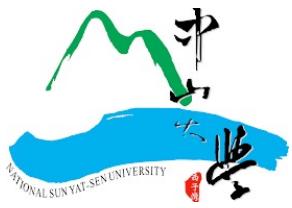


Genomics- Human Genome Project

基因組學 - 人類基因體計劃及其應用

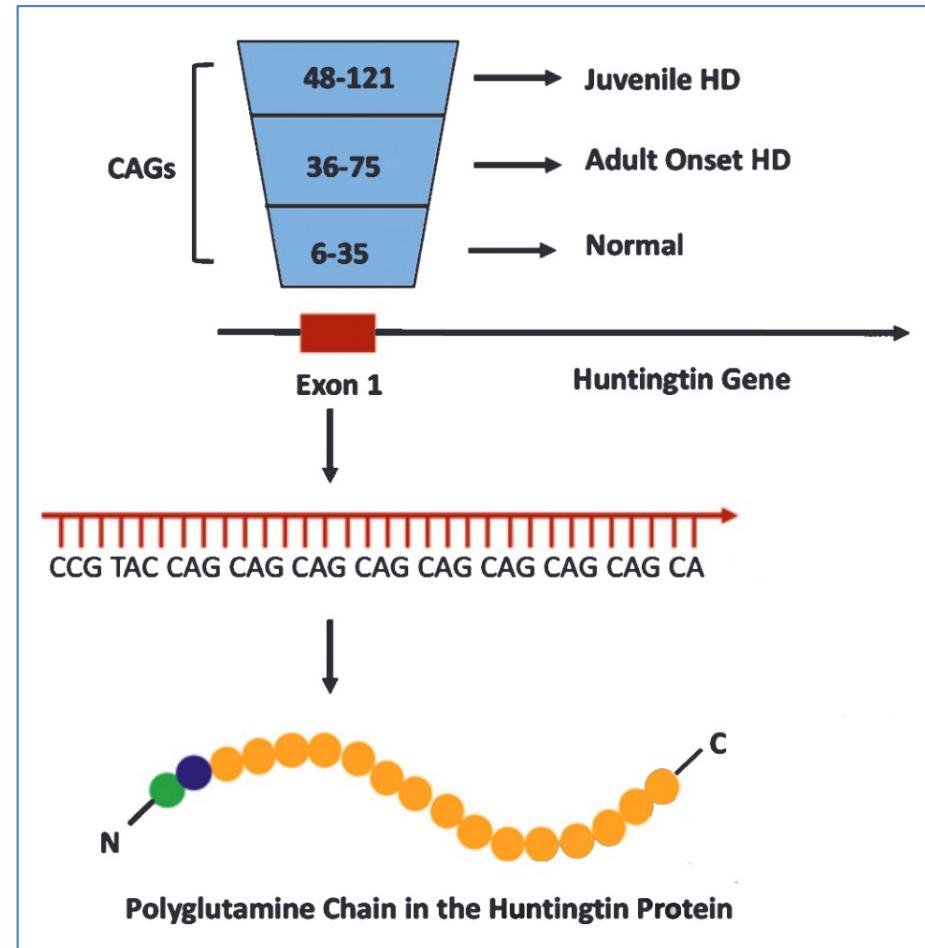
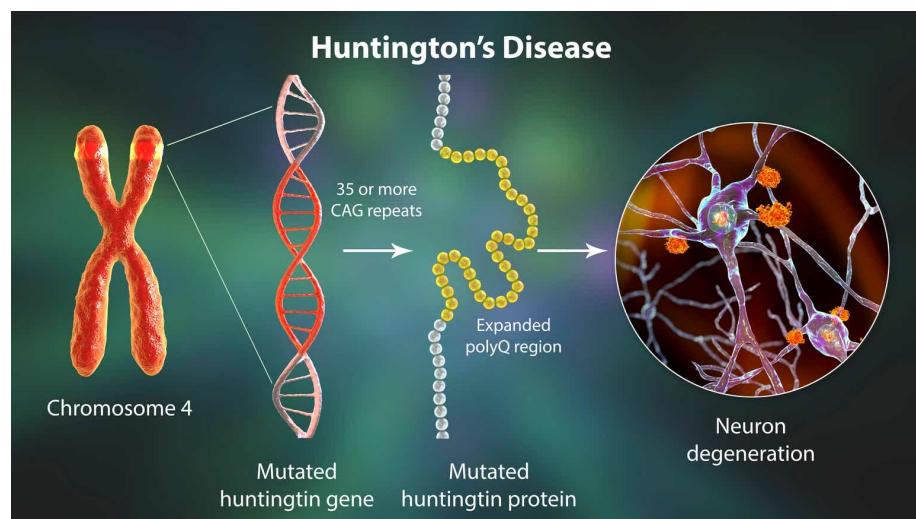


國立中山大學 生物科學系 黃明德

在人類基因組解序之前科學家要如何尋找致病基因？

- 亨丁頓舞蹈症-Huntington's Disease

- 患者會不自主運動、臉部輕微的抽搐
- 發生率約1/10,000 (U. S. A.)
- 1983年HD基因定位於第四號染色體
- 1993年HD基因成功被選殖
- 患者HD基因具有較多CAG重覆於第一外顯子
- 目前尚不明瞭HD蛋白之功能
- 患者發病年齡和CAG重覆有關



Dr_Microbe/iStock via Getty Images

如何透過遺傳學尋找致病基因

I. 遺傳族譜分析

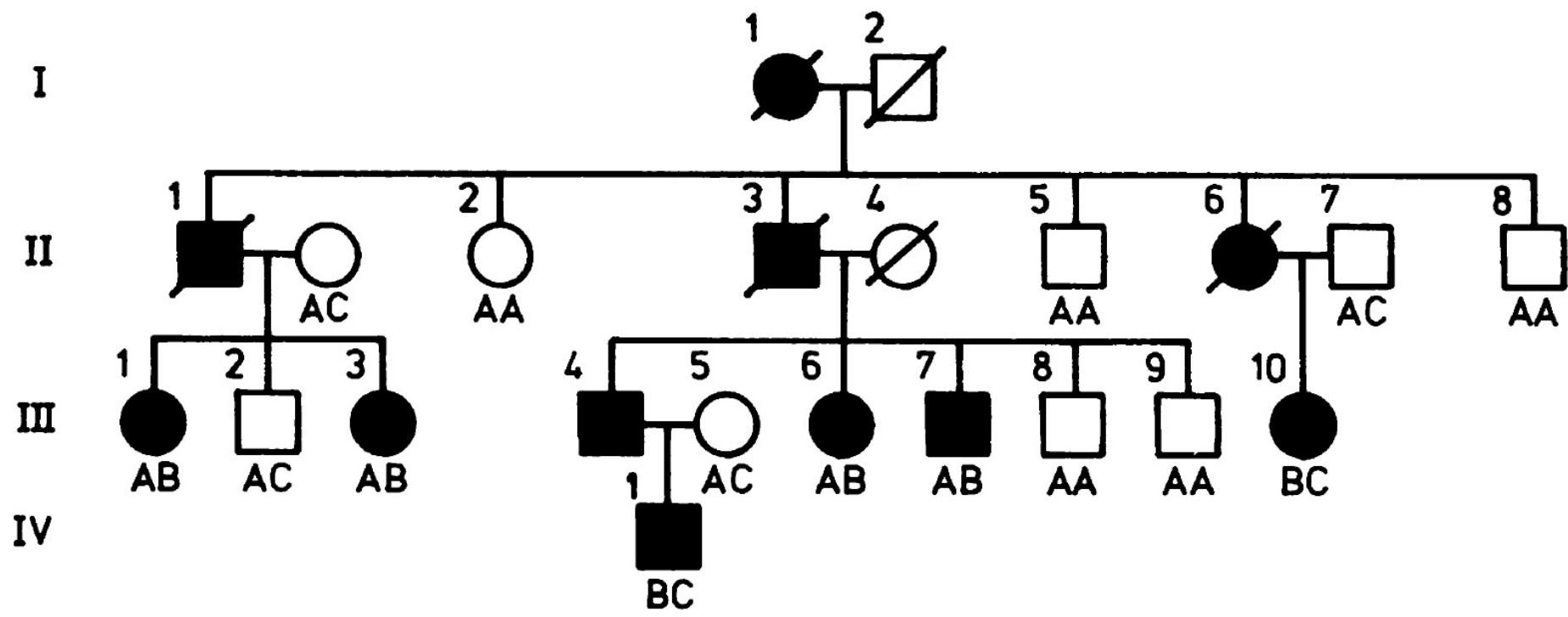
- 顯隱性分析
- 體染色體或性聯遺傳

II. 染色體定位

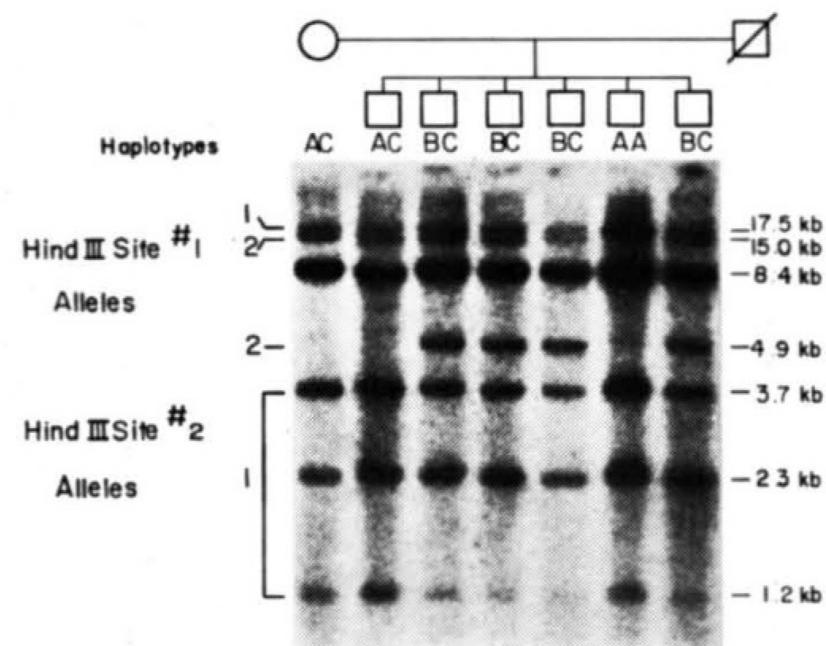
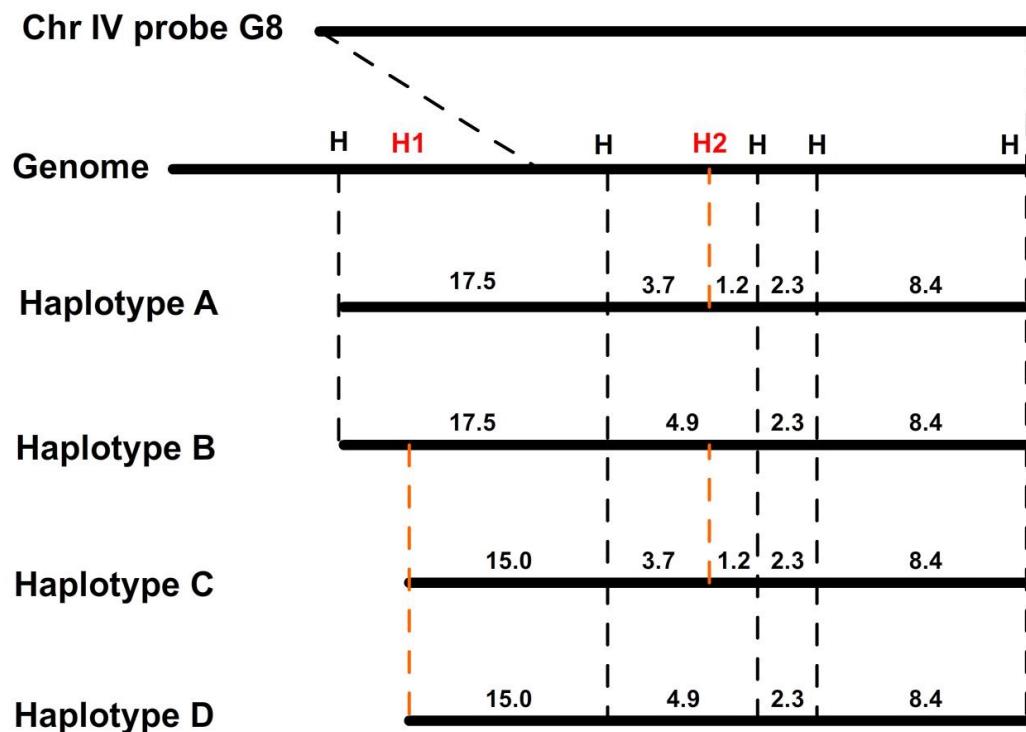
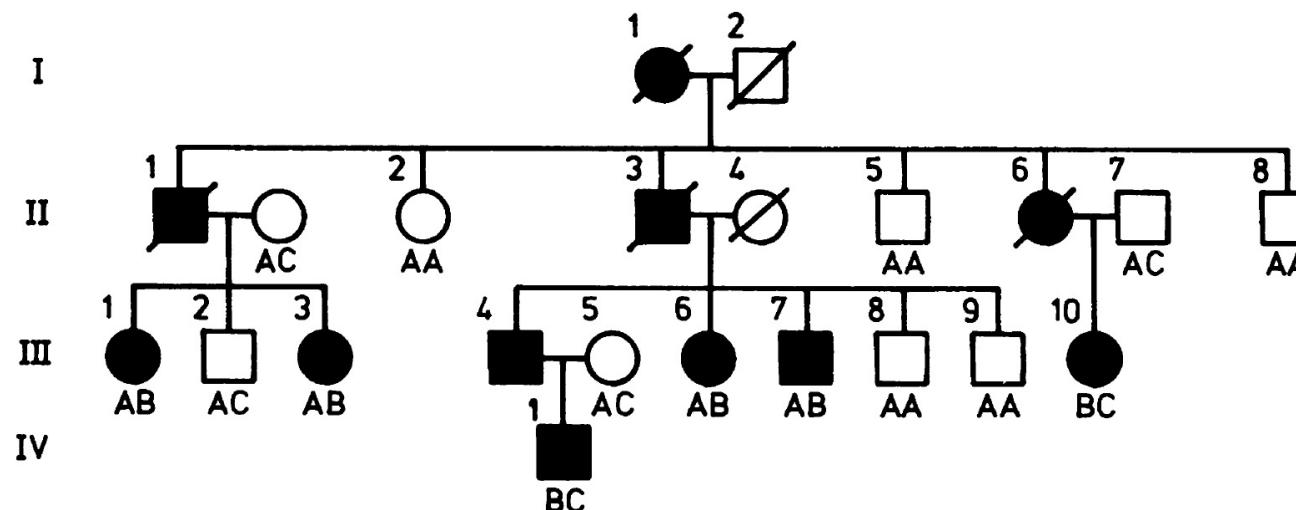
- 分子標誌建立 & 分析

III. 基因序列確認

- 分子選殖



Genetic linkage between Huntington's disease and the DNA polymorphism G8



- 利用性狀或分子標誌(molecular marker)標定基因在染色體上的位置



性狀連結

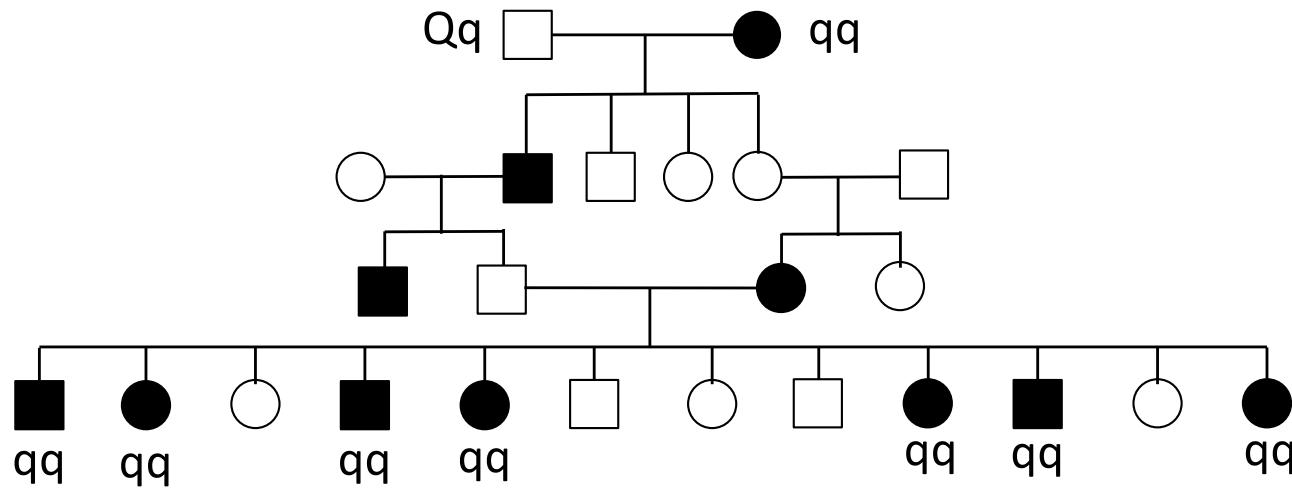
分子標誌：

在不同個體間具有多型性的DNA序列，可透過PCR或Southern blotting檢測

分子標誌

Codon Position	1 123	2 123	3 123	4 123	5 123	6 123	7 123	8 123
Sequence #1	AGC	CTG	ACG	CTA	CGT	CGT	ATA	CCT
Sequence #2TCT.
Sequence #3	A..
Sequence #4A	..TC	...
Sequence #5AC	...
Sequence #6	..T	..AG
Sequence #7	..TA
Sequence #8A.
Sequence #9	..T	..AG
Sequence #10	..T	..A

家族圖譜與基因定位

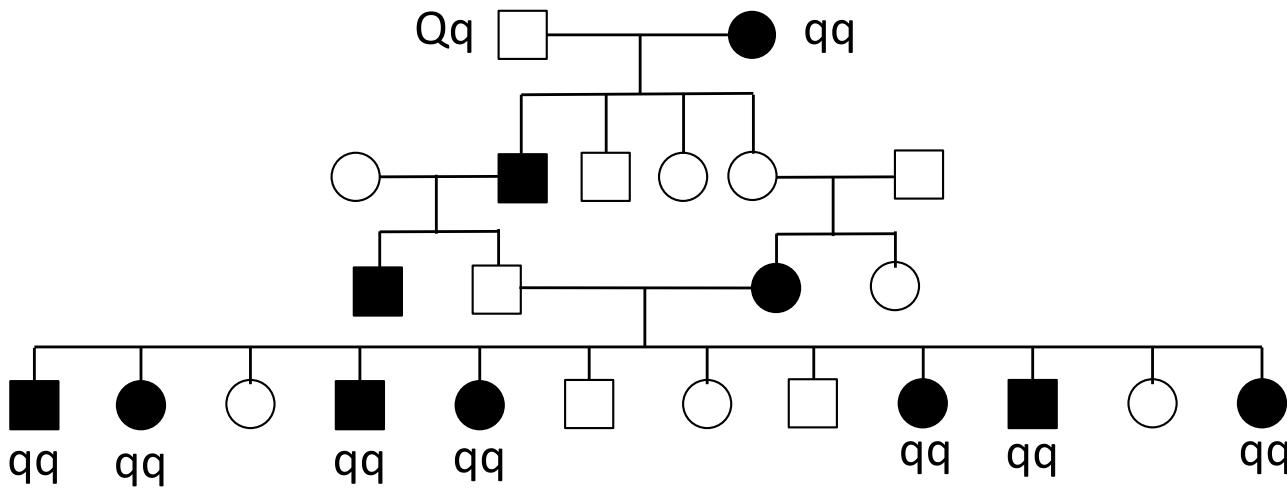


#1	aa						
#2	Bb						
#3	cc						
#4	dd						
#5	ee						
#6	ff						

假設: 曾祖父(Qq)，其餘基因座為顯性（大寫）。曾祖母(qq)，其餘基因座為隱性（小寫）

Q: 基因座Q和那個基因座高度連鎖？

家族圖譜與基因定位



#1	aa						
#2	Bb						
#3	cc						
#4	dd						
#5	ee						
#6	ff						

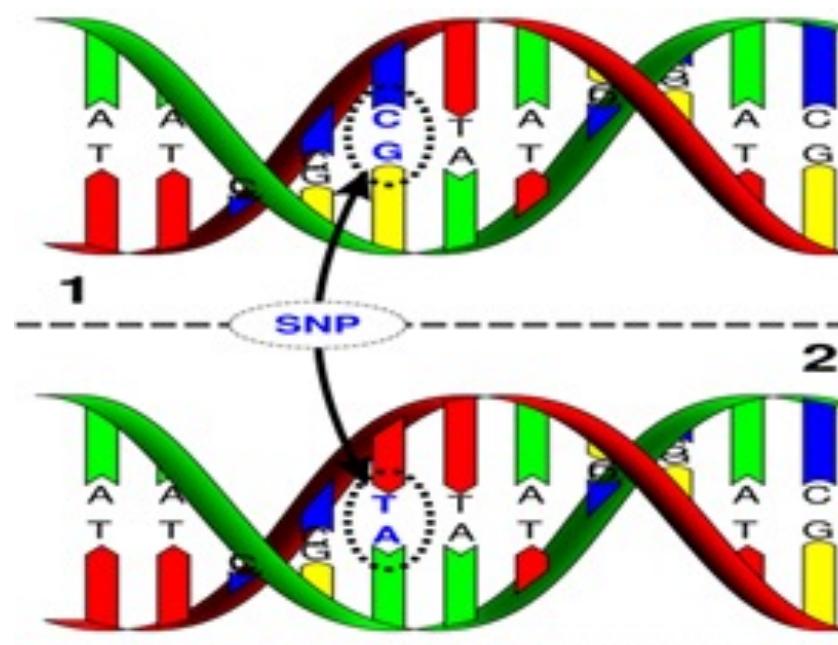


Molecular markers

- DNA polymorphisms are the different DNA sequences among individuals, groups, or populations and can serve as a **molecular marker** for its own location in the chromosome, including SNP, satellites DNA, DNA insertion or DNA deletion and transposon, etc.
- Method to detect molecular markers
 - Morphological
 - Biochemical
 - enzyme activity detection
 - Molecular (PCR or hybridization)
 - SNP (single nucleotide polymorphism)
 - Satellite DNA
 - AFLP (amplified fragment length polymorphism)
 - RFLP (restriction fragment length polymorphism)
 - RAPD (random amplified polymorphic DNA)

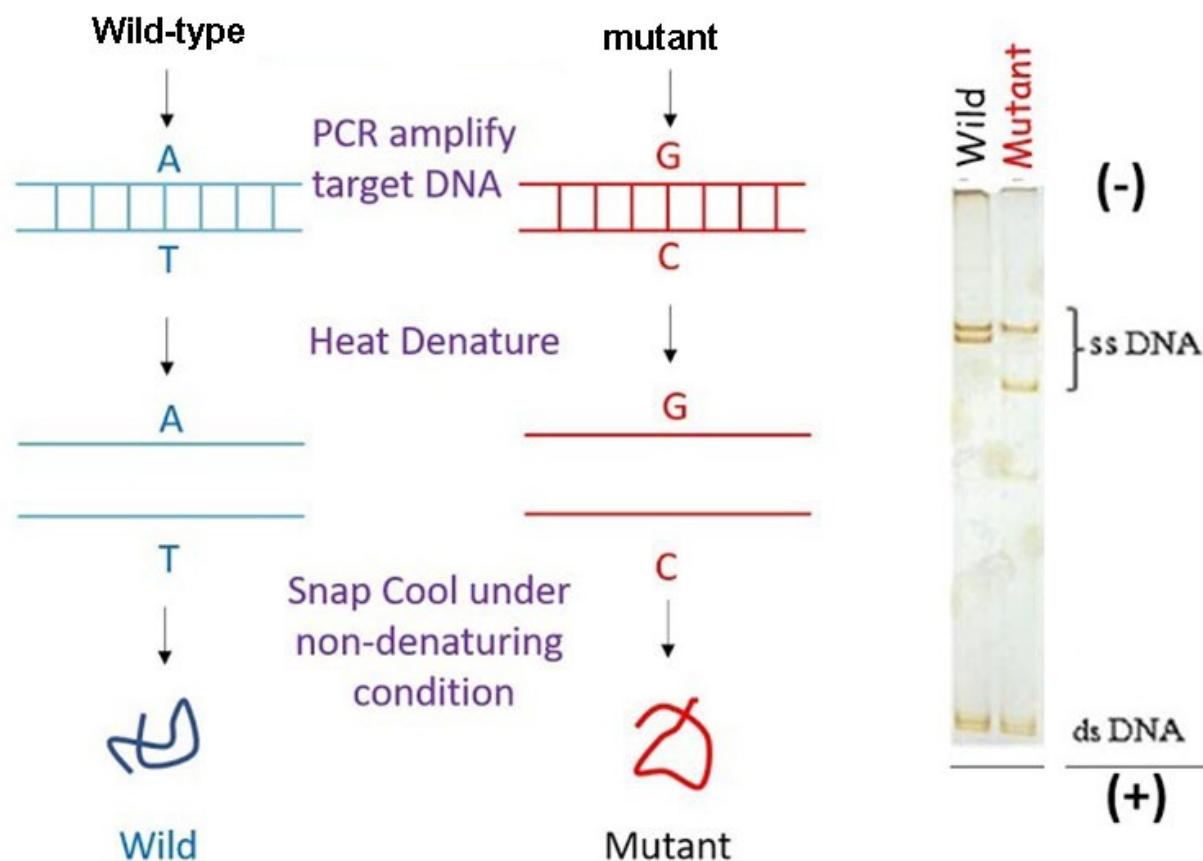
Single nucleotide polymorphism (SNP)

- detected by
 - sequencing
 - SSCP (Single-strand conformation polymorphism)
 - DHPLC
 - DNase cel digestion

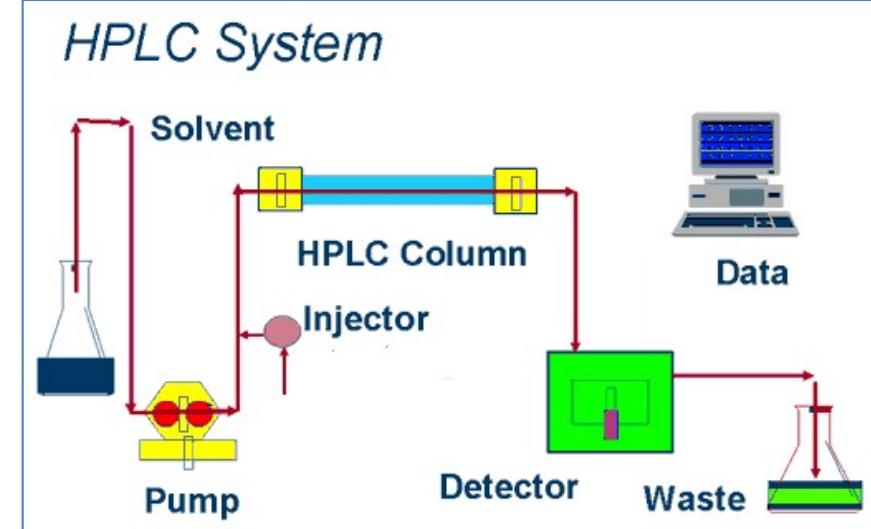
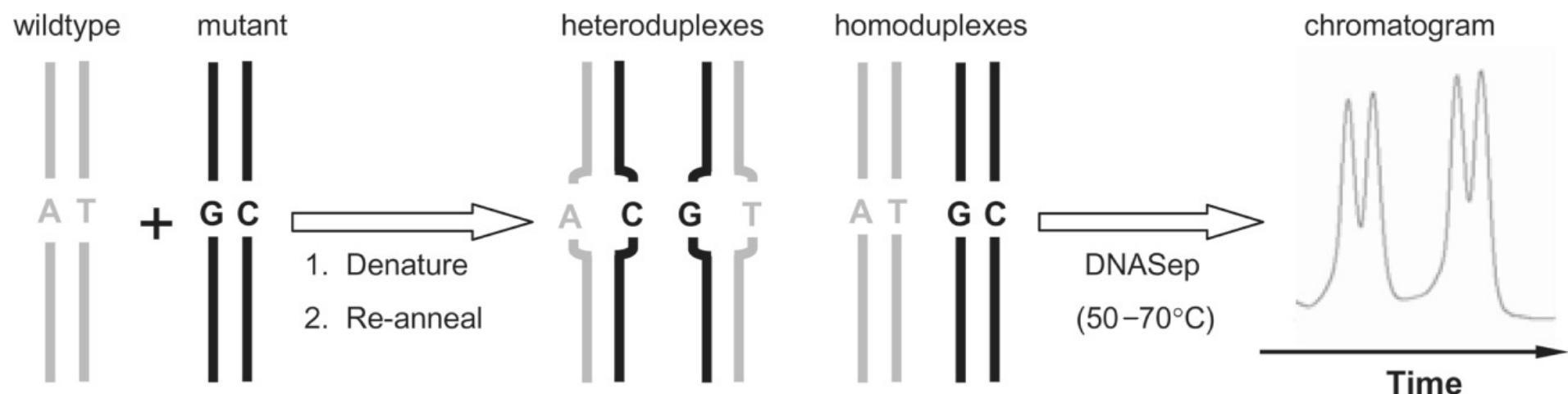


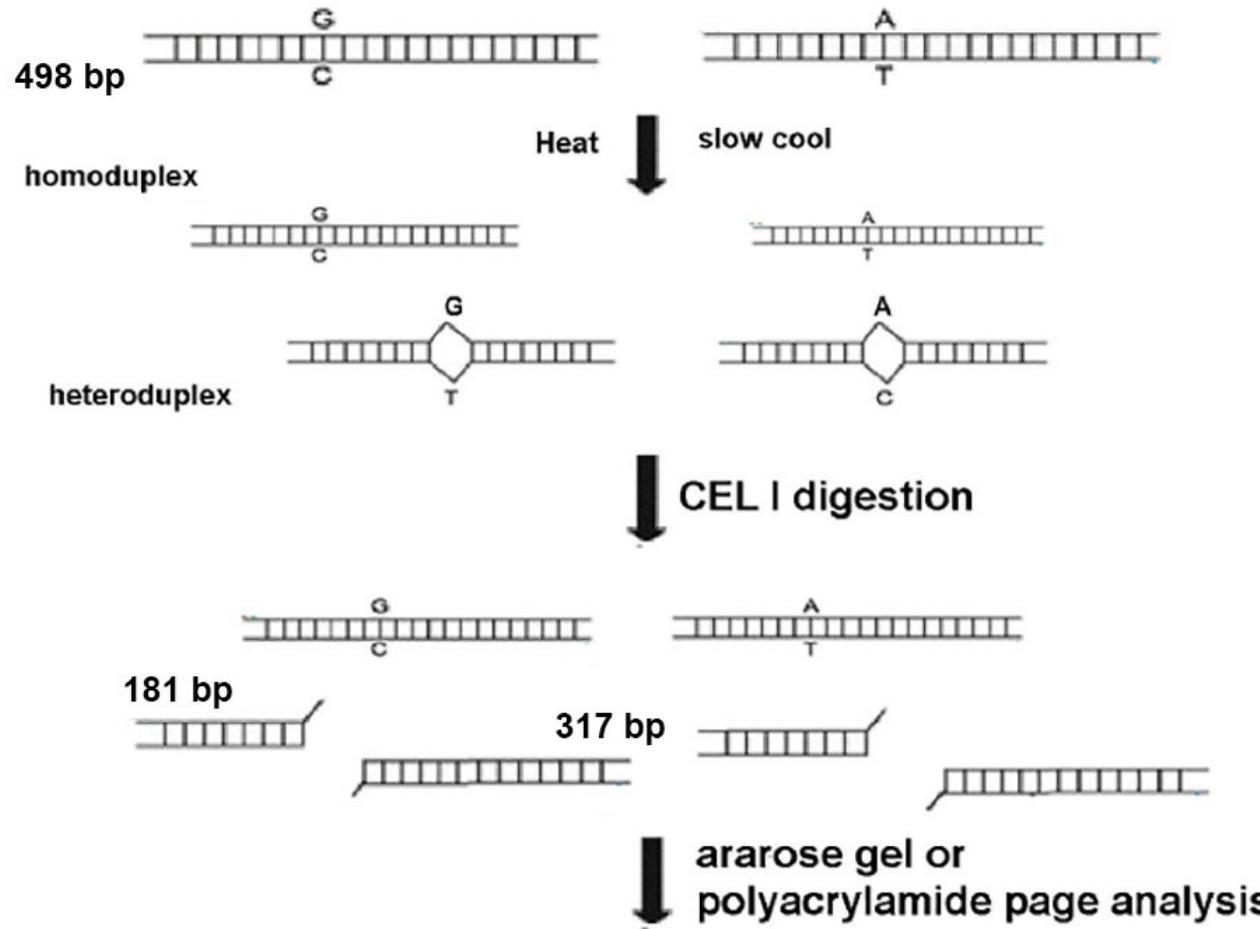
SSCP (Single-strand conformation polymorphism)

- SSCP is the electrophoretic separation of single stranded nucleic acids based on subtle difference in sequence (often a single base pair) which results in a different secondary structure and a measurable difference in mobility through a gel.



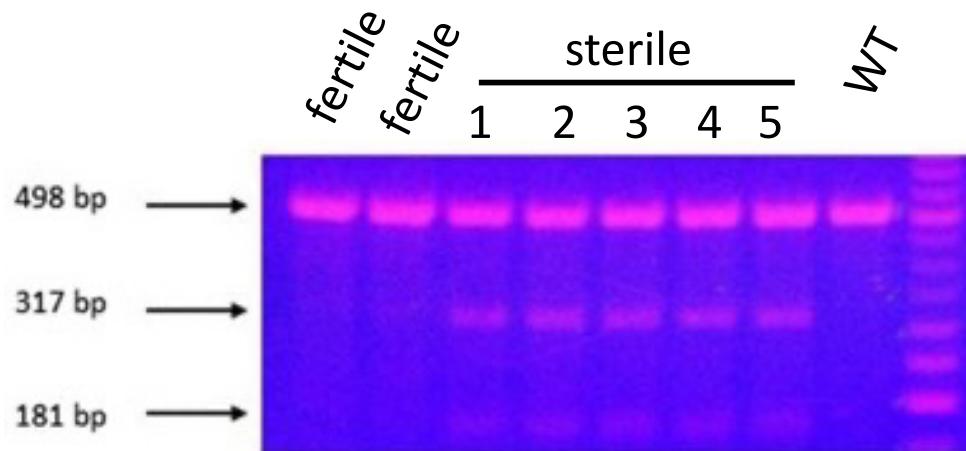
Denaturing high performance liquid chromatography (DHPLC)





Cel nuclease

- **Cel I nuclease**
- From cereley stalks
- Cut at mismatch

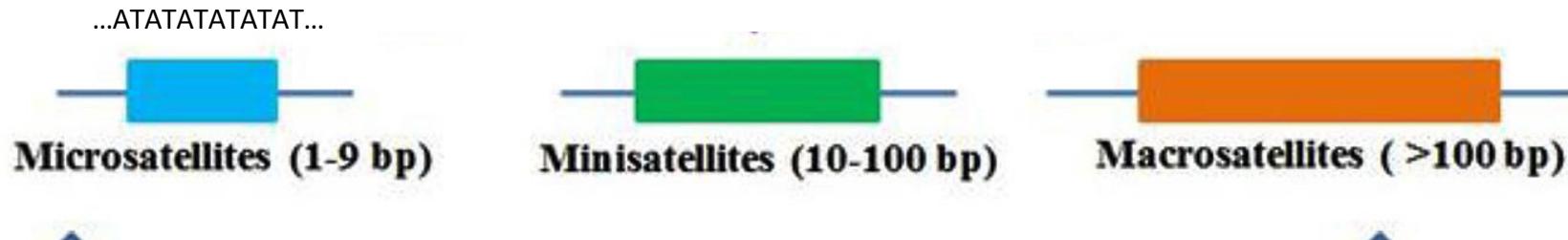


De Guzman et al., 2019

Satellite DNA

- **Satellite DNA**

- noncoding repetitive DNA which are tandem repeats (≥ 5 times)
- microsatellites (1-9 bp repeats) ,
detected by PCR
- minisatellites (10-100 bp repeats) ,
detected by Southern hybridization
- marcosatellites (> 100 bp repeats)
detected by Southern hybridization



Detection of satellite DNA by PCR

- By PCR
 - microsatellite
- By Southern hybridization
 - minisatellite
 - marcosatellites

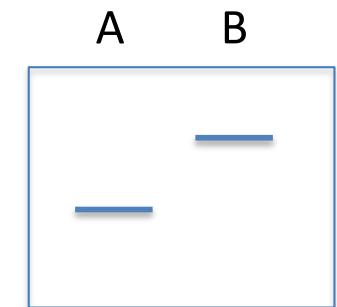
Individual A



Individual B

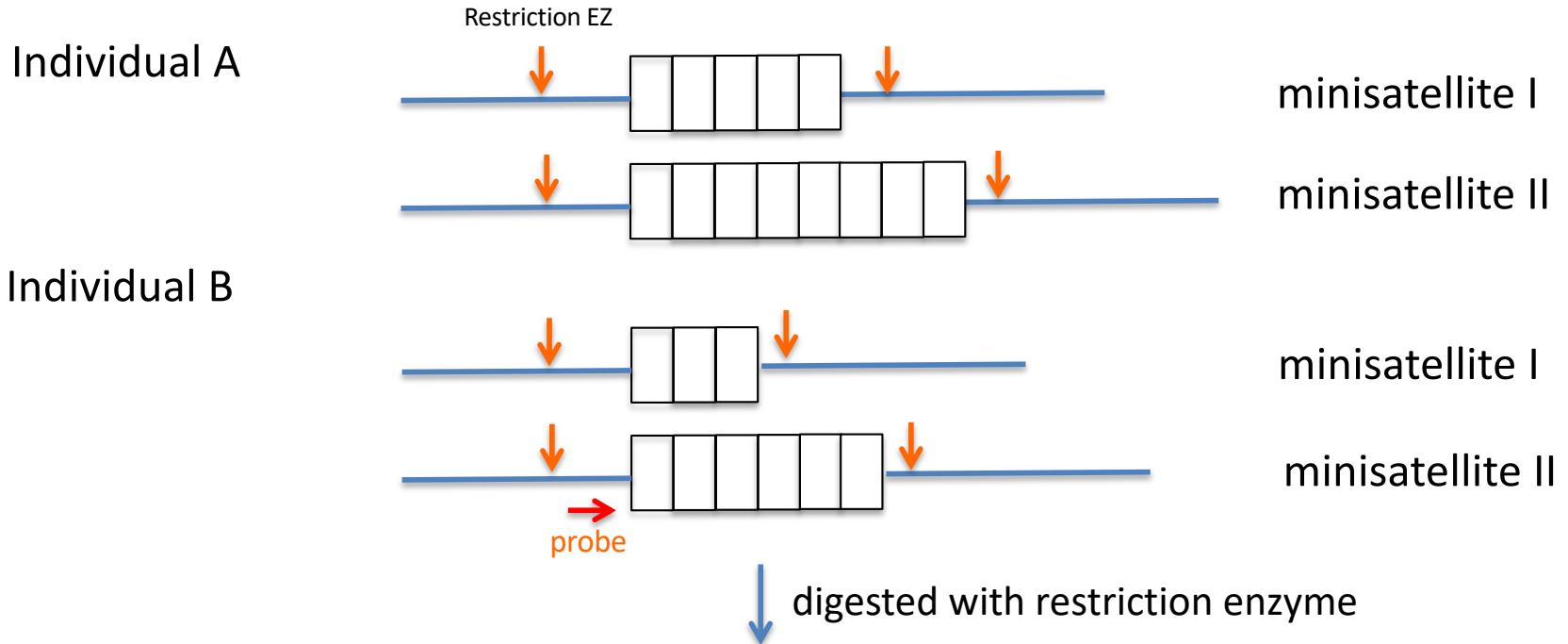


PCR



Analyzed by agarose gel or polyacrylamide page

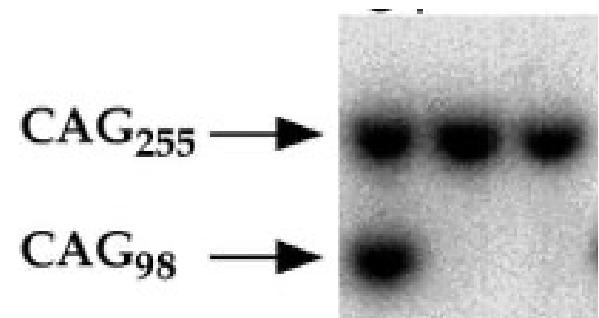
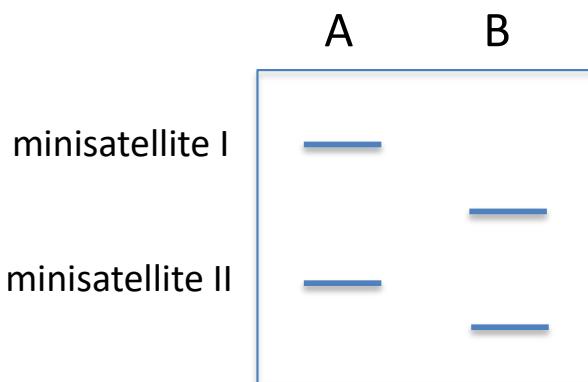
Detection of satellite DNA by Southern hybridization



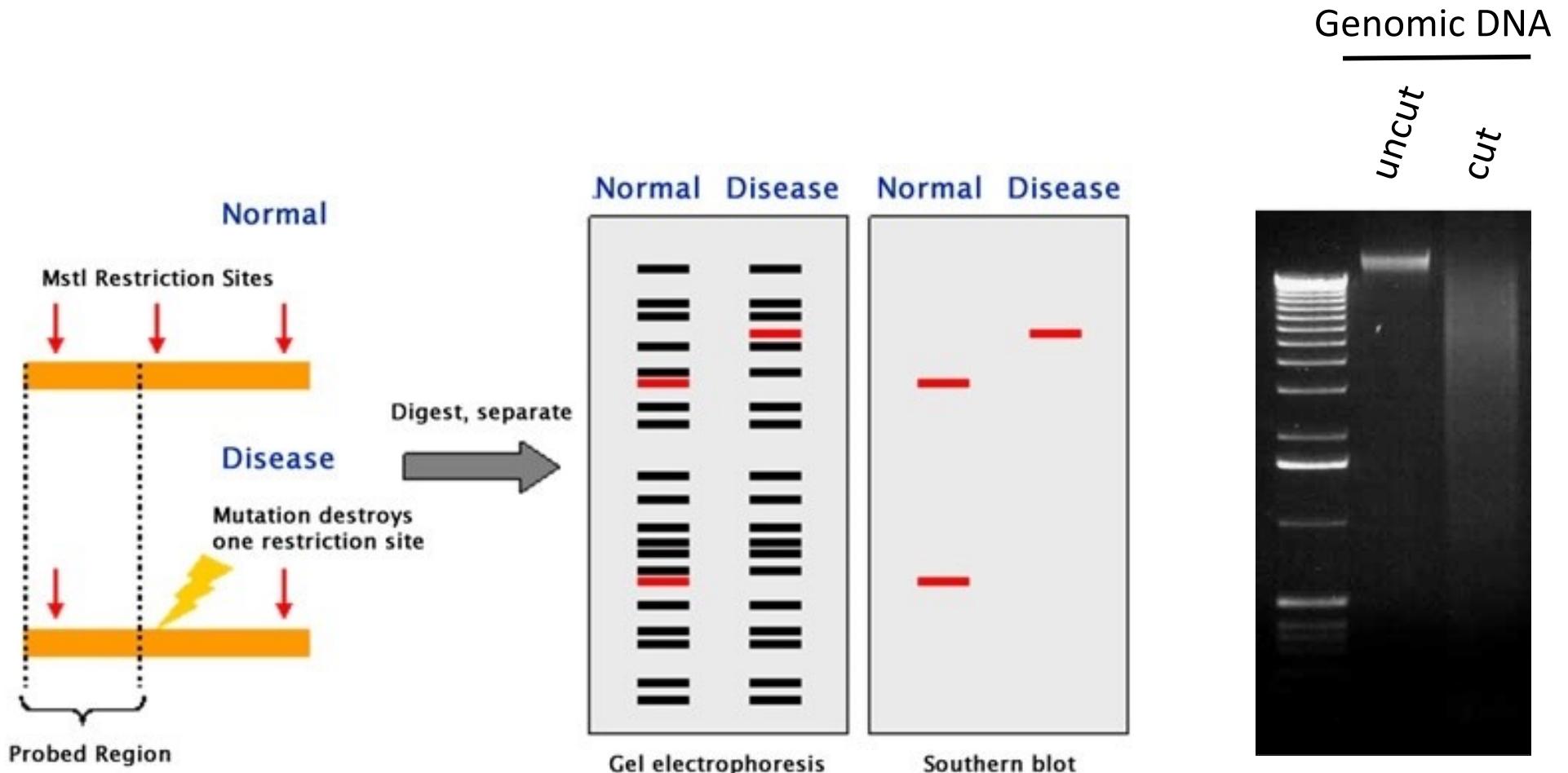
Southern blotting

1. separated with agarose gel
2. detected with probe

 One unit of
Minisatellite
(10-100 bp)



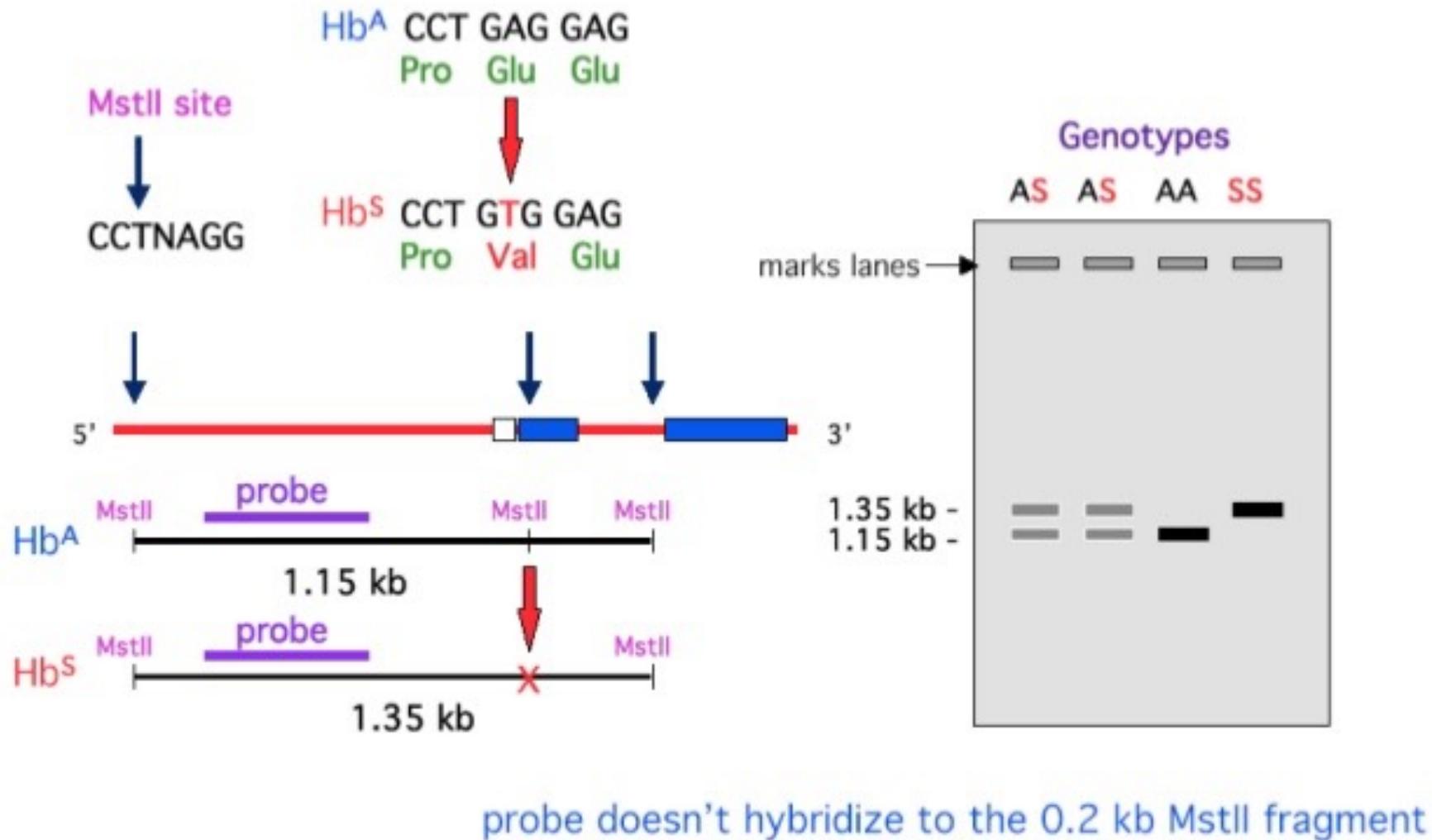
RFLP (restriction fragment length polymorphism)



Example of RFLP

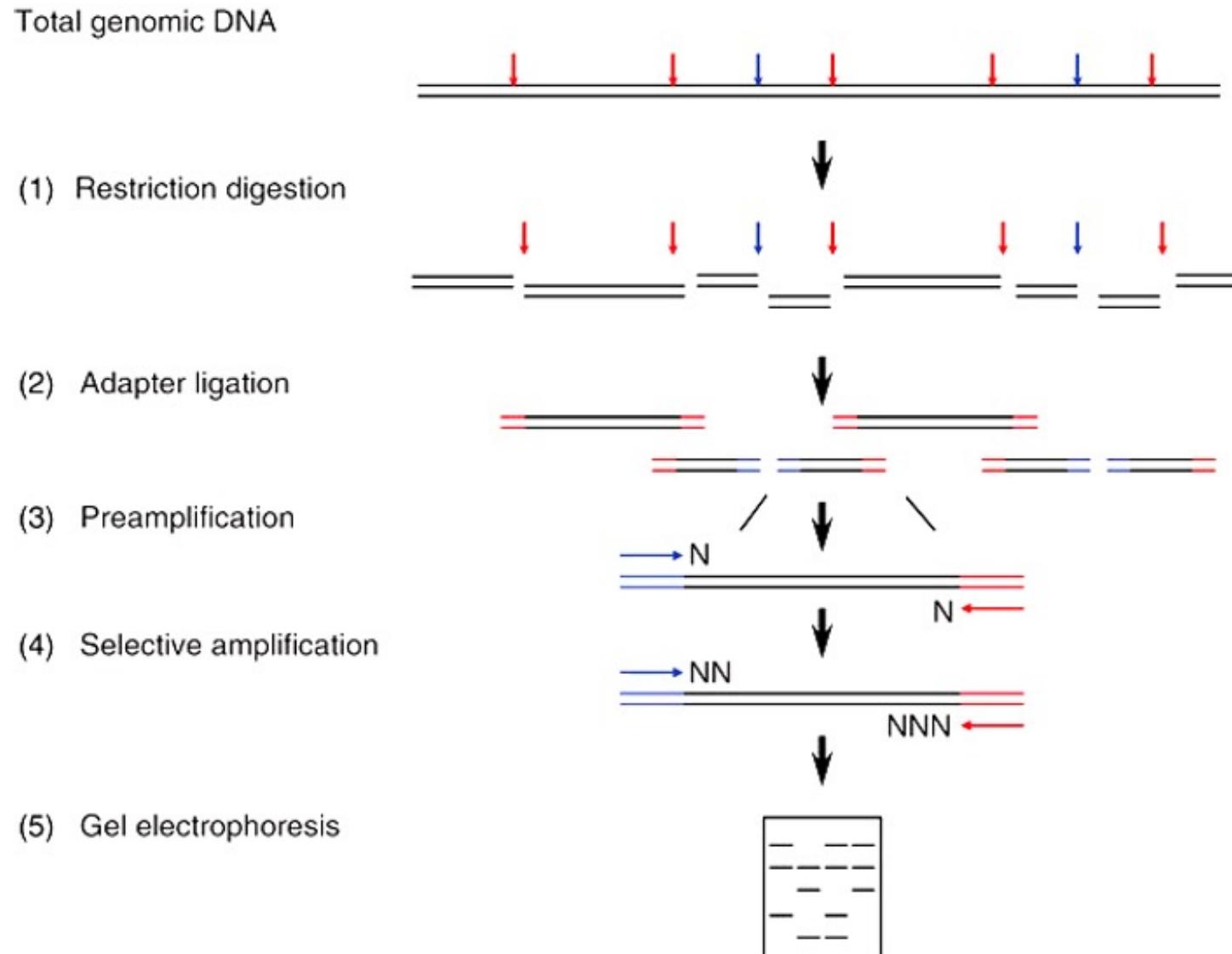


Diagnostic of Sickle Cell Disease/Trait



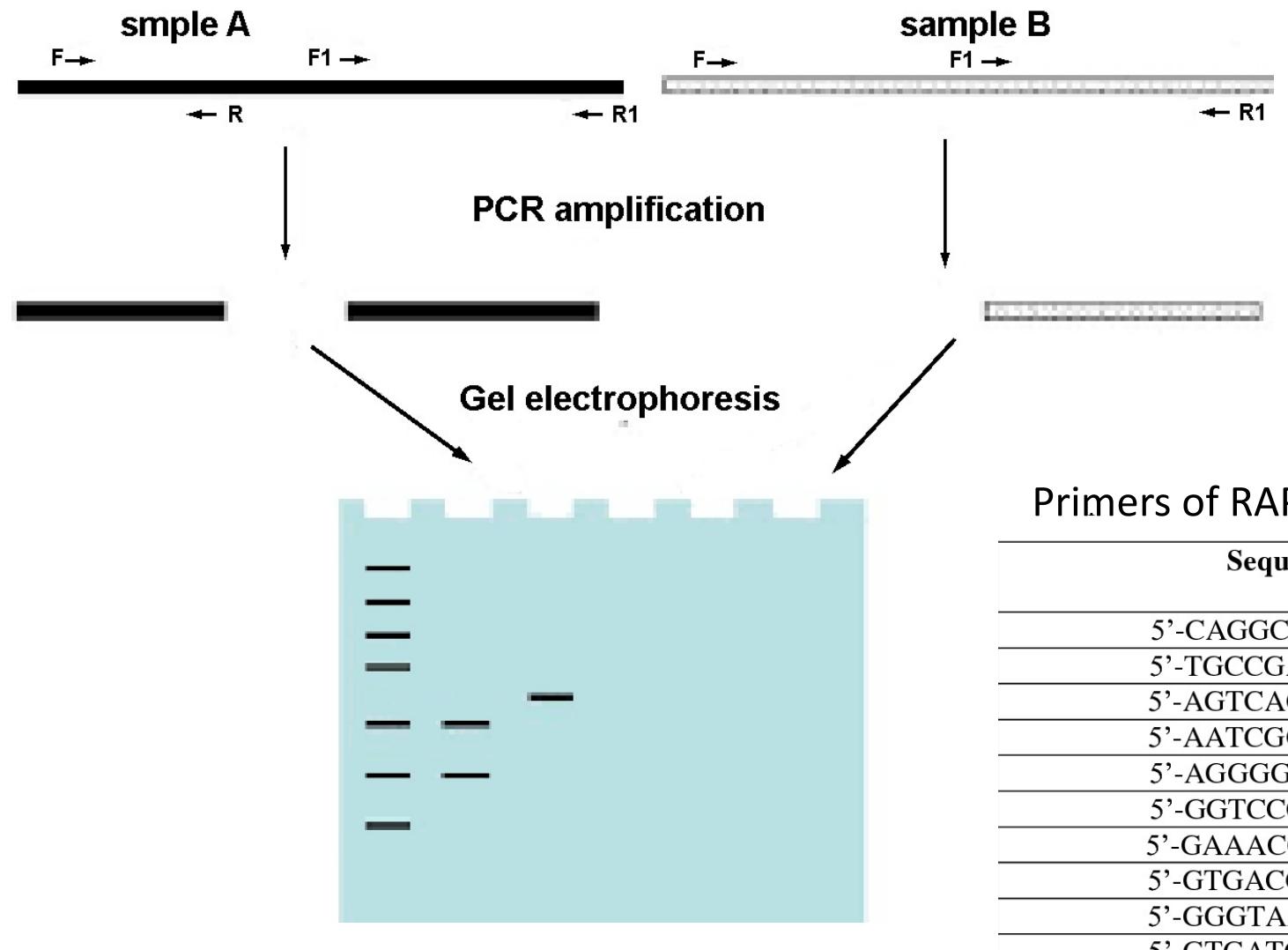
AFLP (amplified fragment length polymorphism)

- Amplified fragment-length polymorphism (AFLP) is a DNA fingerprinting method that employs restriction enzyme digestion of DNA followed by selective amplification of a subset of fragments and separation by electrophoresis on a polyacrylamide gel.



RAPD (random amplified polymorphic DNA)

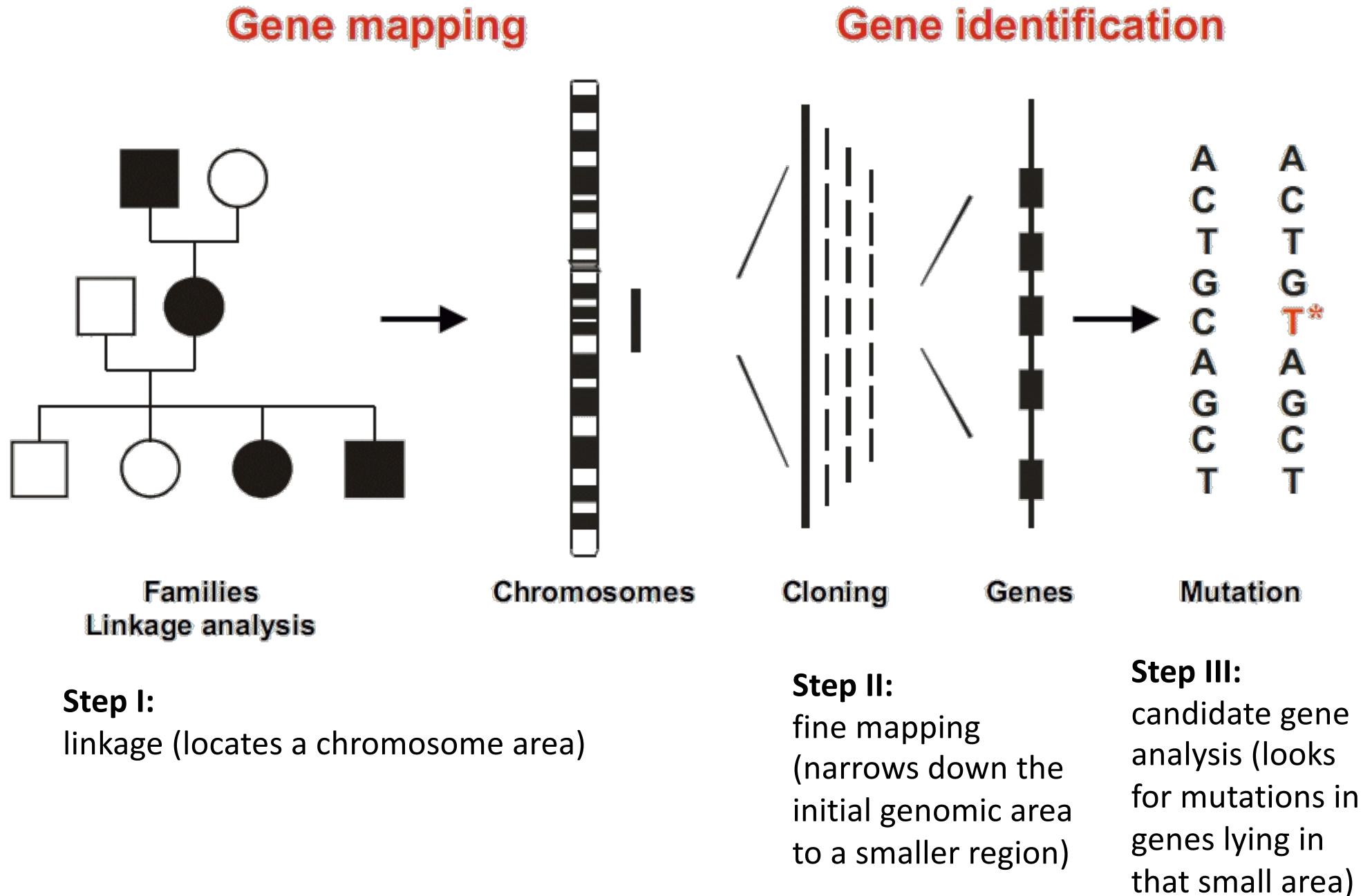
- The DNA are amplified by several arbitrary, short primers (8–12 nucleotides).



Detection of molecular markers

- SNP (single nucleotide polymorphism)
 - by OPCR, SSCP, DHPLC, DNase cel digestion
- Microsatellites
 - by PCR (microsatellite,
 - by Southern blot (minisatellite)
- RFLP (restriction fragment length polymorphism)
 - restriction enzyme, Southern blot
- AFLP (amplified fragment length polymorphism)
 - restriction enzyme, PCR
- RAPD (random amplified polymorphic DNA)
 - PCR

Positional cloning (map-based cloning)



Solomon islands blonde (索羅門群島金髮)



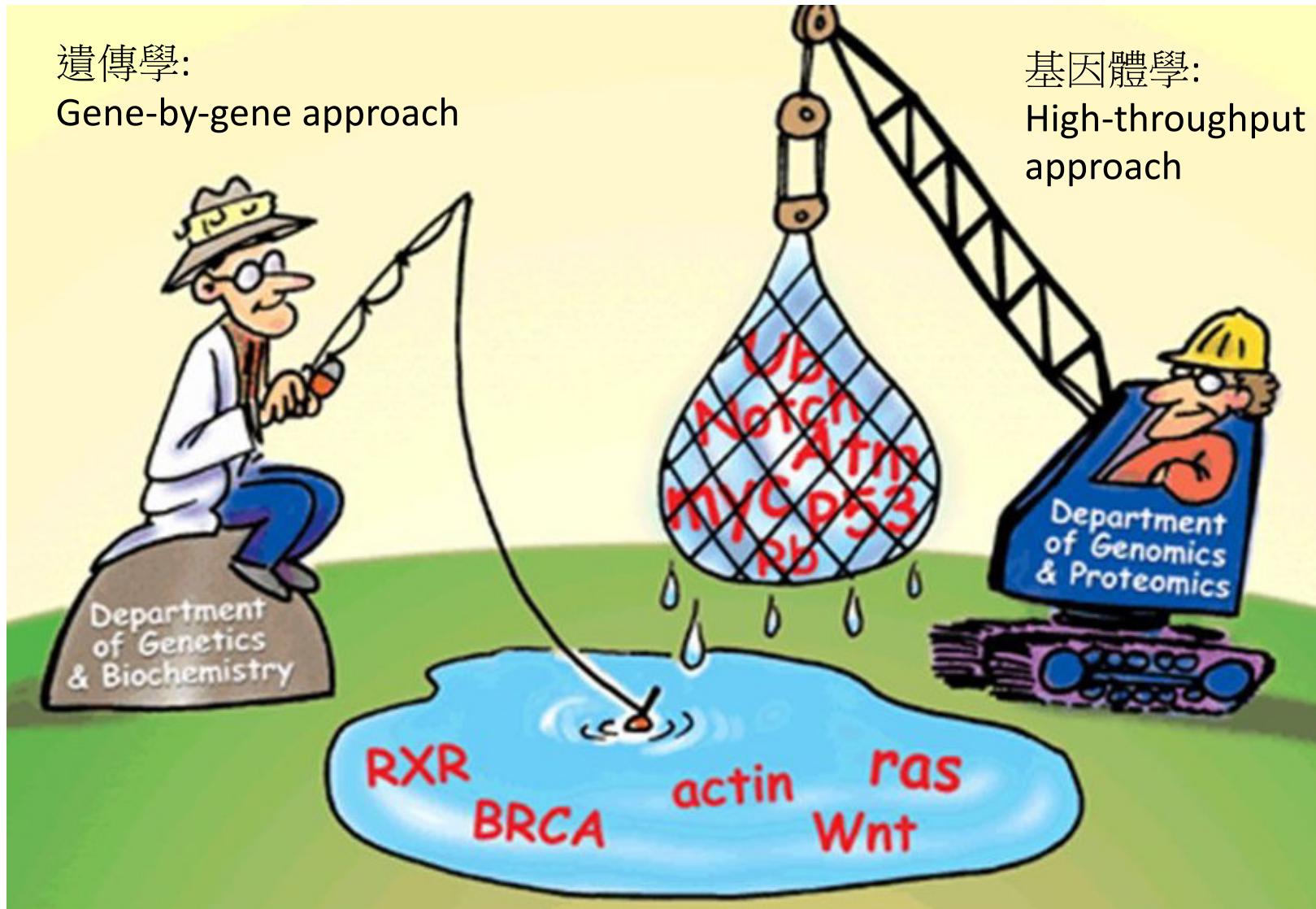
Q: 當金髮不再和膚色性狀相連結，你當如何研究？

基因組學 (Genomics)

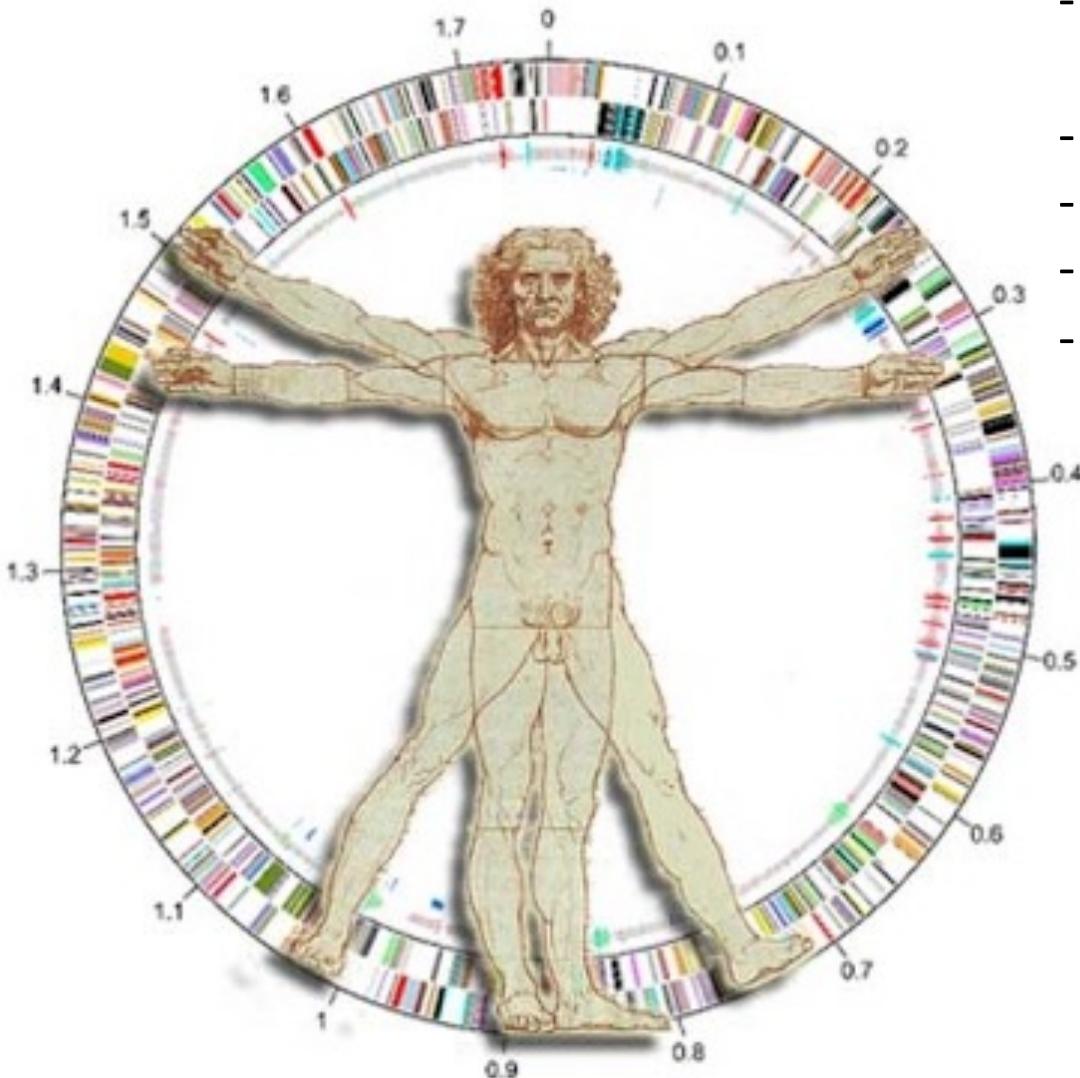
- 基因組(genome)：細胞內所有的DNA，包含核DNA(nuclear DNA)、葉綠體DNA、粒腺體DNA
- 基因組學(基因體學)
 - 目的：研究基因組結構、基因功能、基因演化
 - 工具：
 - 基因組定序 - 基因組完整定序
 - 生物資訊學 - 序列組裝及分析
 - 遺傳學 - 基因功能分析
 - 流程：基因組定序 -> 基因註解 -> 基因功能分析

後基因體時代

- 基因體學與遺傳學相較，為以高通量(high-throughput)策略研究基因功能



Human genome project



- 目的：將人類基因組序列完全定序並註解所有基因
- 1990計劃啟動
- 2003公佈草圖
- 總經費 \$3,000,000,000美元
- 共18個國家參與

人類基因組計劃/Human genome project (HGP)

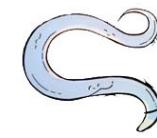
- 1984 – 科學家於美國能源部會議提出構想
- 1986 – 與會科學家再次強調該計劃重要性並討論
- 1988 – 與會科學家一致同意該計劃重要性並準備著手進行
- 1990 - 提出初步構想(為期15年, 經費美金 \$3,000,000,000, 採用階層式定序法)
- 1992 - 發布低解析度基因組草圖(genome map)
- 1998 – Celera公司宣佈將以霰彈槍定序法於五年內完成基因組定序，經費 \$300,000,000，完成後將註冊所有基因
- 1999 – 第一條染色體公布 (chromosome 22)
- 2000 - Celera公司宣佈已完成 ~97%
- 2003 - 人類基因組計劃完成 (99%)
- 2022 - 人類基因組計劃完成 (100%)

基因組大小

- 單位(bp / base pair)

$1 \text{ bp} = 1 \text{ bp}$, $1 \text{ kb} = 1,000 \text{ bp}$, $1 \text{ MB} = 1,000,000 \text{ bp}$,
 $1 \text{ GB} = 1,000,000,000 \text{ bp}$

- 基因組大小

Species	<i>Porcine circovirus</i>	<i>Escherichia coli</i>	<i>Caenorhabditis elegans</i>	<i>Drosophila melanogaster</i>	<i>Homo sapiens</i>	<i>Amoeba dubia</i>
Genome Size	1759 bp	4.6 MB	100 MB	130 MB	3.2 GB	670 GB
Common Name	 Virus	 Bacteria	 Nematode	 Fruit fly	 Human	 Ameoba

可表現蛋白
基因數目

3

4288

19,000

13,600

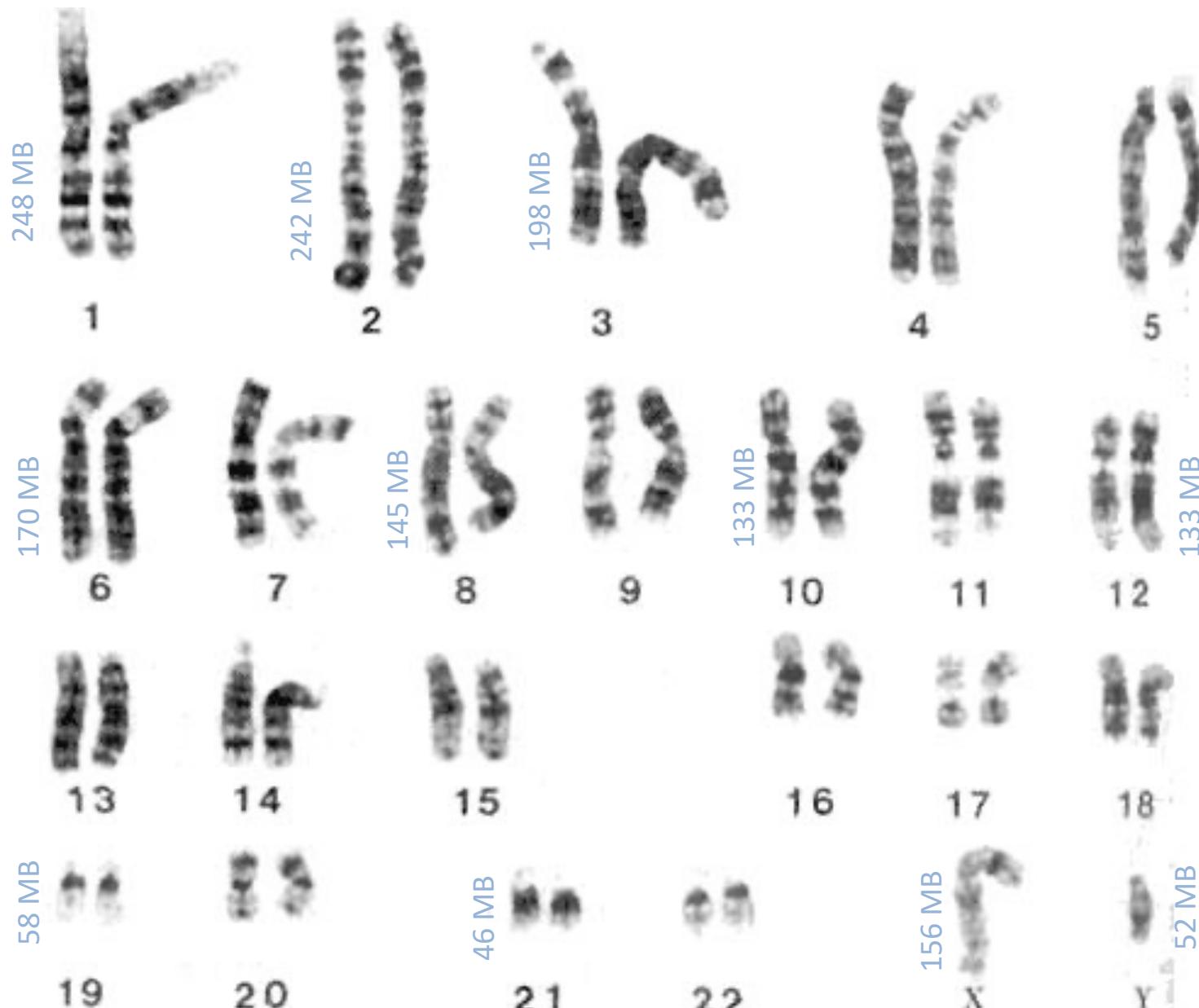
~ 20,000

?

C值謎(C-value enigma): 生物的C值（基因組大小）並不與生物複雜程度相關的現象

G值謎(G-value paradox)：生物的G值（基因數量）並不與生物複雜程度相關的現象

人類染色體數目及其大小

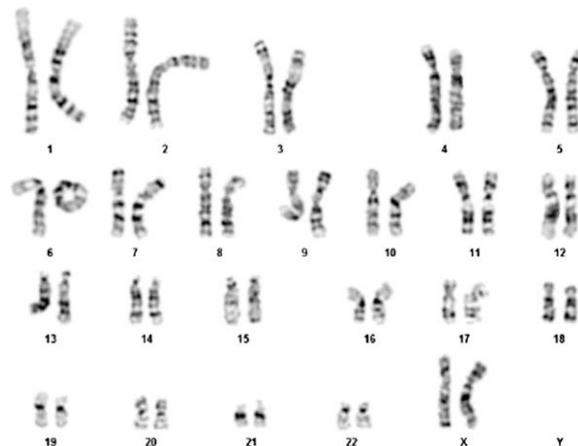


Total: 3,234.83 Mb

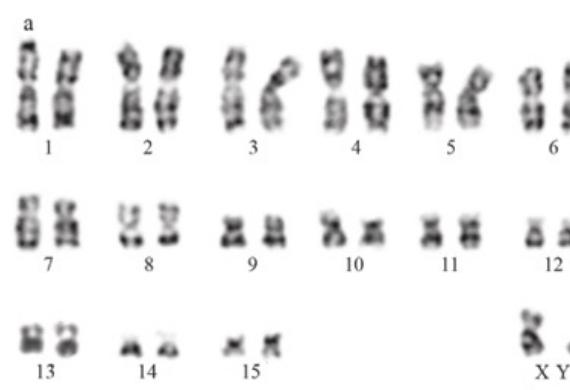
染色體條帶技術

- 利用染劑使染色體呈現各自獨特條帶形態，藉以區別染色體的不同

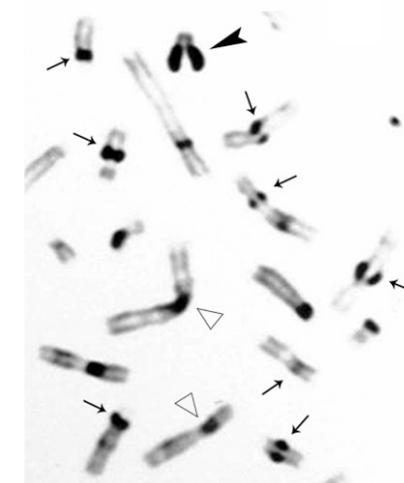
技術	方法	亮帶	暗帶
G 帶	胰酶 + Giemsa	GC rich	AT rich
R 帶	塩處理 + Giemsa	AT rich	GC rich
Q 帶	Quinacrine	GC rich	AT rich
C 帶	$\text{Ba}(\text{OH})_2$ + Giemsa	著絲點以外	著絲點



G-banding

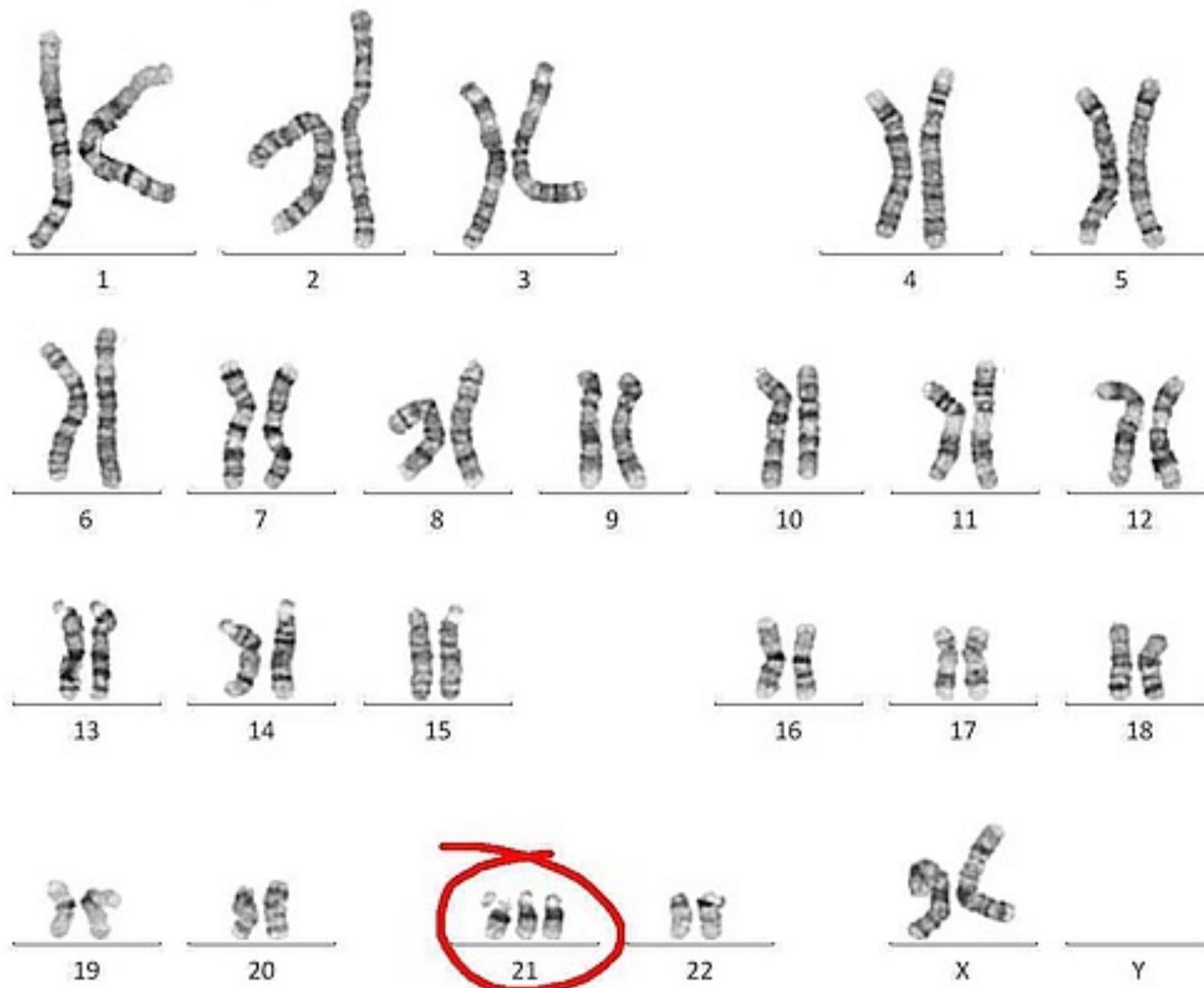


R-banding



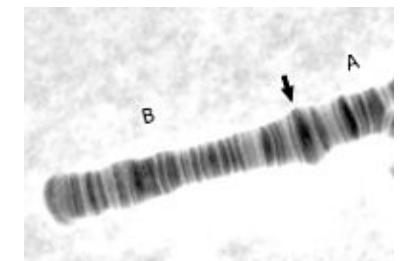
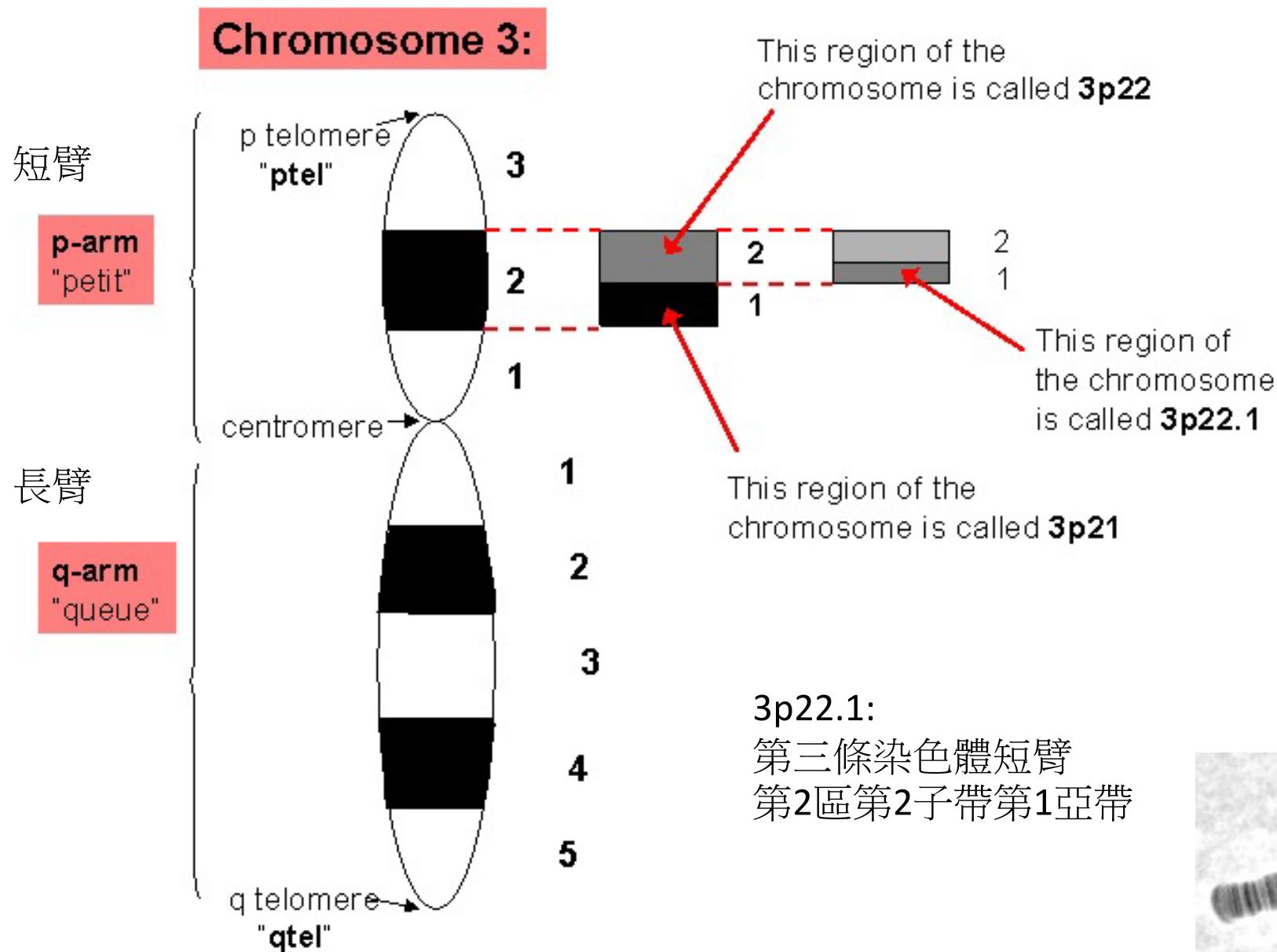
C-banding

唐氏症-第21號染色體異常



<https://www.quora.com/How-does-having-an-extra-chromosome-cause-Down-syndrome>

細胞遺傳圖譜 (cytogenetic map)

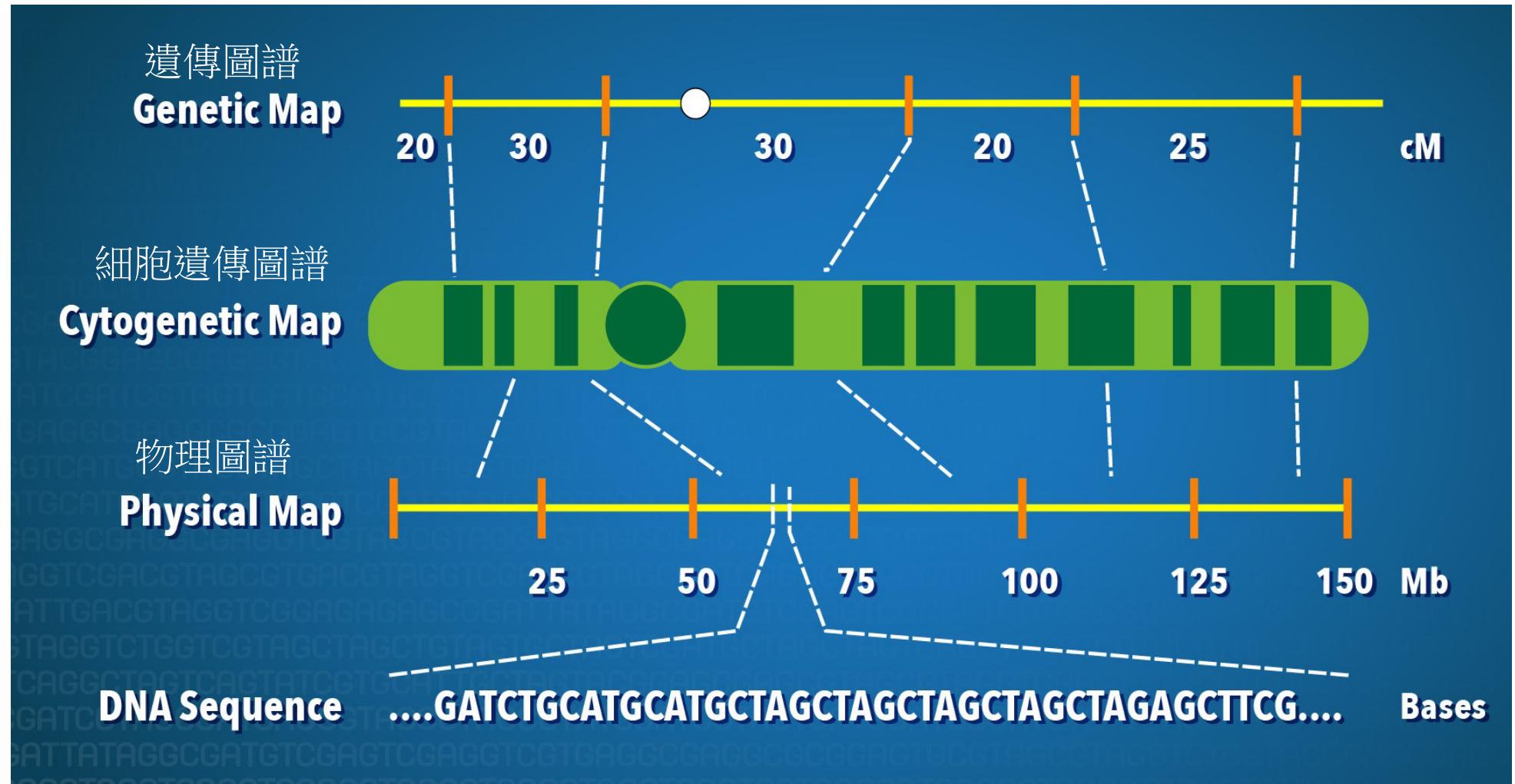


Genome map (基因組圖譜)

Cytogenetic map: 由染色體染色而來，沒有單位，以區域劃分

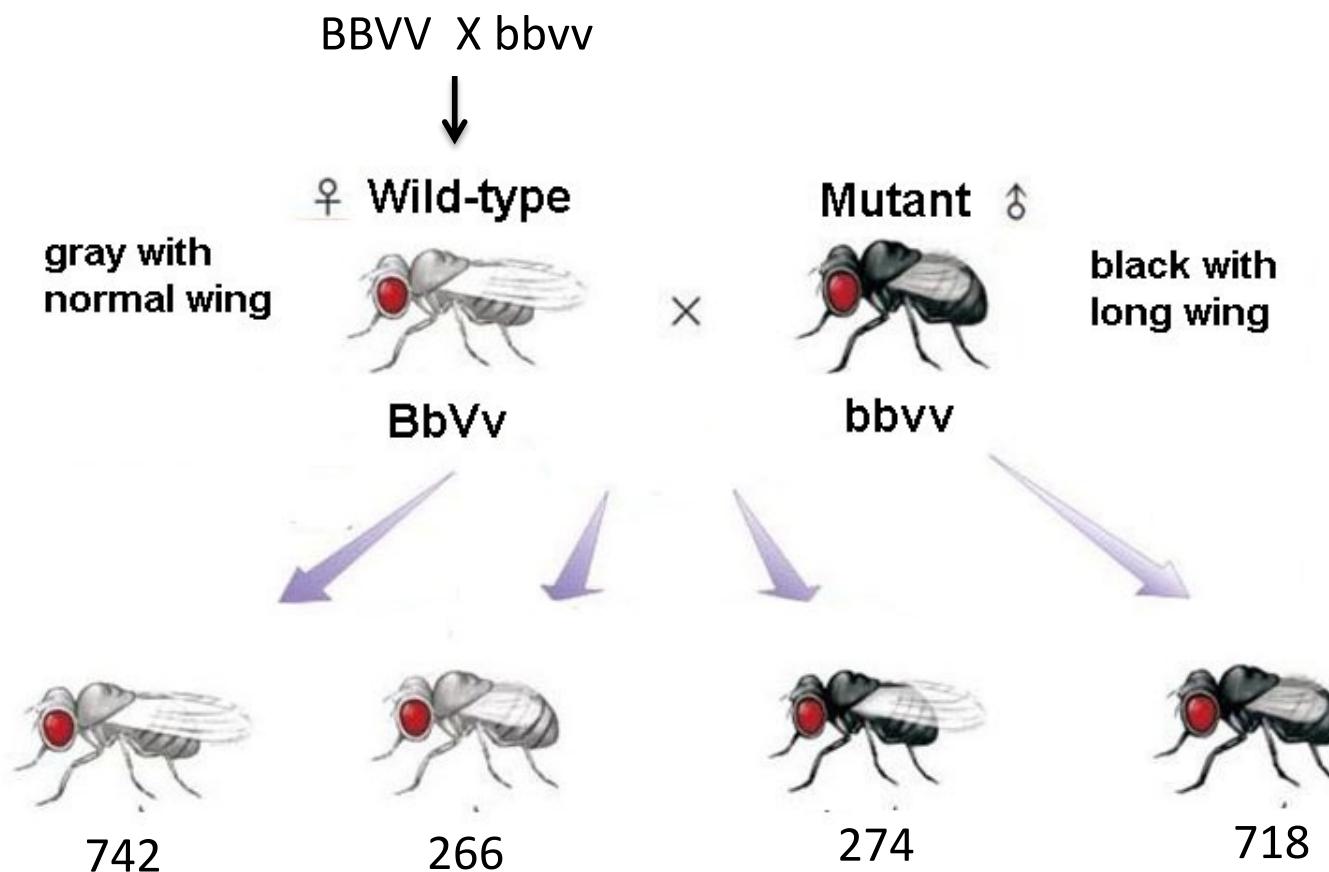
Genetic map: 由互換率計算而來，單位cM (centimorgan)

Physical map : 由序列定序而來, 單位bp



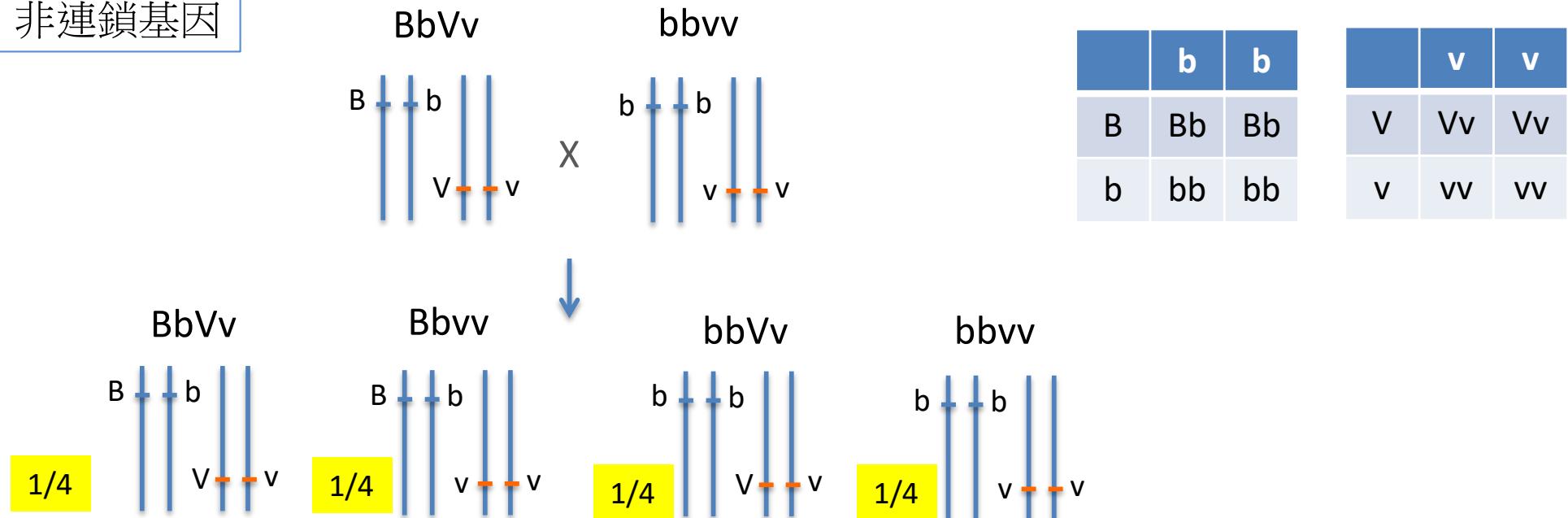
Chromosome recombination (染色體重組)

野生型果蠅(灰身長翅)和突變果蠅(黑身短翅)交配產生2000隻後代，其中742隻為灰身長翅、266隻為灰身短翅、274隻為黑身長翅，以及718隻為黑身短翅。請問控制體色(B基因)及翅膀長度(v基因)兩基因是否為連鎖？如果是的話，請問其相距多少centimorgan？

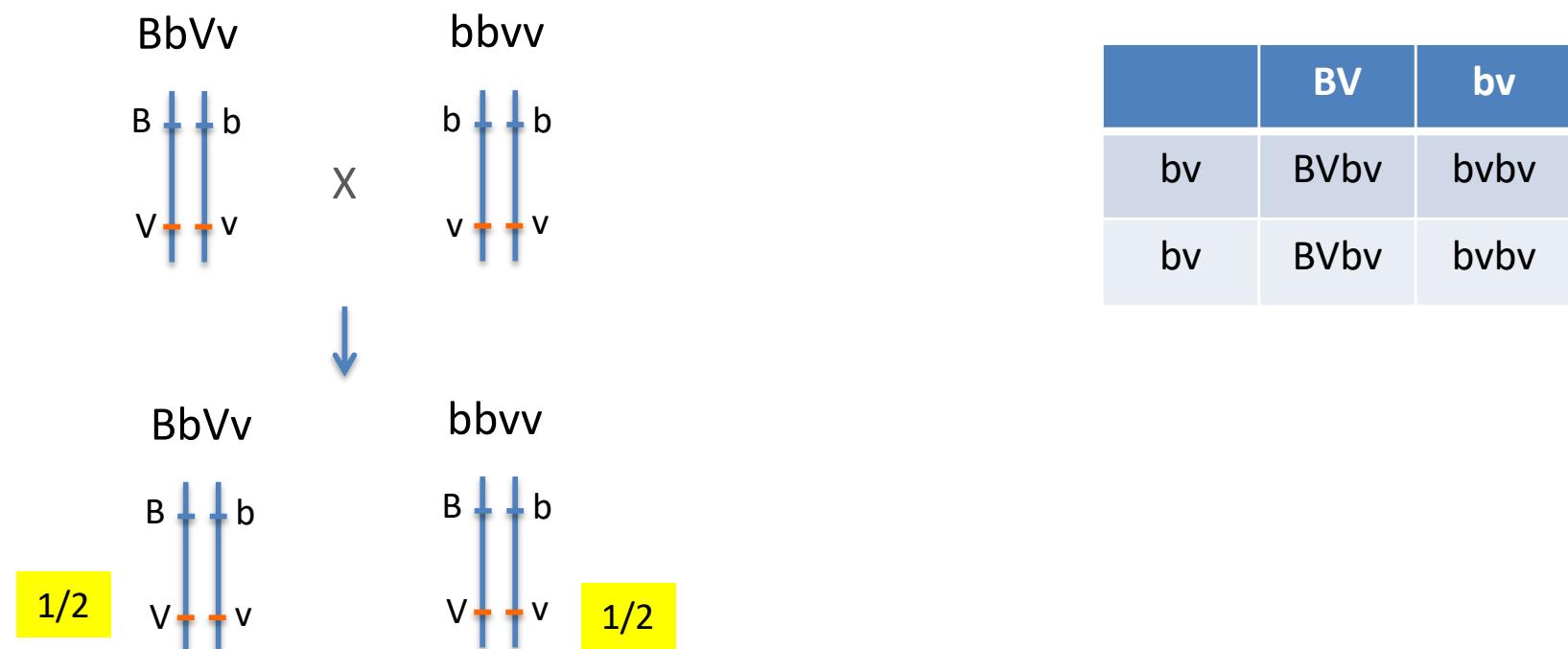


假設：母本BbVv來自於基因型為 BBVV 及 bbvv的後代

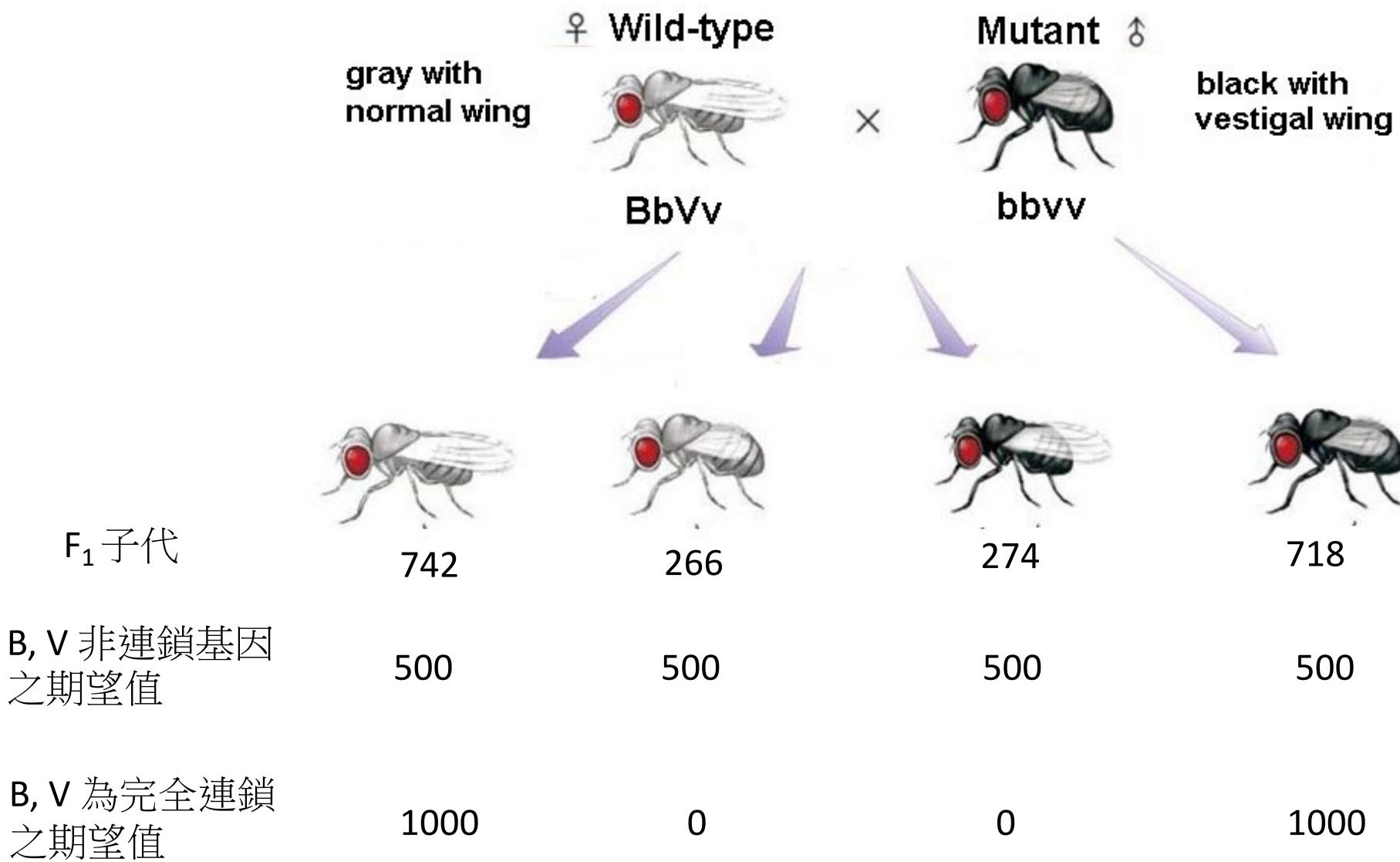
非連鎖基因



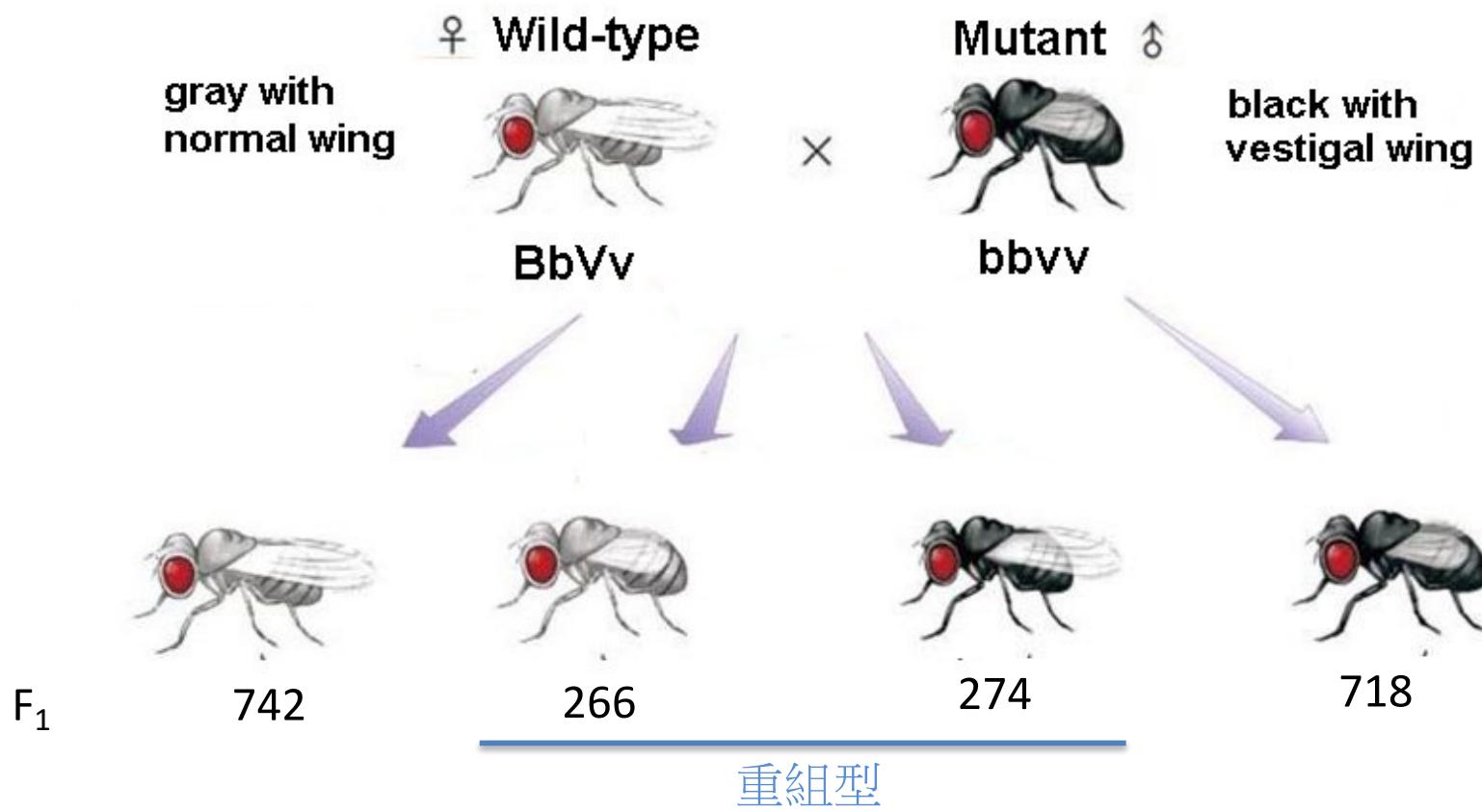
基因連鎖



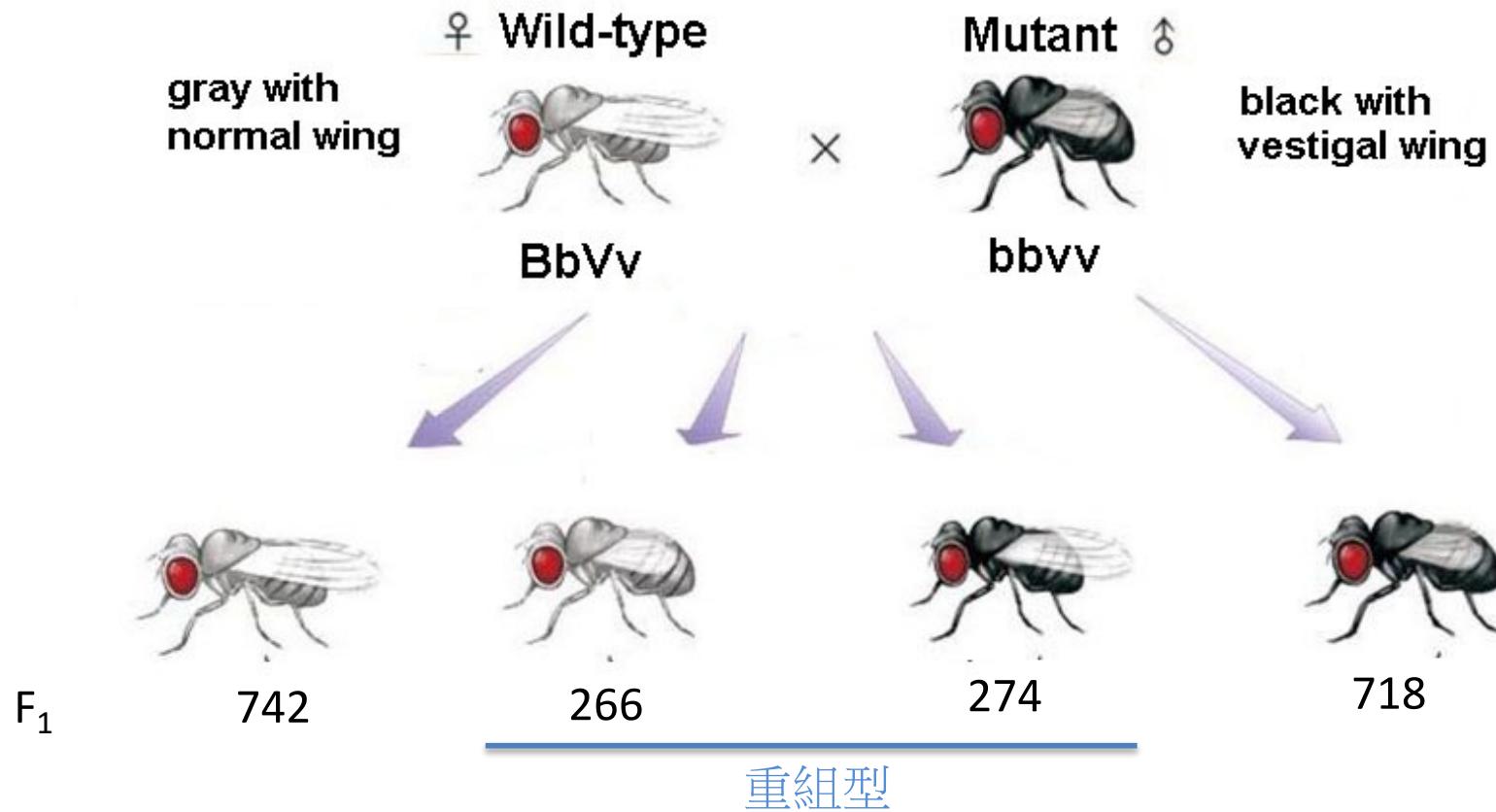
Chromosome recombination (染色體重組)



Recombination rate (重組率)



Recombination rate (重組率)



重組率 = 重組型 / 所有子代

$$= (266 + 274) / 2000 = 0.27$$

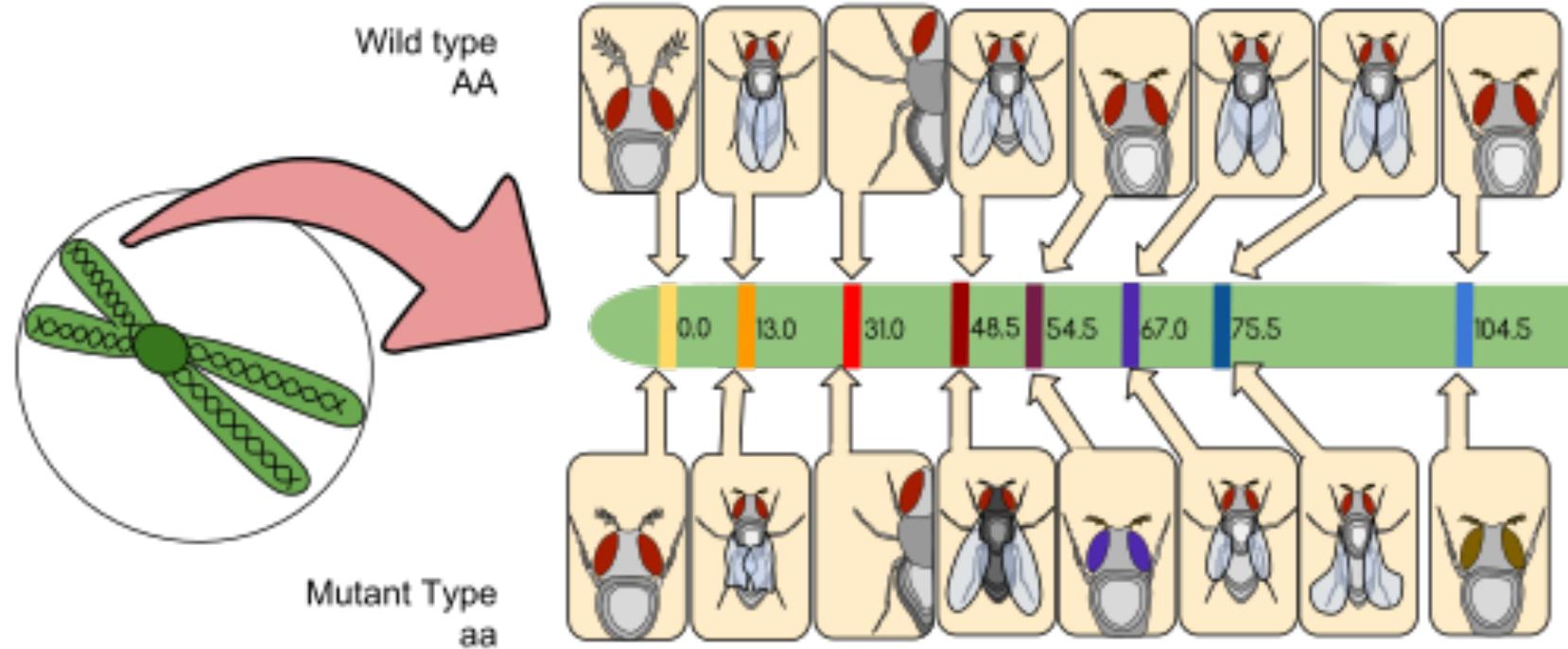
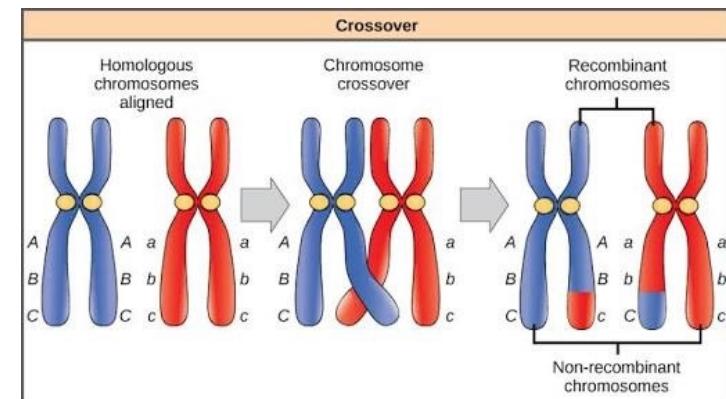
兩基因距離為 27 centimorgan (cM)

*重組率最高為 50%

*重組率的精準度和樣品數量成正相關

遺傳圖譜

- 1% 重組率 = 1 cm (centimorgan)
- 重組率 < 50% --> “連鎖”
- 重組率 = 0% --> “完全連鎖”
- 在人類細胞中 1 cM 約 1 Mb

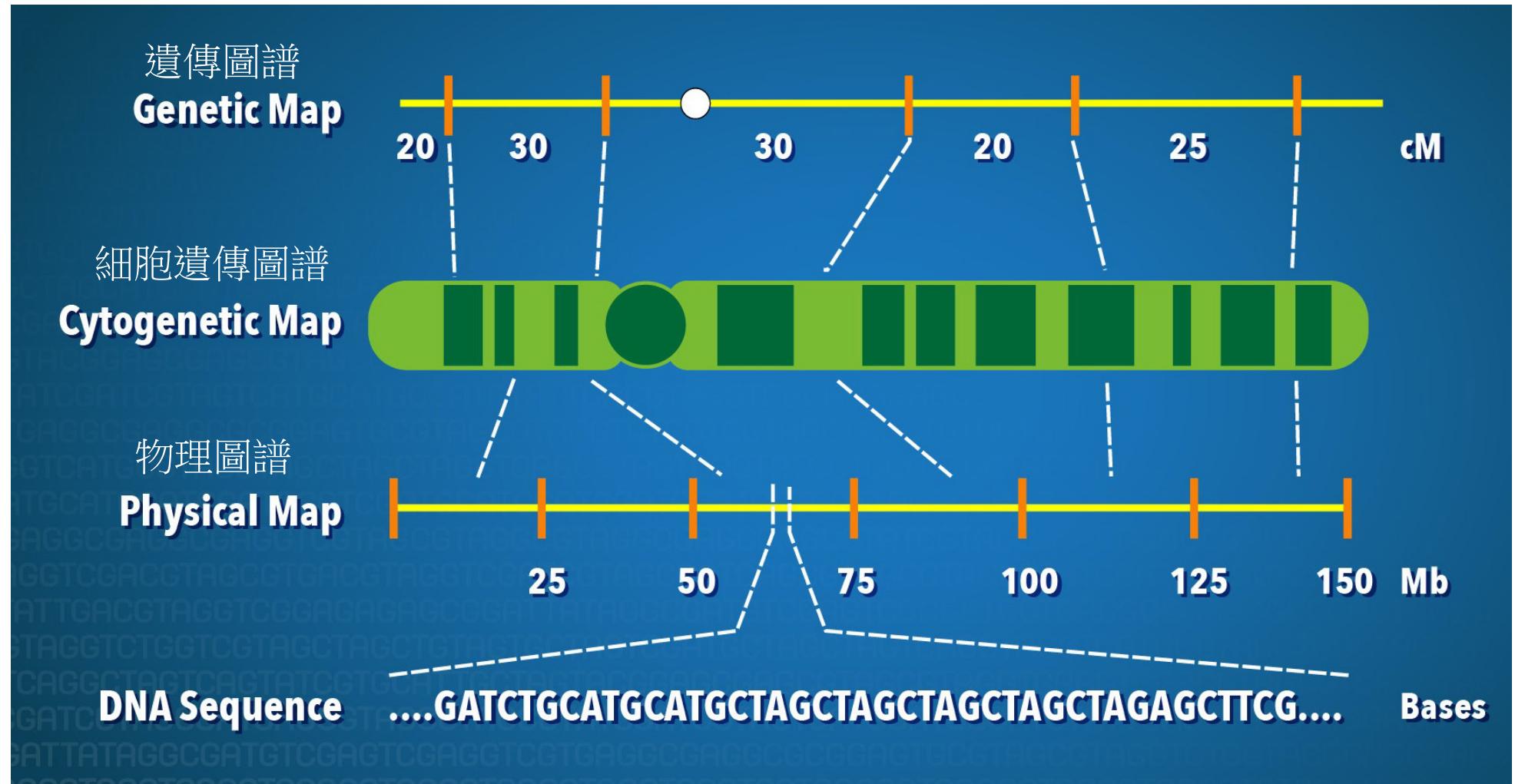


Genome map (基因組圖譜)

Cytogenetic map: 由染色體染色而來，沒有單位，以區域劃分

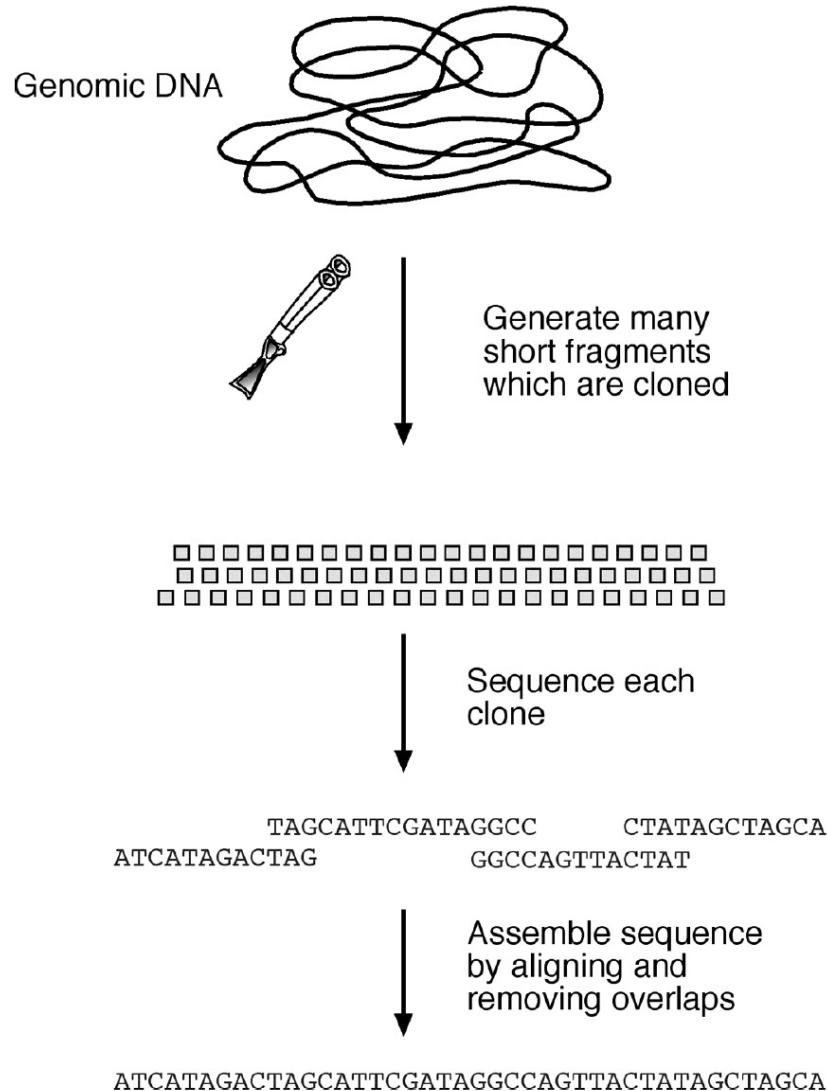
Genetic map: 由互換率計算而來，單位cM (centimorgan)

Physical map : 由序列定序而來, 單位bp

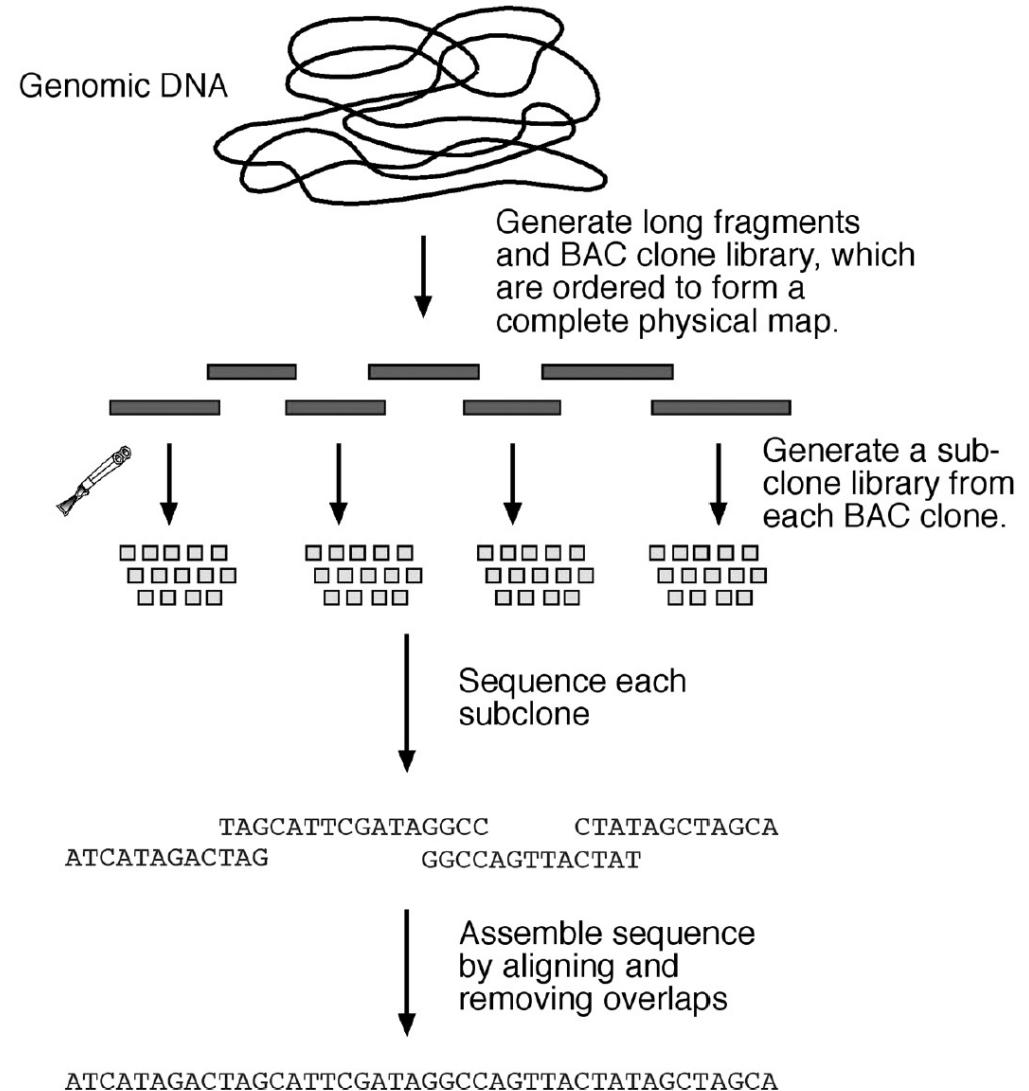


全基因組定序策略

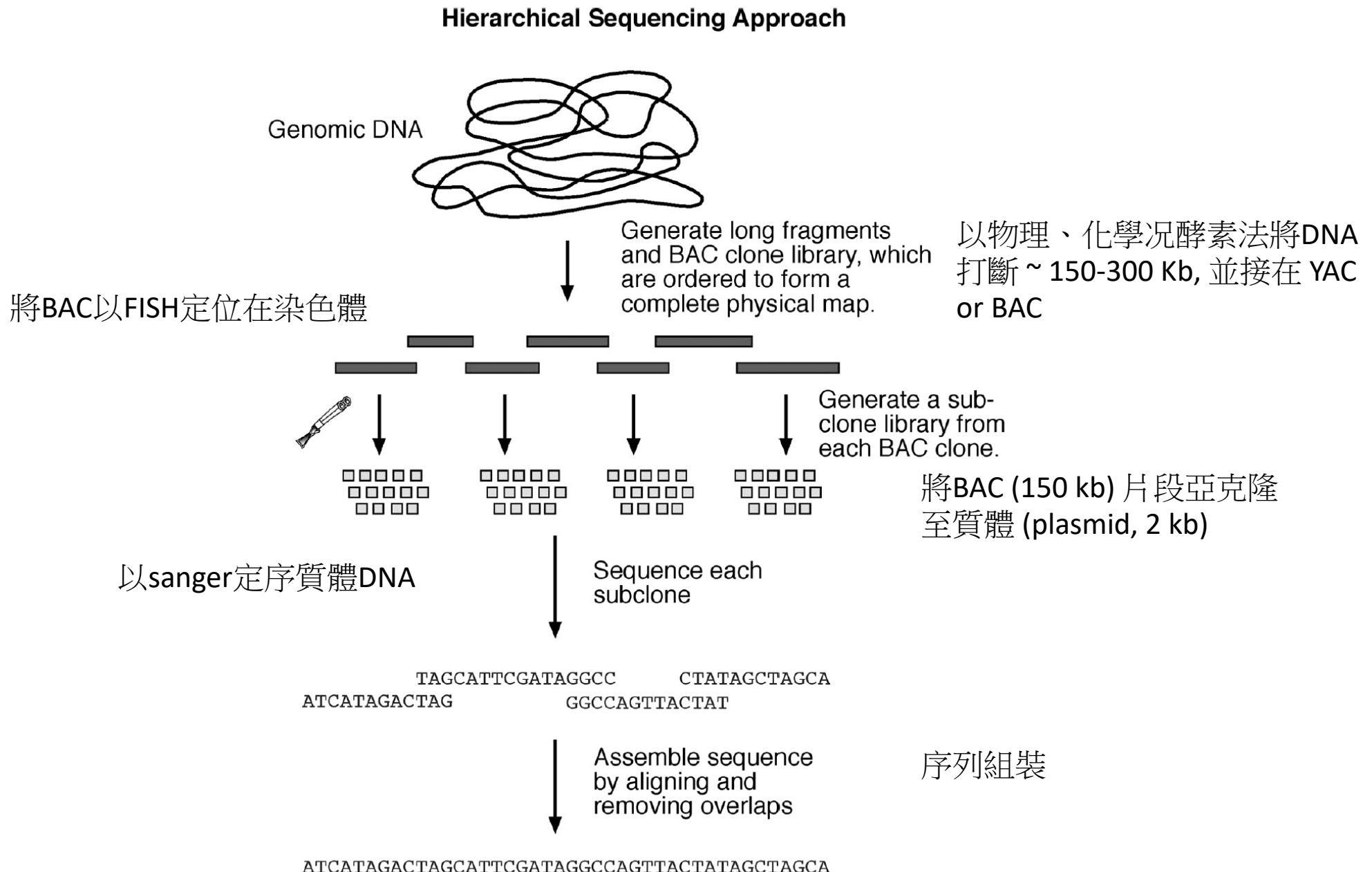
Shotgun sequencing approach (霰彈槍定序法)



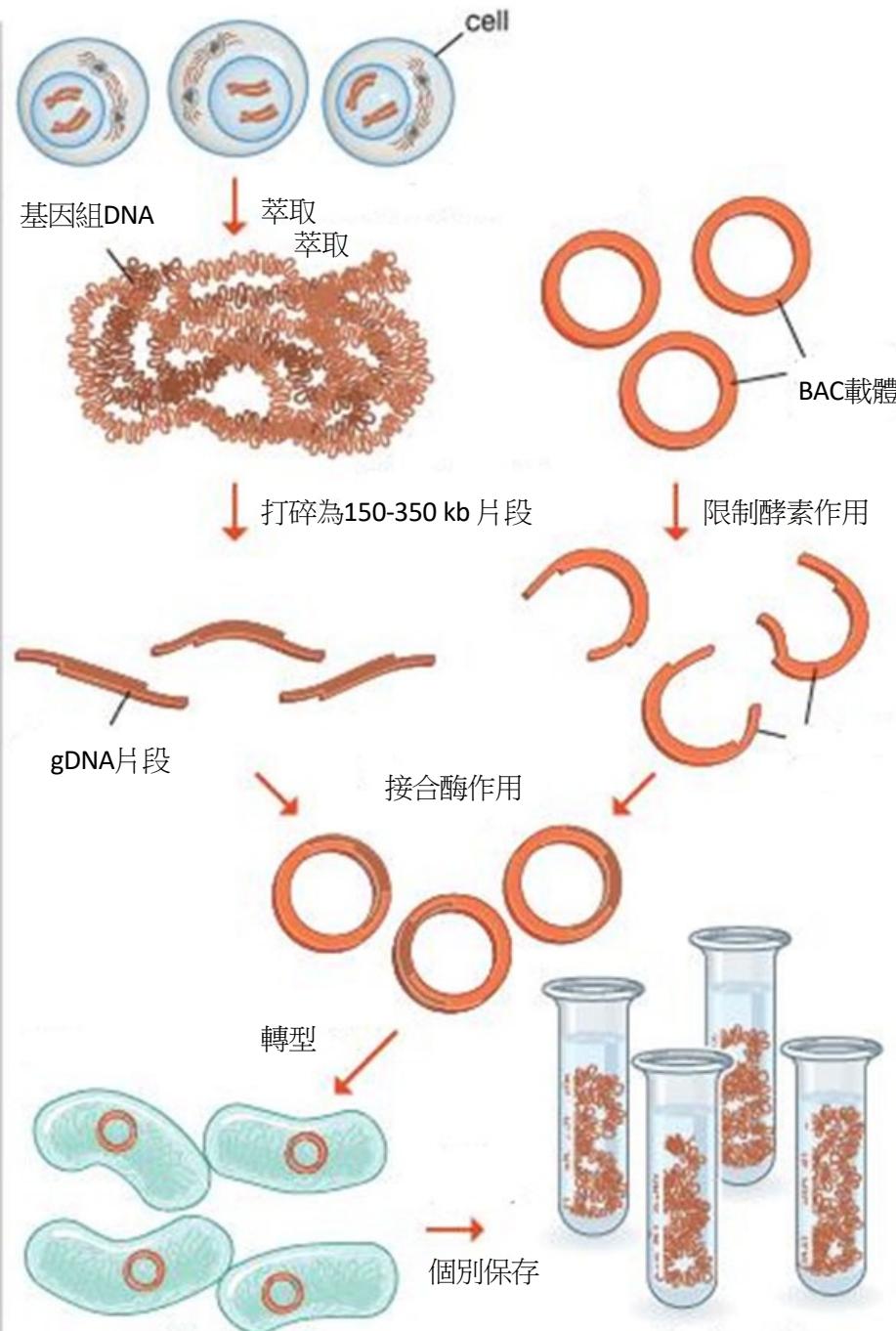
Hierarchical sequencing approach (階層式定序法)



Hierarchical sequencing approach (階層式定序法)



基因組庫建構 (Genomic Library Construction)



如果一個BAC可承載150kb,對於人類基因組3.2 Gb而言，9,142個BAC為一倍的覆蓋率，而若要達到90%以上的覆蓋率，至少要6倍的BAC數目(54,852個)

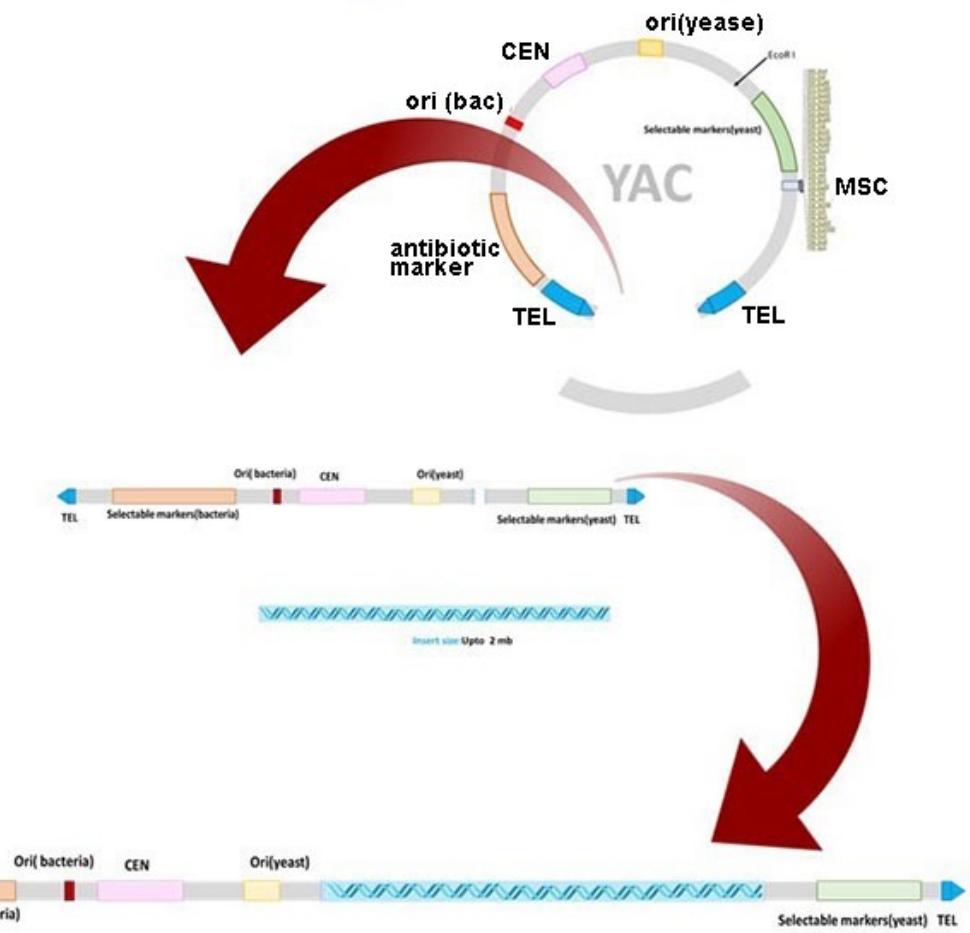
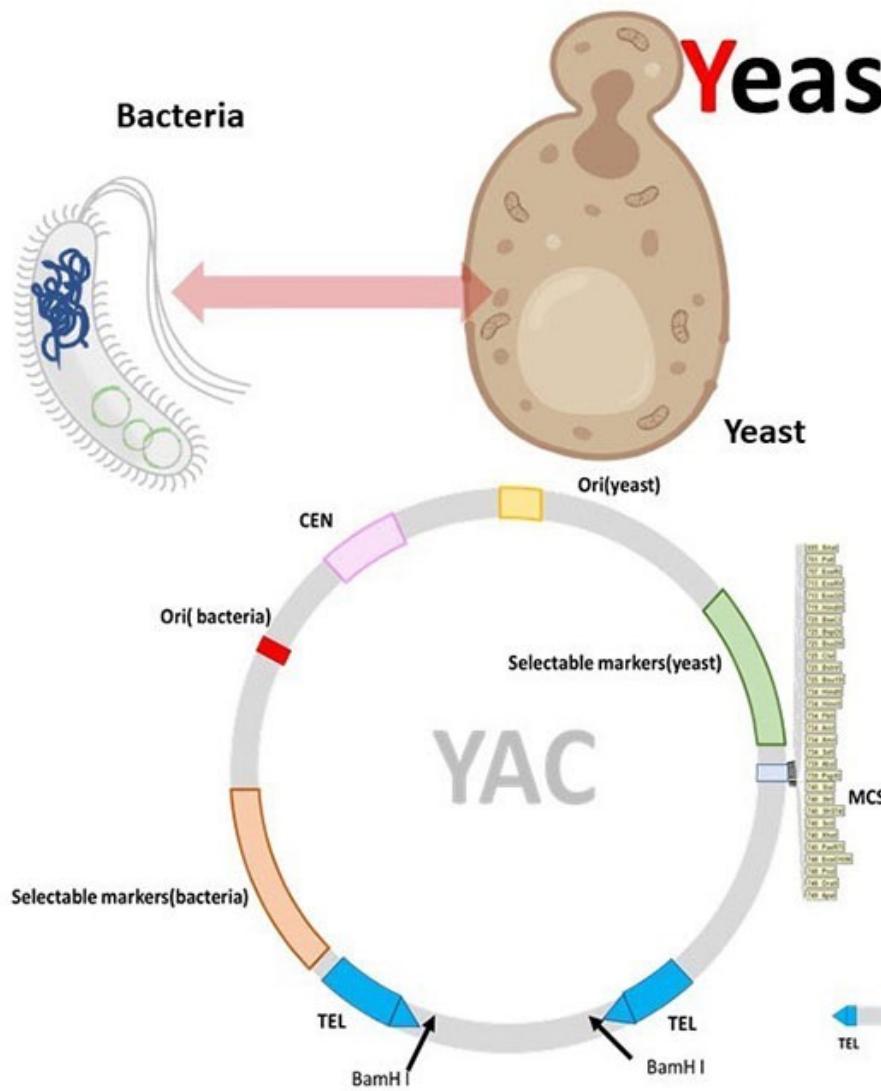
DNA 載體

載體	承載量	宿主細胞
人類人造染體 (HAC)	6000 - 10000 Kb	human cell
酵母人造染色體 (YAC)	100 - 3000 kb	Yeast
細菌人造染色體(BAC)	150 ~ 350 kb	E. Coli
噬菌體載體 (PAC)	100- 300 kb	E. Coli
黏質體 (Cosmid, 噬菌體載體/質體之複合體)	35-45 kb	E. Coli
質體 (plasmid)	<= 15kb	E. Coli

如果以一倍的覆蓋率計算，人類基因組 (3,200,000 kb) 需要
320 HAC
1,066 YAC
9,142 BAC
10,666 PAC
71,111 Cosmid
213,333 plasmid

* 通常一個基因組庫的要求為6倍覆蓋率以上

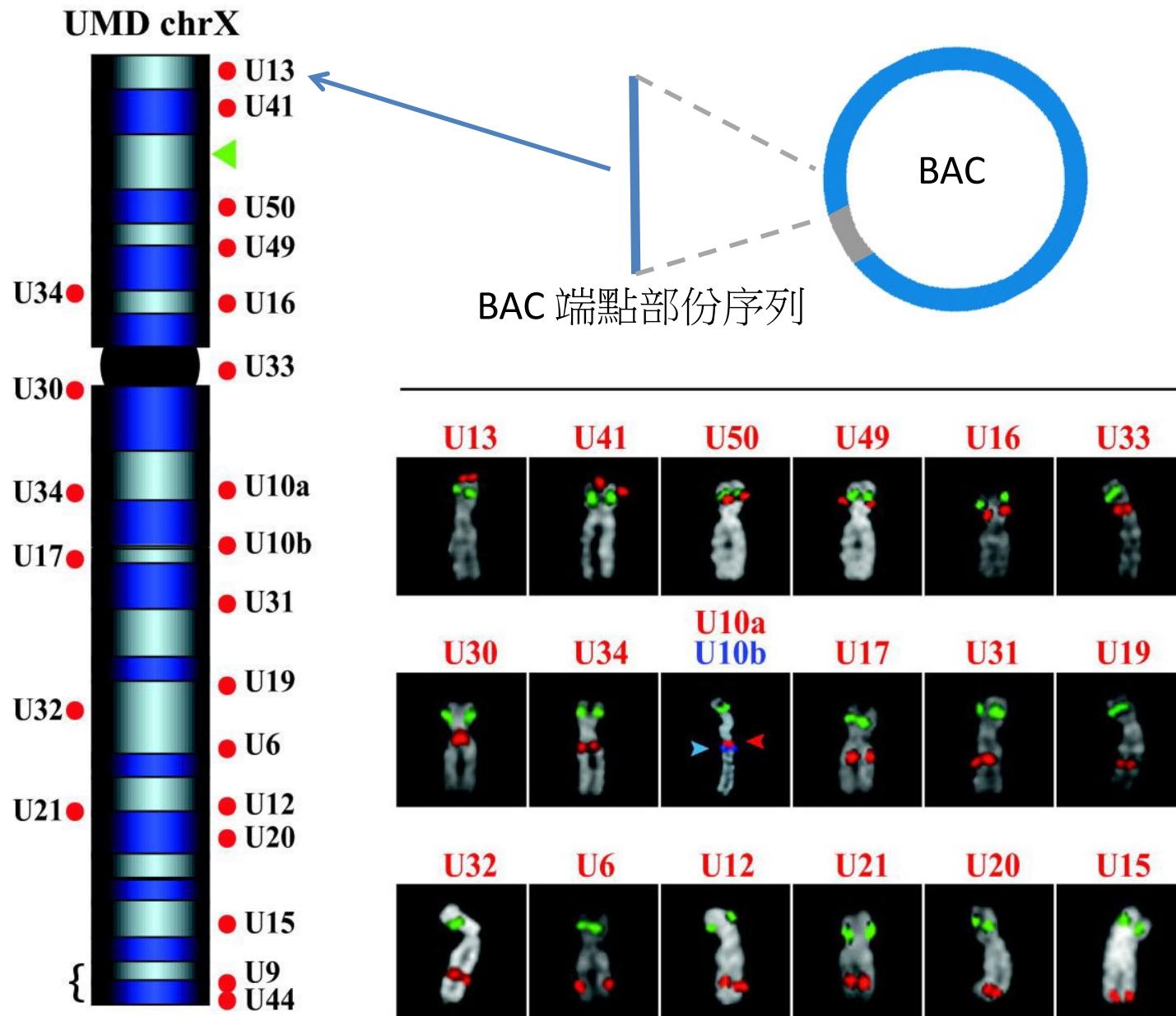
Yeast Artificial Chromosome



Capacity: 100 -3,000 kb

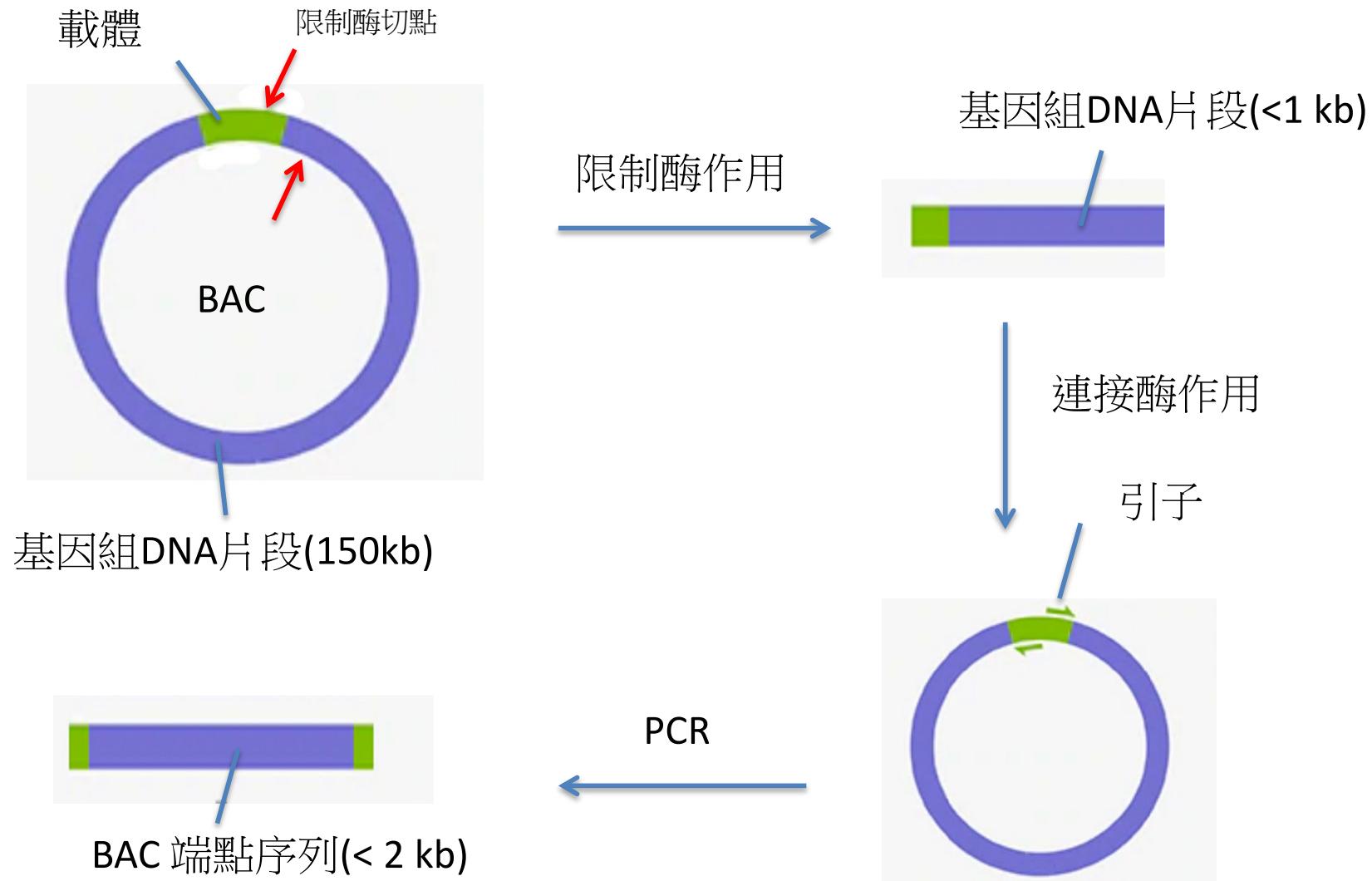
TEL: telomere
CEN: centromere
MSC: multiple cloning site
Ori: origin of replication

以螢光原位雜合法(FISH)將BAC定位在染色體上定位



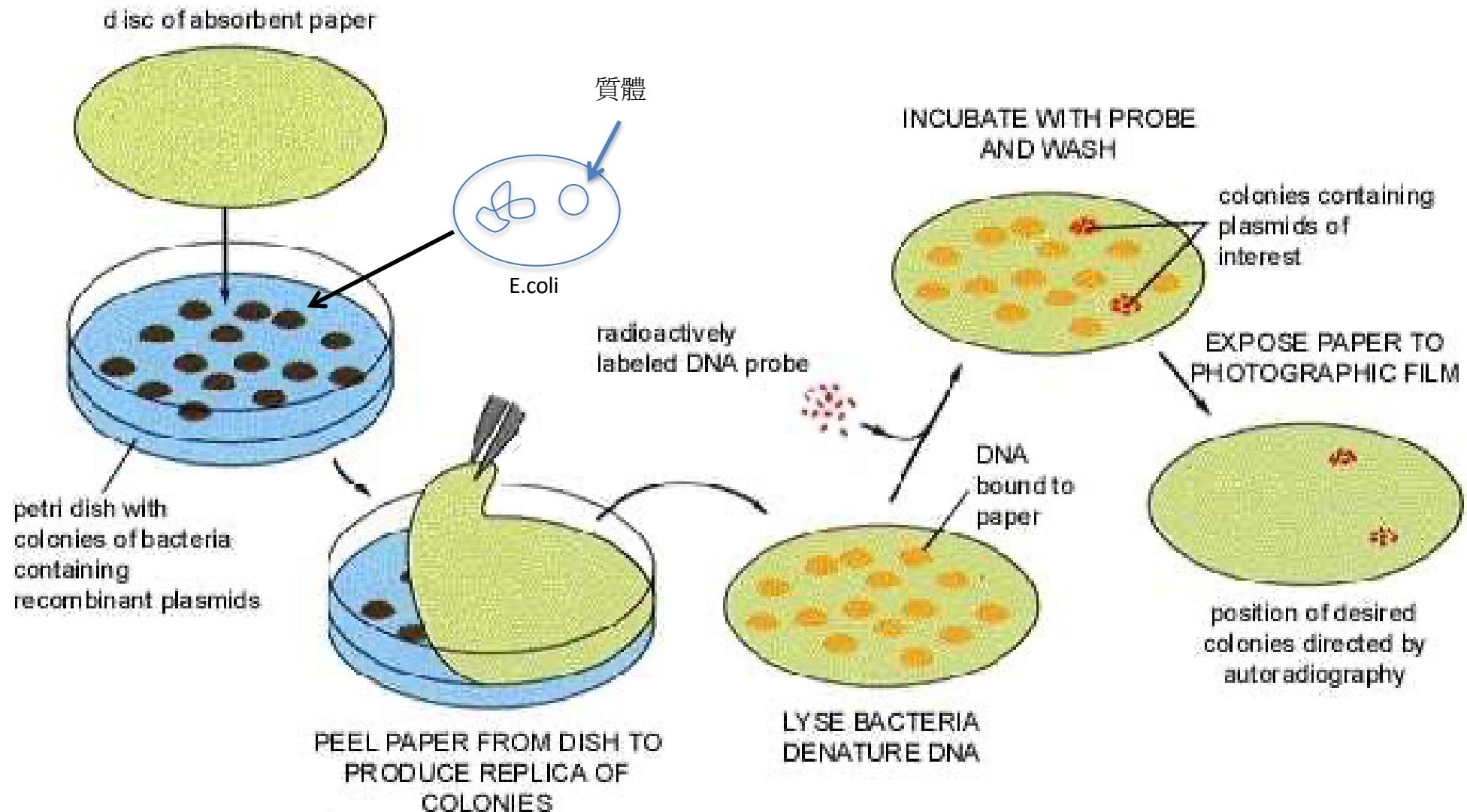
Inverse PCR (反向PCR)

- 最好聚合酶最長可做到30kb，但很貴，成功率也很低。大部份聚合酶僅能做到~ 2 kb
- 基因組DNA片段為未知，無法設計引子。引子只能設計在載體序列
- 目的為擴增BAC端序列，以進行FISH或BAC library screening

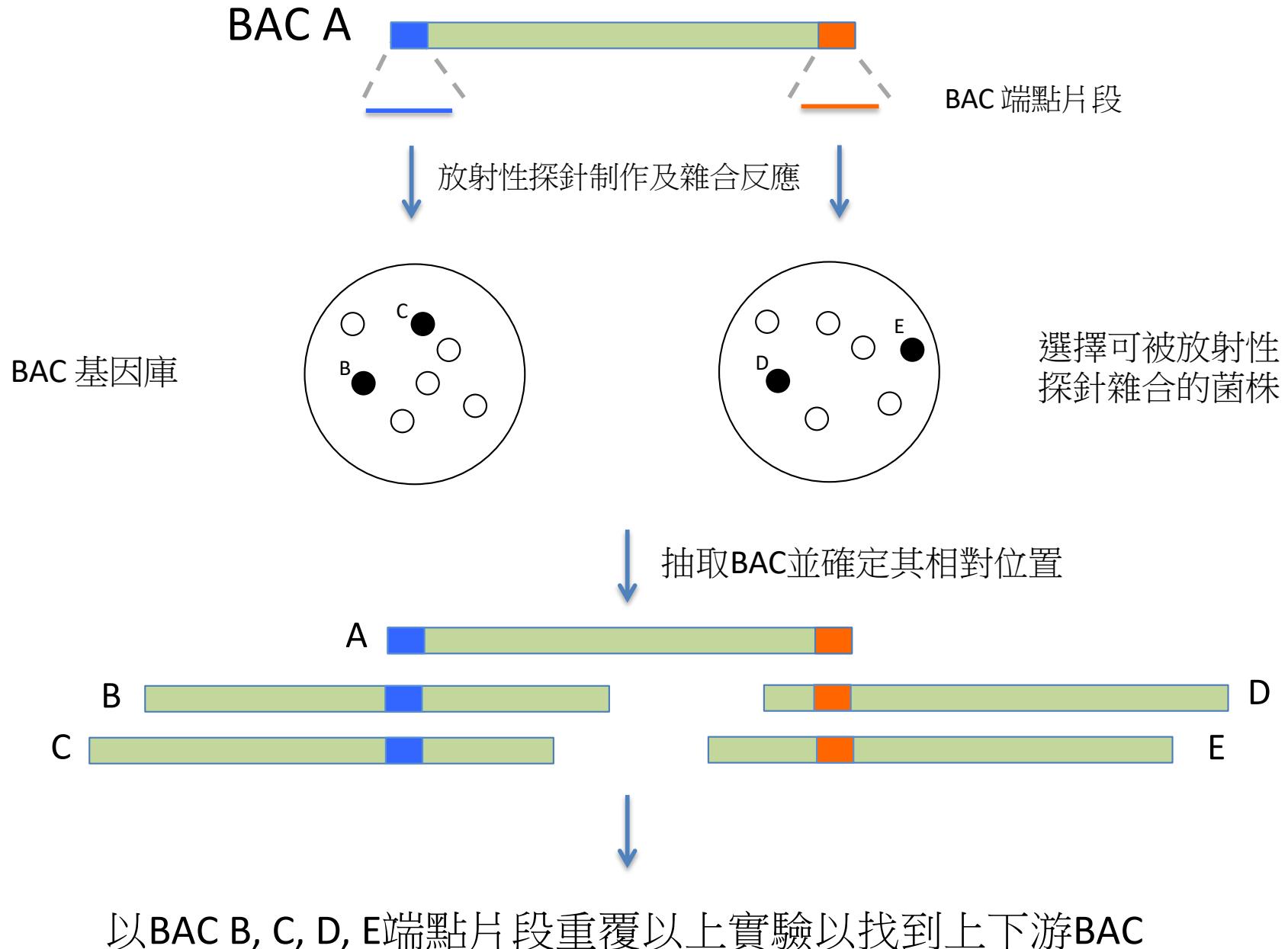


Colony hybridization(菌落雜交)

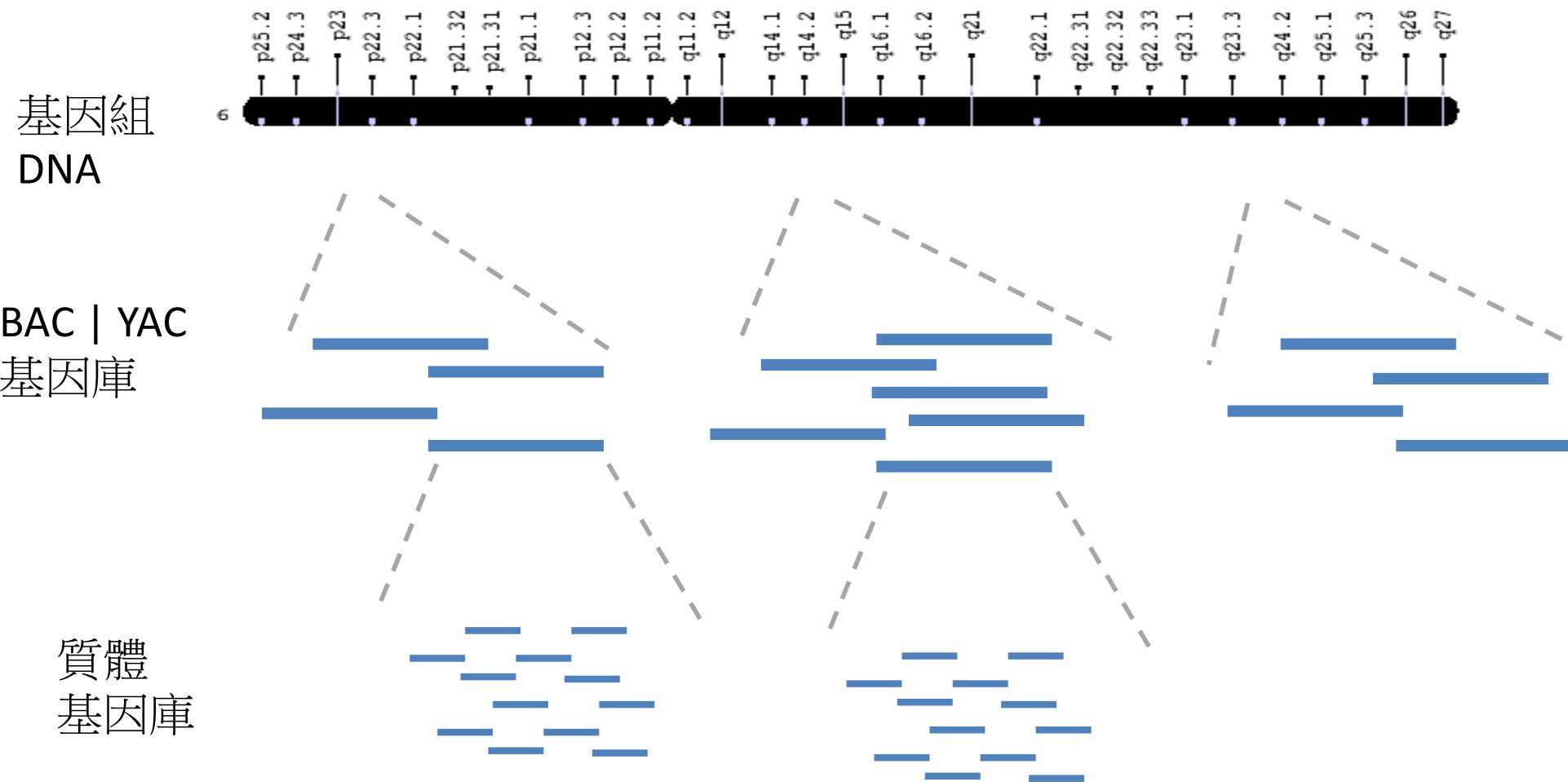
- 主要原理為單股DNA可和放射性探針(probe,單股DNA)結合
- 目的為尋找含有特定序列的菌落
- 培養菌落 → 拓印至硝化纖維膜 → 鹼破壞打破細胞並使DNA變性 → 放射性探針雜合 -> 訊號偵測



BAC基因庫篩選

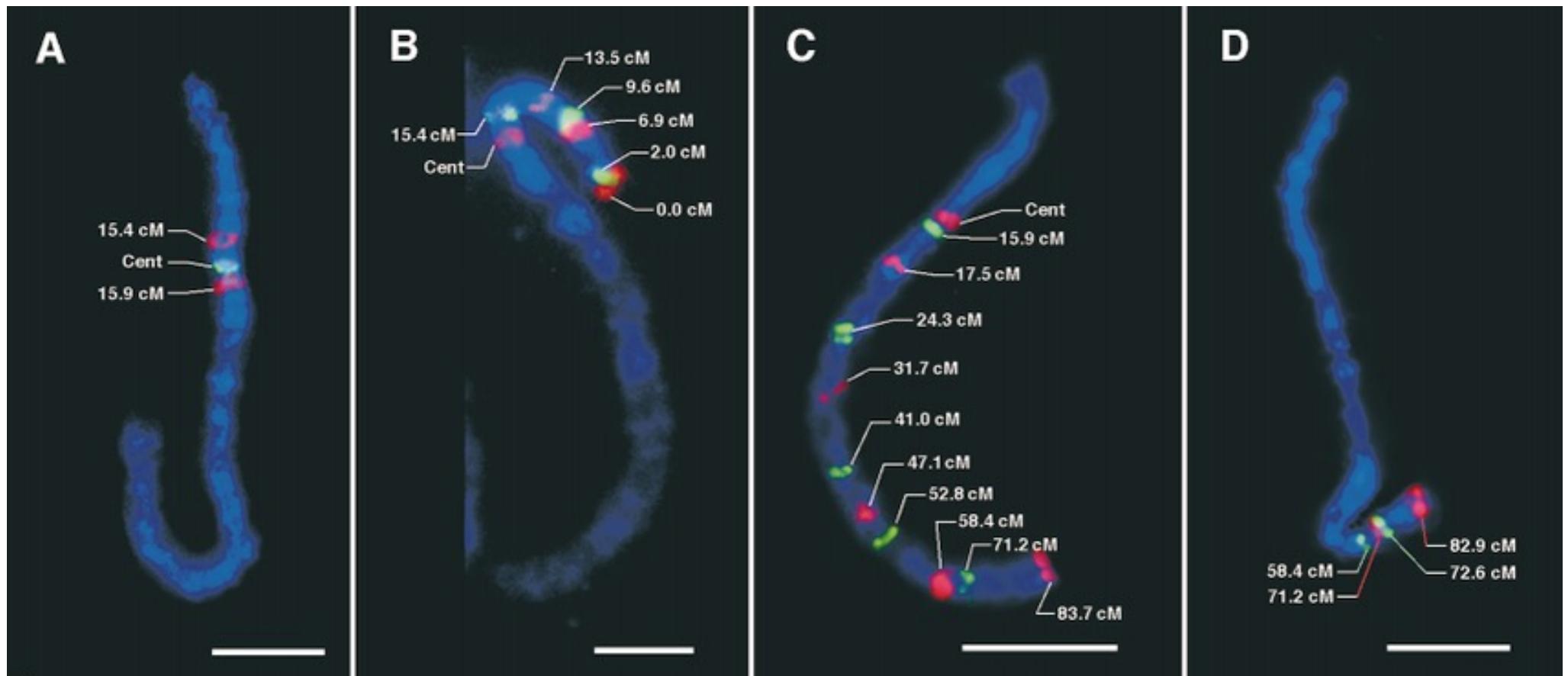


階層式定序



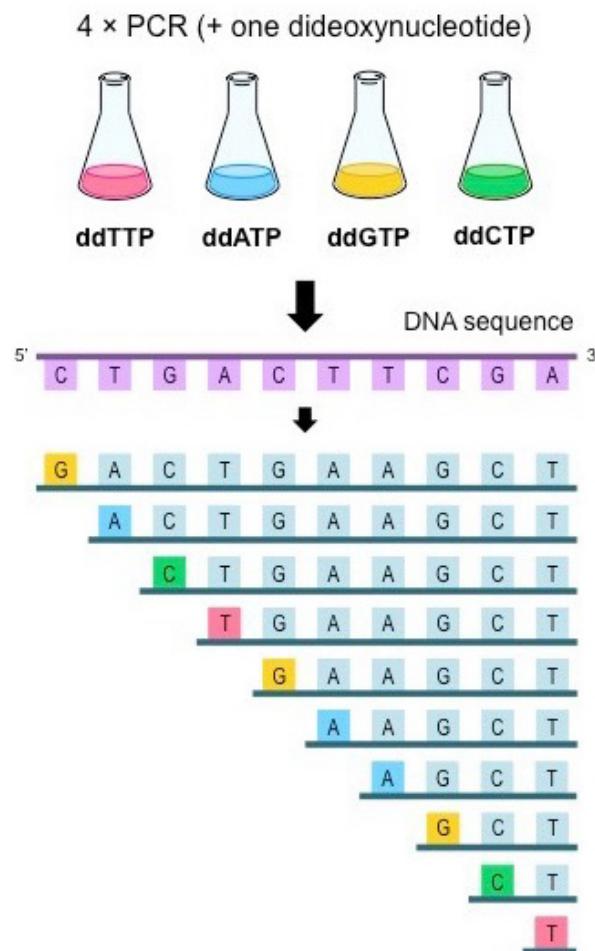
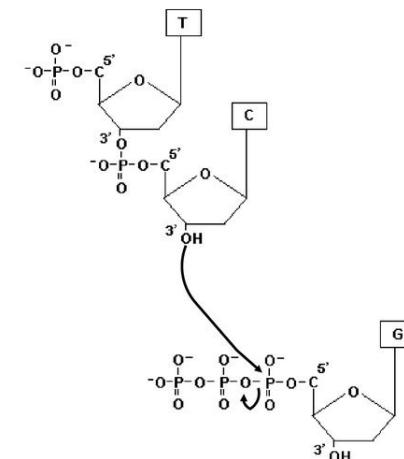
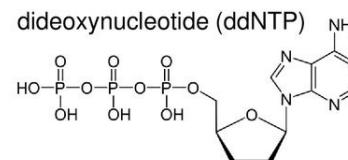
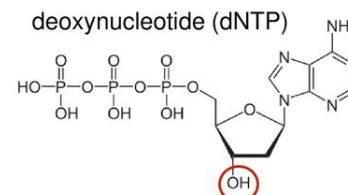
螢光原位雜合FISH (Fluorescence in situ hybridization)

- 主要原理為單股DNA可和螢光標定探針(probe,單股DNA)結合
- 目的為確認目標序列在染色體上的位置
- 細胞固定於玻片 → 以formamide將染色體變性 → 融光標定探針雜合 → 融光顯微鏡觀察

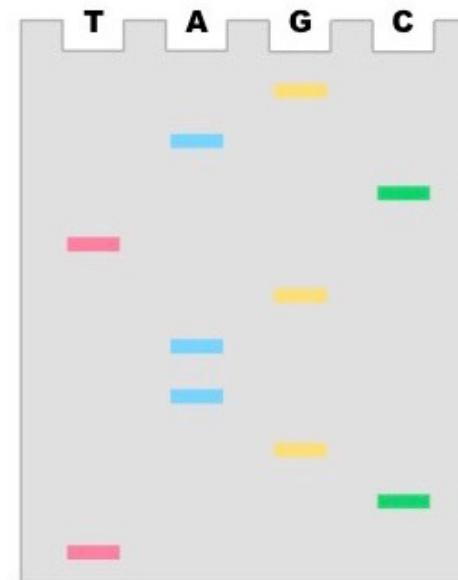
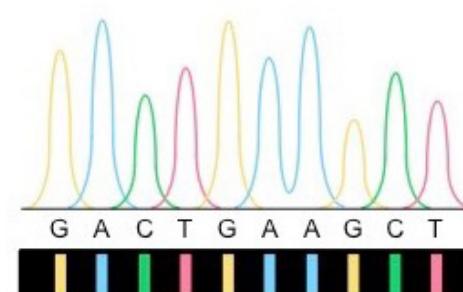


Sanger sequencing

- 藉由螢光標定ddNTP(雙脫氧核苷酸,五碳糖缺乏3端OH基)使PCR反應停止
- 定序長度 500-1200 bp



Use a sequencing machine
Separate with a gel

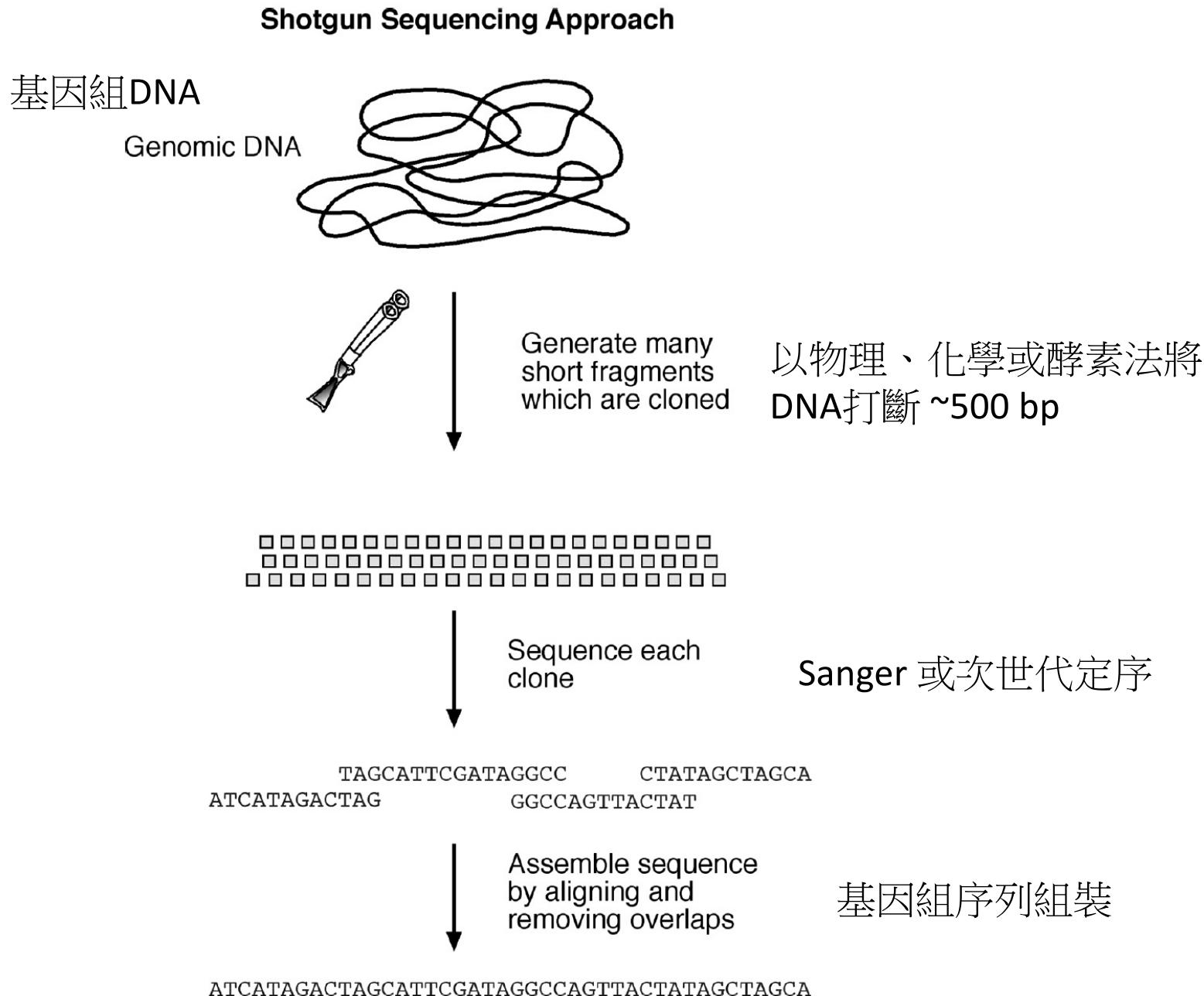


定序序列
5' ————— 3'
引子
5' ← 3'

PCR

跑膠，螢光分析

霰彈槍定序法 (Shotgun sequencing approach)



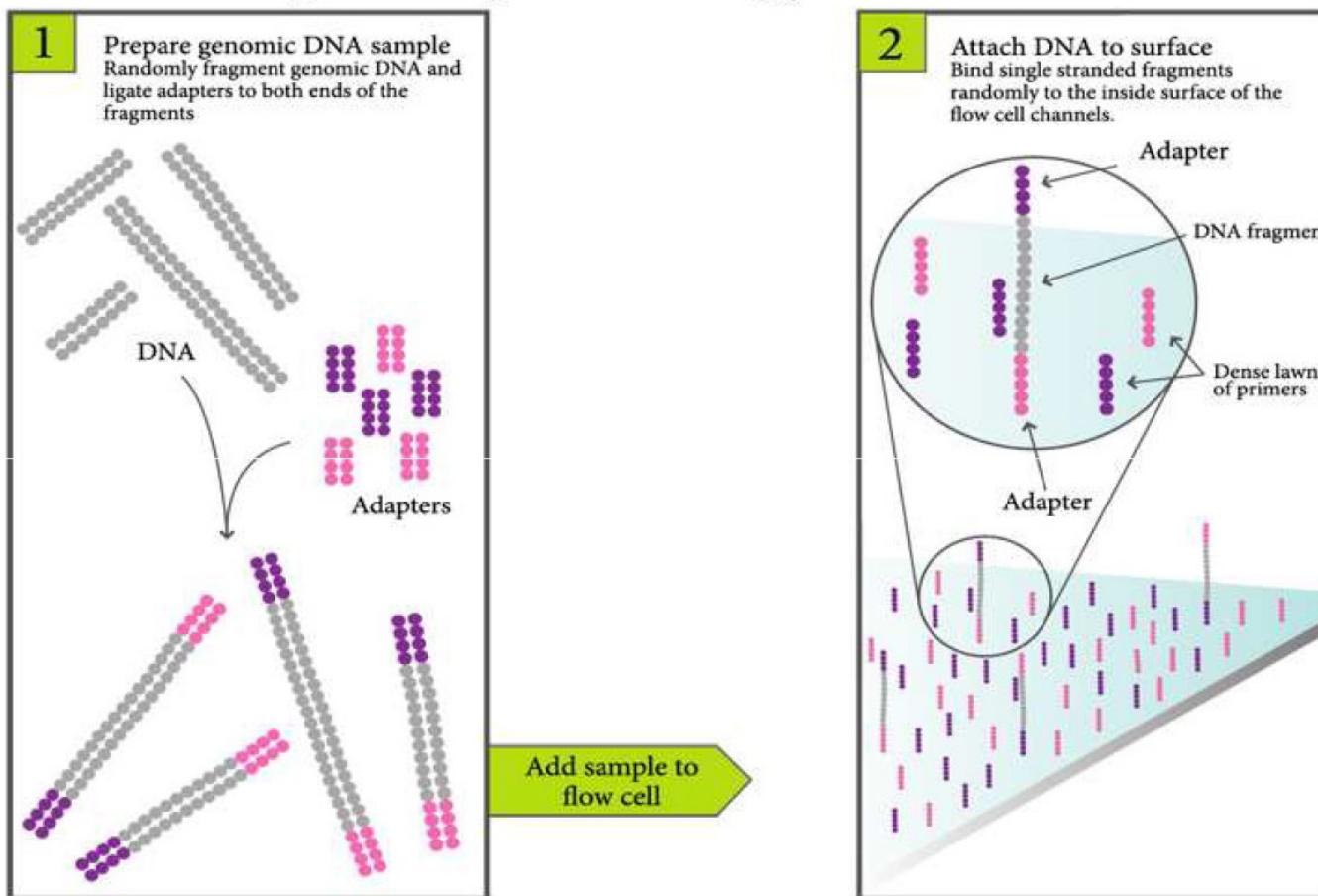
DNA定序 (DNA sequencing)

- 傳統定序法 Sanger sequencing
 - 定次定序長度 ~ 500bp - 1 kb
 - 每次上機 96 samples
 - 1977年發明
- Next generation sequencing (次世代定序, NGS)
 - 又稱大規模並行測序 (Massive parallel sequencing)
 - ~ 100 - 500 bp /read (依定序技術不同)
 - 每次上機可獲得 > 100 Gb 資料量
 - 不同公司技術不同，目前以 illumina 開發的 SBS 技術為主流 (2006年發明)
 - 可用於基因體、轉錄體、基因甲基化及環境微生物定序

SBS (sequencing by synthesis)

- 將Genomic DNA打斷為約 ~500 bp, 並接上adapter
- 將DNA變性並固定在含有引子的定序盤
- 引子的序列為根據adapter設計

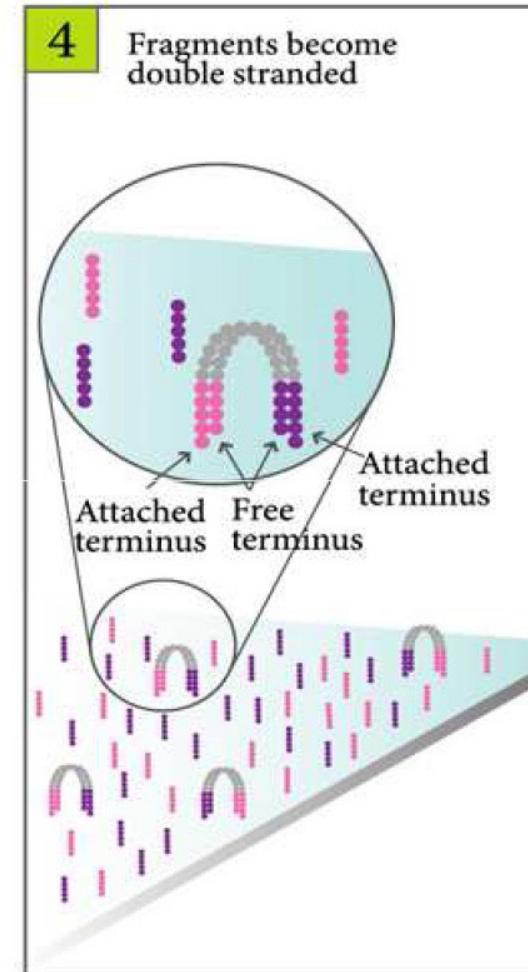
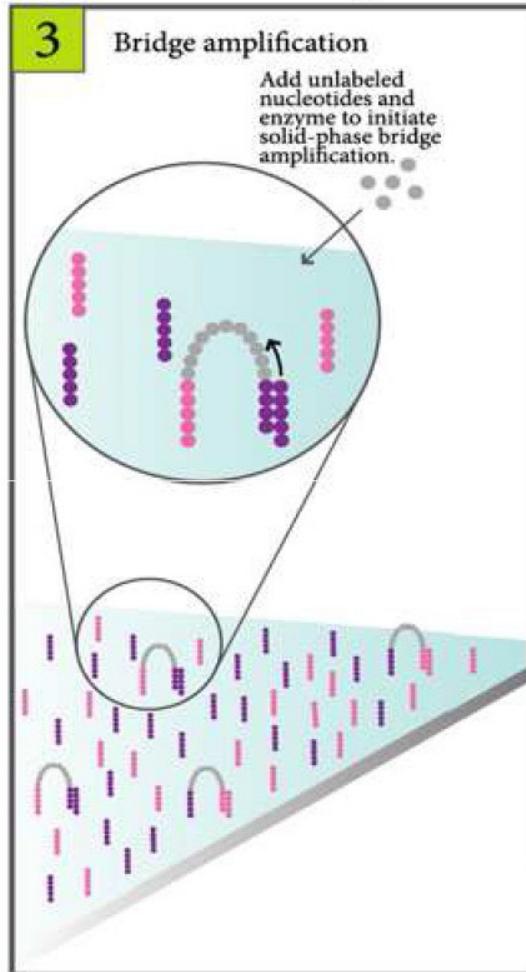
Sequencing Technology Overview



SBS

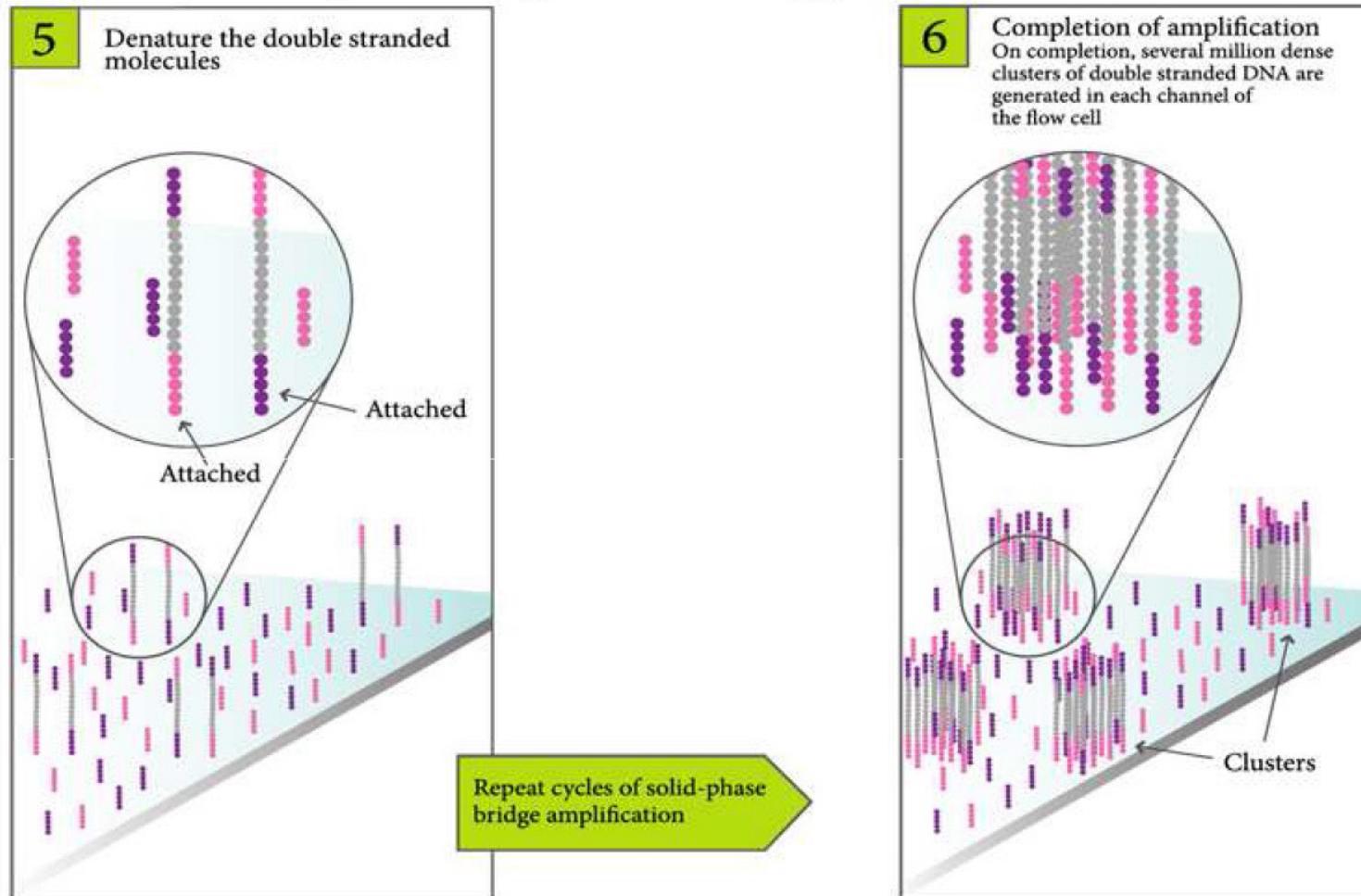
- 利用引子進行PCR，將單股的標的序列複制為雙股

Sequencing Technology Overview



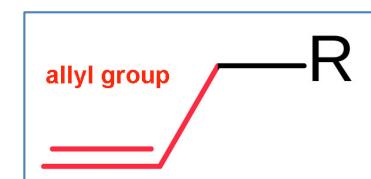
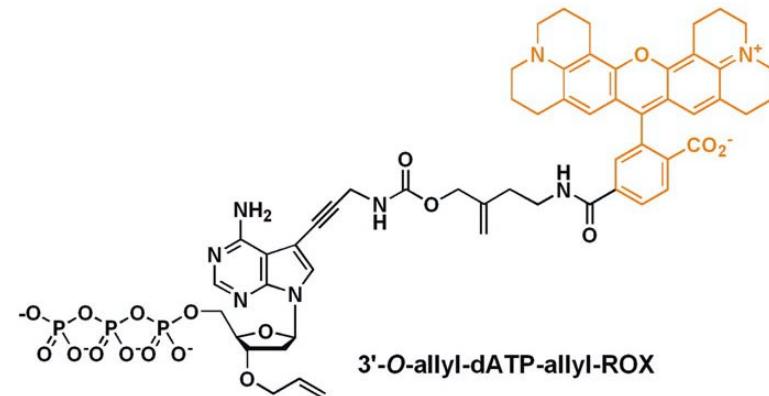
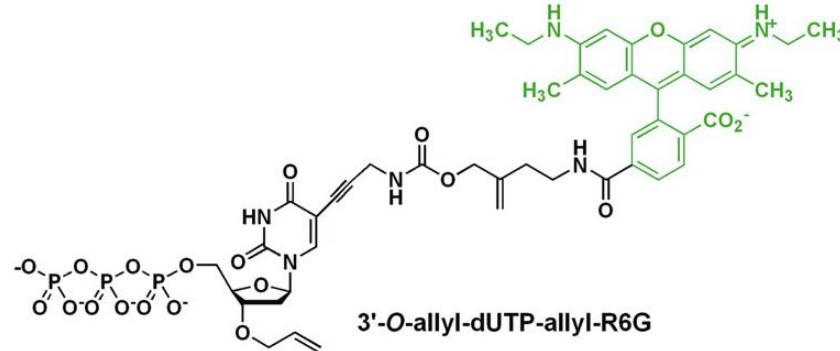
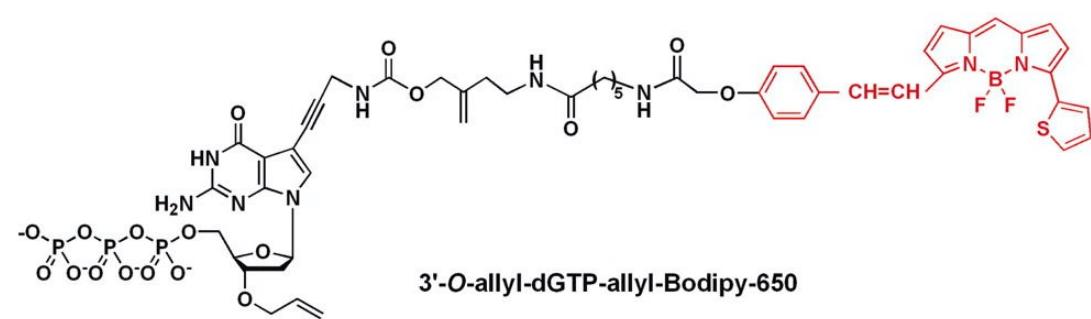
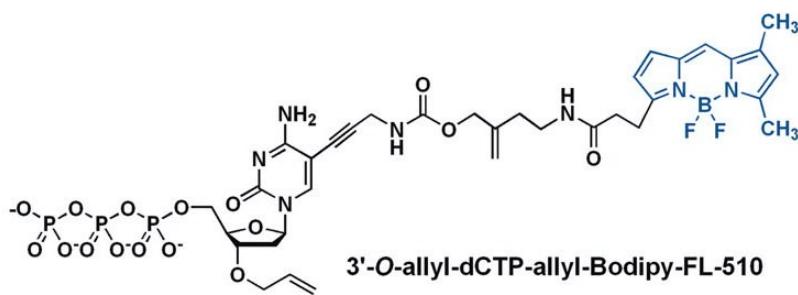
- 將複制的雙股變性成單股，並重覆PCR步驟，序列數目將以2次方成長
- 複制的目的為擴大定序訊號

Sequencing Technology Overview

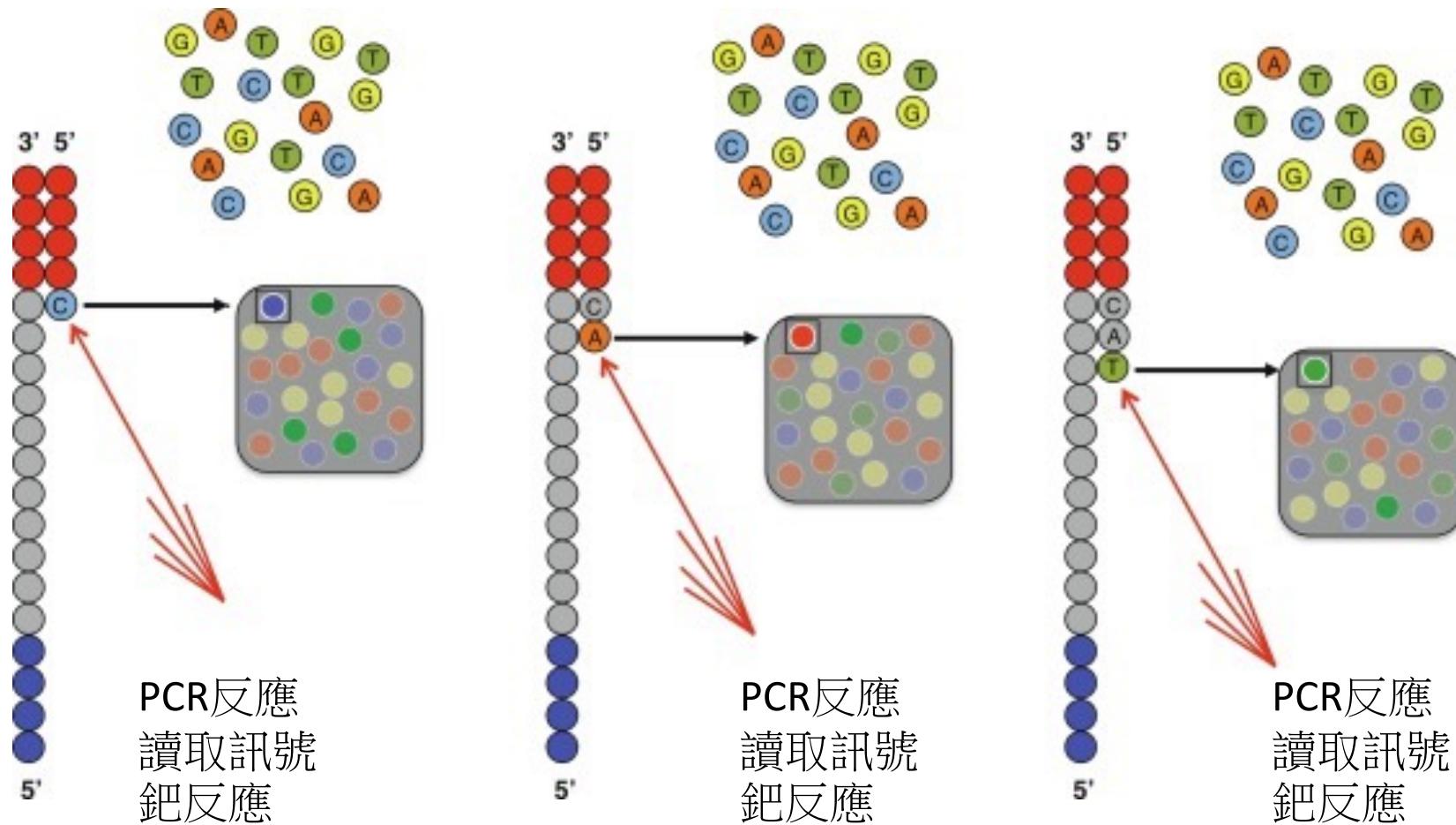


SBS – 反應用特殊核苷酸

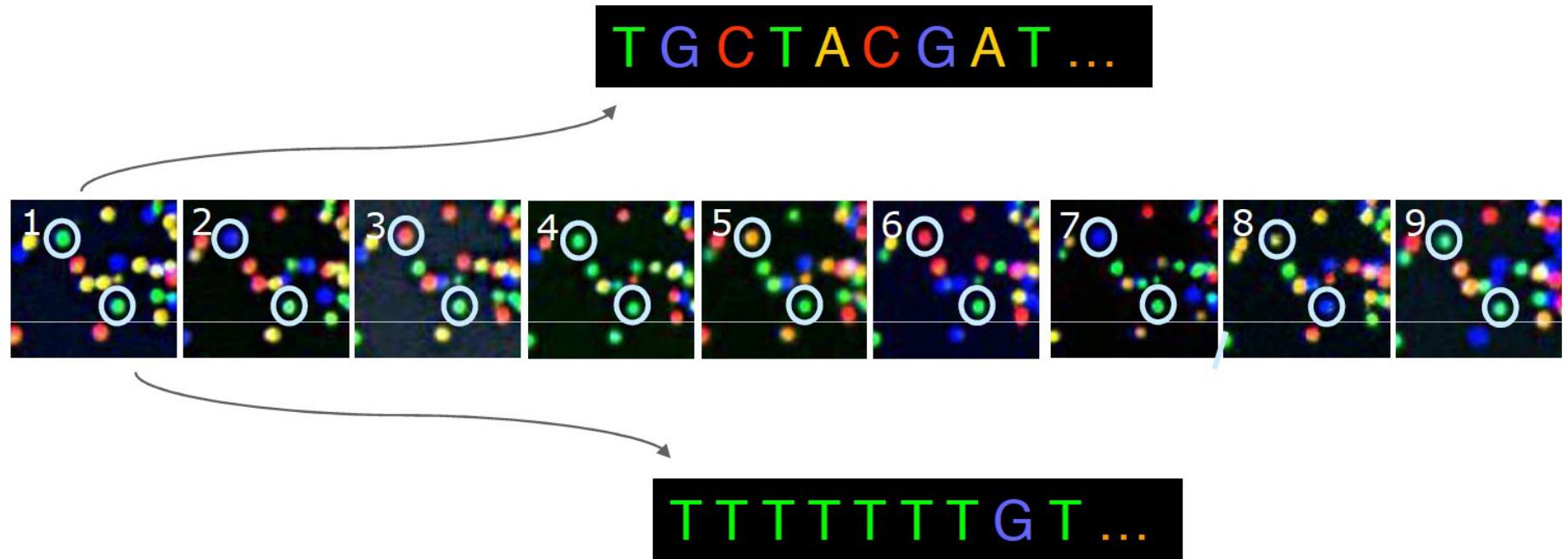
- SBS所使用的核苷酸為特殊鹼基，上頭帶有發色基，可在激光後發出不同顏色
- 核苷酸五碳糖3端上帶有烯丙基(allyl group), 可使PCR反應無法進行
- 發色基及五碳糖3端的烯丙基可以鉀(pd)反應將其移除，使PCR反應繼續



SBS – PCR 反應及訊號讀取



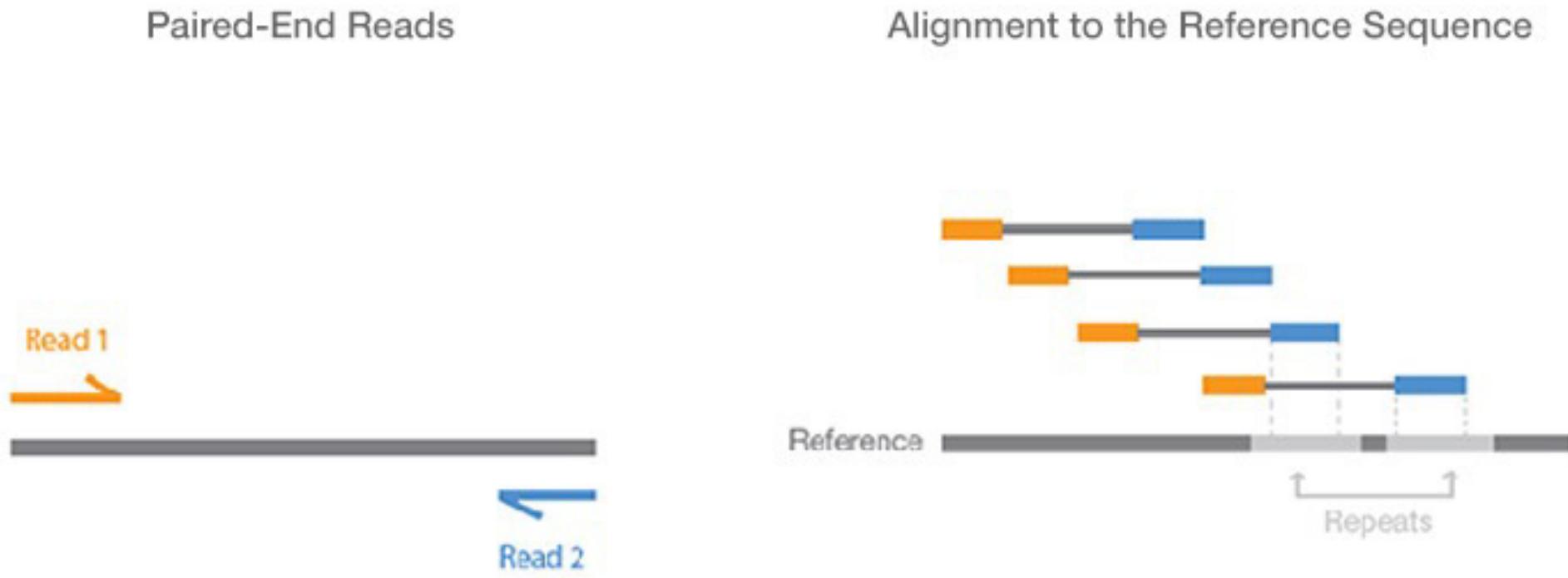
SBS—訊號解讀 (base calling)



The identity of each base of a cluster is read off from sequential images.

Genome assembly (基因組序列組裝)

- 定序時目標序列為 ~500 bp，每次由其兩端讀取 150 bp，這些定序小片段稱之為**read**
- 序列的中間部份為未被定序區域
- 同一目標序列上會有二條**read**，可幫助序列組裝正確性
- 組裝原理為具有重覆序列的片段即可能位在基因體同一位置



定序組裝(assembly)的難題

- 人類基因組為3G(30億個鹼基)，假設定序時每次只能讀 150 bp，那麼至少要有2百萬條reads，若考慮重疊性，那至少要一千二百萬條reads才能完整涵蓋基因組。
- 每條序列都如同一塊小拼圖，如何把六百萬條序列完整地對到染色體？尤其是染色體不同位置但序列卻即為相似或序列多次重覆之位置。
--> 將染色體分為大片段，並得知每大片段是由染色體那一區域而來

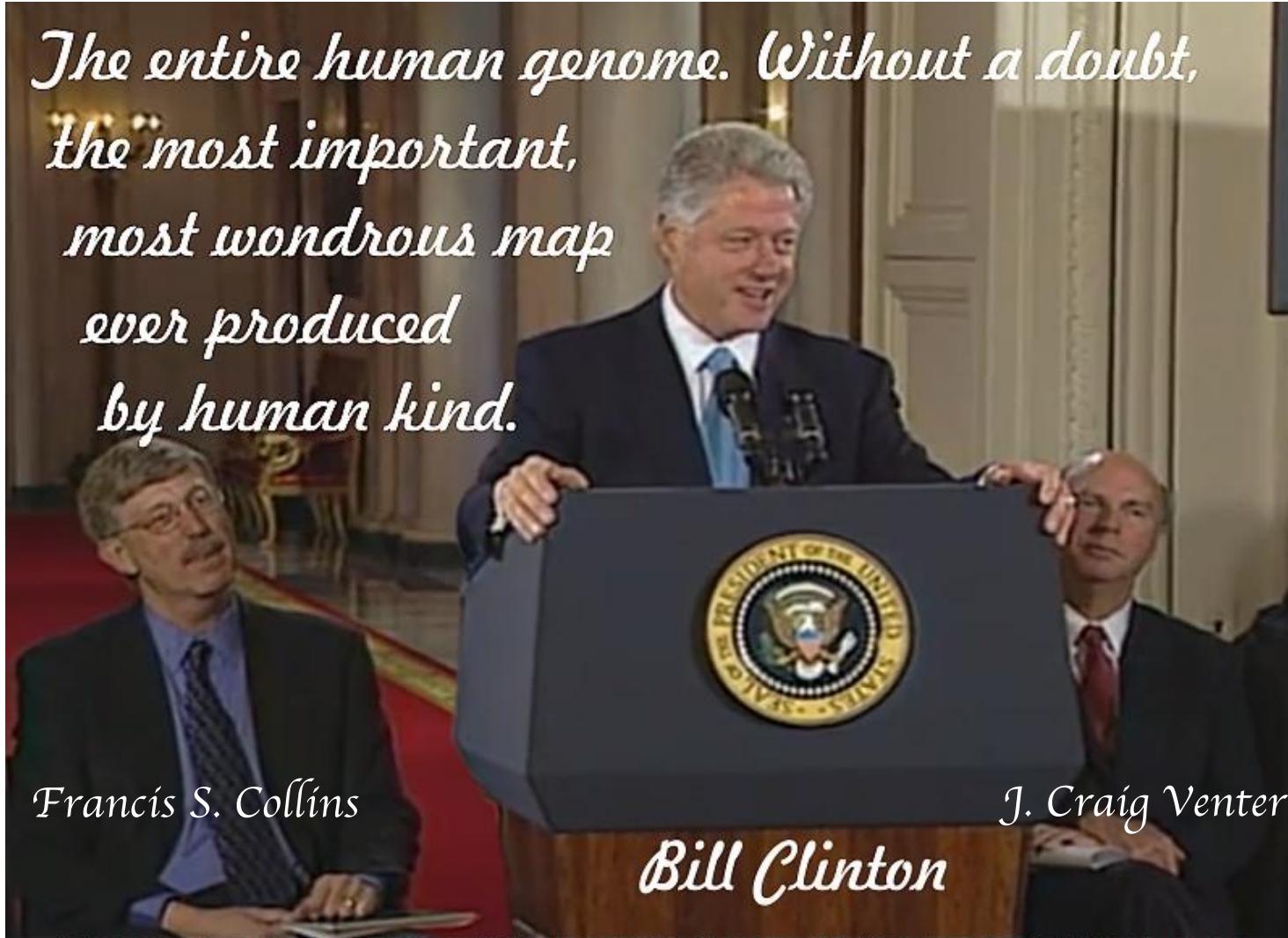


霰彈式定序及階層式定序比較

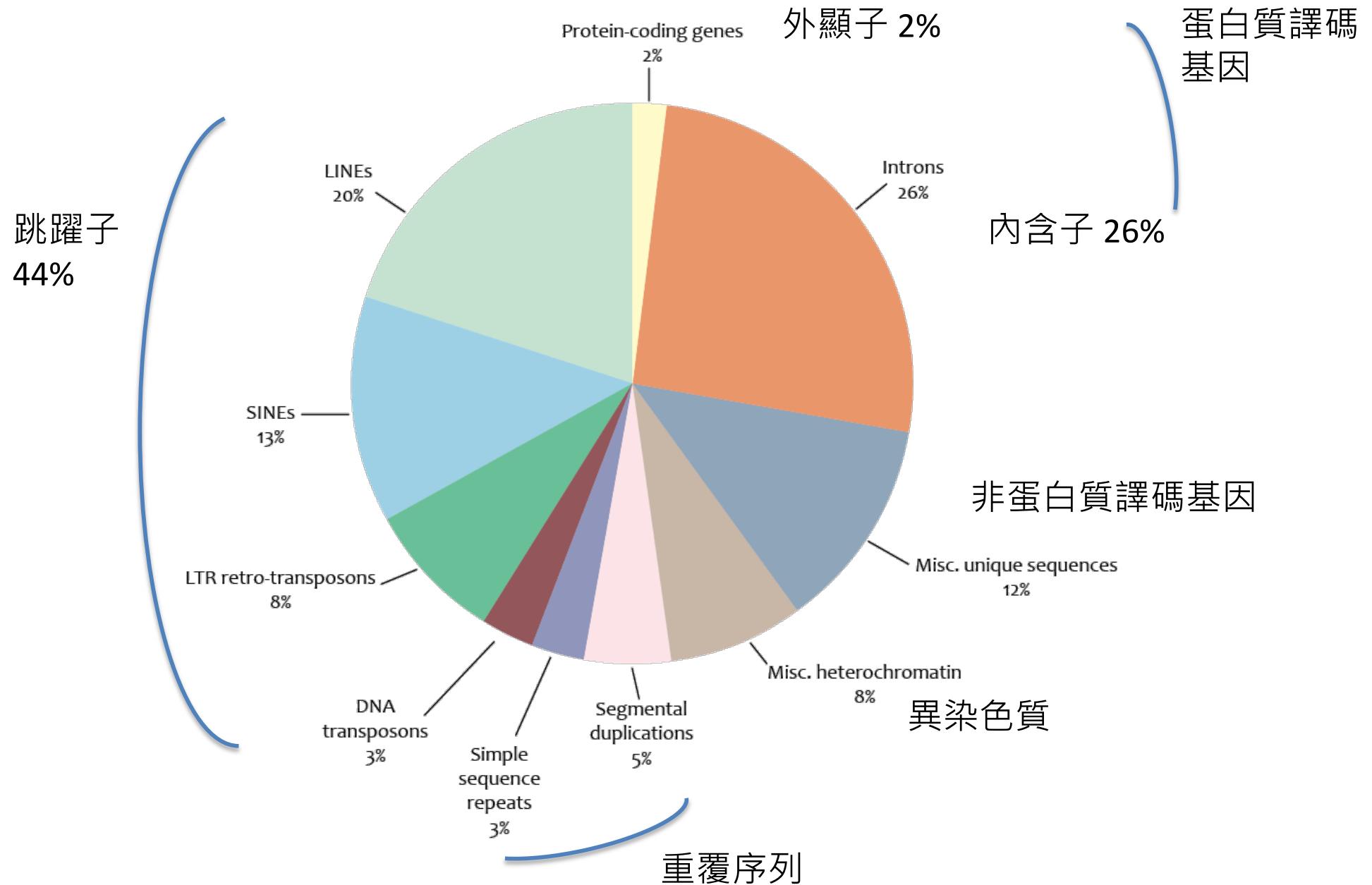
	霰彈式定序	階層式定序
時間	短	長
經費	少	多
人力	少	多
空缺區域(gap)	大	小

- 當第一個物種以階層式定序出來後，可作為其它物種基因組參考圖譜，因此其它物種即使使用霰彈式定序法，也可以得到較精確的基因組圖譜。
- 目前已有公司號稱只要1,000美元即可完成個人定序

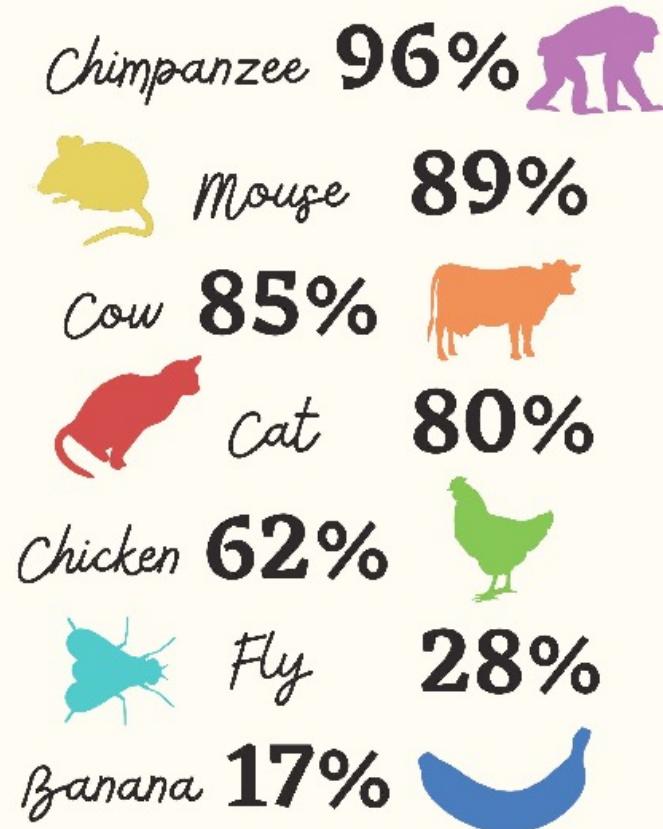
人類基因體計劃共識－所有基因為人類共享，不可註冊



人類基因組序列分析



What portion of genes do humans share with...?



- the portion of the human genome that contains a unique ortholog in the comparison species, considering only protein coding genes
- calculated using omadb, a python wrapper for the OMA database
- inspired by Natasha Glover's blog post "The Banana Conjecture"

定序後的挑戰—基因註解

Q:人類基因組約3.2G，含有約20,000-25,000個蛋白質譯碼基因，如果以每個基因平均長度10 kb 計算，那麼這些基因約佔基因組8% (250,000 kb)。那麼我們要如何知道基因的位置？

1. 轉錄體定序(RNA-seq)：定序RNA的序列，並將其對回至基因組相對位置，該方法較準確，但部份表現量較低的基因不一定會被定序

2. 軟體預測：由基因的特質預測，如果有參考基因庫會較準確

a. ORF finder: 尋找開放讀序框架

對原核生物較有用，因為沒有內含子(intron)

缺點：對真核較無用

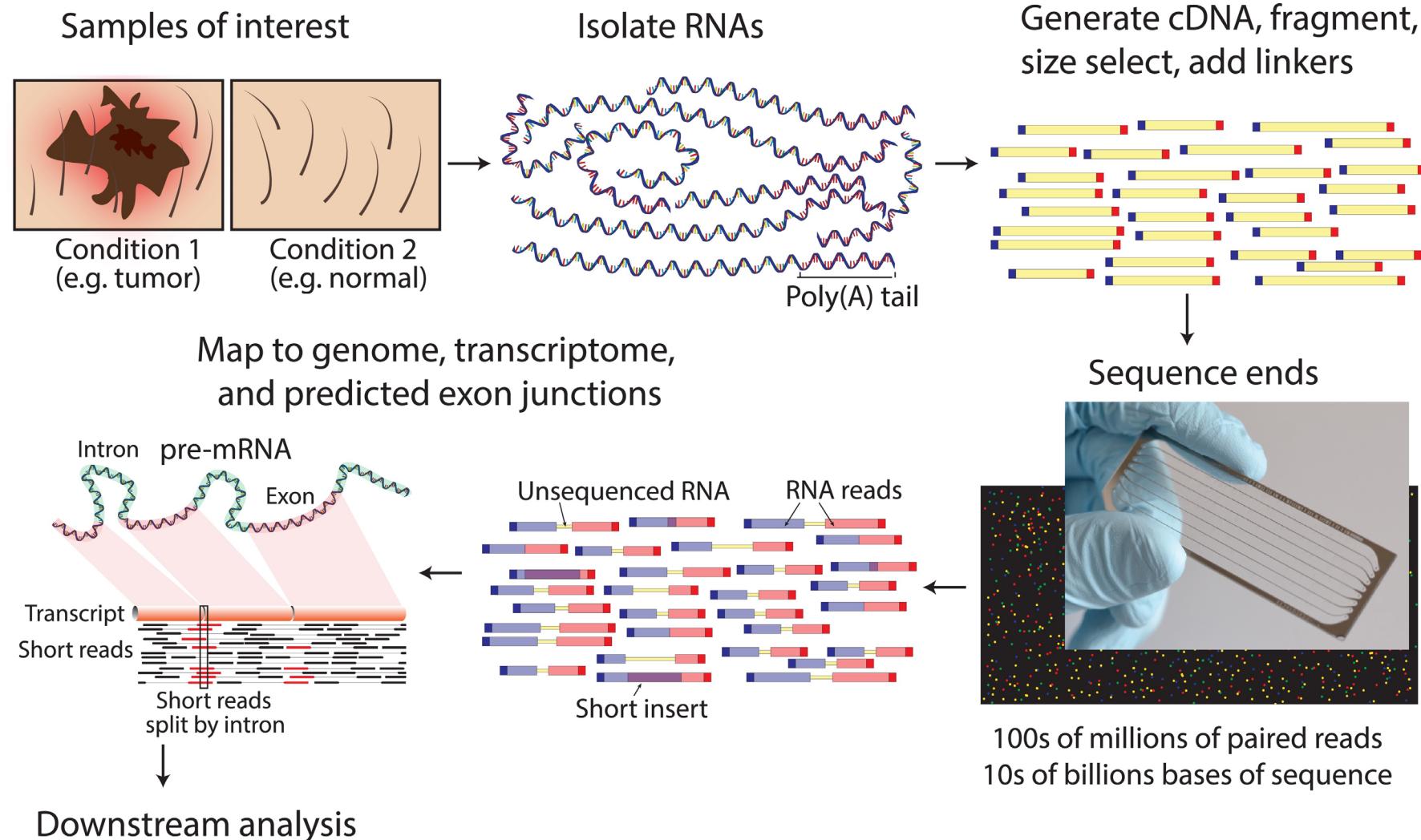
b. 比對已知基因，並對現有基因進行預測

Softberry – FGENESH+

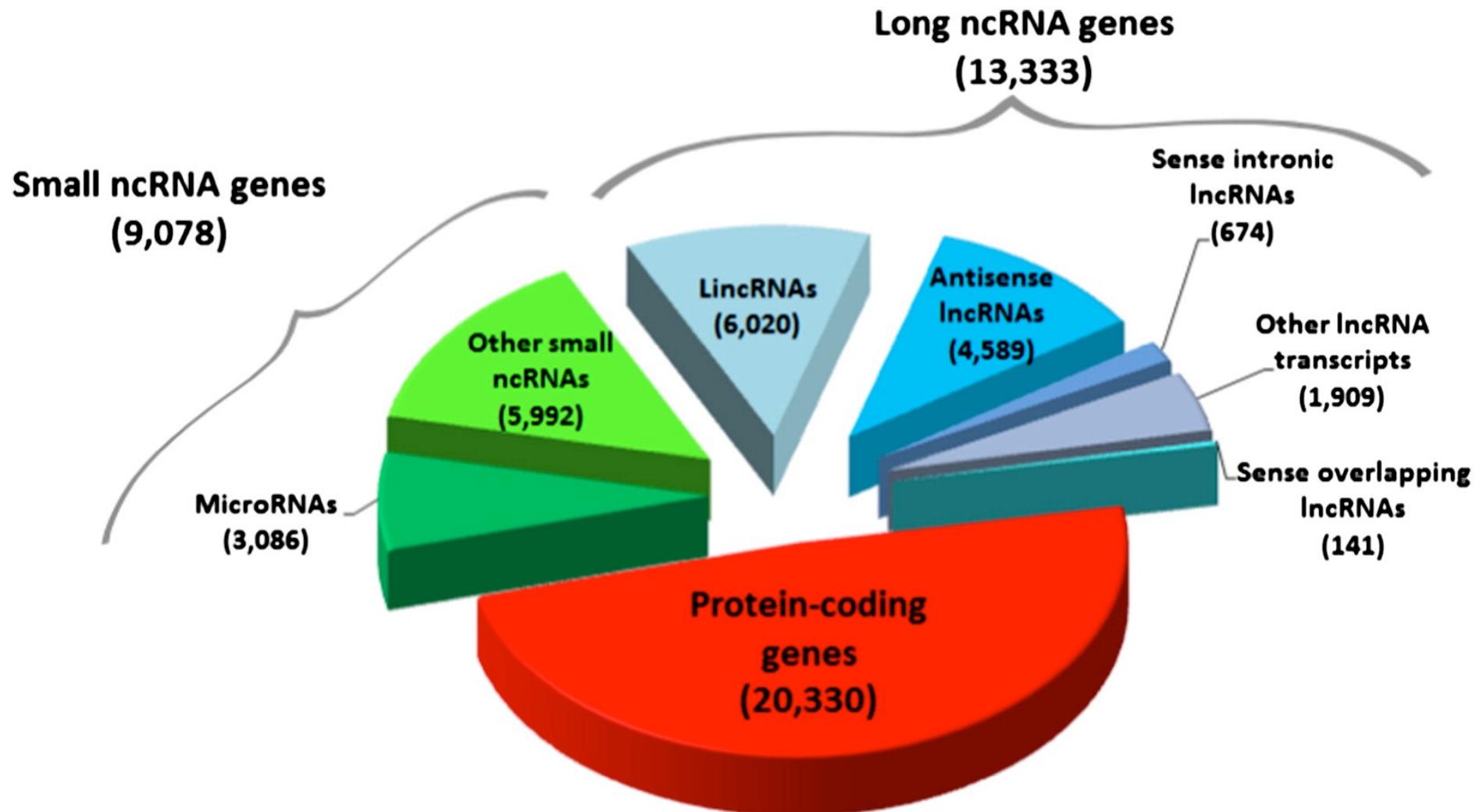
缺點：對變異較大的基因較難進行預測

轉錄體定序 (RNA-seq)

抽取RNA → 轉成雙股DNA並切成小片段 → 穗彈式定序 → 對應到基因組相對位置



人類基因組基因註解



ORF預測軟體 – translate tolls in Expasy

Q：假設你有一條來自大腸桿菌序列如下，請問你要如何以尋找ORF的方式預測該片段是否含有基因片段？

aaacggctcatcgtaaaggcgtattgccatgctaaatctggtaccggcaagcagttat
gtgaaaacgctggaacatctgattcgtagaaggatgtccaggaatagaaaaaatacatcag
cgacattgacagttatgtcaagagcttgctgttagcaaggttagcctattacatgaacaatatg
aacgttaattattgccgatgaccatccgatagtcttgcgttattcgcaaactactgagcaa
attgagtgggtgaatgttgcggcgaattgaagactctacagcactgatcaacaacctgccc
gaaactggatgcgcatgtgtgattaccgatctctccatgcctggcgataagtacggcgatg
gcattacctaatacaagtacatcaagcgccattcccaagcctgtcgatcattgttctgactat
gaacaacaaccggcgattcttagtgcggattggatctggatatcgaaggatcgtgctga
aacaaggcaccgaccgatctgccgaaagctctgcccgctgcagaaaggaaagaaat
ttacccggaaagcggttcgcctgtggaaaaaatcagtgcgtggttacggtgacaag
cgtctctcgccaaaagagagtgaagttctgcgcctgtttcggaaggcttcgtggatcaccga
gatcgctaaaaagctgaaccgcagtattaaaccatcagtagccagaagaaatctgcgatg
atgaagctgggtgtcgagaacgatcgccctgctgaattatctcttcagtgacctaagt
ccggcagataaagactaatcacctgttaggccagat

Open reading frame(ORF): 開放讀序框架

- 基因組中能被轉譯成氨基酸序列的DNA片段
- 一條DNA序列共含有六種可能性的讀取框架(reading frame)
- 框架讀取：每3個鹼基為一組
- 若讀取框架中，含有起始(ATG, Methionine)及終止碼(TAA, TAG, TGA)，且該讀取框架夠長，則其可能為一個基因
- 若讀取框架片段沒有包含起始或終止碼，則該片段則有可能是某基因片段

AAACCATGCTAAATCTGGTAGGCAAGCAGTTGTGAAAAA

Frame 1 -> K P C - I W - A S S C E

Frame 2 -> N H A K S G R Q A V V K

Frame 3 -> T M L N L V G K Q L - K

 G H - I Q Y A L L Q S F <- Frame 4

 V M S F R T P L C N H F <- Frame 5

 F W A L D P L C A T T F <- Frame 6

ORF預測結果

5'3' Frame 1

KRLIVLKAYLPC-IWYPASSYVKRWNI-FVRRMFQE-KNTSATLTVM**S**RACCSKVAYY**MNNM**NVIIADDHPIVLFGIRKSLEQIEWNVVGEFEDSTAL
INNLPKLD~~A~~HVLITDLS**M**PGDKYGDGITLIK**I**KRHFPSLSIIVLT**MNNN**PAILS~~A~~VLDLDIEGIVLKQGAPTDLPKALAALQKGKKFTPESVSR~~L~~EK
ISAGGYGDKRLSPKESEVRLFAEGFLVTEIAKKLNRSIKT~~IS~~SQKKSA**MM**KLGVENDIALLN~~Y~~LSVTLS PADKD-S PVGQ

5'3' Frame 2

NGSSS-R RICHAKSGTRQAV**M**-NAGTSDS-EGCSRNKR~~I~~HQRH-QLCQELAVAR-PIT-TI-T-LLPMT**I**R-SCSVFANHLSKLSG-MLSANLKT**LQH**-
STTCRNW**MRMC**-LPISPCLAI~~STA~~**M**ALP-S STSSAISQACRSLF-L-TTTRRF~~L~~VRYWIWISKGSC-NKVHRPICRKLS~~PR~~CRKGRNL~~PR~~KAFLACWKK
SVLVVTVTSVRQKRVKFCACLRKASW-**PR**SLKS-TAVLKPSVARRNLR--S WVSRTISPC-I ISLQ-P-V RQIKTNHL-**AR**

5'3' Frame 3

TAHRLKGVF**M**LNLVPGKQLCETLEH~~L~~IREKDVP~~G~~IEKYISDIDSYVK~~SLL~~-QGS~~I~~LLHEQYERNYCR-P SD~~S~~LVRYSQIT-AN-V GECCRRI-**RLYSTD**
QQPAETGCACVDYRSLHAWR-VRRWHYLNQVHQ~~A~~PFPKPV~~D~~HCSDYEQ~~Q~~PGDS-CGIGSGYRRDRAETRCTDRSAESSRRAE~~E~~EIYPGKRFSPVGKN
QCWWLR-Q ASLAKRE-S SAPVCGR~~L~~PGDRDR-KAE~~P~~QY-NH~~Q~~-P E~~E~~ICDDEAGCRERYRP~~A~~ELSLFSDLKSGR-RLITCRPD

3'5' Frame 1

IWP~~T~~GD-SLSAGLK~~V~~TEER-F SRAISFSTPSFI IADFFWLL**M**V LILRF~~S~~FLAISVTRKPSANRRRT~~S~~LGERRLSP-P~~P~~ALIFSNRRETLSGVNFFP
CSAARAFGRSGAPCFSTIPSISRSNTALRIAGLLFIVRT**M**IDRLGKWL**MY**LIK**V**MPSPYLS~~P~~GMERSVINTCASSFG~~R~~LLISAVESSNSPTTFHSI
CSSDLRIPNK~~T~~IGWSSAIITFIL**F**M-A TLLQ~~Q~~ALDITVNADVFFYSWNILLTNQ**MF**QRF~~T~~-L~~LAG~~YQI-HGKYAFKT**MSR**

3'5' Frame 2

SGLQVISLYLPDLRSLKRDNSAGRYSRHPASSSQISSGY-WF-YCGSAF-RSRSPGSLPQTGAELHSLLARDACHRNHQH-F FPTGEKRFPG-I SSLS
AARREL~~S~~ADRSVHLVSARSLRYPDP~~I~~PH-ESPGCCS-SEQ-STGLGNGA-CT-LR-CH~~R~~RTYRQAWRDR-S THAHPVSAGC-SVL-SLQIRRQHSPTQF
AQVICEYTRLSDGHRQ-LRSYCSCNRLPCYSKLLT-LS**MSL**MYFSI~~P~~GTFSRIRC~~S~~VSHNCLPGTRFS**M**ANTPLRR-AV

3'5' Frame 3

LAYR-LV~~F~~ICRT-GH-REIIQ~~Q~~GDIVLDTQLHHRRFLLATDG~~F~~NTAVQLFSDLGHQEAFRKQAQNFTLFWR~~E~~LT~~V~~T~~T~~STDF~~F~~QQARNAFRGKFLPFL
QRGESFRQIGRC~~T~~LFQHD~~P~~FDI~~Q~~IQYRTKNRRVVHSQNNDRQ~~A~~WE**M**ALDVLD-GNAIAVLIAR~~H~~GEIGNQH**M**R~~I~~QFRQVVDQCCRVFKFADNIHPLNL
LK-FANTEQDYRM**V**IGNNYVHVHIGYLATASS-HNCQCR-CI FLFLEHPSHESDVP~~A~~F~~H~~ITACRV~~P~~D~~L~~AWQ~~I~~R~~L~~-DDEPF

由已知基因進行基因預測

- 真核生物基因含有內含子，較難以真核生物基因含有內含子，較難以ORF finder的方式預測，因此可以已知蛋白或RNA序列進行預測
- 範例：人類胰島素基因片段 1120 bp, 蛋白質序列 110 a.a.

ATG GCCCTGTGGATGCGCCTCCTGCCCTGCTGGCGCTGCTGGCCCTCTGGGGACCTGACCCAGCCGCAGCCTTGTGAACCAAC
ACCTGTGCGGCTCACACCTGGTGGAAAGCTCTACCTAGTGTGCGGGGAACGAGGCTTCTTACACACCCAAGACCCGCCGGGA
GGCAGAGGACCTGCAGGGTGAGCCAAC TGCCCATTGCTGCCCTGGCCGCCCGAGCCACCCCTGCTCCTGGCGCTCCCACCCA
GCATGGGCAGAAGGGGGCAGGGAGGCTGCCACCCAGCAGGGGGTCAGGTGCACTTTTTAAAAAGAAGTTCTTGGTCACGTC
CTAAAAGTACCAGCTCCCTGTGGCCCAGTCAGAACATCTCAGCCTGAGGACGGTGTGGCTCGCAGCCCCGAGATAACATCAGAG
GGTGGGCACGCTCCTCCACTCGCCCTCAAACAAATGCCCGCAGCCCATTCTCCACCCCTCATTGATGACCGCAGATTCAA
GTGTTTGTAAAGTAAAGTCTGGGTGACCTGGGTACAGGGTCCCCACGCTGCCTGCCTGGCGAACACCCCATCACGCC
CGGAGGAGGGCGTGGCTGCCTGCCTGAGTGGGCCAGACCCCTGTCGCCAGGCCTACGGCAGCTCCATAGTCAGGAGATGGGG
AAGATGCTGGGACAGGCCCTGGGAGAAGTACTGGGATCACCTGTTAGGCTCCACTGTGACGCTGCCGGCGGGGA
AGGAGGTGGGACATGTGGCGTTGGGCCTGTAGGTCCACACCCAGTGTGGGTGACCCCTCCCTAACCTGGGTCCAGCCGGC
TGGAGATGGTGGGAGTGCACCTAGGGCTGGCGGGCAGGCAGGACTGTGTCTCCCTGACTGTGTCCCTGTGTCCCTGC
CTGCCGCTGTTCCGGAACCTGCTCGCGGCCAGTCCTGGCAGTGGGCAGGTGGAGCTGGCGGGGCCCTGGTCAGGC
AGCCTGCAGCCCTGGCCCTGGAGGGTCCCTGCAGAAGCGTGGCATTGTGGAACAATGCTGTACAGCATCTGCTCCCTTACCA
AGCTGGAGAACTACTGCAACTAG

MALWMRLPLLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAEDLQVGQVEL
GGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN

若以ORF預測軟體預測人類胰島素基因，會發現無法預測出完整基因片段

5'3' Frame 1

MALWMRLPLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAEDLQGEPTAHCCWPAPPATPCSWRSHPAWAEGGRRLPPSRSGSGALF-KEVLL
VTS-K-PAPCGPVRISA-GRCWLQRQPRDTSEGHHAPPSTRPSNKCPAAHFSTLI--PQIQVFC-VKS梧TWGHHRVPHAAICLWANTPSRPEEGVAACLSGPDPPCRQASRQL
HSQE-MGKMLGTGPGEKYWDHLFRPLI-RCPGAGEGGGTCRWGL-VHTQCG-PSL-PGSSPAGDGWECDLGLAGRRALCLPDCVLLCPSASPLFRNLLCAARPGSGAGGA
GRGPWCRQPAALGPGGVPAEAWHCGTMLYQHLLPLPAGEELLQL

5'3' Frame 2

WPCGCASCPCWRCWPSGDLTQPQPL-TNTCAAHTWWKLST-CAGNEASSTHPRPAGRQRTCRVSQLPIAAPGRPQPPPAPGAPTQHGQKGAGGCHPAGGQVHFFKKFWS
SRPKSDQLPVAQSESQPEDGVGFSGSPEIHQRVGTLPLAPQTNAQPISPPSFDDRRFKCFVK-SPG-PGVTGCPTLPASGRTPHHARRRAWLPA-VGQTPVARPHGSS
IVRRWGRCWGQALGRSTGITCSGHCDAAPGRGKEVGHVGVGACRSTPSVGDPSENGLPARLEMGGSAT-GWRAGGHCVSLTVSSCVPLPRRCSGTCARHVLAVGQVEL
GGGPGAGSLQPLALEGSLKRGIVEQCCTSICSLYQLENYCN-

5'3' Frame 3

GPDAPPAPAGAAGPLGT-PSRSLCEPTPVRLTPGGSSLPSVRGTRLLLHTQDPPGGRGPAG-ANCPLLPLAAPSHPLLALPPSMGRRGQEATQQGVRCTFLKRSSLG
HVLKTSSLWPSQNLSLRTVLASAAPRYIRGWARSSLHSPLKQMPRSPFLHPHLMTADSSVLLSKVLGDLGSQGAPRCLPLGEHPITPGGGRGCLPEWARPLSPGLTAAP
-SGDGEDAGDRPWGEVLGSPVQAPTVTLPREGGRRWDMWALGPVGPVWVTLPLTWVQPGWRWVGVRPRAGGQAGTVSP-LCPPVSLCLAAVPEPALRGTSWQWGRWSW
AGALVQAACSPWPWRGPCRSVALWNNAVPASAPSTSRTTAT

3'5' Frame 1

LVAVVLQLVEGADAGTALFHNTILLQGPLQQQLQAACTRAPAPAQLHLPHCQDVPRAGSGTAARQRDTGGHSQGDTVPACPPALGRTPTHLQPGWTQVRGRVTHTGCGPT
GPNAMSHLLPPPRGSVTVGA-TGDPSTSPQGLSPASSSPSPDYGAAVRPGRGLAHSGRQPRPPPGVMGCSPRGRQRGAPCDPRSPRTLLNKTEAVIK-GWRNGLRGI
CLRGEWREERAHPLMYLGAAEANTVRLRLF-LGHRELVTFR-PRELLFKKVHLTCPWVAASCPLLPMLGGSARSRGWLGAARGSNGQLAHPAGPLPPGGSWVCRSLVP
RTLGELEPPGVSRGVGSQRLRLGQVPRGPAAPAGAGGASTGP

3'5' Frame 2

-LQ-FSSW-REQMLVQHCSTMPrFCRDPSRAKGCRLPAPGPPSSTCPTARTCRAEQVPEQRRGRGTQEDTVRETQCAPPARQP-VALPPISSRAGPRLEGGSPTLGVDLQ
APTPTCPTSFPRPGAAASQWEPEQVIPVLLPRACPQHLPHLLTMEPLP-GLATGVWPTQAGSHALLRA-WGVRPEAGSVGHPVTPGHPGLYLTKHNLNRSSNEGGE-MCGAF
V-GASGGRSVPTL-CISGLPKPTPSSG-DSDWATGSWSLLGRDQENFFLKKCT-PAGWQPPAPFCPCWVGAPGAGGGWGRPGAAMGSWLTQVLCLPAGLGCVEEASFP
AH-VESFHQV-AAQVLVHKCGCWRSPEGQQRQQQGEAHPQGH

3'5' Frame 3

SCSSSPAGRSRCWYSIVPQCHASAGTPPGPRAAGCLHQGPRAPPAPLPGRAAQSRFRNSGEAEGRRTQSGRHSARLPASPRSHSHPSAGLDPG-REGPHWVWTYR
PQRPHVPPPSAPGQRHSGSLNR-SQYFSPGPVPSIFPIS-LWSCREWQSGGPLRQAATPSSGRDGFAQRQAAWGLT-PQVTQDFT-QNT-ICGHQMRVEKWAAGHL
FEGRVEGGACPPSDVSRCRSQHRPQAEILTGPQGAGHE-DVTKRTSF-KSAPDPLLGGSLLPPSAHGRWERQEQGVAGGGQGQQWAVGSPCRSSASRRVLGV-KKPRSP
HTR-RASTRCEPHRCWFTKAAAGSGPQRASSASRGRRRIHRA

FGENESH+ of Sofeberry

- 利用已知蛋白序列對未知基因組序列進行基因預測

The screenshot shows the FGENESH+ web interface on the Softberry platform. The left sidebar lists various bioinformatics tools: Home, Gene finding in Eukaryota, Gene finding with similarity (highlighted in dark blue), Operon and Gene Finding in Bacteria, Gene Finding in Viral Genomes, Next Generation, Alignment (sequences and genomes), Genome visualization tools, Search for promoters/functional motifs, Deep learning recognition, Protein Location, RNA structures, Protein structure, Pathway prediction, Protein/DNA 3D-Visual Works, and Manipulations with sequences.

The main content area is titled "Services Test Online" and features the "FGENESH+" tool. It includes a reference section: "Reference: Solovyev V.V. (2007) Statistical approaches in Eukaryotic gene prediction. In Handbook of Statistical genetics (eds. Balding D., Cannings C., Bishop M.), Wiley-Interscience; 3d edition, 1616 p." Below this is a text input field for nucleotide sequences, which contains a sequence of DNA bases. A red box labeled "待預測基因組序列" (Target genome sequence) has an arrow pointing to this input field.

Below the nucleotide sequence input is a text input field for protein sequences, containing a sequence of amino acids. A red box labeled "人類胰島素蛋白序列" (Human insulin protein sequence) has an arrow pointing to this input field.

At the bottom, there is a dropdown menu set to "Human (Homo sapiens)" and two buttons: "search" and "Reset". To the right, a note states: "Total 539 genome-specific parameters are available for genefinders of FGENESH suite".

FGENESH+預測結果

G	Str	Feature	Start	End	Score	ORF	Len	Derry	Derry	Derry	Derry	
1	+	1 CDSf	1	-	187	286.95	1	-	186	186	1	62 100
1	+	2 CDS1	975	-	1120	172.67	977	-	1120	144	64	111 100

Predicted protein(s):

```
>FGENESH:[mRNA] 1 2 exon (s) 1 - 1120 333 bp, chain +
ATGGCCCTGTGGATGCGCCTCCTGCCCTGCTGGCGCTGCTGCCCTCTGGGACCTGAC
CCAGCCGCAGCCTTGTGAACCAAACACACTGTGCGGCTCACACCTGGTGGAAAGCTCTAC
CTAGTGTGCGGGGAACGAGGCTTCTTCTACACACCCAAGACCCGCCGGAGGCAGAGGAC
CTGCAGGTGGGGCAGGTGGAGCTGGCGGGGCCCTGGTGCAGGCAGCCTGCAGCCCTG
GCCCTGGAGGGTCCCTGCAGAAGCGTGGCATGTGGAACAATGCTGTACCAAGCATCTGC
TCCCTCTACCAGCTGGAGAACTACTGCAACTAG
>FGENESH: 1 2 exon (s) 1 - 1120 110 aa, chain +
MALWMRLLPLLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAEQ
LQVGQVELGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN
```

外顯子1

外顯子2

結果：RNA 序列

結果：蛋白質序列

人類胰島素基因序列 1120 bp

ATG GCC CTGT GGAT GCG CCTC CTGCCCTGCTGGCGCTGCTGCCCTCTGGGACCTGACCCAGCCGAGCCTTGTGAACCAAC
ACCTGTGCGGCTCACACCTGGTGGAAAGCTCTACCTAGTGTGCGGGGAACGAGGCTTCTACACACCCAAGACCCGCCGGGA
GGCAGAGGACCTGCAGGGTGAGCCAAGTGCCTGCCATTGCTGCCCTGGCCGCCAGCCACCCCTGCTCCTGGCGCTCCCACCCA
GCATGGGCAGAAGGGGCAGGAGGCTGCCACCCAGCAGGGTCAGGTGCACTTTTAAAAAGAAGTTCTTGGTCACGTC
CTAAAAGTACCAGCTCCCTGTGGCCCAGTCAGAATCTCAGCCTGAGGACGGTGTGGCTCGCAGCCCCGAGATAACATCAGAG
GGTGGGCACGCTCCTCCACTCGCCCTCAAACAAATGCCCGCAGCCCATTCTCCACCCCTCATTGATGACCGCAGATTCAA
GTGTTTGTAAAGTAAAGTCCTGGGTGACCTGGGTACAGGGTCCCCACGCTGCCTGCCTGGCGAACACCCCATCACGCC
CGGAGGAGGGCGTGGCTGCCTGCCAGAGTGGCCAGACCCCTGTCGCCAGGCCACGGCAGCTCCATAGTCAGGAGATGGGG
AAGATGCTGGGGACAGGCCCTGGGAGAAGTACTGGGATCACCTGTTAGGCTCAGGTGGACCCCTCCCTTAACCTGGTCCAGCCGGC
AGGAGGTGGGACATGTGGCGTTGGGCCTGTAGGTCCACACCCAGTGTGGGTACCCCTCCCTTAACCTGGTCCAGCCGGC
TGGAGATGGGTGGGAGTGCACCTAGGGCTGGCGGGCAGGCCACTGTGTCCCTGACTGTGTCCCTGTGTCCCTGC
CTGCCGCTGTTCCGGAACCTGCTCTGCCAGTGGCAGTGGGCAGGTGGAGCTGGCGGGGCCCTGGTCAGGC
AGCCTGCAGCCCTGGCCCTGGAGGGTCCCTGCAGAACGCGTGGCATTGTGGAACAATGCTGTACCATGCTCCCTTACCC
AGCTGGAGAACTACTGCAACT**TAG**

外顯子1

內含子

人類胰島素RNA序列 333 bp

外顯子2

ATG GCC CTGT GGAT GCG CCTC CTGCCCTGCTGGCGCTGCTGCCCTCTGGGACCTGACCCAGCCGAGCCTTGTGAACCAAC
ACCTGTGCGGCTCACACCTGGTGGAAAGCTCTACCTAGTGTGCGGGGAACGAGGCTTCTACACACCCAAGACCCGCCGGGA
GGCAGAGGACCTGCAGGTGGGCAGGTGGAGCTGGCGGGGCCCTGGTCAGGCAGCCTGCAGCCCTGGCCCTGGAGGG
GTCCCTGCAGAACGCGTGGCATTGTGGAACAATGCTGTACCATGCTGC TCCCTTACCAAGCTGGAGAACTACTGCAACT**TAG**

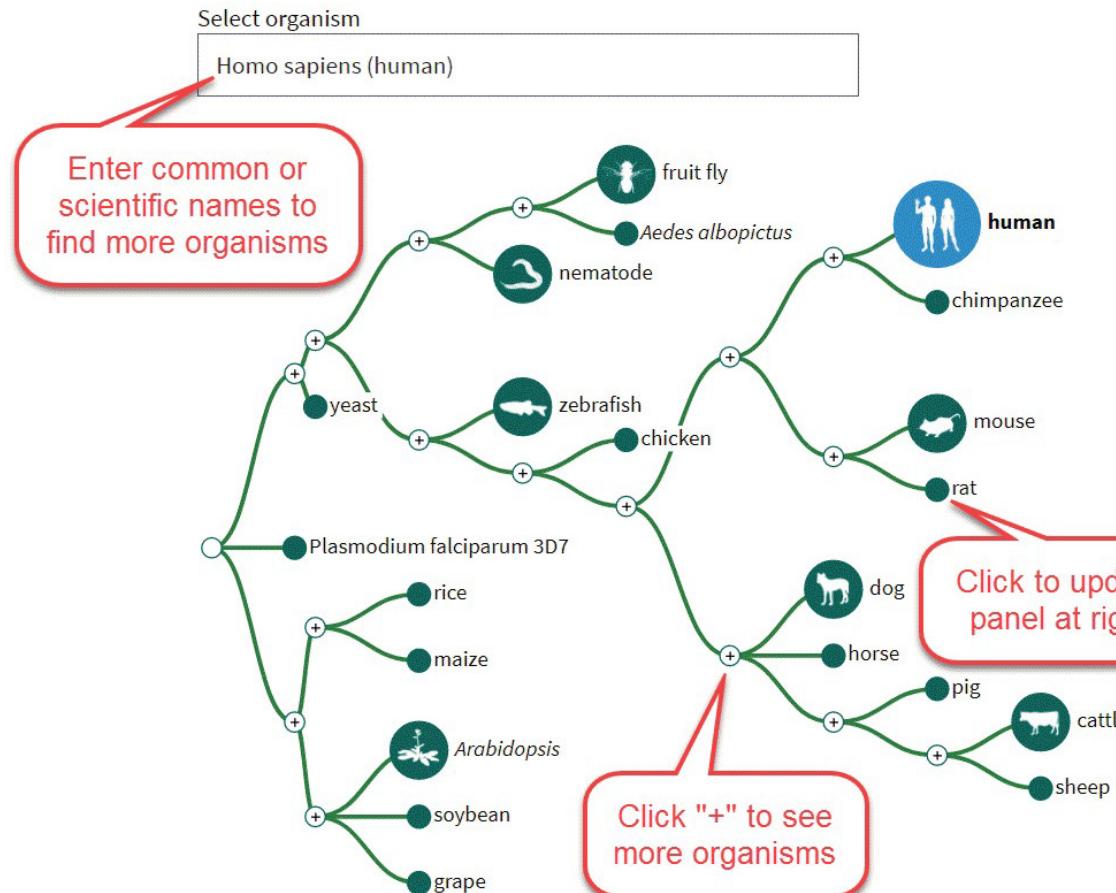
人類胰島素蛋白質序列 110 bp

MALWMRLPLLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAEDLQVGQVEL
GGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN

NCBI Genome Data viewer

Genome Data Viewer

GDV is a genome browser supporting the exploration and analysis of more than 540 eukaryotic RefSeq genome assemblies. [i](#)



Homo sapiens (human) genome

Search within selected assembly

Search in genome
Location, gene or phenotype [Q](#)

Examples: TP53, chr17:7667000-7689000, rs334, DNA repair

Assembly
GRCh38.p11 [▼](#)

[Browse genome](#) [BLAST genome](#)

Assembly details

Name	GRCh38.p11
RefSeq accession	GCF_000001405.37
GenBank accession	GCA_000001405.26
Download via FTP	RefSeq , GenBank
Submitter	Genome Reference Consortium
Level	Chromosome

Annotation details

Annotation Release 108
Release date

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 X Y

Homo sapiens (human) genome



Search in genome



Search “Aquaporin”

Genes

Other

Name	Location
AQP4	Chr18: 26.85M - 26.87M
AQP1	Chr7: 30.91M - 30.93M
AQP2	Chr12: 49.95M - 49.96M
AQP3	Chr9: 33.44M - 33.45M
AQP5	Chr12: 49.96M - 49.97M
AQP9	Chr15: 58.14M - 58.19M
AQP8	Chr16: 25.22M - 25.23M

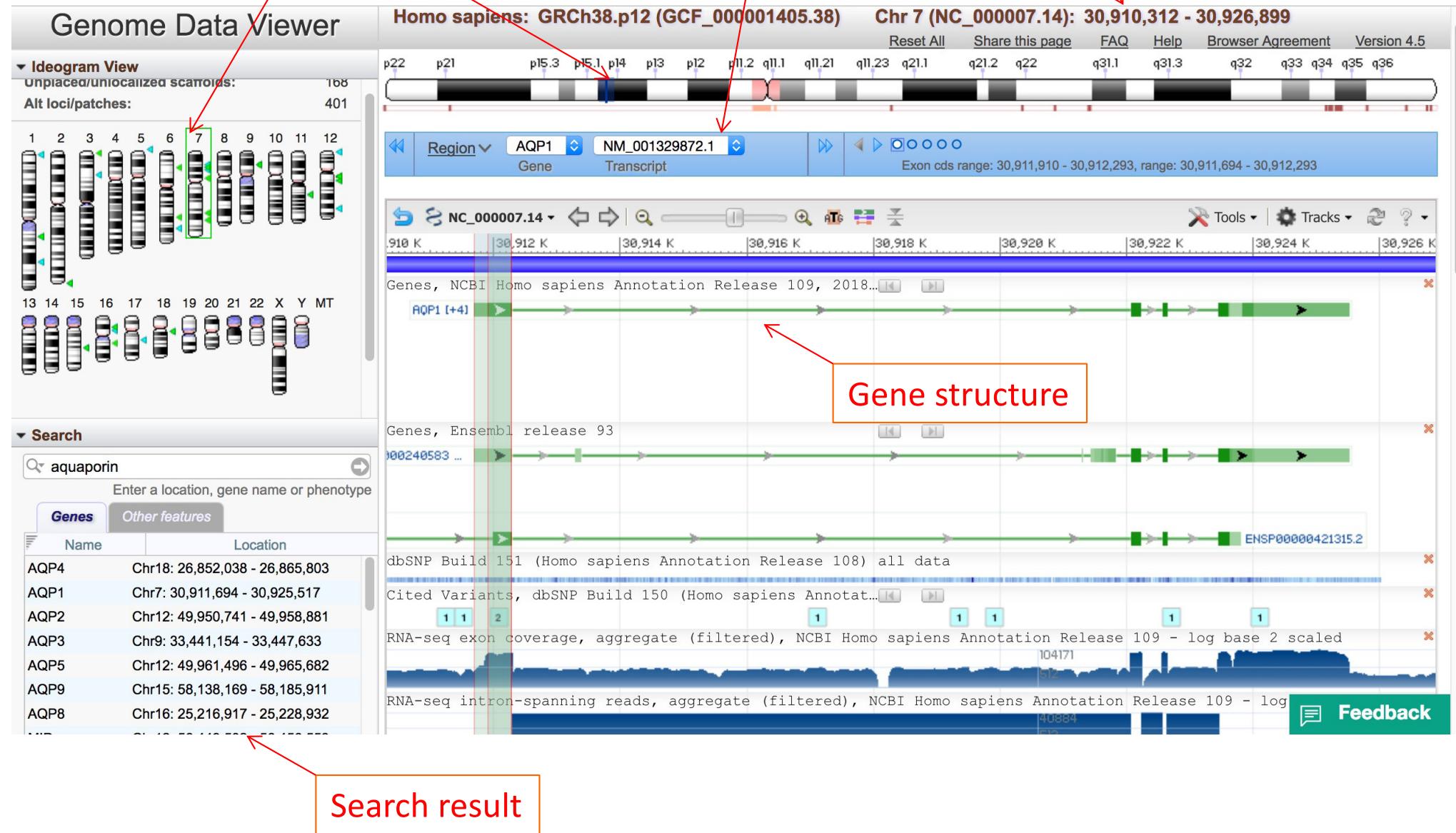


result

Examples: [TP53](#), [chr17:7667000-7689000](#), [rs334](#), [DNA repair](#)

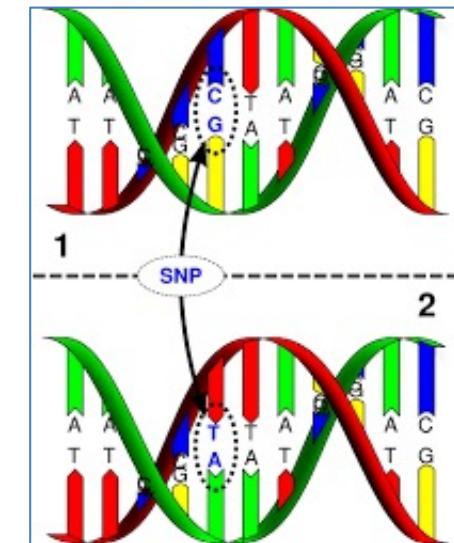
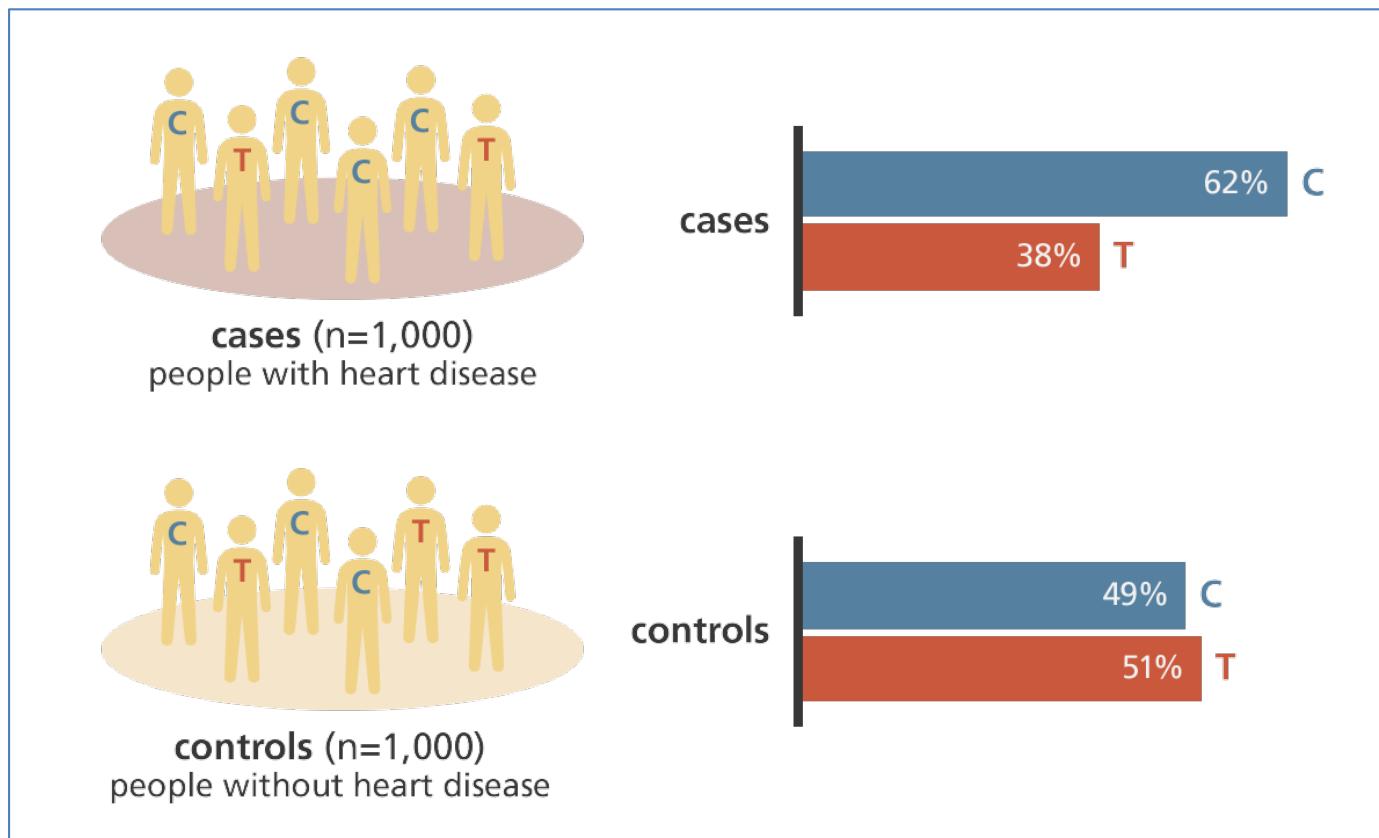
Assembly





人類基因體的應用

1. 基因功能分析
2. 醫藥開發 - 精準醫學
3. 基因尋找 – GWAS
4. Metagenomics – 微生物菌相分析
5. Epigenomics – DNA 甲基化分析
6. 其它



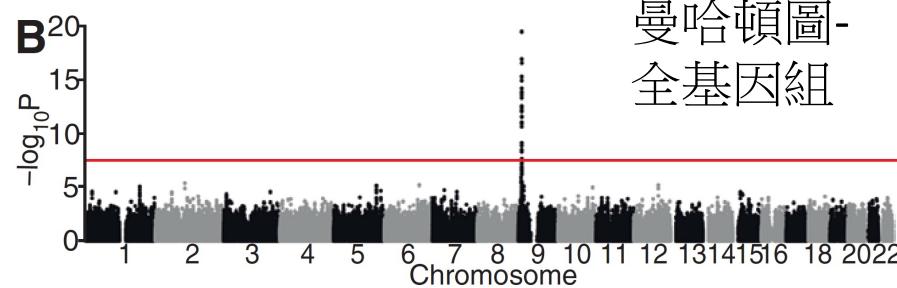
SNP (單一核苷酸
多型性)

GWAS分析

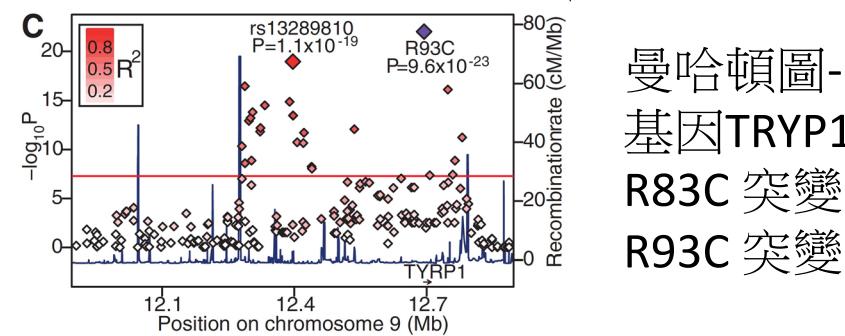
GWAS (全基因組關聯分析)

- 由全基因組中找出變異序列或基因，個體間的遺傳背景不宜差異太大
- 金髮個體：42位， 黑髮個體：43位
- 全基因組定序 → GWAS分析 → 候選基因 → 基因功能分析

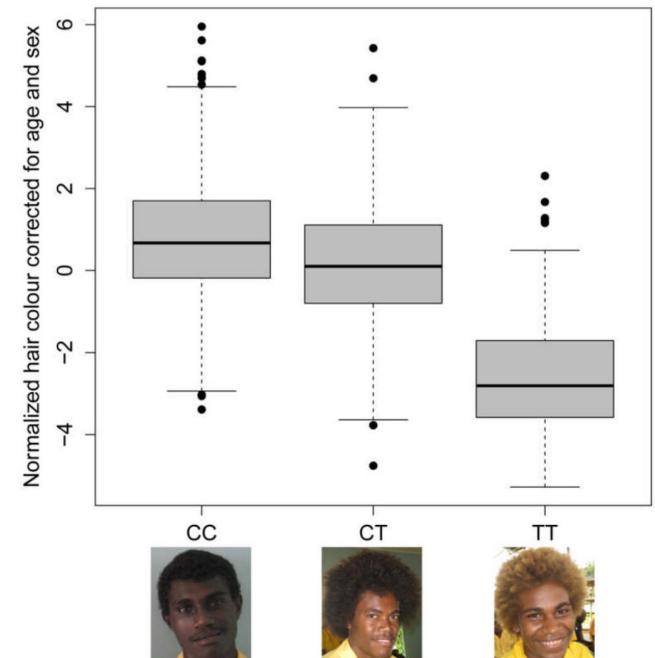
42位



43位



R93C DNA突變



Thanks for your attention