
Des environnements de calcul sécurisés pour les chercheurs et datascientists

Gouvernance des données, protection des données et centre de calcul

23 janvier 2024



Kamel Gadouche, directeur du Centre d'Accès Sécurisé aux Données, kamel.gadouche@casd.eu







Données



Secure use file



FPR ou Scientific use files



Open data ou public use files

Données les plus détaillées possibles

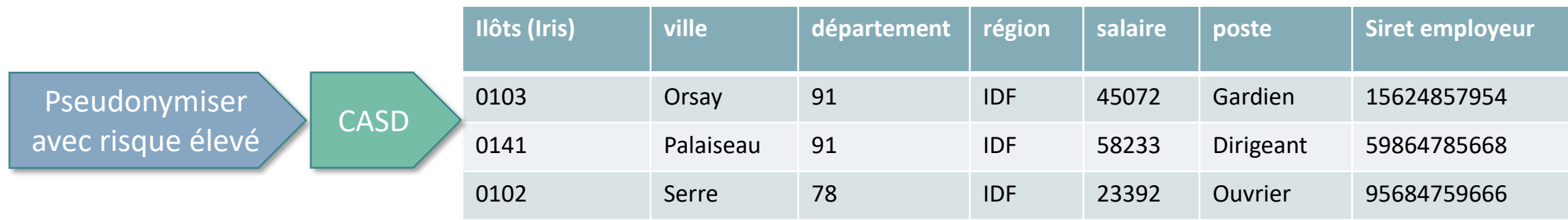
**Données à très faible risque
de réidentification**

**Données anonymisées
(risque théoriquement nul
de réidentification)**

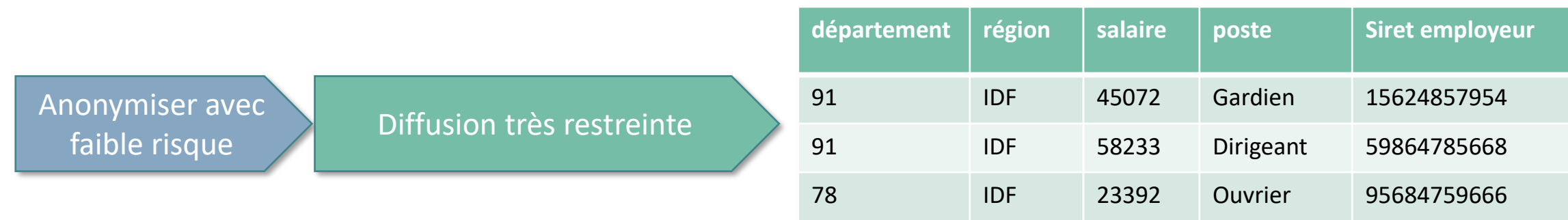
Fichier brut

nom	prénom	rue	ville	département	région	Tel	salaire	poste	Siret employeur
Dupont	Paul	Albert Camus	Orsay	91	IDF	6025212	45072	Gardien	15624857954
Durant	Julie	des sentiers	Palaiseau	91	IDF	6959849	58233	Dirigeant	59864785668
Morel	Pierre	Marie Curie	Serre	78	IDF	6446827	23392	Ouvrier	95684759666

Fichier « Secure Use File »



Fichier FPR « Scientific use file »



Fichier open data « Public use file »





Gouvernance

—



Loi 51-711

articles

1

le SSP = Insee + SSM

Statistiques issues des données d'enquêtes et administratives

3

Les données transmises ne peuvent faire l'objet d'aucune communication de la part du service dépositaire.

6

Transmission si décision de l'administration des archives, prise après avis du comité du secret statistique et relative à une demande effectuée à des fins de statistique publique ou de recherche scientifique ou historique.

6

Les renseignements ne peuvent en aucun cas être utilisés à **des fins de contrôle fiscal ou de répression économique.**

6bis

Les bénéficiaires des communications de données prises après avis du comité du secret statistique s'engagent à ne communiquer ces données à quiconque.

7bis

Obligation de transmission de données administrative à des fins d'établissements de statistiques avec interdiction de les céder pour le service dépositaire à un autre organisme.

7ter

Le comité du secret statistique peut autoriser l'accès aux données issues du 7bis pour des besoins de recherche scientifique ou de réalisation d'étude.

Le service statistique public (SSP) – Art. 1

Définition du Service Statistique Public

Gouvernance des données :

- Conseil national de l'information statistique
- Autorité de la statistique publique



Entreprises



Citoyens

Données brutes Art. 1

Article 3bis

Collecte

Article 1bis



Insee

SSMs

DARES

Agri-SSP



Produisent

Statistiques et études

Diffusent en respectant le secret statistique (7bis)

Citoyens, média, pouvoirs publics, ...

Autres organismes publics

DGFip



cnav



Agriculture



Personnel soumis au secret professionnel



Chercheurs et statisticiens (hors SSP)

Le service statistique public (SSP) – Art. 1

Collecte des données issues d'autres organismes.

Personnel soumis au secret statistique.

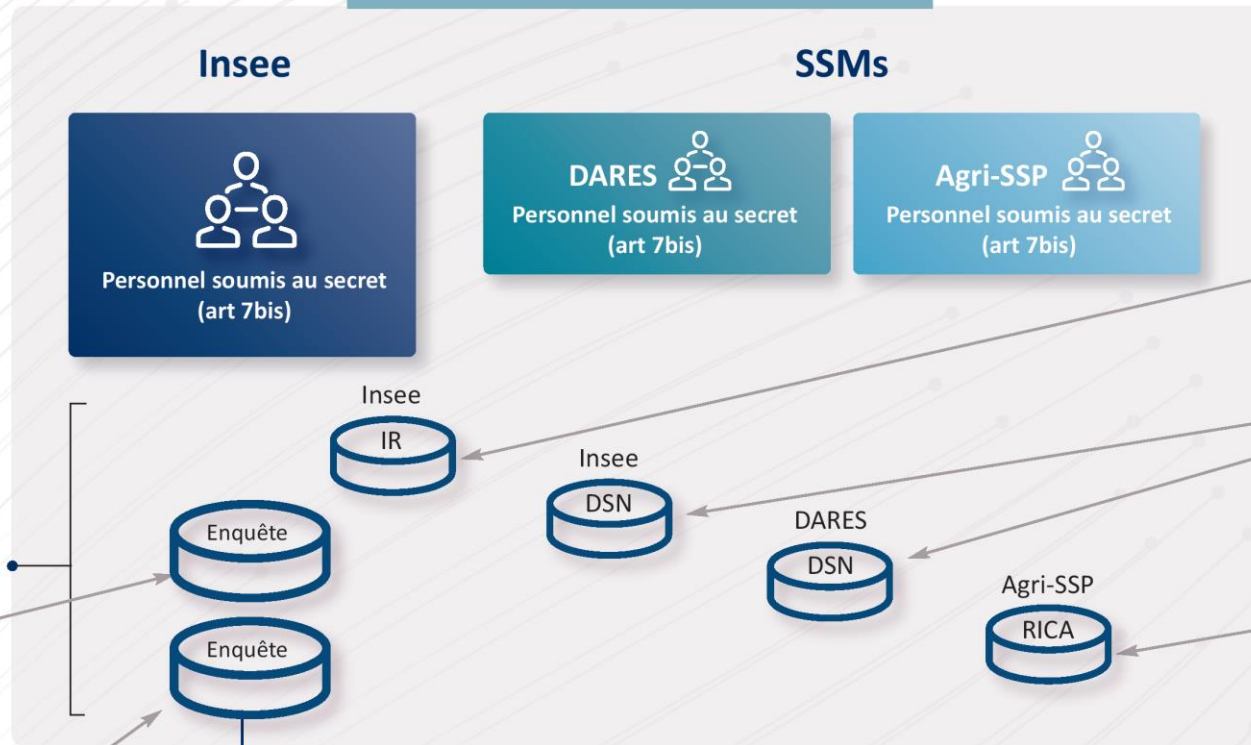
 Non cession des données.



Entreprises



Citoyens



Données brutes Art. 1

Article 3bis

Collecte

Article 1bis

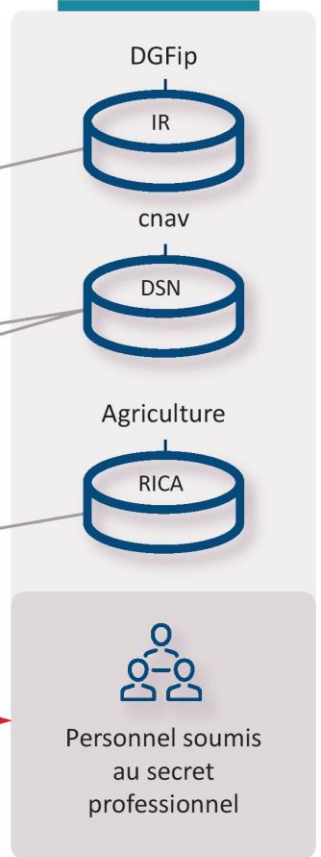
Produisent

Statistiques et études

Diffusent en respectant le secret statistique (7bis)

Citoyens, média, pouvoirs publics, ...

Autres organismes publics




Chercheurs et statisticiens (hors SSP)

Le service statistique public (SSP) – Art. 1

Comité du secret statistique.

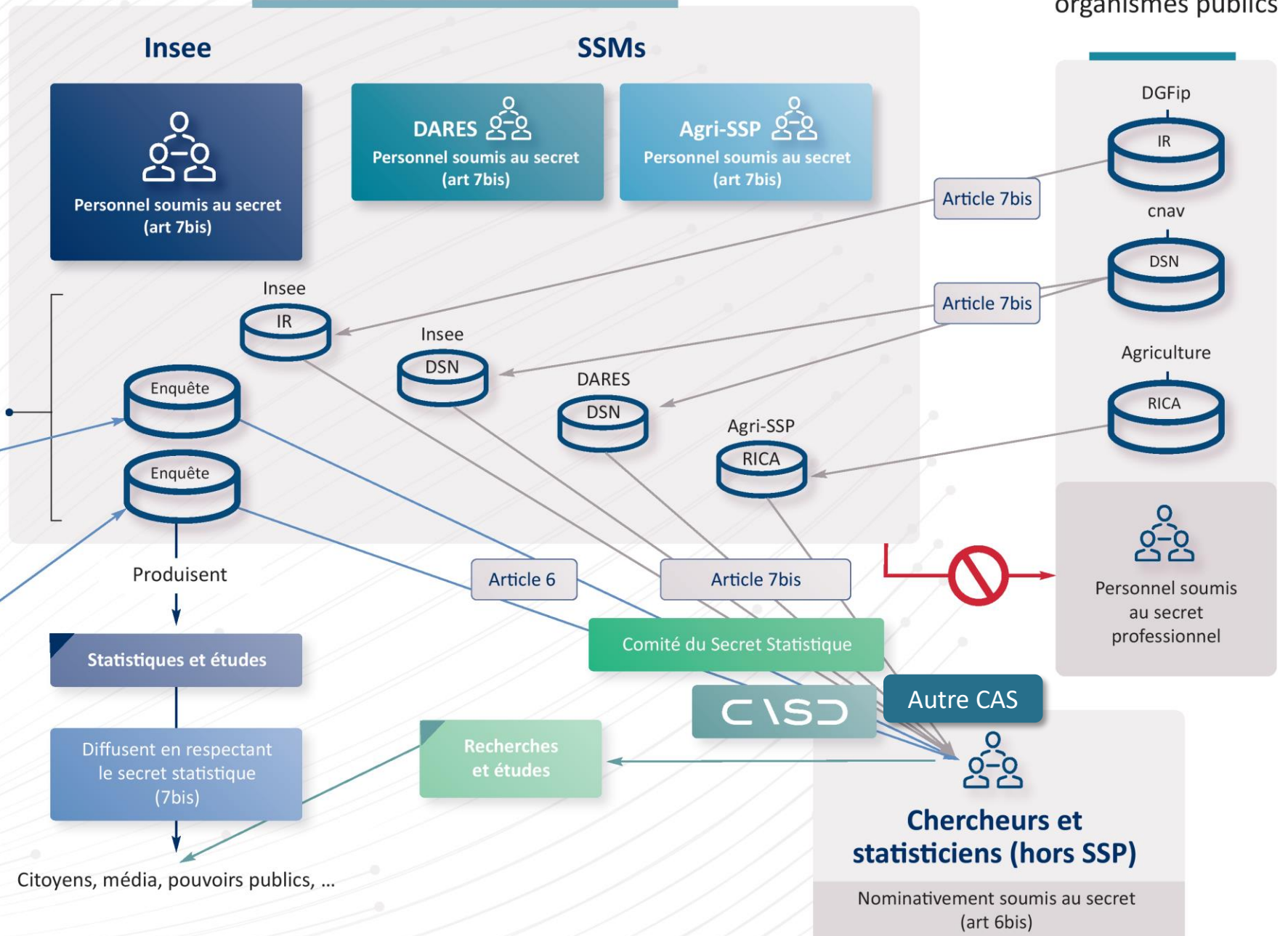
Accès aux données pour les chercheurs et statisticiens par le CASD et sont soumis au secret.



Entreprises



Citoyens



Citoyens, média, pouvoirs publics, ...

Chercheurs et statisticiens (hors SSP)

Nominativement soumis au secret (art 6bis)

Le CASD en bref

- Groupement d'intérêt public (GIP), consortium regroupant l'Insee, le Genes, la Banque de France, le CNRS, HEC Paris et l'École polytechnique
 - A but non lucratif
 - Facturation pour la recherche à marge négative (contributions des membres)
- 30 personnes
 - Service IT et Datascience
 - Service Data Management
 - Service Project Management
- 500+ sources de données mises à disposition
 - Plusieurs centaines de To de données
- 1000 institutions utilisatrices



Offre

Environnement
de calcul pour
la data science

Puissance de calcul paramétrable : 1
serveur à des clusters de serveurs

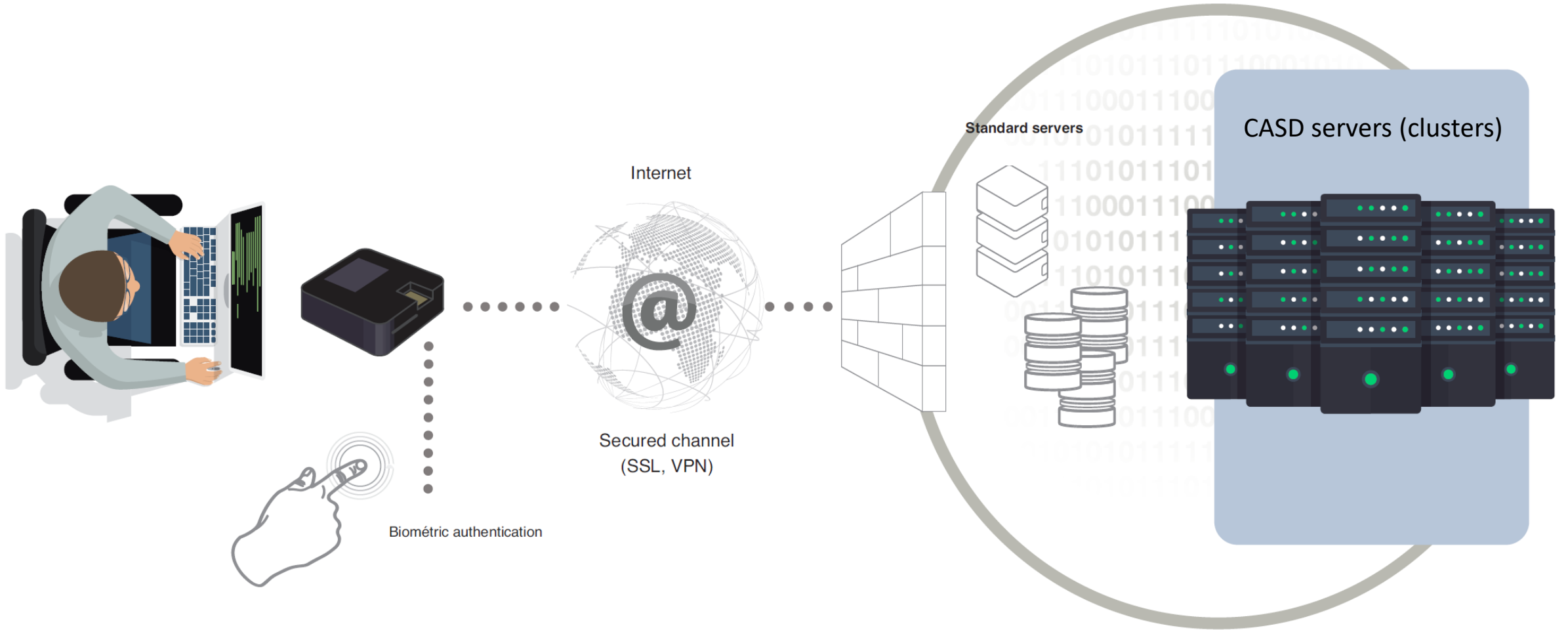
Large offre logicielle : plus de 80
logiciels en catalogue avec
notamment les packages pour R,
python et stata

Support technique, sauvegarde,
localisation multisite

Contrôle des sorties,
documentation, accompagnement
pour les démarches, appariements

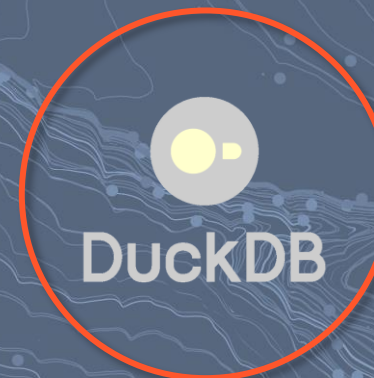


usages



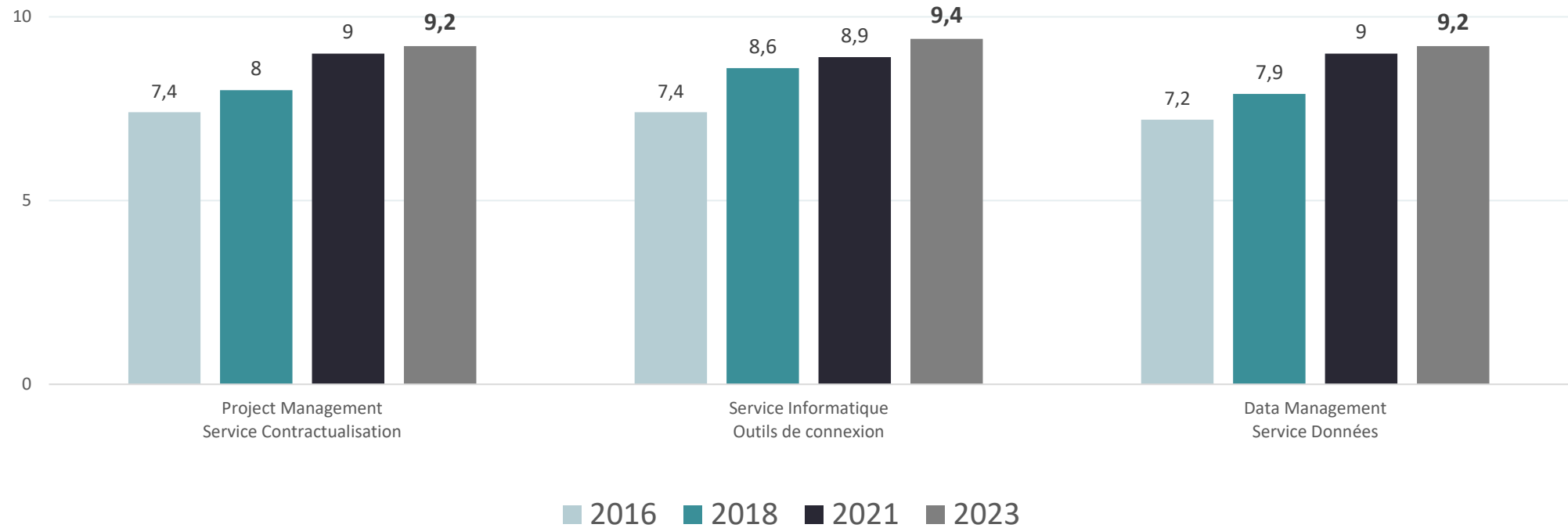


Datascience



Enquêtes de satisfaction

NOTATIONS PAR SERVICE DU CASD
EN 2016, 2018, 2021 ET 2023



Sécurité



Conformité **RGPD** par le Bureau Veritas et **autorisation CNIL** (n°2014-369)



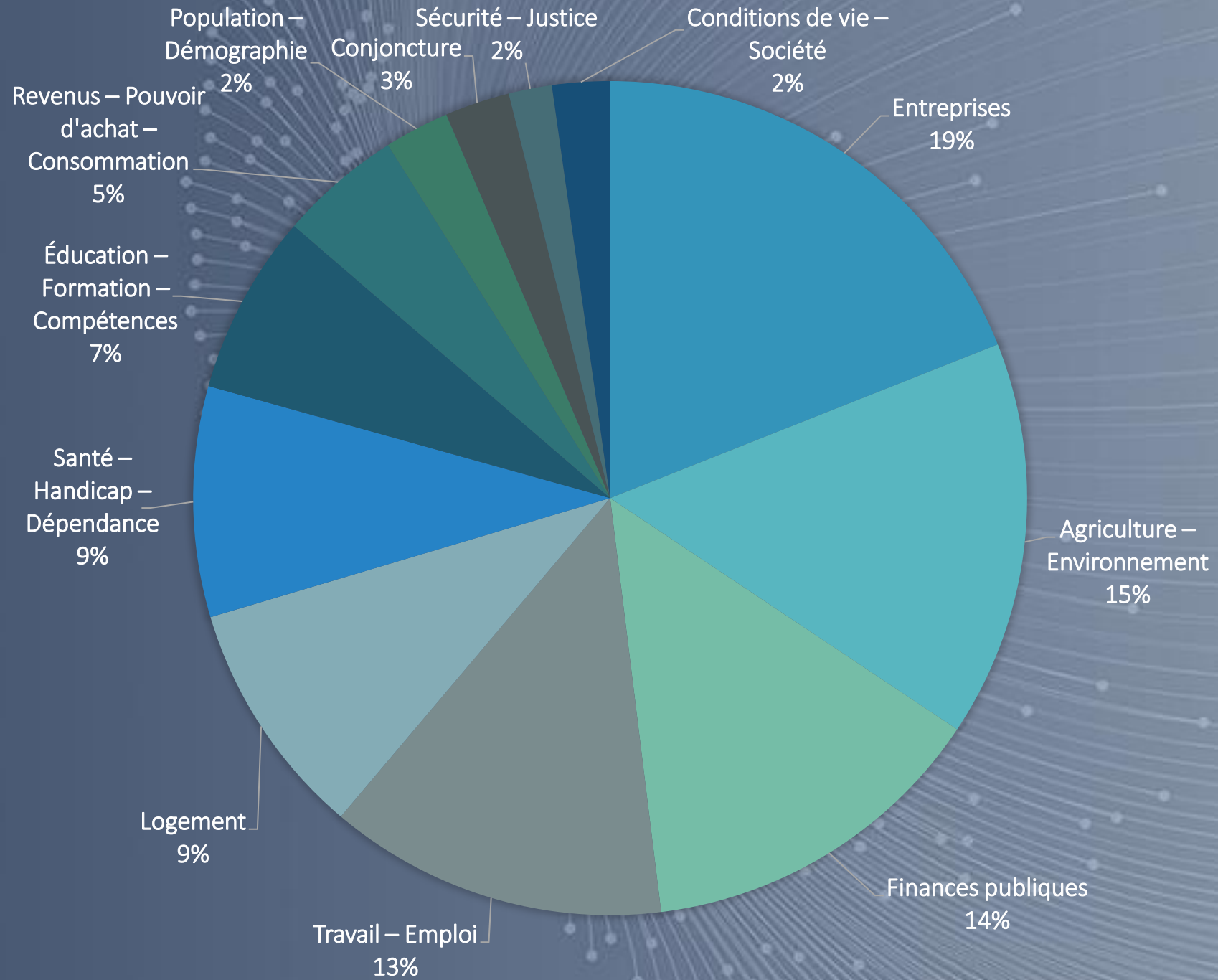
Certifications **ISO 27001** et **ISO 27701** (RGPD)
Certification « **hébergeur de données de santé** », homologué SNDS, EDS

Audits techniques réguliers de sécurité réalisés





Données



Plusieurs dizaines de producteurs dans de multiples domaines : insee, fiscalité, santé, travail, enseignement supérieur et recherche, agriculture, justice, sécurité intérieure, banque, etc.

Une procédure d'accès unifiée au bénéfice des chercheurs et datascientists, un atout essentiel

2/3 des 700 projets menés sur le CASD s'appuient sur des données de plusieurs producteurs

10% des projets concernent des études localisées dans un ou plusieurs territoires (Conseil régionaux ...)



Intersectoriel



Conclusions

Un appui initial et continu des producteurs et du CNIS, en particulier le comité du label.

Des progrès considérables permis à la gouvernance mise en place au fil des années.

Un dispositif intersectoriel qui découle de l'organisation de la gouvernance de la statistique publique et permet de prévenir les silos.

Un dispositif donc assez unique au monde à cette échelle qui s'explique en grande partie par la mise en place de lieux d'échange, de concertation et de décision : le comité du secret statistique, le CNIS, l'ASP, l'INSEE en tant que coordinateur...