Dec. 9th. 2022 ABC symposium

## Introduction to

## the Korea National Genome Project

## **Ji-Hwan Park**

#### **Korea Bioinformation Center (KOBIC)**

Korea Research Institute of Bioscience & Biotechnology (KRIBB)





#### Why do we need the National Genome Project?

Drastic advances in high-throughput sequencing and data analysis and sharing technologies



#### A new wave of population genomics & precision medicine





## Why do we need the National Genome Project?

#### National genomic projects have been launched in > 41 countries





Stark *et al.*, *AJHG* (2019) Kovanda et al., *BMC Human Genomics* (2021)



**Our Mission** 

# **The National Genome Project**

Providing large-scale genomic & clinical data to scientific & industrial community in the purpose of studying precision medicine,

while protecting the participants' privacy





#### Workflow of the NGP



## Cohorts in the NGP, pilot phase (2020 – 2022)

|                | genome                                       | The pilot phase of the NGP provides and clinical information from up to <b>25,000 samples</b> :   |
|----------------|--|---|
| Rare           | Rare Disease (RD)                            | Total 15,000 participants (predominantly trio)  |
| nal<br>Ithy    | Normal Participants<br>(KoGES)               | Total 5,000 participants  |
| Norn<br>/Heal  | Normal Participants<br>(KGP)                 | Total 2,383 participants  |
| ently<br>rable | Alzheimer's Disease<br>(AD)                  | Total 500 Alzheimer's disease patients and 500 normal participants                                |
| Curre          | Autism Spectrum<br>Disorder (ASD)            | Total 849 participants (>200 patients & their parents; >149 families)                             |
| Cancer         | Colorectal Cancer<br>(CRC)                   | Total 300 patients (600 tumor and paired normal adjacent tissues)                                 |
|                | Non-Smoking Lung<br>Adenocarcinoma<br>(NSLA) | Total 84 never-smoker lung adenocarcinoma patients (168 tumor and paired normal adjacent tissues) |
| KKK            | X  | Kobic 국가생명연구자원정보센터     Korea Bioinformation Center  |

- Clinical data from 23,833 participants (= 25,000 samples) have been collected (from 15,000 RD, 300 CRC, 84 NSLA, 849 ASD, 1,000 AD, 1,600 KGP, & 5,000 KoGES cohorts)
- 20,000 whole-genome sequencing (WGS) data have been newly generated and 5,000 WGS data were collected from NGP cohort partnership





## Progress of the NGP pilot phase, as of 1<sup>st</sup> Nov. 2022



# Collection of Samples & Data

#### **Collection of samples & clinical data based on patient consent**

#### **Bioethics and BioSafety Act, Republic of Korea**

#### Participants' consent

■ 생명윤리 및 안전에 관한 법률 시행규칙 [별지 제34호서식]

인체유래물 연구 동의서 (앞쪽) 동의서 관리번호 성명 생년월일 인체유래물 주 소 기증자 성별 전화번호 성 명 관계 법 정 대 리 인 전화번호 Institutional Review Board 성 명 연구책임자 전화번호 이 동의서는 귀하로부터 수집된 인체유래물등(인체유래물과 그로부터 얻은 유전정보를 말합니다)을 질병의 진단 및 치료 등의 연구에 활용하기 위한 것입니다. 동의는 자발적으로 이루어지므로 아래의 내용을 읽고 궁금한 사항은 상담

이 중의서는 귀하도구려 구입된 인세류대출당(인세휴대출과 그로구터 일은 유신정모를 절압니다)을 절명의 신단 및 지묘 법 개발 등의 연구에 활용하기 위한 것입니다. 동의는 자발적으로 이루어지므로 아래의 내용을 읽고 궁금한 사항은 상담 자에게 묻고 질문할 기회를 가지고 충분히 생각한 후 결정하시기 바라며, 이 동의서에 대한 동의 여부는 귀하의 향후 검 사 및 치료 등에 어떤 영향도 미치지 않습니다.

- The research participation agreements and plans are approved by multiple Institutional Review Boards (IRBs).
- Blood or saliva samples and clinical information are collected, according to the strictly managed quality control (QC) processes.





#### Approval of research plans

#### **Rare disease recruitment**

The rare disease patients take part in the NGP via the following hospitals:



Seoul National University Hospital (SNUH)



Pusan National University Yangsan Hospital (PNUYH)



Seoul National University Bundang Hospital (SNUBH)



Inha University Hospital (INHAUH)



Jeonbuk National University Hospital (JBUH)



Chungbuk National University Hospital (CBNUH)



Samsung Medical Center (SMC)



The Catholic Univ. of Korea Seoul St.Mary's Hospital (SSMH)



Severance Hospital (YUHS)



Cheju Halla General Hospital (CHH)



Chonnam National University Hwasun Hospital (CNUHH)



Wonju Severance Christian Hospital (WSCH)



Asan Medical Center (AMC)



Kyungpook National University Chilgok Hospital (KNUH)



Ajou University Hospital (AJOUMC)



Inje University Busan Paik Hospital (INJE)



Chungnam University Hospital (CNUH)



#### **NGP** cohort partnership

In cooperation with previous projects (NGP cohort partnership), NGP collected diverse cohort datasets,



Current research & medical staffs in NGP cohort partnership may participate in the next phase of NGP, which aims to collect the data from > 1M participants





#### NGP cohort partnership

The NGP cohort partnership encompasses diverse types of cohorts (cancer and currently incurable disease patients & normal/healthy participants)

Following implications from the cohort partnership:

Establishment of the standard operating procedure (SOP) for diverse types of diseases and healthy participants

|   | ≣ |
|---|---|
| C |   |

- Ethical, Legal and Social Implications (ELSI)
  - (*e.g.*, Consent forms and research plans regarding the prospective / retrospective studies in the NGP)
- <u>Strategies for sample and data collection</u>
   (e.g., sampling and banking solid tumors rather than blood)
- Data standardization or strategies for update of clinical information (*e.g.*, trajectory analysis of clinical information & outcomes)





### Standardization & QC of clinical information

The Research Environment Platform provides structured clinical data:



**KKIRK** 

 Kobic 국가생명연구자원정보센터 Korea Bioinformation Center

# **Genome sequencing**

## Whole-genome sequencing (WGS) process



KODIC 국가생명연구자원정보센터

KKIKK

#### As of 1<sup>st</sup> Nov. 2022, a total of **20,000 samples were sequenced**

Rare Disease (RD)

15,000 participants

Normal Participants (KoGES)

5,000 participants

#### For each sample, 138 GBase (>30X) on average was generated

WGS data were generated by Macrogen, DNALink, TheragenBio, and LabGenomics





#### Sequence data QC

For each sample, **138 GBase (>30X) on average** was generated

| Phred Q30         | 91.6% |
|-------------------|-------|
| De-duplicate read | 87.8% |
| Mappable read     | 99.8% |

| Genome Coverage | ≥ 1X  | 94.9% |
|-----------------|-------|-------|
| Genome Coverage | ≥ 10X | 94.3% |
| Genome Coverage | ≥ 30X | 86.7% |







# Data processing & Quality control

## WGS data analysis & QC

**kobic** Data processing and QC pipelines for the WGS data have been established

to call germline or somatic variants with high consistency and confidence.









## Data processing and QC pipelines for the WGS

Data processing and QC pipelines for the WGS data have been established to call germline or somatic variants with high consistency and confidence.



Germline mutation

**GATK-Spark**: MapReduce-based distributed parallel computing GATK-Spark pipeline performs <u>3~4 times faster</u> than Java-based GATK

#### Somatic mutation

**FPGA-based DRAGEN platform**, the Hardware accelerated variant analysis, enables somatic variant analysis **20 times faster** than GATK-Spark (60 min per sample)

KODIC 국가생명연구자원정보센터



#### **Research Environment Platform (REP)**

The platform provides processed WGS data and de-identified clinical information, enabling the researchers to conduct a multitude of integrative data analyses in the secured workspace.



# **Rare disease analysis**

#### Rare disease data analysis (example)

Sample\_ID: 405XXXYYYY Sex: Male Age: <7 Diagnosis : Neurodevelopmental disorders (Delayed speech and language development)

| Analysis Name (Proj<br>405XXXYY                | ject)<br>YY                          |   |  | Age<br>-    |   |                    |            |                               | Et -                         | hnicity          |                 |
|--|--------------------------------------|---|--|-------------|---|--------------------|------------|-------------------------------|------------------------------|------------------|-----------------|
| Phenotype: X-                                  | -linked alpha-thalassemi 🖣           | Age of Onset D<br>Birth - 1 Year (i) 20     | isease Prevalence Mode of Inheritance<br>00 Individuals (i) X-Linked |             | The sym<br>- bone d   | ptoms o<br>eformit | of thalass | semia<br>e <b>cially in t</b> | he face                      |                  |                 |
| Gene<br>ATRX C<br>Transcript(s)<br>NM_138270.5 | Variant<br>c.2671G>C<br>p.E891Q loss | Population Frequenc<br>Genotype:<br>Impact: | y: 0% gnomAD<br>Hom (100.00% Allele Fraction)<br>Missense            |             | <ul> <li>- dark urine</li> <li>- delayed growth and development</li> <li>- excessive tiredness and fatigue</li> </ul> |                    |            |                               |                              |                  |                 |
| Open   | < Previous                           | Next > Use Classification                   | View Bibliography  |             | - yellow  | or pale            | skin       |                               | _                            |                  |                 |
| Filter Settings 🕶                              | Search                               |   | 4913405 variants   |             |   |                    |            |                               |                              |                  | View Settings - |
| Gene   |                                      | Alteration                                  | Phenotype  | Proband (i) | Mode of Inheritance   | Function           | Impact     | CADD Score                    | Max Population Frequency     | Variant Findings | HGMD Accession  |
| ATRX   | 1                                    | c.2671G>C<br>p.E891Q                        | X-linked alpha-thalassemia-mental retard                             |             | X-linked  | loss               | Missense   | 14.53                         | 0% gnomAD                    | 74               |                 |
| FECH   | 2                                    | c.315-48T>C                                 | Erythropoietic protoporphyria  |             | recessive   | loss               | -          | <10                           | 34.34% gnomAD<br>(Latino)    | 138              | CS020196 (DFP)  |
| GJB2   | <b>U</b>                             | c.109G>A<br>p.V37I                          | Autosomal recessive deafness type 1A                                 |             | recessive   | loss               | Missense   | 21.7                          | 8.35% gnomAD<br>(East Asian) | 1346             | CM000016 (DM)   |
| NPHP   | 4 !                                  | c.2818-2A>T                                 | Nephronophthisis   |             | recessive   | loss               | Splicing   | 24.4                          | 0% gnomAD                    | 15               | CS057899 (FP)   |
| ΟΤΟΑ   |                                      | c.2122G>T<br>p.E708*                        | Autosomal recessive deafness type 22                                 |             | recessive   | normal             | Stop Gain  | 39                            | 1.03% gnomAD<br>(East Asian) | 18               | CM199073 (DM)   |
| PRMT   | 3                                    | c.1337A>G<br>p.N446S                        | Cancers and Tumors   |             | -   | loss               | Missense   | 23.1                          | 0% gnomAD                    | 2                | -               |

Rare Disease (RD)

#### 8,000 participants





## **Genome Sequencing for Undiagnosed Genetic Diseases**

#### Annual Review of Medicine

What Has the Undiagnosed Diseases Network Taught Us About the Clinical Applications of Genomic Testing?

#### David R. Murdock,<sup>1</sup> Jill A. Rosenfeld,<sup>1</sup> and Brendan Lee<sup>1,2</sup>

<sup>1</sup>Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, Texas 77030, USA; email: blee@bcm.tmc.edu

<sup>2</sup>Texas Children's Hospital, Houston, Texas 77030, USA

Genetic testing has undergone a revolution in the last decade, particularly with the advent of next-generation sequencing and its associated reductions in costs and increases in efficiencies. The Undiagnosed Diseases Network (UDN) has been a leader in the application of such genomic testing for rare disease diagnosis. This review discusses the current state of genomic testing performed within the UDN, with a focus on the strengths and limitations of whole-exome and whole-genome sequencing in clinical diagnostics and the importance of ongoing data reanalysis. The role of emerging technologies such as RNA and long-read sequencing to further improve diagnostic rates in the UDN is also described. This review concludes with a discussion of the challenges faced in insurance coverage of comprehensive genomic testing as well as the opportunities for a larger role of testing in clinical medicine.



#### Table 1 Comparison of reportable variant types detected by different genetic testing assays

| Genetic<br>test | SNVs/<br>Indels | CNVs    | Mosaic variants                  | Repeat<br>expansions | Balanced and<br>complex SVs | mtDNA       | Genes tested                                 | VUS<br>potential <sup>a</sup> | Cost <sup>a</sup> |
|-----------------|-----------------|---------|----------------------------------|----------------------|-----------------------------|-------------|--|-------------------------------|-------------------|
| Sanger          | Yes             | No      | Limited                          | No                   | No                          | If included | Single-few                                   | +                             | \$\$              |
| CMA             | No              | Yes     | Yes (CNV only)                   | No                   | No                          | No          | Up to ~20,000                                | ++                            | \$\$              |
| NGS panel       | Yes             | Yes     | Yes                              | No                   | No                          | If included | Few-hundreds                                 | ++                            | \$                |
| WES             | Yes             | Limited | Limited (depends<br>on coverage) | No                   | No                          | If included | ~20,000 (coding<br>regions only)             | +++                           | \$\$\$            |
| WGS             | Yes             | Yes     | Limited (depends<br>on coverage) | Yes                  | Yes                         | Yes         | ~20,000 (coding<br>and noncoding<br>regions) | ++++                          | \$\$\$\$          |

#### Table 2 Comparison of published RNA-seq studies and their contrasting approaches and results

| Reference | UDN<br>study | Phenotypes                           | Tissue   | Events detected              | Analysis method <sup>a</sup> | RNA-seq<br>diagnostic rate |
|-----------|--------------|--------------------------------------|--|------------------------------|------------------------------|----------------------------|
| 35        | No           | Neuromuscular                        | Muscle   | Splicing                     | Candidate + outlier          | 35%                        |
| 34        | No           | Neuromuscular                        | Muscle, fibroblasts,<br>fibroblast-derived<br>myotubes | Expression,<br>splicing, ASE | Outlier                      | 36%                        |
| 29        | Yes          | Multiple (neurologic<br>most common) | Whole blood  | Expression,<br>splicing, ASE | Outlier                      | 7.5%                       |
| 30        | Yes          | Multiple (neurologic<br>most common) | Whole blood,<br>fibroblasts, muscle,<br>bone marrow    | Expression,<br>splicing, ASE | Candidate                    | 18%                        |
| 23        | Yes          | Multiple (neurologic<br>most common) | Whole blood, fibroblasts                               | Expression,<br>splicing, ASE | Outlier                      | 17%                        |

# Long-read sequencing is necessary for the undiagnosed rare diseases



Murdock (2022) Ann Rev Med



# The Korean Pan-genome Project

## The complete sequence of a human genome







#### **T2T consortium to the Human Pangenome Project**







#### Needs for a pan-genome (genome graph)

genetics

# Fast and accurate genomic analyses using genome graphs

TECHNICAL REPORT

https://doi.org/10.1038/s41588-018-0316-4

KRIBB



Rakocevic (2020) Nat Genet

Annual Review of Genomics and Human Genetics The Need for a Human Pangenome Reference Sequence

#### Karen H. Miga<sup>1</sup> and Ting Wang<sup>2</sup>

<sup>1</sup>UC Santa Cruz Genomics Institute and Department of Biomedical Engineering, University of California, Santa Cruz, California 95064, USA; email: khmiga@ucsc.edu

<sup>2</sup>Department of Genetics, Edison Family Center for Genome Sciences and Systems Biology, and McDonnell Genome Institute, Washington University School of Medicine, St. Louis, Missouri 63110, USA; email: twang@wustl.edu



KODIC 국가생명연구자원정보센터

#### **Construction of the human pan-genome**



2. Construction of complete reference genome (template)

1. De novo assembly of

5. Pan-genome based alignment

3. Identification of variants

(SNVs, small INDELs, & SVs)



Science (2021), Genome biology (2020)



#### Roadmap of the Korea pan-genome project



During the phase 1, we performed *de novo* assembly of the human WGS data, generated from the long-read sequencing technique (PacBio HIFI)







 The Korea National Genome Project (NGP) enables the researchers to explore clinical and genomic data from 23,833 participants (25,000 samples) in the secure cloud service (Research Environment Platform)

KOBIC launched the Korean Pan-genome Project
 for the better performance to identify the uncharted structural variations
 and integrate numerous genetic variants from Korean population





### Acknowledgement

|   | Dr. Seon-Young Kim<br>Jin Ok Yang<br>Ik Su Byeon<br>Dr. Jong-Hwan Kim<br>Dr. Kiwon Jang<br>Jae Ho Lee<br>Gun Woo Park | Dr. Byungwook Lee<br><b>Dr. Sang-Ok Kim</b><br>Bang-Hyuck Lee<br>Dr. Jongbum Jeon<br>Dr. Wonyong Jeong<br>Young Mi Sim<br>Dr. Seon-Kyu Kim | Dr. Pan-Gyu Kim<br>GunHwan Ko<br>Jong-Cheol Yoon<br>Dr. Jaeeun Jung<br>Taeyeon Hwang<br>Dongmin Jang<br><b>Dr. Soobok Joe</b> |
|---|---|--|---|
| Kişti<br>www.kisti.re.kr                            | Dr. Jun-Hak Lee   | Dr. Hyojin Kang  | Dr. Yukyung Jun   |
| KHIDI   | Dr. Kwan Ik Lee   | Dr. Misook Kwak  | Dr. Jong-suk Park   |
| <b>Keit</b><br>한국산업기술평가관리원                          | Dr. Kang-Woo Lee  | Dae Sung Kim   |   |
| 질병관리청<br>KDCA                                       | <b>Dr. Hyun-Young Park</b><br>Dr. Myungguen Chung<br>Namhee Kim   | Dr. Hee Youl Chai<br>Dr. Min Jin Go  | Dr. Eugene Kim<br>Dr. Jung-Eun Kim  |
| <mark>사 울 대 학 교</mark><br>SEOUL NATIONAL UNIVERSITY | Dr. Murim Choi  | Dr. Jun Kim  | Jeongeun Lee  |
| www.korea UNIVERSITY                                | Dr. Joon-Yong An  | Lizzy Choi   | Ganghee Lee   |
|   | NGP co  | hort Partnership   |   |
| SE  |   | <b>YONSEI</b><br>UNIVERSITY CHOSUN<br>UNIVERSITY   | 국립암센터   |
|   | Support   | ted & Funded by:   |   |
| Ministry  | of Health Ministry of S   | cience and ICT Ministry of T   | rade.   |

and Welfare

**KKIBB** 



Industry and Energy



# Thank you!



