

Introduction

1. Les Objectifs

- Sur un même jeu de données:
 - Étudier différents **formats** de données
 - NetCDF / Zarr / Kerchunk
 - Étudier différents **protocoles** d'accès à la donnée
 - POSIX / S3 / FTP / HTTP / OPeNDAP
 - Étudier différentes opérations
 - Chargement / Lecture / Subsetting
- Gitlab: https://gitlab.ifremer.fr/odatis/benchmarks/data_access

2. Le Dataset

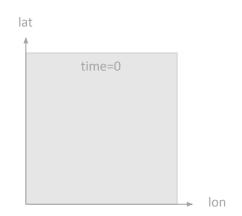
OSISAF: Global Metop Sea Surface Temperature (fichiers de 12h) (time, longitude, latitude).



1. NetCDF (Network Common Data Form)

NetCDF est un format de fichier largement utilisé pour stocker des données scientifiques **multidimensionnelles**, il est auto-descriptif contenant les données mais aussi les métadonnées.



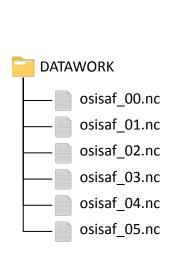


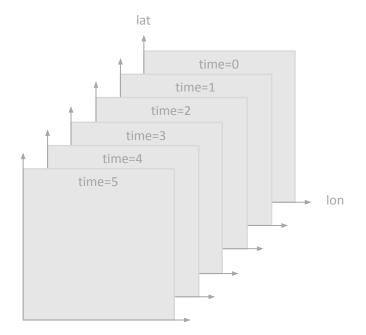




1. NetCDF (Network Common Data Form)

NetCDF est un format de fichier largement utilisé pour stocker des données scientifiques **multidimensionnelles**, il est auto-descriptif contenant les données mais aussi les métadonnées.







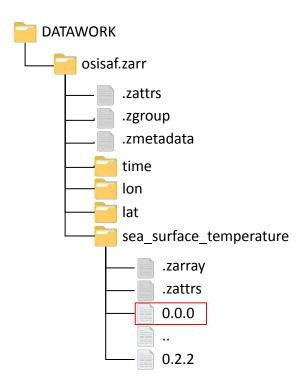
2. Zarr

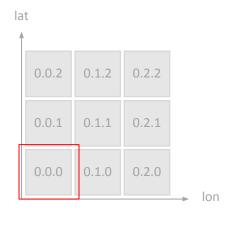
Zarr est un format moderne et flexible pour le stockage de données scientifiques **multidimensionnelles**. Contrairement à NetCDF, qui repose sur des fichiers binaires monolithiques, Zarr est conçu pour un accès distribué et parallèle aux données.





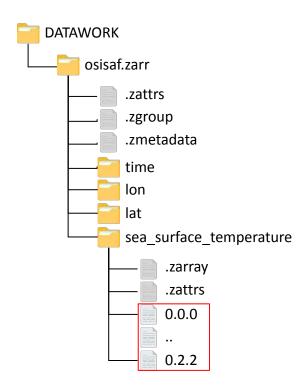
2. Zarr

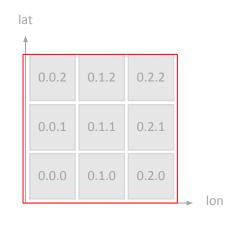






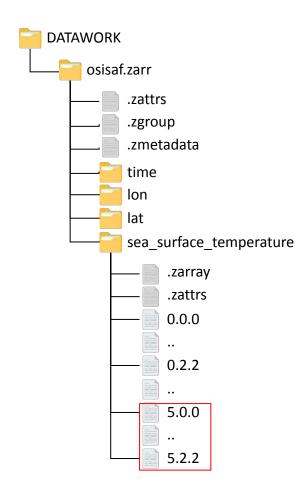
2. Zarr

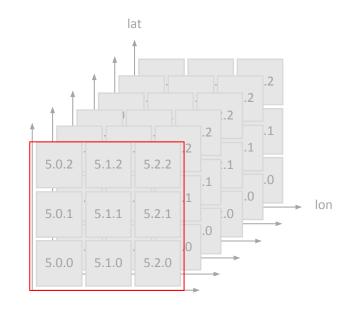






2. Zarr







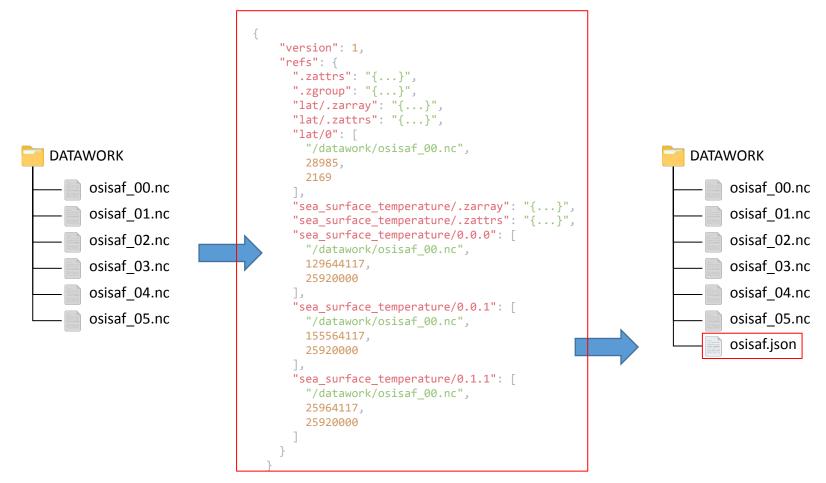
3. Kerchunk

Kerchunk est une bibliothèque conçue pour créer une émulation Zarr à partir de fichiers NetCDF.

Il fonctionne en créant un index léger qui pointe vers les données réelles, sans avoir à les dupliquer.



3. Kerchunk





Les protocoles

Voici les protocoles étudiés:

- POSIX (local)
- S3
- FTP
- HTTP
- OPeNDAP

Les opérations

Voici les opérations effectuées:

- Chargement (Lecture des métadonnées)
- Lecture (Lecture des données -> calcul de la moyenne)
- Subsetting (Lecture d'un sous ensemble de la donnée -> 10%)



Benchmark (Résultats)

1. NetCDF - Les différents protocoles

Configuration				
Cache	False			
Nombre de coeur	5			
Chunks	{'time':1, 'lon': 3600, 'lat': 1800}			
Slices	{'lon': slice(0,7200), 'lat':slice(0, 360)}			
Protocole	POSIX			

netcdf	al de la companya de				format
s3	opendap	http	ftp	posix	filesystem
1.50	0.75	1.17	7.09	1.00	load (9.4s)
27.46	11.16	7.96	3.48	1.00	read (3.6s)
27.76	17.49	22.69	2.88	1.00	subset (0.9s)

Valeurs normalisées (par NetCDF en local)

Benchmark (Résultats)

2. Local VS Distant

Configuration				
Cache	False			
Nombre de coeur	5			
Chunks	{'time':1, 'lon': 3600, 'lat': 1800}			
Slices	{'lon': slice(0,7200), 'lat':slice(0, 360)}			

s3			posix			filesystem
rr kerchunk	zarr	netcdf	kerchunk	zarr	netcdf	format
72 0.65	0.72	1.50	0.69	0.73	1.00	load (9.4s)
06 1.48	3.06	27.46	1.13	0.99	1.00	read (3.6s)
79 3.72	6.79	27.76	1.51	1.56	1.00	subset (0.9s)

Valeurs normalisées (par NetCDF en local)

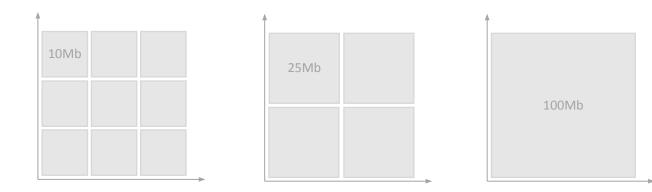


- Taille du dataset (30 fichiers de 200Mb = 6gb)
- Avec ou sans cache (local / serveur)
- Différentes méthodes pour accéder à la données (fsspec / xarray)
- Taille des chunks (de 10 à 100Mb)





- Taille du dataset (30 fichiers de 200Mb = 6gb)
- Avec ou sans cache (local / serveur)
- Différentes méthodes pour accéder à la données (fsspec / xarray)
- Taille des chunks (de 10 à 100Mb)



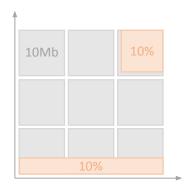


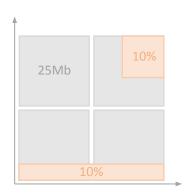


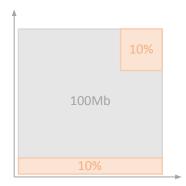
- Taille du dataset (30 fichiers de 200Mb = 6gb)
- Avec ou sans cache (local / serveur)
- Différentes méthodes pour accéder à la données (fsspec / xarray)
- Taille des chunks (de 10 à 100Mb)
- Ressources (CPU / thread) (2 et 5)
- Position du subset (10% 1 seul chunk VS plusieurs chunk)



- Taille du dataset (30 fichiers de 200Mb = 6gb)
- Avec ou sans cache (local / serveur)
- Différentes méthodes pour accéder à la données (fsspec / xarray)
- Taille des chunks (de 10 à 100Mb)
- Ressources (CPU / thread) (2 et 5)
- Position du subset (10% 1 seul chunk VS plusieurs chunk)









1. De nombreux paramètres

- Taille du dataset (30 fichiers de 200Mb = 6gb)
- Avec ou sans cache (local / serveur)
- Différentes méthodes pour accéder à la données (fsspec / xarray)
- Taille des chunks (de 10 à 100Mb)
- Ressources (CPU / thread) (2 et 5)
- Position du subset (10% 1 seul chunk VS plusieurs chunk)
- Nombre d'itérations (entre 1 à 10)
- ...

2. De nombreux obstacles

- Processus très itératifs (de part le nombre de paramètres)
- Erreurs / Bug / Problèmes réseaux



Benchmark (Conclusion)

1. Les protocoles

Protocoles	Avantages	Inconvénients
POSIX	Rapidité (+++), possibilité de chargement et subsetting	
S3	Rapidité (++), possibilité de chargement et subsetting, tout format de fichier	Difficultés de prise en main
FTP	Dans certaine condition préférable de télécharger une fois, pour travailler en local.	Pas de chargement ou de subsetting, nécessite de télécharger TOUTE la donnée en locale.
HTTP	Rapidité (+), possibilité de chargement et subsetting (HTTP-RANGE?)	
OPeNDAP	Rapidité (+), possibilité de chargement et subsetting	

Benchmark (Conclusion)

2. Les formats

Protocoles	Avantages	Inconvénients
NetCDF	Chargement, subsetting	Moins optimisé pour le cloud
Zarr	Lecture/écriture en parallèle, chargement, subsetting	Contient de nombreux petit fichiers
Kerchunk	Avoir les avantages du zarr sur des fichiers NetCDF. Très utile pour les collections historiques au format NetCDF	Le fichier de référence JSON peut devenir très lourd s'il y a de nombreux chunk (possibilité d'utiliser un fichier parquet)

















