

CONFÉRENCE

TIC et mer: nouveaux défis et solutions

Les technologies de l'information au service de la recherche marine



Gérer des bases de données de plus en plus grandes et complexes

Partage et interactions de bases de données

Méthodes de fouille et d'analyse

9H45: ACCUEIL

10H00-16H00: PRÉSENTATIONS

MATHIAS HERBERTS, JEAN-FRANÇOIS PIOLLÉ,
STÉPHANIE MAHÉVAS, GUILLAUME MAZE,
THOMAS LOUBRIEU, GILBERT MAUDIRE,
PHILIPPE LENCA, RONAN FABLET

16H00: TABLE RONDE AVEC RENÉ GARELLO

Thomas Loubrieu

(Ifremer)

“Small to Big Data”

26 Novembre 2013, Ifremer, Brest



<http://wwz.ifremer.fr/bigdata>

Small to Big data

IFREMER/IDM/ISI

T. Loubrieu

Résumé

A partir d'expériences en gestion de données du département Informatique et Données Marines d'IFREMER, on passera en revue différentes problématiques liées à la volumétrie des données. Sur ces problématiques, on passera en revue les solutions possibles d'une part avec les technologies traditionnelles (formats, indexation, interoperabilité et partage de données) et d'autre part avec les systèmes dématérialisés. Par cette revue, on tentera d'explicitier la complémentarité des solutions informatiques traditionnelles et big data.

Plan:

- Solutions “performance” dans l'informatique traditionnelle
- Vers le big data

Besoin traité ici

Traitement scientifique de données (vs gestion de centre de données):

- plutôt lecture
- plutôt exploratoire (ni opérationnel, ni répétitif)
- pas de mission d'archivage long terme

Problèmes rencontrés

- problème de **capacité de stockage**
- problème de **rapidité d'accès** aux données stockées (e.g sur les disques)
- problème de **capacité de traitement**
- problème de **rapidité de traitement**

En transversal:

- problème de **temps de développement** des solutions.

Solutions traditionnelles (1/5)

Volume de stockage:

- encodage (xml/csv/binaire)
- compression (**attention** à la rapidité d'accès)
- éviter la redondance d'information
- mutualiser les copies entre les utilisateurs (**attention** stockage des petites différences) dans un centre de données:
 - Confiez vos produits à un centre de données.
 - Faites confiance aux centres de données pour gérer vos entrées (inputs).

Solutions traditionnelles (2/5)

Rapidité d'accès sur disque:

- indexation: via encodage (ascii vs netcdf)
- indexation via SGBD (mysql, oracle, ...)

Solutions traditionnelles (3/5)

Capacité de traitement (RAM):

- ne lire que ce qui est strictement utile au traitement
- ne pas conserver en mémoire les copies d'objets
- traiter la donnée par morceaux (attention rapidité)

Solutions traditionnelles (4/5)

Rapidité de traitement (lecture disque, CPU):

- choix du langage (matlab -> fortran)
- pré-traiter la donnée (la modéliser et la stocker dans le modèle approprié au traitement).
- systèmes parallèles High Performance Computing (caparmor):
 - possibilité de traitements parallèles inter-dépendants (MPI).
 - Cluster de 300 machines sur réseau infiniband
 - Stockage sur disques parallèles lustre.

Solutions traditionnelles (5/5)

Limiter les temps de développement:

- homogénéiser les formats de données
- ré-utiliser des bibliothèques de fonctions, pour lire les données (e.g. python/pyNIO) ou traitements (e.g. statistiques R)

Exemples de solutions (1/3)

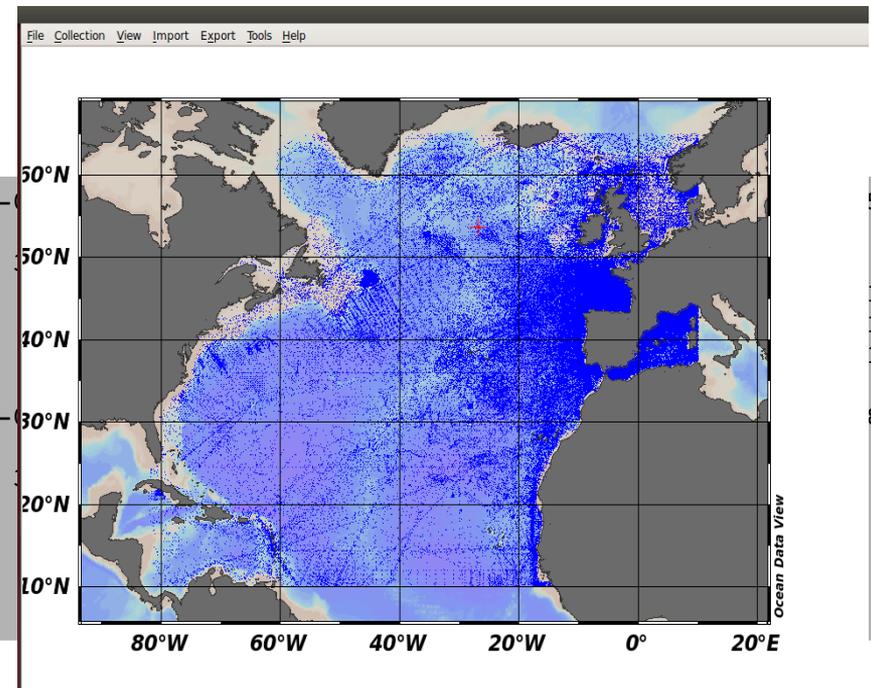
- Encodage (ASCII vs binaire): format ODV ascii → ODV collection en binaire (700 000)

```
0401C 1 652/1 1 C 1 2001-03-07T15:50:00.000 1 358.00931 1 56.24580
43 1 114 1 38 1 1033356 1 652/1 1 PSAL,RFVL,TEMP 1 2011-04-05 15:
07T15:50:00.000 1 0 5.6170 1 34.51900 1 0
2001-03-07T16:10:00.000 1 0 5.6170 1 34.51900 1
2001-03-07T16:30:00.000 1 0 5.6170 1 34.51900 1
2001-03-07T16:50:00.000 1 0 5.6170 1 34.51900 1
(... ?)
0401C 1 654/2 1 C 1 2001-03-08T15:50:00.000 1 358.74701 1 56.24850
43 1 114 1 38 1 1033358 1 654/2 1 PSAL,RFVL,TEMP 1 2011-04-05 15:
08T15:50:00.000 1 0 6.0030 1 34.88200 1 0
2001-03-08T16:10:00.000 1 0 5.9800 1 34.90600 1
2001-03-08T16:30:00.000 1 0 5.9800 1 34.90600 1
2001-03-08T16:50:00.000 1 0 5.9800 1 34.82400 1
```

Exemples de solutions (1/3)

- Encodage (ASCII vs binaire): format ODV ascii

```
0401C 1 652/1 1 C 1 2001-03-07T15:50:00.000 1
43 1 114 1 38 1 1033356 1
07T15:50:00.000 1 0 5.6170 1
2001-03-07T16:10:00.000 1
2001-03-07T16:30:00.000 1
2001-03-07T16:50:00.000 1
(... ?)
0401C 1 654/2 1 C 1 2001-03-08T15:50:00.000 1
43 1 114 1 38 1 1033358 1
08T15:50:00.000 1 0 6.0030 1
2001-03-08T16:10:00.000 1
2001-03-08T16:30:00.000 1
2001-03-08T16:50:00.000 1
```

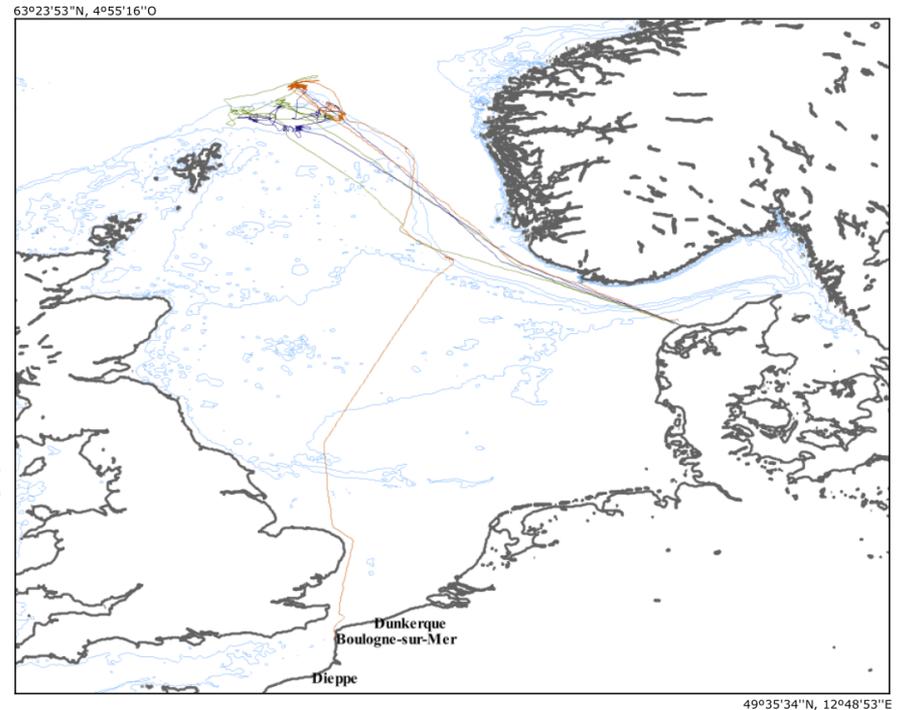


Exemples de solutions (2/3)

Position des navires de pêche, VMS : 16 M de positions par année depuis > 10 ans:

- index paginé par année
- index spatial

On sait utiliser, afficher rapidement les positions d'une zone ou d'un navire sur une carte.



Exemples de solutions (3/3)

Modèle de données (Harmonie, Q2), ms access → oracle, traitement SACROIS:

- 4 sources de données hétérogènes (marées déclarées, marées reconstituées par les positions navires, ventes, calendrier d'activité)
- On produit mensuellement des synthèses :
 - ♦ 33 000 marées,
 - ♦ 1M de captures

Vers le big data ? (1/2)

Dans les cas suivants:

- limite vraie de capacité matérielle
- partage des données au delà du cadre local (thématique et géographique)

Cadre “architecture” INSPIRE

<http://inspire.jrc.ec.europa.eu/index.cfm/pageid/48>

- INSPIRE is based on a number of common principles:
- 1)Data should be collected only once and kept where it can be maintained most effectively.
- 2)...

Cela signifie que:

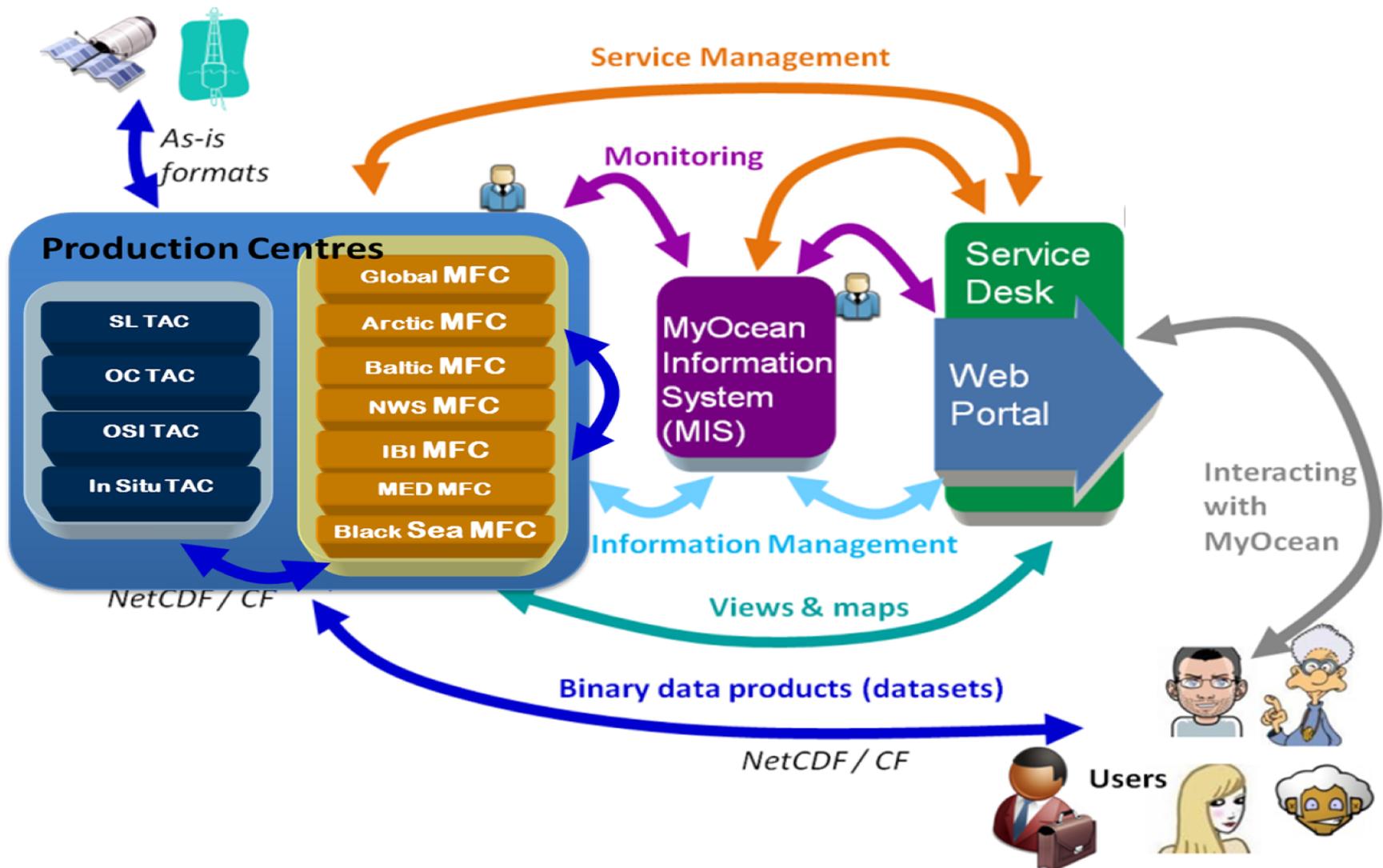
La copie maîtresse des données doit être proche de la "compétence" qui a créé l'information et là où les infrastructures permettent de la conserver et la diffuser sur le long terme.

Propositions techniques (dans ce cadre)

Système de documentation/diffusion de données distribués dans les centres de données marines:

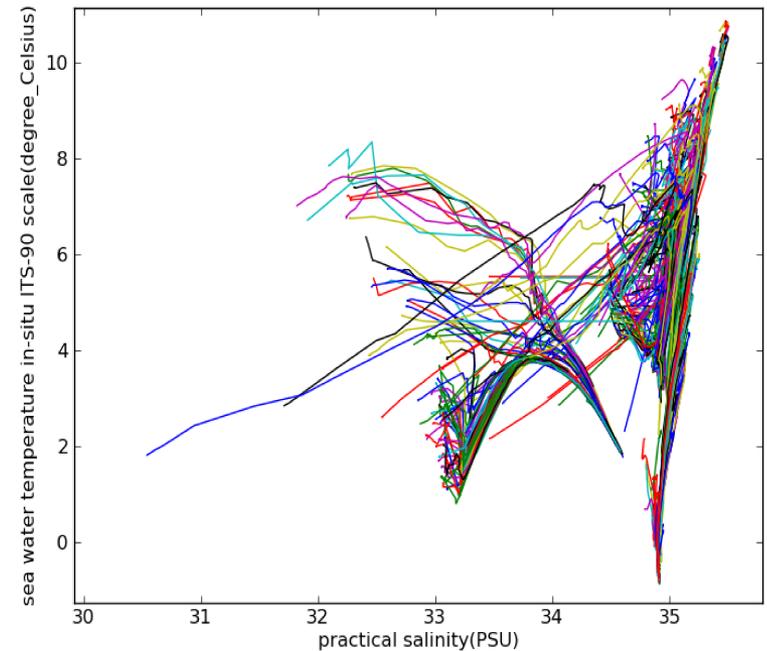
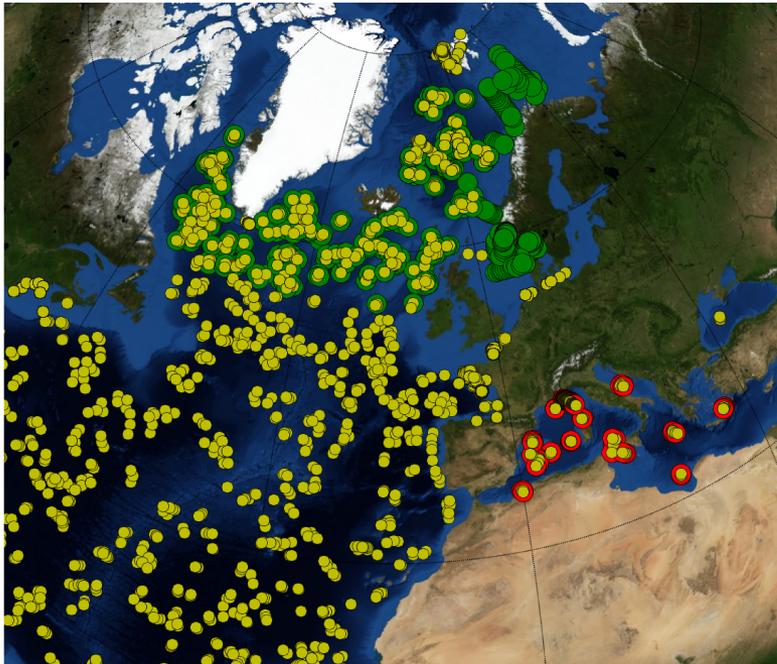
- les ressources 'données' sont documentées et référencées via des **catalogues Nationaux** (Sextant), **Européens** (myocean, seadatanet) ou **internationaux** (geoss, ocean data portal). Les ressources de données sont **normalisées** (e.g. Open GIS Consortium).
- la **donnée est modélisée de façon homogène** partout (e.g. oceanotron).
- l'utilisateur peut accéder (visualisation/telechargement) à des sous-ensembles de données (**filtre phenomenon + 4D**), (e.g. THREDDS Data Server).
- intégration directe de ces données dans les **algorithmes de traitement**: python, matlab, idl, ... (e.g. OPenDAP)

Exemple d'architecture myOcean



Oceanotron

Accès aux observations in-situ océanographie physique distribuée dans un modèle commun, subset-able, script-able.



THREDDS Data Server

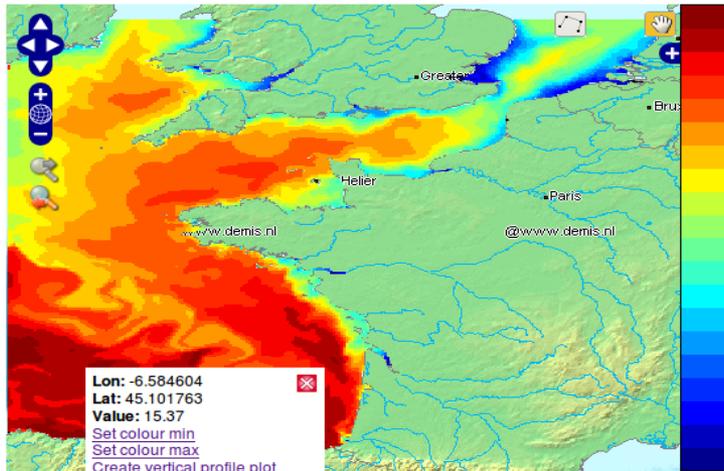
Accès aux grilles structurées de données 4D (analyses, modèles, ici 300Go par an):

Layer: Ifremer Thredds Data Server > PREVIMER F1 MARS3D
MANGA4000 FORECAST > sea_water_temperature
Units: degree_Celsius
Depth (level): 0.01666666753590107
Date/time: 25 Nov 2013 00:00:00 UTC [first frame](#) [last frame](#)

November, 2013						
Today						
Sun	Mon	Tue	Wed	Thu	Fri	Sat
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30

Select date

[Fit layer to window](#)

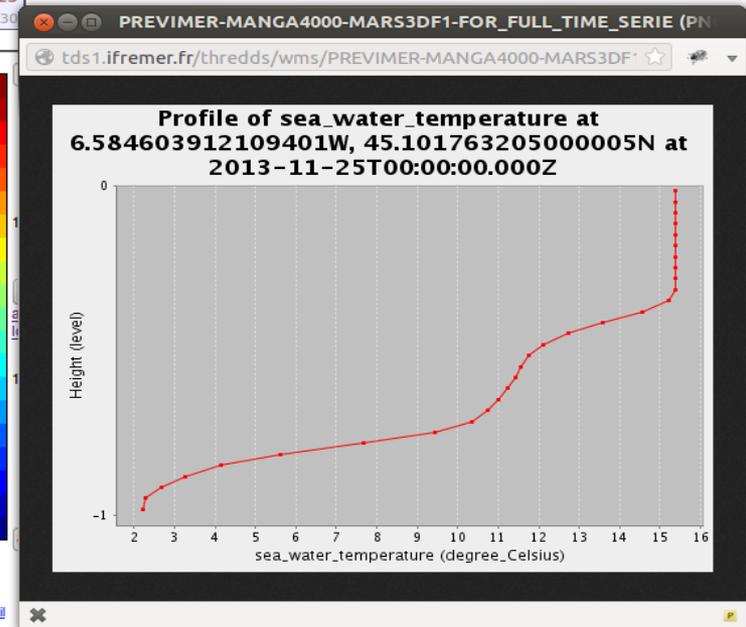


[test image](#)

Overlay opacity: 100%

Powered by [OpenLayers](#) and [OGC](#) standards

[Permalink](#) | [email](#)



Limites actuelles de ces solutions

Accès aux données par web services:

- **accès internet** (traitement éloigné de la donnée)
- **suivi de modifications** de données (reproductibilité des résultats).
- atteindre les **utilisateurs ciblés** (communication, formation)
- **maturité**

Pourquoi le big data ?

- limite vraie de capacité matérielle
- partage des données au delà du cadre local (thématique et géographique)
- lorsque les ressources informatiques d'infrastructure (liés au matériels:disque/ram/cpu) se développent très vite grâce à la virtualisation
- et relâchent les exigences sur la conception/optimisation logicielle (en amont des flux de données)

Perspectives big data

Pour un centre de données:

- **cloudifier** les applications de partage de données de l'ère "INSPIRE" pour dépasser leur limite de performance.
- mettre à disposition vos données dans des environnements **big data dédiés au domaine marin** (et alors gérer leur documentation, leur mise à jour):
 - pour la science traditionnelle (e.g. à IFREMER: données météo-france, mercator-océan, ...)
 - pour la gestion des données brutes des capteurs
- vous permettre de pousser vos données dans des environnements de **big data partagés** (pour la "data science")

Challenge

Aujourd'hui on sait que

1. certains centres de données marines n'ont pas les compétences logicielles ou d'infrastructures suffisantes.
2. des nouvelles infrastructures sont disponibles (sans compétence données marine):

Possibilité de **redistribution des rôles**, de l'architecture pour la gestion des données.