**National Genomics Data Center**

# CNCB-NGDC Progress Report

**Yiming Bao**

Director
National Genomics Data Center
Beijing, China

中国科学院北京基因组研究所（国家生物信息中心）
BEIJING INSTITUTE OF GENOMICS  CHINESE ACADEMY OF SCIENCES / CHINA NATIONAL CENTER FOR BIOINFORMATION

# CONTENTS

1 Updates

2 Development

3 Future directions

# The Team

☐ **Steering Advisors**          ☐ **Professors**



**67 students**

**53 Staff**

# The growing of capability



2016

2017

2018

2019

2020

2021

2022

CNCB-NGDC

# GSA Data Growth



**18.7 PB as of 2022-12-09**

CNCB-NGDC

# Grants that support data deposited in GSA



| Total Holdings | Sample Types | Organisms | Platforms | Organizations | Downloads | Journals | Articles | Agencies | Global Visits |

## Grants

Growth of grants

| Agencies | Grants | GSAs | Experiments | Runs |
|---|---|---|---|---|
| Ministry of Science and Technology of the People's Republic of China (MOST) | 729 | 507 | 27685 | 30577 |
| National Natural Science Foundation of China (NSFC) | 4552 | 1773 | 86747 | 93631 |
| Chinese Academy of Sciences (CAS) | 401 | 424 | 24197 | 27787 |
| Others | 5022 | 467 | 37995 | 38729 |

CNCB-NGDC

# GSA Supported Publications

Total Journals: 457

Total Articles: ~2000

# GSA for Human
## Genome Sequence Archive

GSA for Human    |    find a GSA-Human accession    🔍

e.g., HRA000087; HRA000150

Home    Submit    Browse    Search    DAC    Documentation    Policy    👤 Login    🔒 Register

## Genome Sequence Archive for Human

The Genome Sequence Archive for Human (GSA-Human), as a part of GSA in the National Genomics Data Center, is a data repository specialized for human genetic related data derived from biomedical researches. Aside from basic data archive services, GSA-Human features:

- Specializing in human related omics data archives.
- Supplying controlled-access data management services.
- Providing secure online data request services.

### Submit
Submit data to GSA for Human

### Browse
View meta-informations about the released data

### Request Data
Download data after get the access permission

## Data Statistics



## 🌐 Statistics

Studies: 2220

Individuals: 179600

Samples: 308392

Experiments: 352979

Runs: 443714

## ✍ Help & Support

If you have any question or would like to give us any suggestion/comment or report a bug, please feel free to contact us.

**Email:** gsa@big.ac.cn

**QQ group:** 548170081

We highly appreciate your comments and suggestions for further improvements.

## 📁 Latest Released GSA-Human

| Accession | Description |
|---|---|
| HRA003325 (2022-11-09) | SnoRNA-sequencing |
| HRA001912 (2022-11-08) | PAIso-seq2 for human oocyte and embryo |
| HRA001911 (2022-11-08) | PAIso-seq for human oocyte and embryo |
| HRA001288 (2022-11-08) | human oocyte and embryo PAIso-seq |
| HRA001289 (2022-11-08) | human oocyte and embryo PAIso-seq2 |

CNCB-NGDC

# Human Data Backup & Registration Protocol



Backup Center

Registration Center

Researchers

CNCB-NGDC GSA-Human

Backup

Backup ID

Registration ID

Backup ID

Submission

Accession #s

# Genome Warehouse - GWH

## Genome Warehouse

The Genome Warehouse (GWH) is a public repository housing genome-scale data for a wide range of species and delivering a series of web services for genome data submission, storage, release and sharing.

**Submit**
Deposit meta-information into GWH databases

**Download**
Transfer GWH data to your computer

**Browse**
View genome information about the released data

**Documentation**
Find help documents to learn more about GWH

## Statistics

- **Integrated Animals**: 61 genomes, 62 assemblies
- **Integrated Plants**: 77 genomes, 88 assemblies
- **Release of Direct Submissions** (Total: 13348): 2430 Animals; 1206 Plants; 24 Protists; 152 Fungi; 1129 Bacteria; 122 Archaea; 1525 Viruses; 6730 Metagenomes; 30 Others
- Direct Submissions (Total: 27977): 9199 Animals; 3239 Plants; 25 Protists; 164 Fungi; 4983 Bacteria; 130 Archaea; 2988 Viruses; 6746 Metagenomes; 503 Others

## Data Growth



## Virus Resources



*RCoV19*    *MPoxVR*

## GWH-supported Deposition

Data submissions to GWH have been reported by multiple journals, including:

# Multi-omics association knowledgebases



GWAS Atlas
A curated resource of genome-wide variant-trait associations

- 3432 Studies
- 15 Species
- 830 Publications
- 1444 Traits
- 278109 Associations
- 462 Causal variants
- GeneFinder
- LeadSNPFinder
- GWAS Atlas
- MHPlotter
- Submissions
- QQPlotter

*Nucleic Acids Research*, 2020, 2022

EWAS Atlas @ EWAS Open Platform
A knowledgebase of epigenome-wide association studies

- Associations: 642,544
- Traits: 717
- Cohorts: 3,497
- Tissues/Cells: 197
- Studies: 1,586
- Publications: 991

*Nucleic Acids Research*, 2019

TWAS Atlas

Transcriptome-Wide Association Studies (TWAS) Atlas is a curated knowledgebase of transcriptome-wide association studies, integrating trait-associated transcriptome signals from TWAS publications and constructing TWAS knowledge graph to provide reliable and practical resource for researchers.

Please enter 1 or more characters

e.g. Lung Cancer; EFO:0004339; RBM6; ENSG00000004534.14

| Traits | Associations | Genes | Tissues | Publications |
|--------|-------------|-------|---------|--------------|
| 257 | 401,266 | 22,247 | 135 | 200 |

*Nucleic Acids Research*, 2023

# RCoV19

# Machine learning detection of high-risk SARS-CoV-2 variants



**Delta detected about 3 months ahead of the WHO announcement**

# Early-warning of high-risk variants

# Bioinformatics toolkit - BIT

# Sequence Alignment Tools



https://ngdc.cncb.ac.cn/bit/seqaln

# IT Infrastructure

**Public Data Storage**

- **24.5PB** Storage capacity
- **9.6PB Newly added storage**

**Application Clusters**

- **67** Application servers
- **11** Newly added servers

**Backup**

- **15.6PB** Archive backup capacity
- **Tape library system 2 sets**
- **2.4PB Newly added capacity**

**Internet Bandwidth**

- **4Gb** Data center Internet bandwidth
- **The average bandwidth usage is more than 70%**

**HPC**

- **210** Computing Nodes
- **8900** CPU computing cores
- **269** TFlops CPUs' computational ability
- **82** TFlops GPUs' computational ability
- **8 Newly added 512GB memory Nodes**

**HP Storage**

- **8.2PB** Storage capacity

# CONTENTS

1 **Updates**

2 **Development**

3 **Future directions**

# Monkeypox Virus Resource

# TWAS Atlas

## a curated knowledgebase of transcriptome-wide association studies



*Nucleic Acids Research* (2023)

# ASCancer Atlas



**Key features**

- A collection of 2,006 experimentally validated cancer-associated splicing events(CASE)

- Each CASE records the complete information of upstream splicing regulators, splicing event annotations, downstream oncogenic effects and potential treatment targets

- The most complete splicing regulatory network so far

- An interactive splicing visualization tool and a suit of multi-dimensional online splicing analysis tools

*Nucleic Acids Research* (2023)

CNCB-NGDC

**https://ngdc.cncb.ac.cn/ascancer/**

Brain Catalog: a One-Stop Shop for Brain-related Traits

# CCAS

# ProPan

## a comprehensive database for profiling prokaryotic pan-genome dynamics



*Nucleic Acids Research* **(2023)**

# Homologous Gene Database - HGD



*Nucleic Acids Research* (2023)

**Features:**

- Integrated **112,383,644** homologous pairs from multiple resources

- Homologous genes with various annotations

  - **16,909** homologs with traits

  - **276,607** homologs with variants

  - **398,573** homologs with expression

  - **536,852** homologs with Gene Ontology

- Support a comparison function of homologs across multiple species

# Cell Taxonomy



*Nucleic Acids Research* **(2023)**

# GSA's Application for the GCBR-Selection (the Global Core Biodata Resource Selection)

- Initiated by the Global Biodata Coalition, aiming to define Global Core Biodata Resources across biological, life science, and biomedical data resources (biodata resources) worldwide.

- GSA had passed the first round of the selection, and was invited to submit the full application.

- Waiting for the final result on 12th, December, 2022.

## KEY DATES

**21st March 2022** Submission of expressions of interest opened.
**22nd April 2022** Deadline for submission of expressions of interest.
**8th August 2022** Deadline for Submission of full applications.
**Week commencing 12th December 2022** Announcement of the initial GCBR list.

GLOBAL BIODATA COALITION

GLOBAL CORE BIODATA RESOURCE SELECTION

https://globalbiodata.org/scientific-activities/gcbr-selection/

# INSDC meeting

# NGDC – EBI Discussions



**Database Commons**

**Search Engine**

# NCBI Curation Training – 7X

# INSDC Data Mirroring

# BioSample, BioProject and SRA Data



Metadata information has been updated to **2022-11-08**

The data files have been downloaded every day since **2022-04-20**

Data Files： **1.6 PB**

# Integration of SRA Data

# GenBase



- **GenBank/ENA** type database, **GenBank Release 250.0** has been integrated

- **In total**: **592,226** Species, **236,128,977** Nucleotides, **185,121,255** Proteins

- **Direct submission**: **1300** Nucleotides, **1073** Proteins

CNCB-NGDC

**https://ngdc.cncb.ac.cn/genbase/**

# Submission Portal

# Source  Modifiers Table



1. https://ngdc.cncb.ac.cn/genbase/standards
2. https://www.ncbi.nlm.nih.gov/WebSub/html/help/genbank-source-table.html#modifiers

CNCB-NGDC

# Excel structure



➤ Guidelines describe simple help messages, both English and Chinese messages are provided

➤ Color Code mark each modifiers, and classify them to required, optional (recommended), optional

➤ Each cell is identified by Line ID and Column number, so that users can focus on them faster

# Validators

**Structure Validators**

Validate table structure for parsing data correctly.
- Line ID validator: Line ID in the first column is required, if not , fatal error will be thrown.

**Column Validators**

Validate data by column.
- Column_unique_value_validator:  columns with duplicate values are not allowed

**Row Validators**

Validate data by row.
- Row_co_occur_validator: check whether specified values are co-occurred
- Row_mutex_validator:  specified values are mutually exclusive

**Cell Validators**

Same validators are applied to all cells in one column.  If errors are detected, the cell will be mark red, and comment will be add.

- Control vocabulary  validators
- Format validators

**(sourceModifiersValidator.py)**

# Validation Example



*Processed by* **sourceModifiersValidator.py**

# Feature Table

## Three optional formats for submitting sequence features

### Plain text tbl or gff3: >10 sequences

### Excel tbl (with hints ): <10 sequences

# Feature table with INSDC standards



We embed the INSDC standard into excel, and assist users to fill in the qualifier value by prompting for *formats*, *definitions* and *examples*.

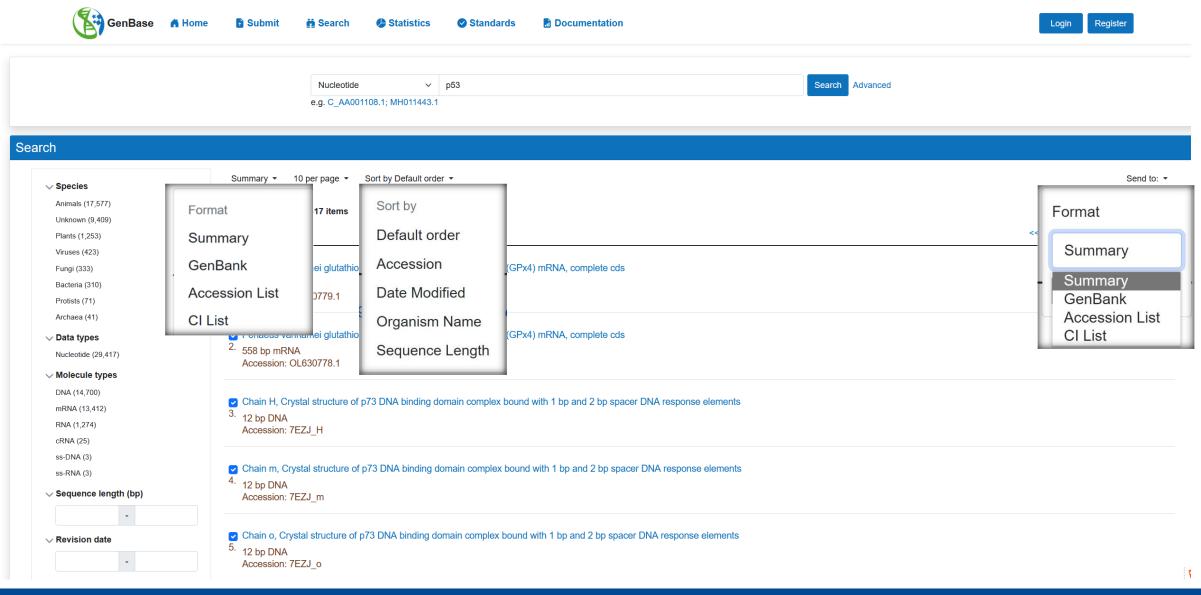Different features correspond to their respective qualifiers, and also correspond to different hints.

We will feed back to user in real time the data that does not conform to the given format, and mark it in red.

User only need to re-submit the corrected excel again until the verification passes. Formatted tbl format will then be used for subsequent steps.

# Search and Download

# Display of Flat Files



GenBank ▾

Serratia marcescens SM39 DNA, complete genome.

**GenBank:** AP013063.1

```
LOCUS       AP013063              5225577 bp    DNA     circular BCT 27-JAN-2017
DEFINITION  Serratia marcescens SM39 DNA, complete genome.
ACCESSION   AP013063
VERSION     AP013063.1
DBLINK      BioProject: PRJDB1121
            BioSample: SAMD00061009
KEYWORDS    .
SOURCE      Serratia marcescens SM39
  ORGANISM  Serratia marcescens SM39
            Bacteria; Proteobacteria; Gammaproteobacteria; Enterobacterales;
            Yersiniaceae; Serratia.
REFERENCE   1
  AUTHORS   Iguchi,A., Nagaya,Y., Pradel,E., Ooka,T., Ogura,Y., Katsura,K.,
            Kurokawa,K., Oshima,K., Hattori,M., Parkhill,J., Sebaihia,M.,
            Coulthurst,S.J., Gotoh,N., Thomson,N.R., Ewbank,J.J. and Hayashi,T.
  TITLE     Genome evolution and plasticity of Serratia marcescens, an
            important multidrug-resistant nosocomial pathogen
  JOURNAL   Genome Biol Evol 6 (8), 2096-2110 (2014)
   PUBMED   25070509
   REMARK   DOI:10.1093/gbe/evu160
REFERENCE   2  (bases 1 to 5225577)
  AUTHORS   Hayashi,T., Iguchi,A. and Ogura,Y.
  TITLE     Direct Submission
  JOURNAL   Submitted (27-MAY-2013) Contact:Tetsuya Hayashi Kyushu University,
            Department of Bacteriology, Faculty of Medical Sciences; 3-1-1
            Maedashi, Higashi-ku, Fukuoka 812-8582, Japan
COMMENT     ##Genome-Assembly-Data-START##
            Assembly Method       :: phred/phrap/consed package
            Genome Coverage       :: 12.6x
            Sequencing Technology :: Sanger
            ##Genome-Assembly-Data-END##
FEATURES             Location/Qualifiers
     source          1..5225577
                     /organism="Serratia marcescens SM39"
                     /mol_type="genomic DNA"
                     /strain="SM39"
```

GBFF ▾                                                          Send to: ▾

Ipomoea batatas ocimene synthase 1 (OS1) mRNA, complete cds.

**GenBase:** C_AA001108.1

FASTA

```
LOCUS       C_AA001108              1644 bp    mRNA    linear   PLN 31-AUG-2022
DEFINITION  Ipomoea batatas ocimene synthase 1 (OS1) mRNA, complete cds.
ACCESSION   C_AA001108
VERSION     C_AA001108.1
KEYWORDS    .
SOURCE      Ipomoea batatas (sweet potato)
  ORGANISM  Ipomoea batatas
            Eukaryota; Viridiplantae; Streptophyta; Streptophytina; Embryophyta;
            Tracheophyta; Euphyllophyta; Spermatophyta; Magnoliopsida;
            Mesangiospermae; eudicotyledons; Gunneridae; Pentapetalae; asterids;
            lamiids; Solanales; Convolvulaceae; Ipomoeeae; Ipomoea.
REFERENCE   1  (bases 1 to 1644)
  AUTHORS   Xiao,Y., Qian,J., Hou,X., Zeng,L., Liu,X., Mei,G. and Liao,Y.
  JOURNAL   Unpublished
REFERENCE   2  (bases 1 to 1644)
  AUTHORS   Xiao,Y., Qian,J., Hou,X., Zeng,L., Liu,X., Mei,G. and Liao,Y.
  TITLE     Direct Submission
  JOURNAL   Submitted (31-AUG-2022) Center of Economic Botany, Core Botanical
            Gardens, South China Botanical Garden, Chinese Academy of Sciences,
            Changxing, Guangzhou 510630, China
FEATURES             Location/Qualifiers
     source          1..1644
                     /organism="Ipomoea batatas"
                     /mol_type="mRNA"
```

# Grants Awarded for International Collaboration

| Funding Agency | Project Title | Collaborators | Amount |
|---|---|---|---|
| IUBS | Open Biodiversity and Health Big Data | Multiple countries | Euro 10,000 |
| CAS | Global Genomics Data Sharing | USA | RMB 800,000 |
| ANSO | Precision warning method for high-risk variants of emerging infectious diseases | Brazil, France, Pakistan | RMB 1,300,000 |
| ANSO | Whole genome sequencing and miRNA biomarkers for an enhanced understanding of mechanism of tuberculosis infection in cynomolgus macaques (Macaca fascicularis): A translational knowledge to clinical study | Thailand, USA | US$ 150,000 |

# CONTENTS

1 **Updates**

2 **Development**

3 **Future directions**

# Future Directions

- Progress of CNCB infrastructure

- Collaboration with DDBJ and KOBIC

- Global Core Biodata Resource

- Partnership with INSDC

- Talent recruitment

**National Genomics Data Center**

# Acknowledgement

**NGDC Members** https://ngdc.cncb.ac.cn/people

CNCB-NGDC