

PIA3: orchestration des flux tout au long de la vie des données

<https://frama.link/ifb-ag20-mudis4ls>

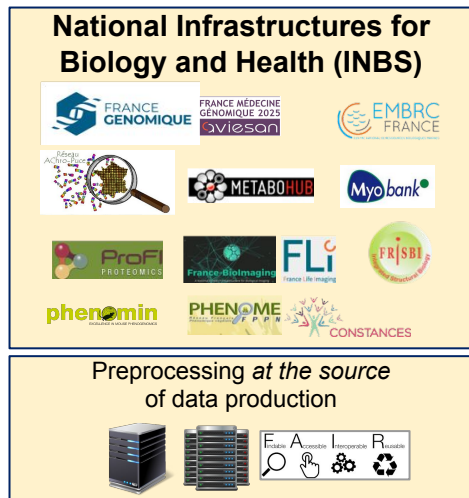
Noms



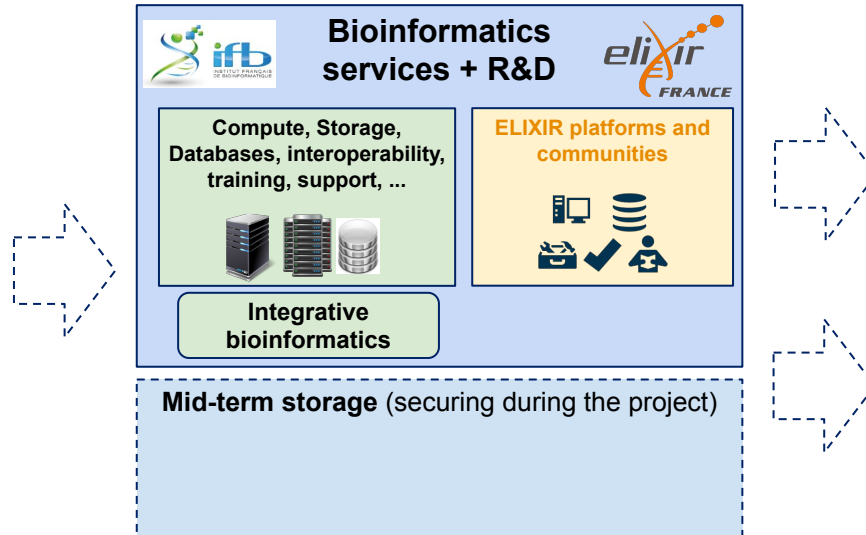
Orchestration des flux de données

Life cycle of the data – The typical stages

1. Production (weeks - months)



2. Analysis (typically 6 months - 4 years)



3. Conservation (5, 10, 20 years)



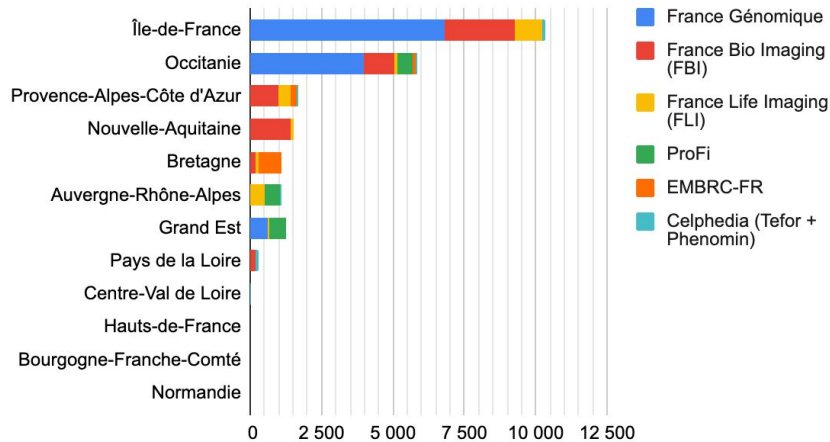
Cartographie des besoins de stockage des INBS

Nous avons demandé à chaque infrastructure nationale de biologie-santé **INBS** d'estimer la répartition par région

- Des besoins de stockage internes (numérique accolé aux machines de production ou autre raison)
- Délégués à l'IFB (notamment pour intégration avec autres types de données)

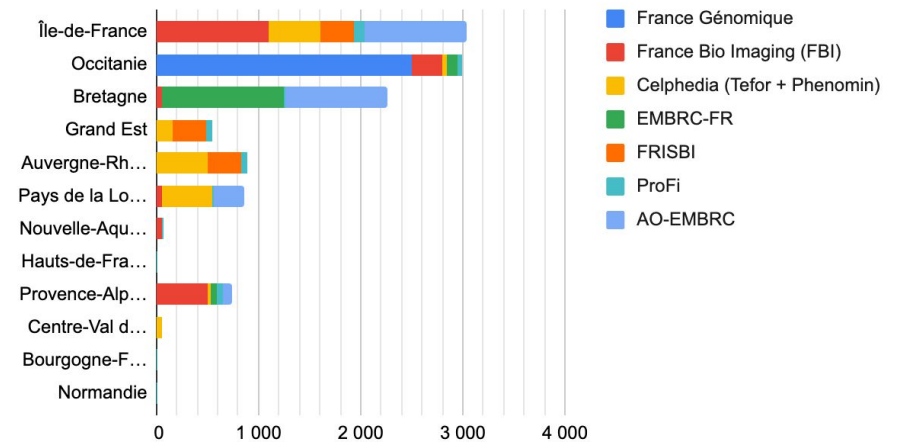
Besoins de stockage internes

Internal needs in mass storage for data-producing National Infrastructures in Life Sciences and Health



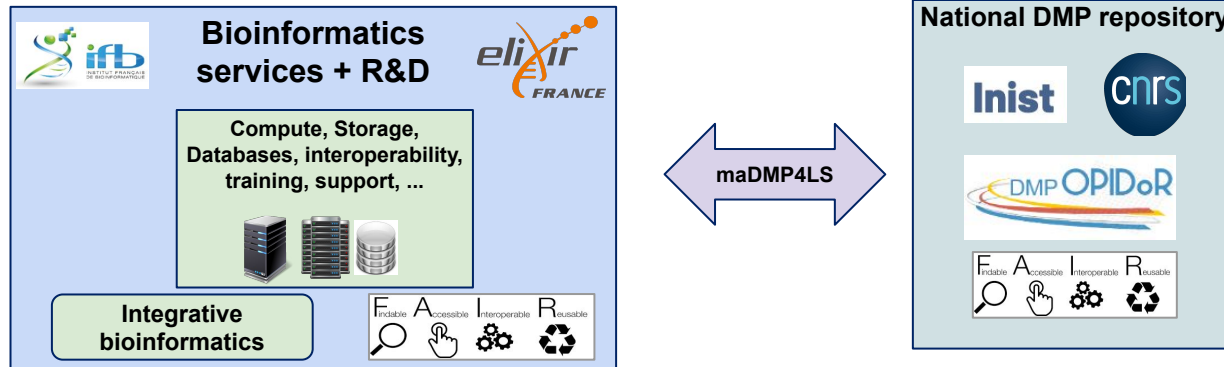
Besoins de stockage délégués à l'IFB

Mass storage delegated to IFB by other Research Infrastructures for life science and Health



Orchestration du flux de la donnée tout au long du cycle de vie



- Projet-flash ANR **machine-actionable DMP for Life Sciences (maDMP4LS)**,
 - allocation des ressources (espace-projet, volume, ressource calcul, accès collaborateurs) par accès programmatique aux DMP déposés sur OPIDoR (INIST)
 - mise à jour conjointe du DMP et des ressources au fil de l'évolution du projet
- **Gestion à chaque phase** des données et des métadonnées : production (collaboration avec autres infrastructures), analyse (serveurs IFB), conservation (centres de stockage), accès (dépôts nationaux et internationaux).
- Programmation des flux de données via le maDMP (extension du projet maDMP4LS) : des sites de production aux serveurs d'analyse, puis de ceux-ci vers les centres de stockage et vers les dépôts.



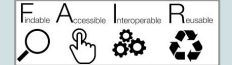
De la production à l'analyse des données

- Les infra productrices de données nécessitent du calcul et stockage local
- Offre IFB: environnements spécialisés pour l'analyse + intégration + interopérabilité des données
- Description, dès la conception du DMP, des modalités de transfert des plateformes de production/preprocessing vers les plateformes IFB pour l'analyse
- Automatisation de ces transferts (maDMP)


National DMP repository


F Indisible A Accessible I Interoperable R Reusable




Data-producing infrastructures



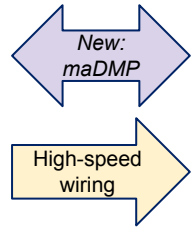
Preprocessing at the source of data production




Hosting

Hosting


Local storage + data centers



Bioinformatics services + R&D




Compute, Storage, Databases, interoperability, training, support, ...

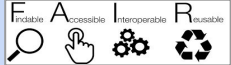


Integrative bioinformatics

ELIXIR platforms and communities




F Indisible A Accessible I Interoperable R Reusable



Hosting

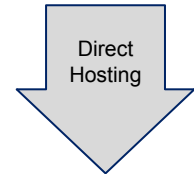
Hosting via Mesocenters (some platforms)



Hosting

Direct Hosting

Data centers



Conservation et accessibilité des données

■ Conservation des données:

- ❑ ni dans nos missions, ni dans nos projets
- ❑ collaborer avec les data centres régionaux et/ou nationaux
- ❑ Rôle de l'IFB: FAIRifier les données conservées
 - Formats standards
 - Annotation des données (métadonnées)

■ Dépôts

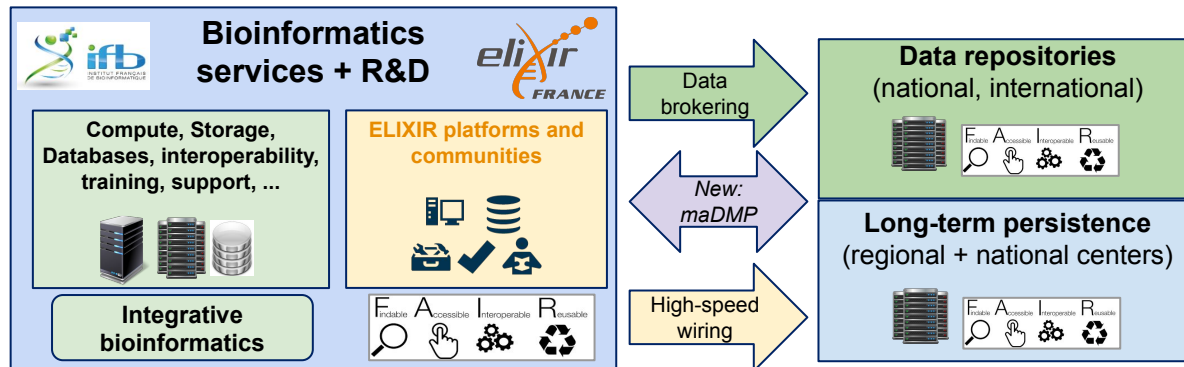
- ❑ Carence actuelle de solution au niveau national
- ❑ Existence de quelques solutions institutionnelles (INRAE, CIRAD, IRD, CEA)
- ❑ Dépôts internationaux spécifiques d'un type de données (séquences NGS, images, ...)

■ Data brokering

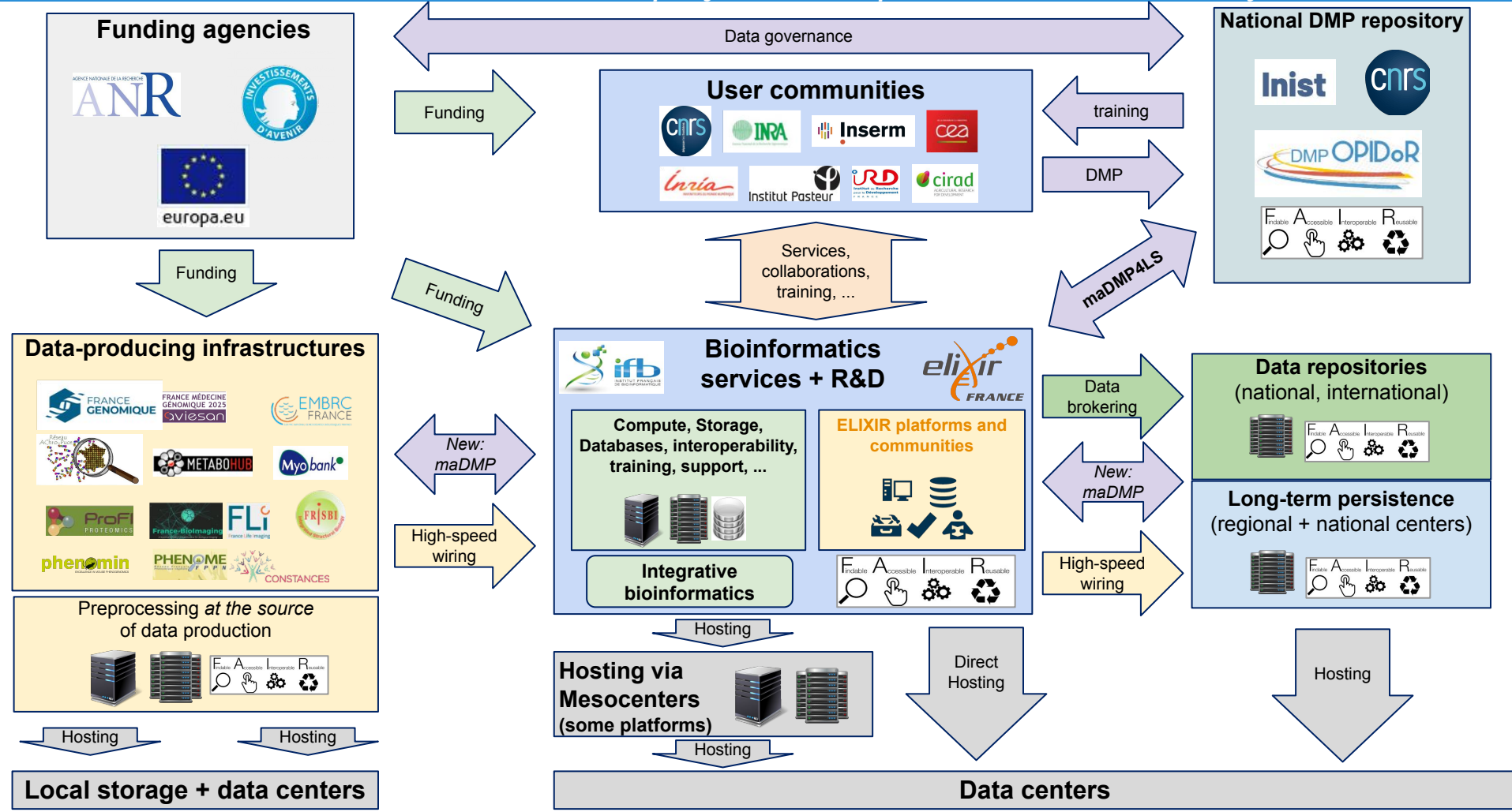
- ❑ projet-pilote avec European Nucleotide Archive (EBI) pour que l'IFB devienne courtier des données françaises. Extension envisagée à d'autres dépôts.

■ Elargissement et défis

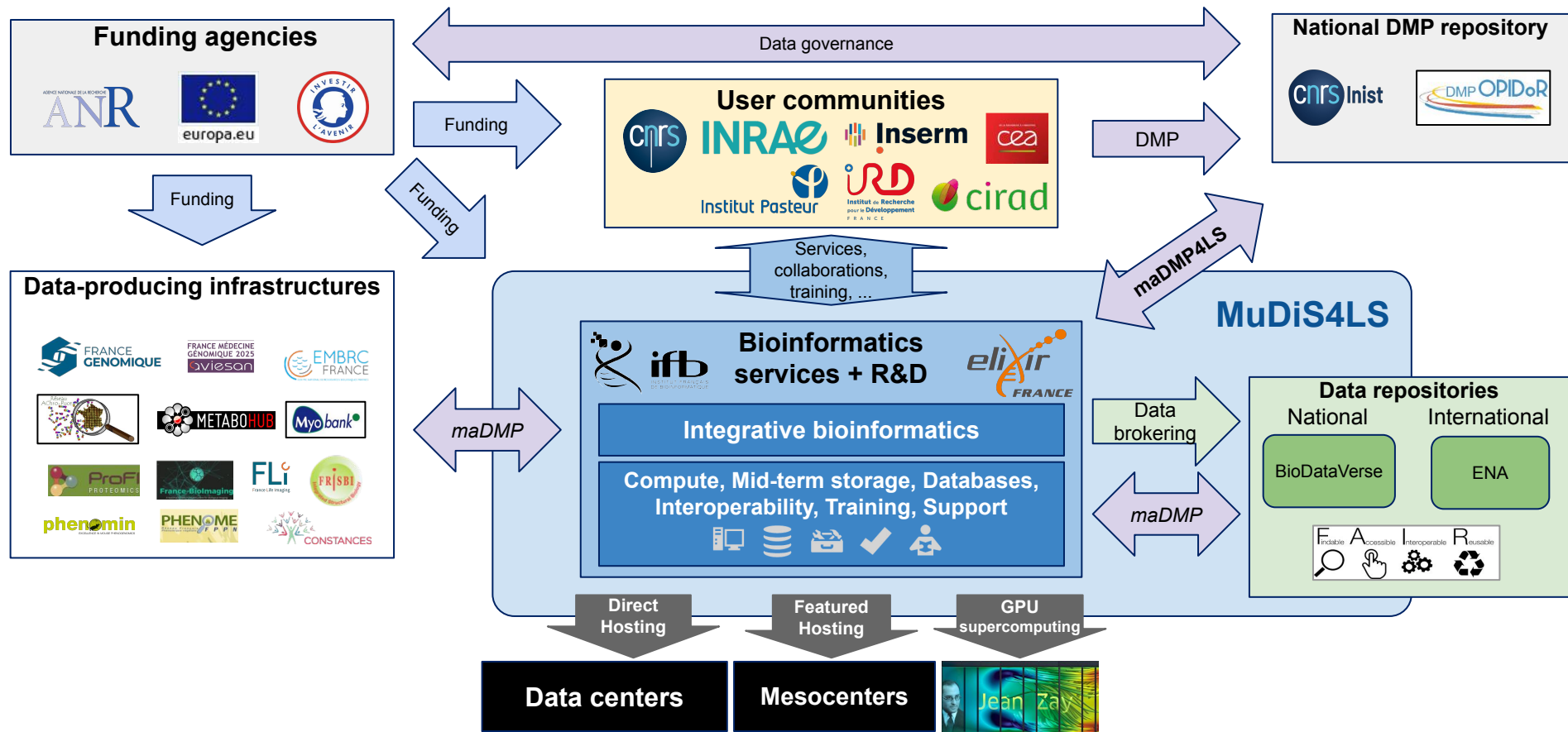
- ❑ Extension aux autres dépôts internationaux ?
- ❑ Traitement des données d'imagerie ?



At the cross-road of data flows – from project conception to results delivery



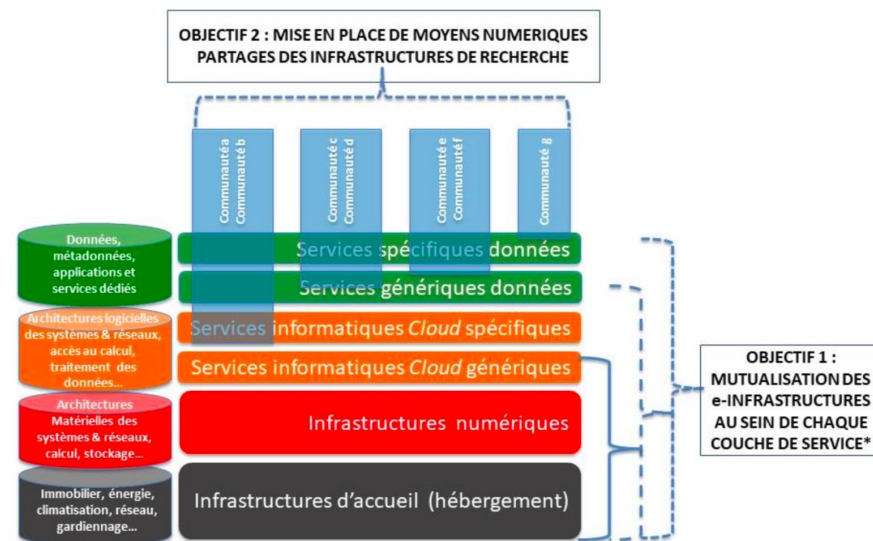
Orchestration des flux de données



Statistiques et figures

- **Rationalisation** des moyens numériques nationaux
 - ❑ Actuellement : 30.000 salles de calcul en France (labos, instituts, mésocentres, centres nationaux)
 - ❑ Demain : 30 data centres nationaux et régionaux labellisés
- **Mutualisation** des moyens numériques et humains
- **Engagement des établissements partenaires** → recrutement de personnel permanent
- **Structuration des services**
 - ❑ Infrastructure d'accueil (bâtiments, climatisation, électricité, réseau)
 - ❑ **Infrastructures numériques (moyens de stockage et calcul)**
 - ❑ **Services informatiques**
 - ❑ **Services données**

REPRÉSENTATION SCHÉMATIQUE DE L'ACCÈS DES INFRASTRUCTURES DE RECHERCHE AUX DIFFÉRENTES COUCHES DE SERVICE DES INFRASTRUCTURES NUMÉRIQUES (ou e-INFRA)

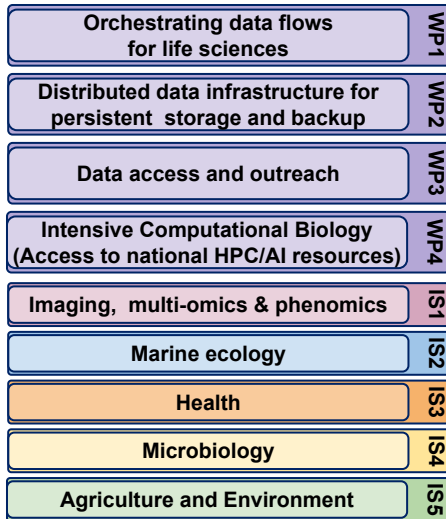


* Plus la couche concernée est basse, plus la mutualisation doit être démontrée.

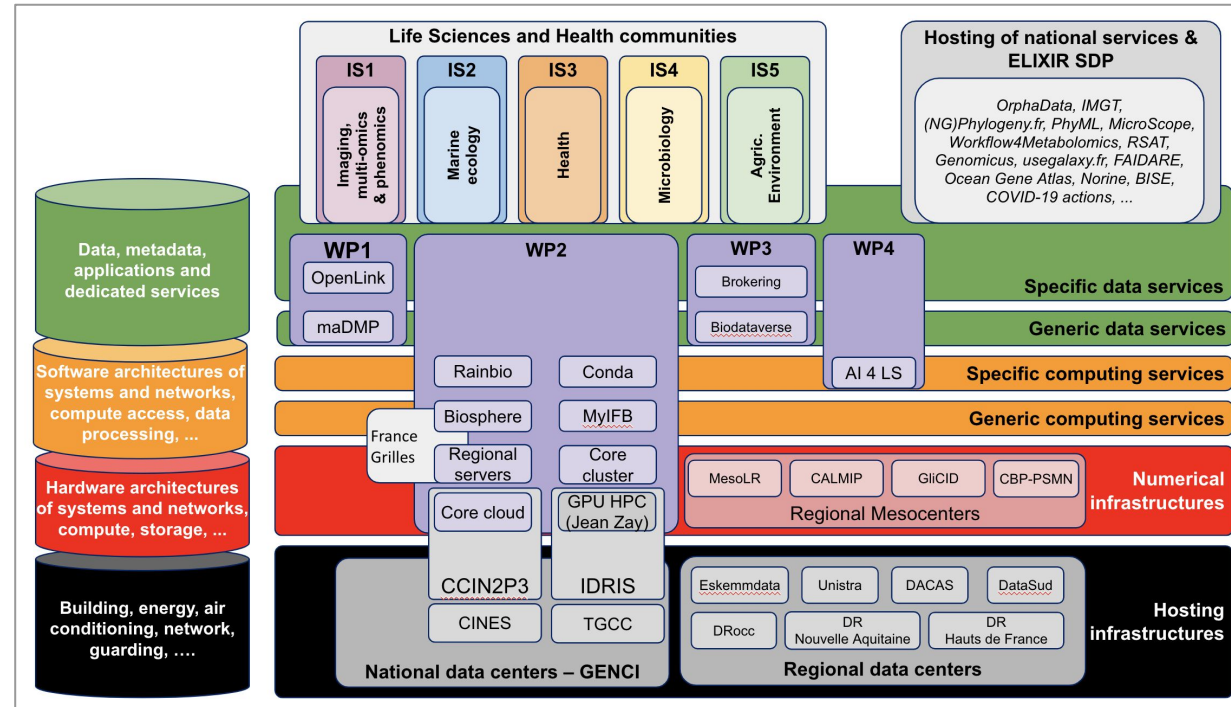
6.2 Représentation schématique de l'accès des infrastructures de recherche aux différentes couches de service des infrastructures

Key elements of the MUDIS4LS project

- 39 partner teams
- 14 partner organisms
- 4 national + 7 regional data centers
- 4 partner mesocenters
- 170 persons, 39 FTE;
→ 2.49 person•centuries
- Requested funding: 20M€
 - 13.70 M€ equipment
 - 3.53 M€ functioning
 - 1.28 M€ personnel
 - 1.81 M€ overheads
- 4 technological work packages (WP)
- 5 thematic implementation studies (IS)
- Engagement des tutelles : 7 postes permanents



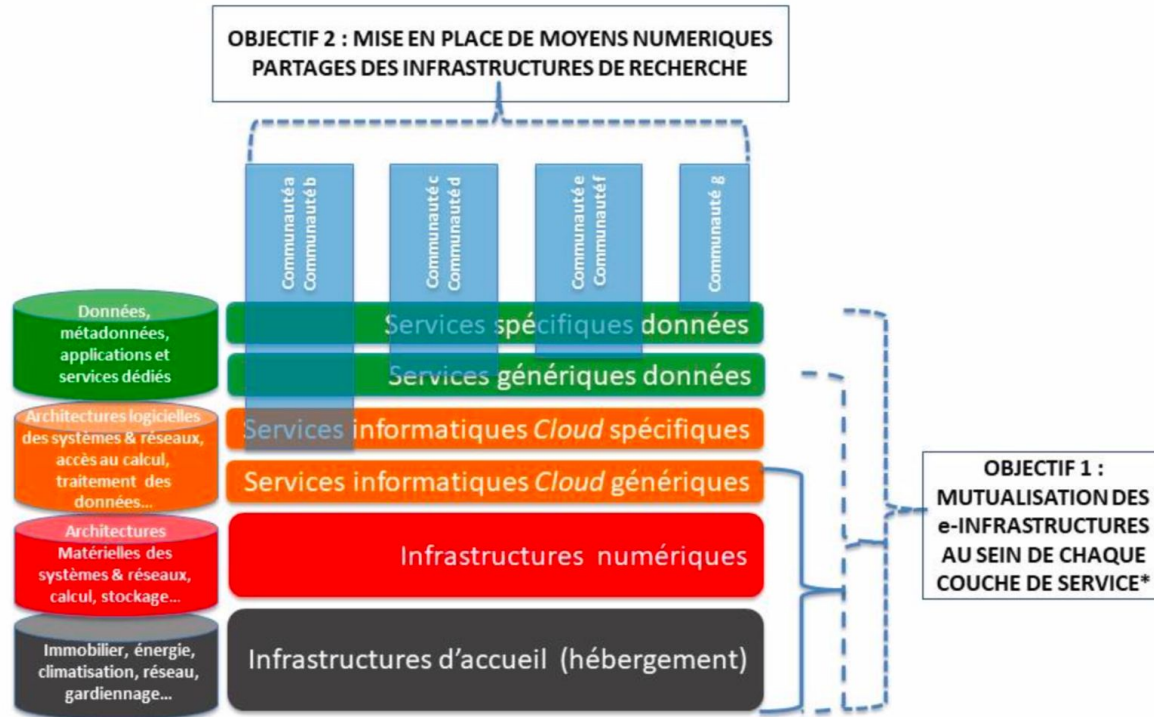
Mutualised Digital Space for Life Sciences (MuDiS4LS)



Récapitulatif total des demandes financières par destination

Description	Coût total	Aide demandée	Apport
Equipement	13 704 588.20 €	13 704 588.20 €	0.00 €
Personnel	14 285 017.04 €	1 278 609.60 €	13 006 407.45 €
Fonctionnement	3 533 523.90 €	3 533 523.90 €	0.00 €
Facturation interne	0.00 €	0.00 €	0.00 €
Frais de gestion	1 481 337.74 €	1 481 337.74 €	
Frais d'environnement	8 942 244.16 €		8 942 244.16 €
Total	41 946 711.04 €	19 998 059.44 €	21 948 651.61 €

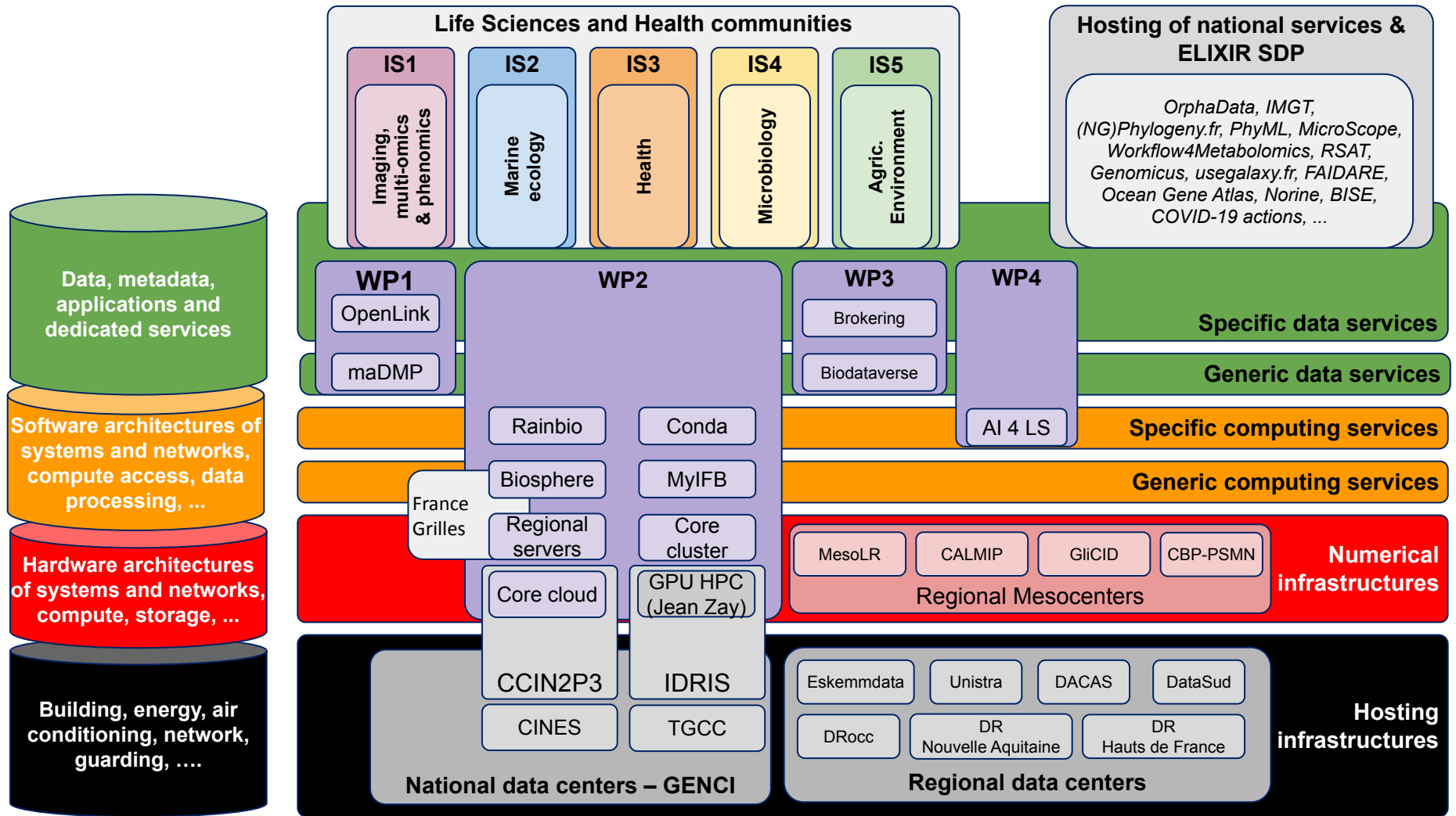
REPRÉSENTATION SCHEMATIQUE DE L'ACCES DES INFRASTRUCTURES DE RECHERCHE AUX DIFFÉRENTES COUCHES DE SERVICE DES INFRASTRUCTURES NUMÉRIQUES (ou e-INFRA)



* Plus la couche concernée est basse, plus la mutualisation doit être démontrée.

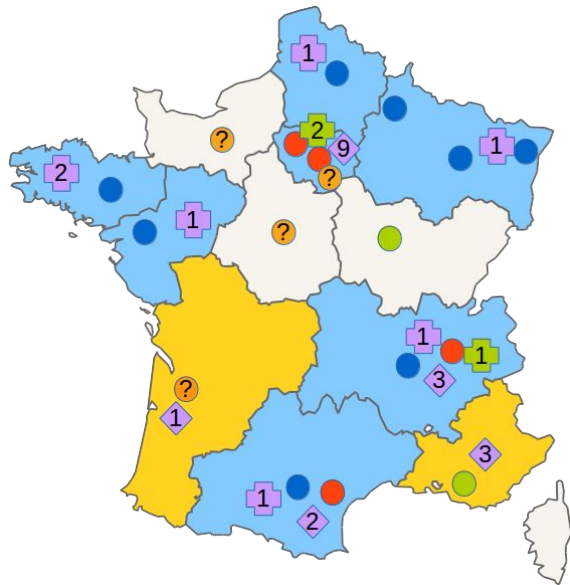
6.2 Représentation schématique de l'accès des infrastructures de recherche aux différentes couches de service des infrastructures

EQUIPEMENTS STRUCTURANTS POUR LA RECHERCHE /EQUIPEX+

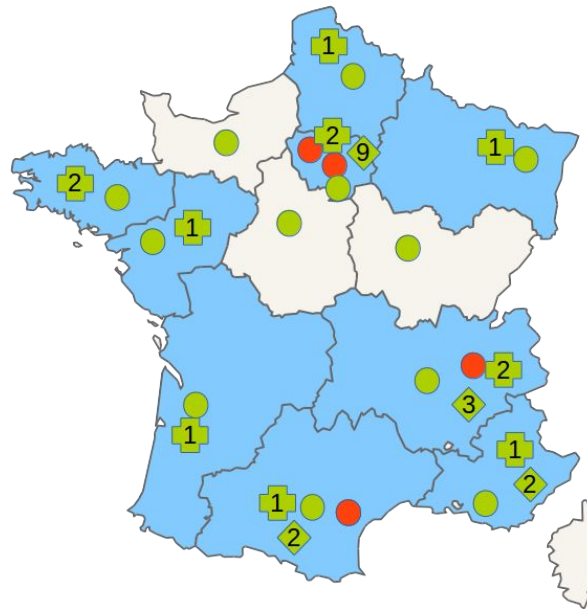


NNCR anchoring in national and regional infrastructures

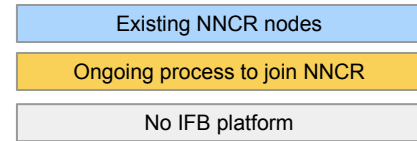
2020 status



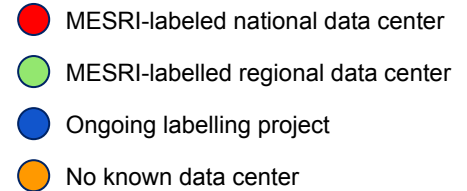
2025 target



Regions



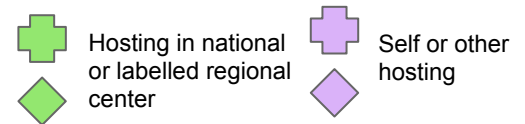
Regional status of data centers



NNCR status of IFB platform



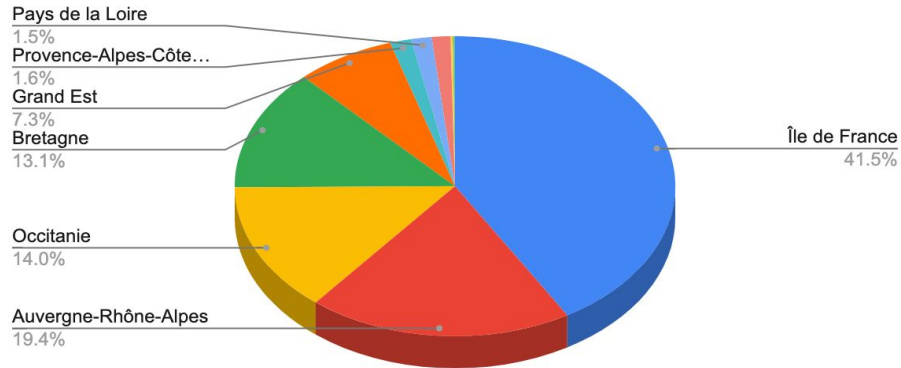
Anchoring status of IFB platforms



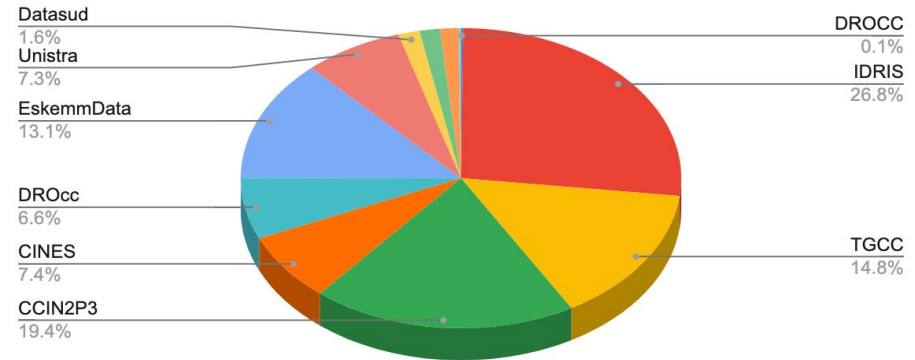
Regional distribution of equipment + hosting costs

Total cost (equipment + hosting costs): 17,471 k€

Cost per French region
(equipment + associated functioning)



Cost per hosting data center
(equipment + associated functioning)

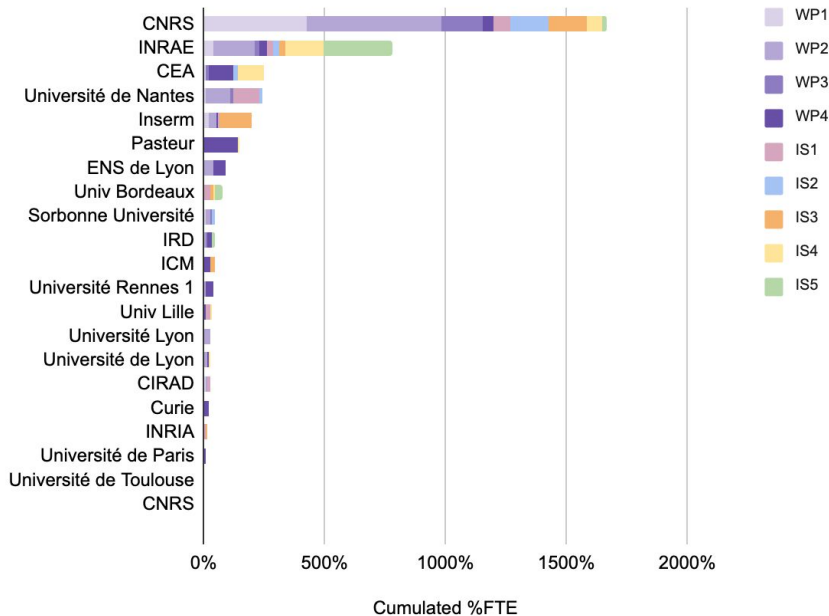


Distribution of the personnel by organism and WP/IS

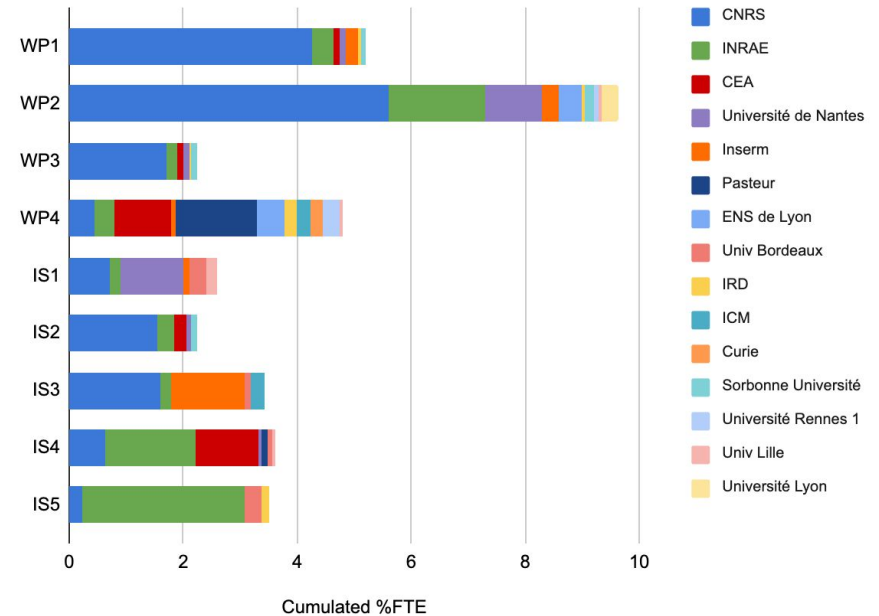
Total human effort

- 170 persons
- 39 FTE
- 2999.28 person•months = 2.49 person•centuries

Personnel per partner organism and per WP/IS



Personnel per WP/IS and per partner organism



Répartition régionale du personnel

Total human effort

- 170 persons
- 39 FTE
- 2999.28 person•months = 2.49 person•centuries

Personnel involvement per region / WP

