# NGS를 이용한 게놈분석

#### Jong Bhak 박종화 게놈연구재단

2012.05.02



## Let's do more sequencing

# 시퀀싱 더 하자!

By Jong Bhak

결론 2.

# Sequencing 7 billion people on Earth as cheaply as possible

# 70억명 해독

By Jong Bhak

## 감사의 말

- 연구자체에 열정을 가지고 정확한 데이터를 내는 많
   은 국내외 과학자들
- -과학기술인을 지원하는 국민들
- -게놈관련 상품을 사주는 고객들
- 성신여대 와 김상태교수님
- 테라젠 및 게놈재단의 많은 동료들
- -재정적 지원(기부)를 해준 테라젠 (고진업대표)



### Gene, Environe, Phene (유전자, 환경자, 형질자)

- Gene: single genetic factor (genetic atom)
- Environe: single environmental factor (temp.)
- Phene: single phenotype (trait: hair color)
- Genome, Envirome, Phenome (Traitome)



#### Single Gene . Environe . Phene Variation

![](_page_7_Figure_1.jpeg)

#### **Environe Variation**

Jong Bhak. Under BioLicense: public domain

# **Genome Envirome and Phenome**

- Genome = gene types + their variome
- Envirome = environe types + their variome
- Phenome = phene types + their variome

## GenoEnviroPheno Unpredictability Graph

• Genome

Phenome

![](_page_9_Figure_3.jpeg)

![](_page_9_Picture_4.jpeg)

Jong Bhak. Under BioLicense: public domain

## Genecomplex $\leftarrow \rightarrow$ Phenecomplex

• GeneComplex

![](_page_10_Figure_2.jpeg)

#### PheneComplex

Jong Bhak. Under BioLicense: public domain

![](_page_11_Figure_0.jpeg)

#### **Biomatrix: Putting structure in omics** BiO Ome Interactome Textome Transcript Proteome **Functome** Genome **Matrix** file **Pipeline** Public DB Auto Update & Dump to MySOL OITEK Sequence CDI UniPr at MEDLINE PSIbase InterPi o OMIM GO Resource External I Uniaene NCBI Structure EC SCOL dbEST dbSNU DDBJ POB РЕал GEO Locu Link . . . . . . . . . **Expression** Enst MBL SMAR Pipeline 분석 Pathway Data Full length cONA Link Reference ntemal NGIC Preproce ss Regulation PPI prediction **Materials BioEngine** Annotation D В Network Domain Assignment BioDiversity Public Service 소재은행 Info Type Portal 2ndary DB 2ndary Info

# 게놈학이란?

 Genomics is the bioinformatic study of genes of individual organisms, populations, and species.

http://genomics.org

# DNA 해독 and DNA 타이핑

해독 → 궁극적

![](_page_14_Picture_2.jpeg)

• 타이핑 → cheap and efficient

![](_page_14_Picture_4.jpeg)

# 차세대 해독기

- 일루미나
- 라이프테크 (구 AB)
- 로슈
- IBS
- CompleteGenomics
- PacBio
- Helicos
- NanoPore
- Genia

# **Genomic T**:

![](_page_16_Figure_1.jpeg)

#### **Genome Diversity and Genome Variation**

![](_page_17_Figure_1.jpeg)

#### nsSNVs in proteins

![](_page_18_Figure_1.jpeg)

#### The first Korean Genome Sequence

Downloaded from genome.cshlp.org on August 11, 2009 - Published by Cold Spring Harbor Laboratory Press

![](_page_19_Picture_2.jpeg)

### The first Korean genome sequence and analysis: Full genome sequencing for a socio-ethnic group

Sung-Min Ahn, Tae-Hyung Kim, Sunghoon Lee, et al.

*Genome Res.* published online May 26, 2009 Access the most recent version at doi:10.1101/gr.092197.109

#### Data publicized: 2008. Dec. ftp://bioftp.org

## The first Korean Genome (SJK)

- First analyzed by Gacheon medical school LCDI and KOBIC, KRIBB in 2008 (Joint effort among LCDI, KOBIC, and 국가참조표준센터)
- First annotated and made public on 4<sup>th</sup> Dec. 2008 (through web and ftp)
- SNP, CNV, indels were analysed
- Automated phenotypic association study was done
- Non-syn. Analysis
- Phylogenetic study of mtDNA, Y Chr And autosomes showed Korean relationship to Chinese and Japanese.
- First intra-Asian genome comparison (Chinese and Korean)
- Analyzed at: 7.8, 17.3, 23.5 and 28 x folds
- By Jan. 23.5 fold sequenced and analyzed
- Openfreely Available from: <u>http://koreagenome.org</u>

## The Karyogram of the donor DNA

No obvious chromosomal abnormalities!

![](_page_21_Figure_2.jpeg)

#### **Classification and number of intra-genic SNPs**

![](_page_22_Figure_1.jpeg)

#### **Comparison of individual SNPs**

![](_page_23_Figure_1.jpeg)

# Size distribution and classification of short indels found in SJK

![](_page_24_Figure_1.jpeg)

Using MAQ, we identified 342,965 short indels  $\rightarrow$  We found that only 247 (0.1%) were validated,113,287 (33.0%) non-validated, and 229,431 (66.9%) indels were not found in dbSNP

#### Indels in SJK genic regions

	Indel				
Index	Indel number	Homozygous	Heterozygous	Gene	
				number	
5'UTR	27	9	18	26	
CDS	49	16	33	40	
3'UTR	319	114	205	247	
Intron	127,516	45,430	82,086	12,421	
Total	127,911	45,569	82,342	12,734	

# Homo– and heterozygous deletions in SJK genome

![](_page_26_Figure_1.jpeg)

(A) Homozygous 2.3 kb genomic deletion and (B) Heterozygous 5 kb genomic deletion.

#### Detection and identification of structural variants

- We found structural variants by using paired-end reads.
  1.2920 deletions (100bp ~ 100kb)
  2.415 inversions (100bp ~ 100kb)
  3.963 insertions (175bp ~ 250bp)
- We found deletion SVs in 21 coding genes.
   → All heterozygous deletions

# Repeat composition in SJK deletion variants

![](_page_28_Figure_1.jpeg)

#### Soybean 유전체 분석

#### **Overview**

- Sequencing: *G.max* 품종 1, 품종 2 - Illumina GAIIx
- Reference : G.max Williams 82 (Glyma1)
- Tools : bwa, samtools, breakdancer
- 분석항목
  - 1차분석
    - 레퍼런스 맵핑 통계
  - 2차분석
    - SNP
    - Small Indel
    - SV
  - 3차분석
    - Comparison
    - Pattern analysis
    - Pathway mapping

![](_page_31_Figure_0.jpeg)

#### **SNV** density profile

![](_page_32_Figure_1.jpeg)

![](_page_33_Figure_0.jpeg)

![](_page_34_Figure_0.jpeg)

#### 2. Bioinformatics Challenge :

### still massively mapping

#### How to map/compute 6 billion X 6 billion matrix?

![](_page_35_Figure_3.jpeg)

#### Adding one more dimension?

#### How to map/compute **RNA** expressions In relation with bio-function?

![](_page_36_Figure_2.jpeg)

#### Adding even more dimension?

#### How to map/compute Phenome?

![](_page_37_Figure_2.jpeg)

#### How to map/compute epigenome?

![](_page_38_Figure_1.jpeg)

#### How to map/compute Microbiome?

![](_page_39_Picture_1.jpeg)

# PGP and KPGP (게놈연구재단 프로젝트)

- "한국인게놈프로젝트"
- Doing Biology with Sequencing

### **Personal Genome Project (PGP)**

#### Public Open Source Genome Project

>Volunteers from the general public working together with res earchers to advance personal genomics.

 Led by Prof. George Church in Harvard Medical School
 100,000 informed participants from the general public (US Citizen).

► Research Data freely available to the public

![](_page_41_Picture_5.jpeg)

![](_page_41_Figure_6.jpeg)

#### **Mission**

Personal Genome Project is to encourage the development of personal genomics technology and practices that:

- > are effective, informative, and responsible
- > yield identifiable and improvable benefits at manageable lev els of risk

> are broadly available for the good of the general public

![](_page_42_Picture_0.jpeg)

![](_page_42_Picture_1.jpeg)

# KOREAN PERSONAL GENOME PROJECT

PERSONAL GENOMICS INSTITUTE

![](_page_42_Picture_4.jpeg)

![](_page_43_Figure_0.jpeg)

이용약관 | 개인보호정책 | QnA | 오시는길

TEL:031-888-9317 FAX:031-888-9314 E-MAIL:info@kpgp.kr

#### Korean Personal Genome Project (KPGP)

- > Extension of Harvard PGP Project in Korea
- > Led by Jong Bhak in Personal Genomics Institute in Korea
- > Period : 2007 -2017

#### ≻Plan

- 1단계 2007년 ~ 2009년, 1명 (SJK): Korean Reference
- 2단계 2010년 ~ 2011년, 100명
- 3단계 2012년 ~ 2013년, 3000명
- 4단계 2014년 ~ 2017년, 10000명

![](_page_44_Picture_9.jpeg)

![](_page_44_Picture_10.jpeg)

![](_page_44_Picture_11.jpeg)

![](_page_44_Picture_12.jpeg)

![](_page_44_Picture_13.jpeg)

#### Korean Genome Project Vision & Mission

Sequencing every single Korean on Earth

![](_page_45_Picture_2.jpeg)

Personalized Medicine using Genomics

![](_page_45_Figure_4.jpeg)

#### Mission

1. Korean standard genome information DB

2. Open and sharing genome project

3. New technology development for personal genomics and medicine.

4. Community formation for genomic issues in society (legal, ethical)

#### KT 한국인 게놈 프로젝트 참여

![](_page_46_Picture_1.jpeg)

KT, 한국인 게놈 연구에 클라우드 컴퓨팅 제공

![](_page_46_Picture_3.jpeg)

![](_page_46_Picture_4.jpeg)

#### KPGP–20 Korea Telecommunication (KT) & Theragen

![](_page_47_Figure_1.jpeg)

#### **Sequencing & Analysis Summary**

FTP directory /BiO/Distrib	ute/Open-KPGP/ at bioftp.org - Window	v			
🕞 🔵 🗢 🙋 ftp://bioff	tp.org/BiO/Distribute/Open-KPGP/				
☆ Favorites		• Sample	Sample : Blood Genomic DNA		
FTP director	y /BiO/Distribute/Op	• Sequenc	ing Platform	: Illumina HiSeq 2000	
To view this FTP site in Windows Explorer, click P		• Data ar	Data analysis tool information		
Up to higher level di	rectory	Raw dat	a	fasta	
09/16/2011 01:40AM	Directory TGP2010D0001		a		
09/16/2011 01:40AM	Directory <u>TGP2010D0005</u> Directory TGP2010D0006	<ul> <li>mapping</li> </ul>	g result	bwa(0.5.9)	
09/16/2011 01:40AM	Directory TGP2010D0007	<ul> <li>SNIV</li> </ul>		samtools(0116)	
09/16/2011 01:40AM	Directory TGP2010D0009				
09/16/2011 01:40AM	Directory TGP2010D0010	• INDEL		samtools(0.1.16)	
09/16/2011 01:40AM	Directory TGP2010D0012	• \$\/		breakdancer (11)	
09/16/2011 01:40AM	Directory TGP2010D0013	50			
09/16/2011 01:40AM	Directory <u>IGP2010D0014</u> Directory <u>IGP2010D0015</u>	<ul> <li>NSSNV</li> </ul>		GRF's pipelin	
09/16/2011 01:40AM	Directory TGP2010D0016		CC		
09/16/2011 01:40AM	Directory TGP2010D0017	• 51A11511	C2	GRES pipeline	
09/16/2011 01:40AM	Directory <u>IGP2010D0018</u> Directory <u>IGP2010D0020</u>	<ul> <li>Reference</li> </ul>			
09/16/2011 01:40AM	Directory TGP2010D0021	• NEIEIEIIC			
09/16/2011 01:40AM	Directory TGP2010D0022				
09/16/2011 01:40AM	Directory <u>TGP2010D0023</u> Directory TGP2011D0010				
09/16/2011 01:40AM	Directory TGP2011D0011	<ul> <li>Open Da</li> </ul>	atabase :		
09/16/2011 01:40AM	Directory TGP2011D0012	ftma / /la	after and /D:	O/Distribute (Onese KDCD	
09/16/2011 01:40AM	Directory TGP2011D0013	<u>ttp://b</u>	lottp.org/Bl	<u>O/Distribute/Open-KPGP</u>	
09/16/2011 01:40AM	Directory TGP2011D0014				
09/16/2011 01:40AM	Directory TGP2011D0018				
09/16/2011 01:40AM	Directory TGP2011D0019				
09/16/2011 01:40AM	Directory TGP201100020				
09/16/2011 01:40AM	Directory TGP2011D0022				
09/16/2011 01:40AM	Directory TGP2011D0023				
09/16/2011 01:40AM	Directory <u>TGP2011D0024</u>				

## **KPGP–20 Admixture view**

![](_page_49_Figure_1.jpeg)

#### **KPGP–20 Cumulative Variation Analysis**

![](_page_50_Figure_1.jpeg)

### **Sample Size Requirements for Detection**

Number of People (N)	Population Frequency ( <i>f</i> )	Probability of Detection (p)
	0.1	0.9852
	0.05	0.8715
50	0.05	0.9941
50	0.02	0.8674
400	0.02	0.9824
	0.01	0.8661
200	0.01	0.9821
200	0.005	0.8653

# PAPGI (PASNP 2.0): 범아시아 집단 유전체 프로젝트

게놈연구재단

**Pan Asian Population Genomics Initiative** 

#### **Useful links**

- <u>http://pgi.re.kr</u>
- http://www4a.biotec.or.th/PASNP
- <u>http://papgi.org</u>
- <u>http://opengenome.net</u>
- http://genomics.org
- <u>http://personalgenome.org</u>
- <u>http://personalgenome.net</u>