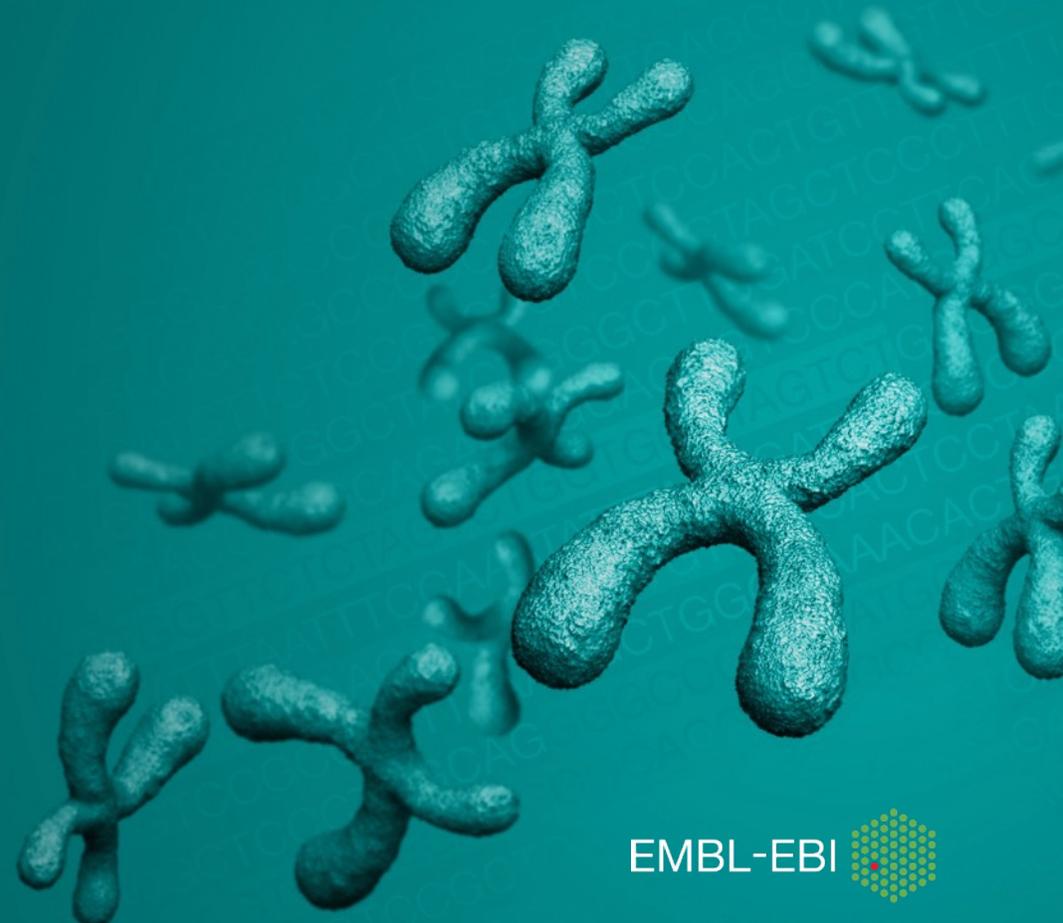


# Management of digital sequence information

Guy Cochrane



CCTAACCAAGTACTACTTAAAATCAGTAAACTGCCACTCCCTCCCCACTACCCTATATGA  
AAATAATTAAGCTCCTACGCTTCTTTTTTTGACCGCCCTCTACCTCTTTTCACAAA  
GGCTATCAACTCGAATTTTCTATTGCTTTATTTTATCAAAAAAACATAAAAATTTTGT  
ATCAAAAAATATCATTCAAGAGGGGGGATTTTCACCCCCACCGCCGGCACCAAAGCCAA  
AATTCTTGGAATCAAACCTCCTTGTCCCTGGTTTTCCCTAATGTACGAGTAACGTGGCAA  
CATATTATGCCTATCGTGATCCAACCTATCTGCCACACGACTATTTTTTAGTACTCTTCG  
GAGTGTAAAGCCGAACACTTAGGGCGAGAAACCACCAACCCGCTCCTGACGATACGATGAC  
CAGATCTGAGGACTTTTATTGTAGAGTGCCTTACTTCCCTTGAGGCGCCACTGGTTAAAA  
TCTATGGACACGACTCGAAGATTCATTCATCAATTGGATCGAACGGGTACCTGGCGGCTG  
CTTGATAACTAATCAGCCCATGATCCCTCAGCCTCCTCTAAGCACATCTGGTAATTTTTT  
ATTTTTCTGTGTGGTCAACCAACATTACGGTAATATGTCTGGTACTACACGATCTAAAGC  
TGAACATAACATGCAATTGTTTTATTTGGACCAAATAGAATGAGTAAATTATATGAATGA  
TTAT

# A sequence in isolation...

```
CCTAACCAAGTACTACTTAAAATCAGTAAACTGCCACTCCCTCCCCACTACCCTATATGA
AAATAATTAAAAAGCTCCTACGCTTCTTTTTTTGACCGCCCTCTACCTCTTTTCACAAA
GGCTATCAACTCGAATTTTCTATTGCTTTATTTTATCAAAAAAACATAAAAATTTTGT
ATCAAAAAATATCATTCAAGAGGGGGGGATTTTCACCCCCACCGCCGGCACCAAAGCCAA
AATTCTTGGAATCAAACCTACTCTTGTCCCTGGTTTTCCCTAATGTACGAGTAACGTGGCAA
CATATTATGCCTATCGTGCATCCAACCTATCTGCCACACGACTATTTTTTAGTACTCTTCG
GAGTGTAAGCCGAACACTTAGGGCGAGAAACCACCAACCCGCTCCTGACGATACGATGAC
CAGATCTGAGGACTTTTTATTGTAGAGTGCCTTACTTCCCTTGAGGCGCCACTGGTTAAAA
TCTATGGACACGACTCGAAGATTCATTCATCAATTGGATCGAACGGGTACCTGGCGGCTG
CTTGATAACTAATCAGCCCATGATCCCTCAGCCTCCTCTAAGCACATCTGGTAATTTTTT
ATTTTTCTGTGTGGTCAACCAACATTACGGTAATATGTCTGGTACTACACGATCTAAAGC
TGAACATAACATGCAATTGTTTTATTTGGACCAAATAGAATGAGTAAATTATATGAATGA
TTAT
```

.. cannot be interpreted

.. is unlikely to be informative

<https://www.ebi.ac.uk/ena/browser/view/HM544578.1>

# A collection of sequences..

260 million sequences



.. can be interpreted

.. drives discovery

# International Nucleotide Sequence Database Collaboration (INSDC)



## Values

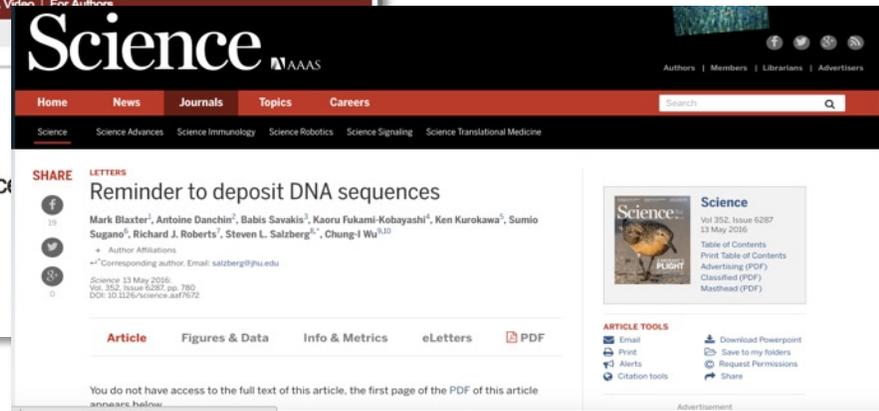
- open access for all
- globally comprehensive
- spanning life science domains
- permanent database of record
- public forum for the scientific process

## Organisation

- established early 1980s
- major ongoing investment
- structure and governance
- model for scientific collaboration



<http://www.insdc.org/>

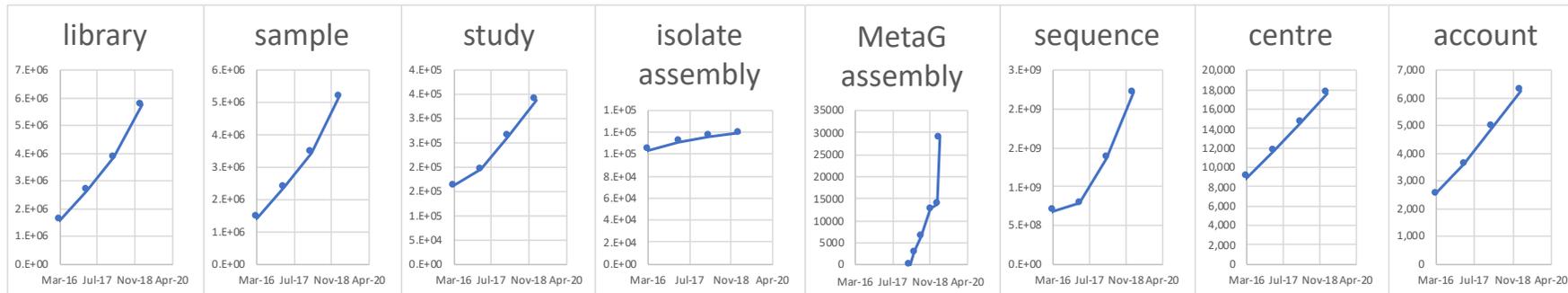


## Instruments

- regular data exchange
- accession scheme
- data standards
- mandatory submission agreement
- services and software (node-level)

# Scale

- **Rate:** 1 new dataset every 6 minutes
- **Data:**  $2 \times 10^9$  sequences and  $1 \times 10^{16}$  base pairs of read data across  $2 \times 10^6$  taxa
- **Usage:** 2,000 submitters; 10x thousands monthly consumers; 10x millions of monthly hits, many times this globally
- **Support:** 46 tickets per day and in-person training delivered to more than 350 users per annum
- **Infrastructure:** includes 100s of petabytes of storage, ~\$50 USD mil./year
- **Growth:**

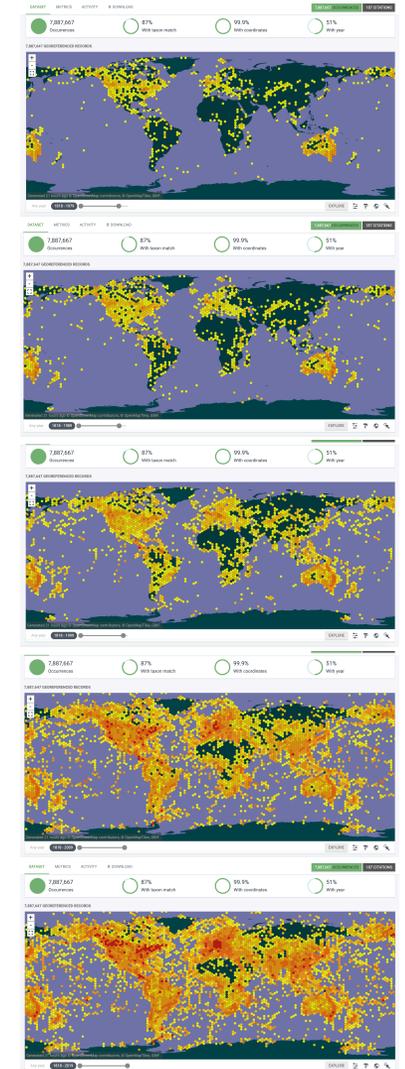


- **E.G. biodiversity data**

- 7.2 million records point to natural history museums, culture collections and biobanks
- 7.9 million sequences with geographical coordinates
- 2.1 million sequences have place name annotation
- 1 million raw sequencing data sets have coordinates

1980

2020

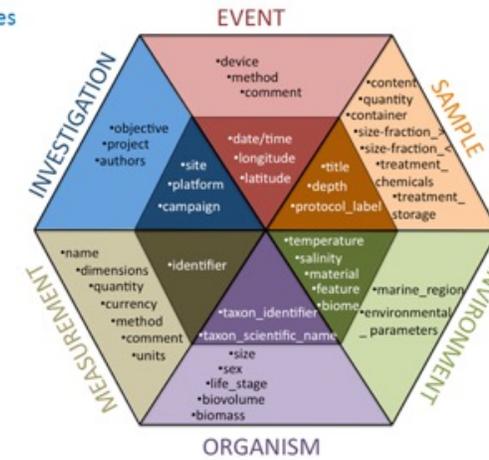


<https://doi.org/10.15468/cndomy>,  
courtesy of GBIF

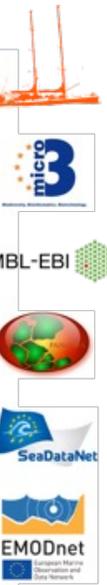
# The INSDC data cycle

- Data provider
  - Validation, organization, adding structure and curation
  - Compliance with standards and established conventions
  - Links to external data (e.g. academic publications)
- Database
  - Curation and integration with the overall corpus
  - Indexing for discovery and reuse
  - Open services, freely available to the world

## Registries & ontologies



Ten Hoopen P, Pesant S, Kottmann R, Kopf A, Bicak M, Claus S, Deneudt K, Borremans C, Thijssse P, Dekeyzer S, Schaap D, Bowler C, Glöckner F.O., Cochrane G. *Data standards for Marine Microbial Biodiversity, Bioinformatics and Biotechnology (M2B3) Standards in Genomic Sciences* 10:20 doi:10.1186/s40793-015-0001-5 (2015).  
<http://www.standardsingenomics.com/content/10/1/20>

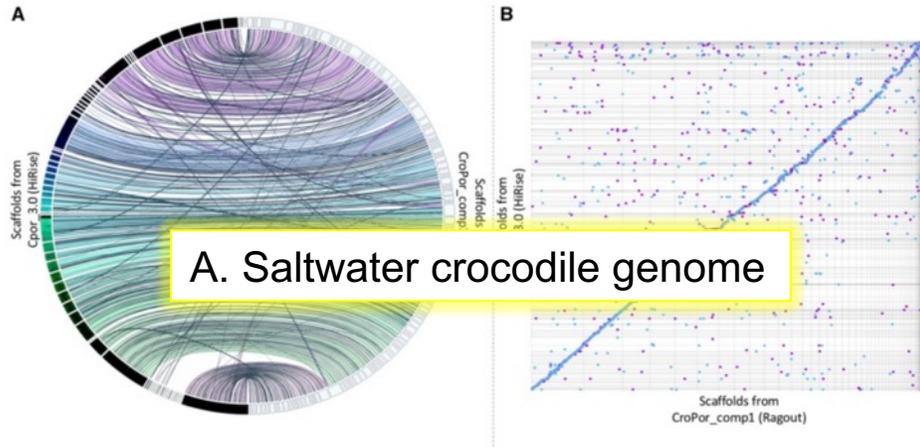


**Submission**

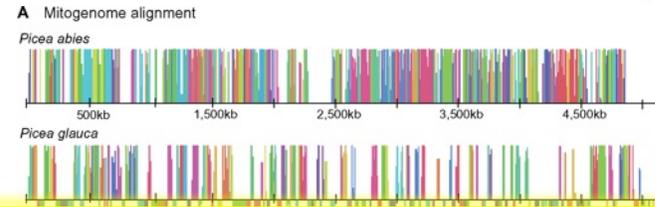
**Access**

# Direct use

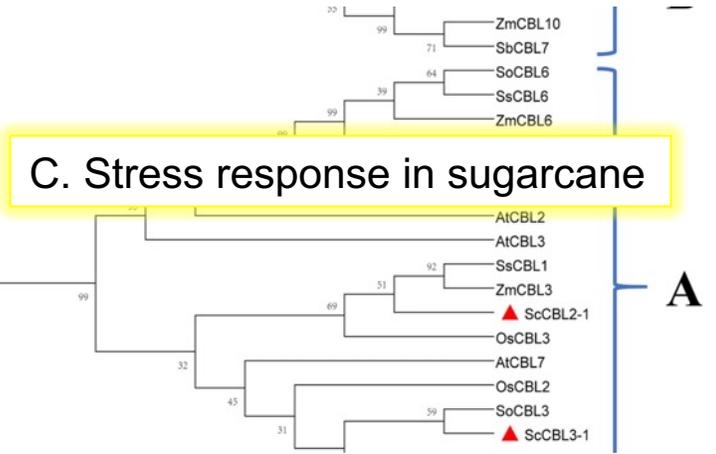
- 1,600 publications citing sequence accessions per month in 2020



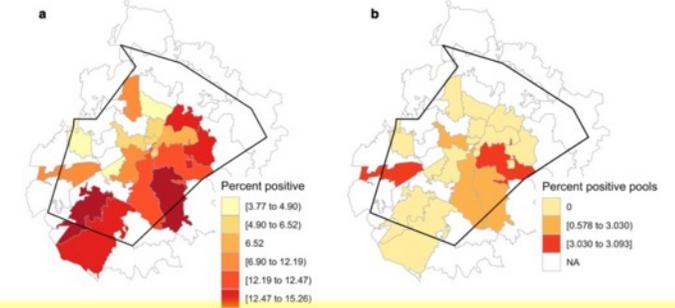
**A. Saltwater crocodile genome**



**B. Norway spruce evolutionary biology**



**C. Stress response in sugarcane**



**E. Mosquito-borne heartworm in dogs**

**A.** Ghosh A et al., Genome Biol Evol. 2020 Jan;12(1) 3635-3646. doi:10.1093/gbe/evz269. PMID: 31821505; © The Author(s) 2020; <http://creativecommons.org/licenses/by/4.0/>

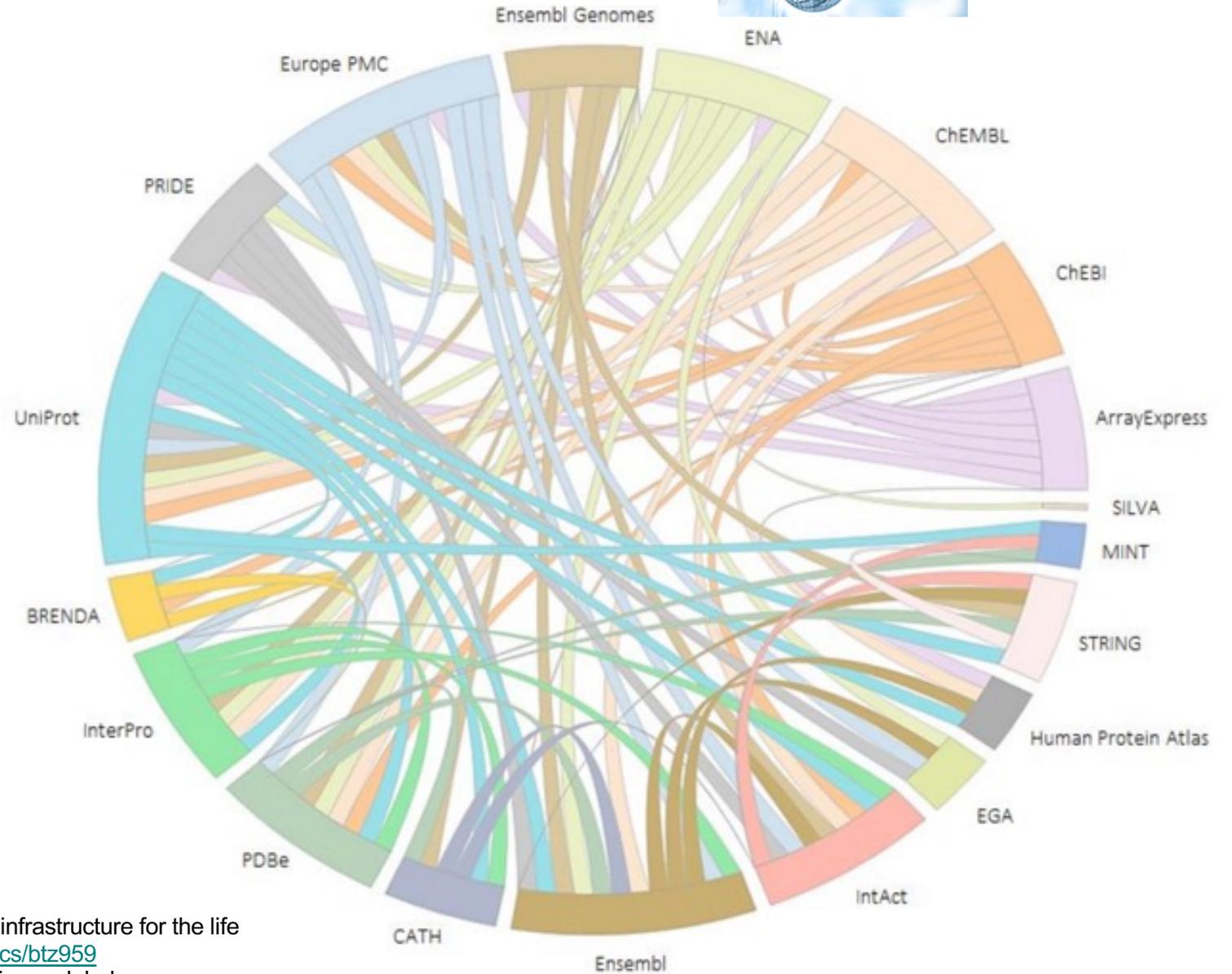
**B.** Sullivan AR et al., Genome Biol Evol. 2020 Jan;12(1) 3586-3598. doi:10.1093/gbe/evz263. PMID: 31774499; © The Author(s) 2020; <http://creativecommons.org/licenses/by-nc/4.0/>

**C.** Su W et al., Sci Rep. 2020 Jan;10(1) 167. doi:10.1038/s41598-019-57058-7. PMID: 31932662; <http://creativecommons.org/licenses/by/4.0/>

**D.** Saeed AM et al., Curr Microbiol. 2020 Jan. doi:10.1007/s00284-020-01874-y. PMID: 31925514.

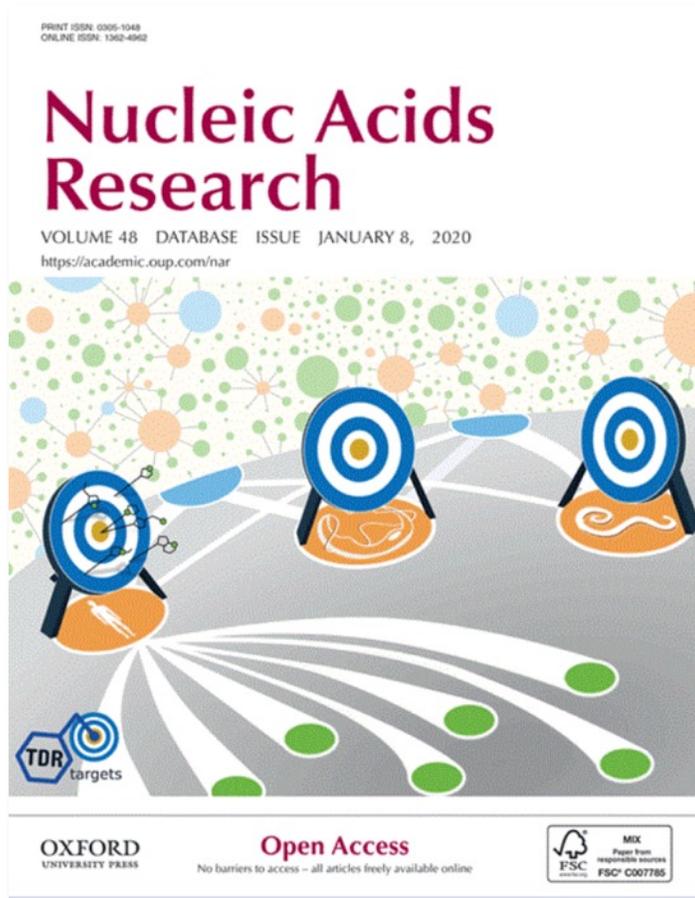
**E.** Spence Beaulieu MR et al., Parasit Vectors. 2020 Jan;13(1) 12. doi:10.1186/s13071-019-3874-0. PMID: 31924253; PMCID: PMC6953185. <http://creativecommons.org/licenses/by/4.0/>; <http://creativecommons.org/publicdomain/zero/1.0/>

# Data reach

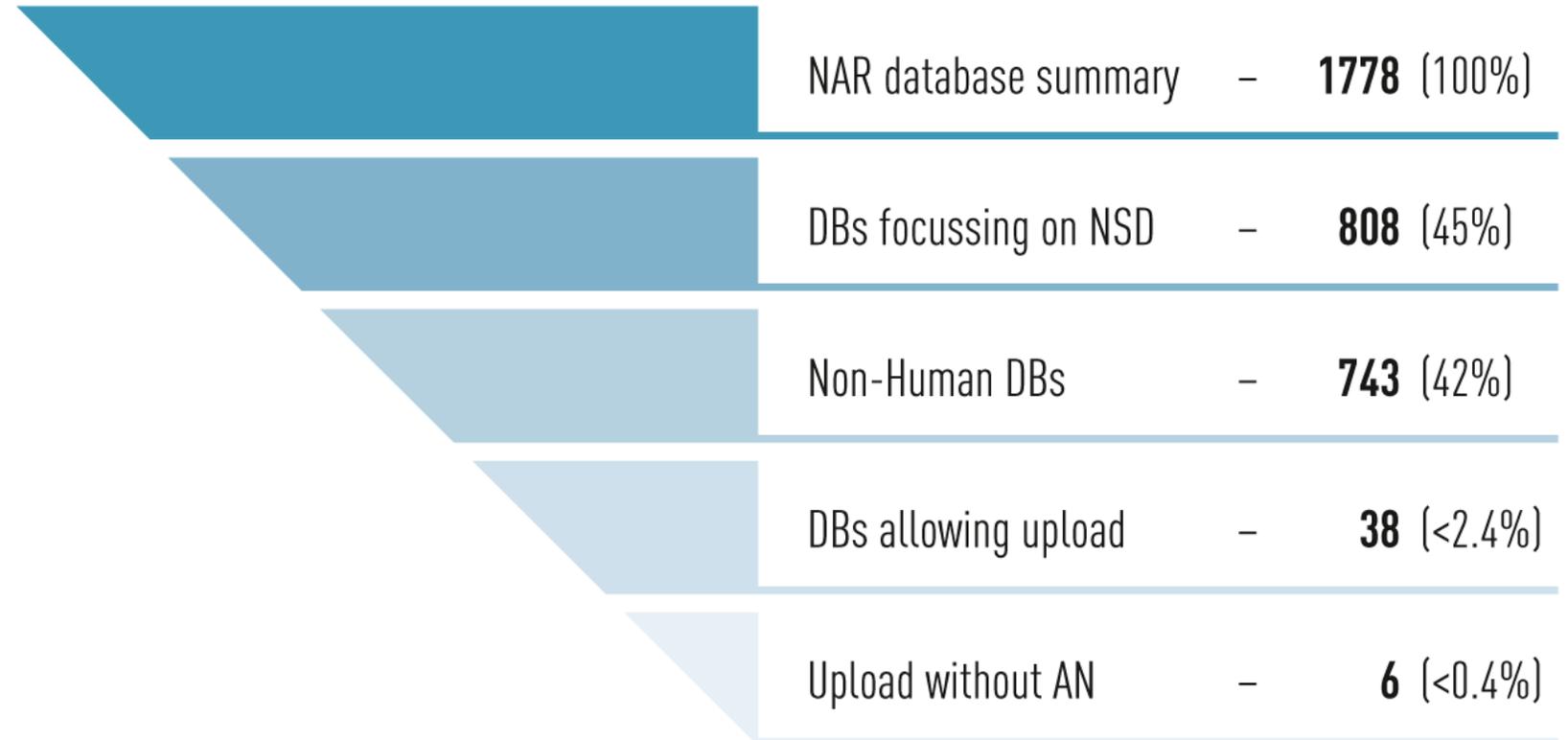


Drysdale *et al.* (2020) The ELIXIR Core Data Resources: fundamental infrastructure for the life sciences. *Bioinformatics*, 2020, 1–7; <http://doi.org/10.1093/bioinformatics/btz959>  
Cook *et al.* (2020) The European Bioinformatics Institute in 2020: building a global infrastructure of interconnected data resources for the life sciences. *Nucleic Acids Research* 48:D17-D23; <http://doi.org/10.1093/nar/gkz1033>

# Further reach



## Public database inventory



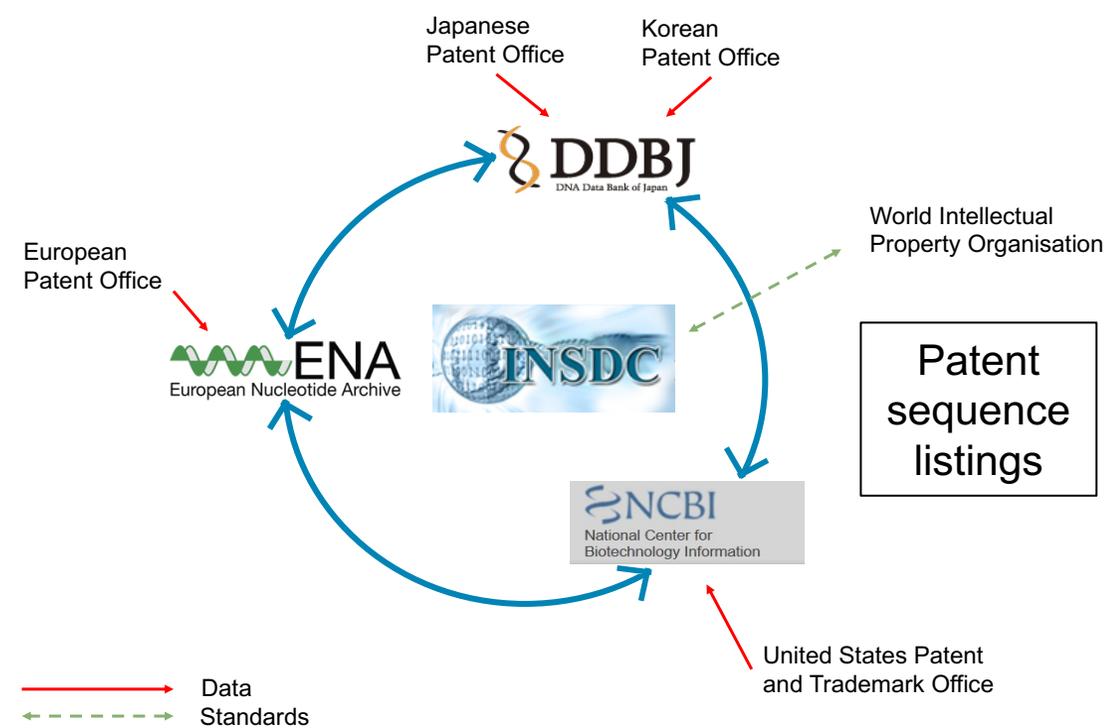
Rohden F, Huang S, Dröge G, Hartman Scholz A, and contributing authors (2019). Combined study in DSI in public and private databases and DSI traceability. <https://www.cbd.int/abs/DSI-peer/Study-Traceability-databases.pdf>

# Terms of use

## Data Services

1. The online data services and databases of EMBL-EBI are generated in part from data contributed by the community who remain the data owners.
2. When you contribute scientific data to a database through our website or other submission tools this information will be released at a time and in a manner consistent with the scientific data and we may store it permanently.
3. EMBL-EBI itself places no additional restrictions on the use or redistribution of the data available via its online services other than those provided by the original data owners.
4. EMBL-EBI does not guarantee the accuracy of any provided data, generated database, software or online service nor the suitability of databases, software and online services for any purpose.
5. The original data may be subject to rights claimed by third parties, including but not limited to, patent, copyright, other intellectual property rights, biodiversity-related access and benefit-sharing rights. For the specific case of the EGA database and human data consented for biomedical research, these rights may be formalised in Data Access Agreements. It is the responsibility of users of EMBL-EBI services to ensure that their exploitation of the data does not infringe any of the rights of such third parties.

<https://www.ebi.ac.uk/about/terms-of-use>



EXAMPLES OF PREPUBLICATION DATA-RELEASE GUIDELINES		
Project type	Prepublication data release recommended	Prepublication data release optional
Genome sequencing	Whole-genome or mRNA sequence(s) of a reference organism or tissue	Sequences from a few loci for cross-species comparisons in a limited number of samples
Polymorphism discovery	Catalogue of variants from genomic and/or transcriptomic samples in one or more populations	Variants in a gene, a gene family or a genomic region in selected pedigrees or populations
Genetic association studies	Genomewide association analysis of thousands of samples	Genotyping of selected gene candidates
Somatic mutation discovery	Catalogue of somatic mutations in exomes or genomes of tumour and non-tumour samples	Somatic mutations of a specific locus or limited set of genomic regions
Microbiome studies	Whole-genome sequence of microbial	Sequencing of target locus in a limited number of samples
RNA profiling		Expression profiles of a gene(s) or pathway(s)
Proteomic studies	Mass spectrometry data sets from large panels of normal and disease tissues	Mass spectrometry data sets from a well-defined and limited set of tissues
Metabolic studies	Catalogue of metabolites in one or more tissues of an organism	Analyses of metabolites induced in a perturbed biological system(s)
RNAi or chemical library screen	Large-scale screen of a cell line or organism analysed for standard phenotypes	Focused screens used to validate a hypothetical gene network
3D-structure elucidation	Large-scale cataloguing of 3D structures of proteins or compounds	3D structure of a synthetic protein or compound elucidated in the context of a focused project

Toronto Principles

ABS-related obligations in existing agreements

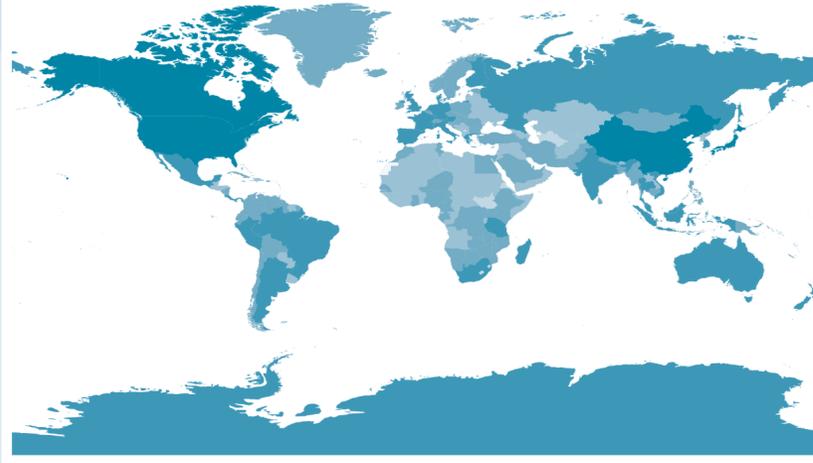
*Nature* volume 461, pages168–170(2009)



# Direction of flow of non-monetary benefit

What is the country of origin for non-human NSD?

1. China	18.23 %
2. United States	17.39 %
3. Canada	9.10 %
4. Japan	7.24 %
5. India	3.46 %
6. Australia	2.66 %
7. Mexico	2.54 %
8. Brazil	2.30 %
9. Germany	1.83 %
10. Spain	1.58 %



Countries of DSI origin

DSI user countries

1. United States	22.69 %
2. China	15.42 %
3. India	6.16 %
4. Japan	3.97 %
5. Germany	3.67 %
6. United Kingdom	3.45 %
7. France	2.84 %
8. Brazil	2.83 %
9. Spain	2.31 %
10. Russian Federation	2.25 %

