

Bridging Innovation and Policy: Leveraging Emerging Data for Healthcare Decisions

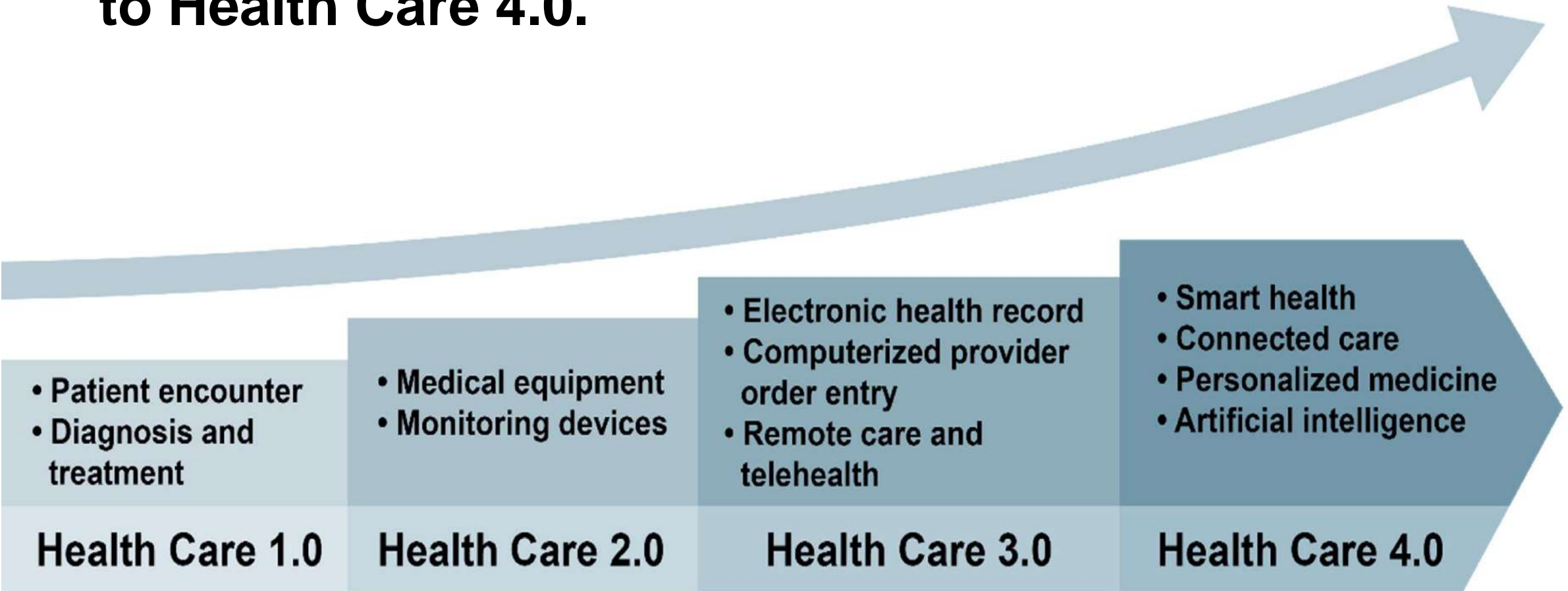
Hyun-Young Park, MD, PhD

National Institute of Health, Republic of Korea

Contents

- Healthcare big data
- Role of healthcare big data for outbreak control
- NHIS and EMR data integration
- Genomic data for precision medicine
- NIH's Bio Big Data Initiative

Historical evolution of health care 1.0 to Health Care 4.0.

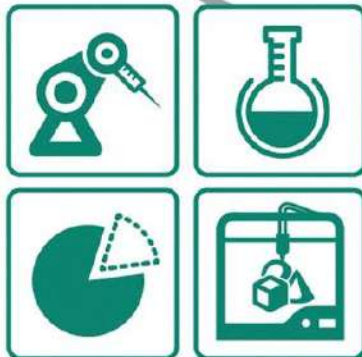


Healthcare 4.0, *Smart Medicine*

The new brain & new hands in Healthcare 4.0

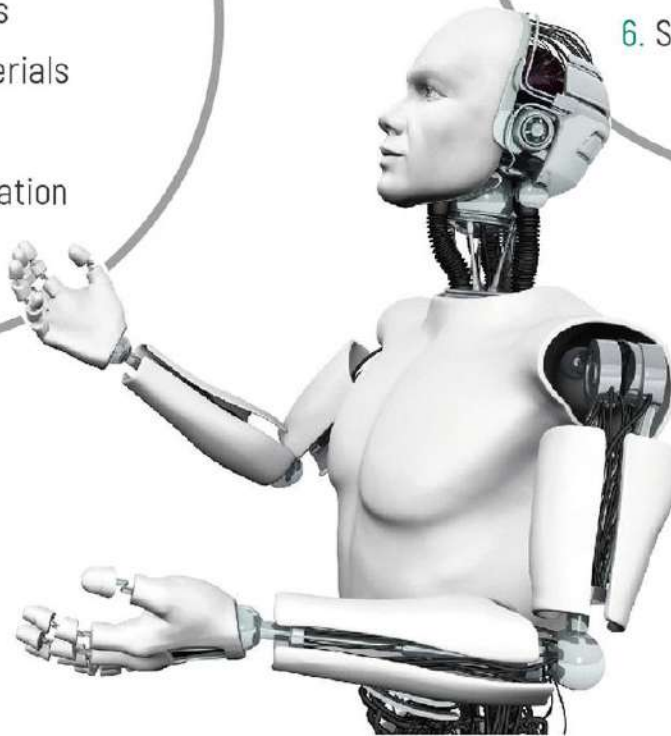
New hands

1. Robot
2. Mini-laboratory
3. Wearable devices
4. Customized materials
5. 3D printing
6. Speed & minimization



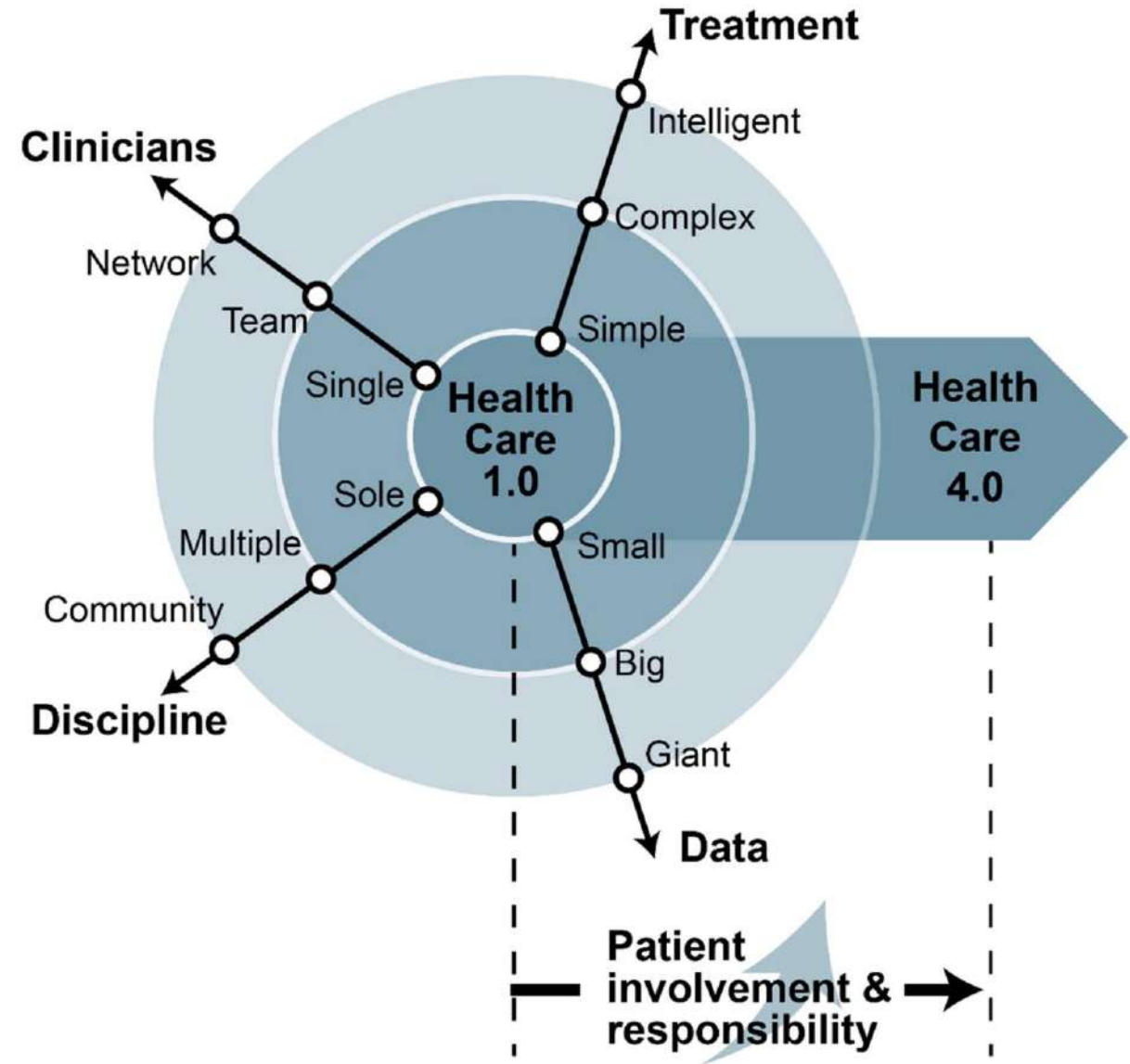
New brain

1. Precision medicine
2. Artificial intelligence
3. Big data
4. Internet of things
5. Telemedicine
6. Shared decision making

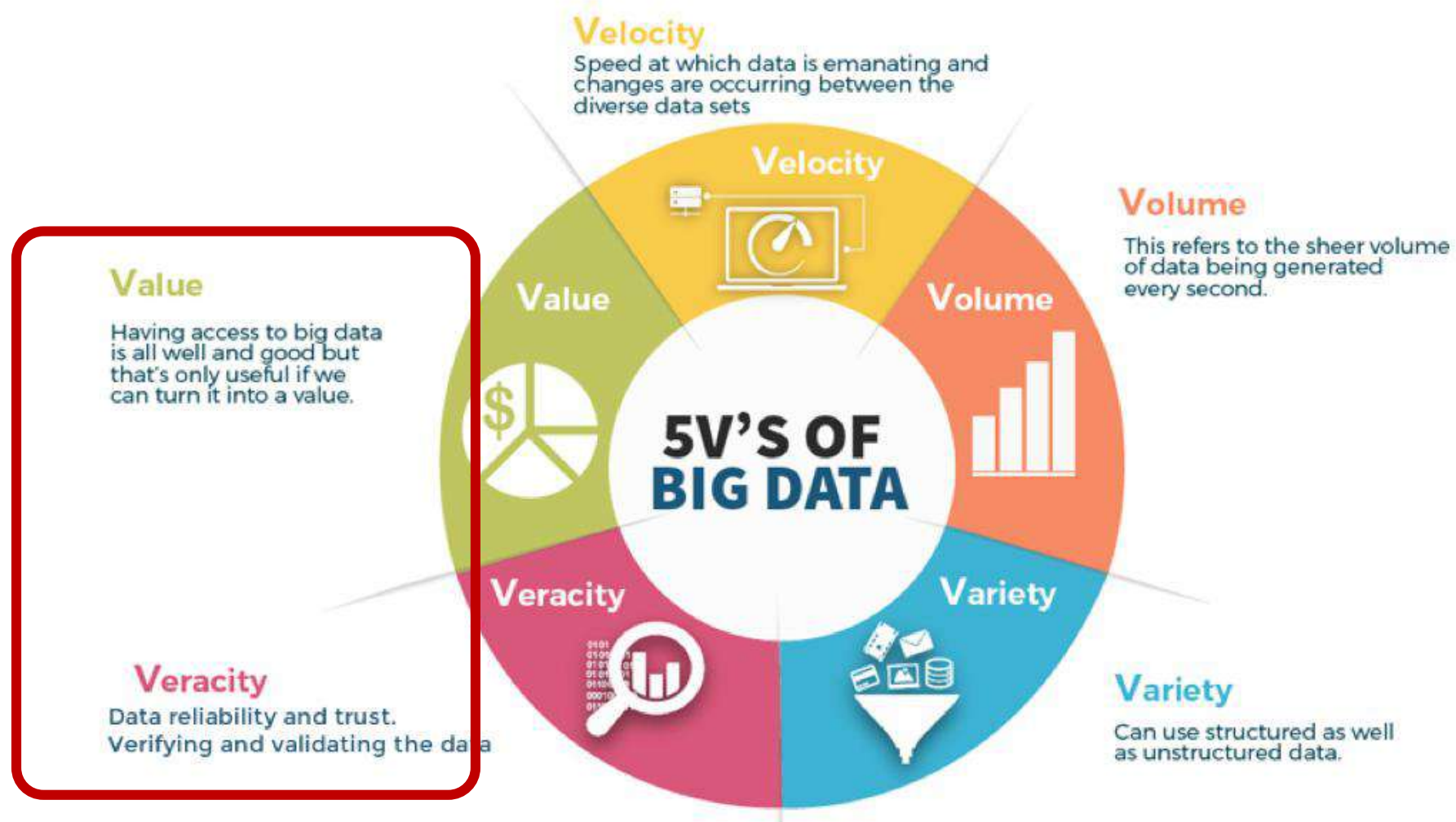


Characteristics

- **Big and giant data**
streams with large variations in dimension, quality, format, and characteristic
- Patients and clinicians are **increasingly involved and share responsibilities** for monitoring their health, reporting symptoms, and participating in shared decision making for treatment and care planning



Big Data in Healthcare



Current Healthcare Big Data

- Electronic Health Records (EHRs): Diagnoses, lab results, prescriptions, vital signs, clinical notes
- Medical Imaging Data
- Claims and Administrative Data: Insurance billing, procedures, medications, healthcare utilization
- Health Screening
- Social Determinants of Health (SDOH): Income, education, employment, environment

What is Emerging Healthcare Data?

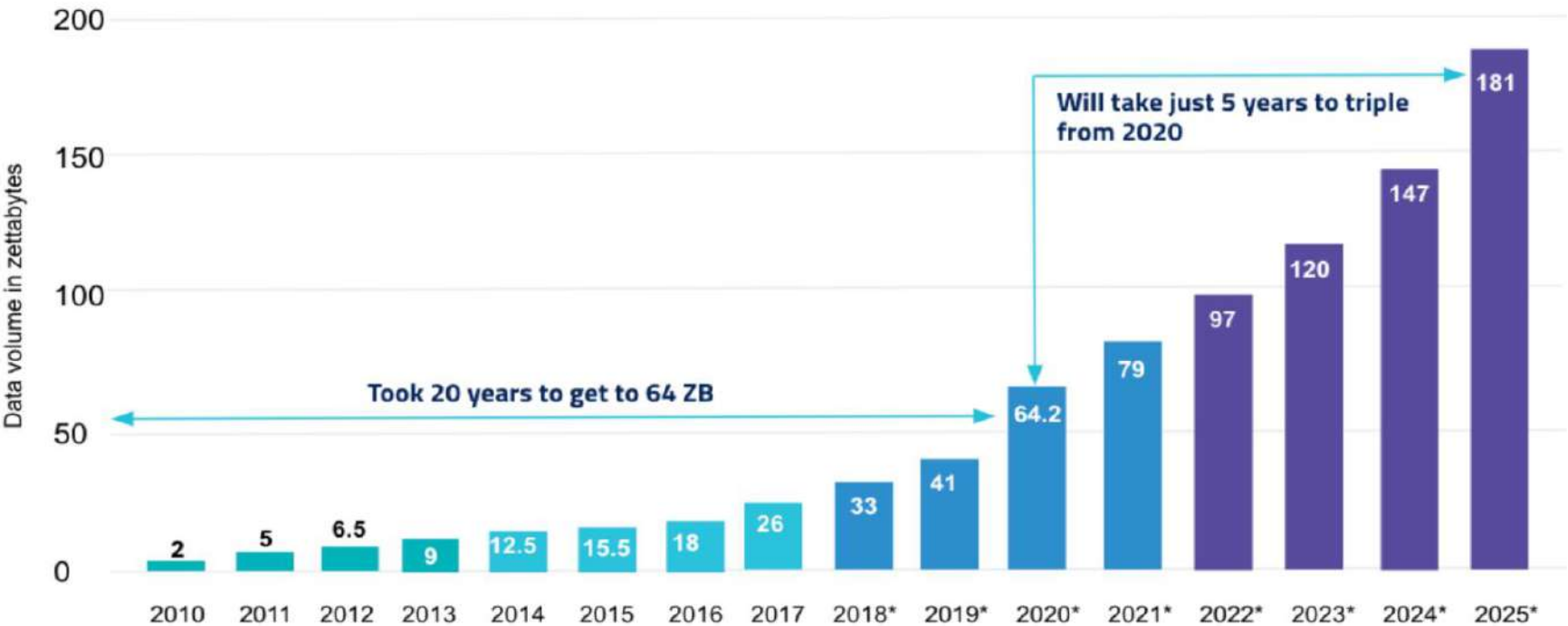
- Genomic data (e.g., whole-genome sequencing)
- Multi-omics data generated by research
- Real-world evidence (RWE) from clinical practice and patient registries
- Social and behavioral health data
- Wearable/IoT health data
- Patient-Reported Outcomes (PROs)
- Artificial intelligence-generated insights

Big Data Grows Ever Bigger

1 ZETTABYTE =
**1 000 000 000
000 000 000 000**
BYTES

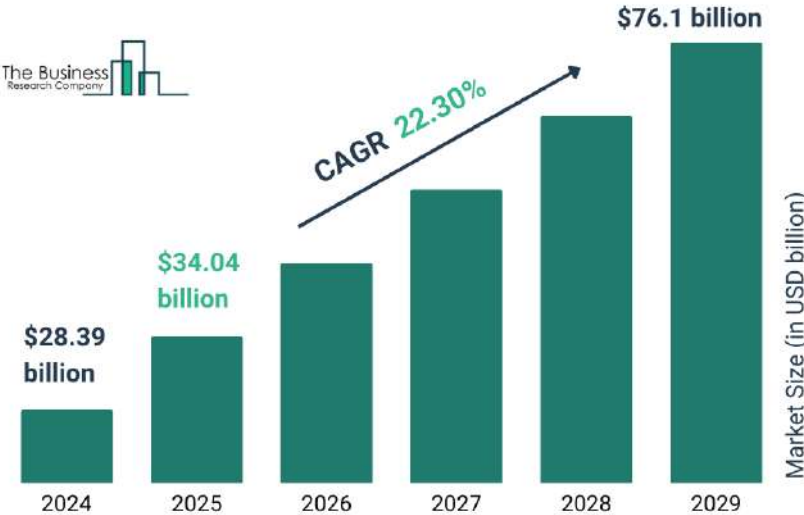
©ComputerHope.com

Volume of data/information created,
captured, copied, and consumed worldwide (Zettabytes)



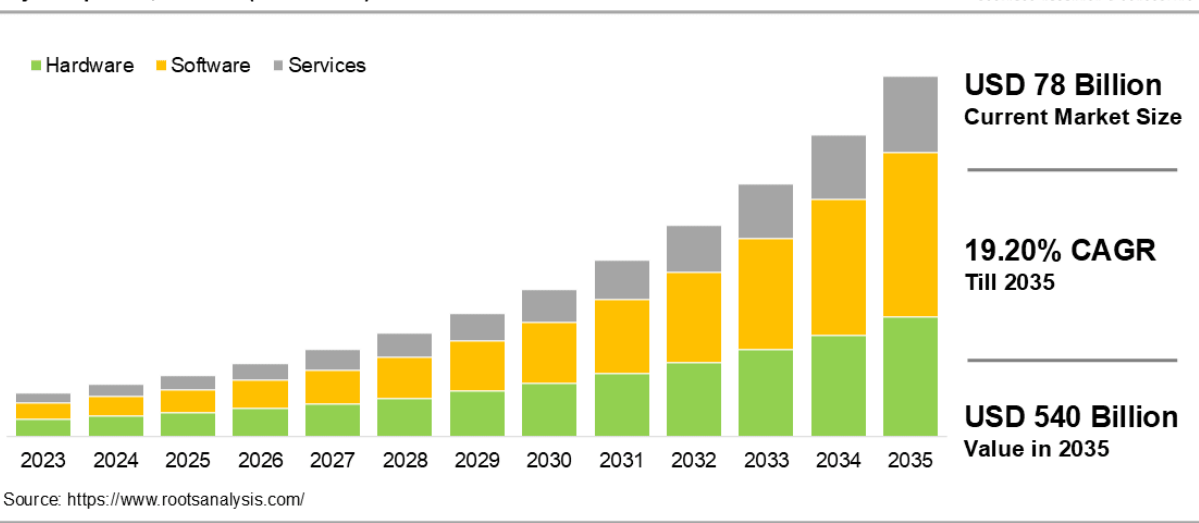
35%
of all data will
be life sciences
+ healthcare by
2025

Big Data Healthcare Global Market Report
2025

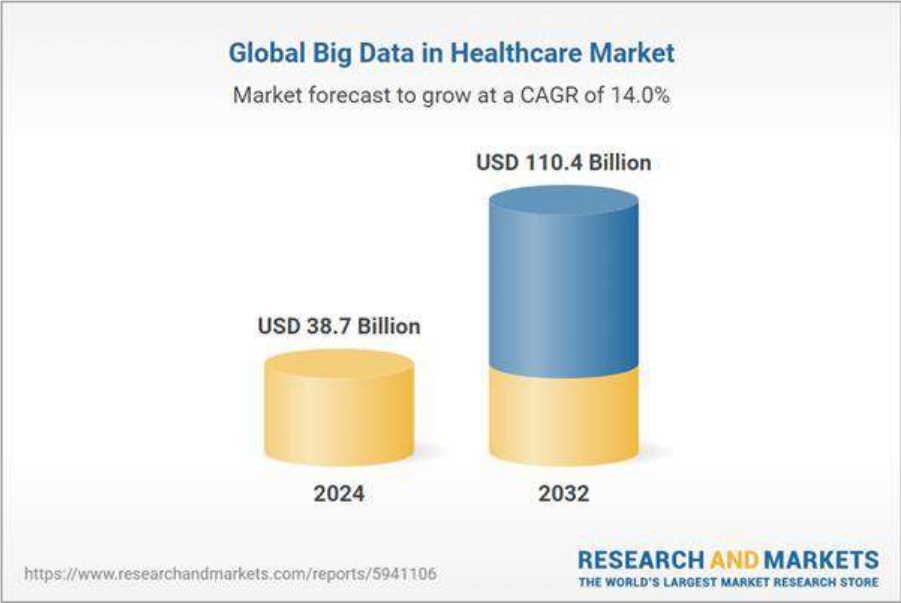


The Era Where Data
Becomes Economic Value

Big Data in Healthcare Market
By Component, Till 2035 (USD Billion)



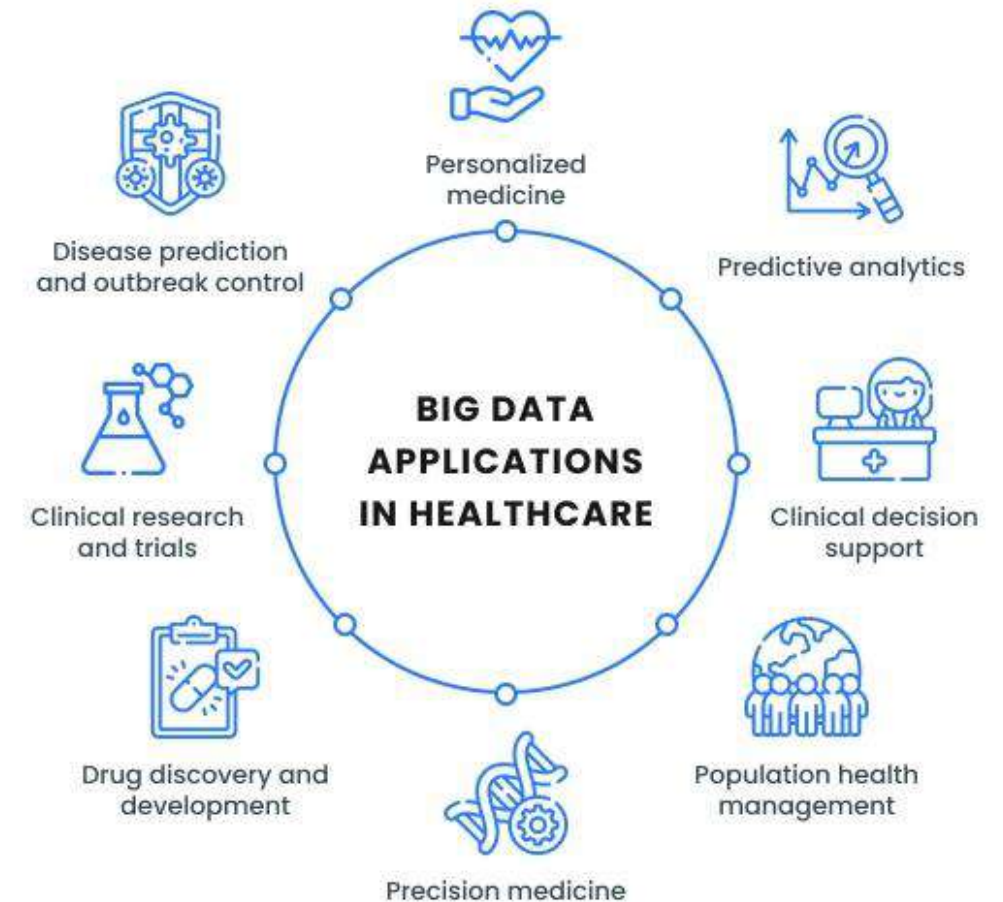
Global Big Data in Healthcare Market
Market forecast to grow at a CAGR of 14.0%



Healthcare Big Data

Opportunities and Impact

- **Clinical decision support**
- **Precision medicine:** Tailored therapies based on genomic and lifestyle data
- **Population health management:** Predictive analytics for early intervention
- **Policy innovation:** Data-driven design of reimbursement, screening, and public health programs
- **Drug discovery and clinical trials:** Accelerated recruitment and targeted trials



Role of Healthcare Big Data for Outbreak Control

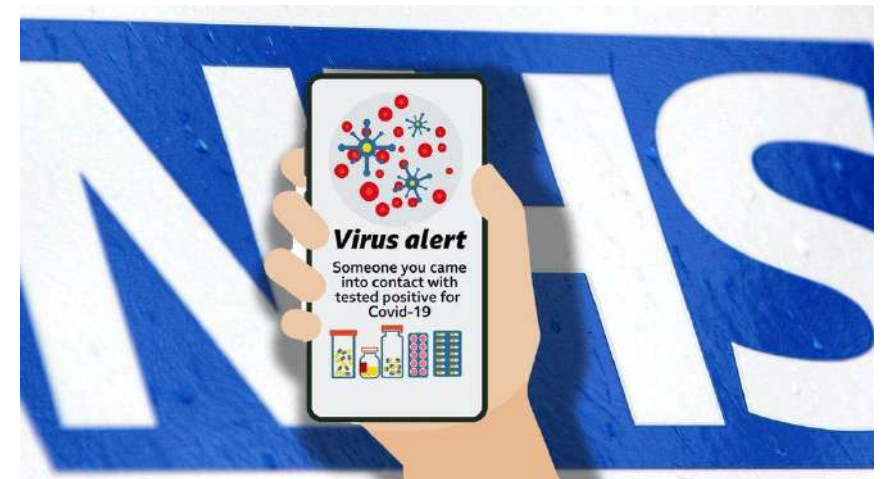
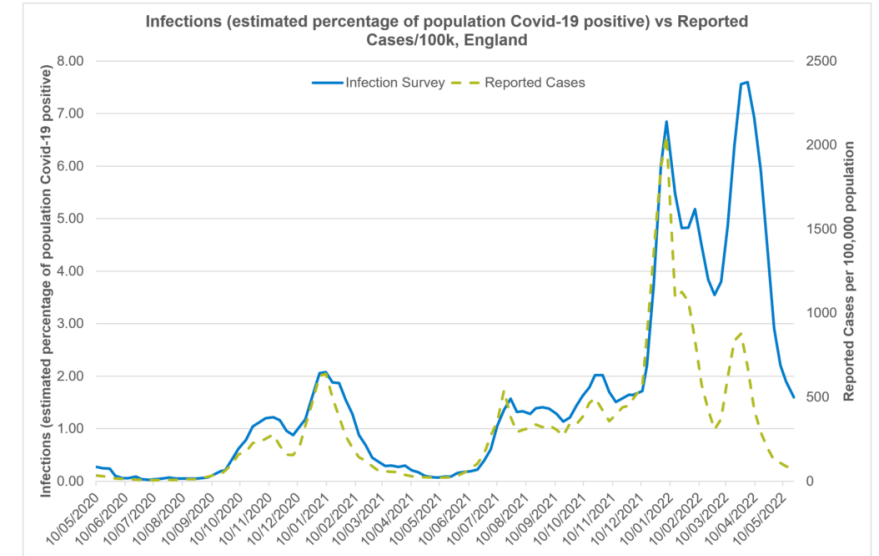
COVID-19 Infection Survey (ONS)

National Health Service (NHS) and Office for National Statistics (ONS) to monitor infections, hospitalizations, and deaths

- Regular testing of random household samples, including asymptomatic individuals
- Used big data analytics to estimate: National and regional infection rates, Infection prevalence by age group, region, and time

Impact:

- Provided real-time community transmission insights
- Guided local lockdown decisions and vaccine rollout
- One of the most accurate surveillance tools globally





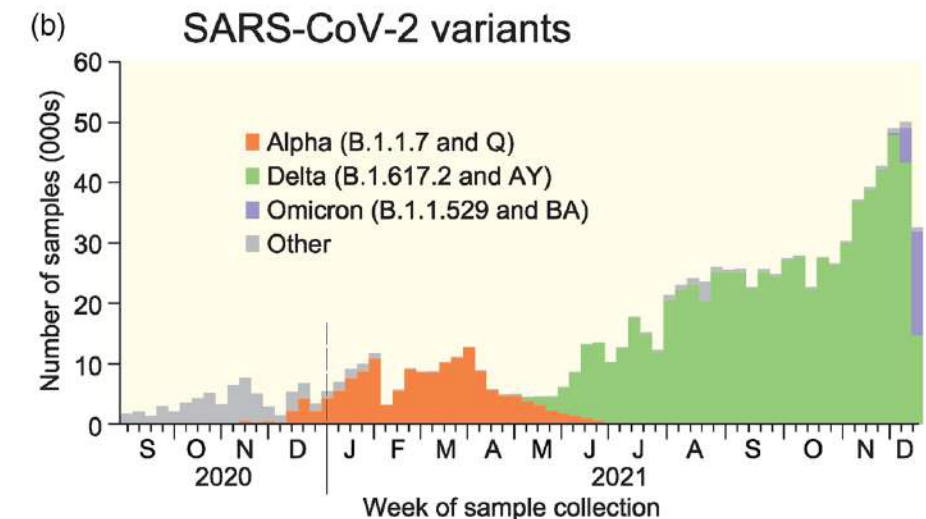
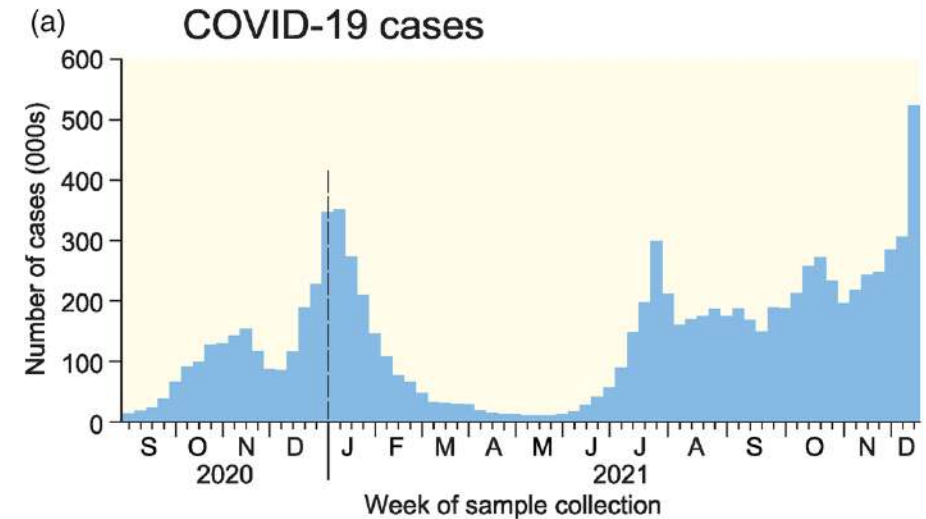
**COVID-19
GENOMICS
UK CONSORTIUM**

COVID-19 Genomics UK (COG-UK) Consortium established in March 2020

- Sequenced over 2 million SARS-CoV-2 genomes
- Linked viral genome data with patient demographics and clinical outcomes

Impact:

- Detected and tracked variants of concern (e.g., Alpha, Delta, Omicron)
- Informed public health strategies and vaccine updates
- A model for integrated pathogen surveillance globally



- Secure, open-source analytics platform created in March 2020
- Designed to enable secure, large-scale analysis of electronic health records (EHRs) for COVID-19 research while maintaining strict privacy protection
- Data: Primary care records (diagnoses, medications, demographics), Linked to hospital data, death registrations, and COVID-19 test results
- Impact
 - 1) COVID-19 Risk Factors identification
 - 2) Vaccine Safety and Effectiveness: Monitored vaccine uptake and safety in real time, Evaluated rare adverse events
 - 3) Health Inequalities: disparities in COVID-19 mortality among ethnic minorities, Access to care and shielding recommendations, Geographic variation in vaccine uptake

Big Data and Public Health Crisis

Integrated big data and genomics can turn a reactive public health response into a proactive and precision-driven strategy.

KDCA Covid-19 NHIS cohort



KDCA Covid-19 NHIS cohort

Data sharing

K-COV-N sharing website (<http://nhiss.nhis.or.kr>)

The screenshot shows the NHISS National Health Insurance Data Sharing Service website. At the top, there is a navigation bar with links for 'Introduction', 'Research DB', and 'Index of Medical utilization', along with a 'KOREAN' language toggle. The main header features the NHISS logo and the text 'National Health Insurance Data Sharing Service'. Below this, a large banner reads 'We lead in evidence-based health and medical policy and support of academic research'. The main content area is divided into three columns: 'Research DB' (National health insurance data, Sample Research DB, Customized Research DB), 'Index of Medical utilization' (Monitoring index of disease Provisional guide), and 'NHIS' (National Health Insurance Sharing Service). Each column has a corresponding icon and a button labeled 'PROVISIONAL GUIDE >' or 'INTRODUCTION >'.

Key achievements using K-COV-N

This block shows the header of a journal article from Elsevier. It includes the Elsevier logo, the journal title 'Asian Journal of Psychiatry', and the journal homepage URL: www.elsevier.com/locate/ajp. It also mentions that contents lists are available at ScienceDirect.

Short communication

COVID-19 vaccination, incidence, and mortality rates among individuals with mental disorders in South Korea: A nationwide retrospective study

Brief Communication
Infectious Diseases,
Microbiology & Parasitology



COVID-19 Vaccination Rates in Patients With Chronic Medical Conditions: A Nationwide Cross-Sectional Study

RESEARCH

Open Access

Excess mortality during the Coronavirus disease pandemic in Korea



- 224 studies were approved, 36 papers published (April 2022~2025)



NHIS Data

Korea's National Health Insurance Service (NHIS) data serves as a representative example of how large-scale health data can inform and shape healthcare policy decisions.

Overview of NHIS Data

- Operated by the **National Health Insurance Service**
- Covers **over 98%** of the Korean population
- Contains longitudinal data: eligibility, healthcare utilization, prescriptions, health screening, death
- Linked to other datasets: national cancer registry, cause-of-death registry, etc.

Example

Article

Effectiveness of the Korean National Cancer Screening Program in Reducing Colorectal Cancer Mortality

Hyeon Ji Lee ¹, Kyeongmin Lee ², Byung Chang Kim ³, Jae Kwan Jun ^{1,2}, Kui Son Choi ^{1,2} and Mina Suh ^{1,2,*} Cancers 2024, 16, 4278.

- Objective: To assess the impact of the Korean National Cancer Screening Program (KNCSPP) on colorectal cancer (CRC) mortality.
- Methodology: A nested case-control study utilizing cohort data from the KNCSPP, encompassing 5,944,540 individuals aged ≥ 50 years as of 2004. The study linked this data with the Korea Central Cancer Registry (KCCR) and death certificate data from Statistics Korea.
- Findings: Individuals who underwent CRC screening using the fecal immunochemical test (FIT) had a **26% lower risk of CRC-specific mortality** compared to those who were never screened (Odds Ratio: 0.74; 95% Confidence Interval: 0.71–0.76). The reduction in mortality was more pronounced with increased frequency of screening.

NHIS Database

Pros and Cons

Strengths

Nationally representative

Longitudinal, large-scale

Real-world utilization and claims data

Useful for RWE and health policy evaluation

Limitations

Limited clinical detail

Potential for miscoding

No patient-reported outcomes or genomics

Access constraints and time lag

*NHIS data is a **powerful tool for public health research and policymaking**, but should be **complemented** with clinical, genomic, and patient-reported data*

Key Obstacles in Linking Healthcare Big Data in Korea

Lack of Interoperability

- Heterogeneous data formats
- Limited adoption of international standards

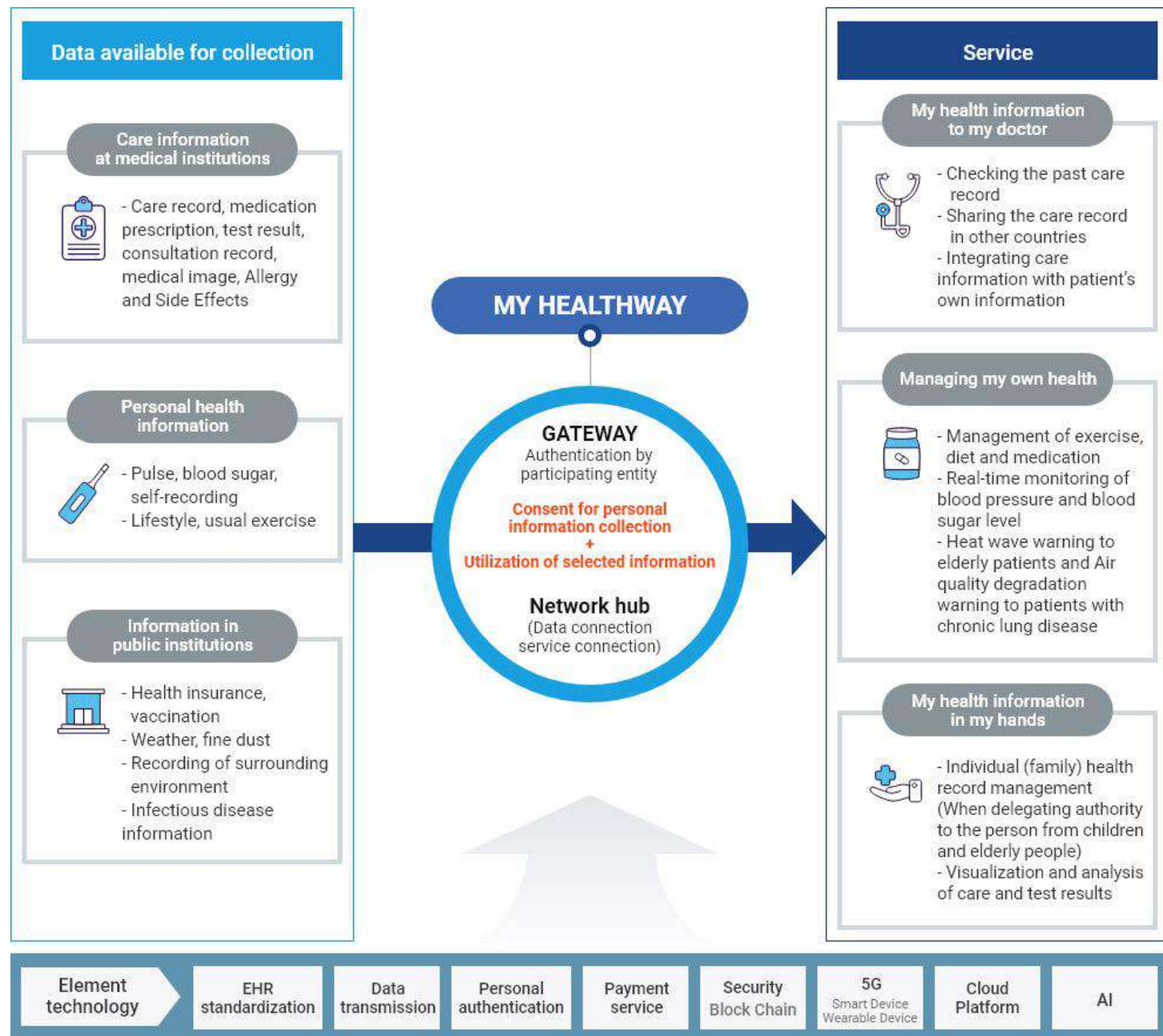
Data Ownership and Governance Issues

- Unclear data ownership between healthcare providers, public institutions, and patients obstructs coordinated data use.
- Absence of a unified governance framework for managing cross-sector data access and accountability.

Privacy and Public Trust Concerns

- High public sensitivity regarding health data misuse or commercial exploitation undermines support for industry-driven uses.
- Lack of **trusted intermediary institutions** to manage anonymization, consent, and ethical oversight weakens social acceptance.

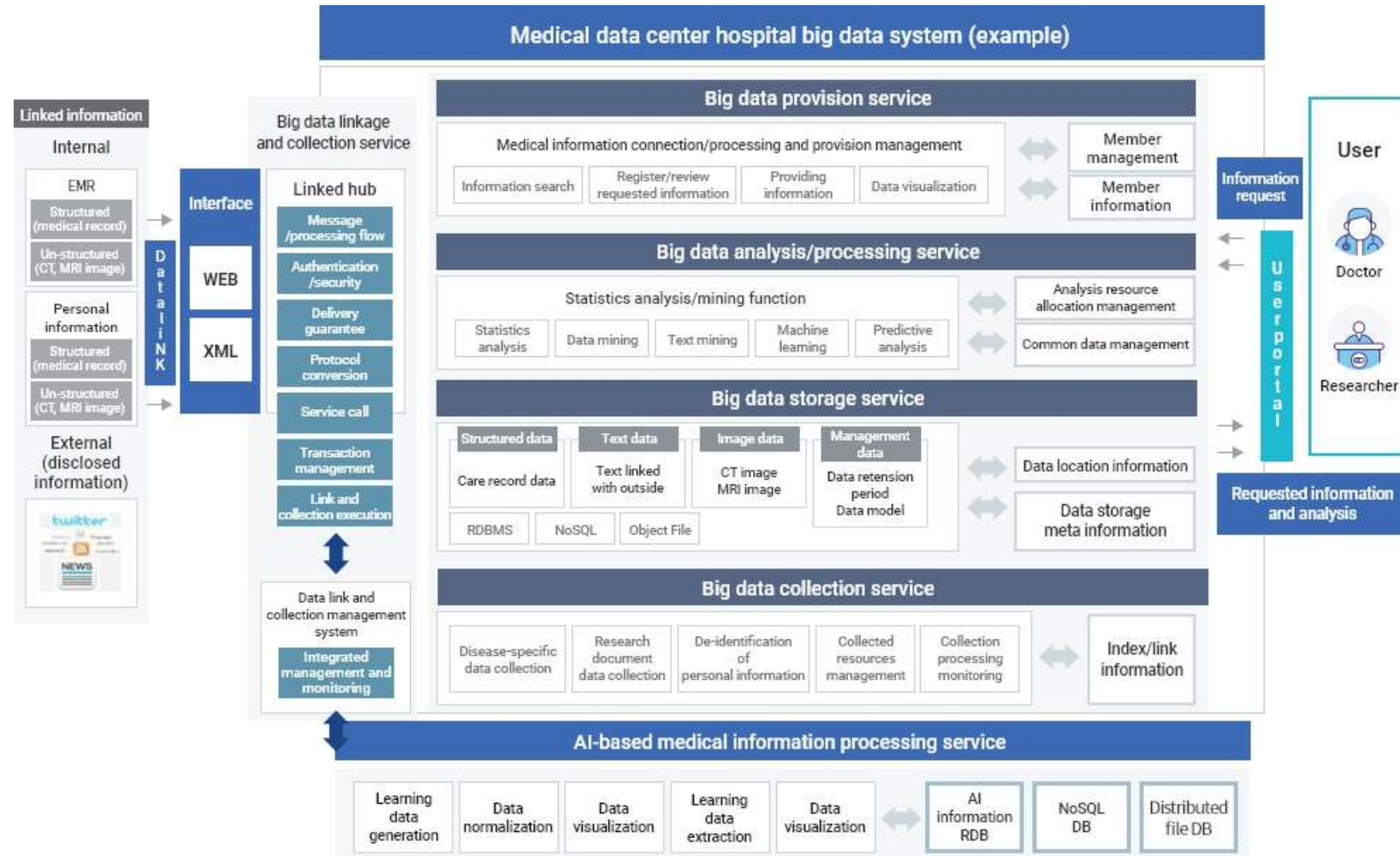
My Healthway



Healthcare data center hospitals

Promoting research on utilization of medical information

- Definition and construction of medical data set for promotion of submitted application scenario research
- Contributing to the development of new medical technologies and promoting research for the development of new drugs, medical devices, and AI by using medical data



Healthcare data center hospitals



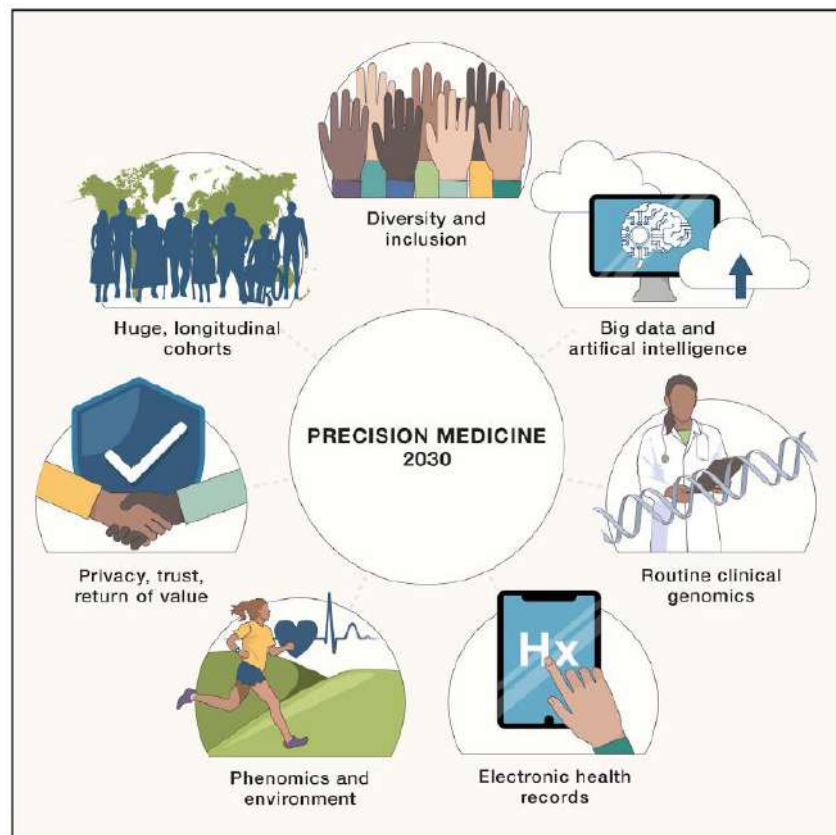


Genomic Data for Precision Medicine

Importance of Genomic Data

- **Precision medicine:** personalized drug selection, pharmacogenomics.
- **Rare diseases:** faster and more accurate diagnoses.
- **Cancer genomics:** tumor profiling, immunotherapy targeting.
- **Public health:** genomic surveillance, population-level risk stratification.

Seven opportunities for precision medicine by 2030

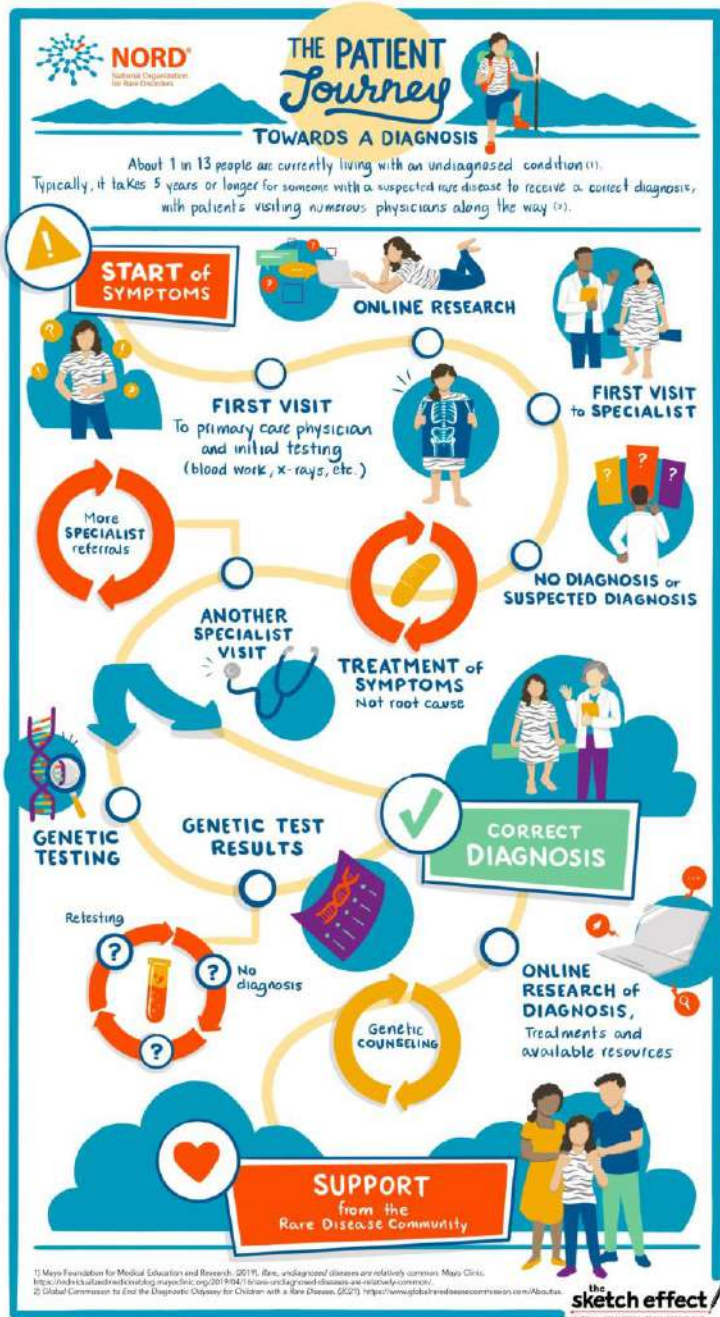


<https://doi.org/10.1016/j.cell.2021.01.015>

Table 1. Envisioning how precision medicine will affect clinical medicine and research in the next decade

	Where we are today	Where we will be in 2030
<i>Clinical applications</i>		
Genomics for disease	Primarily limited to rare disease and select cancers.	Genomics is routine. Genetic causes and targeted therapies are discovered for many “common” diseases. Microbiome measures are routinely included.
Pharmacogenomics (PGx)	Common in cancer and within select applications of older medications at select sites.	Genome-aware EHRs make PGx easy and automatically update rules from central guidelines. New PGx associations discovered from clinical data.
Genomics for healthy individuals	In research, whole-genome sequencing and search for mutations in one of the ACMG59 genes, present in about 3% of people. Variant interpretation is hard.	ACMG59 grows to > 200, variant interpretation improved by huge, diverse sequenced populations. Cell-free DNA becomes a mainstay of cancer screening
EHRs	Episodic capture from healthcare without robust genomics support. EHR data is essentially not portable.	Genome- and device- enabled. Data can be easily moved between EHRs and to participant apps.
Environmental influences on health	Patient-reported habits and exposures	Geocode-based exposure linkage Real time monitoring of multiple environmental exposures Precision nutrition
Wearable sensors	Ad hoc use of activity monitors	Continuous monitoring of physical activity, sleep, metabolic parameters
<i>Research applications</i>		
Population demographics	>80% European ancestry	>50% non-European ancestry
Routinely available data	Surveys of health conditions, lifestyle, behavior, and diet. GWAS data, lab assays, structured EHR data, and geocoded exposure linkages.	Whole genomes, lab assays, surveys, full EHRs, environmental, genomic and sensor data. Includes imaging, narrative, geocoded, and continuous monitoring approaches to clinical care, activity, precision nutrition, and environment.
Size of cohorts used in analysis	Up to 500K, data downloaded and manually harmonized to sets of several million	>100M using cloud-based federated analyses facilitated by common standards
Largest genomic studies performed on a trait	>1M (GWAS)	>50M (GWAS) >2M (WGS)
Cost of a whole genome	\$500	\$20*

*Sequencing costs have often fallen faster than Moore’s law. Using Moore’s law, sequencing costs would be 1/32 of US \$500, or \$15.63.



Rare Diseases

Diagnostic Odyssey

- Multiple specialists
- Diagnostic uncertainty
- Impact on patient quality of life
- Financial burden

Rare Disease Diagnosis in Korea

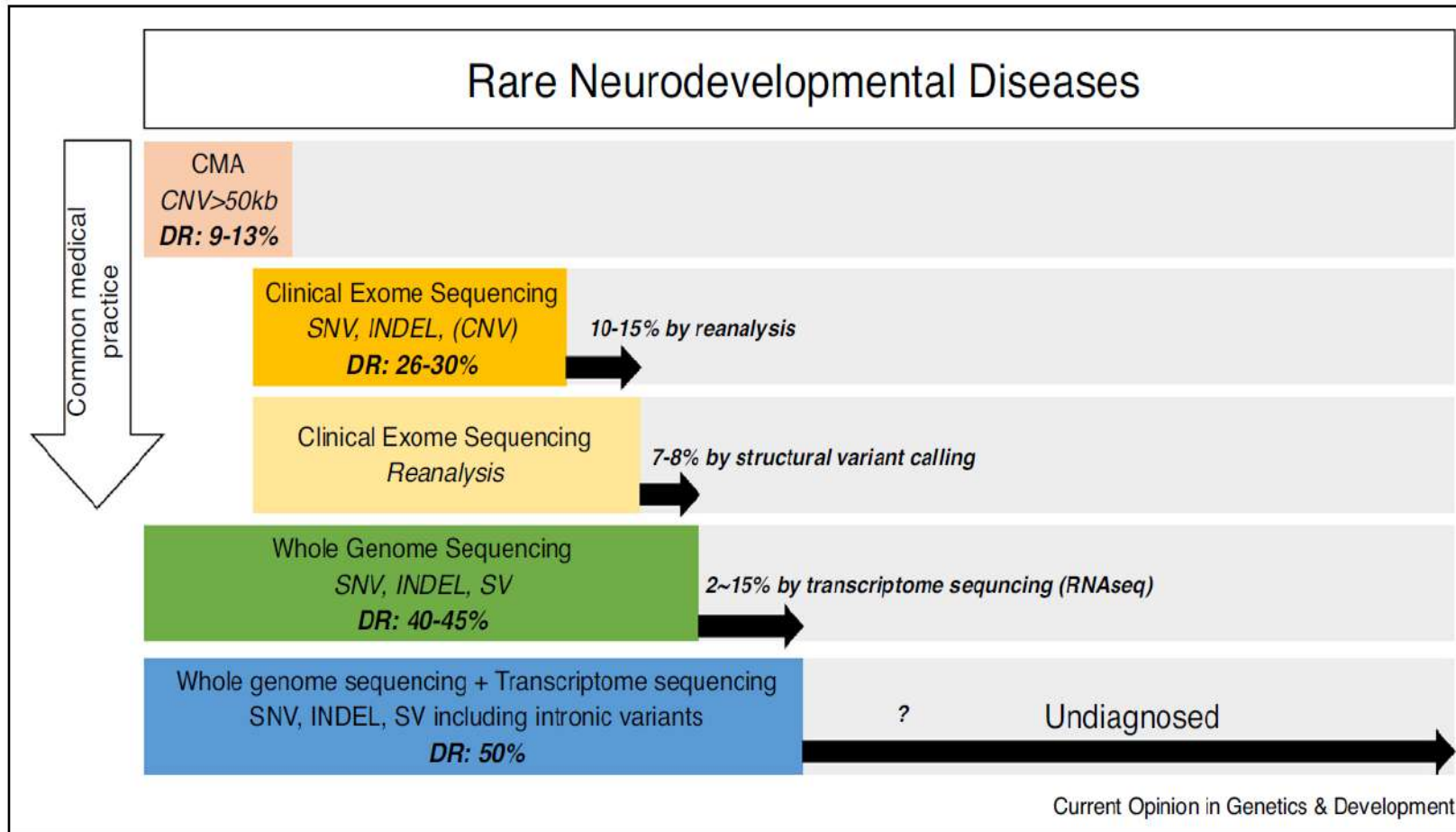
- Patients undergo multiple **single-gene or small panel tests**
- Each test targets specific genes, often requiring **dozens to hundreds** of sequential tests
- Covered by National Health Insurance (NHI),
Patient out-of-pocket cost: ~10%
For panel test, patient have to pay up to 80%
- Access to WGS is mainly available through **government-supported research or diagnostic aid programs.**

Diagnostic Yield of Genetic Test

- **Traditional genetic testing (e.g., microarrays, gene panels):**
Diagnostic yield: **10–25%**
- **Whole Exome Sequencing (WES):**
Diagnostic yield: **25–40%**
- **Whole Genome Sequencing (WGS):**
Diagnostic yield: **40–60%**

Diagnostic Yield of Genetic Test

Genomic Tests for Diagnosing Rare NDD

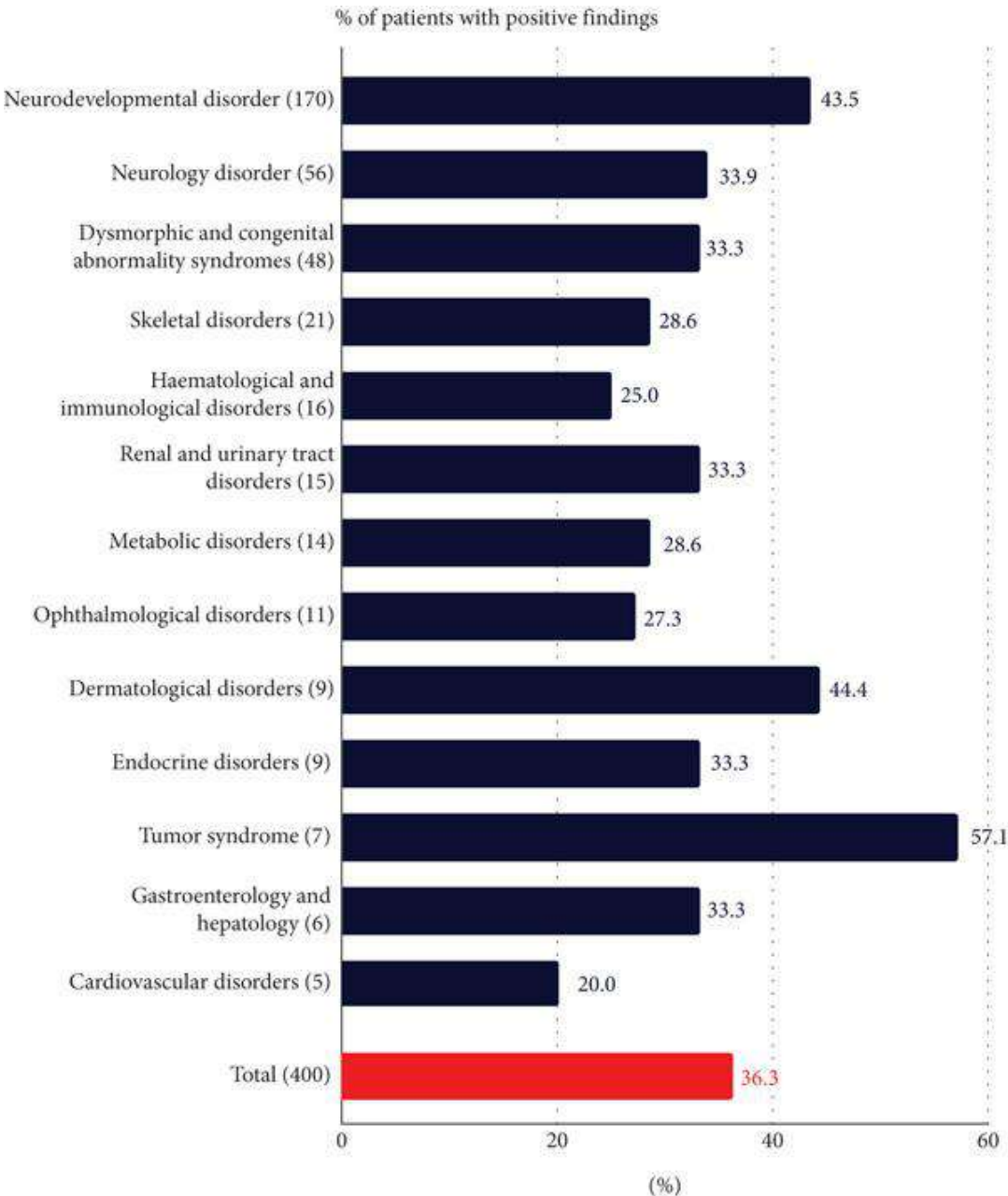
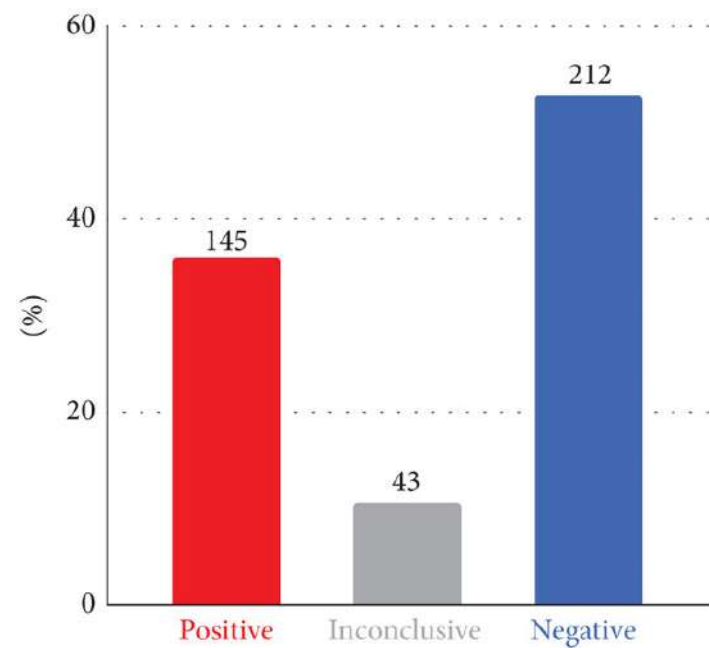


Research Article

Genome Sequencing of Rare Disease Patients Through the Korean Regional Rare Disease Diagnostic Support Program

<https://doi.org/10.1155/humu/6096758>

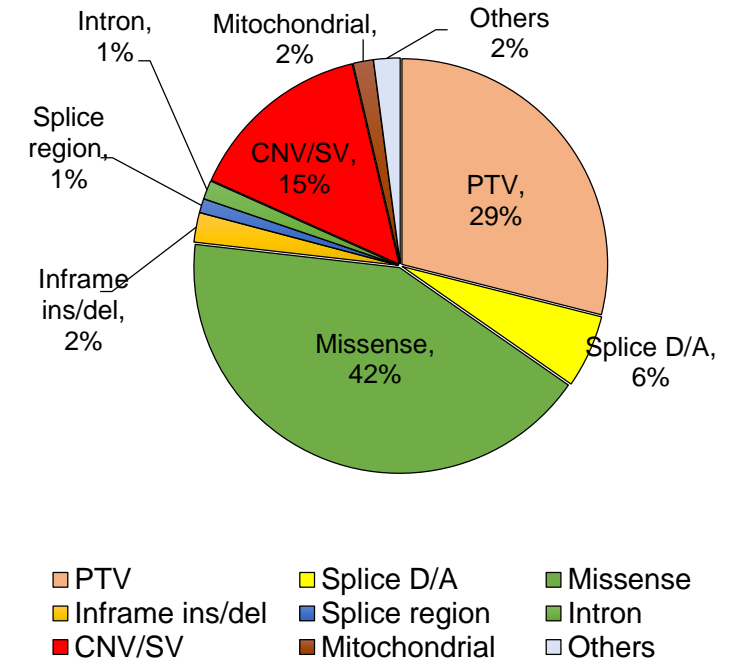
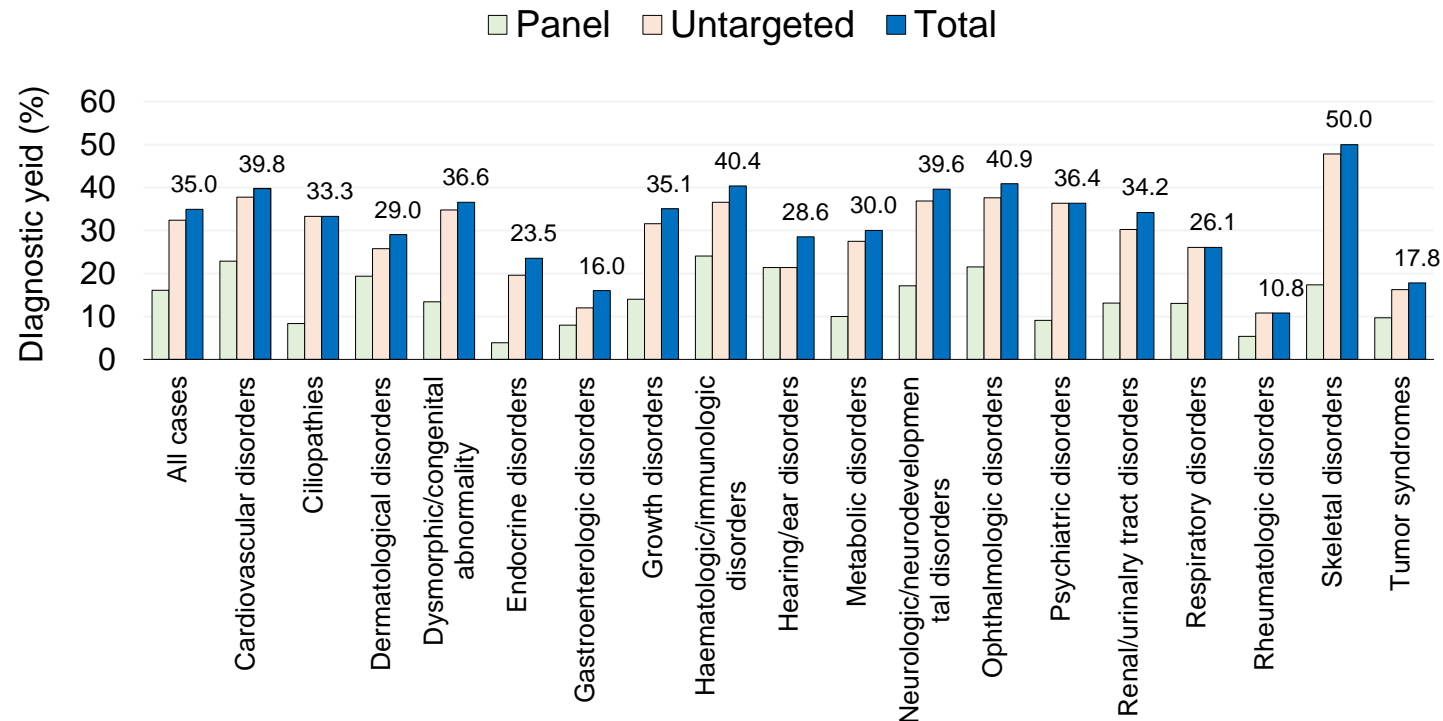
overall diagnostic yield was 36.3% (145/400), with 4.8% (7/145) of the diagnosed patients being reported with variants that could not have been identified



Pilot stage of National Bio Big Data Project

Diagnostic Yield in Rare Diseases

Panel-based analysis vs. WGS



- Total diagnostic yield: 35.0%
- 19 disease categories, 491 diagnostic genes

Challenges in Applying WGS for Rare Disease Diagnosis

1. Interpretation of variants

- Many results fall into Variants of Uncertain Significance (VUS).
- Difficulty in linking novel or rare variants to clinical symptoms.
- Requires continuous updates from global genomic databases (e.g., ClinVar, gnomAD).

2. Lack of standardized guidelines

- Interpretation and classification of variants may vary across labs and institutions.
- Need for consistent application of international standards (e.g., ACMG/AMP guidelines).

3. Limited clinical actionability

- Not all findings are immediately clinically actionable.
- In many cases, diagnosis may not lead directly to treatment options.
- Raises ethical questions about disclosure and patient communication.

Challenges in Applying WGS for Rare Disease Diagnosis

4. Data integration and infrastructure

- WGS data must be integrated with clinical records (EHRs) for full value.
- Requires robust data storage, analysis pipelines, and security infrastructure.

5. Workforce and expertise gaps

- Shortage of trained clinical geneticists, bioinformaticians, and genetic counselors.
- Need for professional training to interpret WGS results in a clinical context.

6. Cost and reimbursement issues

- In Korea, WGS is not yet reimbursed in standard care settings.
- Need for policy development on cost-effectiveness and insurance coverage.

7. Ethical, legal, and social issues (ELSI)

- Management of incidental or secondary findings.
- Ensuring informed consent, data privacy, and patient rights.
- Long-term data governance for secondary use and research.

ACMG PRACTICE GUIDELINE

Exome and genome sequencing for pediatric patients with congenital anomalies or intellectual disability: an evidence-based clinical guideline of the American College of Medical Genetics and Genomics (ACMG)

Kandamurugu Manickam^{1,2}, Monica R. McClain³, Laurie A. Demmer⁴, Sawona Biswas⁵, Hutton M. Kearney⁶, Jennifer Malinowski⁷, Lauren J. Massingham^{8,9}, Danny Miller¹⁰, Timothy W. Yu^{11,12}, Fuki M. Hisama¹³ and ACMG Board of Directors^{14*}

PURPOSE: To develop an evidence-based clinical practice guideline for the use of exome and genome sequencing (ES/GS) in the care of pediatric patients with one or more congenital anomalies (CA) with onset prior to age 1 year or developmental delay (DD) or intellectual disability (ID) with onset prior to age 18 years.

METHODS: The Pediatric Exome/Genome Sequencing Evidence-Based Guideline Work Group ($n = 10$) used the Grading of Recommendations Assessment, Development and Evaluation (GRADE) evidence to decision (EtD) framework based on the recent American College of Medical Genetics and Genomics (ACMG) systematic review, and an Ontario Health Technology Assessment to develop and present evidence summaries and health-care recommendations. The document underwent extensive internal and external peer review, and public comment, before approval by the ACMG Board of Directors.

RESULTS: The literature supports the clinical utility and desirable effects of ES/GS on active and long-term clinical management of patients with CA/DD/ID, and on family-focused and reproductive outcomes with relatively few harms. Compared with standard genetic testing, ES/GS has a higher diagnostic yield and may be more cost-effective when ordered early in the diagnostic evaluation.

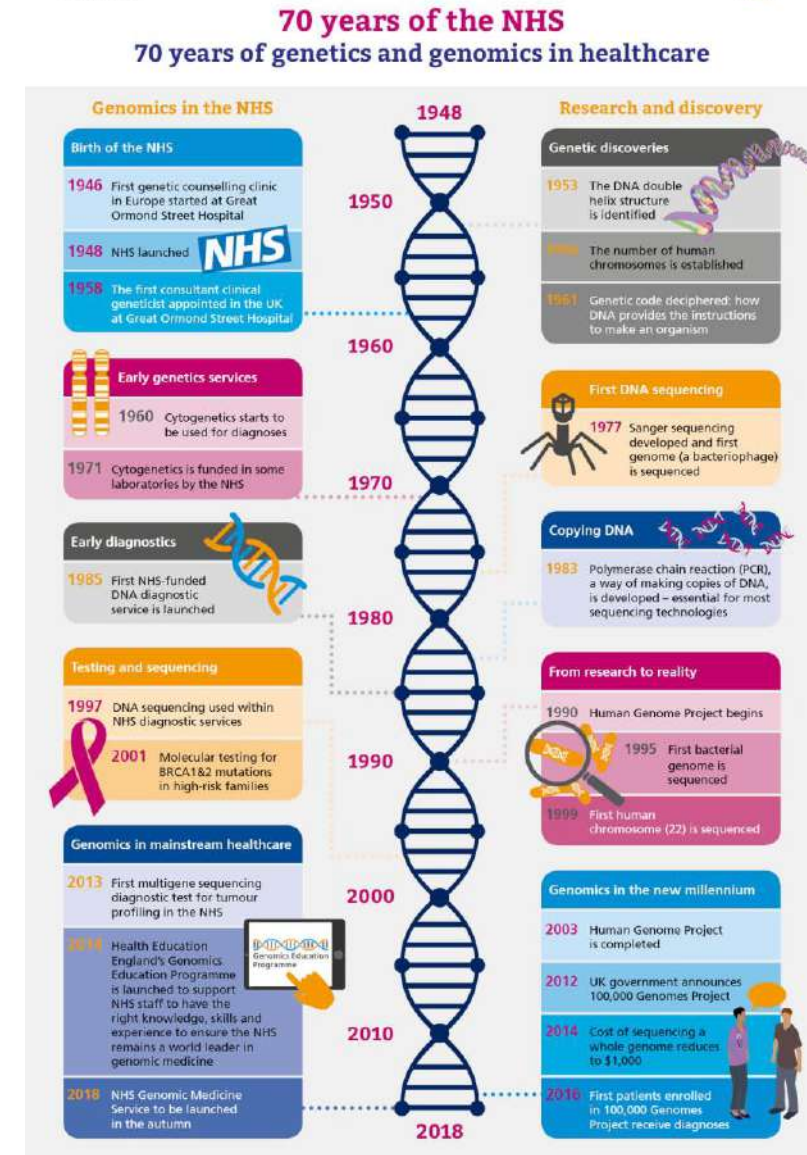
CONCLUSION: We strongly recommend that ES/GS be considered as a first- or second-tier test for patients with CA/DD/ID.

Genetics in Medicine (2021) 23:2029–2037; <https://doi.org/10.1038/s41436-021-01242-6>

NHS England established the NHS GMS in 2018

Strategies

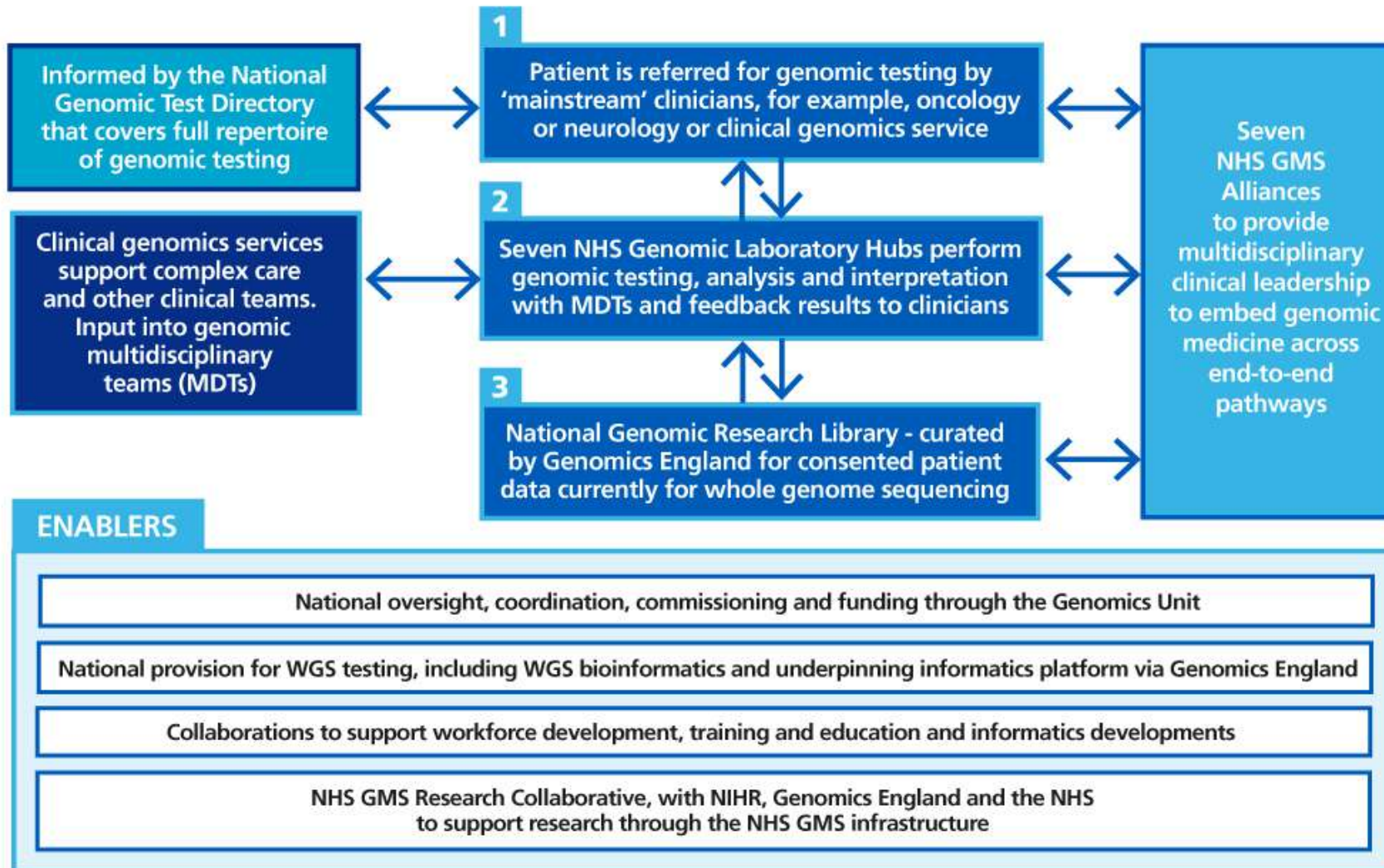
- Embedding genomics across the NHS, through a world leading innovative service model from primary and community care through to specialist and tertiary care.
- Delivering equitable genomic testing for improved outcomes in cancer, rare, inherited and common diseases and in enabling precision medicine and reducing adverse drug reactions.
- Enabling genomics to be at the forefront of the data and digital revolution, ensuring genomic data can be interpreted and informed by other diagnostic and clinical data.
- Evolving the service through cutting-edge science, **research and innovation to ensure that patients can benefit from rapid implementation of advances.**



Find out more about genomics and the future of healthcare

www.genomicseducation.hee.nhs.uk

Infrastructure of NHS Genomic Medicine Service



Therapeutic Application

Zolgensma for SMA Gene Therapy



A Caregiver's Guide to ZOLGENSMA

I'll always remember the day we received the one-time-only dose for SMA

Malachi, treated at ~4 months and pictured at 4 years, was diagnosed with SMA Type 1.

Indication

ZOLGENSMA® (onasemnogene abeparvovec-xioi) is a prescription gene therapy used to treat children less than 2 years old with spinal muscular atrophy (SMA). ZOLGENSMA is given as a one time infusion into a vein. ZOLGENSMA was not evaluated in patients with advanced SMA.

Important Safety Information

ZOLGENSMA can cause acute serious liver injury. Liver enzymes could become elevated and may reflect acute serious liver injury in children who receive ZOLGENSMA. Patients will receive an oral corticosteroid before and after infusion with ZOLGENSMA and will undergo regular blood tests to monitor liver function. Contact the patient's doctor immediately if the patient's skin and/or whites of the eyes appear yellowish, or if the patient misses a dose of the corticosteroid or vomits it up.

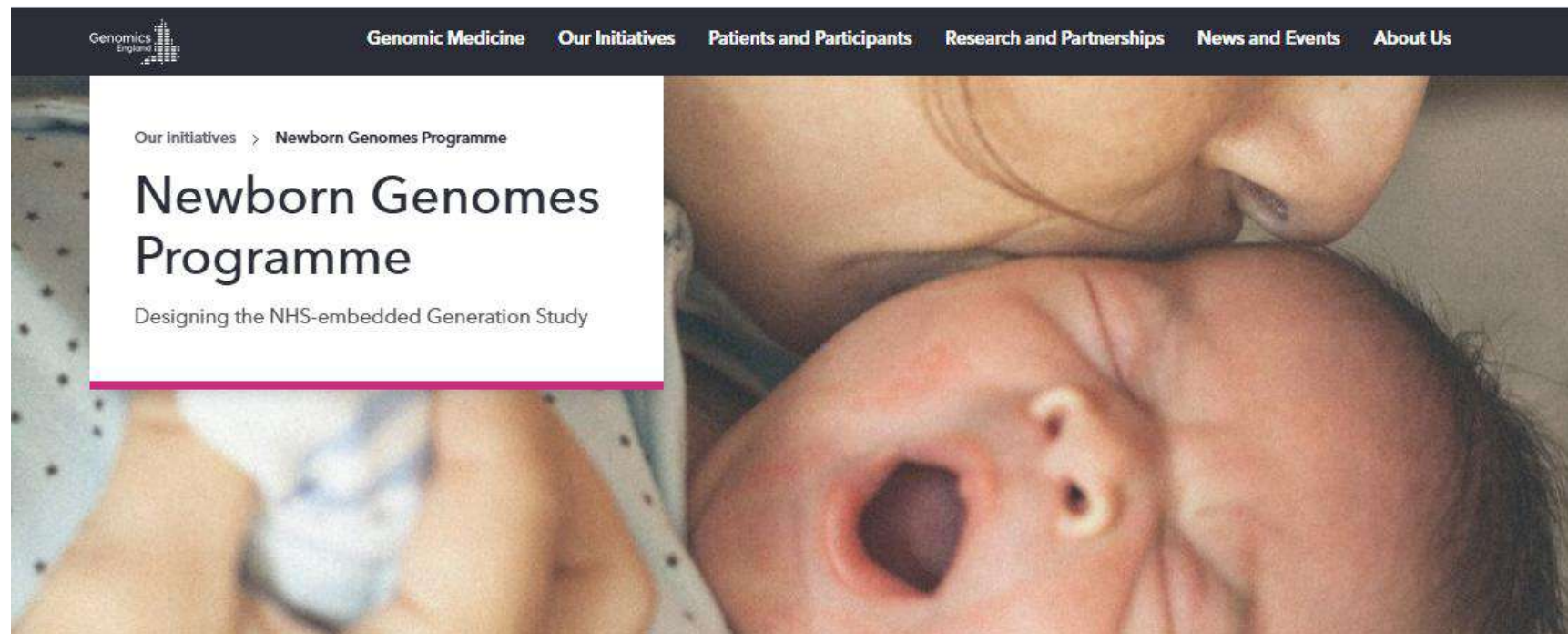
Please see additional Important Safety Information on page 17 and the accompanying Full Prescribing Information.

Table 1. Approved gene therapies for rare diseases in Europe

Drug name (Company)	Type	Clinical indication	Gene	Inheritance pattern	Year of approval in Europe
Zolgensma (Novartis)	AAV vector-based gene therapy	Pediatric patient (<2 years) spinal muscular atrophy	SMN1	Autosomal recessive	2020
Upstaza (PTC Therapeutics International)	AAV vector-based gene therapy	Adults and children (from 18 months) with severe aromatic L-amino acid decarboxylase deficiency	DDC	Autosomal recessive	2022
Strimvelis (Orchard Therapeutics)	Gamma-retroviral vehicle for ex-vivo stem-cell therapy	Patients with severe combined immunodeficiency with no available human leukocyte antigen-matched related stem-cell donor	ADA	Autosomal recessive	2016
Roctavian (BioMarin International)	AAV vector-based gene therapy	Severe hemophilia A	F8	X-linked recessive	2022
Luxturna (Novartis)	AAV vector-based gene therapy	Adults and children with loss of vision due to inherited retinal dystrophy	RPE65	Autosomal recessive	2018
Libmeldy (Orchard Therapeutics)	Lentiviral vehicle for ex-vivo stem-cell therapy	Children with metachromatic leukodystrophy	ARSA	Autosomal recessive	2020

Abbreviation: AAV, adeno-associated virus.

Newborn Screening



[Newborn Genomes Programme](#) [How we work](#) [Engagement](#) [Ethics](#) [How we choose conditions](#) [Evaluation](#)

The Generation Study

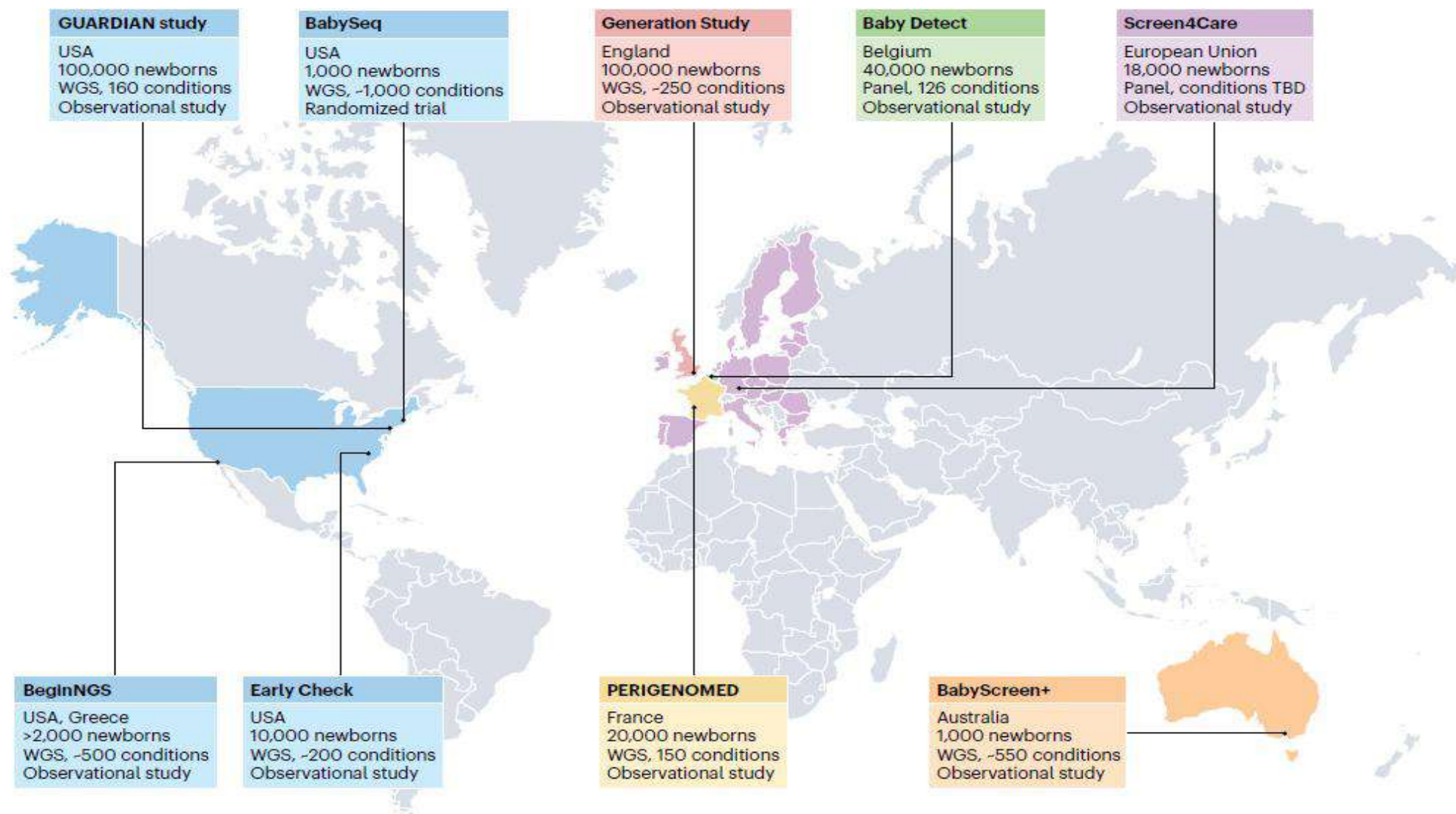
Every year hundreds of babies are born in the UK with rare genetic conditions. Early intervention can enhance the health and quality of life of many of these babies. But these conditions can be hard to diagnose, leading to delays in care.

The Generation Study is a groundbreaking research study which will sequence the genomes of 100,000 newborn babies. We are running our study in partnership with the NHS to understand whether we can improve our ability to diagnose and treat genetic conditions.



Learn more about how the Generation Study works by watching this video.

IC^{NS} International Consortium on Newborn Sequencing



Pilot Project for Genomic Newborn Screening in Korea

1. Target conditions

- Identify early-onset, actionable genetic conditions that may benefit from prompt intervention
- Define which genes and conditions will be screened
- Examine the impact of returning secondary findings (e.g., adult-onset disorders, carrier status)

2. Develop protocols for informed consent and clinical decision support

3. Data privacy and long-term storage

- WGS data is sensitive and must be protected throughout the individual's life
- Requires long-term secure storage infrastructure and policies on secondary use
- Consider re-consent models as the child matures

4. Assess integration with national health records and long-term follow-up systems

Rapid genome sequencing for critically ill neonates

Genetic diagnoses were achieved in 10 patients (total: 20 trio, diagnostic yield: 50%)

Sample #	Phenotype	TAT	Gene	Variant	Inheritance	Variant type
rWGS-03	polycystic kidney dysplasia	5day 4h	<i>AGT</i>	NM_001384479.1:c.292G>A	Homozygous	Missense
rWGS-04	(unconjugated) hyperbilirubinemia * maternal hemolytic anemia	5day 2h	<i>ANK1</i>	NM_000037.4:c.2394_2397 del	Heterozygous (Maternal)	Frameshift
rWGS-09	Fetal thrombotic vasculopathy	5day 14h	<i>NSD1</i>	NM_022455.5:c.5885T>C	De novo	Missense
rWGS-10	Renal vein thrombosis, both., IVC thrombus	5day 1h	<i>SERPINC1</i>	NM_000488.4:c.235C>T	Homozygous	Missense
rWGS-11	Asymptomatic hyperammonemia	7day	<i>OTC</i>	NM_000531.6:c.513G>C	X-linked	Missense
rWGS-12	Hyperglycemia, Type I diabetes mellitus	5day 1h	<i>INS</i>	NM_000207.3:c.265C>T	De novo	Missense
rWGS-14	ARPKD, Potter syndrome	4day 13h	<i>PKHD1</i>	NM_138694.4:c.6840G>A NM_138694.4:c.6602T>A	Compound heterozygous	Stop gained, stop gained
rWGS-17	thrombocytopenia with growth retardation, feeding difficulties	3day	<i>SAMD9</i>	NM_017654.4:c.2414A>G	De novo	Missense
rWGS-18	Hyperammonemia, lactic acidosis	4.5day	<i>OTC</i>	NM_000531.6:c.841T>G	X-linked	Missense
rWGS-19	Hydrops fetalis, Congenital thrombocytopenia, Lung hypoplasia	3.5day	<i>LZTR1</i>	NM_006767.4:c.742G>A	De novo	Missense

Pharmacogenetics

Strong guideline frameworks

- CPIC (Clinical Pharmacogenetics Implementation Consortium)
- DPWG (Dutch Pharmacogenetics Working Group)
- FDA labeling includes PGx info for over 400 drugs.

Implementation Models

- Preemptive Testing: Genetic data obtained before prescribing, stored in EHR (e.g., St. Jude, Mayo Clinic)
- Reactive Testing: Ordered at the time of prescribing based on need
- Integration into Clinical Decision Support (CDS) is key for usability.

Genomic medicine initiatives

- USA: Widespread implementation in academic health centers ; PGx panels offered by commercial labs
- Europe: National strategies (e.g., Netherlands, UK) embedding PGx into care pathways
- Asia: Japan and Singapore have PGx in routine drug safety

Implementation of Genomic Medicine in Clinical Practice

Current limitations:

- No reimbursement pathway

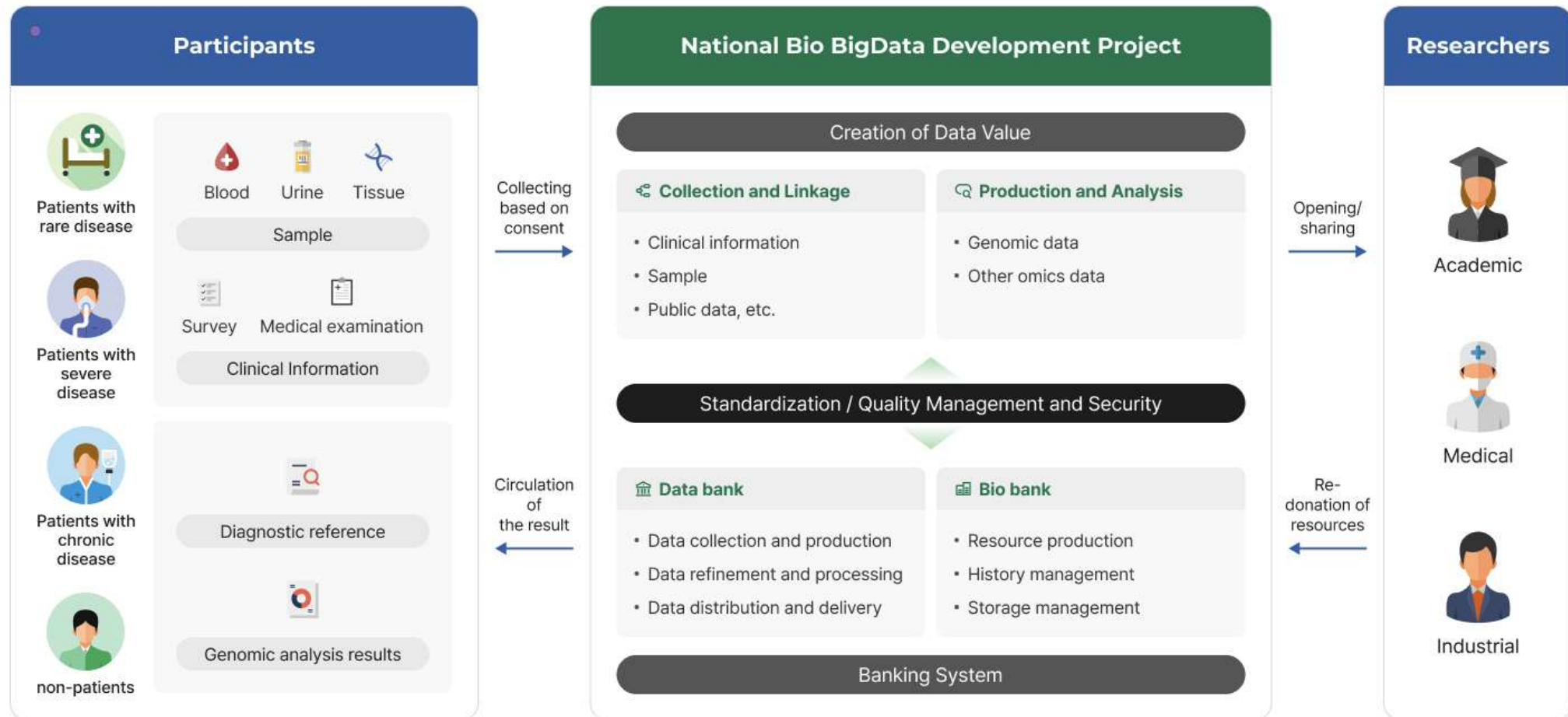
- Limited scalability and clinical feedback loop

Policy opportunity:

- Transition WGS from research-only to **standardized clinical application**, supported by ethical, regulatory, and reimbursement frameworks.

National Bio Big Data Project

Integrated Bio Big Data of 1,000,000 people

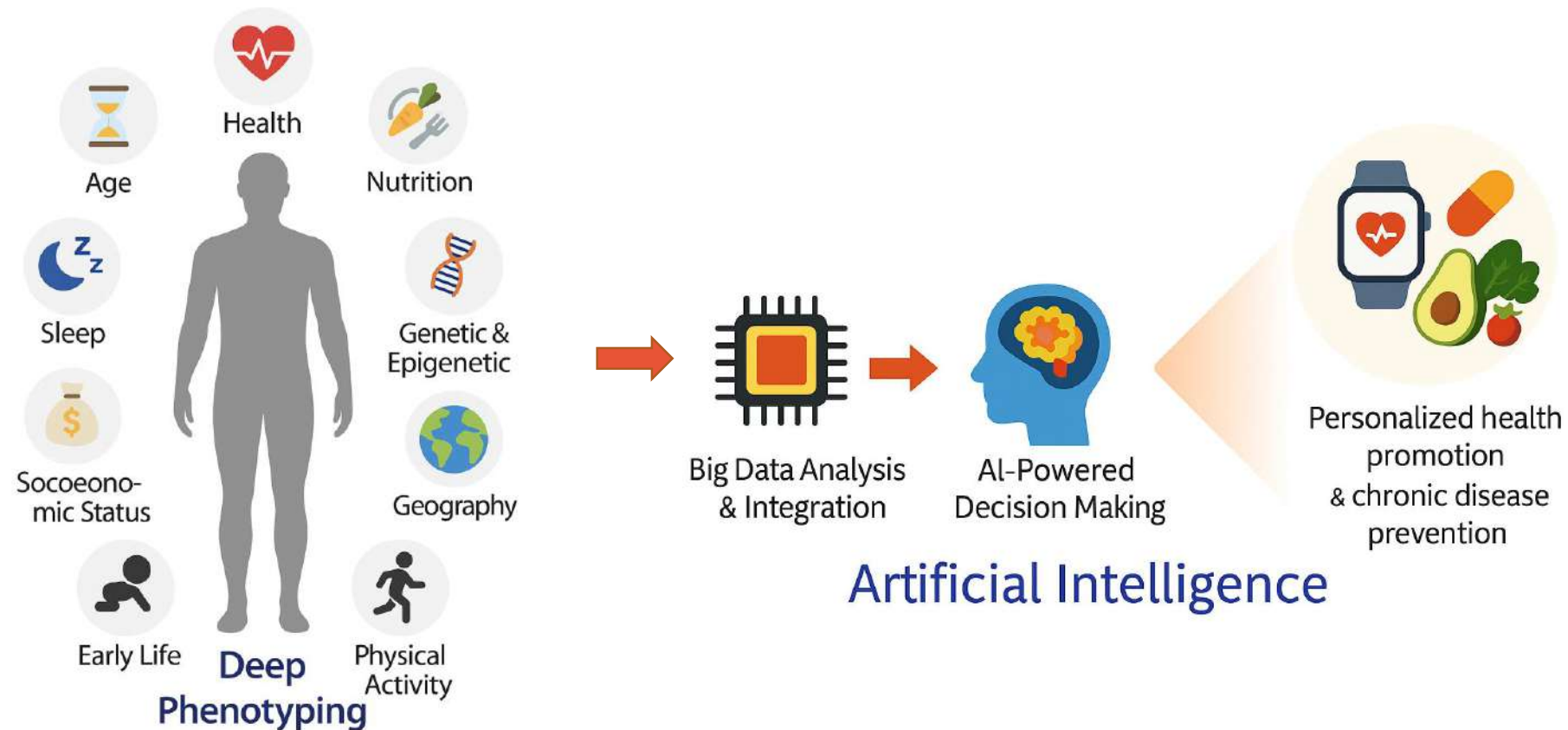




**Biomedical Research Data
beyond Genomics**

KNIH's Bio Big Data Initiative

Deep Phenotyping is essential for Precision Health



Deep phenotyping and artificial intelligence for health promotion and chronic disease prevention.

KNIH's Bio Big Data Initiative



Large-scale
Cohort Studies



Korea Biobank
Project



Korean Genome
Analysis Projects



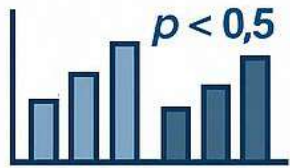
Biomedical data
beyond the genome



Data sharing
CODA



Large-Scale Cohorts as Key Drivers of Precision Health Research



Capturing Genetic and Environmental Diversity

- Enables analysis of diverse populations across age, ethnicity, lifestyle, and regions
- Helps uncover gene-environment interactions and rare variant effects

Powering Statistical Validity

- Large sample sizes improve statistical power for discovering disease biomarkers and risk factors
- Facilitates robust subgroup analysis (e.g. sex-specific, age-specific effects)

Longitudinal Insights for Disease Progression

- Tracks health trajectories over time to understand natural history of diseases
- Supports early detection and predictive modeling of chronic diseases

Foundation for Data Integration

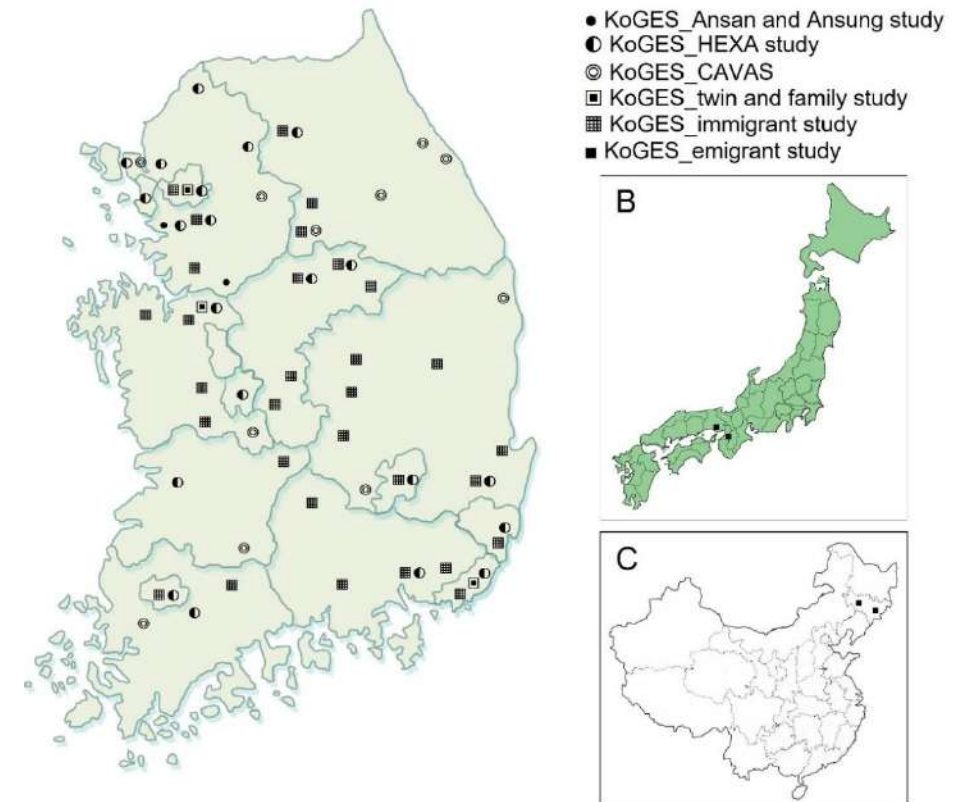
- Allows linkage with genomic, clinical, behavioral, and digital health data
- Supports AI-driven precision diagnostics and personalized treatment pathways

KoGES

Korean Genome and Epidemiology Study

- Started in 2001, **6 sub-cohorts** representing diverse population groups.
- **KoGES-Ansan & Ansung study is still ongoing with over 20 years of F/U.**
- **240,000** participants have been recruited by the end of 2014
- **Major target diseases** : T2DM, metabolic syndrome, hypertension, cardiovascular disease, osteoporosis, cancer

KoGES						
Population-based studies			Gene-environment model studies			
	KoGES_ Ansan and Ansung study	KoGES_ HEXA study	KoGES_ CAVAS	KoGES_ twin and family study	KoGES_ immigrant study	KoGES_ emigrant study
Baseline	10,030 (2001-2002)	173,195 (2004-2013)	28,337 (21,715)* (2005-2011)	3,202 (2005-2013)	7,191 (Immigrants 2,875, Korean spouses 1,911, Children 2,405) (2006-2014)	3,556 (1,062 in Japan, 2,493 in China) (2005-2011)
Follow up						
1st	8,603 (2003-2004)	65,608 (2012-2016)	12,463* (2007-2014)	2,030 (2008-2014)	1,824** (2012-2014)	773*** (2008-2013)
2nd	7,515 (2005-2006)		11,399* (2008-2016)	940 (2009-2014)		549*** (2010-2011)
3rd	6,688 (2007-2008)		6,423* (2011-2016)			520*** (2012-2013)
4th	6,665 (2009-2010)					
5th	6,238 (2011-2012)					
6th	5,906 (2013-2014)					
7th	6,318 (2015-2016)					
8th	6,157 (2017-2018)					
9th	5,856 (2019-2020)					
10th	5,511 (2021-2022)					



NIH's Cohort or Registry

Currently operating more than 30 cohorts or registries

General Population / Special Group Cohorts

Cohort Name	Start year	Number of Participants
KoGES Rural (Ansung)	2001	5,018
KoGES Urban (Ansan)	2001	5,012
KoGES Rural (CAVAS)	2004	28,337
Korea Nurses Health Study (KNHS)	2013	20,613
Women's Health Study	2014	4,684
Korea Neonatal Network (KNN)	2013	22,651
Korea Frailty and Aging Cohort Study (KFACS)	2016	3,011
Korea Urban and Rural Elderly Study (KURE)	2011	3,517
Cardiovascular & Metabolic Etiology Research Cohort (CMERC)	2013	11,375
Korean Transplant Registry (KOTRY)	2014	45,396

Disease-focused Cohorts

Cohort Name	Start year	Number of Participants
Korean Stroke Registry, KSR	2017	25,515
Korean childhood Asthma Study	2016	961
Korea COPD Subgroup Study	2016	2,952
Resistant hypertension cohort	2018	1,432
Community-based Dementia Cohort	2019	4,789
Parkinson Diseases Registry	2021	749
Hospital-based Dementia Cohort	2021	646
Early-onset Dementia Cohort	2021	406
Korean Severe Asthma Registry	2022	591
Korea HIV/AIDS Cohort	2006	1,644
Korea HCV Cohort	2007	3,858
Korea HBV Cohort	2014	3,021
Tuberculosis Cohort	2019	2,001

Cohorts for Aging Research



KoGES

Korean Genome and Epidemiology Study

Started in 2001

Study population: aged 40~69

Collaboration: Aju University, Korea University, Hanyang University



KURE

Korean Urban and Rural Elderly cohort

Started in 2012

Study population: aged 65 and older

Collaboration : Yonsei University



KFACS

Korean Frailty and Aging Cohort Study

Started in 2016

Study population: aged 70~84

Collaboration : Kyung Hee University

Korean Centenarian Study (2025~)

Starts in 2025

Study population: aged 90 and older

Focusing on healthy aging and longevity

Biomedical Data Generation



Omics data

- Genome: K-chip(DNA microarray chip) & imputed data, WGS
- Epigenome: DNA methylation profiles using Illumina arrays (e.g., 850K, 935K)
- Metabolome & Proteome



Brain MRI

- ~5,700 participants in KoGES- Ansan and CAVAS
- Structural & functional analysis of brain imaging(T1, T2, Flair, DTI, fMRI)
- Regional brain volume & morphometric measures



Microbiome

- A total of 800 participants in KoGES-Ansung & KFACS
- including:
 - Fecal DNA samples
 - Microbiome sequencing data
 - Related clinical and epidemiological information



Air pollution

- Cohort data have been linked to air pollution exposure data
 - PM10, PM2.5
 - SO2, NO2, CO, O3
- Meteorological data: Humidity, Wind speed, Precipitation, Cloud cover, Solar radiation, etc



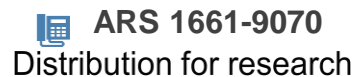
Link with Health Data

- Death data(National statistics)
- Cancer registry (National Cancer Center)
- National Health Insurance data



Korea Biobank project
National Biobank of Korea

Korea Biobank Network = NBK + 47 regional biobanks



Others

Korea Biobank Project



국립중앙인체자원은행
NATIONAL BIOBANK OF KOREA

0.47

million participants

- **Population-based cohorts or Surveys**

- Korean Genome Epidemiology Study (KoGES, n=249K+)
- Korea National Health & Examination Survey (KNHANES, n=108K+)

- **Infection related studies**

- National immunity surveys, etc. (n=48K+)

- **Other biobanking studies**

- CMERC, KOTRY, BICWALZS etc. (n=52K+)

blood DNA Urine

10.5 mil vials stored in a central biobank (NBK)

Epidemiological & clinical data



380K participants

Multi-omics data

(K-chip, WES, WGS, RNA-seq, methyl, absolute IDQ, etc)



170K participants



한국인체자원은행네트워크
Korea Biobank Network

0.76

million participants

- **Specialized Biobank Subnetwork**

- 10 subnetwork (10 Hub Biobanks + 20 collaborative Biobanks)
- Establishment and operation of human biobank specialization subnetwork

- **Innovative Biobanking Consortium**

- 4 consortium (Chronic cerebrovascular diseases & Alzheimer's diseases, Sarcoma, developmental disorder)
- Establishment of clinician-led biobanks and healthcare R&D Ecosystem

Tissue blood DNA

9.5 mil vials stored in across the regional biobanks (n=47)

Clinical data

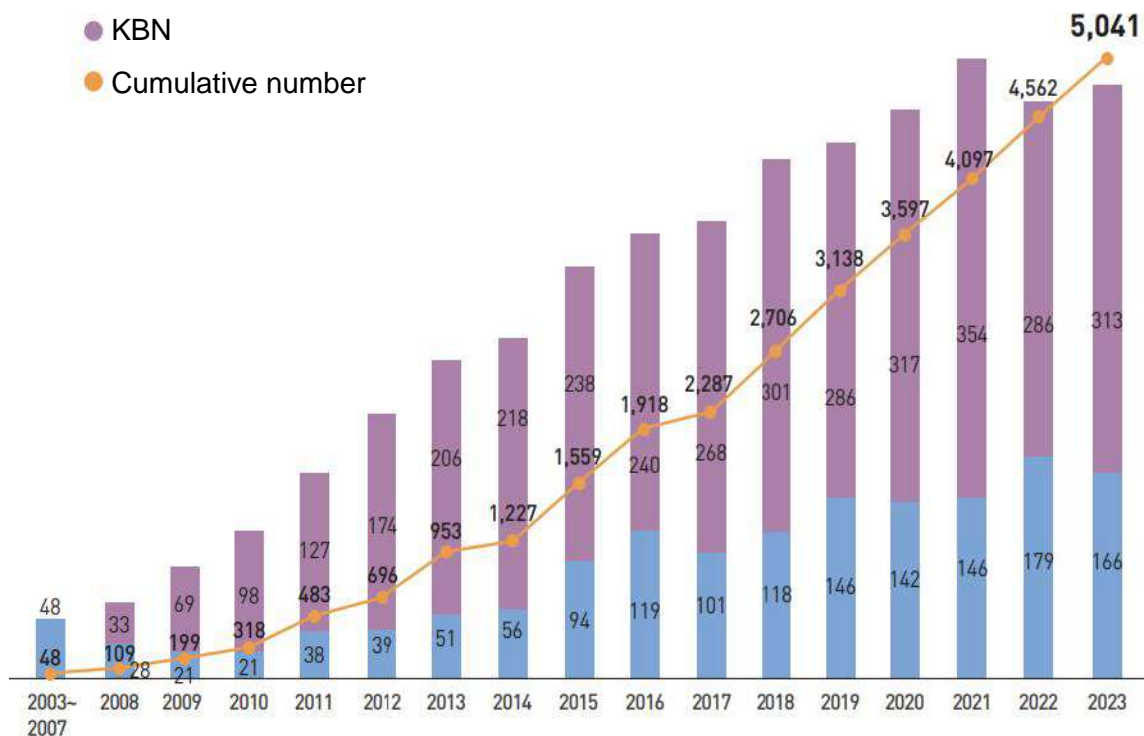


Total 0.76 million (annotated with 11 clinical variables)

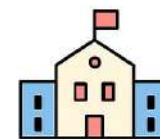
Utilization of Human Biological Materials

(Unit: Project #)

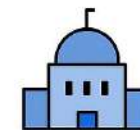
- National Biobank of Korea
- KBN
- Cumulative number



Human biological materials have been used in over 5,000 studies



University/
hospitals
4,361



Government
institutions
283



Biotech
companies
254



Other
institutions
143

2053 outcomes have been derived from human resources

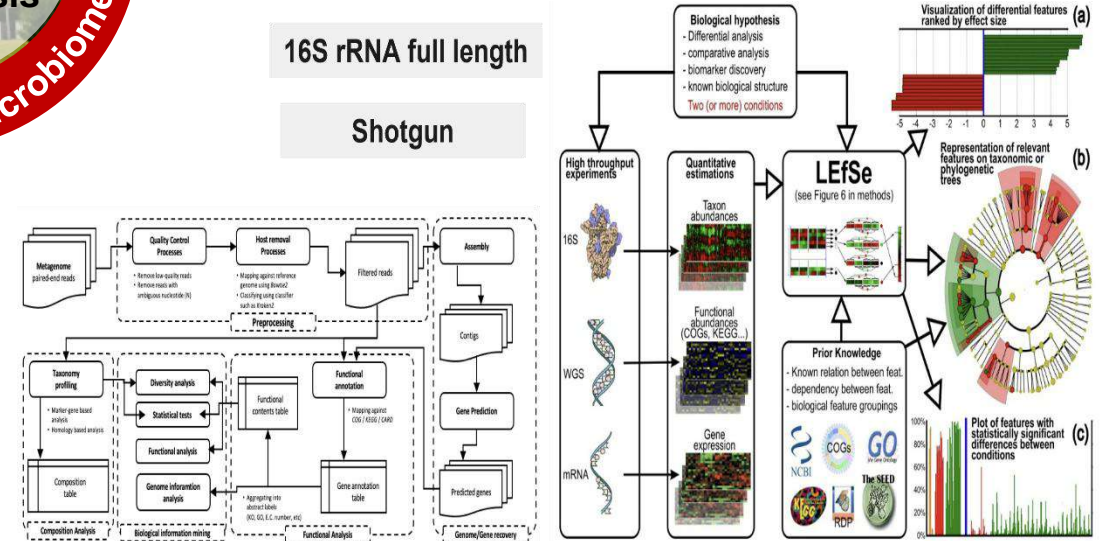
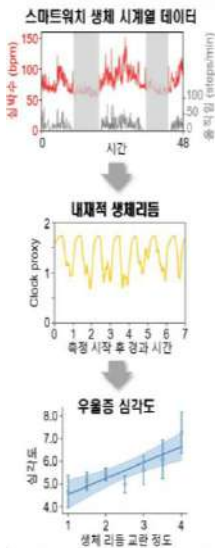
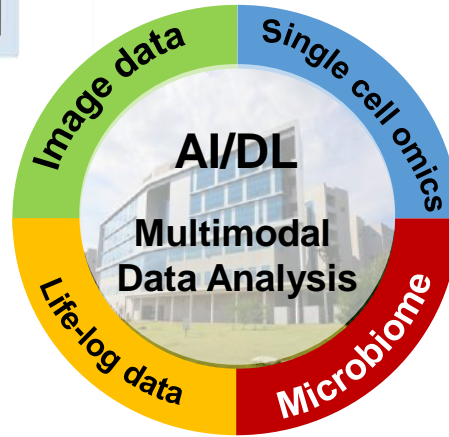
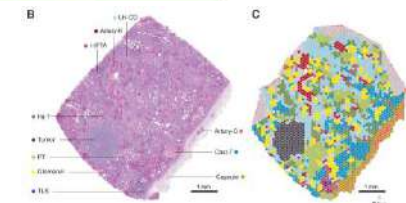
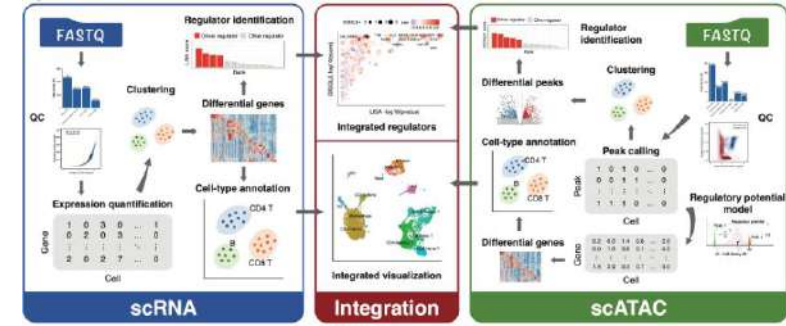
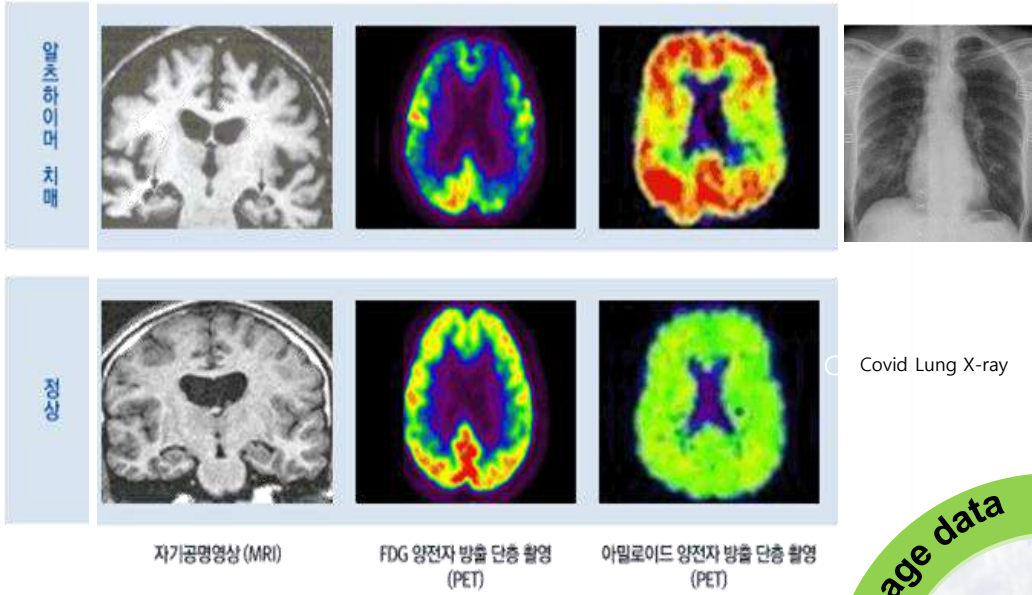


Papers
1,876



Patents
177

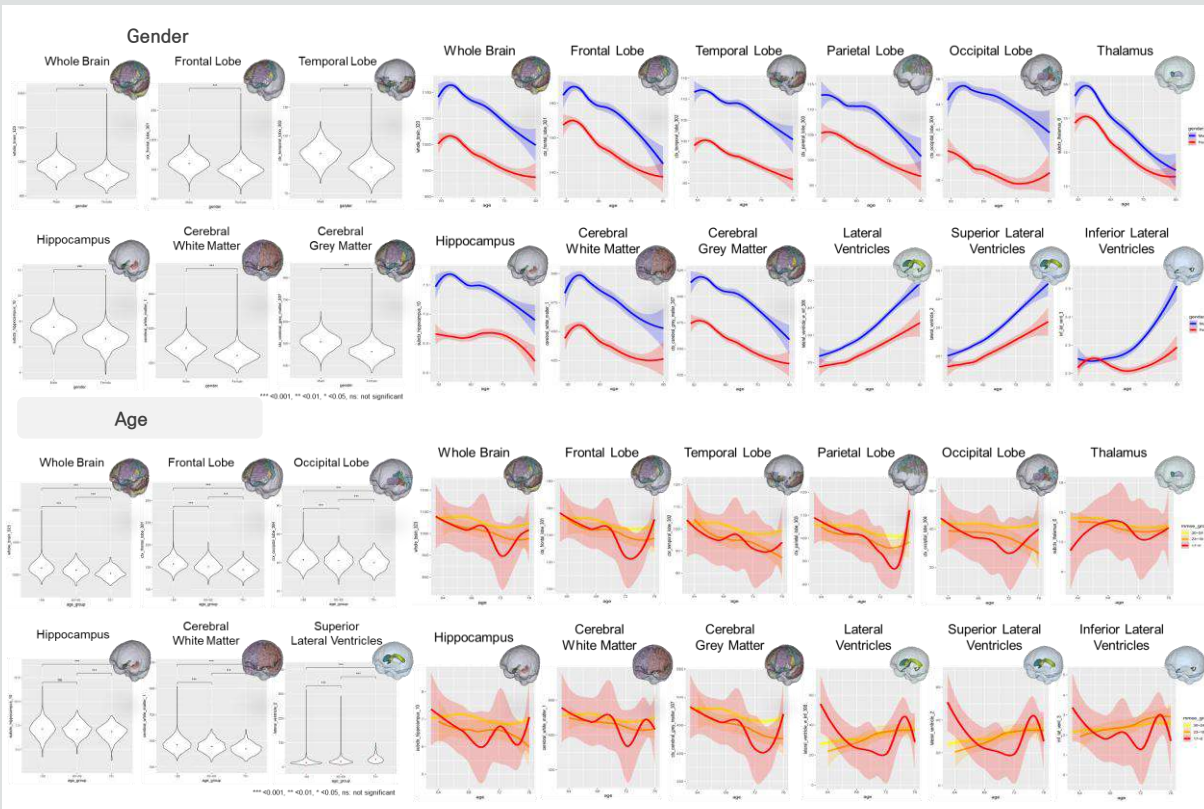
As of the end of 2023



Curation and database development of neuroimaging data

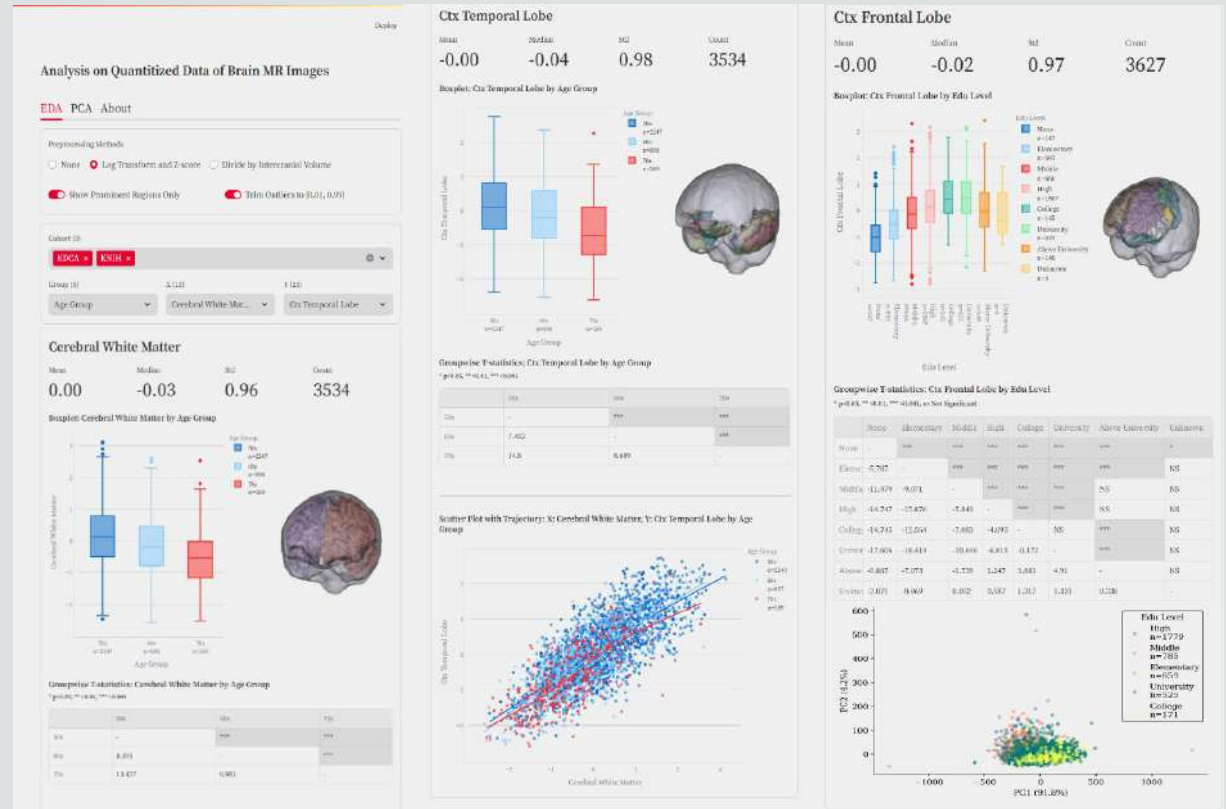
Establishment of a Korean Brain MRI Reference Database

Establishment of a reference dataset of brain regional volume values based on sex, age group, cognitive function, and disease presence/severity (approx. 5,000 individuals)



Development of an App for Visualizing Brain Imaging Data

Development of an interactive web-based visualization tool for multi-center cohort-based brain MRI quantification results and clinical/epidemiological data sets



Central Data Repository

CODA

Research data deposited in CODA
can be used for your research



Search the data you want by keyword.



Public Resources



Project
181



Participants
831,664



Approved Access&Sharing
159

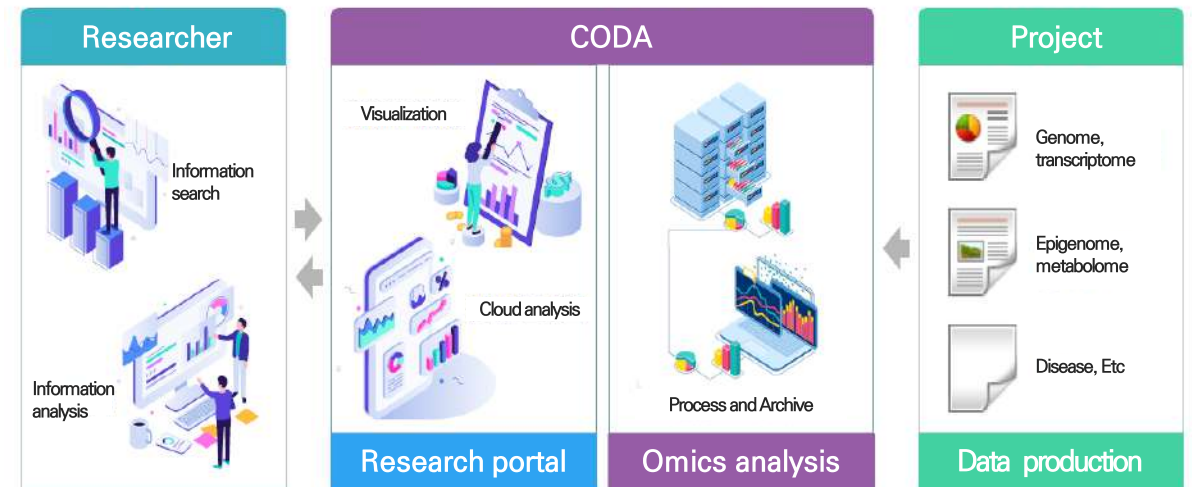


Omics Data
438,923

CODA

Clinical & Omics Data Archive

- ✓ It collects **clinical/epidemiological** and **omics data** (microarray, whole-exome, whole-genome, transcriptome, metabolome)
- ✓ Collect, share, and utilize biomedical data (resource) derived from national R&D projects or voluntary data deposit
- ✓ Established in 2016 at NIH
in accordance with Article 8(2)(3) of the Enforcement Decree of the Bio-Research Resources Act.
- ✓ The system enables researchers to search for and request necessary data, and to perform analyses using a built-in analysis pipeline within the cloud infrastructure.



Open data platform

- ✓ Improve user accessibility
- ✓ Provision of analytical infrastructure via cloud-based operations



데이터 조회

MY DB관리

자료실

공지사항

OPEN KoGES 소개

OPEN KoGES(KoGES 공유 플랫폼)는 한국인유전체역학조사사업(Korea Genome and Epidemiology Study; KoGES)에서 생산한 21만명의 코호트자료(지역사회(안상·안산)/도시/농촌)를 통합검색부터 활용(분석) 인프라까지 지원하는 원스톱 서비스입니다.

OPEN KoGES는 일반연구진단을 대상으로 한 대규모 코호트 통합 자료를 제공하여, 연구자 필요시 맞춤형DB를 설계하고 생성할 수 있을 뿐만 아니라, 임상·역학 및 유전체 데이터의 다양한 분석 환경(GUI, 파이썬, 터미널 등)을 구축·제공하여 대용량 데이터를 빠르게 탐색하고, 다양한 시각적인 방법으로 분석 활용하실 수 있습니다.

OPEN KoGES는 보건의료 연구데이터의 연구 효율성을 제고하고, 데이터의 활용가치를 극대화하여 연구자들에게 대규모 데이터를 활용할 수 있도록 지속적인 서비스를 지원하고 개발하겠습니다.

주요 질환 진단 여부(기반조사)



데이터 현황

		지역사회 (안상/안산)				농촌				도시			
		조사기간	대상자수	유전체 (800개좌)	변수(개)	조사기간	대상자수	유전체 (800개좌)	변수(개)	조사기간	대상자수	유전체 (800개좌)	변수(개)
3월	01-02	10,000	7,522	1,852		05-11	28,337	18,925	1,028	04-13	173,195	121,048	1,795
11월 추적	03-04	8,000	6,000	1,572		07-14	12,403	9,016	850	07-10	65,806	52,467	1,049
22월 추적	05-06	7,515	5,811	2,348		08-16	11,309	8,235	898				
33월 추적	07-08	6,088	5,186	1,789		11-15	6,423	4,551	857				
43월 추적	09-10	6,665	5,164	2,047		14-15	1,449	1,083	206				
53월 추적	11-12	6,288	4,823	2,123									
63월 추적	13-14	5,909	4,515	1,943									
73월 추적	15-16	6,318	4,858	1,635									
83월 추적	17-18	6,197	4,739	1,144									
93월 추적	19-20	5,854	4,508	1,205									

데이터 분석 유형

임상역학분석

GUI 기반의 분석 및 그래프를 지원하는 클라우드 서비스
#기술용계량 #그래픽스 #동계분석 #변수변환

[더보기 →](#)

임상역학 및 유전체 분석

Jupyter 기반 분석을 지원하며, 연구자를 위한
유전체 역학 분석 유틸리티를 제공한다.
#GWAS #PRS #SNP Heritability #QC

[더보기 →](#)


```

Jupyter bmi.log Last Checkpoint: 8 days ago

File Edit View Settings Help

1 PLINK v2.00a5LM 64-bit Intel (28 May 2023)
2 Options in effect:
3 --allow-no-sex
4 --bfile /home/rex/yjkin/data/KBA_160K_Retractor_masking
5 --covar /home/rex/yjkin/data/linear/bmi/except_BMI_covar.txt
6 --covar-name CT,WC, AS,SEX,AGE
7 --glm hide-covar
8 --out /home/rex/yjkin/data/linear/check/bmi
9 --pheno /home/rex/yjkin/data/linear/bmi/except_BMI_pheno.txt
10 --pheno-name BMI
11
12 Hostname: koges-54
13 Working directory: /
14 Start time: Tue Jan 16 11:28:32 2024
15
16 Note: --allow-no-sex no longer has any effect. (Missing-sex samples are
17 automatically excluded from association analysis when sex is a covariate, and
18 treated normally otherwise.)
19 Random number seed: 1705371632
20 206390L MIB RAM detected, ~1866894 available; reserving 1831950 MIB for main
  
```

GWAS 수행 로그화면

The Need for Integrated Data Ecosystems

- To enable comprehensive and accurate decision-making
- To improve patient outcomes
- To accelerate innovation and research
- To support health system efficiency
- To respond effectively in public health crises

The Need for Integrated Data Ecosystems

Data	Description	Data source
Clinical Data	Real-time medical data from healthcare providers	EMRs, lab results, prescriptions
Public Data	National-level health, claims, and surveillance	NHIS, HIRA, KDCA, mortality data
Research Data	Discovery-driven data from cohort & omics studies	Biobanks, Cohort/Registry, National Research Initiatives

Barriers to Overcome

- **Technical:** Lack of standards, interoperability gaps
- **Legal:** Regulatory silos between health and research data
- **Ethical:** Consent frameworks for secondary use
- **Institutional:** Siloed data ownership, weak incentives for sharing

***"No single dataset is sufficient.
Smarter healthcare requires smarter data
integration."***

By bringing together clinical, public, and research data, Korea can lead the way in evidence-based, personalized, and equitable healthcare decisions.

Korea
National
Institute of
Health

국립보건연구원
Korea National Institute of Health

인류와 미래세대를 위한 질병보건연구

국립보건연구원

혁신적인 보건연구를 선도하는 글로벌 리더

A Global Leader Pioneering Innovative Health Research

Thank you



National Institute of Health
Republic of Korea