

Summary

: GWAS, Phylogenetic analysis

Sung Wook Chi

Division of Life Sciences, Korea University

Overview of genome-wide association study (GWAS)

all type variation
phenotype per associated SNP

Sample design / Collection



WGS



GWAS



Associated Loci
QTL



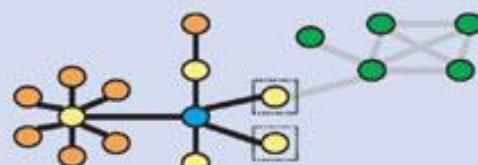
Variation & Gene



Causative variation ?



Functional interpretation validation



Population resources – trios or case-control samples

Whole-genome genotyping

Genome-wide association
with statistical test

Fine mapping

Gene mining

Gene sequencing & polymorphism identification

Identification of causative SNPs

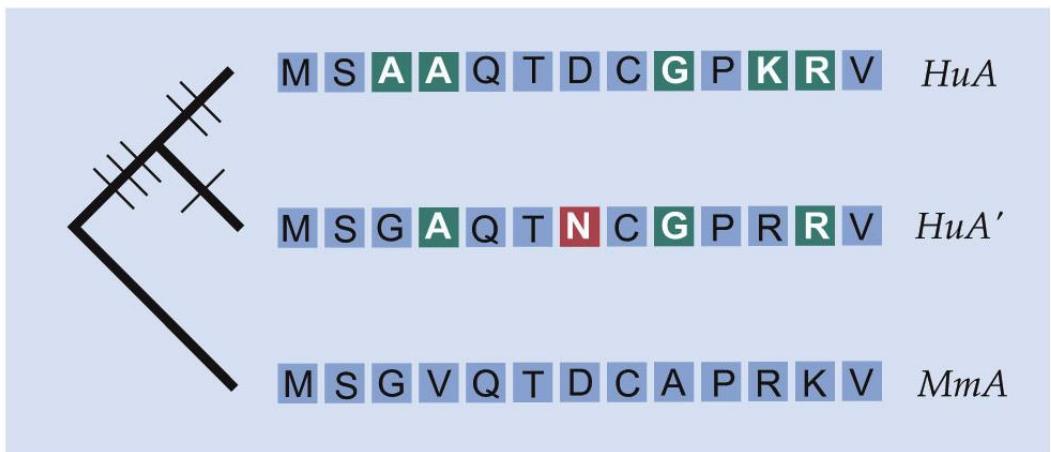
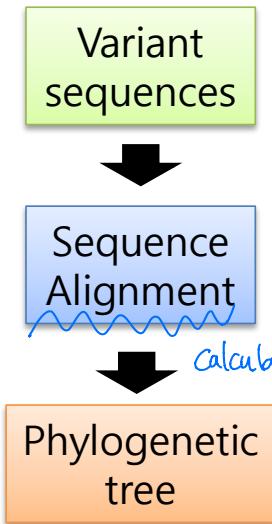
Pathway analysis & target identification

GWAS:

An examination of genetic variation across a given genome whether it associated with phenotypes (quantitative traits, diseases)

Interpretation of phylogenetic tree: Orthologs and paralogs

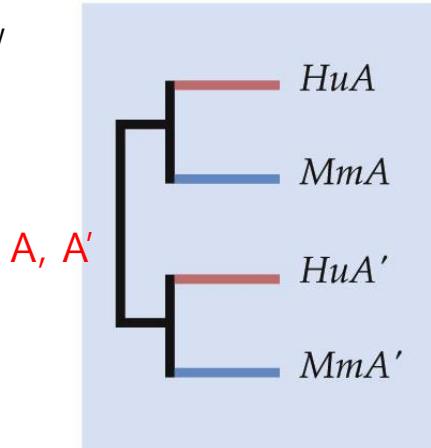
How these diversity were generated through evolution event.
(A)



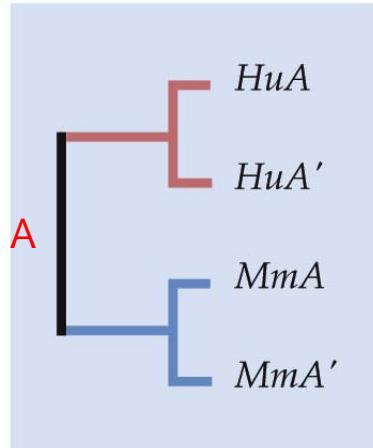
- Recalculate / Infer
- Distance or similarity b/w sequences
- Topology (order)
- Length (evolutional time)

Evolutional relationship / interpretation

(B)



(C)



: Whole genome sequencing (NGS)

1. Whole genome sequencing (WGS by NGS)
2. Variation Analysis for NGS data
3. Genome-wide association study (GWAS)
4. Phylogenetic analysis

Functional Genomics

기능유전학
기능 유전학

: Exome-Seq and functional genomics

Sung Wook Chi
Division of Life Sciences, Korea University

Functional Genomics

-**Functional Genomics** : Studying functions of genes using massive data

1) genome-wide methods

2) Gene expression (transcription, translation, protein-protein interaction)

*_{the} gene : *functional genetics*

Gene expression

Genome

Phenotype

→
Genome-wide
methods

Variation

→
GWAS

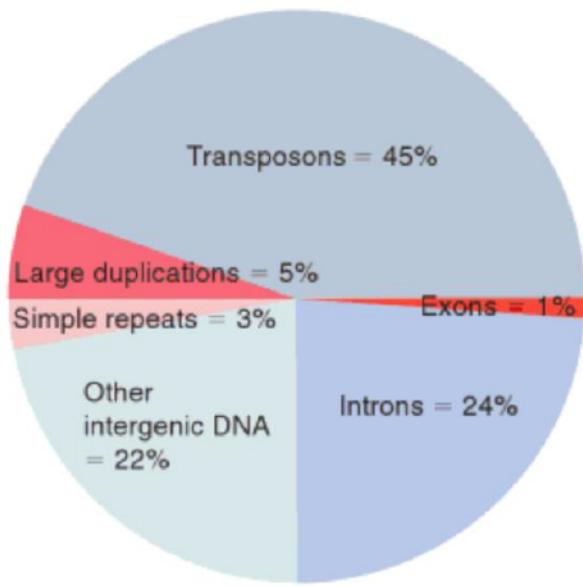
Diseases
Quantitative trait

NGS
- WGS
- **Exome-Seq**

Whole Exome Sequencing (WES)

Whole Exome Sequencing, Why?

- Focuses on the part of the genome we understand best, the **exons** of genes
- Exomes are ideal to help us understand its **relationship to phenotype**.

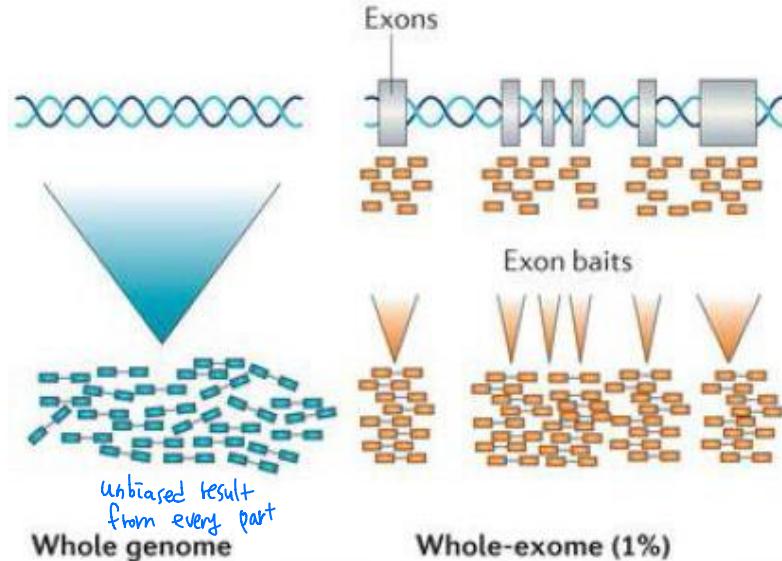


only looking at exon part (~1%)

- Capturing and sequencing the **~1% of the human genome** that codes for **protein sequences** *treatment* *large* *repetitive* *regions* *Exons* *in* *1%*
- On average, whole-exome sequencing identifies **12,000 variants** in coding regions.
- **~90%** are found in publicly available databases.
useful to identify new variation

- A whole exome is 1/6 the cost of whole genome and 1/15 the amount of data

Whole genome vs. Whole exome sequencing



Predominant applications:

- Structural variants
- Point mutations
- Copy number variation

Predominant applications:

- Point mutations
- Copy number variation → duplication

hard to find structure variants

Target Size: 3Gb

30Mb

Scale: 1

0.01 → can have more coverage

File size for human germline WGS (30X)

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

Exome Sequencing Pipeline

Sample DNA Fragmentation

Illumina Library Preparation

Exome Enrichment

⇒ WGS에 있는 단계

Cluster Generation

NGS

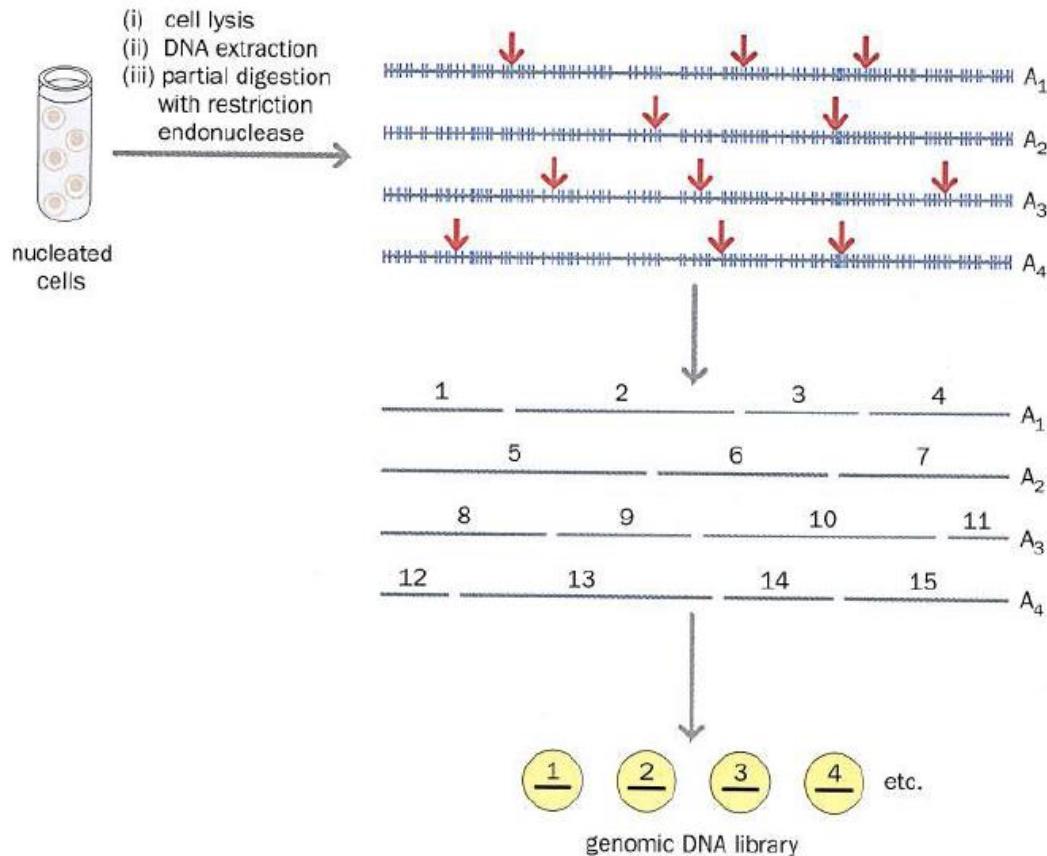


- Analysis of raw data (images)
- Base Calling: Determination of sequences
- Mapping: Alignment to reference genome variation detection
- Variant calling:
 - Substitutions
 - Indels

~30 000 to 50 000 variants /person

Filtering steps (artefacts, known variants,...)

Genomic DNA fragmentation (DNA libraries)

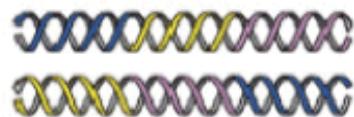


Exome capture, enrichment

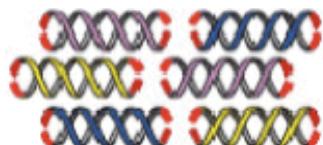
Exome-capture (Agilent)



GENOMIC SAMPLE
(Set of chromosomes)



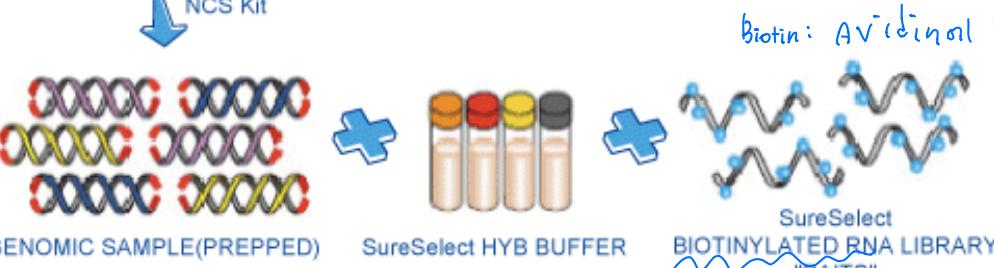
NCS Kit



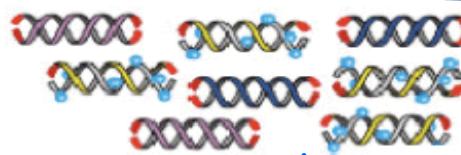
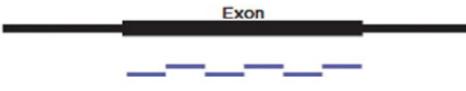
GENOMIC SAMPLE(PREPARED)

only exon part will hybridize
size phase off by 1/4 for purity

SureSelect Target Enrichment System Capture Workflow



Design of capture phases



hybridization only for Exon

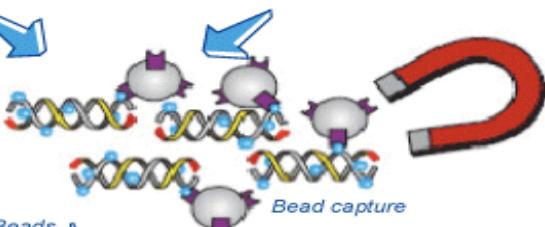


UNBOUND FRACTION
DISCARDED

-Genome-wide capture of all human exome is possible (180,000 human exons)



+ →



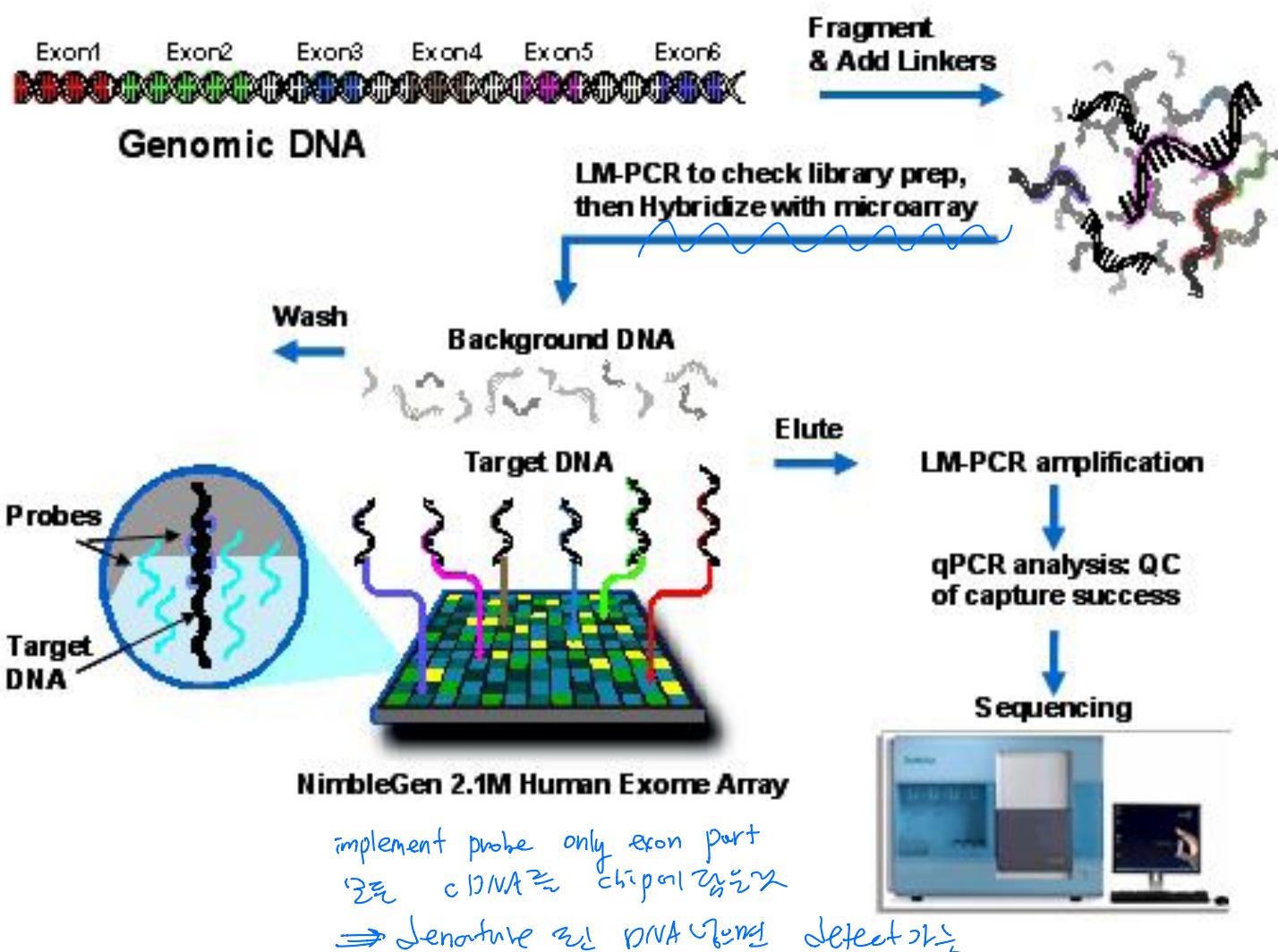
Wash Beads
and
Digest RNA



Amplify

Sequencing

Exome-capture (Nimblegen)



Exome Sequencing



Access

To read this story in full you will need to login or make a payment

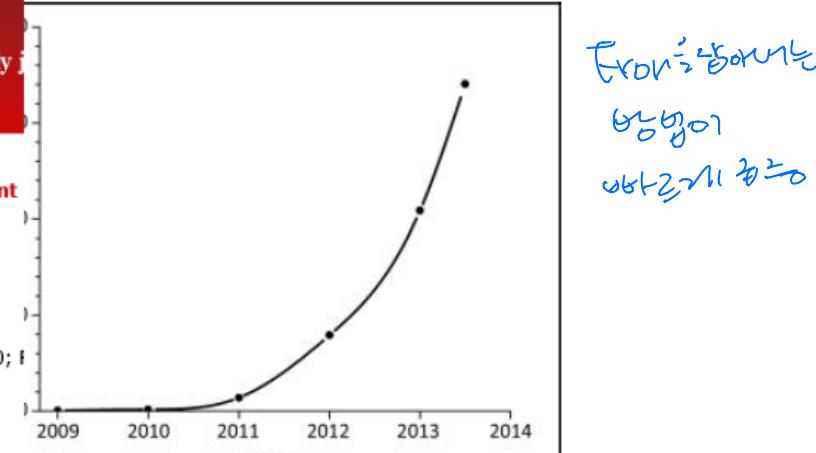
nature.com > Journal home > Table of Contents

Letter

Nature 461, 272-276 (10 September 2009) | doi:10.1038/nature08250; published online 16 August 2009

Targeted capture and massively parallel sequencing of 12 human exomes

Sarah B. Ng¹, Emily H. Turner¹, Peggy D. Robertson¹, Steven D. Abigail W. Bigham², Choli Lee¹, Tristan Shaffer¹, Michelle Wong¹, Bhattacharjee⁴, Evan E. Eichler^{1,3}, Michael Bamshad², Deborah A. Nickerson¹ & Jay Shendure¹



Advantages

- Higher sequence coverage and less raw sequence and cost than WGS.
- Point mutations
- Small indels
- CNAs (need further validation)

Small indels

Avoid noise from contamination

Disadvantages

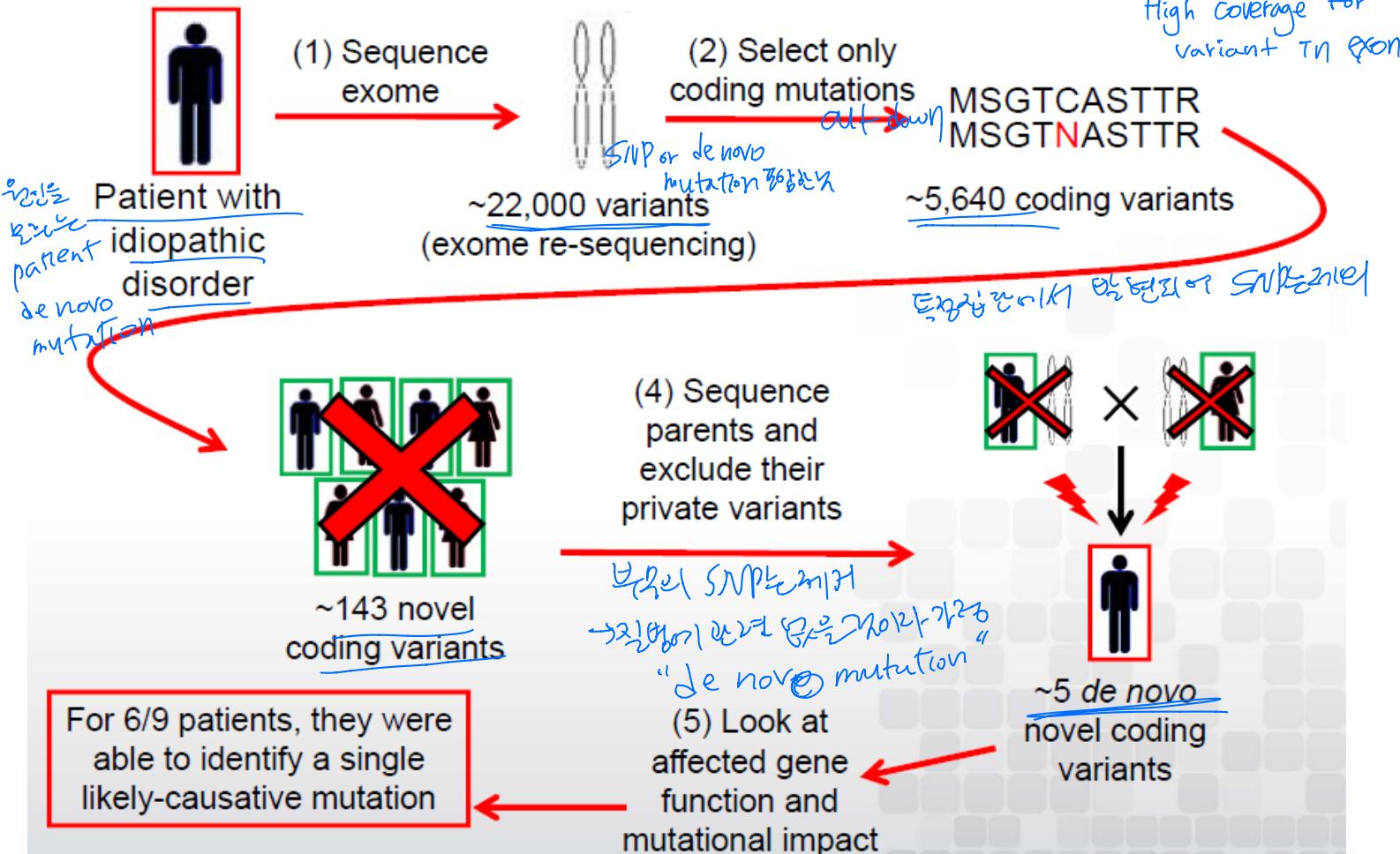
- Incomplete exome capturing
- Limit to mutations in coding regions
- Lack efficient methodology to detect structural variations

proper seq compositional effect
affinity ↑ → 27%
(GC의 양 등에 영향을 미친다)
→ 27% 27% 27% 27%
① 27% 30%
A

Identifying causative de novo mutation

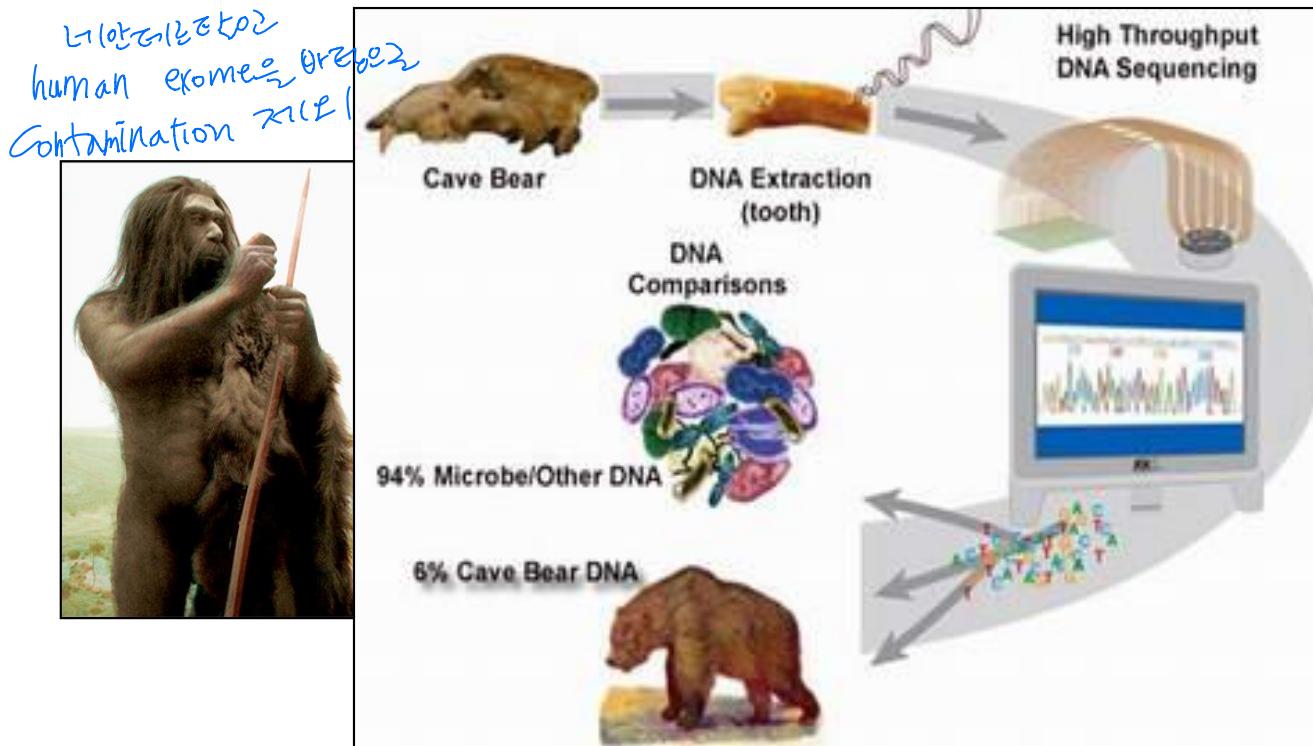
GWA서 유사한 질병

Veltman and colleagues - Nat Genet. 2010 Dec; 42(12):1109-12



Ancient Genomes Resurrected

- Nuclear genomes of ancient remains: cave bear, mammoth,
 - Neanderthal genome (10^6 bp)

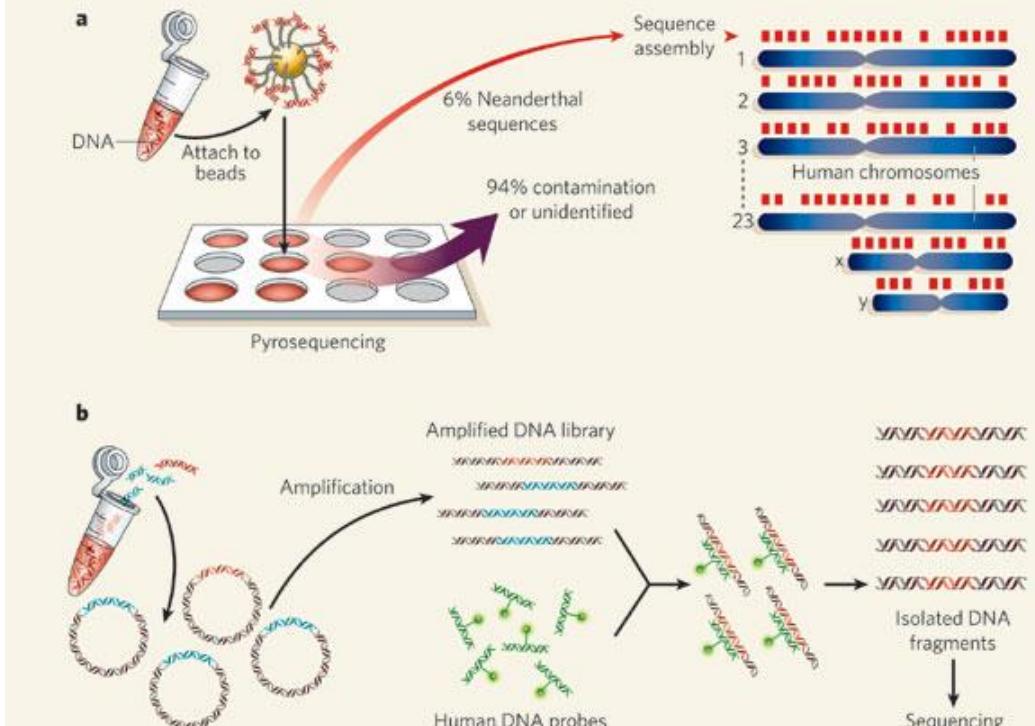


Problems: contamination modern humans and coisolation bacterial DNA
Solution : Exome Capture based on Homo Sapiens

The Neanderthal Genome

- Neandertal genome

- Neandertal genome composed of more than 4 billion nucleotides from three individuals
 - Gene flow has occurred from Neandertals to humans of Eurasian descent, but not to Africans



- sequenced ~14,000 protein-coding positions
 - identified 88 amino acid substitutions that have become fixed in humans
since our divergence from the Neandertals

<https://www.youtube.com/watch?v=zHhaVjzip-o>

Functional Genomics

네안데르탈인은 호모 사피엔스에서는 흔히 발견되지 않으나, 현생 아시아, 유럽인에서 발견되는 것들이 많이 발견
-> 아프리카에서 시작한 호모사피엔스와 섞였을 것으로 추측 (?)

Gene expression

Genome



Phenotype

1. Human genome projects
2. Genome Sequencing
3. Sequence Alignment
4. Genome annotation
4. Genomic variation

NGS

Genetic interaction
RNAi screening

DNA
Chip-Seq / ENCODE

RNA
RNA-Seq / microarray

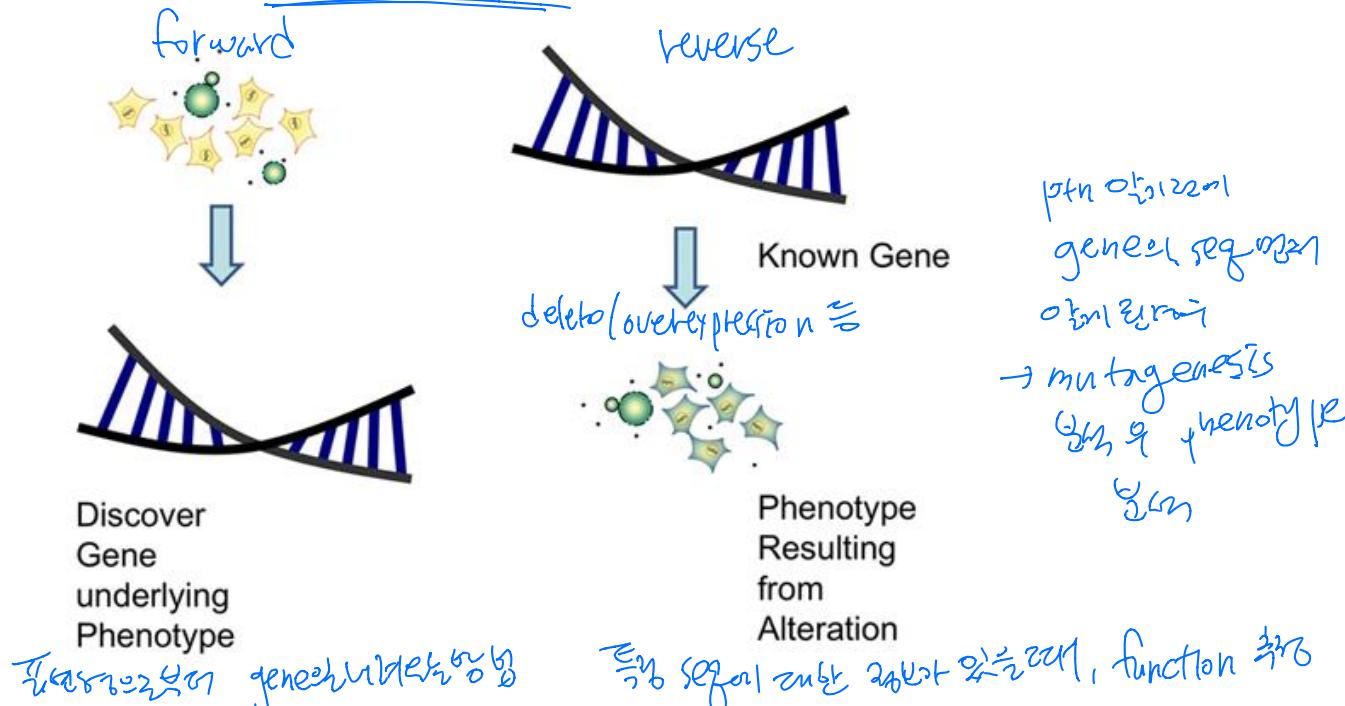
Protein
Translation

1. Biological function
2. Diseases
3. Quantitative trait /population

Forward and Reverse Genetics

Clues to gene function can be inferred from different types of genetic manipulation

- Two major types of genetic analysis have been applied in model systems
 - **Forward genetics** (phenotype → gene) : random mutagenesis
 - **Reverse genetics** (gene → phenotype) : targeted mutagenesis



Selective gene inactivation and modification

- Four different functional classes of mutation can be generated, and different screens can be conducted depending on the desired effect on the target gene

TABLE 12.3 DIFFERENT TYPES OF REVERSE GENETIC SCREEN IN CULTURED MAMMALIAN CELLS

Type of screen	Basis of method
Loss-of-function	usually, the RNA interference (RNAi) pathway is induced to selectively degrade RNA transcripts of a specific target gene (see Figures 12.3 and 12.4)
Gain-of-function	involves transfection of an exogenous gene copy driven by a suitable promoter to overexpress or misexpress that gene in a desired cell type
Dominant-negative	relies on producing a mutant gene product that interferes with the normal product of a specific target gene; similar to loss-of-function screen but suppresses gene activity more efficiently than RNA; often the mutant product is designed to be a truncated version of the wild type that competes with wild-type protein in some way, or the mutant protein becomes incorporated with wild-type proteins in multisubunit complexes, thereby inactivating the complex
Modifier → infer relationship	seeks to identify genes that enhance or suppress a specific phenotype; the initial phenotype is produced by one of the three methods above; thereafter, the mutant cells are subjected to a second genetic modification (which can again be any of the three methods above); if altered expression of a second gene enhances or suppresses the initial phenotype, the second gene is considered to be a modifier of the gene that caused the initial phenotype; it then remains to be worked out how the two genes interact

Genetic interaction

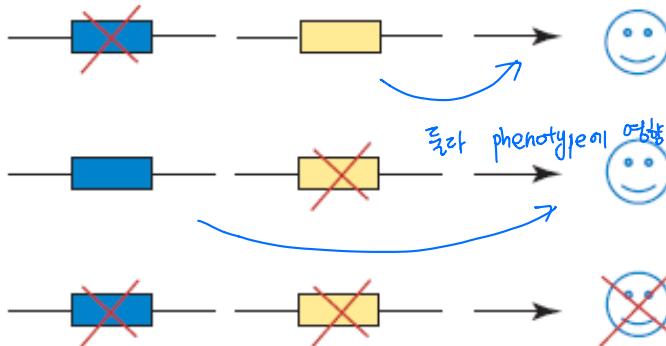
→ introducing loss of function or the control gene

Genetic interaction



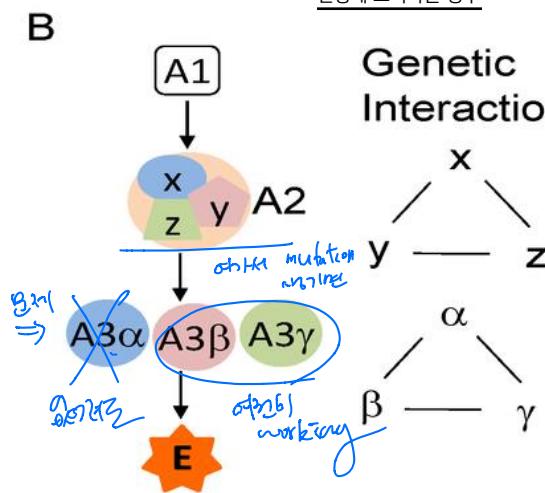
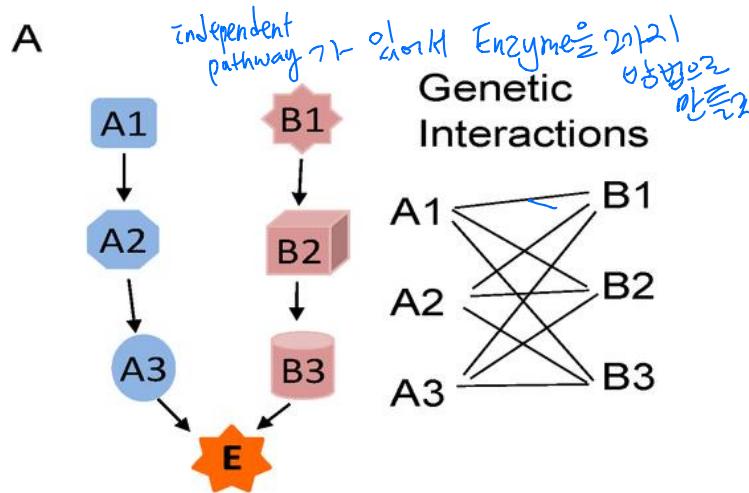
different mutation
same phenotype
or working together

Genetic interaction
(e.g. synthetic sick or lethal interaction)



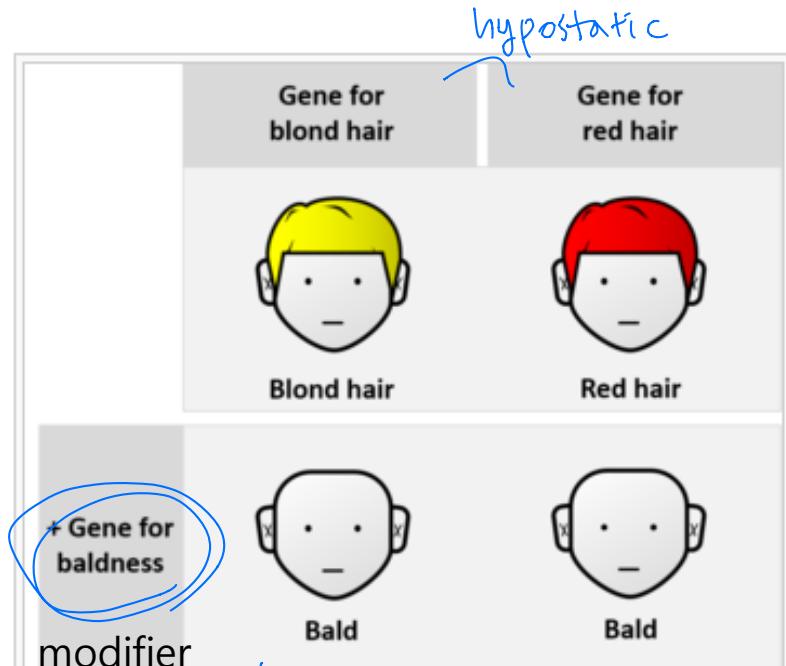
one mutation is not enough
↳ still working

ptn. complex를 만들어 각 complex의 한부
분이 없더라도 나머지가
보완할 수 있는 경우 -> 2개 다 없애야지만 표
현형에 드러나는 경우



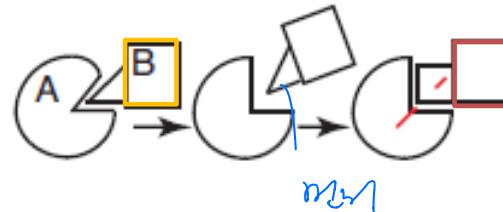
Epistasis

Epistasis is a phenomenon that consists of the effect of one gene being dependent on the presence of one or more 'modifier genes' (genetic background). Similarly, epistatic mutations have different effects in combination than individually

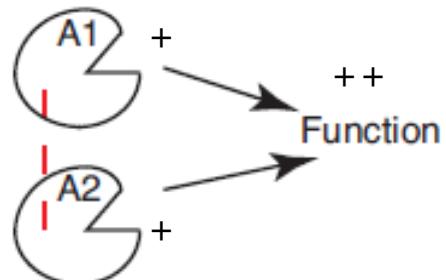


↳ 이 표현을 epistasis라 부른다

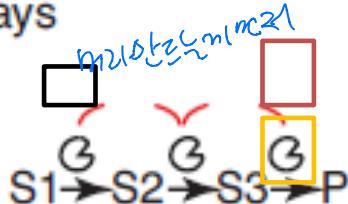
(molecular recognition)
molecular recognition



Redundancy



Positive interactions in linear pathways



Dominant Negative Phenotype

loss of function을 만드는 방법
같은 유전자는 두 가지 유형의 단백질을 만듭니다

homo zygote인 loss-of-ftn을 만드는 것이 어려운데,
dimer을 형성하는 놈은 한 부분만 mutation을 넣어줘도 그 기능을 상실함
-> 이 형태로 loss-of-ftn에서의 표현형을 추측할 수 있음.

- A mutation whose gene product adversely affects the normal, wild-type gene product within the same cell, usually by dimerizing (combining) with it

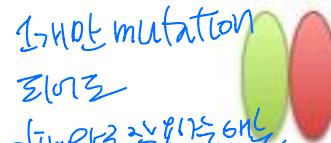
dimer을 만드는 방법

Functional



25%

Non-Functional



loss of function 50%

Non-Functional



25%

Dimer of two normal Factor XI monomers. The dimer has normal function. About one fourth the molecules are such homodimers.

A dimer of one normal (green) and one monomer encoded by the dominant negative gene (red). The dominant-negative gene product inhibits the function of the product of the normal gene rendering the entire molecule functionless. About half the molecules are heterodimers.

A dimer of two dominant negative monomers that is functionless. About one fourth of the molecules are such homodimers.

Selective gene inactivation and modification

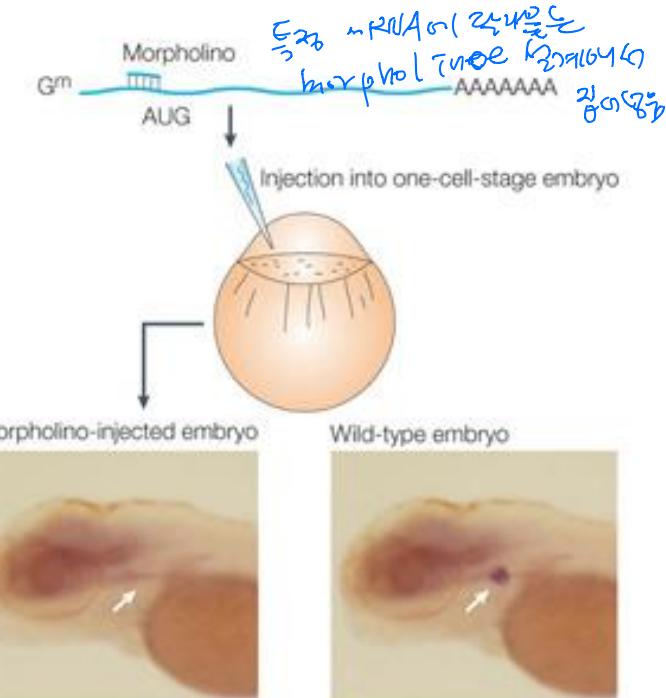
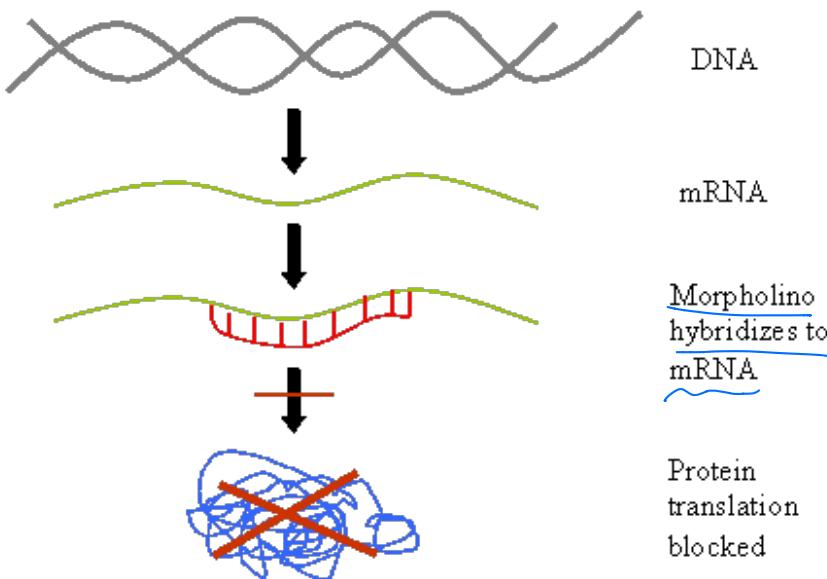
Clues to gene function can be inferred from different types of genetic manipulation

gene 발현 차단을 하는
ex) 아밀로이드 precursor 단백질을 감소.

- Antisense technology: an early general approach to inhibit gene expression using the specificity of base pairing
 - RNA or oligonucleotide constructs are designed to have a sequence that is complementary to that of RNA transcripts from a gene of interest (knockdown)
 - This technology was first developed in the 1980s (antisense RNA)
 - Subsequent way was use of antisense oligonucleotides
 - More recent was use of morpholino antisense oligonucleotides: vastly more stable and robust structure than conventional oligonucleotides, and more consistent in producing significant inhibition of gene expression

Antisense technology : Morpholino

Morpholino ^{consist of oligos ... RNA와 같은} oligos are a class of antisense, a technology used to block access of other molecules to specific sequences within nucleic acids. Morpholinos block small (~25 base) regions of the base-pairing surfaces of RNA and are used as research tools for reverse genetics by knocking down gene function.



RNA interference (RNAi)

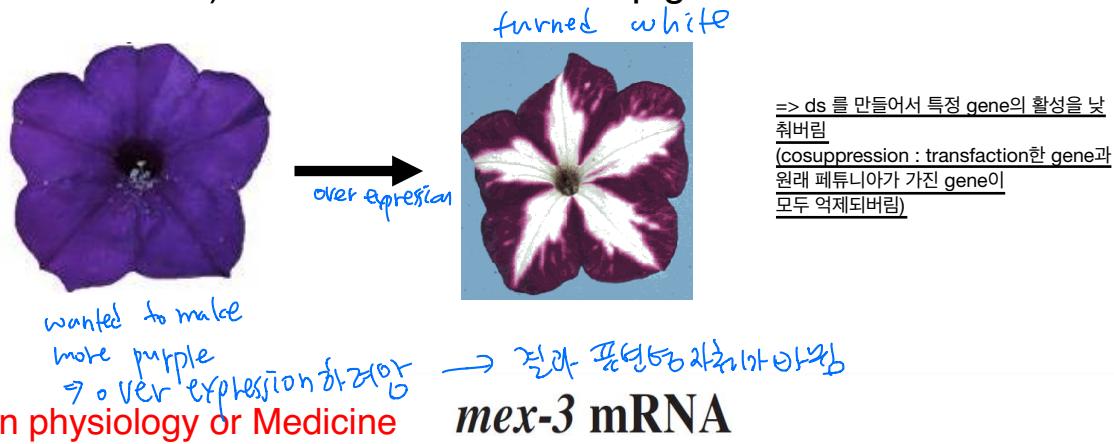
RNAi는 양방향 유전조작의 일환인 것으로 간주된다.

RNA interference (RNAi) is a biological process in which RNA molecules inhibit gene expression, typically by causing the destruction of specific mRNA molecules

double strand mRNA

Phenomena first observed in petunia (1990)

Attempted to overexpress chalone synthase (anthocyanin pigment gene) in petunia (trying to darken flower color) -> Caused the loss of pigment.



2006 Nobel prize in physiology or Medicine

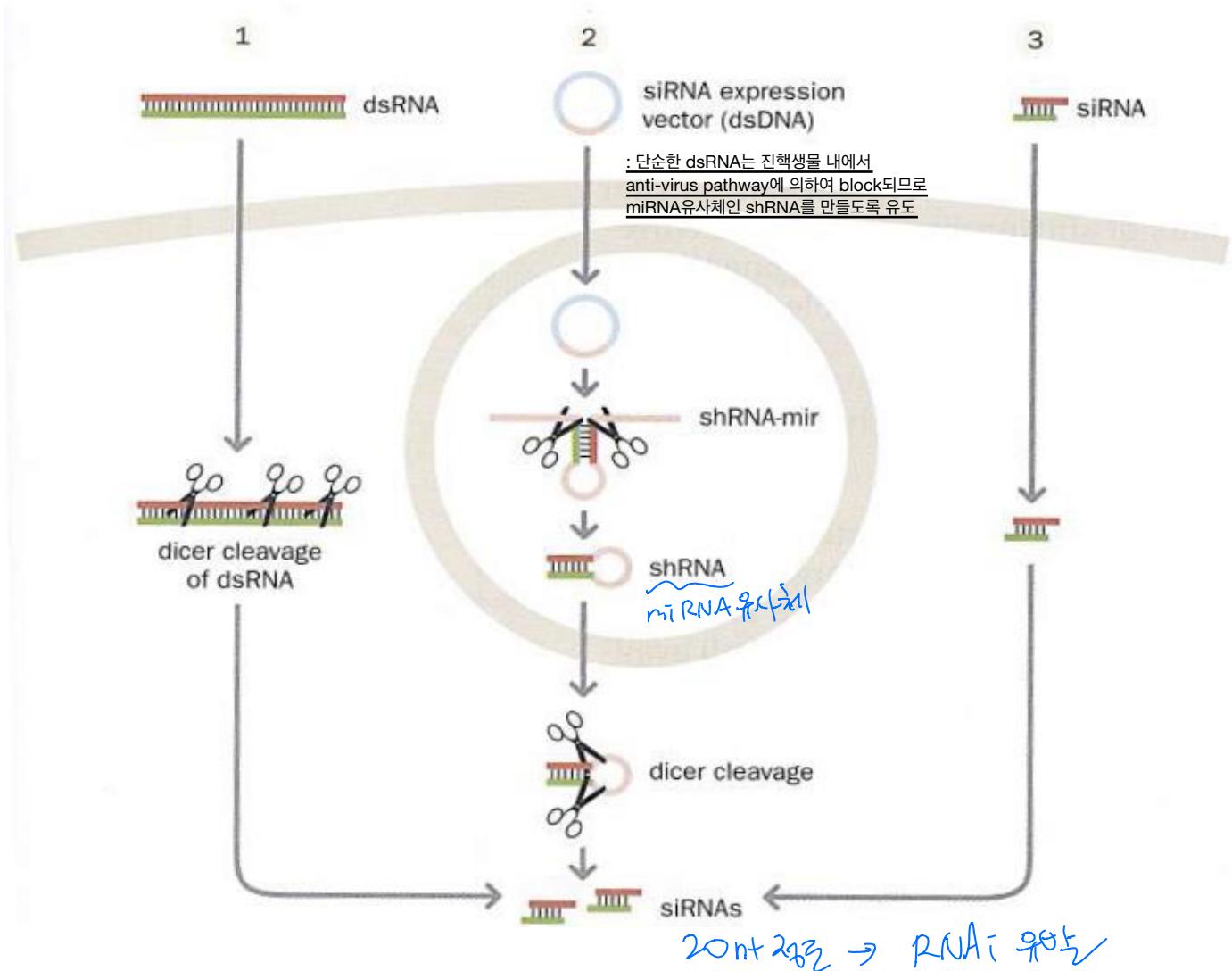
Potent and specific genetic interference by double-stranded RNA in *Caenorhabditis elegans*

anti-sense RNA
sense RNA
gene expression
(1998)
anti → ds RNA인 이유는?



control +dsRNA

RNA interference (RNAi)



RNA interference (RNAi)

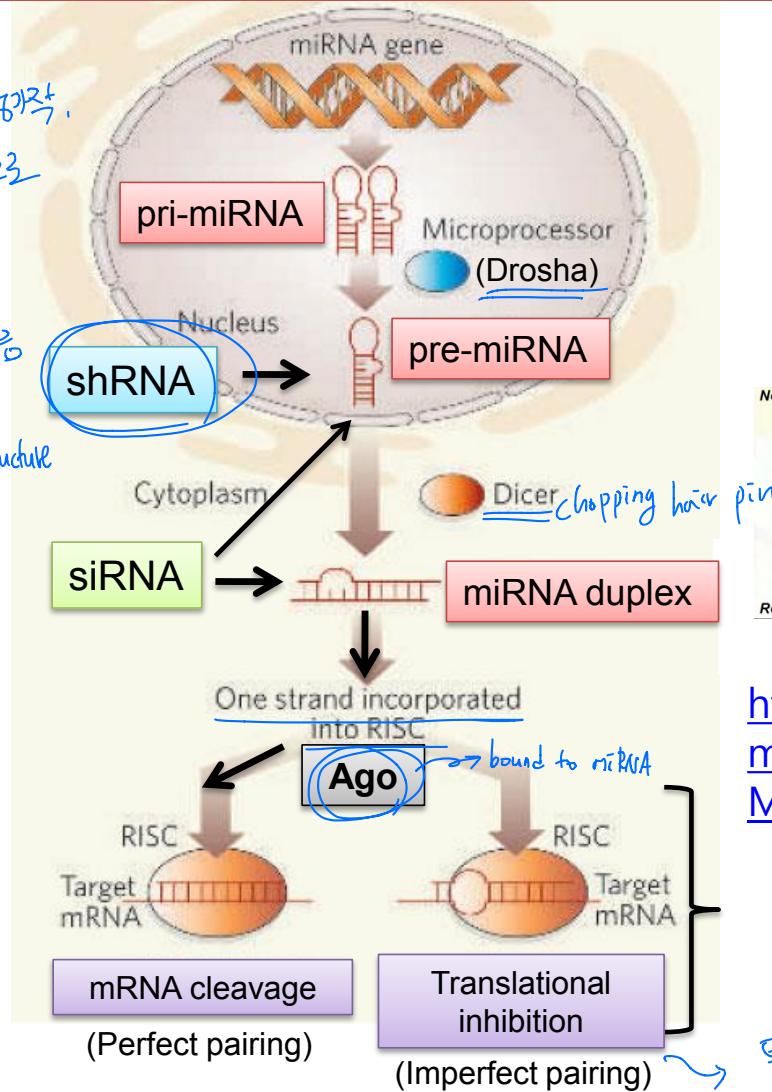


생물학에서 siRNA는?

Q. dsRNA가 mRNA를 block하는가?

⇒ miRNA 풍선과
같은 기작으로 작동

Small hairpin structure



<http://www.youtube.com/watch?v=H5udFjWD-M3E>

특정 mRNA의 번역을 일시정지
 혹은 저하로 통제하는 기작

RNAi screening

design cont siRNA → hybrid to RNA

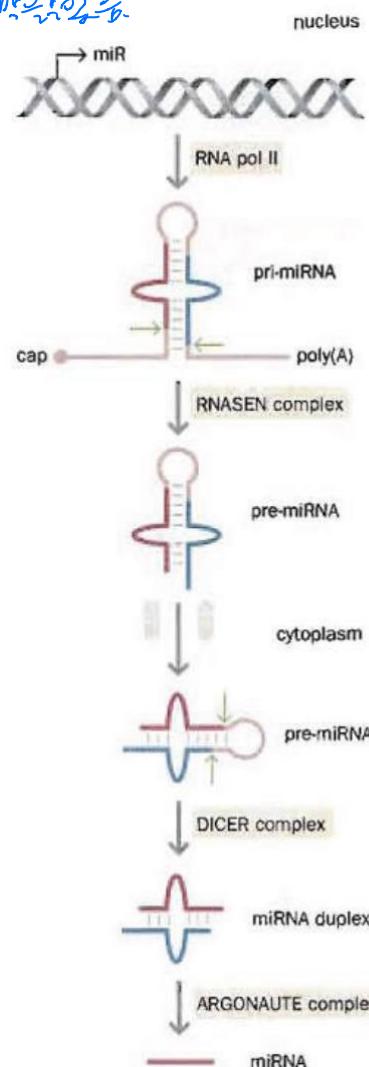
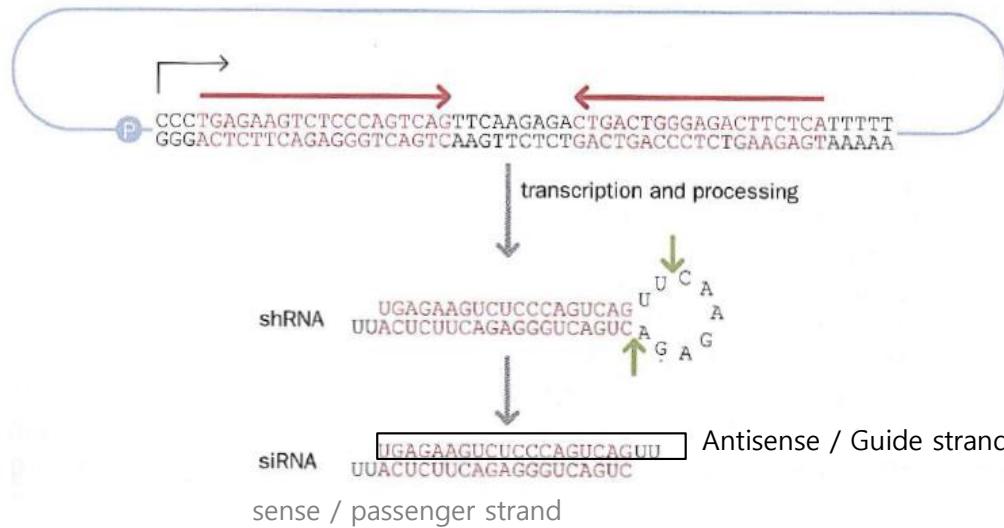
Global RNAi screens provides a systems-level approach to studying gene function in cells

우리 genome에 붙는 여러 종류의 RNA를 디자인 -> 개별적으로 각 gene 발현을 막음
*siRNA는 cis로 작용하기 때문에 단순히 관심있는 gene을 shRNA로 만들 수 있도록 설계하여 클로닝하여 넣어도 효과를 볼 수 있음

- RNAi studies began by attempting to assess the function of individual genes
- In the post-genome era, it became possible to perform large-scale, and eventually ① genomewide, analyses of gene function
- To perform global RNAi screens, suitably ② large nucleic acid libraries (siRNA, shRNA libraries) need to be made
- Short hairpin RNA (shRNA) libraries
 - A short transcript that spontaneously forms a hairpin RNA, which then undergoes cleavage *→ inverted-repeat hairpin structure makes*
 - Pairs of long complementary oligonucleotides are designed in such a way that, when annealed and cloned into a suitable expression vector, a transcript can be produced with inverted repeats that base-pair

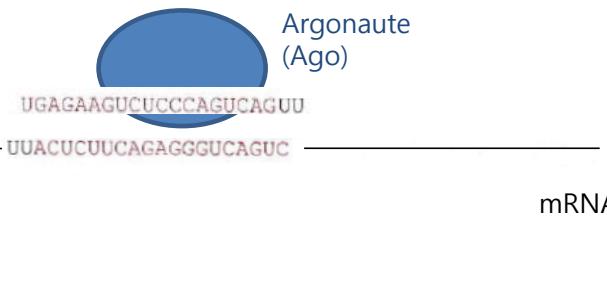
shRNA library

1 probe RNA가 1개의 gene에 8~10개로



RISC (RNA-induced silencing complex)

Target gene

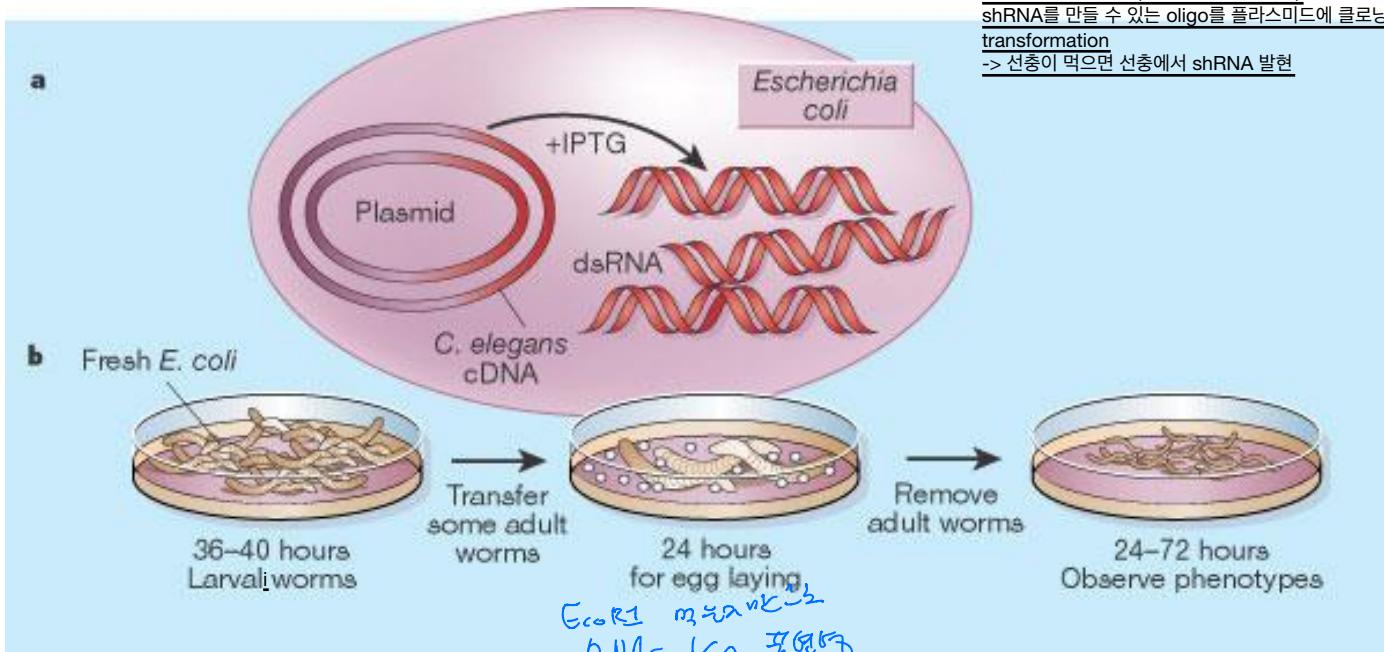


RNAi Screening

②

=> 특정 gene이 표현형에 어떤 영향을 끼치는지 알 수 있음.

선충의 먹이 -> 세균(실험실에서는 E.coli)
shRNA를 만들 수 있는 oligo를 플라스미드에 클로닝 -> E.coli transformation
-> 선충이 먹으면 선충에서 shRNA 발현



Kamath et al. 2003

16,757 strains = 86% of predicted ORFs

Looked for sterility or lethality(Nonv), slow growth (Gro) or defects (Vpep)
1,722 strains (10.3% had such phenotypes)

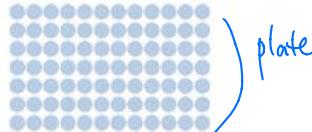
~ 유익한 유전자에 대한 shRNA를 만들고 그에 맞는 유전자는 어떤 defeciton인가?
유익한 유전자를 찾으려면 identity가 높은 유전자를 찾는다.

High-throughput shRNA screening

원

Arrayed Screens

- Hypothesis driven gene sets
- Robust interrogation of each gene
- High Content Screening



→ generally random siRNA
to knock down expression

풀

Pooled Screens

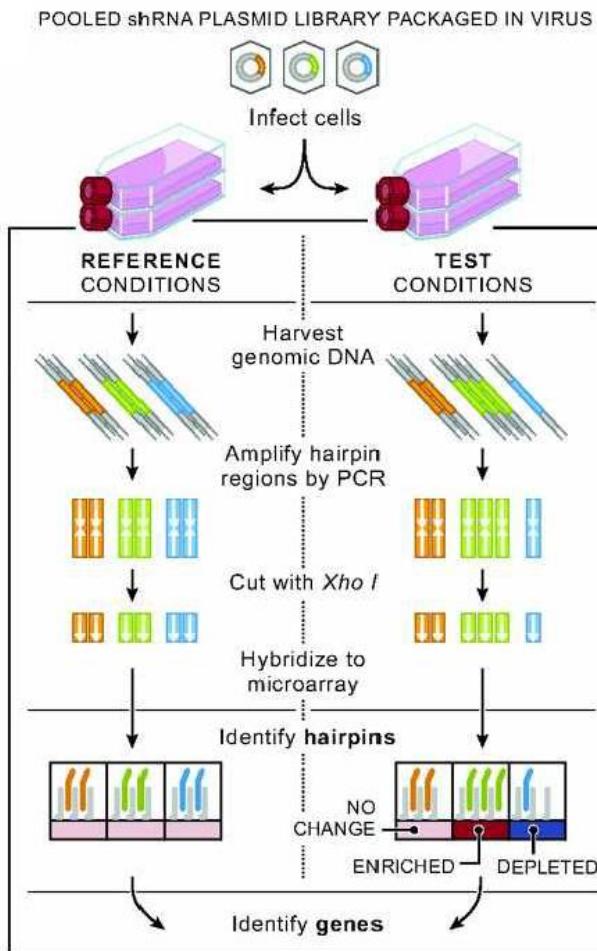
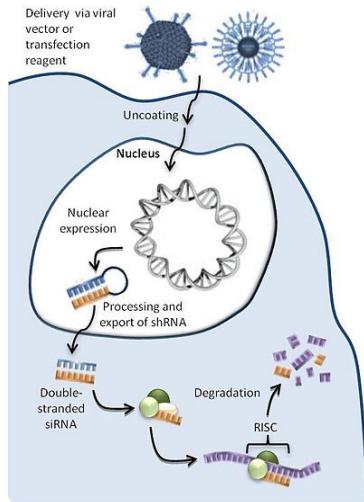
- Genome-scale screening
- Many conditions: drug doses, time points, cell types etc.
- Long time course assays – proliferation and survival
- Lower resolution – readout is differential representation of hairpin measured by microarray or sequencing

DE sh RNA 원고
→ cell 종류와 조건
RNA가 같은 걸 찾는다
제작되는 걸 찾는다
모든 걸 찾는다



High-throughput shRNA screening

Pooled RNAi screening strategy and performance using pools of 45,000 shRNA-expressing viruses.



cancer cell 을 가우고
drug sensitive gene 을 찾고자 함
hairpin 은 단일 유전자 vector로
→ 특정 cell에서만
→ 특정 gene KD
PGK promoter = sense - antisense

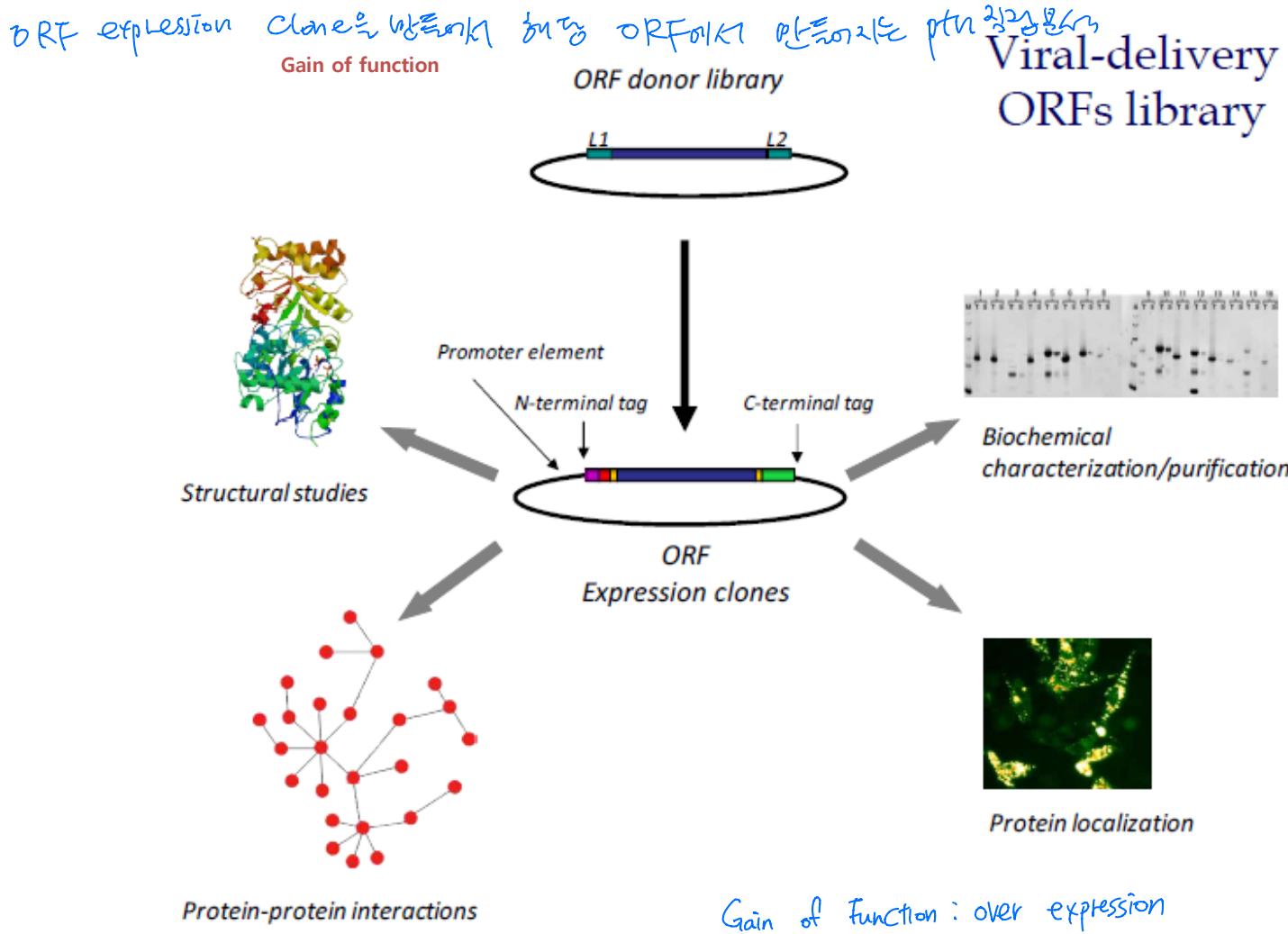
Pool of siRNA
+ / - Treatment (condition)

- Identification of shRNA clones
- sh RNA가 드나드는지 없는지 체크
- Samples vs. controls (Ratio)

- Identification of genes

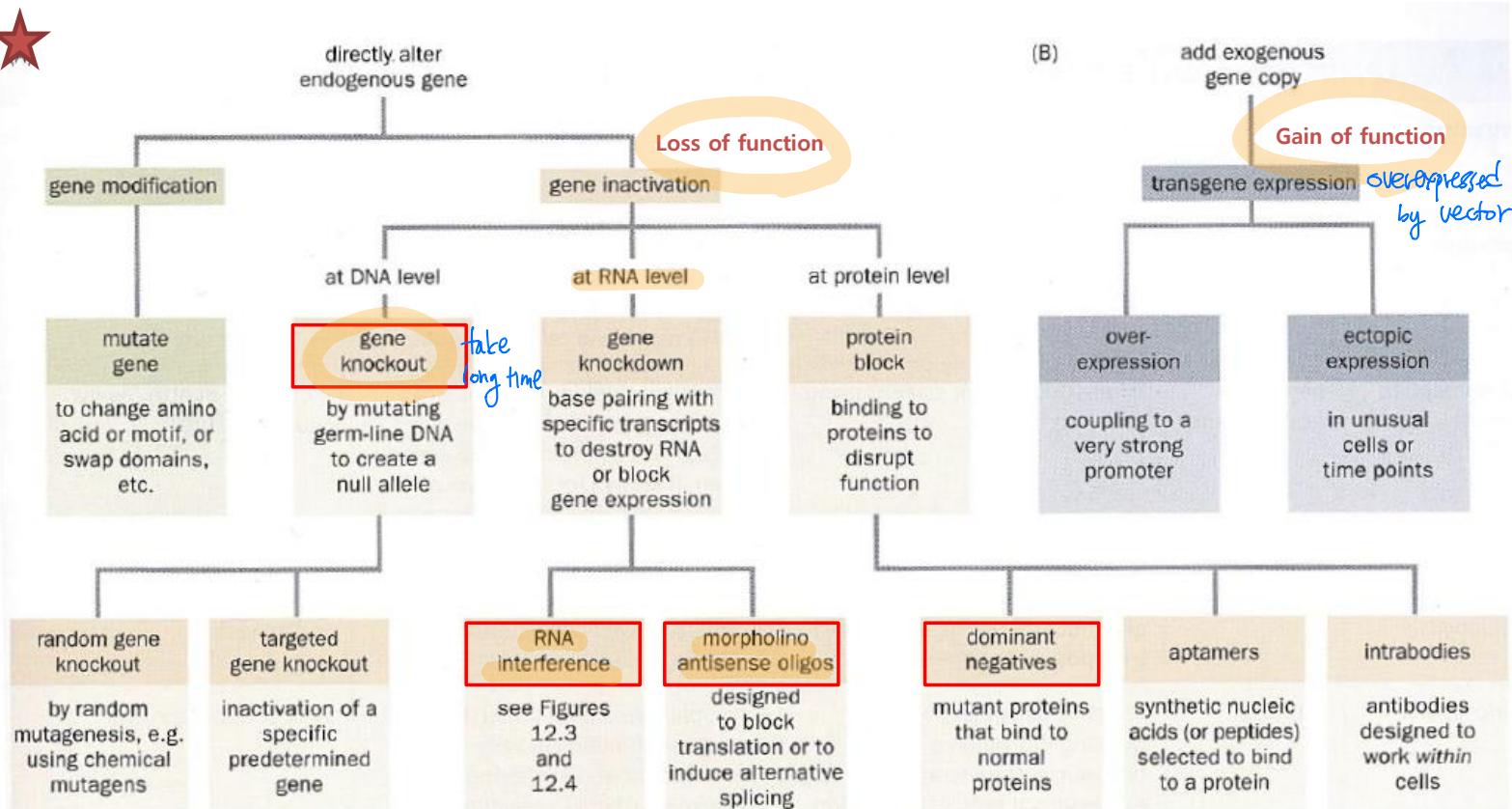
Drug이 DB에 드나드는지 체크
이후에 유전자 찾고자 함
= shRNA가 드나드는지 체크
→ 그에 드나드는 유전자 찾고자 함
Sensitive

Open Reading Frames (ORFs) expression



Selective gene inactivation and modification

- Study of gene function in cultured cells or using cell extracts has limitations
- Defining gene function in this wider context requires the genetic manipulation of model organisms



youtube : Perkin Elmer Screening

