



# Actualités calcul et données au CNRS

**Denis Veynante**

**Direction des données ouvertes de la recherche (DDOR)**

→ JCAD – 4-6 novembre 2024

# Infrastructures numériques au CNRS

# CNRS et infrastructures numériques

## ➤ Opérateur de deux des quatre datacentres d'envergure nationale

### ➤ **IDRIS** (Orsay) Calcul intensif

- Opère le calculateur Jean Zay, financé par GENCI
- Centre de ressources pour la recherche en intelligence artificielle
- Projet CLUSSTER
- Hébergement : mésocentre Paris-Saclay, données CLIMERI, IFB, ...



### ➤ **CC-IN2P3** (Villeurbanne)

- Traitement de données massives pour les activités IN2P3 (LHC, LSST, ...)
- Hébergement : DSI CNRS, HAL, HumaNum, BBEES, ...



## ➤ Deux mésocentres rattachés au CNRS (UAR)

- **CALMIP** (Toulouse)
- **GRICAD** (Grenoble)

+ Demandes d'association d'autres mésocentres

# Jean Zay, supercalculateur GENCI opéré par IDRIS

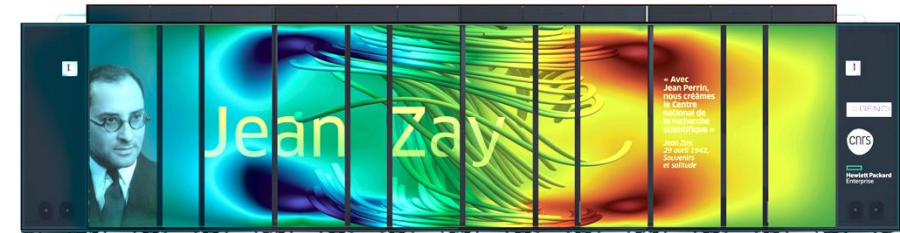
## ➤ Partition scalaire (CPU)

- 720 noeuds de calcul
  - 40 cœurs de calcul
  - 192 Go de mémoire

## ➤ Partitions accélérées (GPU)

- 396 noeuds quadri-GPU NVIDIA V100
  - 126 nœuds 4 GPU V100 – 16 Go
  - 270 nœuds 4 GPU V100 – 32 Go
  - 192 Go de mémoire / nœud
- 31 noeuds octo-GPU NVIDIA V100 – 32 Go
  - 20 nœuds à 384 Go mémoire
  - 11 nœuds à 768 Go mémoire
- 52 noeuds octo-GPU NVIDIA A100 – 80 Go (*extension juin 2022*)
  - 512 Go de mémoire / nœud
- 364 noeuds quadri-GPU NVIDIA H100 – 80 Go (*extension été 2024*)
  - 126 nœuds 4 GPU V100 – 16 Go
  - 270 nœuds 4 GPU V11 – 32 Go
  - 512 Go de mémoire / nœud

125,9 Pfllops crête



# Jean Zay, supercalculateur GENCI opéré par IDRIS

## ➤ Partition scalaire (CPU)

- 720 noeuds de calcul
  - 40 cœurs de calcul
  - 192 Go de mémoire

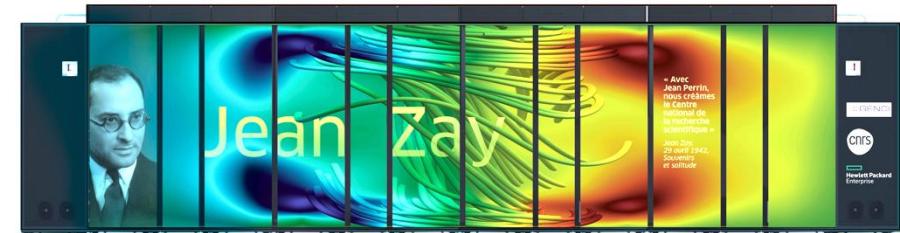
Après décommissionnement de  
808 noeuds (53 %) le 5/2/2024

125,9 Pfllops crête

## ➤ Partitions accélérées (GPU)

- 396 noeuds quadri-GPU NVIDIA V100
  - 126 noeuds 4 GPU V100 – 16 Go
  - 270 noeuds 4 GPU V100 – 32 Go
  - 192 Go de mémoire / noeud
- 31 noeuds octo-GPU NVIDIA V100 – 32 Go
  - 20 noeuds à 384 Go mémoire
  - 11 noeuds à 768 Go mémoire
- 52 noeuds octo-GPU NVIDIA A100 – 80 Go (*extension juin 2022*)
  - 512 Go de mémoire / noeud

Après décommissionnement de  
220 noeuds (36 %) 4 GPU V100 – 16 Go le 5/2/2024



- 364 noeuds quadri-GPU NVIDIA H100 – 80 Go (*extension été 2024*)
  - 126 noeuds 4 GPU V100 – 16 Go
  - 270 noeuds 4 GPU V11 – 32 Go
  - 512 Go de mémoire / noeud

# Extension supercalculateur Jean Zay (été 2024)

- **Commande de l'Etat pour intelligence artificielle**
  - Budget spécifique supplémentaire : + 40 M€
- **Contraintes capacité électrique et de refroidissement du centre**
  - Décommissionnement d'une partie de la configuration initiale
    - 53 % partition CPU, 36 % partition GPU V100
  - Emoi des communautés scientifiques
  - Transferts vers les autres centres (CINES et TGCC)
    - Pas de soucis de ressources globales CPU, toutes machines GENCI confondues
    - Accompagnement des utilisateurs
- **Demandes CPU inférieures aux ressources (novembre 2024)**
  - Explications ?
    - Autocensure ? Récupération des heures PRACE ?
    - ...
- **Ne pas négliger ressources EuroHPC !!!**

# Ouverture des données de recherche

# Motivation : ouverture des données

- **Assurer l'intégrité scientifique** (reproductibilité et validation des résultats)
- **Rendre la recherche plus efficace et non redondante** (pas de duplication inutile)
  - taux de perte des données estimé à 20 % / an
- **Être en capacité de réutiliser les données même sans en être à l'origine**
- **Croiser les données** (nouvelles analyses, voire nouvelles thématiques)
- **Satisfaire le cadre légal d'ouverture des données a priori :**
  - « *Ouvert autant que possible, fermé autant que nécessaire* »
  - *Obligation contractuelle (ANR, Europe, ...)*



# En pratique...

## ➤ **Des communautés très organisées**

- Physique des particules, Astronomie, Sciences de la terre ...

## ➤ **Une offre générique : Recherche Data Gouv**

- Entrepôt et catalogue, 20 ateliers de la donnée, 6 centres de références thématiques, 4 centres de ressources
- Espaces institutionnels (universités, organismes, ...)

## ➤ **... qui ne répondent pas à tous les besoins**

- Communautés qui ne disposent pas d'entrepôts thématiques
- Volumétrie limitée (Recherche data gouv : 50 Go par dépôt, 5 To par organisme)
- Besoin de capacité de traitement à proximité des données volumineuses (limiter les transferts)

## ➤ **Mutualiser et rationaliser infrastructures informatiques et ressources RH**

- *Datacentres labélisés*
- *Optimisation et réduction des coûts et de l'empreinte environnementale*
- *Pas de doublons inutiles, ni de trous*
- *Nouveaux métiers (« data stewardship », ...)*

**Ne pas développer sa propre solution !!!**

# CNRS Research Data

L'espace institutionnel CNRS  
dans l'entrepôt de données  
Recherche Data Gouv



# Objectifs

- Proposer un espace aux scientifiques pour déposer leurs données et les rendre accessibles lorsqu'il n'existe pas d'entrepôts thématique ou institutionnel de référence.
- Poursuivre la contribution à l'écosystème Recherche Data Gouv en créant un espace institutionnel CNRS.
  - Déjà présent via les Centres de Ressources (DORANum, OPIDoR) et les Centres de Références Thématiques (CDS, HUMA-NUM, IFB, etc.)
- Mutualiser les efforts en évitant le déploiement et la gestion d'une instance Dataverse complète propre au CNRS.



# Prérequis

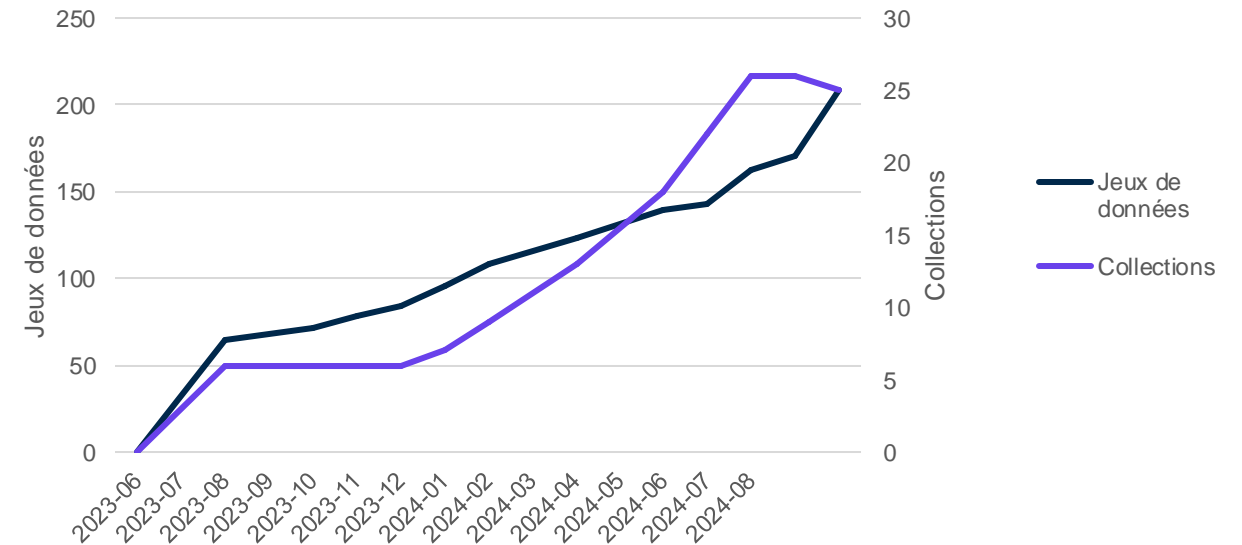
1. Le jeu de données n'a pas déjà été publié.
2. Il n'y a pas
  - D'entrepôt thématique plus pertinent que Recherche Data Gouv
  - D'espace Recherche Data Gouv institutionnel plus pertinent que l'espace CNRS
3. Le déposant a le droit de diffuser les données sur le plan légal :
  - Absence de données sensibles,
  - Respect des réglementations,
  - Accord des coauteurs.



# Statistiques

## Ouverture le 29/06/2023

- 25 collections
  - 8 créées (ICMCB, OSUC, F2D2, IDEES, LNCMI)
  - 14 liées (OSUG, I3S)
  - 3 déplacées (BETA)
- 4 collections à publier (MEDG, Robotex, LOCIE)
- 4 collections en discussion
  
- 209 jeux de données
  - 61 déposés
  - 148 liés
  - 13 en cours de curation
  
- ≈ 100 contacts, 14 demandes orientées vers des entrepôts thématiques
  - 2 PROGEDO, 1 Nakala
  - 8 Data.InDoRES
  - 1 CDS
  - 2 eBRAINS



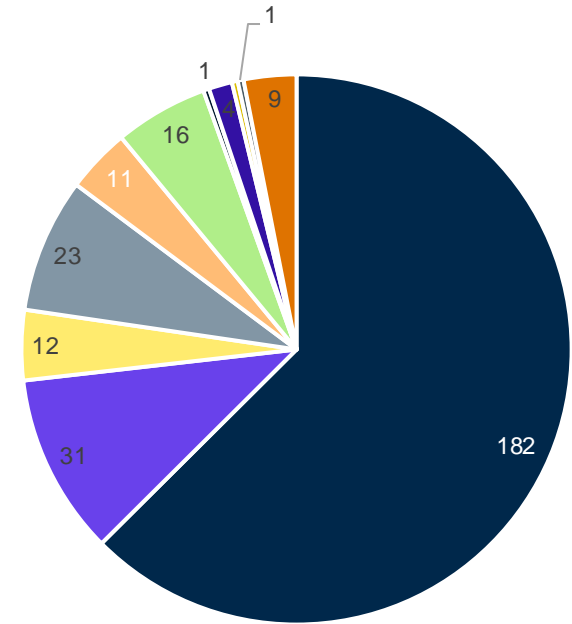
# Statistiques

## Dépôts, téléchargements et types de données

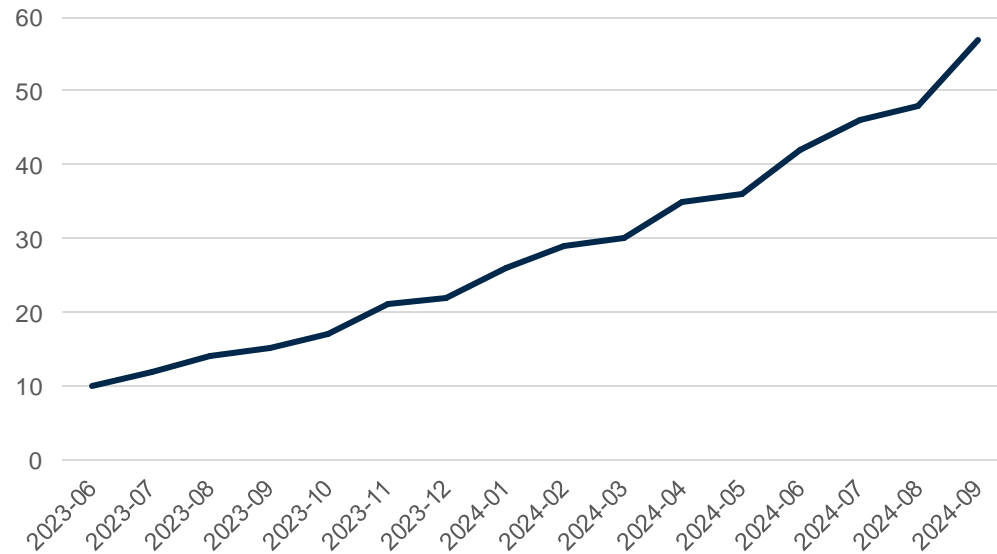


### Types de données

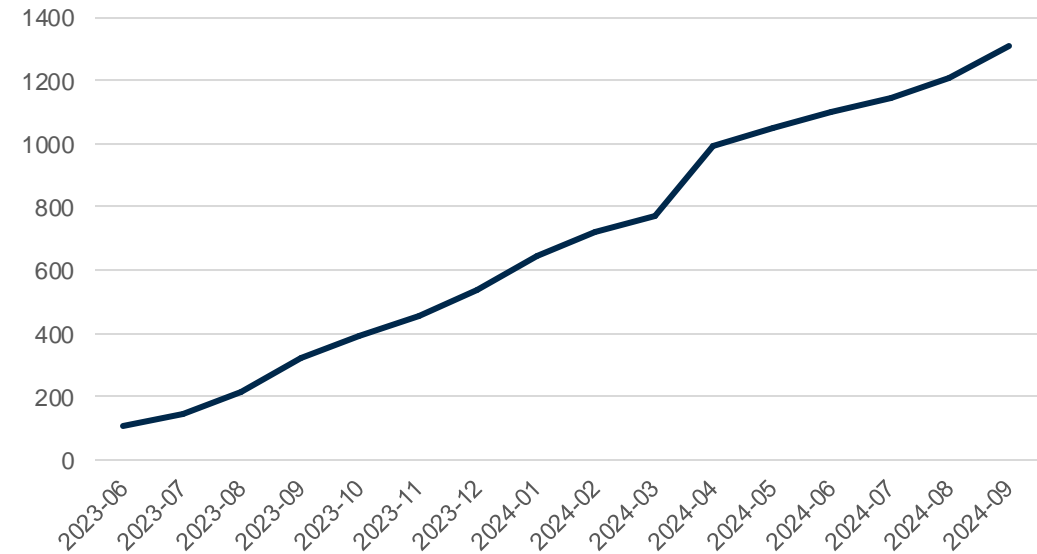
- Dataset
- Image
- Software
- Model
- Text
- Physical object
- Interactive resource
- Audiovisual
- Sound



### Dépôts



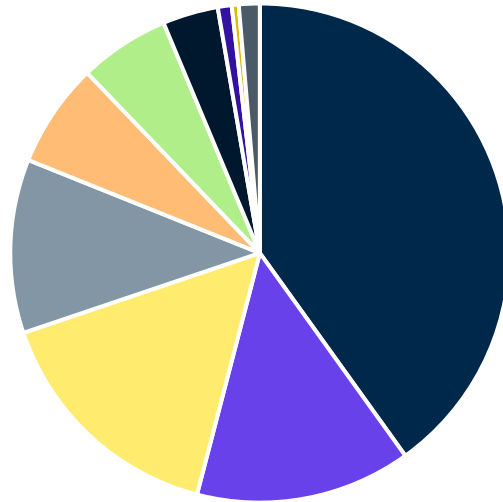
### Téléchargements



# Statistiques

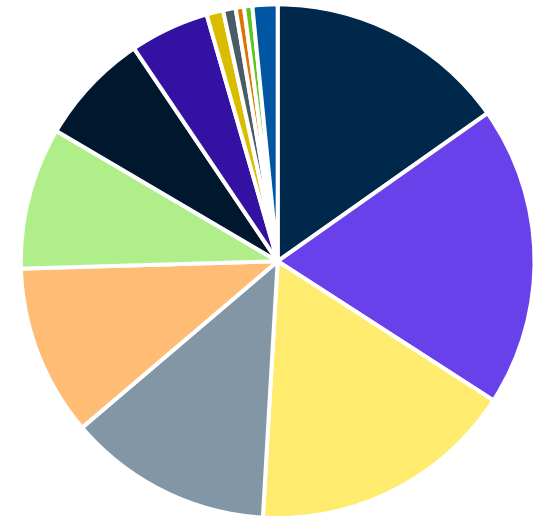
- Experimental data
- Analysis data
- Observational data
- Simulation data
- Survey data
- Computer code
- Aggregate data
- Text corpus
- Audiovisual corpus
- Other

Origine des données



- Chemistry
- Earth and Environmental Sciences
- Physics
- Engineering
- Medicine, Health and Life Sciences
- Computer and Information Science
- Agricultural Sciences
- Social Sciences
- Arts and Humanities
- Astronomy and Astrophysics
- Mathematical Sciences
- Law
- Other

Sujet



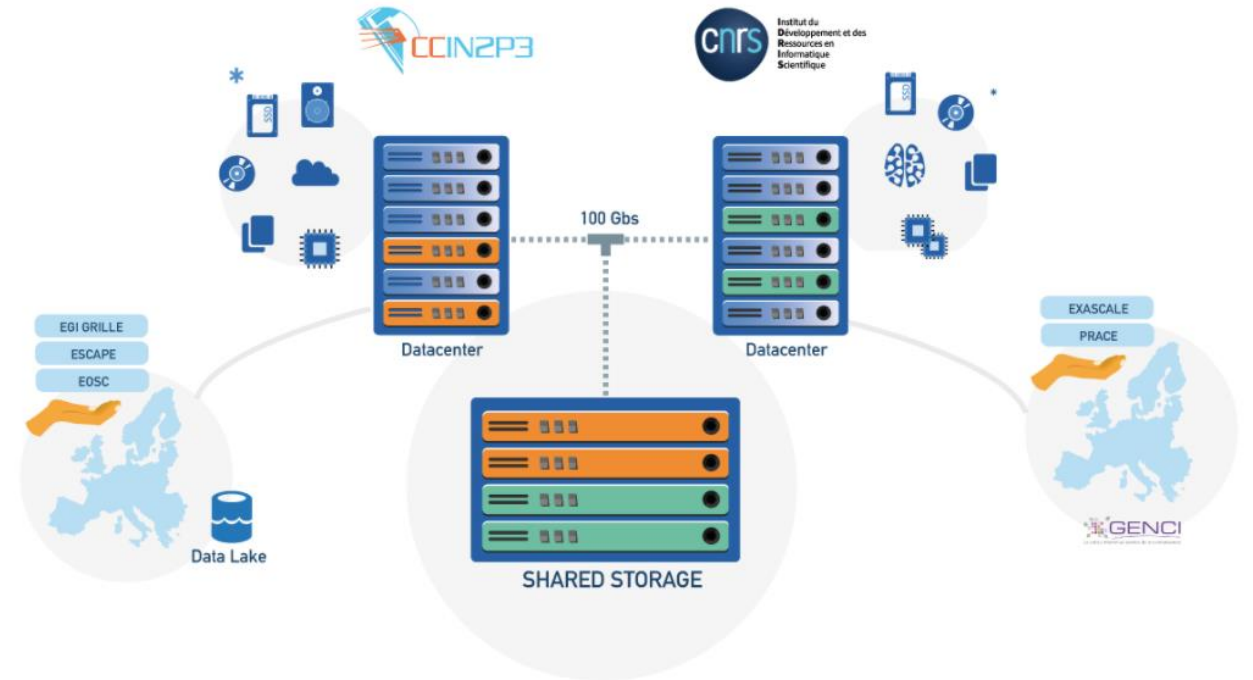
# Projets autour du stockage, traitement et mise à disposition des données de recherche



# Equipex+ FITS (CNRS Federated IT services for Research Infrastructures)

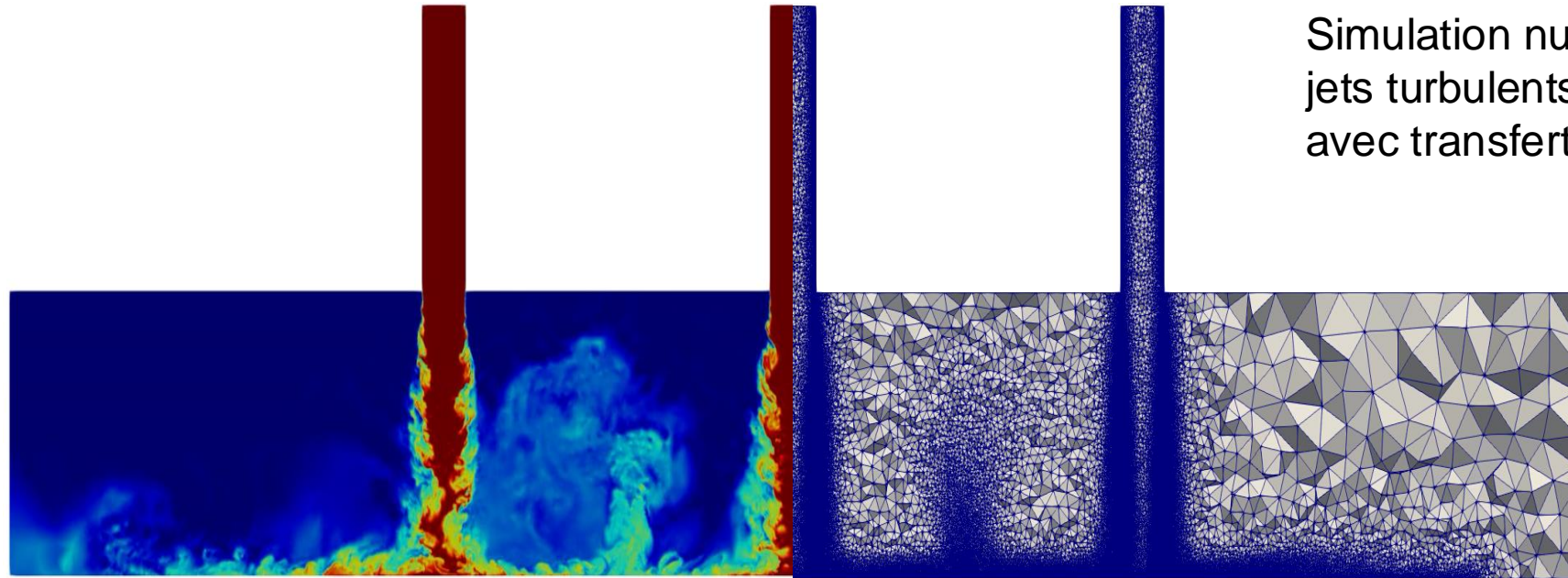
## Infrastructure répartie de stockage, traitement, mise à disposition, diffusion et valorisation des données au service des IR/IR\*

- Basée sur CC-INP2P3 et IDRIS + partenariat GENCI (calculateur Jean Zay)
- Portail unique d'accès aux ressources
- 4 cas d'usage : Soleil, HL-LHC, LSST, IFB
- Extension des capacités des centres
- 15,4 M€ dont 11,4 M€ de travaux
- 8 ans (juin 2021 – juin 2029)
- Mise en place d'un modèle économique
- [www.fits.cnrs.fr](http://www.fits.cnrs.fr)



**Offre réservée aux IR/IR\***

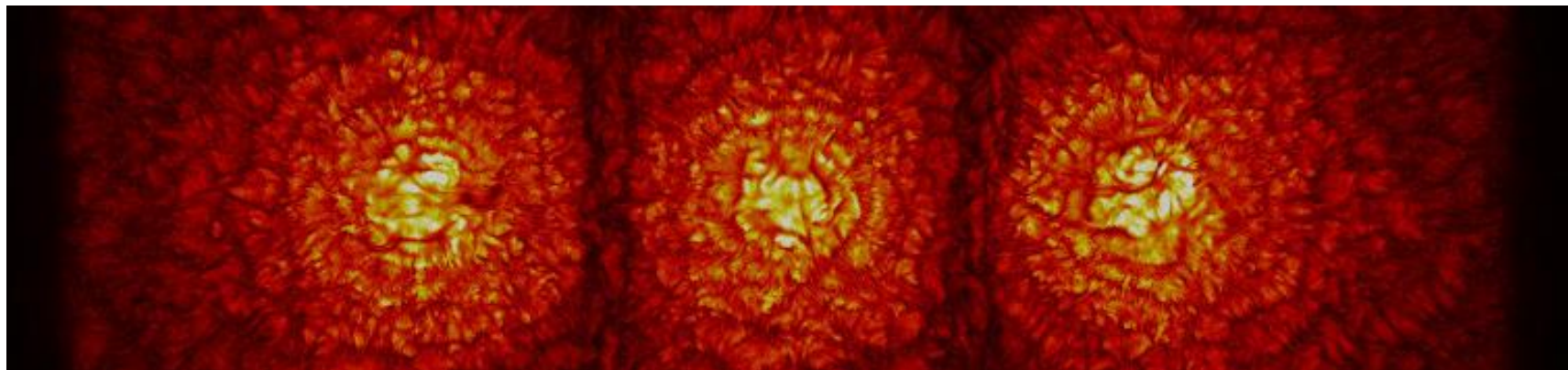
# Un autre besoin : cas d'usage HiFiLES4ML



Simulation numérique directe de jets turbulents impactant une paroi avec transferts thermiques

## Objectifs :

- Compréhension physique
- Validation de modèles



# Projet : services datacentre national à IDRIS

## ➤ Développement d'une offre de service pour les données:

- Trop volumineuse pour relever de Recherche Data Gouv
- Issues de communautés ne disposant pas d'entrepôts thématiques

## ➤ Trois types de services :

- Stockage de données massives (*volumétrie cible de plusieurs Po*)
  - Chaud (*technologie de type disques*)
  - Froid (*bandes magnétiques*)
- Traitement de données avec une puissance cible de plusieurs PFlops
  - CPU (*nœuds classiques et nœuds grosse mémoire*)
  - GPU
- Services d'hébergement de matériels informatiques

## ➤ Rationalisation, structuration et extension de services existants :

- Développés pour répondre à des demandes ponctuelles
  - Hébergement et mise à disposition de données du climat (*IR CLIMERI*)
  - Hébergement de calculateurs (*Mésocentre UPSaclay, IFB, ...*)
- Empreinte environnementale maîtrisée et à l'état de l'art

**En veillant à la cohérence avec  
les services et projets existants  
(FITS, CLIMERI, ...)**

# Positionnement

Typologie des services

## Typologie des utilisateurs

IR/IR\*      Projets nationaux      Communautés      Mésocentres et laboratoires      Autres

Stockage et mise à disposition de données

Recherche Data Gouv

Stockage et mise à disposition de données massives

Traitement de données

Hébergement

FITS

Services Datacentre national  
IDRIS

### ➤ Infrastructure de services extensible et flexible

- Adaptation aux nouvelles demandes et nouveaux besoins
- Hors ZRR

### ➤ Problématiques communes et convergence avec FITS, Clusster, Numpex

- Authentification, portail, cybersécurité
- Mutualisation logicielle et matérielle autant que possible

# Evolution : projet commun

*Déploiement d'une nouvelle génération d'offre de service de stockage, traitement et mise à disposition de données scientifiques  
Data Terra – France-Grilles - IDRIS*

- **Analyse des besoins des infrastructures de service aux données (ISD)**
  - Recommandation du document d'orientations stratégiques du CoSIN
- **Déploiement d'une offre nationale de services :**
  - Stockage de données massives
  - Traitement
  - Hébergement
- **Interconnexion des infrastructures de stockage**
  - IDRIS, mésocentres de Clermont-Ferrand et Strasbourg
- **Analyse des coûts et modèle économique**
  - Assurer la pérennité de l'infrastructure

# Statut du projet

## ➤ Sollicité et soutenu par la DGRI

- Embryon du « cloud stockage et traitement de données »
- Retenu dans le cadre du « fonds d'amorçage CoSIN » à hauteur de 2 M€

## ➤ Co-financé par le CNRS

- Pour 500 k€, acquis et notifiés

## ➤ Calendrier

- A préciser en fonction des ressources humaines disponibles
- Déploiement espéré pour les premiers utilisateurs courant 2025

## ➤ Remarques complémentaires

- Le projet est extensible et pourra intégrer d'autres mésocentres à l'avenir
- L'insertion harmonieuse avec les autres acteurs du « cloud » sera un point d'attention

# Conclusions - résumé

## ➤ Développement d'une offre stockage et traitement de données

- Répondre aux besoins non-couverts par les infrastructures actuelles
- Infrastructure mutualisée, optimisée et à l'empreinte environnementales maîtrisée

## ➤ Deux projets ambitieux et complémentaires...

- FITS
- Offre de services datacentre national

## ➤ ... Qui prendront du temps

- FITS attendu pour juin 2029, au-delà des 4 cas d'usage intégrés
- Déploiement progressif de l'offre datacentre national espéré à partir de mi-2025

## ➤ Modèle(s) économique(s)

- Durée de conservation des données ?
- Ces infrastructures ont un coût : investissements, fonctionnement, jouvence, extension...
  - Seul « l'amorçage » est aujourd'hui couvert
  - Facturations ? Qui paie ?

*Ne pas disperser les efforts  
ni construire de solutions individuelles  
ad hoc !!!*



**Merci de votre attention**