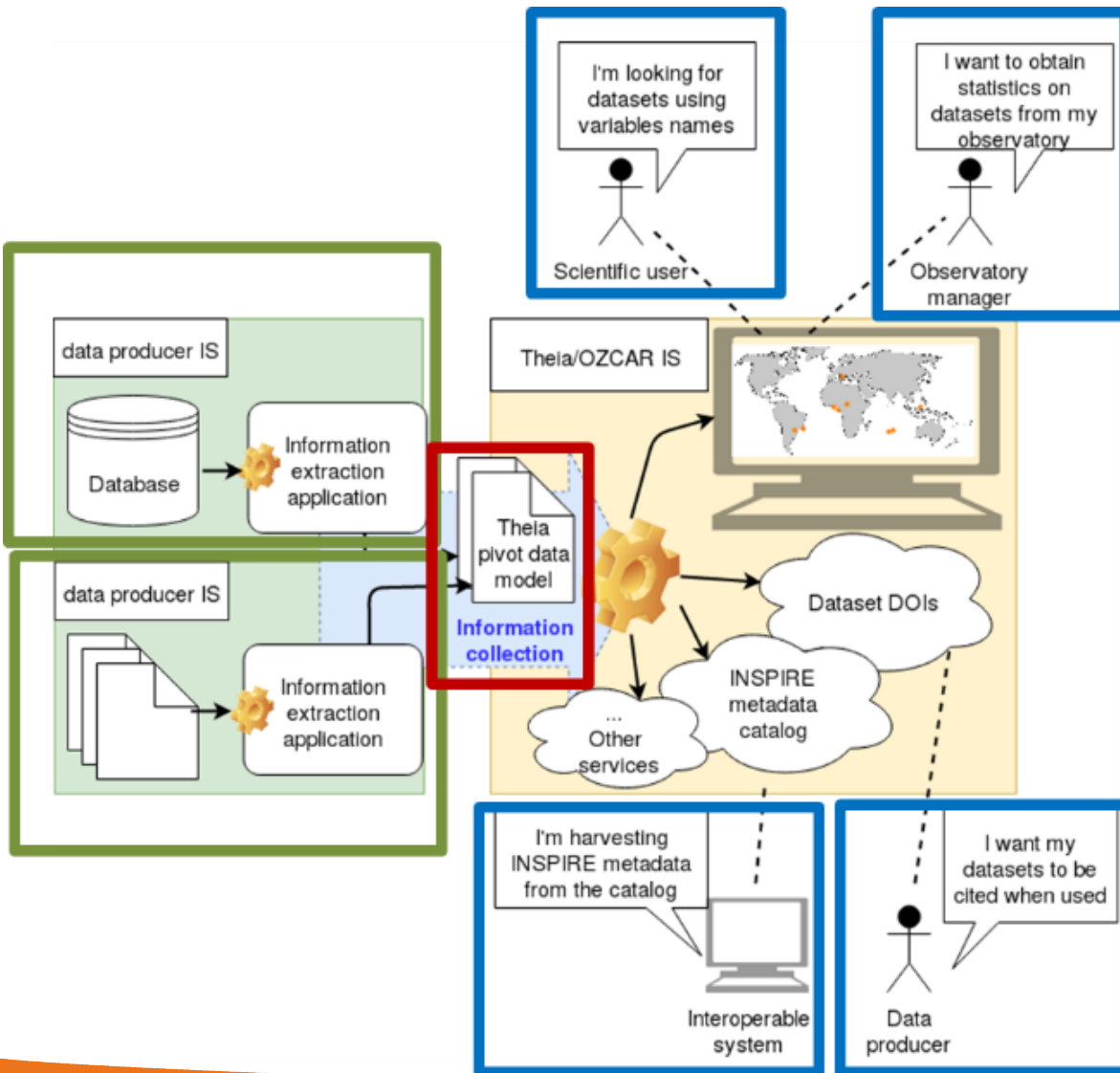
- Finaliser l'interface du portail
- Finaliser l'architecture de production (déploiement continue docker Kubernetes) et mettre en production (déc. 2019)

2020-2022

- Mettre en place des services d'interopérabilité (web service CSW geonetwork, web service SOS)
- Fournir les données dans différents formats NetCDF, csv
- Continuer à interfacier les SI des producteurs
- Vocabulaire : objet d'intérêt (FeatureOfInterest) sampled Feature

Des questions ?

Slides Annexe



Architecture

Construction d'un flux de données :

- 1) Différents producteurs, différents formats
- 2) Le **format pivot** permet :
 - (i) de collecter ces informations,
 - (ii) leur mise à jour en temps réel
- 3) Le SI Theia/OZCAR répond aux requêtes d'utilisateurs humains ou machine

Architecture du système et développement d'un prototype

