

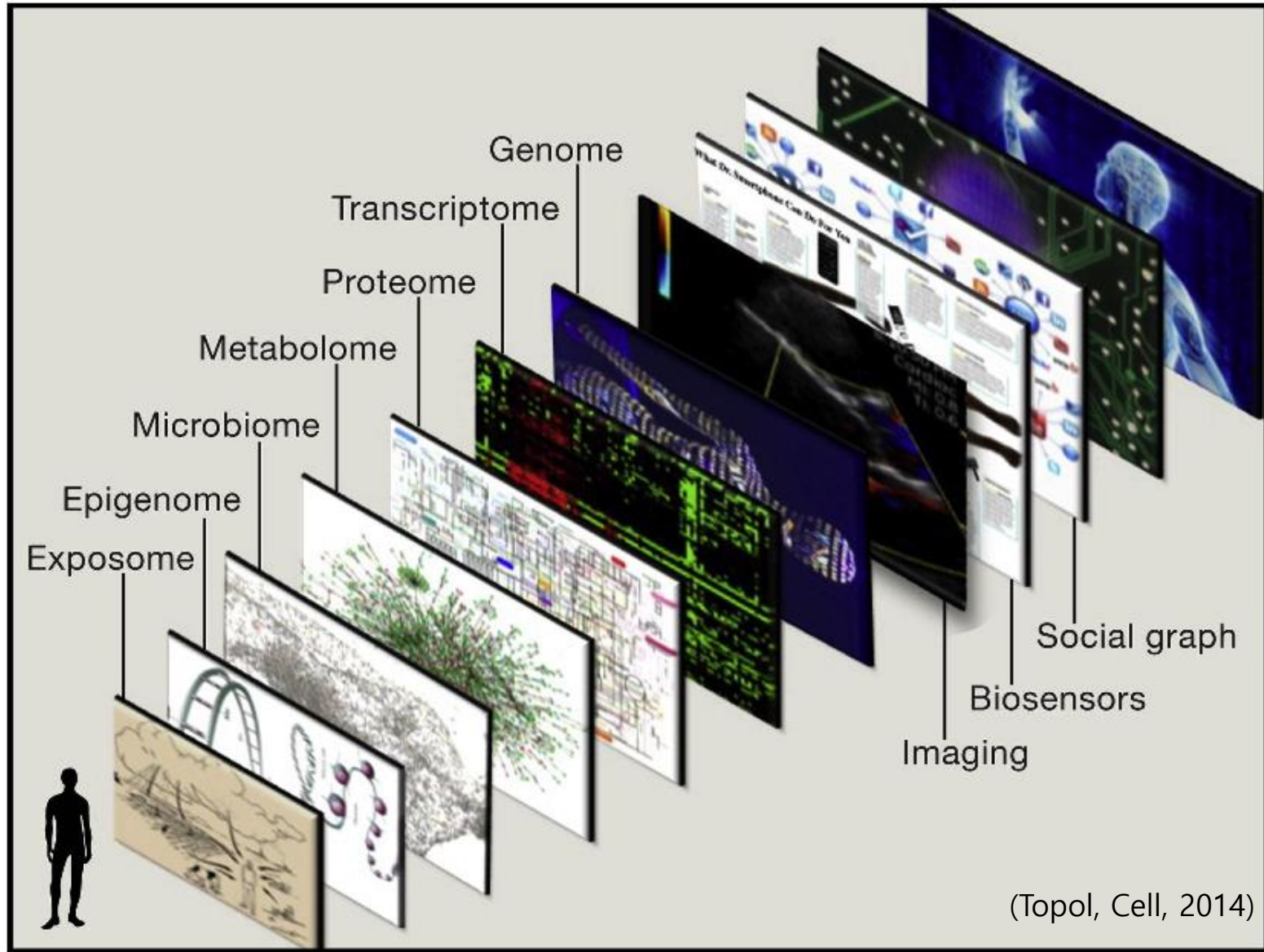


유전체 빅데이터 현황 및 응용사례

차의과학대학교 분당차병원

정제균

Individualized Medicine from Prewomb to Tomb



인간이 평생 만들어내는 데이터의 종류와 크기

Exogenous data

(Behavior, Socio-economic, Environmental, ...)

60% of determinants of health
Volume, Variety, Velocity, Veracity

Genomics data

30% of determinants of health
Volume

Clinical data

10% of determinants of health
Variety



1100 Terabytes
Generated per lifetime

6 TB
Per lifetime

0.4 TB
Per lifetime

Source: "The Relative Contribution of Multiple Determinants to Health Outcomes", Lauren McGover et al., *Health Affairs*, 33, no.2 (2014)

유전체 지도

BIOINFORMATIK - DataBanks in the WorldWideWeb



The ultimate challenge - the human genetic blueprint

바이러스 게놈
3000 염기
1페이지 분량의 정보



박테리아 게놈
3백만 염기
1000 페이지 책 1권 분량

인간 게놈
30억 염기
책 1000권 도서관



DNA Chip vs. NGS

DNA Chip (microarray)



Next Generation **Sequencing machines** (NGS)



IT·BT의 발전



IT



1990 년 2001 년

인간 유전체 프로젝트

\$2.7 billion

10년 소요



BT

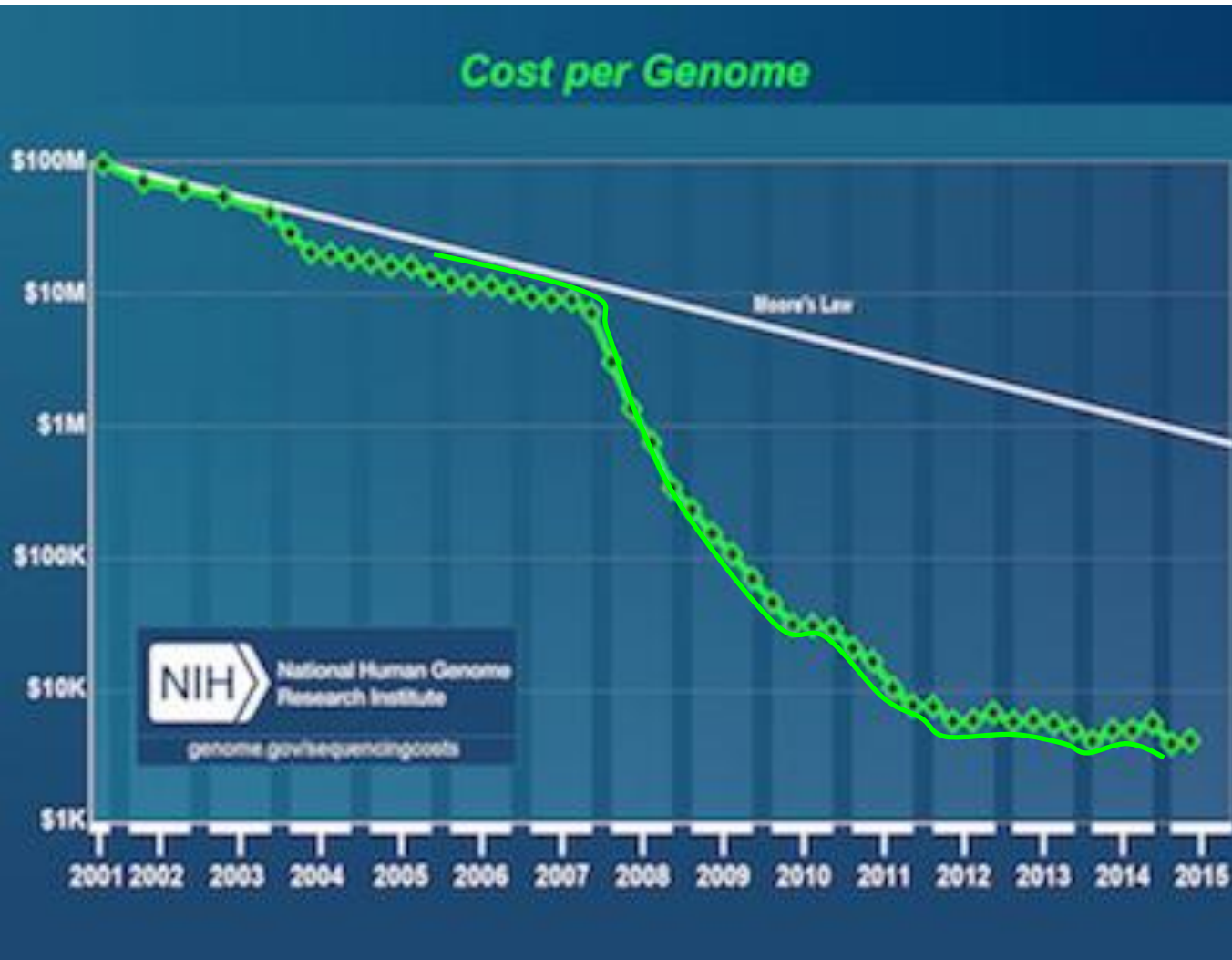
2014 년

\$1,000

수시간 소요



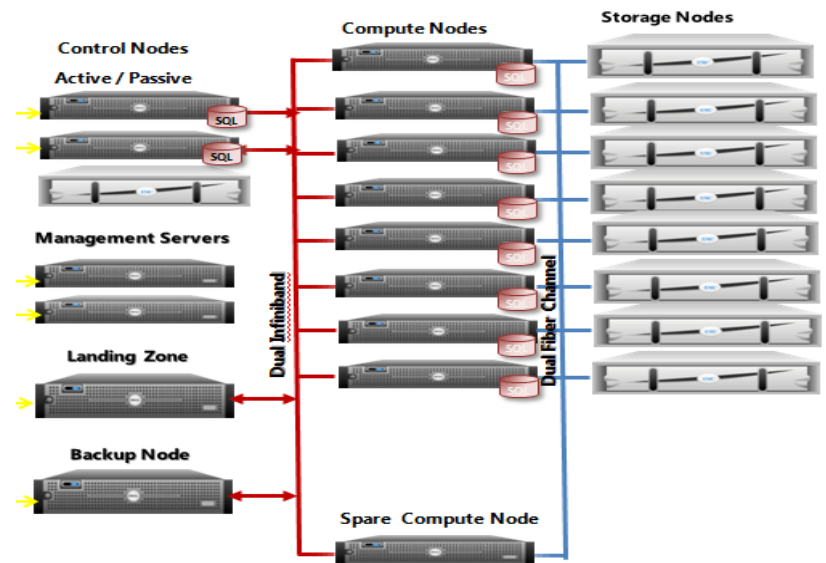
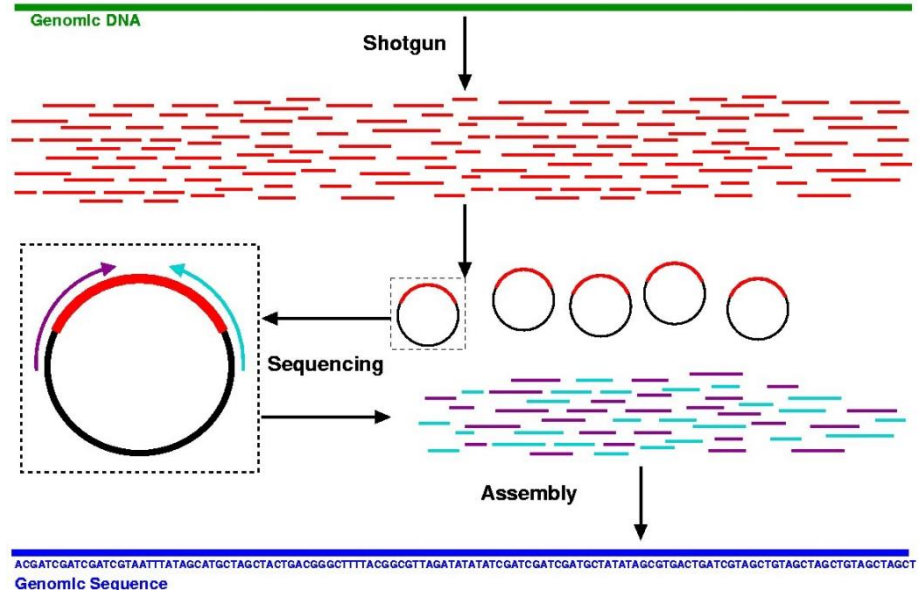
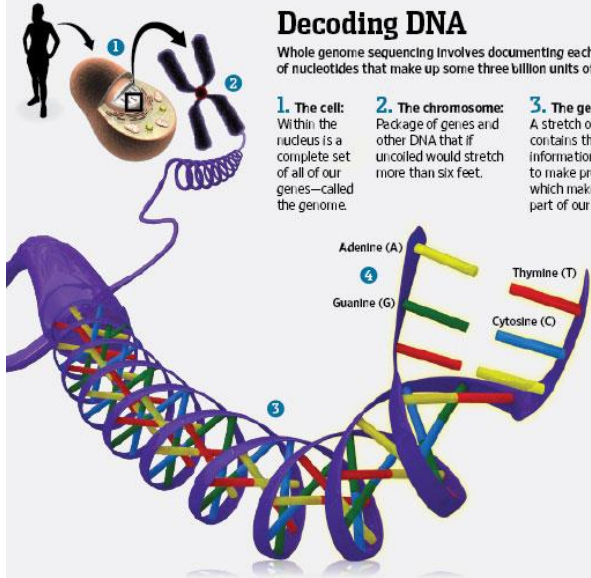
유전체 시퀀싱 가격 추이



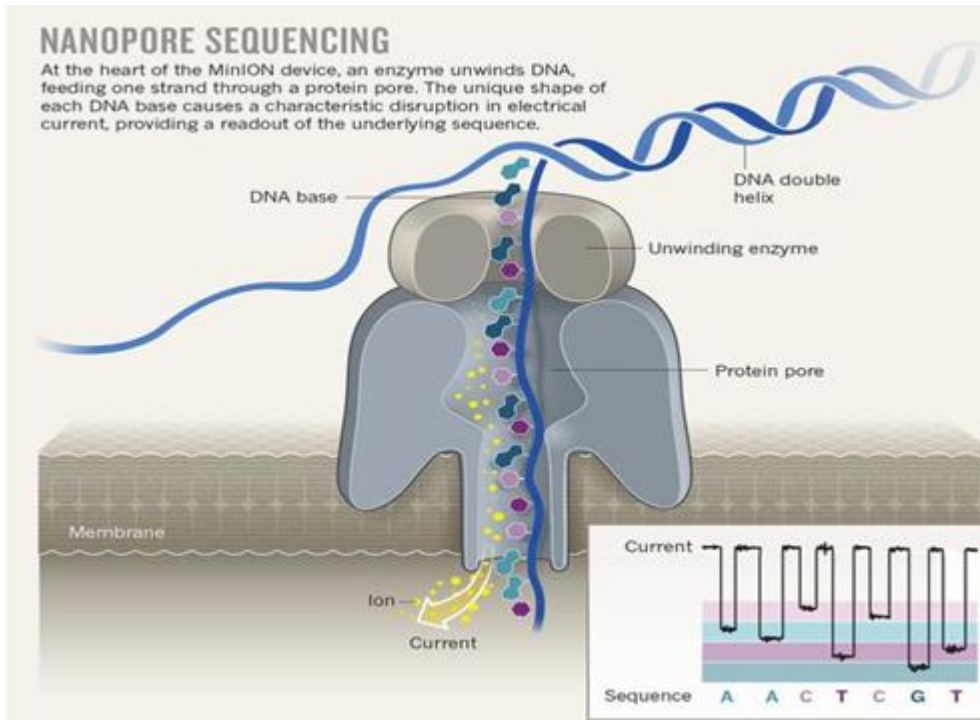
Decoding DNA

Whole genome sequencing involves documenting each of the so-called base pairs of nucleotides that make up some three billion units of DNA in our genetic code.

- 1. The cell:** Within the nucleus is a complete set of all of our genes—called the genome.
- 2. The chromosome:** Package of genes and other DNA that if uncoiled would stretch more than six feet.
- 3. The gene:** A stretch of DNA that contains the information necessary to make proteins, which make up each part of our bodies.
- 4. The bases:** The base pairs always come together in the same way—A with T and G with C. But the sequences along the molecule vary, encoding the genetic information.



나노포어 시퀀서



나노포어 막단백질을 이용 전류 신호 측정을 통해서 ATGC 염기서열을 알아내는 방식

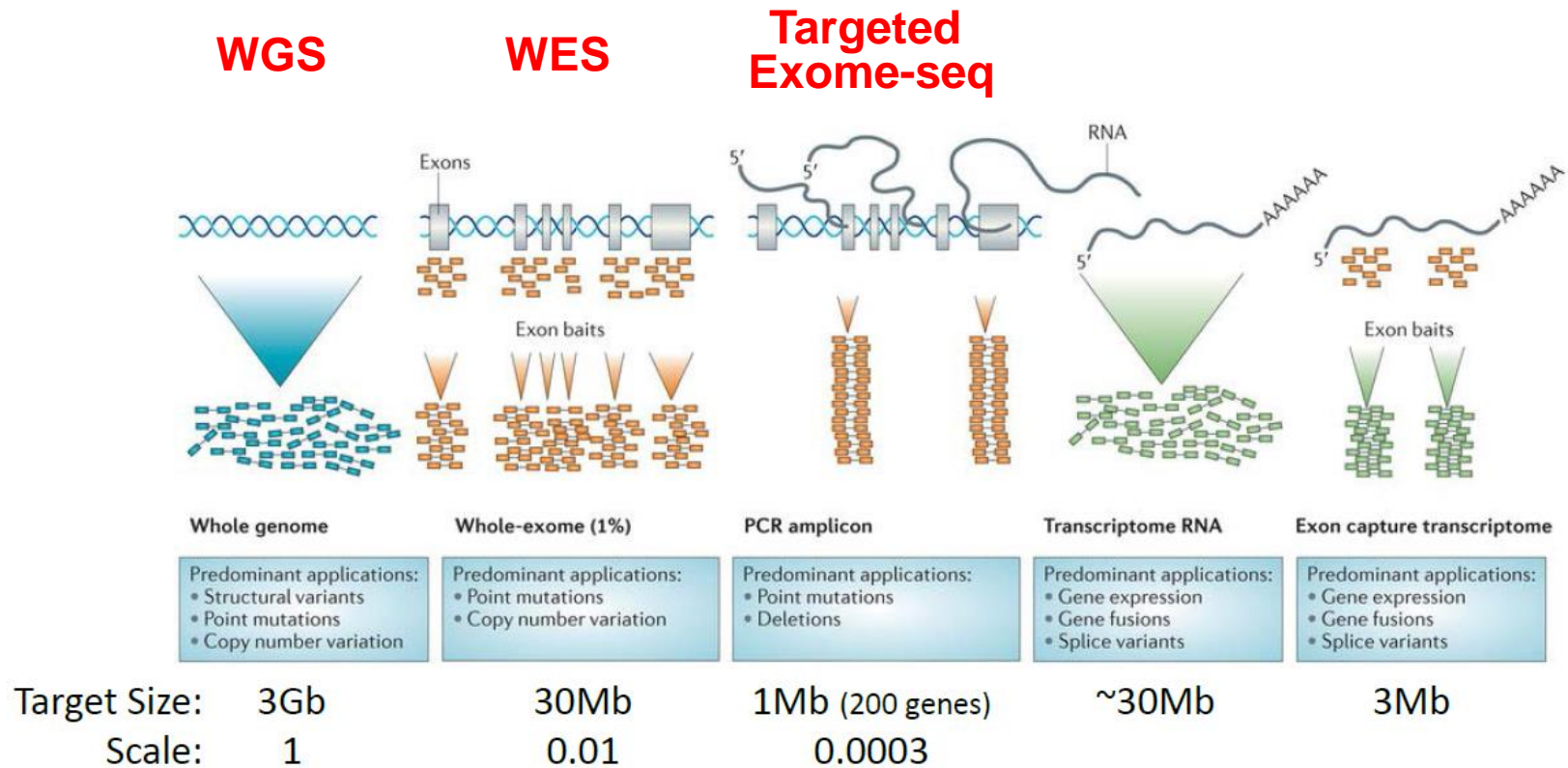


MinION Mk1: portable, real time biological analyses

MinION



연구 대상별 유전체 분석 방법



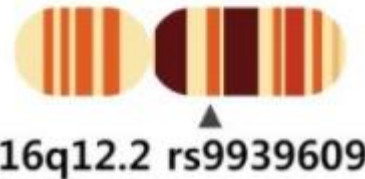
File size for human germline WGS (30X)

- Image Data 16TB
- BaseCall/Quality score data 200GB
- Final Alignment output 1 TB

유전자를 통한 질병 예측?



FTO 유전자

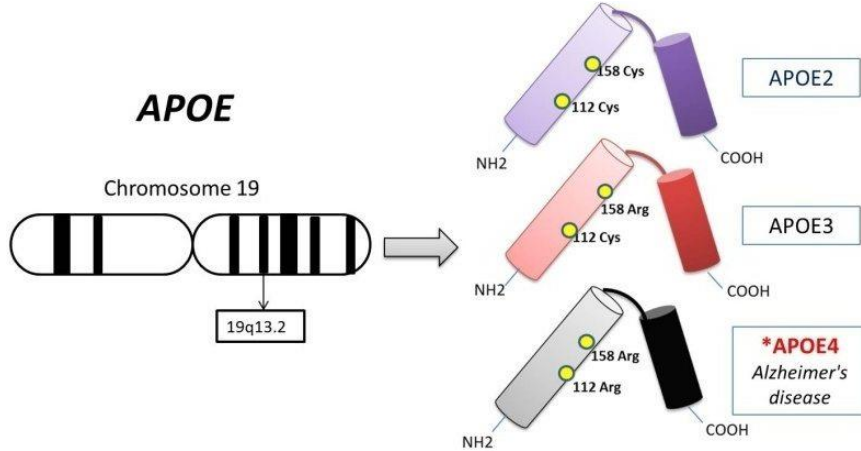
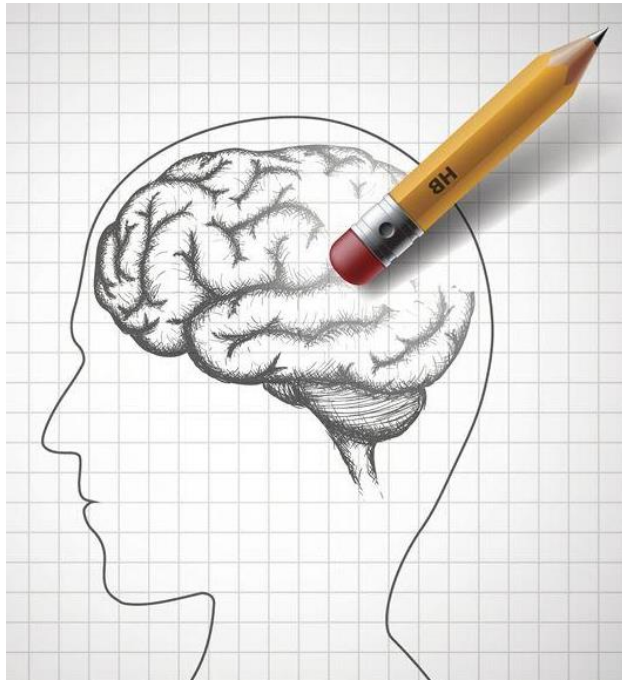


체온조절과 섭식조절을
담당하는
뇌의 시상하부에 위치

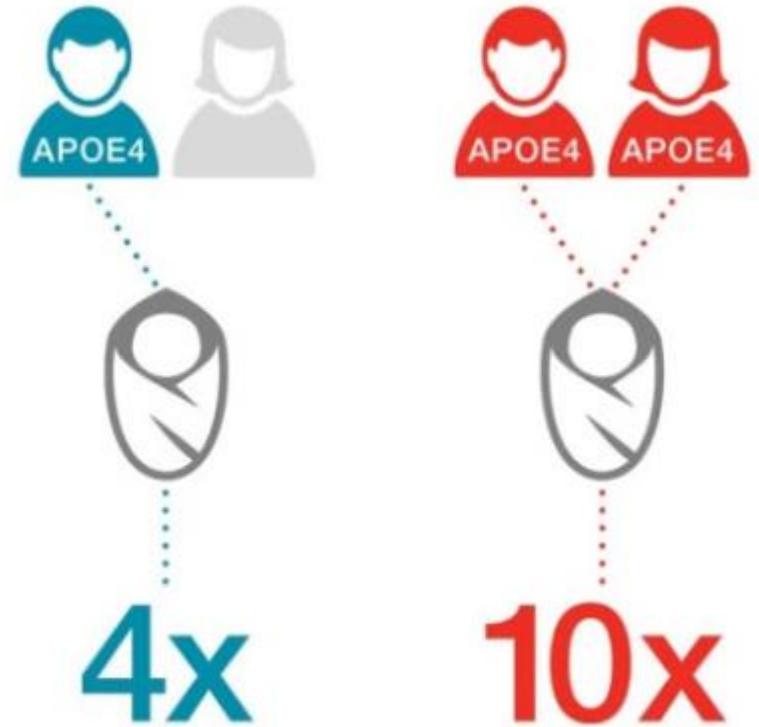
FTO 유전자는 크로모솨 16q12.1에 위치한 유전자입니다.
여러 비만 관련 유전자 중 가장 먼저 발견되어, '비만 유전자'로 불립니다.
BMI(Body Mass Index), 비만위험도, 제2형 당뇨병과 관련이 있습니다.

특정 변이가 생길시
식욕증가, 포만감 감소, 지방세포 에너지 소모 비율이 감소하고
충동적 성격이 강해지며, 지방이 많은 음식을 더 많이 섭취하게 됩니다.

유전자를 통한 질병 예측?

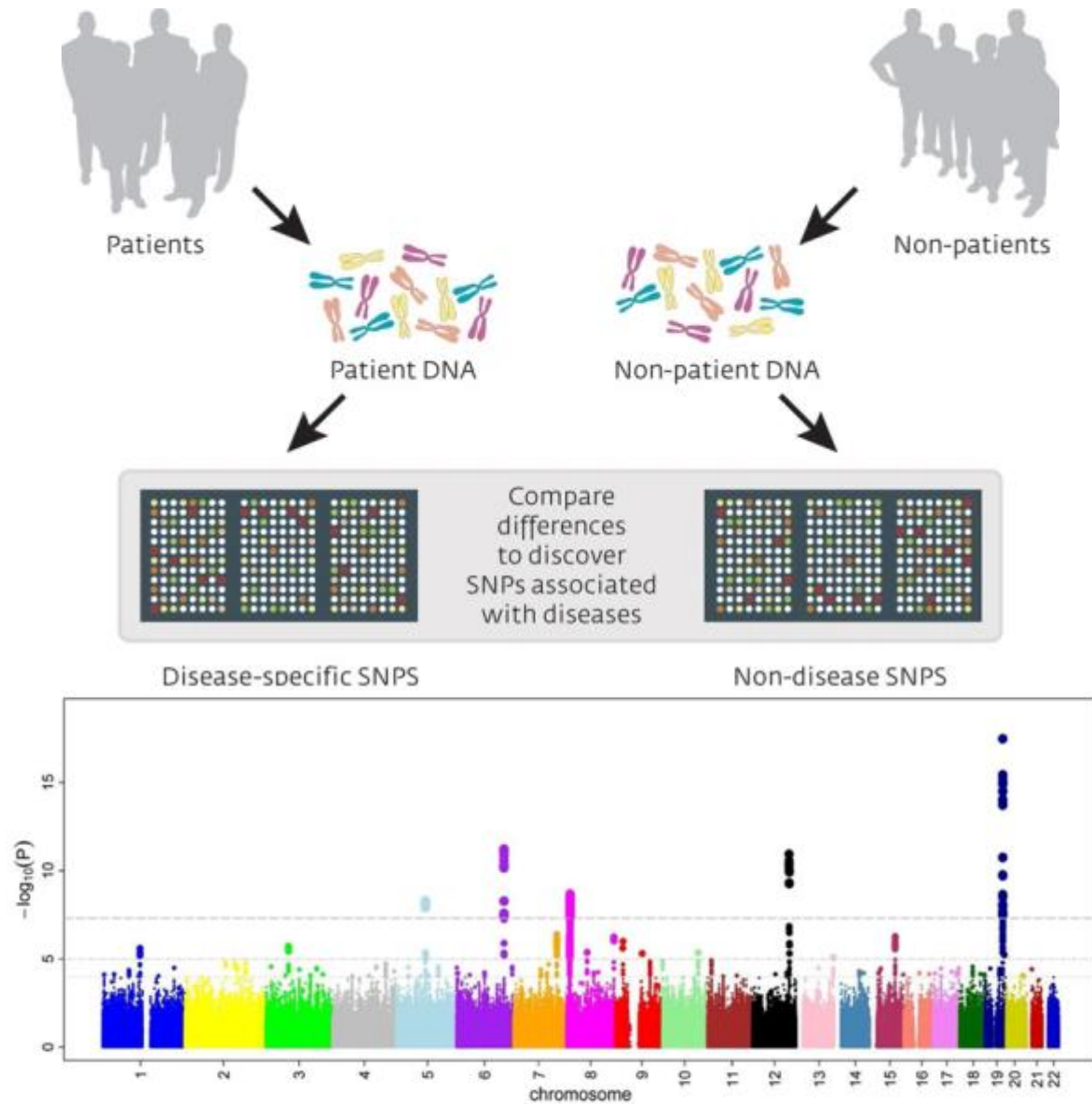


부모 양쪽이 모두 E4변이 일때 무려 10배의 확률이 올라감



(Cure Alzheimer's Fund)

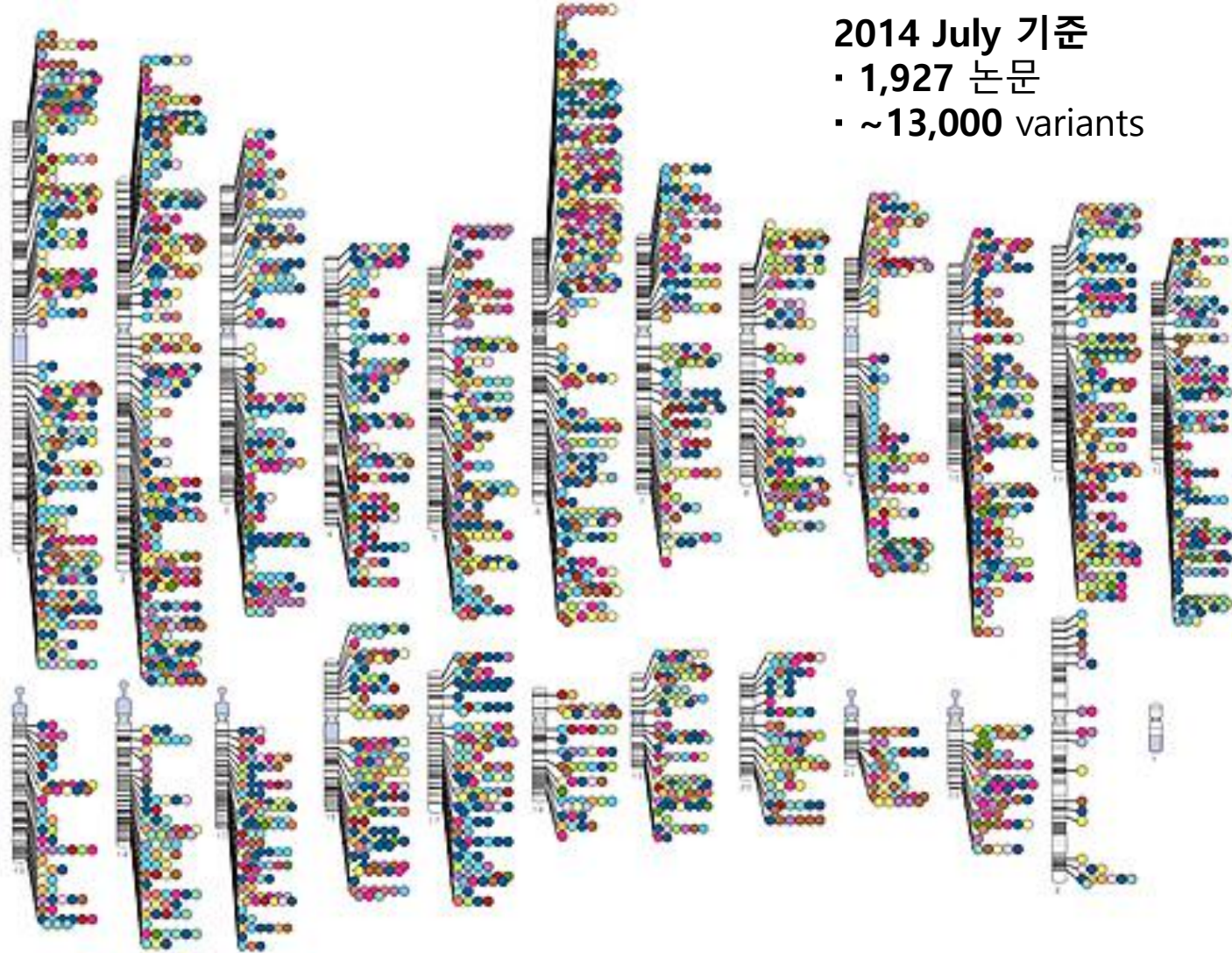
전장유전체 연관분석 (GWAS: Genome-wide association studies)



GWAS (Genome-wide association studies) Catalog

2014 July 기준

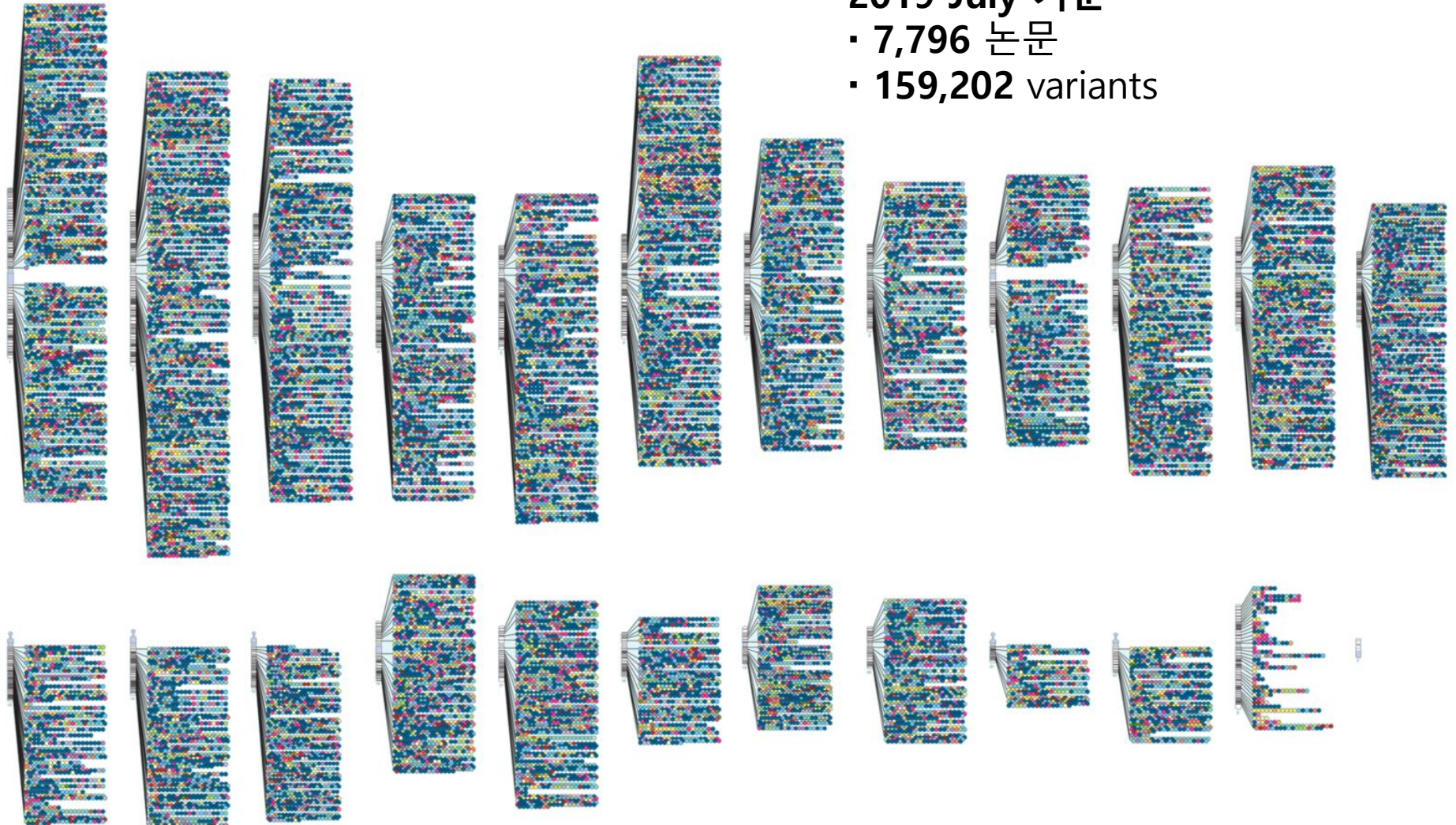
- 1,927 논문
- ~13,000 variants



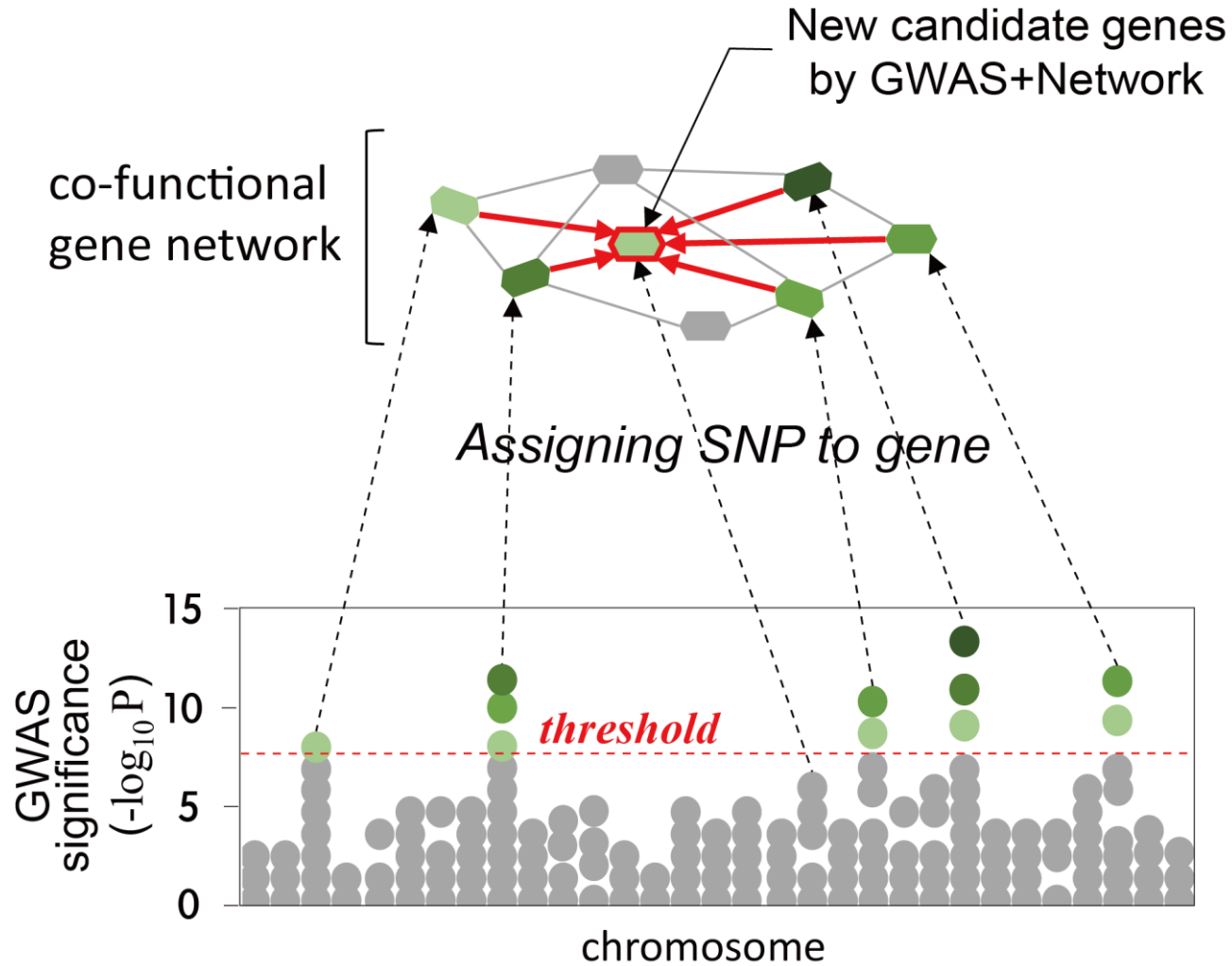
2019 July 기준

· 7,796 논문

· 159,202 variants



Polygenic causes

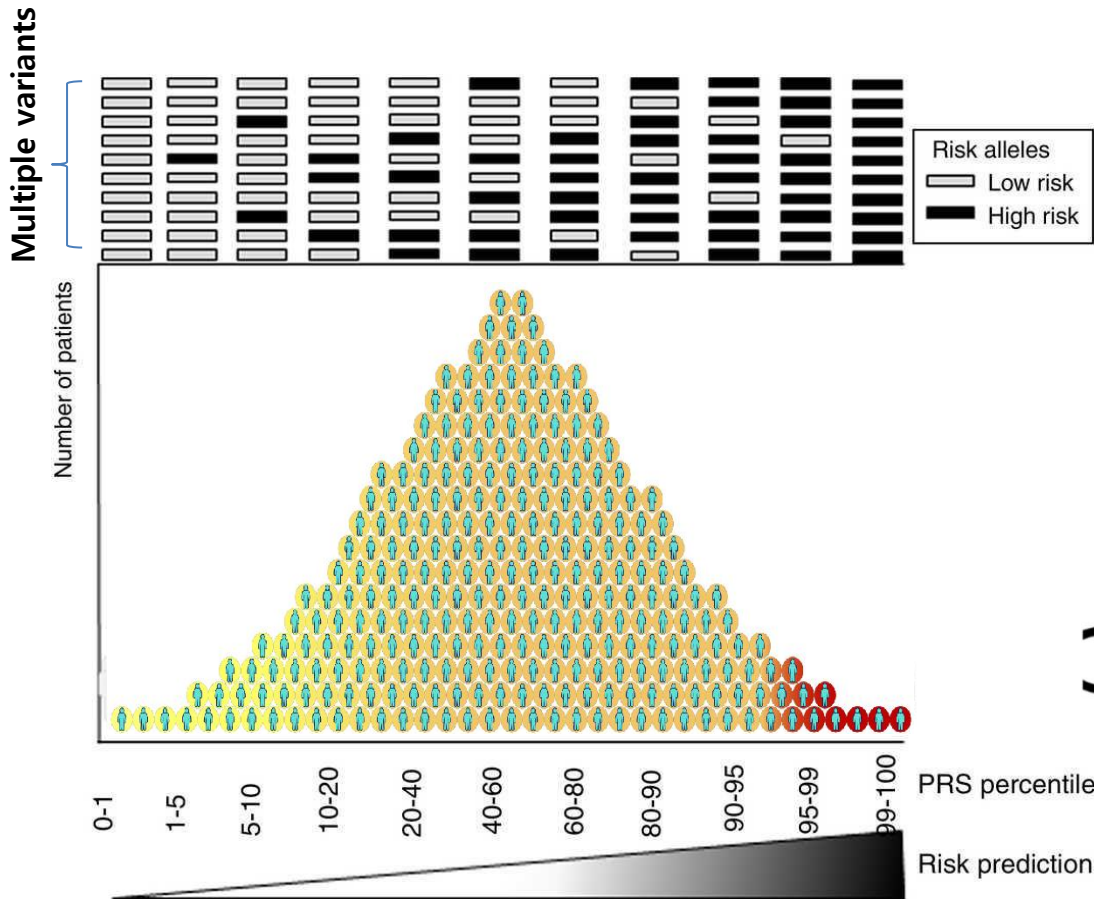
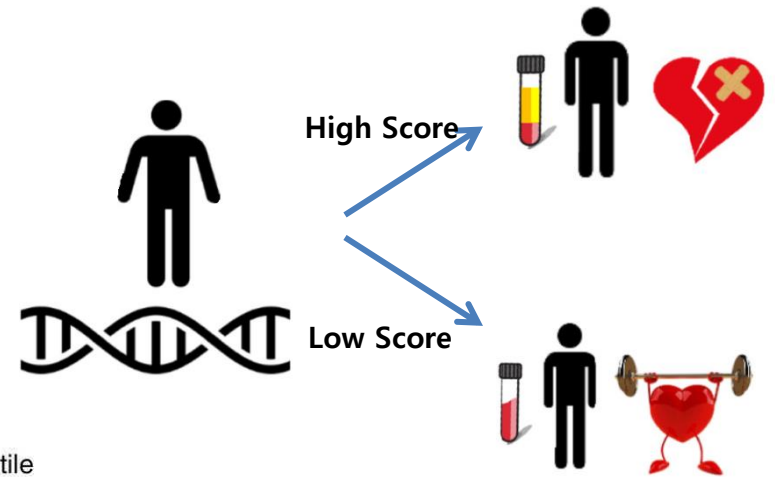


질병 위험도 예측

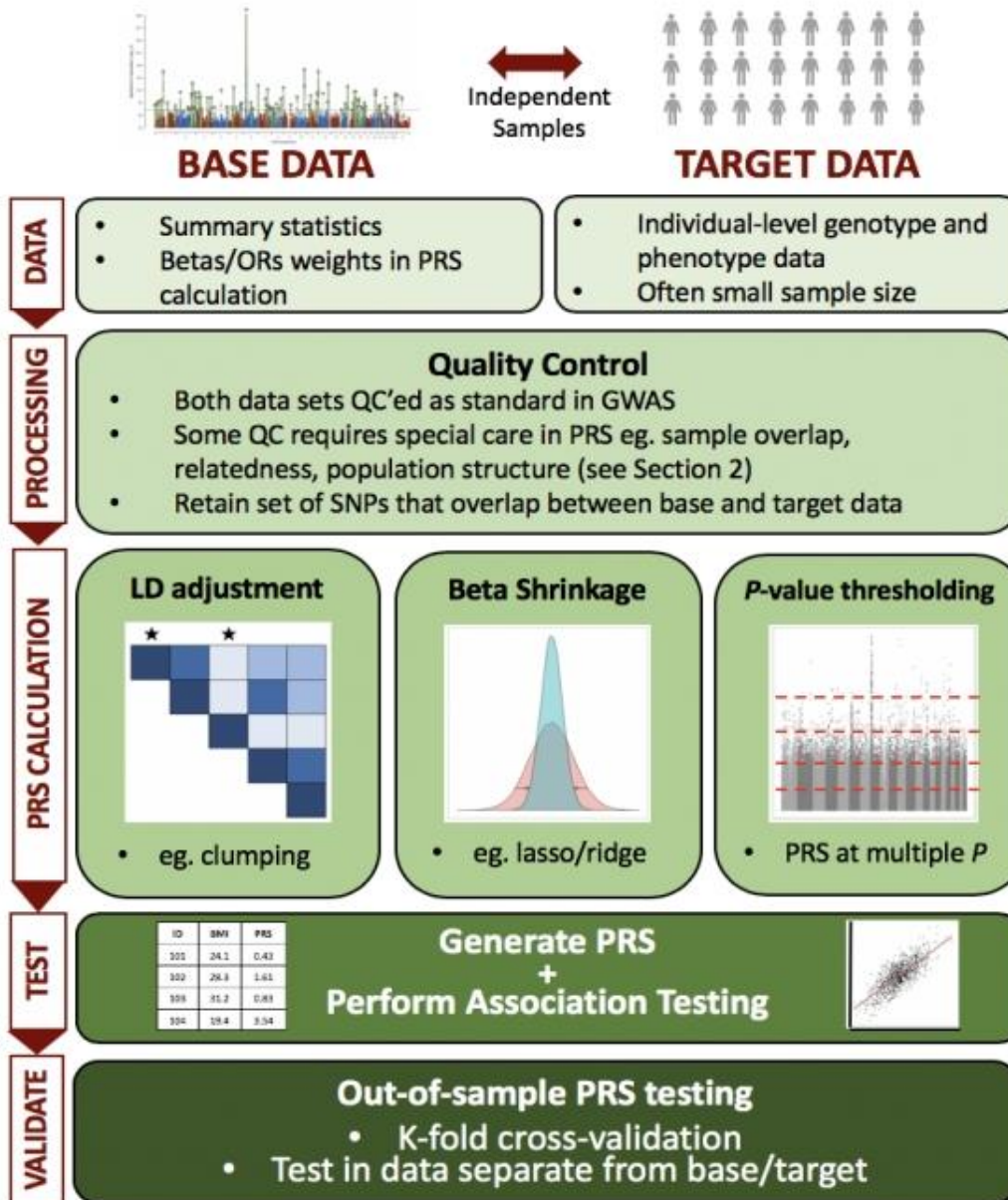
위험도 예측모델 Polygenic risk model

$$\hat{\phi}_j = \sum_{i=1}^M \hat{\beta}_i x_{ij}$$

$\hat{\phi}_j$: PRS
 $\hat{\beta}_i$: Estimated per-allele log odds
 x_{ij} : Genotype (# of variant alleles)



질병 위험도 예측



UK Biobank Project



500,000 whole human genomes

Oversight:



Funding:



wellcome trust



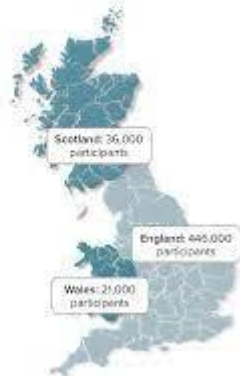
Participants



NHS Genomic Medicine Centres

- Clinical samples and hospital data
- Laboratory processing including molecular pathology
- Broad consent for research and re-contact

UK Biobank's 500,000+ participants



Biorepository



Data

Genomics

Clinical Data

- Identifiable clinical data
- Longitudinal
- Linked to genomic data

Existing Clinical Data

Cancer & RD registries, HES, Mortality data, etc

Sequencing



Genomics

Research Data

- Pseudonymised
- GeCIP and industry partners work within data centre

Data and Analysis Improvement

- Annotation & QC
- Scientists/SMEs
- Product comparison

Fire wall

Clinicians & Academics

Training
NHS Health Education England

Industry

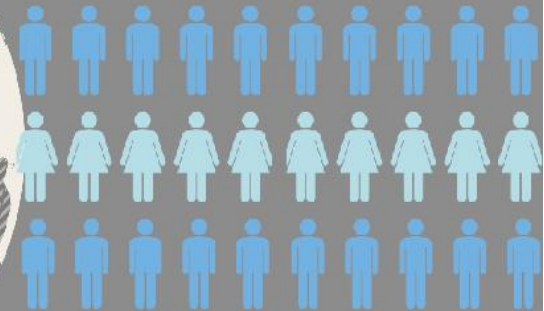
UK Biobank Project

Our data can enable your vision to improve the world's health

OVER 19,000 GLOBAL REGISTRATIONS



1,465 SCIENTIFIC
PUBLISHED PAPERS



UK 22%

INTERNATIONAL

78%



Trans-biobank analysis with 676,000 individuals

The association of polygenic risk scores of complex traits with human lifespan

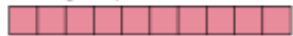


BioBank Japan

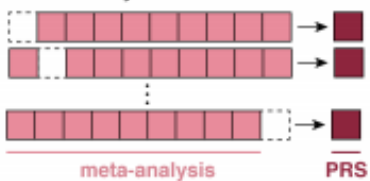


$n = 179,066$

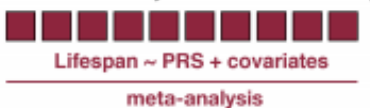
Sub-group GWASs



Meta-analysis and PRS derivation



Survival analysis and meta-analysis



$\text{Lifespan} \sim \text{PRS} + \text{covariates}$

meta-analysis

UK Biobank



$n = 361,194$

When individual-level data available;

Sub-group GWASs



Meta-analysis and PRS derivation



Survival analysis and meta-analysis



$\text{Lifespan} \sim \text{PRS} + \text{covariates}$

meta-analysis

Otherwise;

PRS from public GWAS and survival analysis



$\text{Lifespan} \sim \text{PRS} + \text{covariates}$

Parental lifespan $\sim \text{PRS} + \text{covariates}$

FinnGen



$n = 135,638$

When UKBB sumstats available;

UKBB whole GWAS



UKBB-based PRS construction



$\text{Lifespan} \sim \text{PRS} + \text{covariates}$

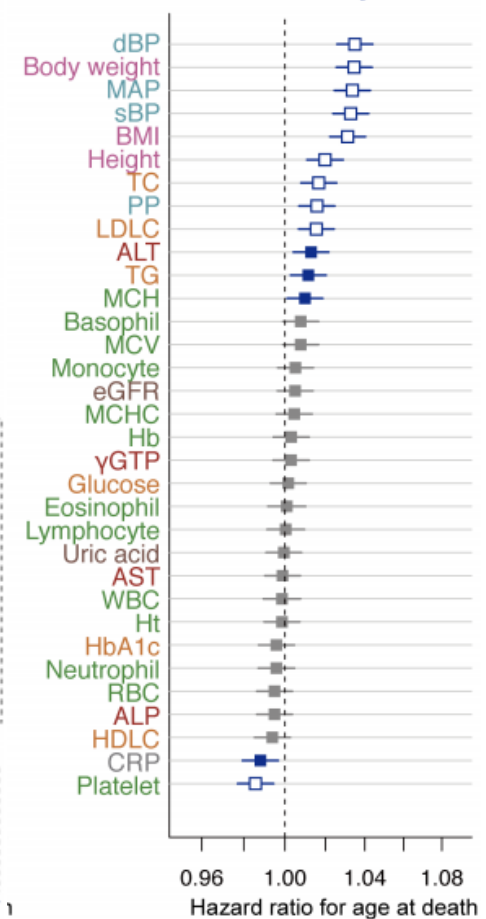
Otherwise;

PRS from public GWAS and survival analysis



$\text{Lifespan} \sim \text{PRS} + \text{covariates}$

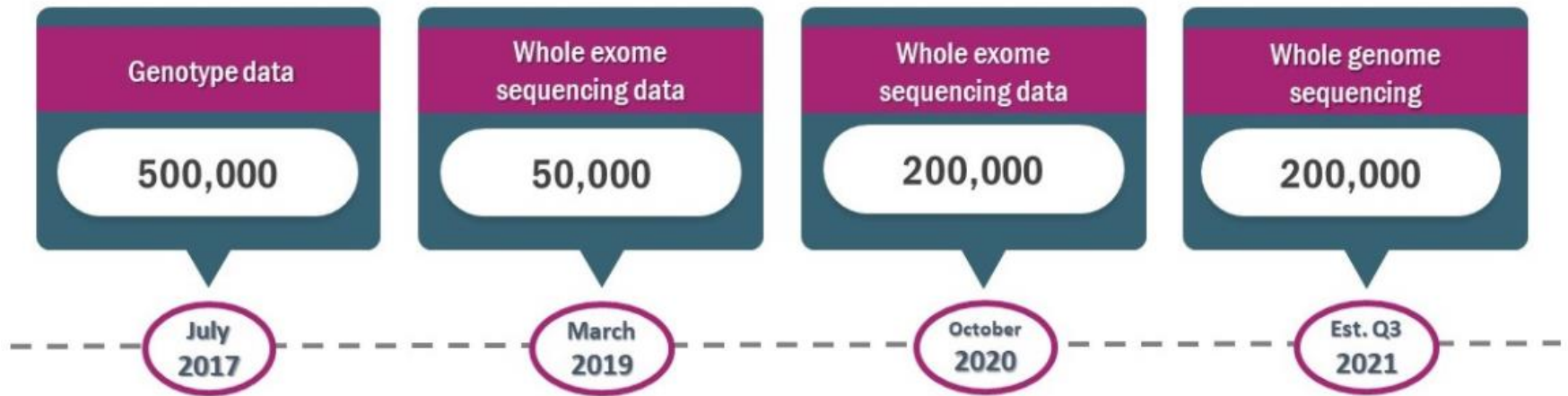
d. Meta-analysis



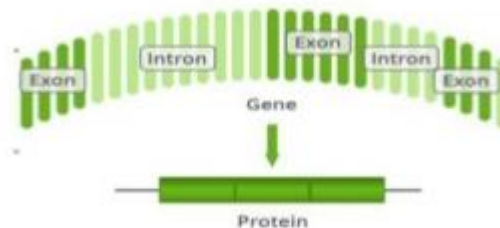
Trans-ethnic meta-analysis

(Nature Medicine, 2020)

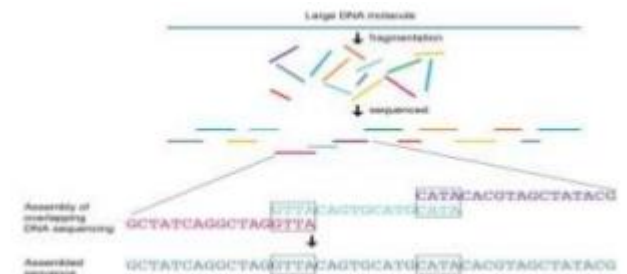
Genetic data release timeline



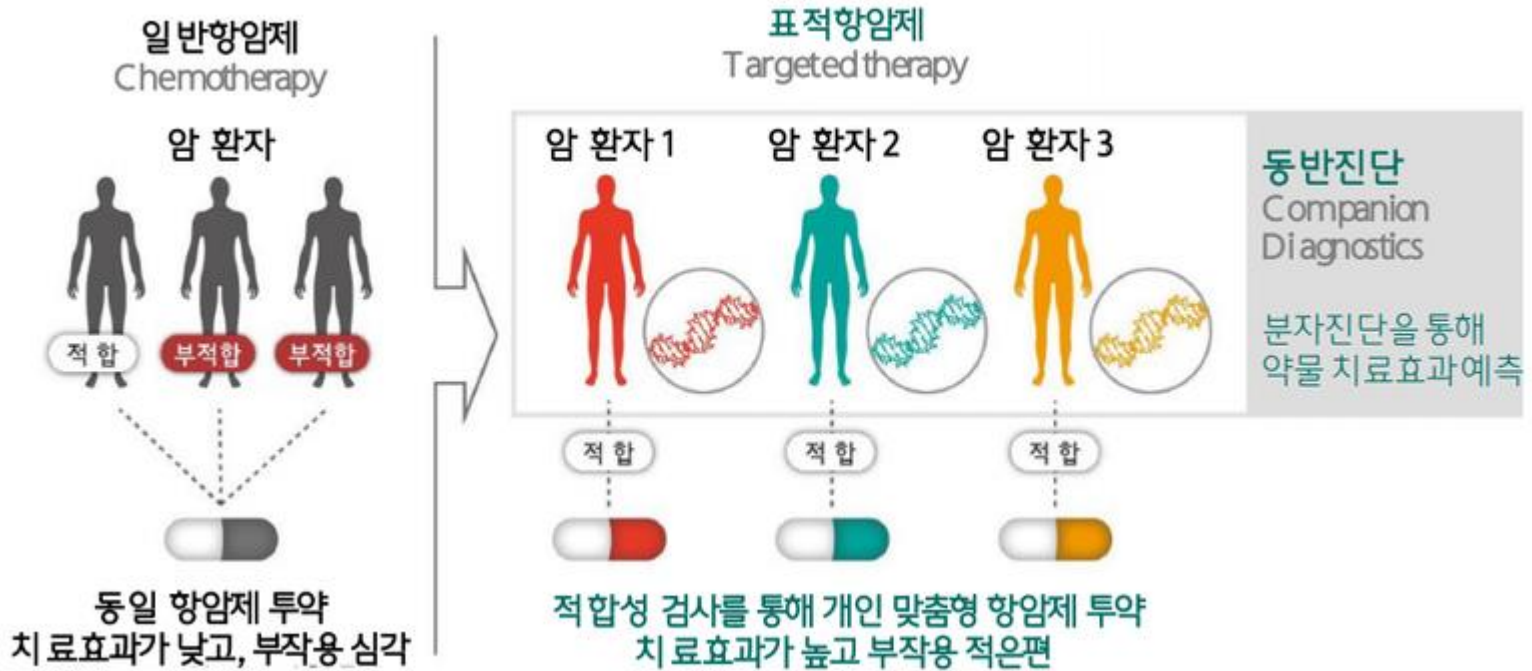
850,000 positions



~2% of genome



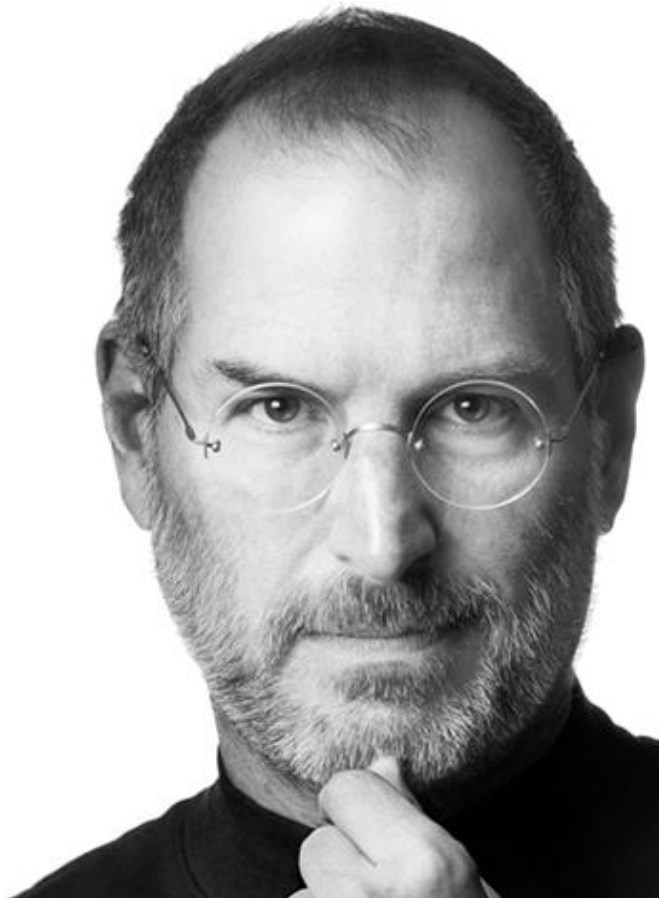
3 billion in genome



- 유전정보, 임상정보, 환경정보
- 효과적인 치료방법(표적항암제 등)을 선택
- 각 개인에 최적화된 진단 및 치료를 적용

출처: 메디파나뉴스

Steve Jobs
1955-2011



종양 세포들에서 유전자 변이들

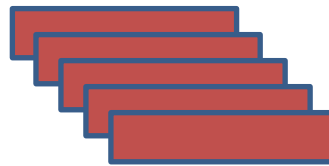


1 포인트 변이



비정상 시그널

2 복제수 변이



증가/삭제 시그널

3 유전자 융합



비정상 시그널

유전자 패널 검사

요양기관명	유전성	고형암	혈액암
서울대병원	○	○	○
인천성모병원	○	○	○
부산대병원	○	○	○
분당서울대병원	○	○	○
아주대병원	○	○	○
서울의과학연구소 용인의원	○	○	○
서울아산병원	○	○	○
이원의료재단 이원의원	○	○	○
서울성모병원	○	○	○
연세대 세브란스병원	○	○	○
랩지노믹스진단검사의학과의원	○	○	○
고려대 안암병원	○	○	
고려대 구로병원	○	○	
삼성서울병원	○	○	
연세대 강남세브란스병원	○	○	

인하대병원
길병원
순천향대 서울병원
녹십자의료재단 녹십자병원
국립암센터
삼광의료재단삼광의료원
분당차병원

건강보험 적용 NGS 검사 대상 (단위:개)

고형암(10종)

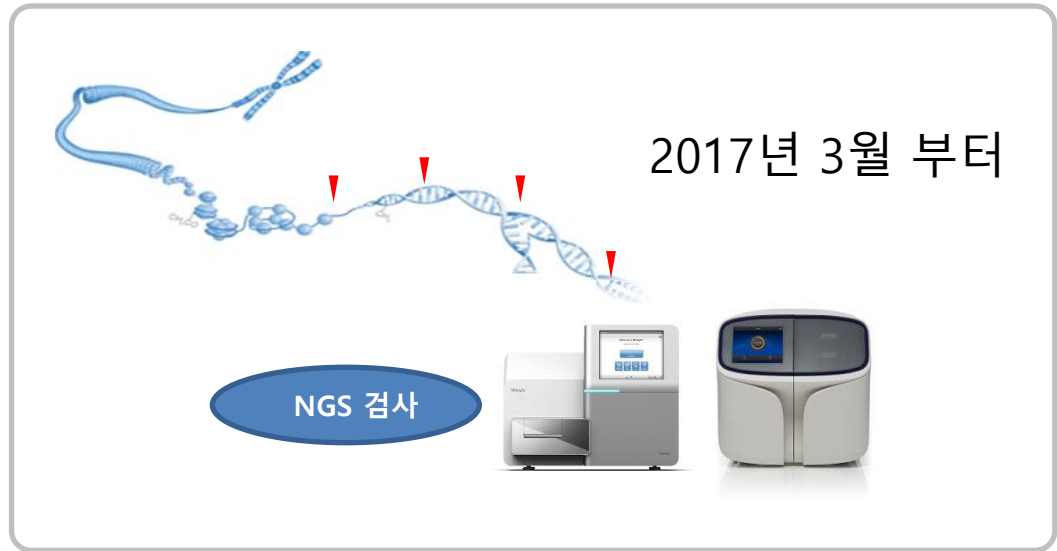
위암, 폐암, 대장암, 유방암, 난소암, 흑색종, 악성뇌종양, 위장관 기저종양, 소아신경모세포종, 원발성 불명암(총 14)

혈액암(5개군)

급성 골수성 백혈병(9), 급성 림프구성 백혈병(5), 악성 림프종(3), 형질세포종(3), 골수형성이상골수증식증양(11)

유전질환(4개군)

유전성 망막색소변성(7), 유전성 난청(4), 샤르코마리투스병(4), 기타(0)



기존 유전자 검사

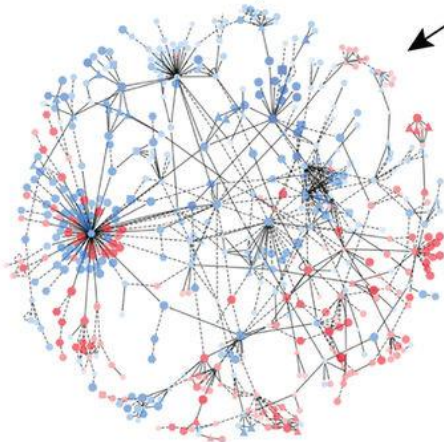


12 tumor types

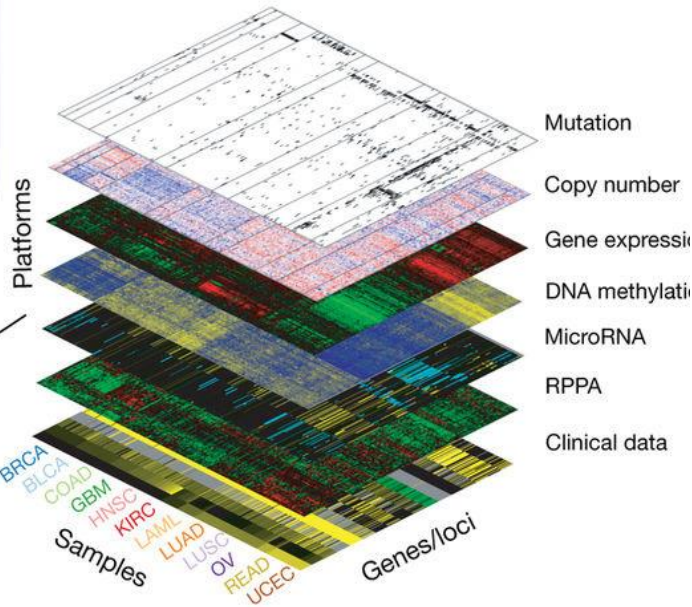
- Leukemia (LAML)
- Lung adenocarcinoma (LUAD)
- Lung squamous (LUSC)
- Kidney (KIRC)
- Bladder (BLCA)
- Endometrial (UCEC)
- Glioblastoma (GBM)
- Head and neck (HNSC)
- Breast (BRCA)
- Ovarian (OV)
- Colon (COAD)
- Rectum (READ)



Thematic pathways



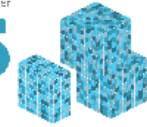
Omics characterizations



NATIONAL CANCER INSTITUTE THE CANCER GENOME ATLAS

TCGA BY THE NUMBERS

TCGA produced over **2.5** PETA-BYTES of data



TCGA data describes **33** DIFFERENT TUMOR TYPES including **10** RARE CANCERS

Based on paired tumor and normal tissue sets collected here **11,000** PATIENTS

To put this into perspective, 1 petabyte of data is equal to

212,000 DVDs

TCGA RESULTS & FINDINGS

- MOLECULAR BASIS OF CANCER**: Improved our understanding of the genomic underpinnings of cancer
- TUMOR SUBTYPES**: Revolutionized how cancer is classified
- THERAPEUTIC TARGETS**: Identified genetic characteristics of tumors that can be targeted with currently available therapies or used to aid with drug development

For example, a TCGA study found the basal-like subtype of breast cancer to be similar to the more aggressive subtype of ovarian cancer on a molecular level, suggesting that despite arising from different tissues in the body, these subtypes may share a common path of development and respond to similar therapeutic strategies

TCGA revolutionized how cancers are classified by identifying tumor subtypes with distinct sets of genomic alterations*

TCGA's identification of targetable genomic alterations in lung squamous cell carcinoma led to NCI's Lung-MAP Trial, which will treat patients based on the specific genomic changes in their tumor.

THE TEAM

20 COLLABORATING INSTITUTIONS across the United States and Canada

WHAT'S NEXT?

The Genomic Data Commons (GDC) focuses TCGA and other NCI generated data sets for scientists to access from anywhere. The GDC also has many expanded capabilities that will allow researchers to answer more clinically relevant questions with increased ease

* The original classification of cancer with a molecular basis, but a lack of cancer type defining, molecularly defining, or defining the molecular basis of cancer.

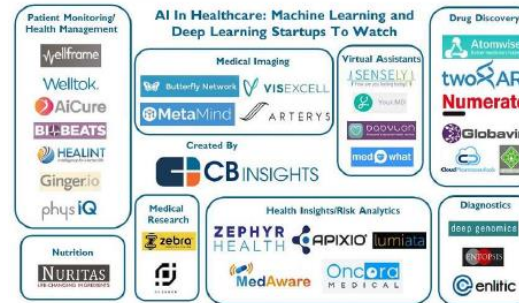
국가 바이오 빅데이터 구축 사업

시범사업 기간 동안 총 25,000명 ('20~'21)

구분	주요 항목	규모
희귀질환	환자 및 부모의 임상 및 유전체데이터 희귀질환 분류, (추정)진단명, HPO 등	15,000명
대장암	환자의 임상 및 유전체 데이터 대장암 진단, 생존 및 치료 정보	400명
자폐증	환자 및 부모의 임상 및 유전체 데이터 행동, 지능지수, 적응능력, 주의력, 인지전환능력, 이상행동, 강박증상, 질병력, 약물투여력 등	500명
일반인	역학 및 유전체 데이터 질병력, 가족력, 생활습관, 신체 및 혈액 검사 및 소변검사 등 국가 일반검진 정보	울산1만명 프로젝트 1,600명 KoGES 5,000명

인공지능헬스케어의 급성장

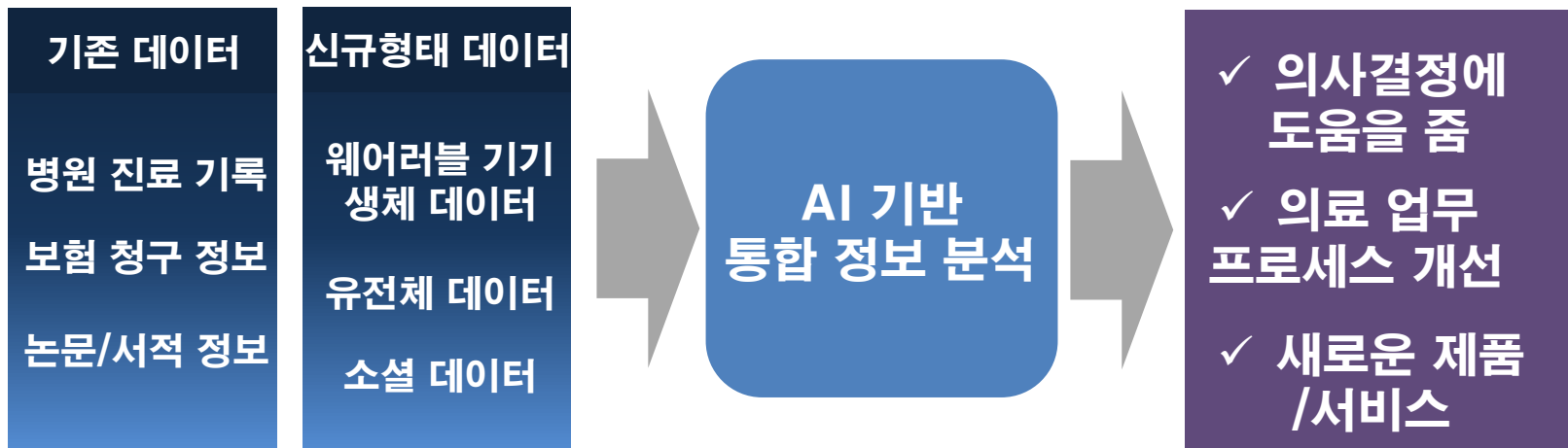
The next major advance in medicine will be the use of AI



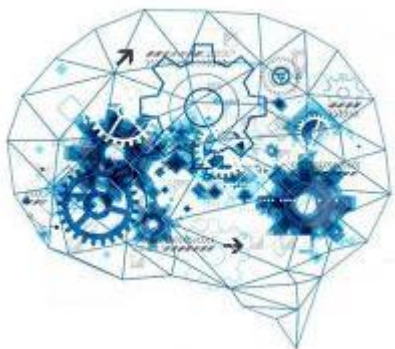
자료: CB insight

90이상의 헬스케어+AI 스타트업에
\$15억(1조6천억투자)

의료 빅데이터



유전체 분야 AI 기술 적용 패러다임



Deep learning
Cognitive computing

...

classifying
prioritizing
interpreting
linking

**AI 알고리즘
적용**

IBM 왓슨
cleveland clinic
resistant to treatments

Geisinger
MyCode clinical genomics

Pathway Genomics
smartphon app of
personal genetic info

Deep Genomics
genetic variant analysis

Mendel.ai
clinical trial matching platform

Sophia Genomics
genome interpretation

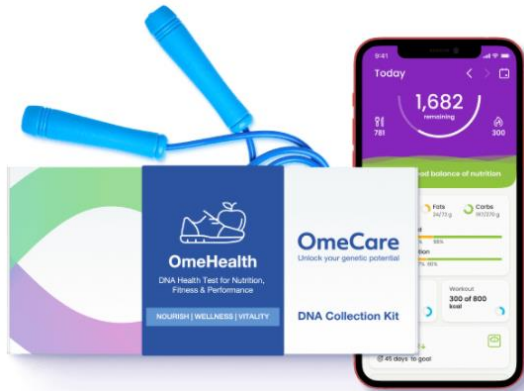
iCarbonX
genomics & health factors

**글로벌 기관
/기업**

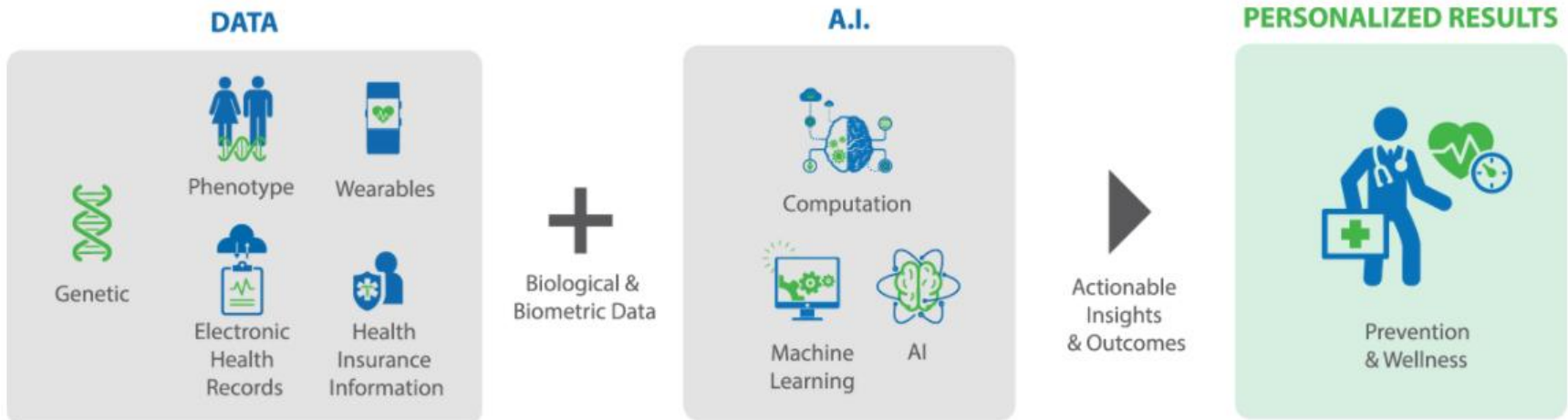


- ✓ **Pathogenic variants**
- ✓ **Drug development**
- ✓ **Drug/treatment response**
- ✓ **Cancer risk**
- ✓ **Clinical trials**

**의료 분야
적용**



- 개인유전체 정보를 이용한 AI 및 딥러닝
- 암, 유전질환, 심혈관질환, 약물반응성 등 예측
- 건강 히스토리를 근거로 개인 건강 지식 제공



The Virtual Medical Assistant

