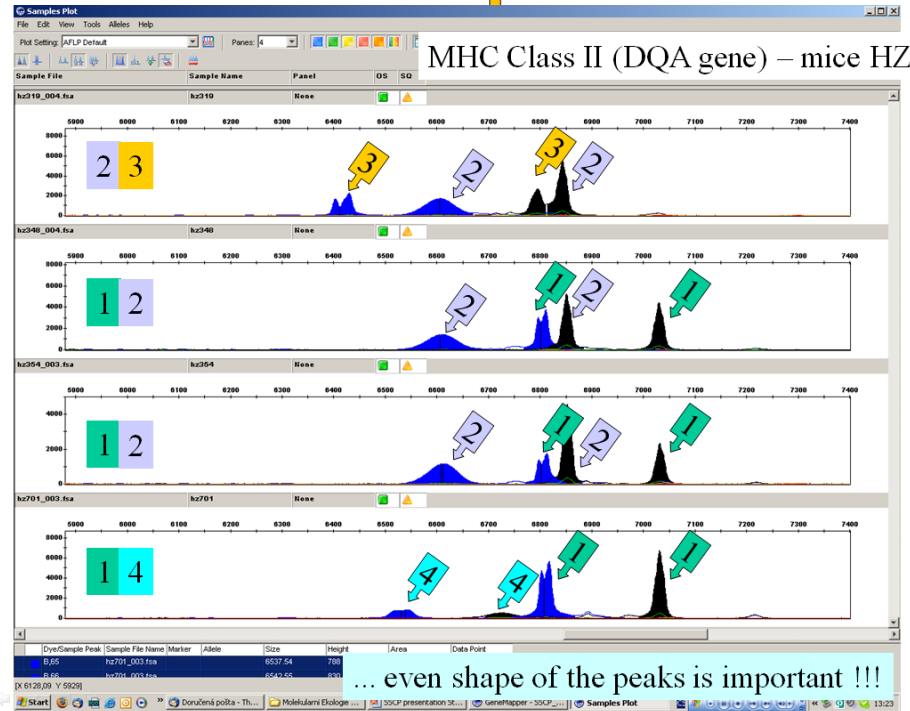
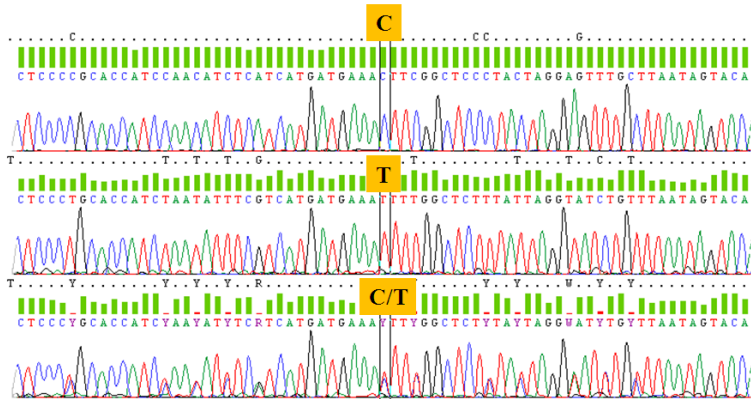


SNPs genotyping - sekvenování? Je drahé a nejasné u heterozygotů

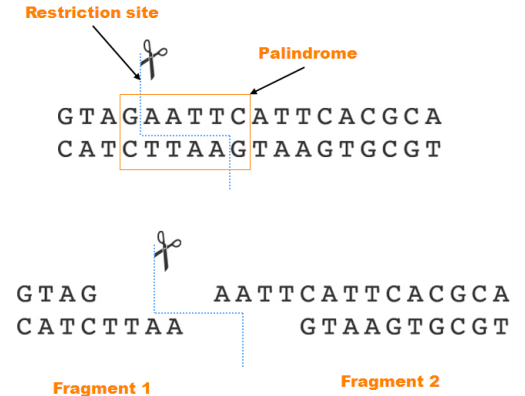


SNP genotyping - old standards

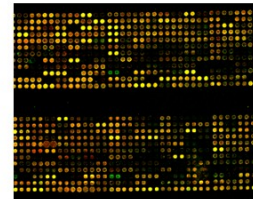
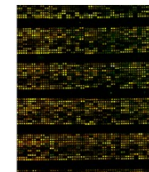
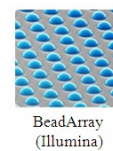
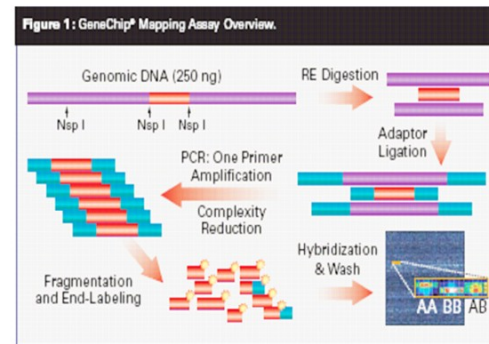
PCR-RFLP
(restriction fragments length polymorphism)

Enzyme Site Recognition

- Each enzyme digests (cuts) DNA at a specific sequence = restriction site
- Enzymes recognize 4- or 6- base pair, palindromic sequences (eg GAATTC)



Detekce: Affymetrix, Illumina



10 – 500 tisíc SNP znaků najednou – „chip technology“

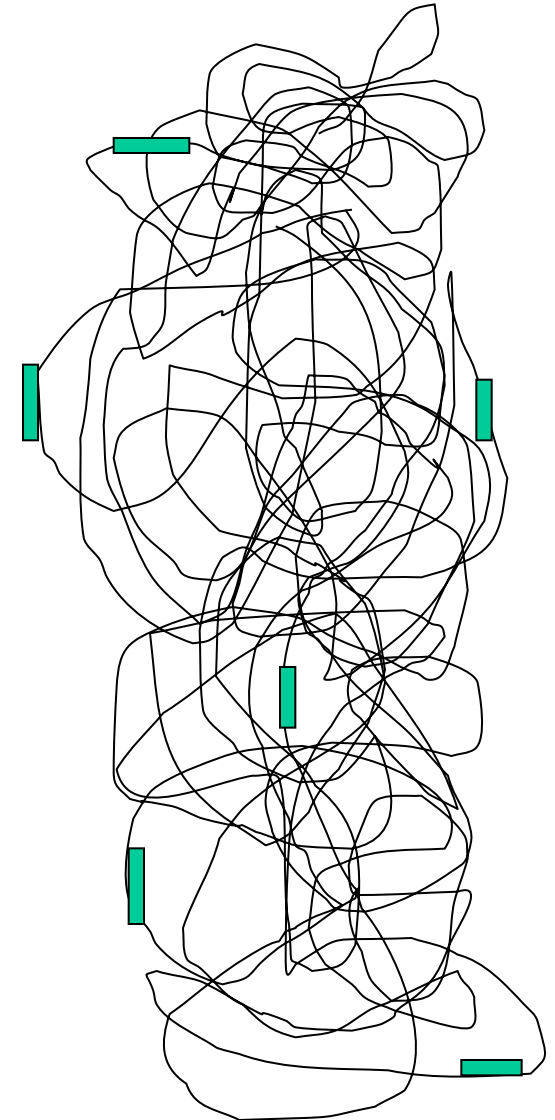
Typy genetických markerů

	Single locus	Codominant	PCR assay	Overall variability
Nuclear multilocus				
Nuclear single locus				
Alozymy	Yes	Yes	No	Low-medium
Mikrosatelite	Yes	Yes	Yes	High
SINE (LINE)	Yes	Yes	Yes	Low
SNPs (sekvence)	Yes	Yes	Yes	Low-high

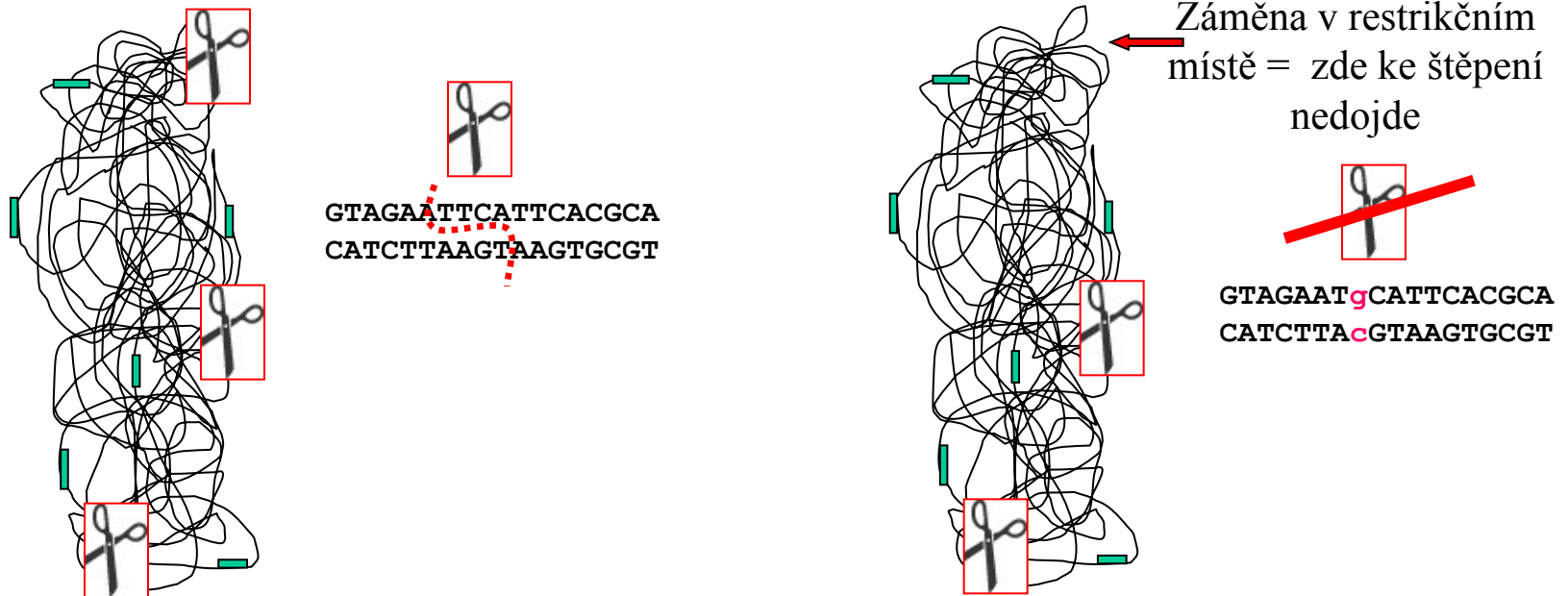
Multi-locus genetic markers

- Mnoho znaků náhodně rozmístěných v genomu - celogenomový scan
- *minisatellite DNA fingerprinting*
- *RAPD* (randomly amplified polymorphic DNA)
- *AFLP* (amplified fragment length polymorphism)
- presence vs. absence **restrikčního místa** (AFLP) či **místa pro dosednutí primerů** (RAPD) = **dominantní znaky** (neodliší heterozygota - proužek na gelu buď je nebo není)
- není nutno znát předem genom studovaného druhu (tj. primery)

Př.: chromozóm 1



Každý jedinec má jedinečný genom



1. Ztráta nebo nabytí restričního místa

Enzyme Site Recognition

- Each enzyme digests (cuts) DNA at a specific sequence = restriction site
- Enzymes recognize 4- or 6- base pair, palindromic sequences (eg GAATTC)

Restriction site

Palindrome

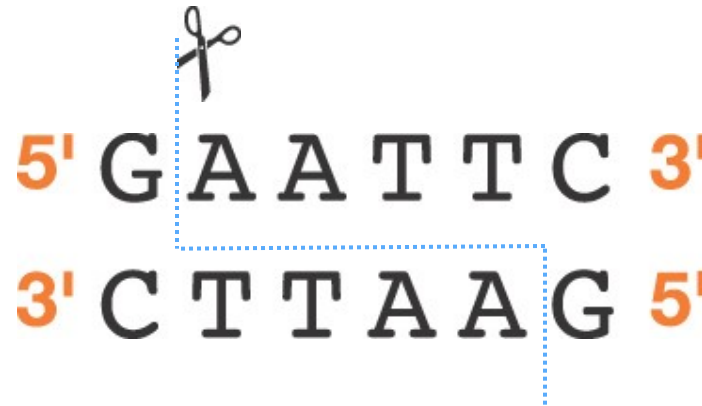
G T A G G A A T T C A T T T C A C G C A
C A T C T T A A G T A A G T G C G T

G T A G A A T T C A T T T C A C G C A
C A T C T T A A G T A A G T G C G T

Fragment 1

Fragment 2

Common Restriction Enzymes



EcoRI

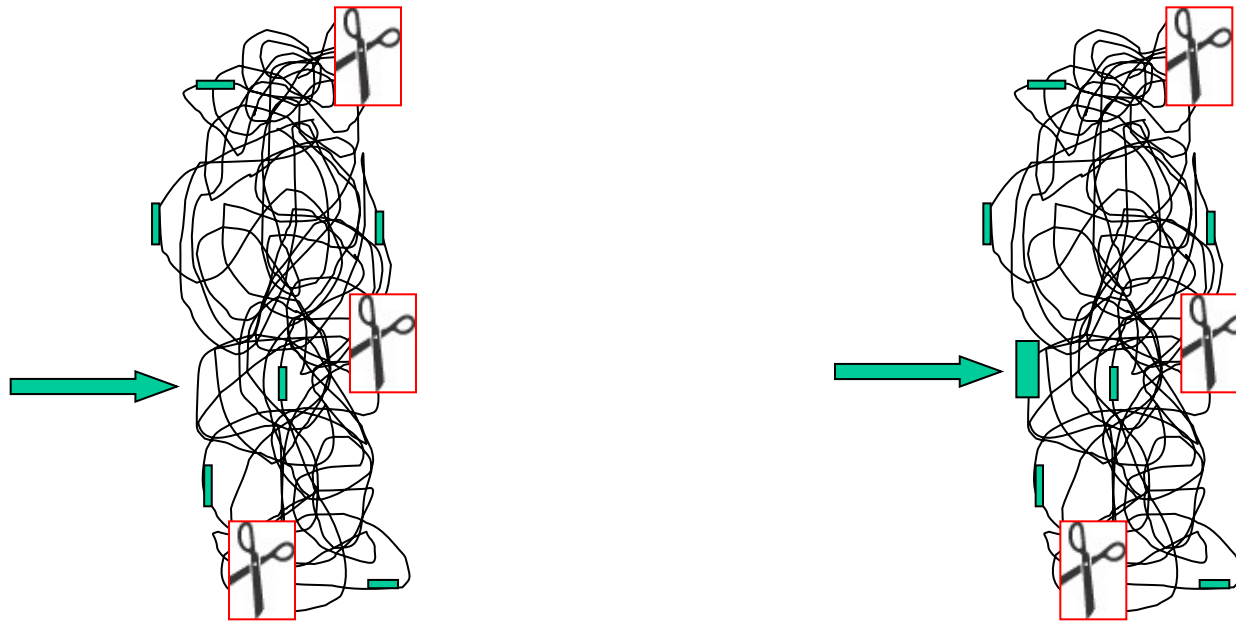
- *Escherichia coli*
- 5 prime overhang



PstI

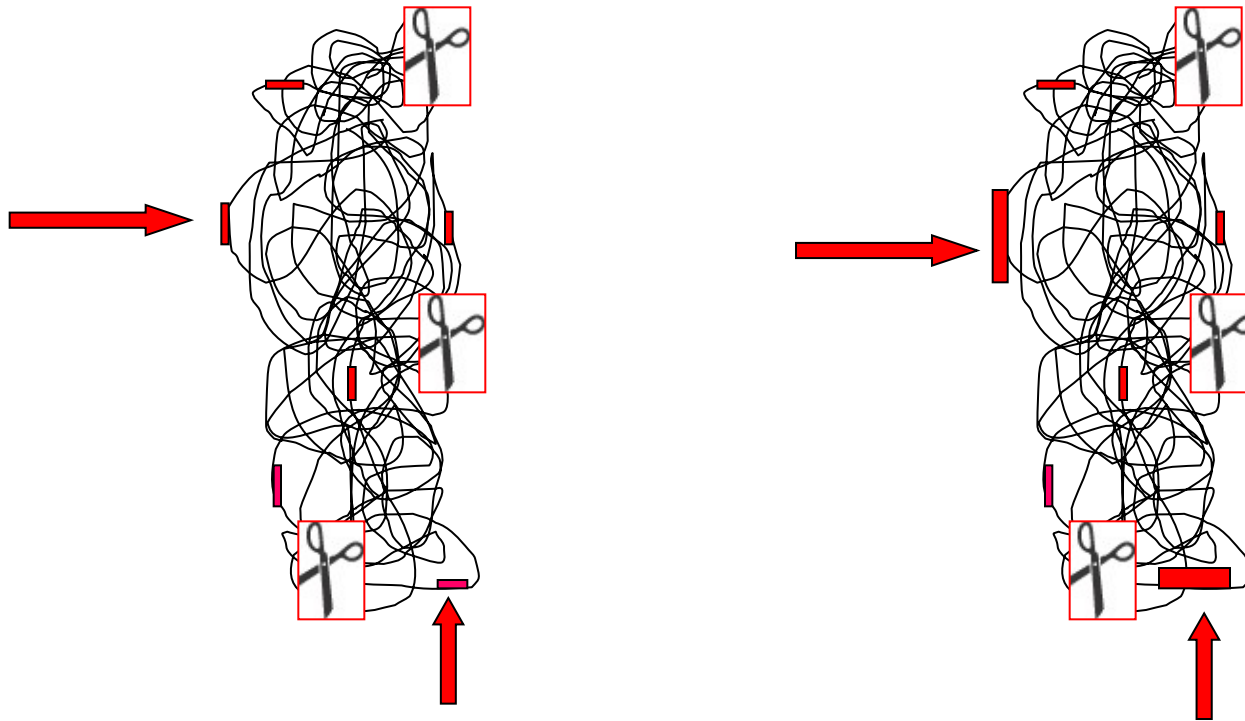
- *Providencia stuartii*
- 3 prime overhang

Každý jedinec má jedinečný genom



2. Ztráta nebo nabytí SINE (např. **Alu** sekvence) nebo LINE

Každý jedinec má jedinečný genom



3. Vysoká mutační rychlost **minisatelitů a mikrosatelitů** -
rozdíly v počtu repeticí, tj. v délce daného úseku

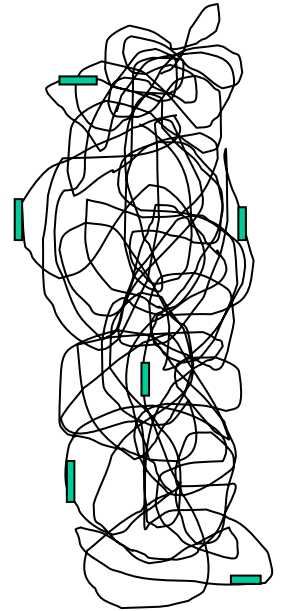
Repetitivní DNA

DNA	Typical sequence length (bp)	Location
Satellites ($>10^6$ repeats/genome)	5-100	Tandem arrays, scattered throughout the genome
Minisatellites ($>10^3$ loci/genome)	20-300	Tandem arrays up to 5 kb in length, scattered throughout the genome
Microsatellites ($>10^4$ loci/genome)	1-6	Tandem arrays up to a few 100 bp in length, scattered throughout the genome
Telomeres	4-8	Tandem arrays up to 1kb in length, at the ends of each chromosome
SINEs ($>10^5$ /genome)	50-500 (100-300)	Interspersed throughout the genome
LINEs ($>10^3$ /genome)	1-5 k	Interspersed throughout the genome

(Minisatellite) DNA fingerprinting

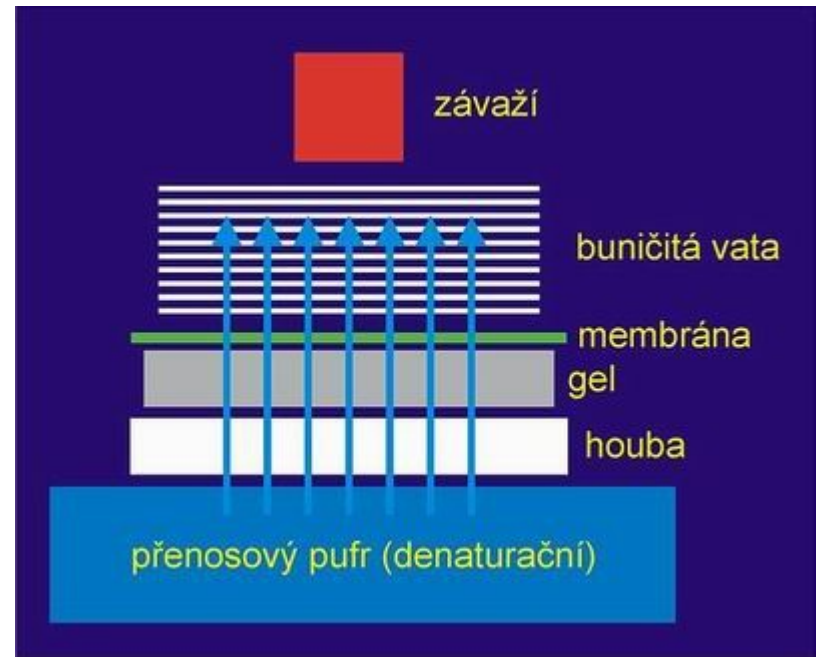
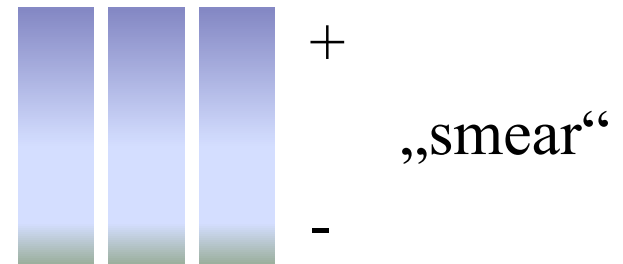
(Jeffreys et al. 1985)

- první celogenomový screening
- restriční štěpení kompletní DNA – sekvenčně specifické **restriční endonukleázy**



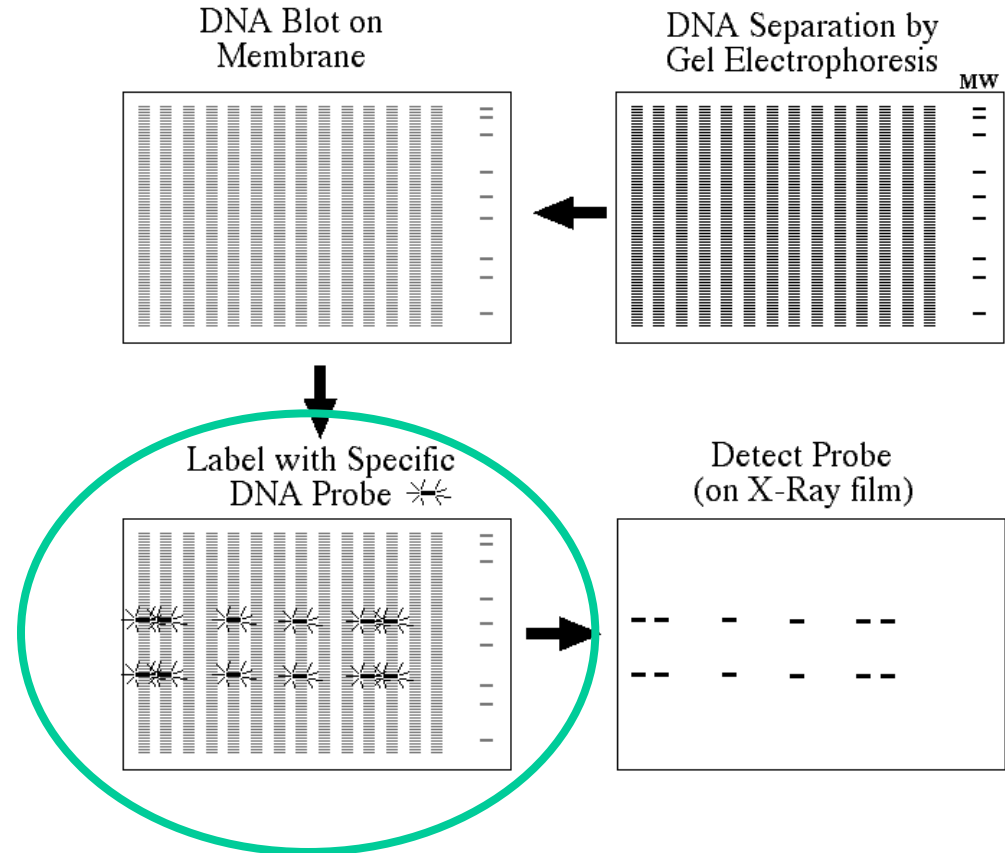
Minisatellite DNA fingerprinting

- elektroforéza rozštěpené DNA
- Southern blotting – přenesení DNA na membránu



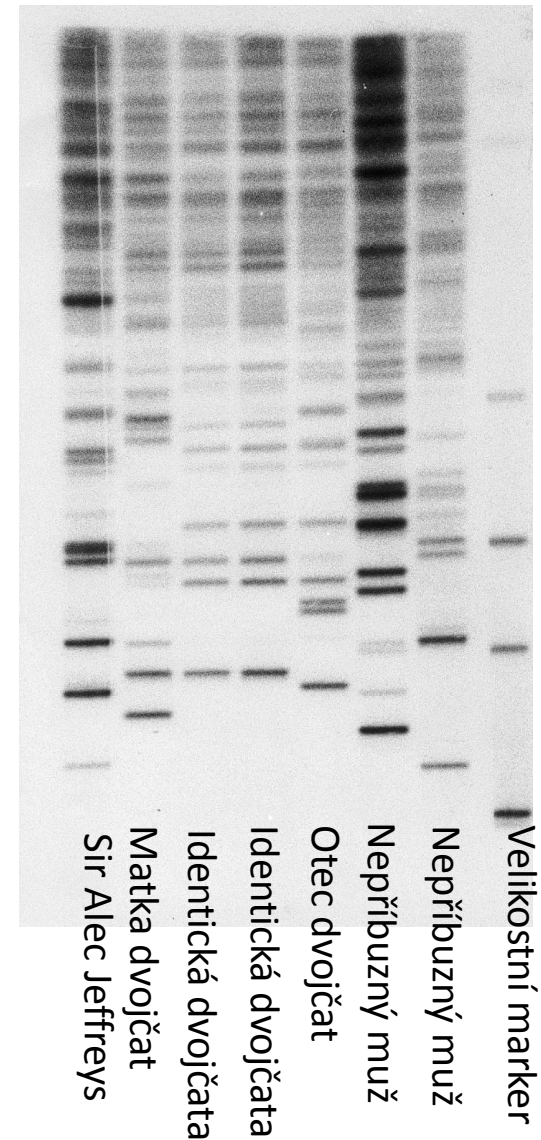
Minisatellite DNA fingerprinting

- elektroforéza
- Southern blotting – přenesení DNA na membránu
- hybridizace se značenou sondou (nejčastěji radioaktivní značení), tj. specifickou sekvencí odpovídající danému minisatelitu (popř. SINE)



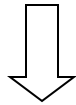
Minisatellite DNA fingerprinting

- elektroforéza
- Southern blotting – přenesení DNA na membránu
- hybridizace se značenou sondou, tj. specifickou sekvencí odpovídající danému minisatelitu
- zásadní objevy např. EPC u ptáků
- v posledních 10-15 letech – přesun k PCR-based metodám



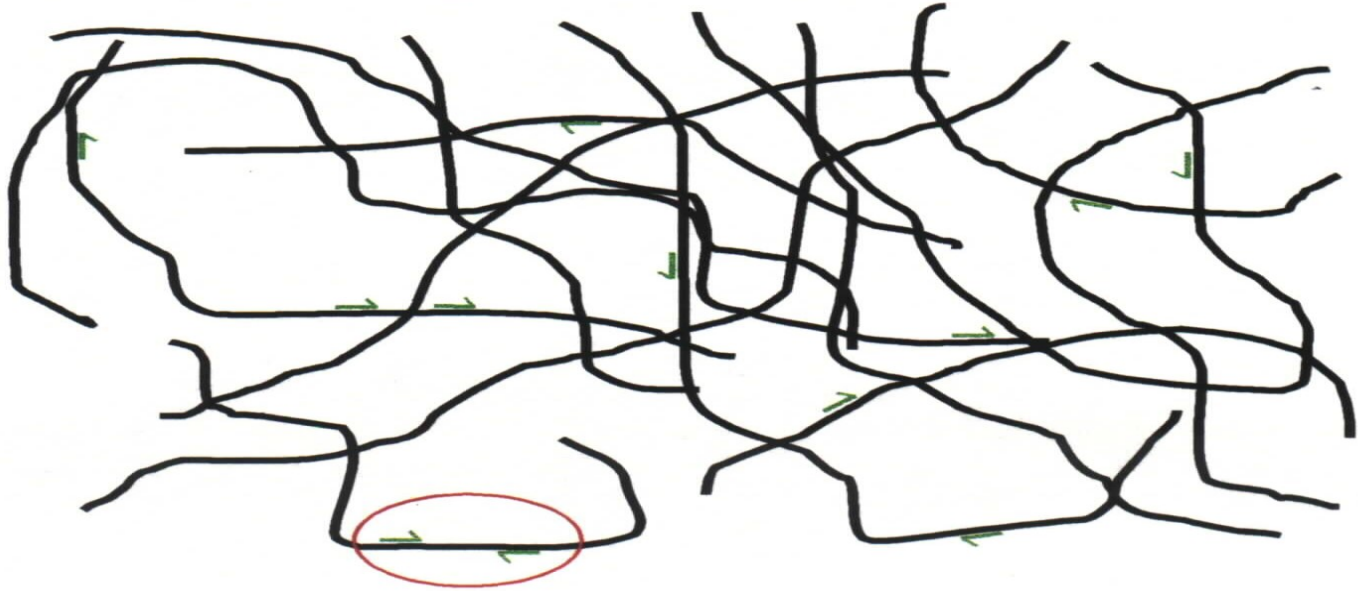
RAPD (randomly amplified polymorphic DNA)

Krátké náhodné oligonukleotidy
(~ 10 bp) jako primery

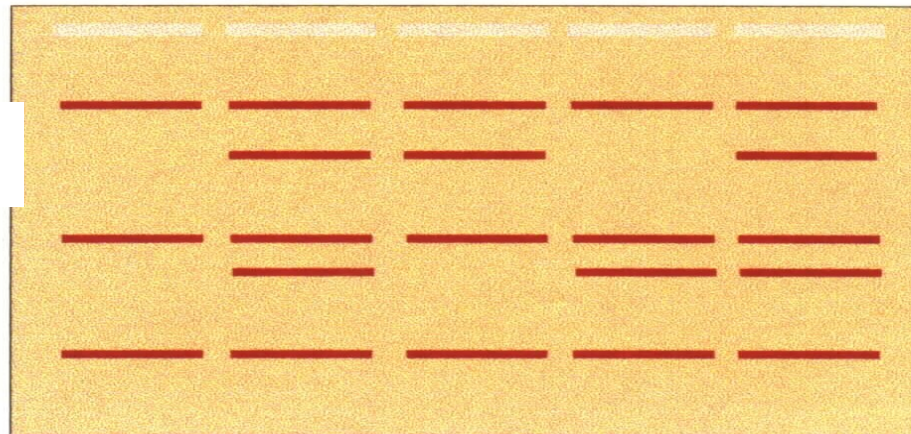


PCR za málo specifických podmínek

genomic DNA



- 1) PCR
- 2) Separation by size on agarose gel

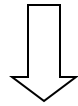


Variabilní DNA detekovaná metodou RAPD je důsledkem ztráty RAPD lokusů v důsledku:

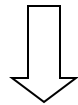
- a) Změna sekvence v místě nasedání primeru
- b) Delece místa nasedání primeru
- c) Velká inzerce mezi dvěma místy nasedání primeru

RAPD - review

Krátké náhodné oligonukleotidy
(~ 10 bp) jako primery

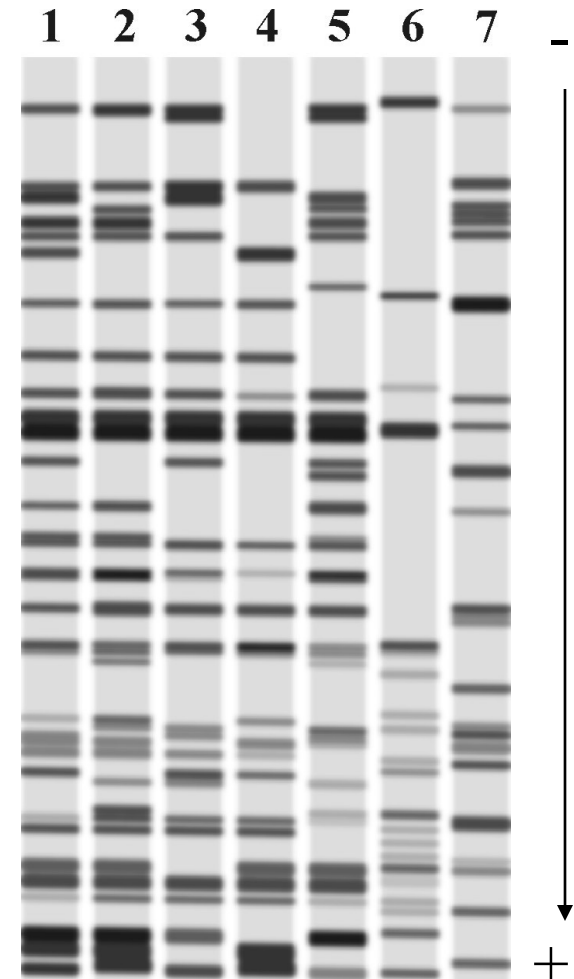


PCR za málo specifických podmínek



Detekce PCR produktů elektroforézou

**Nízká opakovatelnost v důsledku
mnoha faktorů ovlivňujících PCR –
dnes již není akceptována jako
metoda např. pro studium
populační struktury**



AFLP (amplified fragments length polymorphism)

- levná, jednoduchá, rychlá a spolehlivá metoda na generování stovek informativních genetických markerů
- současný screening mnoha různých DNA oblastí distribuovaných náhodně v genomu
- lépe reprodukovatelná než RAPD – obsahuje krok se specifickou PCR
- „genome scan“ – hledání asociací s fenotypovými znaky

Princip AFLP metody („generating AFLP markers“)

(a) AFLP template preparation

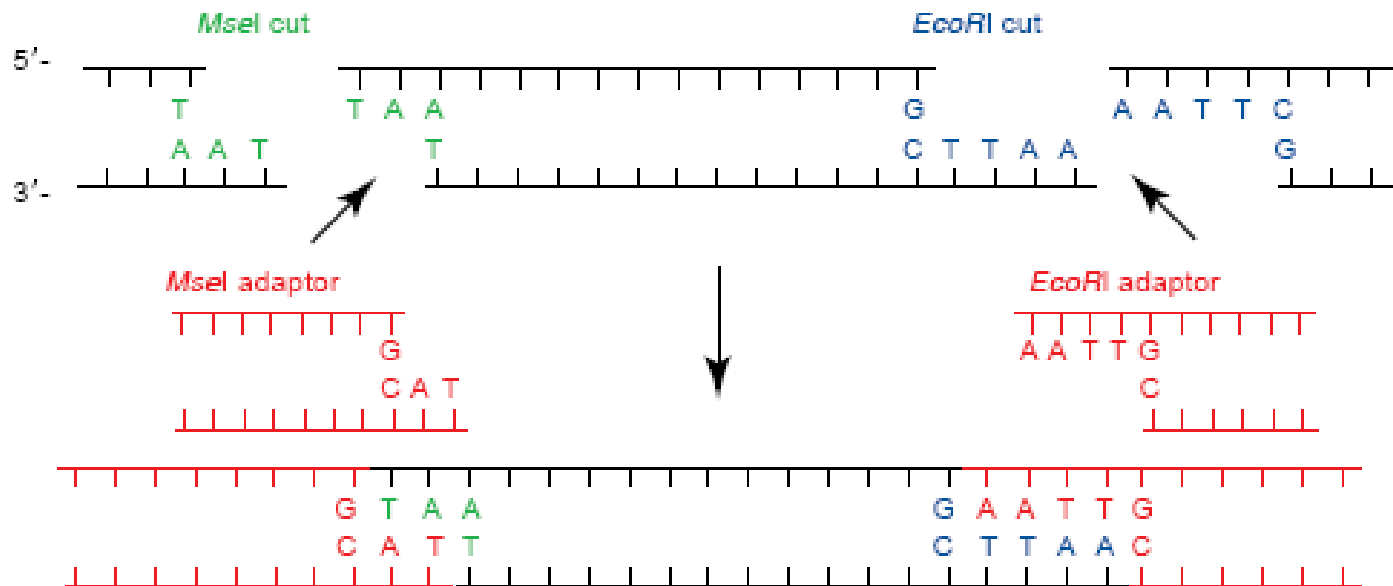
Whole genomic DNA



Restriction enzymes
(*MseI* and *EcoRI*)
and
DNA ligase

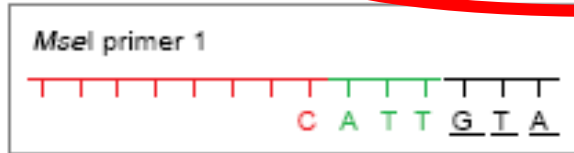


(b) Restriction and ligation

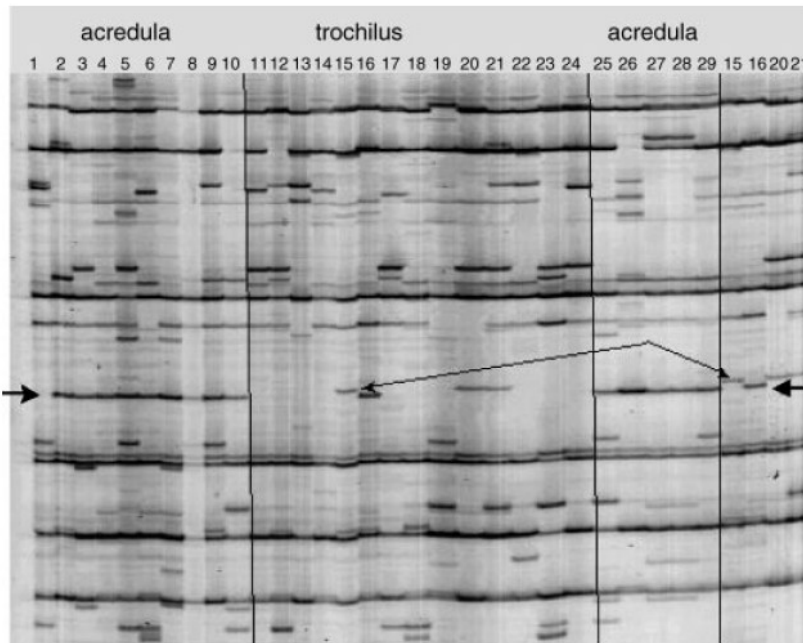
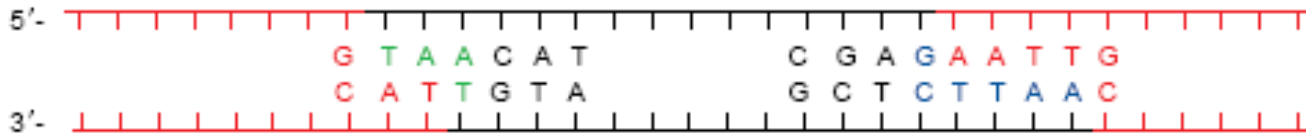


Generating AFLP markers

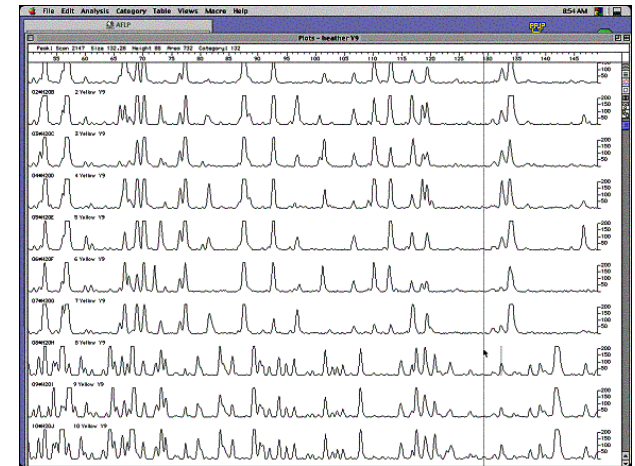
(c) Selective amplification (one of many primer combinations shown)



PCR with primers on adaptors



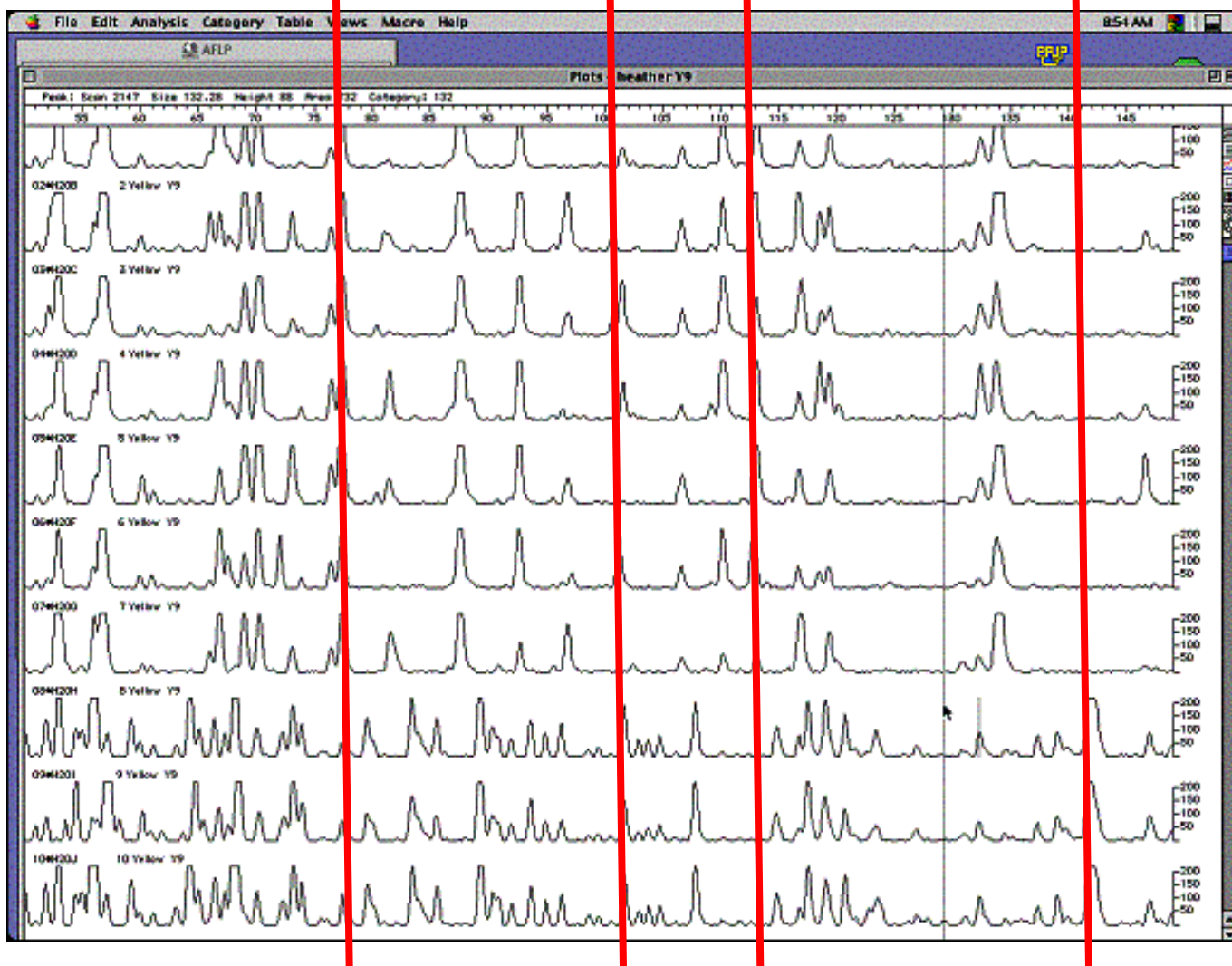
multi-locus genotype



„capillary version“

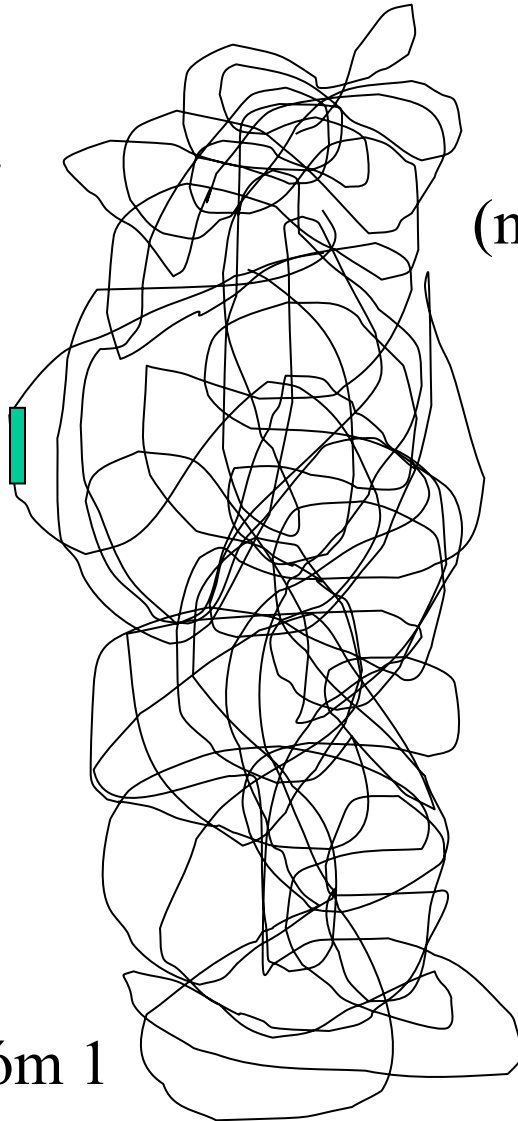
Ex.:
Combination
MseI + EcoRI

Automatizované čtení elektroforetogramu podle
zadaných kritérií (např. pozice a minimální výška píku)

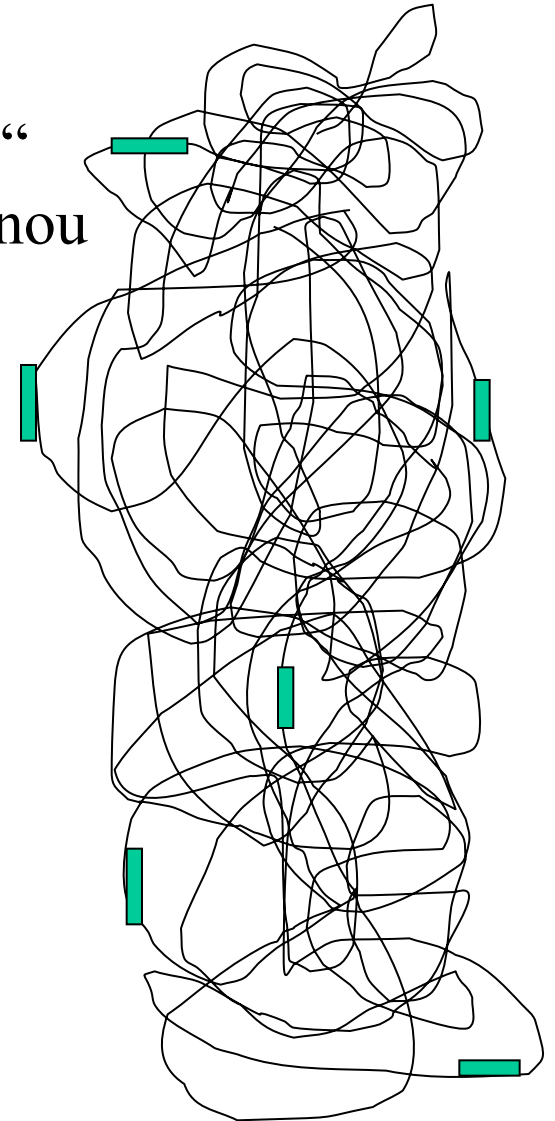


Typy genetických markerů

„single-locus“
(PCR)



„multi-locus“
(neznáme přesnou
lokalizaci
v genomu)



Př.: chromozóm 1

Budoucnost genetických metod v zoologickém výzkumu

1. Nové postupy při sekvenování DNA („genomics“)

Molecular Ecology Resources (2008) 8, 3–17

doi: 10.1111 /j.1471-8286.2007.02019.x

TECHNICAL REVIEW

**Sequencing breakthroughs for genomic ecology and
evolutionary biology**

MATTHEW E. HUDSON

Department of Crop Sciences, University of Illinois, Urbana, 334 NSRC, 1101 W. Peabody Blvd., IL 61801, USA

4-kapilární sekvenátor

=

96 x 500 bp/12 hodin

=

cca 100 000 bp/den

Next-generation sequencing

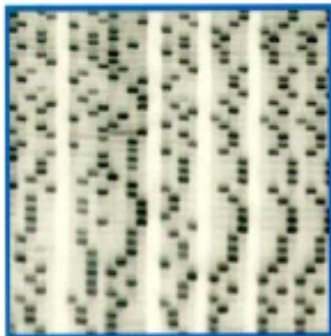
=

cca 1 000 000 000 bp/den

electrophoresis

Evolve Sangerova sekvenování

Pre-1992
“old fashioned
way”



S35 ddNTPs
Gels
Manual loading
Manual base calling

1992-1999
ABI 373/377



Fluorescent ddNTPs*
Gels
Manual loading
Automated base calling*

1999
ABI 3700



Fluorescent ddNTPs
Capillaries*
Robotic loading*
Automated base calling
Breaks down frequently

2003
ABI 3730XL



Fluorescent ddNTPs
Capillaries
Robotic loading
Automated base calling
Reliable*

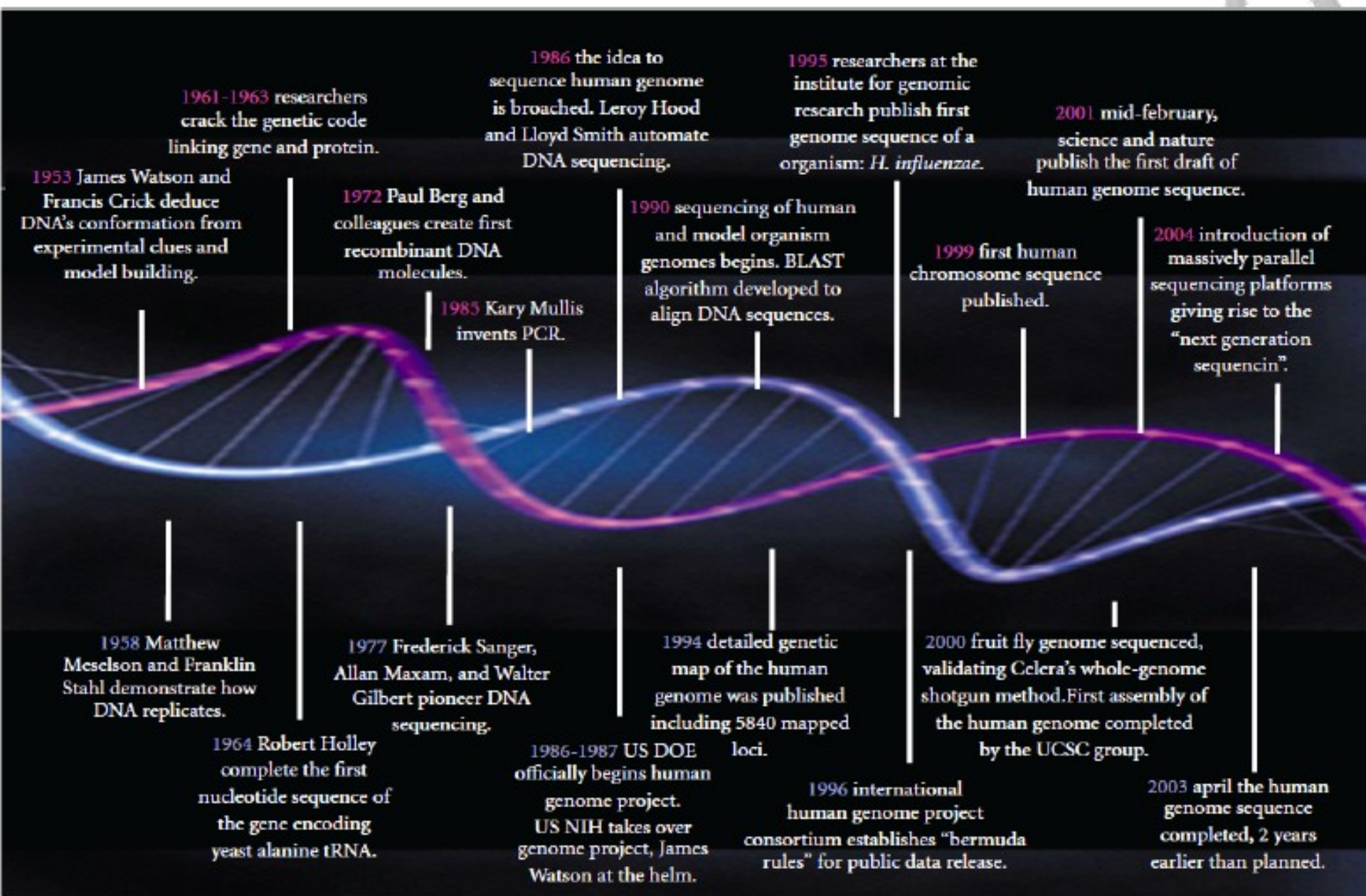
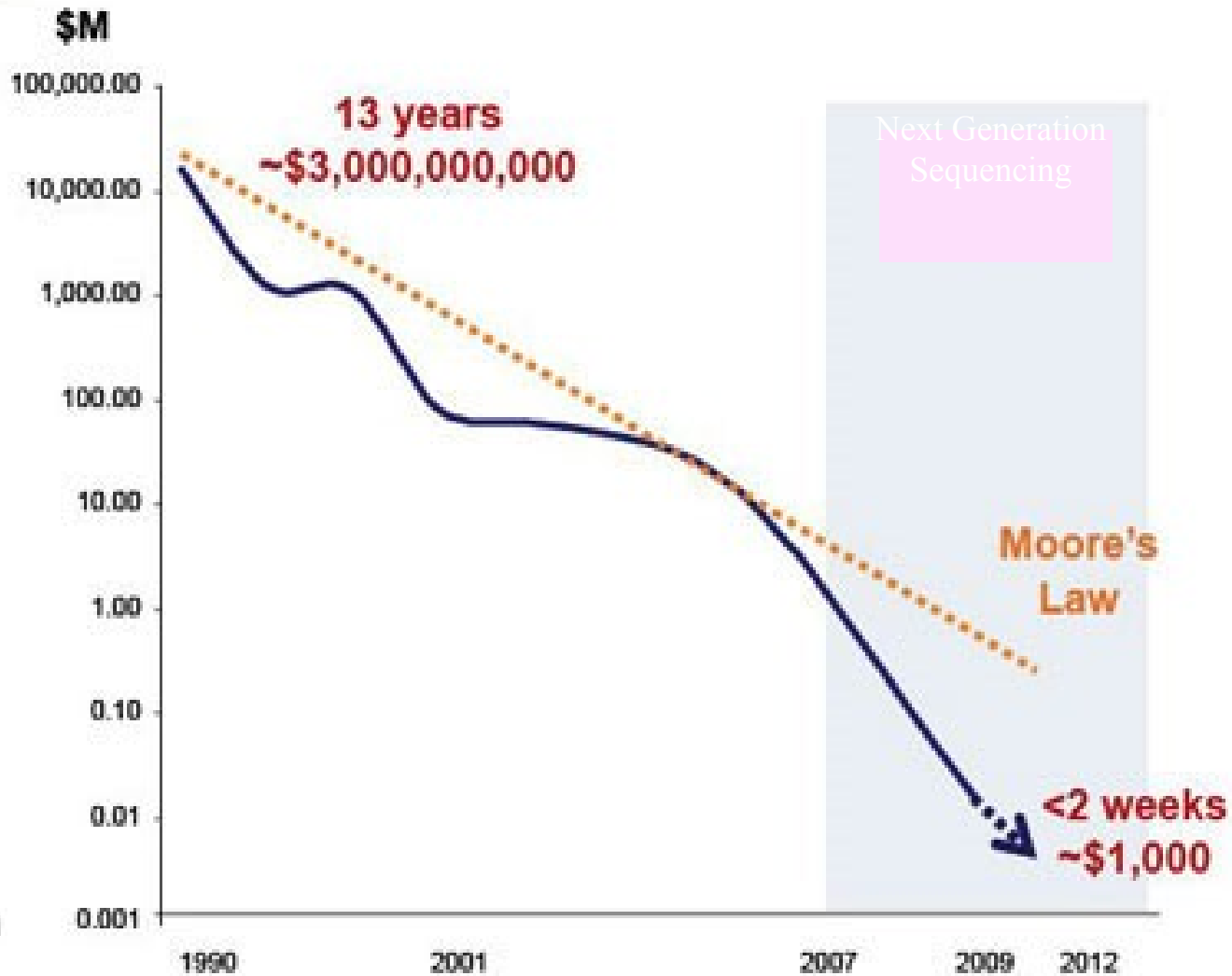


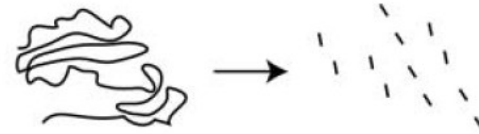
FIGURE 1: Evolution of DNA revolution.

Cost per Human Genome

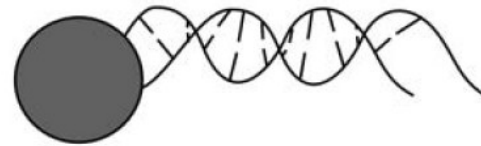


„Next generation sequencing“

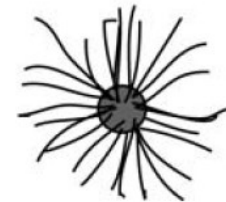
1) Randomly fragment many molecules of target DNA



2) Immobilize individual DNA molecules on solid support



3) Amplify DNA in clonal 'polymerase colony'



„polonies“
(polymerase colonies)

4) Sequence DNA by adding liquid reagents to immobilized DNA colonies



5) Interrogate sequence incorporation *in situ* after each cycle using fluorescence scanning or chemiluminescence



... commercially available since August 2007

Available next-generation sequencing platforms

- **Roche 454**
- **Illumina/Solexa**
- ABI SOLiD
- ABI IonTorrent
- Polonator
- HeliScope
- ...

454 pyrosequencing

- emulzní techniky amplifikace pikolitrové objemy
- simultánní sekvenování na destičce z optických vláken detekce pyrofosfátů uvolňovaných při inkorporaci bází
- První generace GS20 → 200 000 reakcí najednou (zhruba 20 milionů bp) dnes FLX → 400 000 reakcí najednou eukaryotní genom za týden!!!
- Délka jednotlivých sekvencí 100 - 400 (800 bp)



Molecular Ecology (2008) 17, 1629–1635

NEWS AND VIEWS

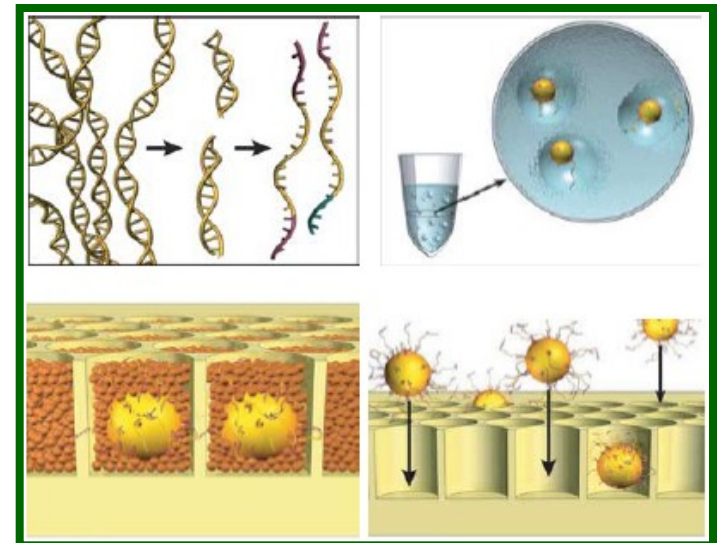
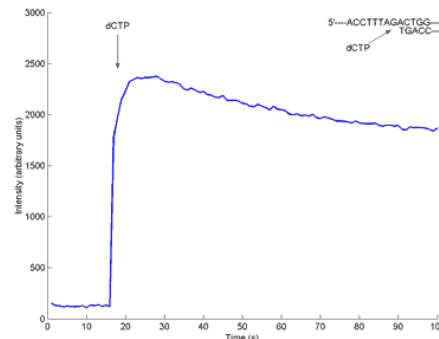
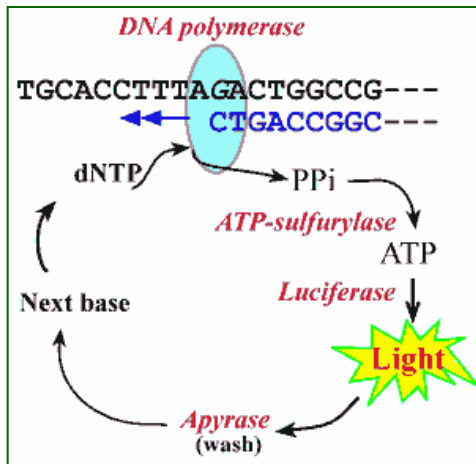
PERSPECTIVE

Sequencing goes 454 and takes large-scale genomics into the wild

HANS ELLEGREN

Department of Evolutionary Biology, Uppsala University, Norbyvägen 18D, SE-75236 Uppsala, Sweden

1 600 000 well plate



Pracovní postup



1

DNA Library Preparation

1. DNA Fragmentation (Nebulization)
2. DNA Fragment Size Selection
3. DNA Sample Quality Assessment (Nebulized or *LMW* DNA Sample)
4. Fragment End Polishing
5. Adaptor Ligation
6. Small Fragment Removal
7. Library Immobilization
8. Fill-In Reaction
9. Single-Stranded DNA Library Isolation
10. DNA Library Quality Assessment and Quantitation

Time: 11 - 72 h

General Laboratory 1



2

Emulsion-Based Clonal Amplification (emPCR)

1. Preparation of the Live and Mock Amplification Mixes
2. DNA Library Capture
3. Emulsification
4. Amplification
5. Bead Recovery
6. DNA Library Bead Enrichment
7. Sequencing Primer Annealing

Time: 11 - 13 h

Controlled Room

Amplicon Room



3

Sequencing / Genome Sequencer FLX Operation

1. The Pre-Wash
2. PicoTiterPlate Device Preparation
3. The Sequencing Run

Time: 11.5 h

General Laboratory 2



4

Data Processing and Analysis

1. Data Processing
 - a) Image Processing
 - b) Signal Processing
2. Data Analysis
 - a) Assembly
 - b) Mapping
 - c) Amplicon Variant Analysis

Time: variable



1. Příprava jednořetězcové DNA knihovny (ssDNA library preparation)

1 DNA Fragmentation (Nebulization):



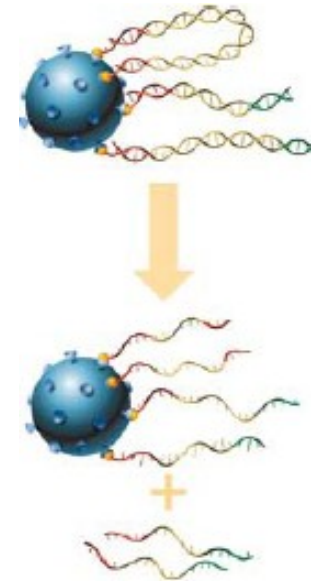
5 Adaptor Ligation:



7 Library Immobilization:



9 ssDNA Library Isolation:



Adaptor A + Adaptor B

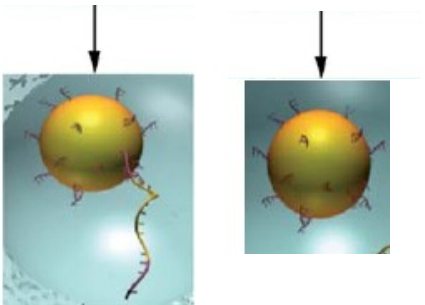
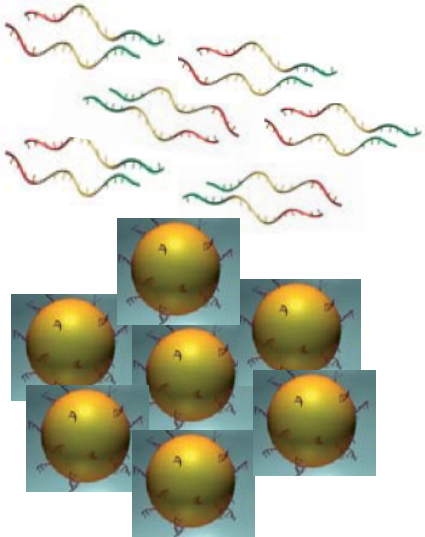
- Slouží jako vazebné místo primerů pro následnou PCR amplifikaci a sekvenování

- Slouží k uchycení na kuličky (na adaptor B je připojen **biotin**)

2. Namnožení každé jednotlivé molekuly pomocí emulzní PCR (emPCR)

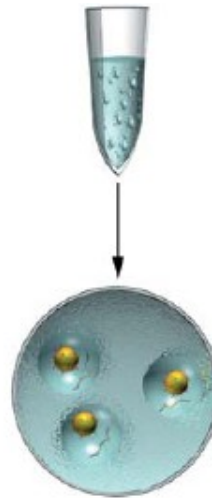
1 DNA Library Capture:

- poměry nastavit tak aby
1 kulička \leq 1 molekula DNA

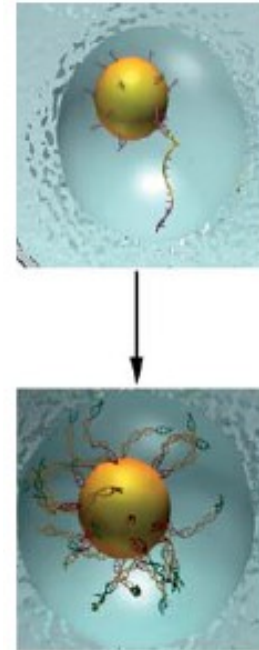


2 Preparation of the Amplific. Mixes

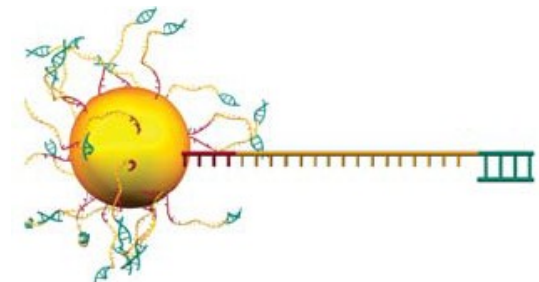
3 Emulsification:



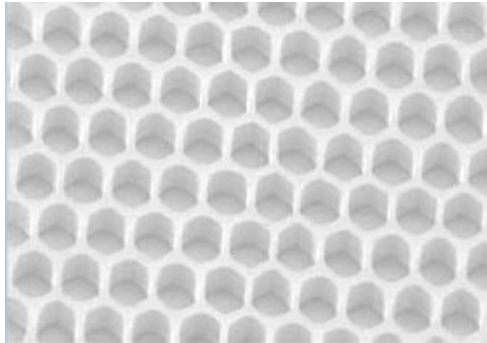
4 emPCR Amplification:



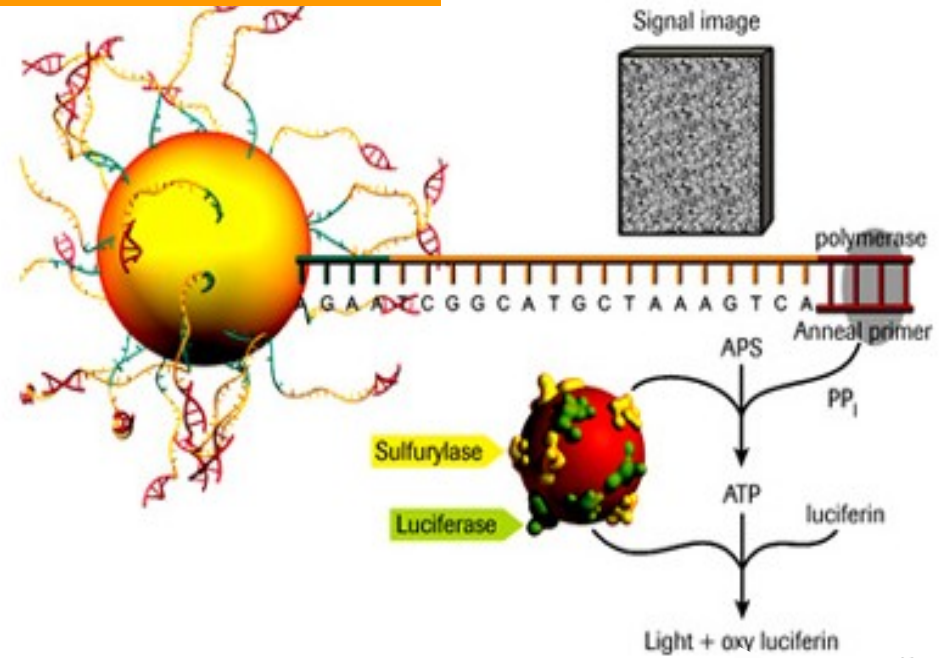
7 Sequencing Primer Annealing:



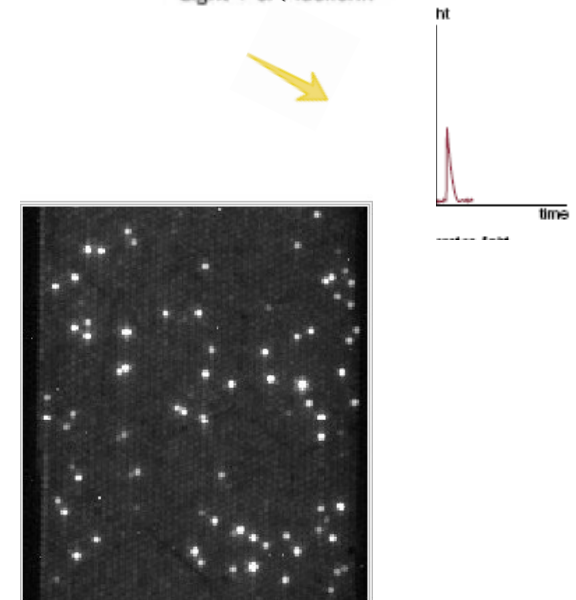
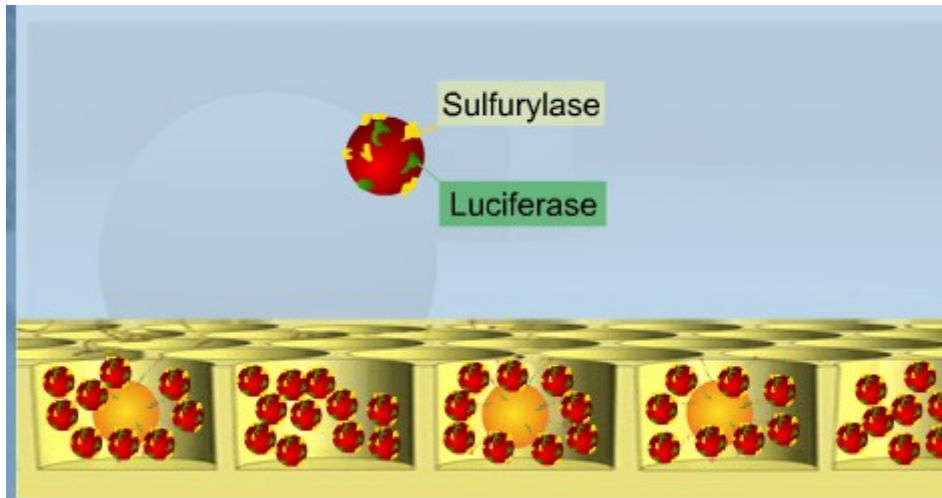
3. Pyrosekvenování („sequencing by synthesis“)



pikotitrační destička



Na jedné destičce 400 000 až 1milión jamek



3. Pyrosekvenovani - detekce signálu

- postupně se přidávají nukleotidy v definovaném pořadí: např. TACG TACG TACG
- po přidání každého nukleotidu a detekci signálu se nukleotid odmyje a přidá se další odmyje

DNA sekvence: **C T C C G**

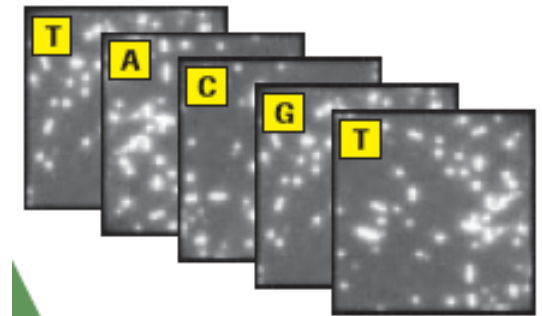
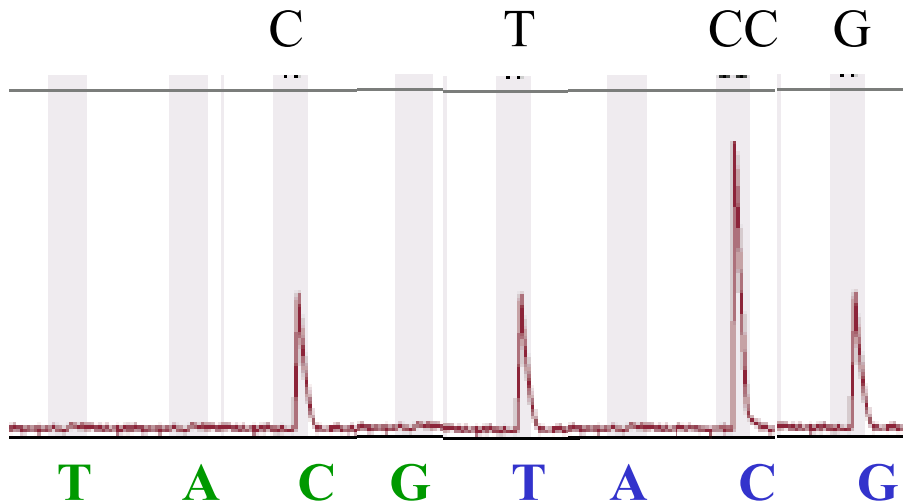


Image Files:
12-15 gigabytes
per run

Problém!!!! Homopolymery např. AAAAAAAAAA

High-throughput - paralelní sekvenování

1 běh (run) = 1 destička:

- 400 000 / 1milión jamek (reads)
- v každé 240 / 400 bp (read length)
- 7.5 / 10 hod

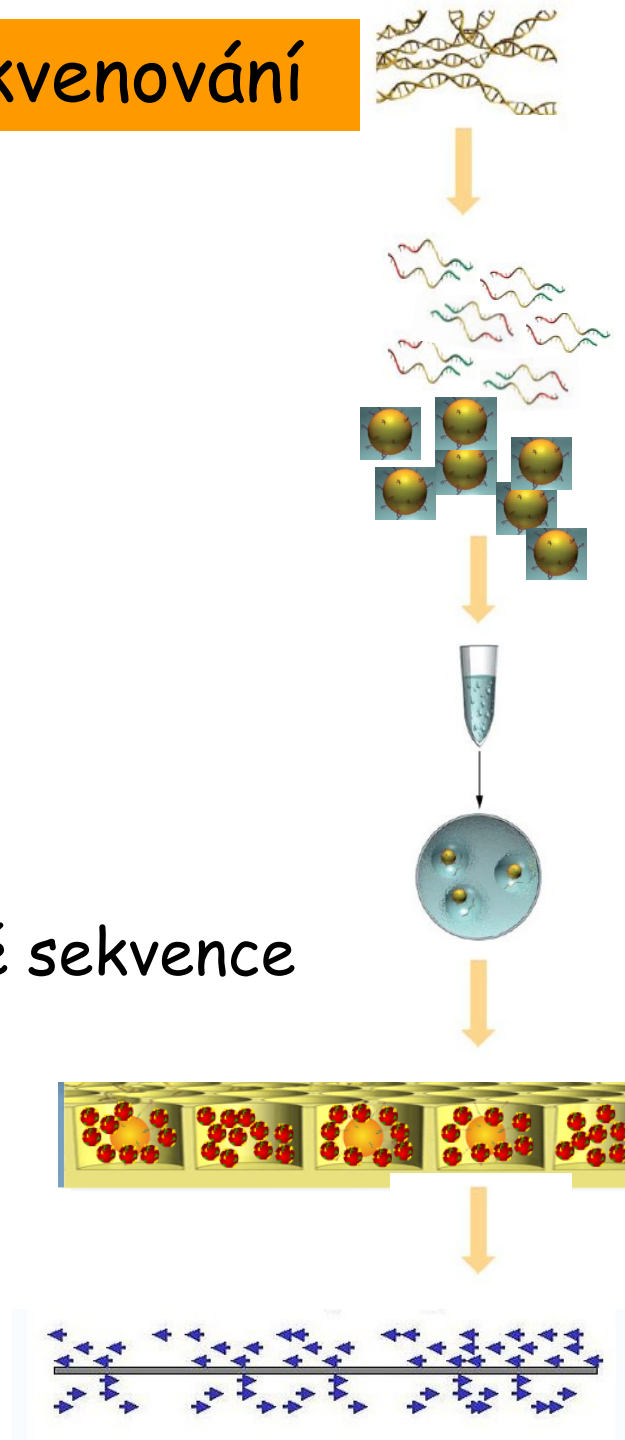
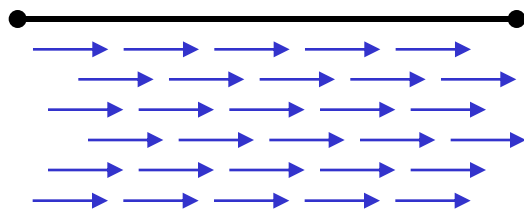
→ 100 Mb / 400 Mb na jednu destičku

→ cena??? 150-350 000 Kč ??? (verze Junior - od 15 tisíc)

!!! Samozřejmě nestačí mít každou bázi osekvenovanou 1x !!!

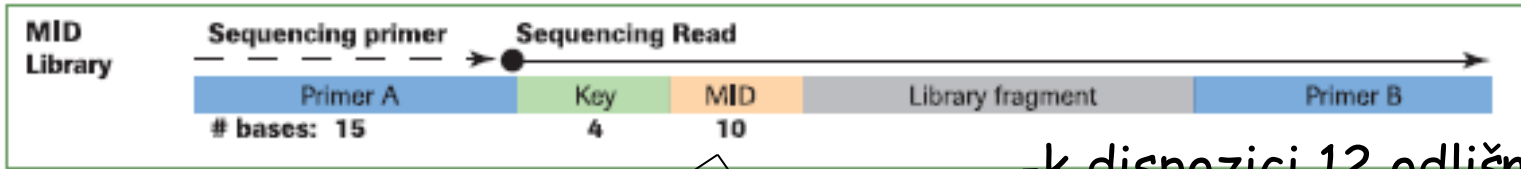
- Pospojování (**reads assembly**) do souvislé sekvence

- Nepřesnosti - pokrytí (**coverage**)

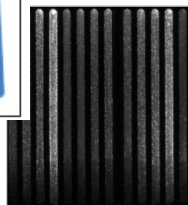
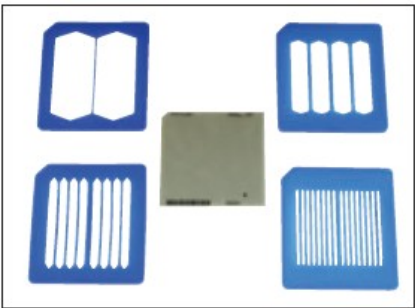
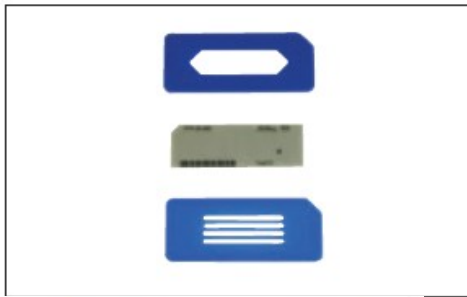


Kapacita destičky **400 Mb**:

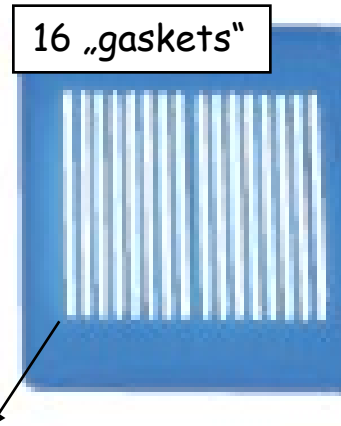
Mus:	2700 Mb	→ 7 run 1x coverage
Caenorhabditis:	100 Mb	→ 1 run 4x coverage
E. coli:	5 Mb	→ 1 run 80x coverage
mitoch. Mus:	0.016 Mb	→ 1 run 25000x coverage
HIV:	0.01 Mb	→ 1 run 40000x coverage



-k dispozici 12 odlišných MID („multiplexing“)



1. CCCCCCCCCC
2. GGGGGGGGGG
- ...
12. CCCCAAAG



$$\begin{array}{r} 12 \text{ MID} \\ \times \\ 16 \text{ gaskets} \\ = \\ \text{max. 192 vzorků} \end{array}$$

V každém max. 12 vzorků
(každý označen svým MID)



454 Genome Sequencers

FLX System

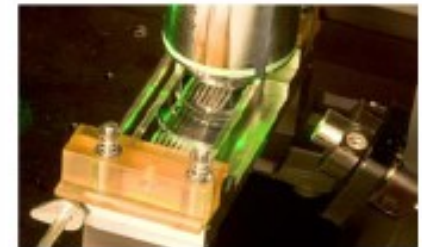
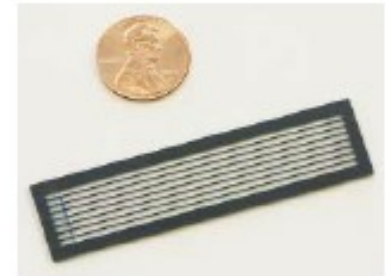
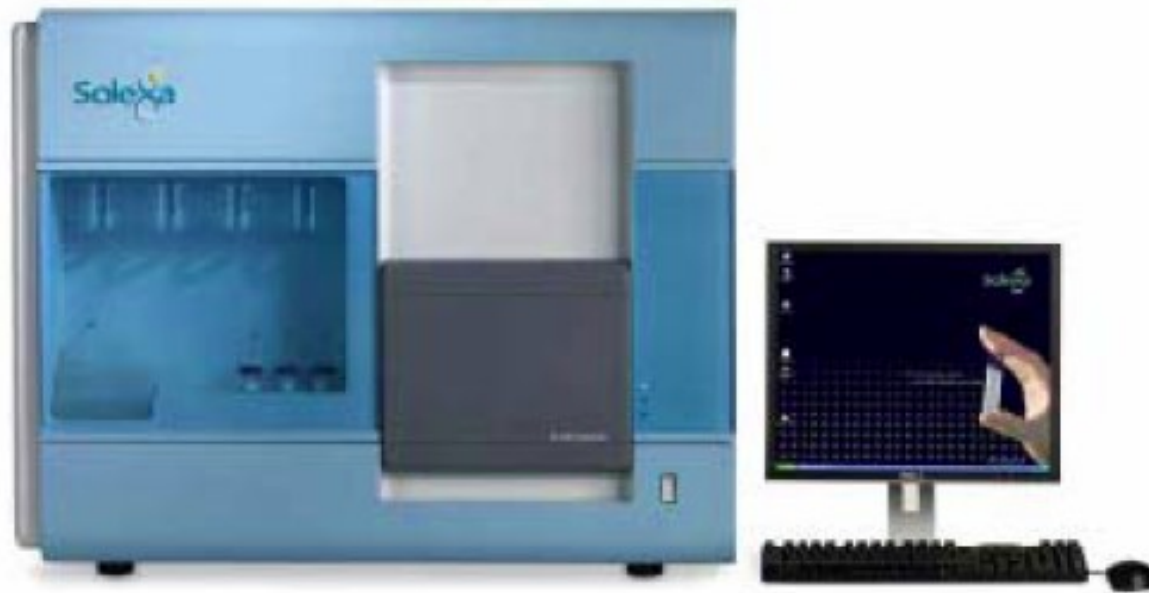
- 1 million of reads/run
- 400-650 bp/read
- 2 přístroje v ČR



GS Junior

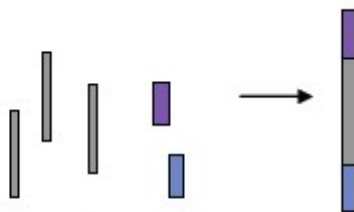
- 0.1 millions of reads/run
- 400 bp/read

Illumina Genome Analyzer Introduction to the Technology



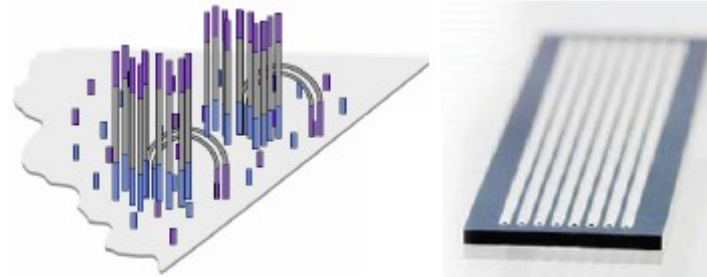
Illumina Sequencing pipeline

1. Sample Prep (1-5 days)



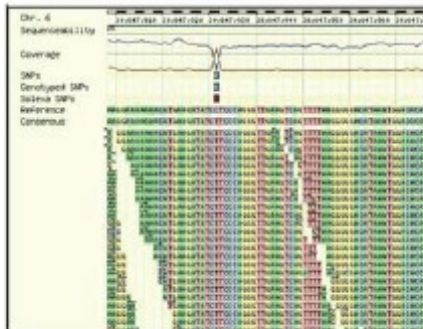
Ligate adapters

2. Cluster generation on flow cell (1.5 day)



Clonal Single molecular Array

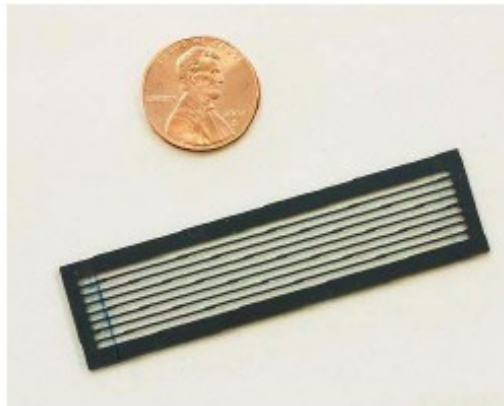
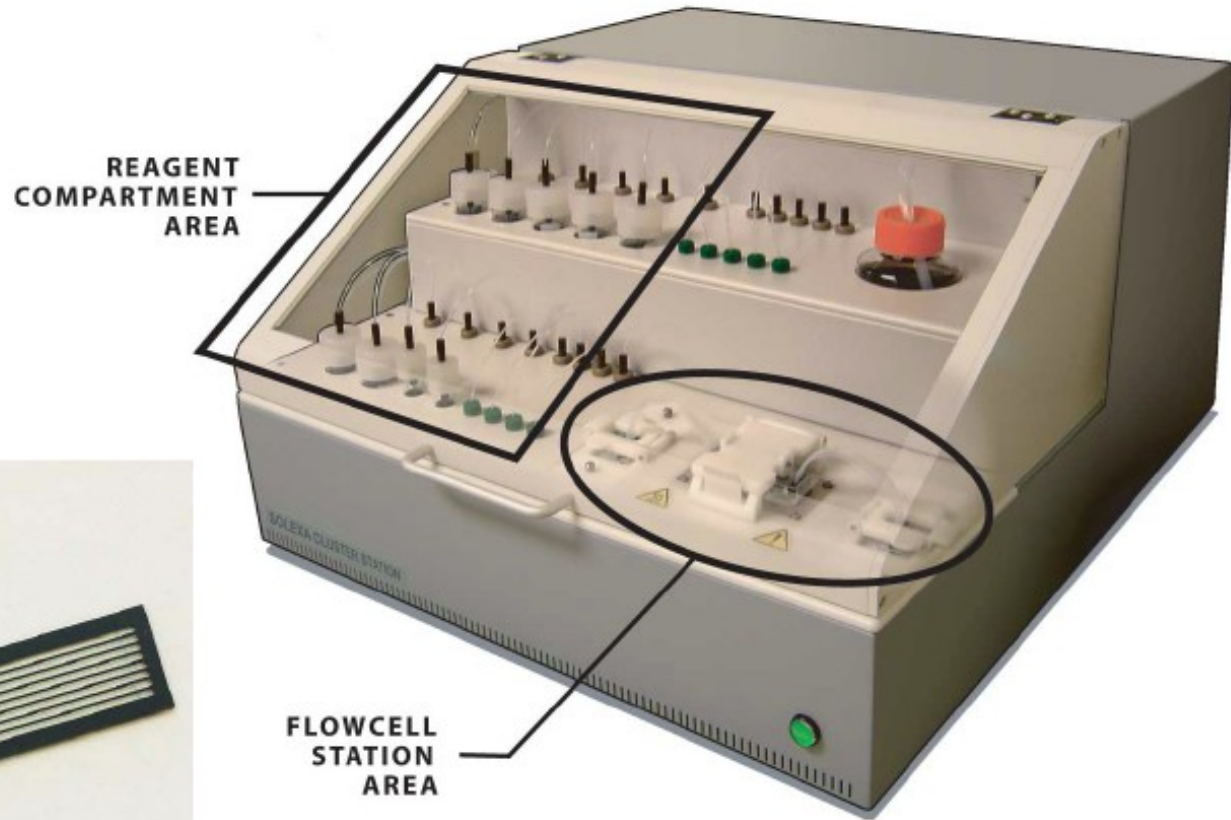
4. Data Analysis (days-months)



3. Sequencing and imaging (2-3 days)

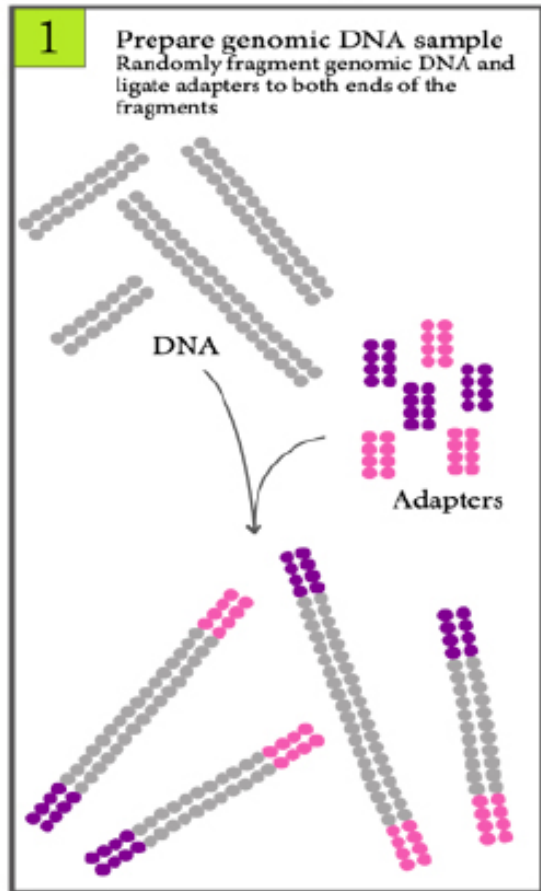


Cluster Generation

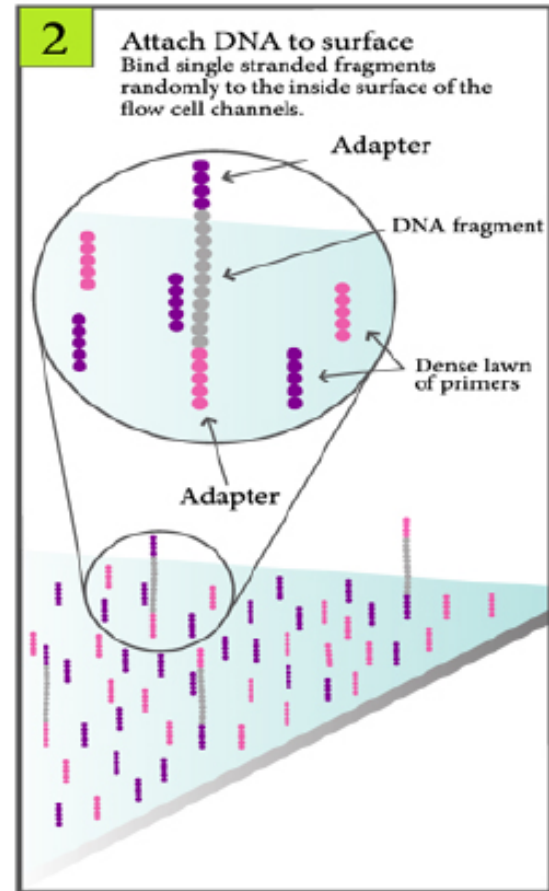


8 channels (lanes)

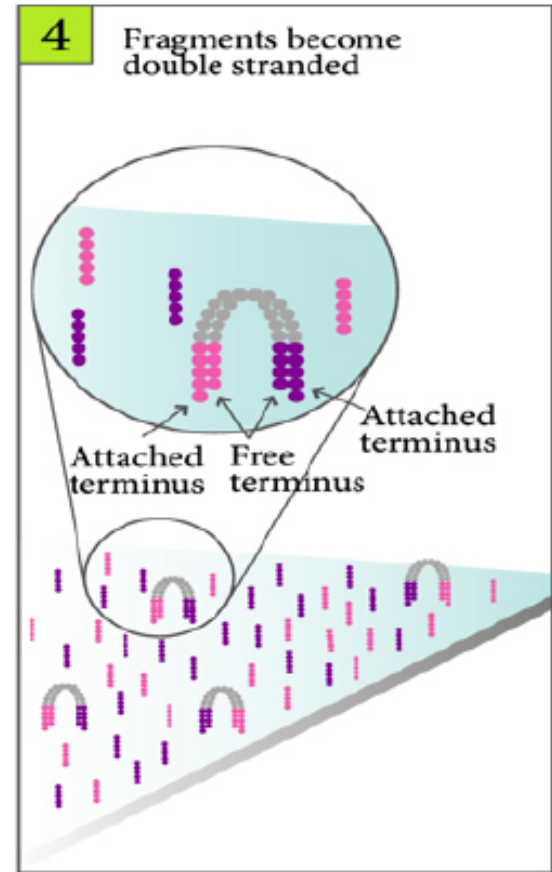
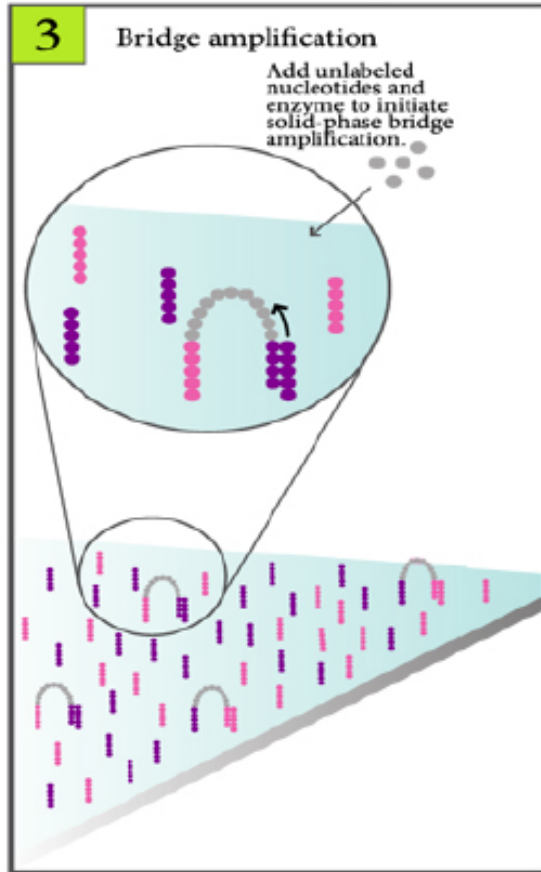
Attach DNA to flow cell



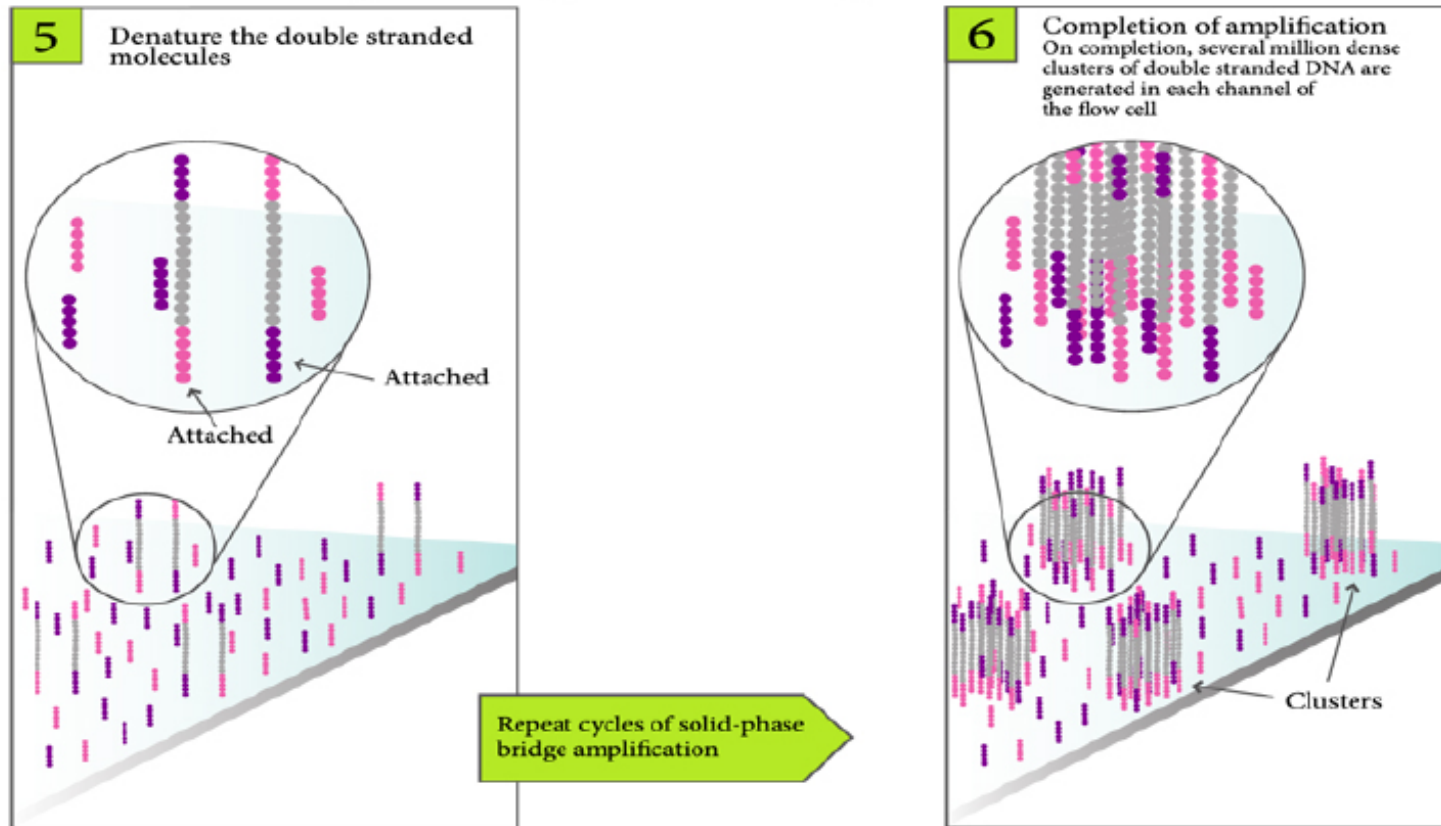
Add sample to flow cell



Bridge Amplification

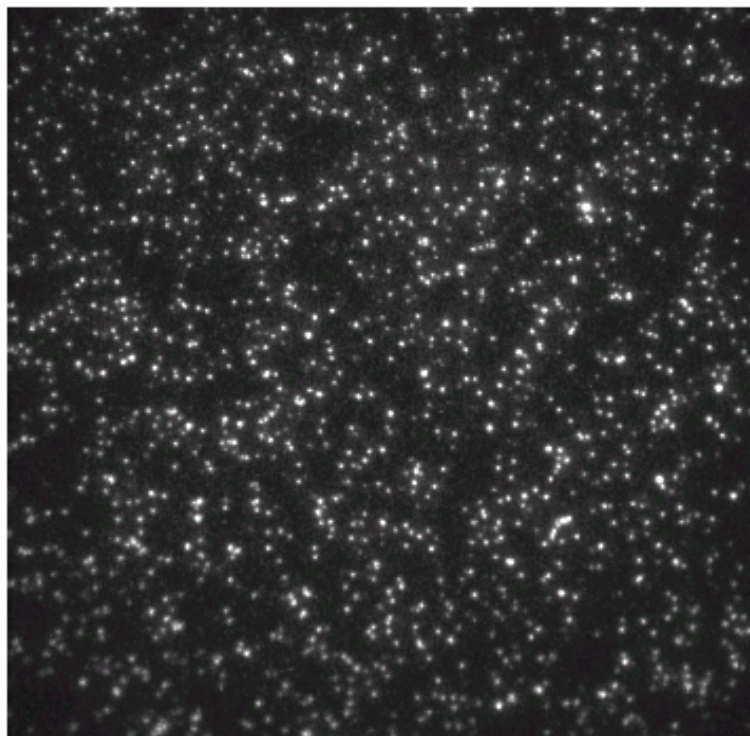


Cluster Generation



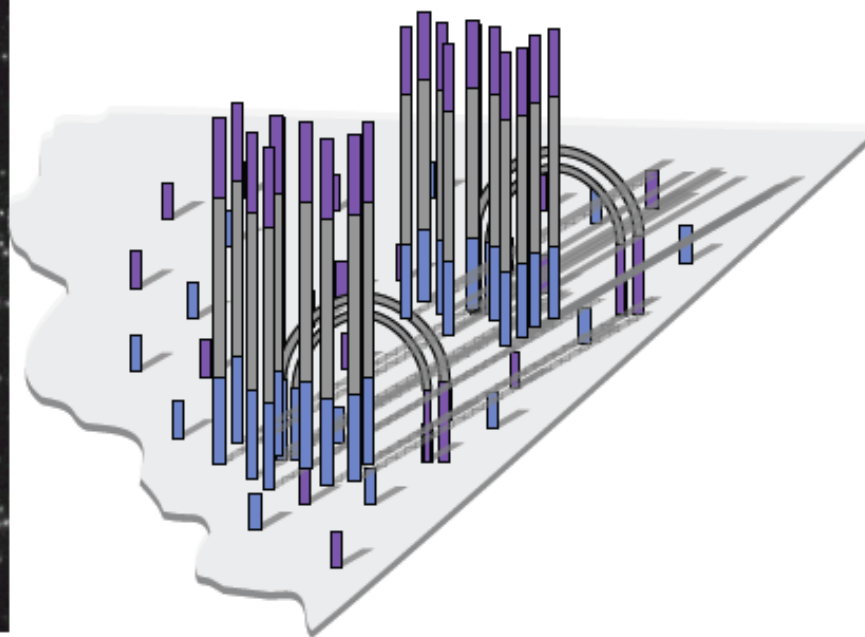
Clonal Single molecular Array

Clonal Single molecule Array



100um

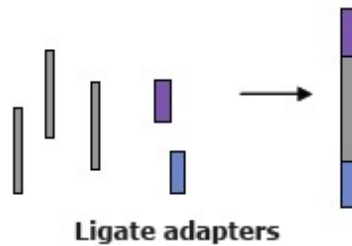
Random array of clusters



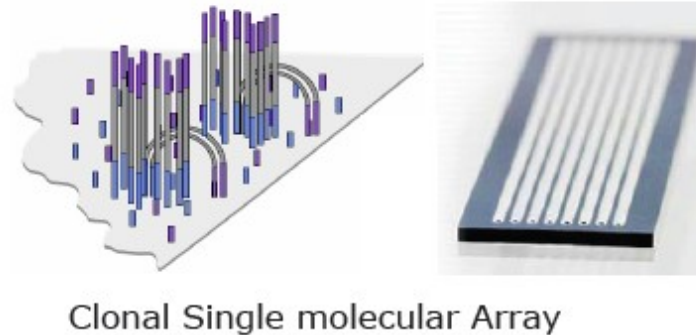
~1000 molecules per ~ 1 um cluster
~20-30,000 clusters per tile
~40 M clusters per flowcell

Illumina Sequencing pipeline

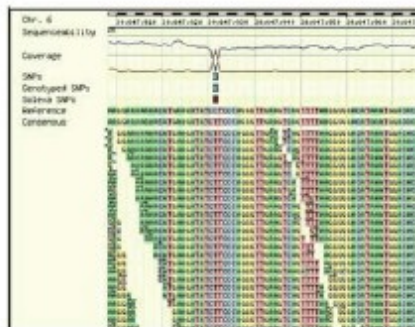
1. Sample Prep (1-5 days)



2. Cluster generation on flow cell (1.5 day)



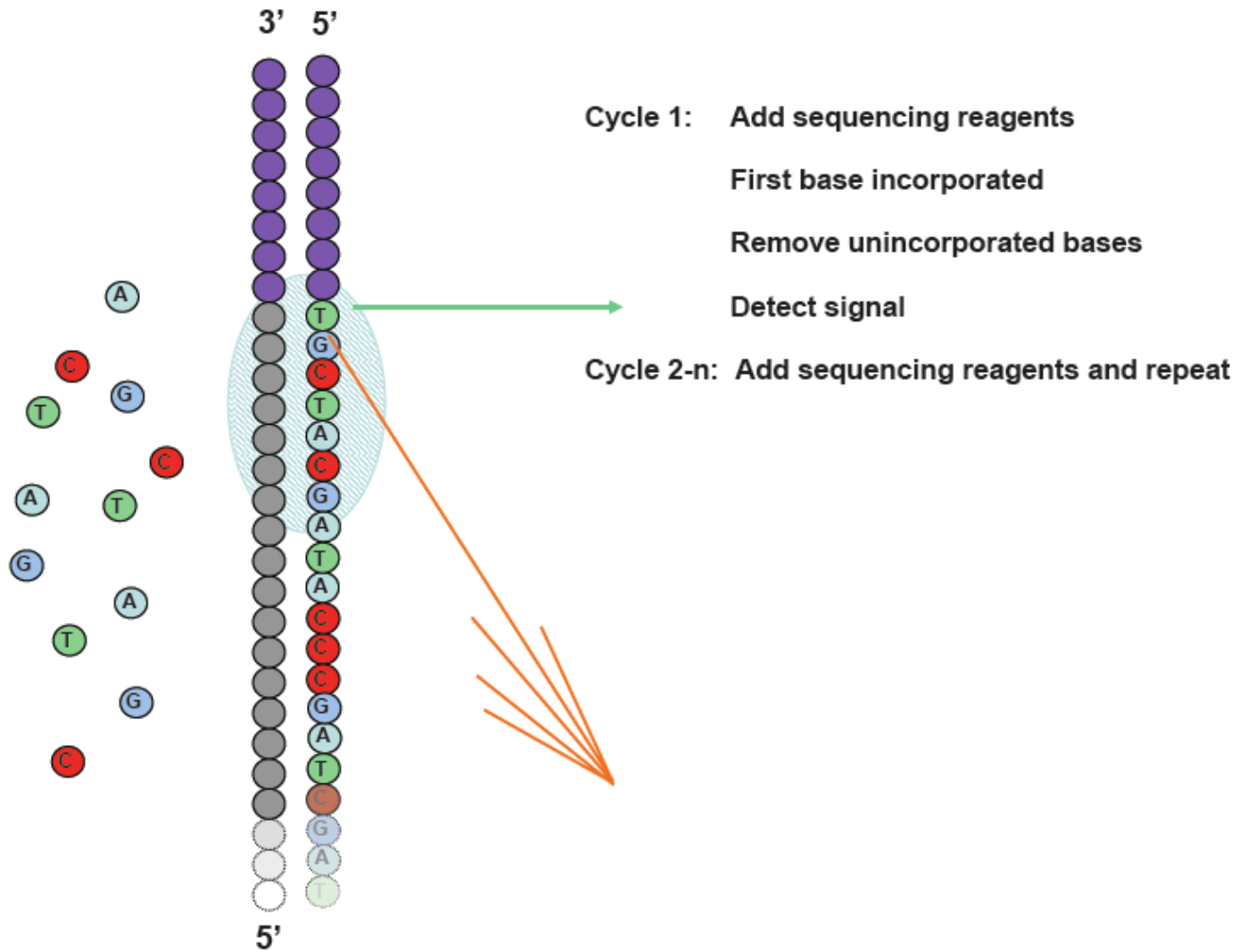
4. Data Analysis (days-months)



3. Sequencing and imaging (2-3 days)



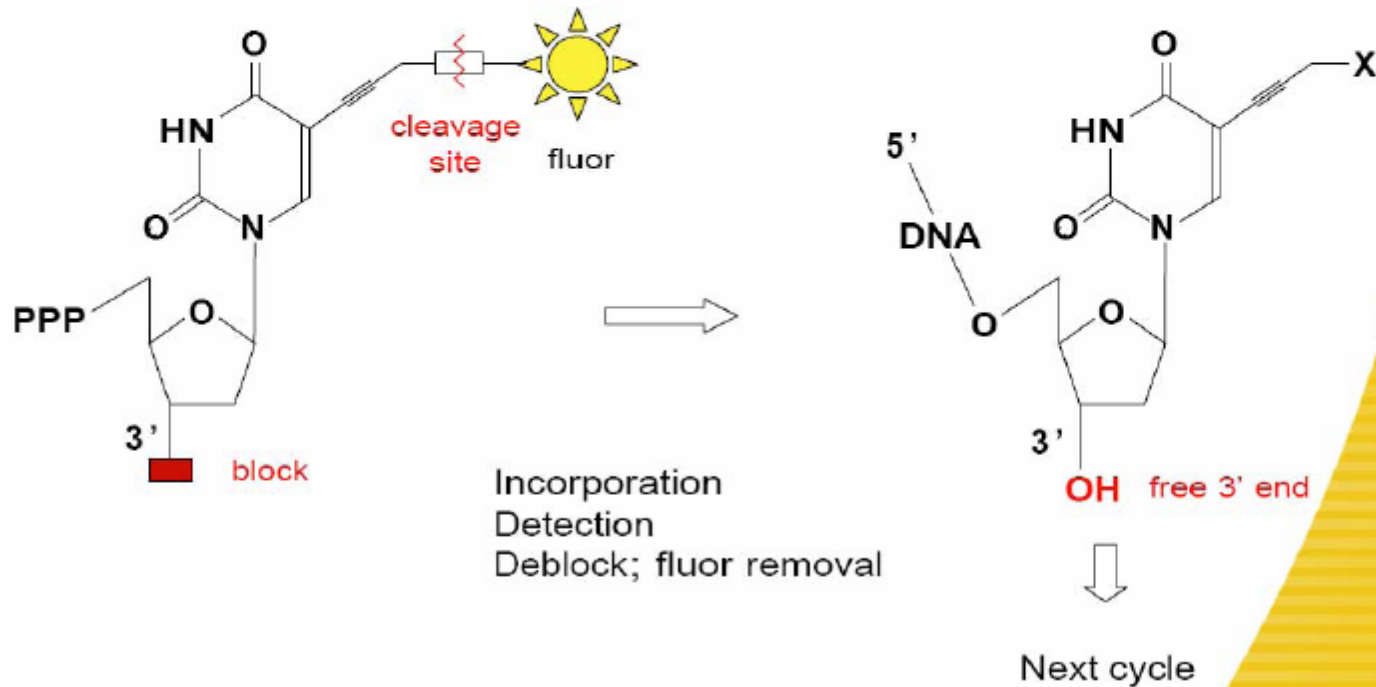
Sequencing By Synthesis (SBS)



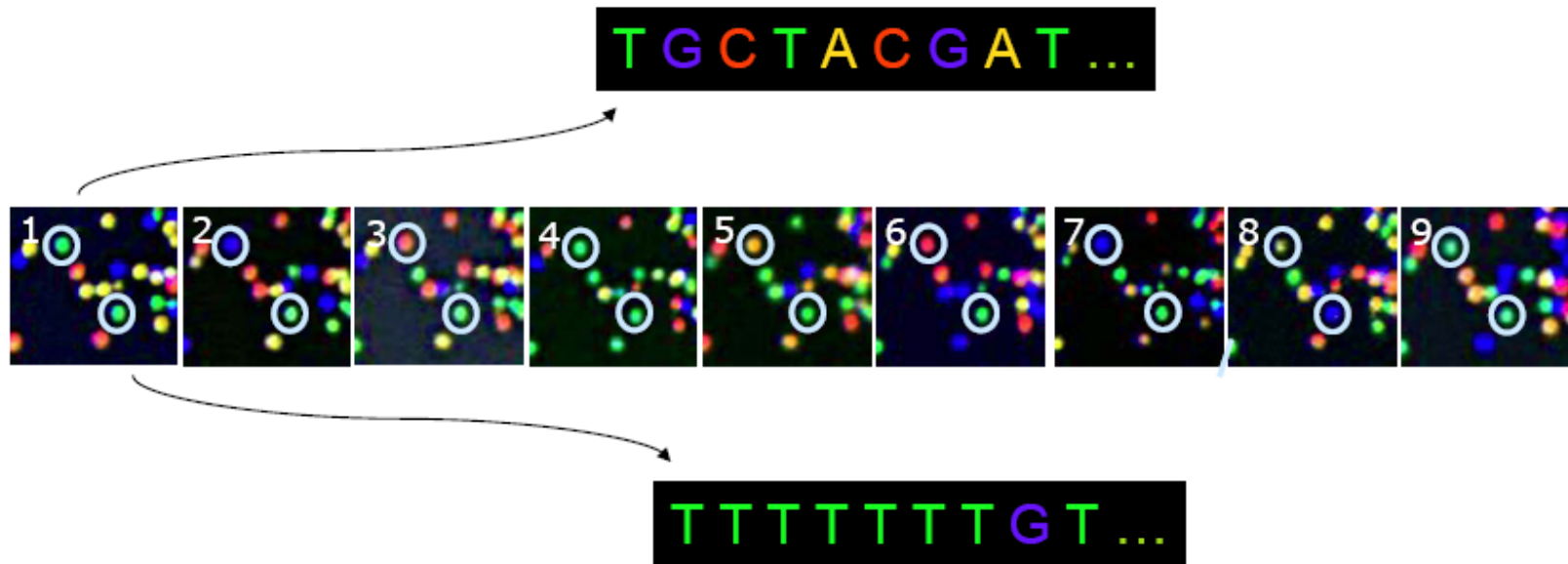
Reversible Terminator Chemistry



- All 4 labelled nucleotides in 1 reaction
- Higher accuracy
- No problems with homopolymer repeats

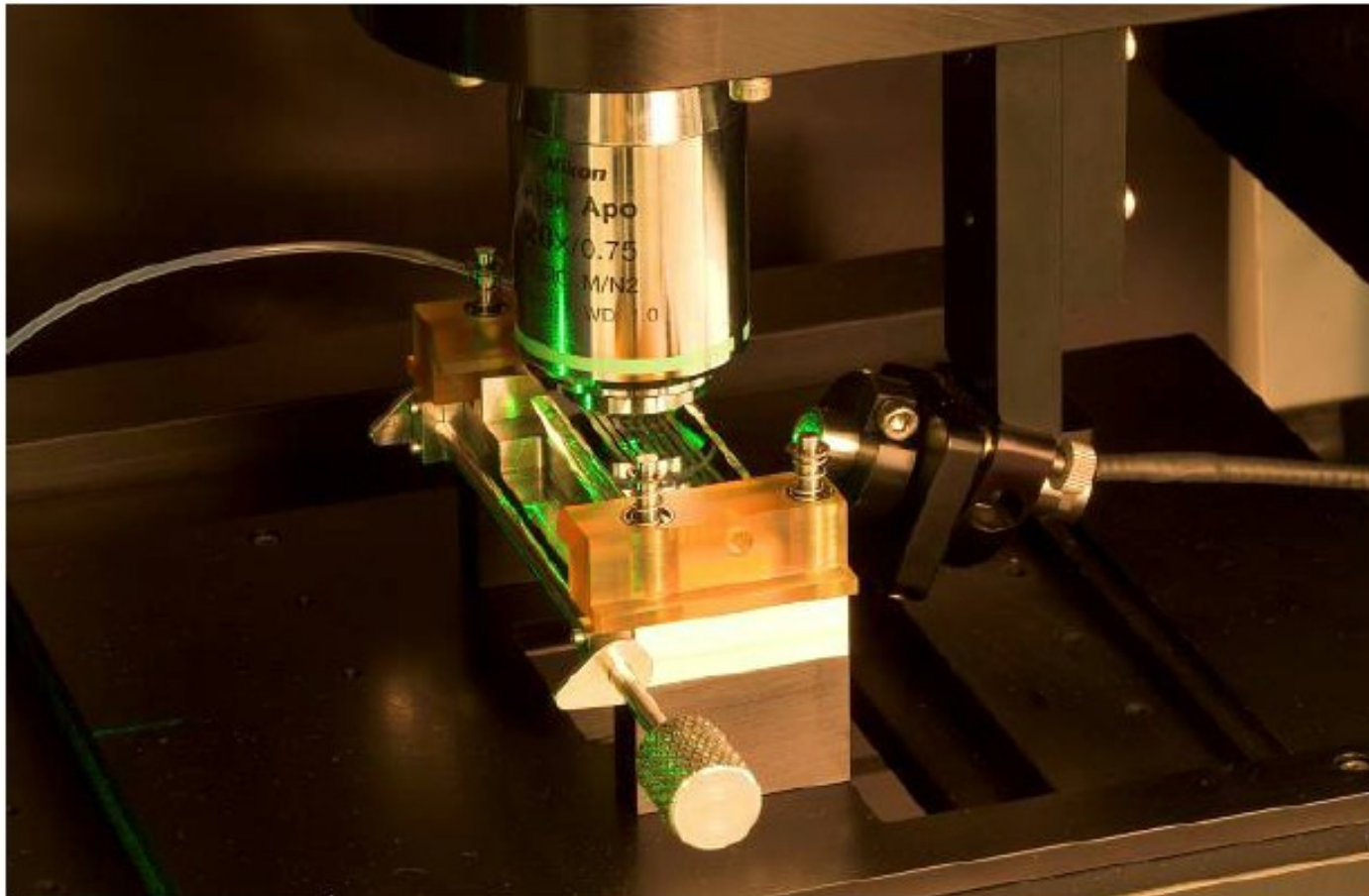


Base Calling From Images



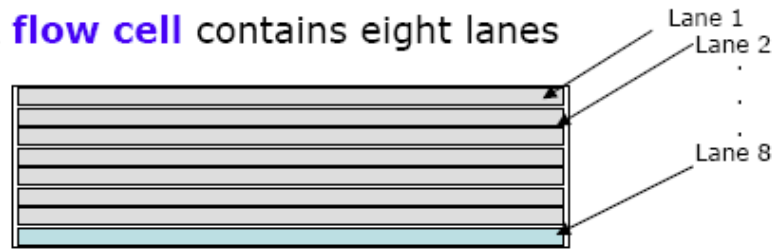
The identity of each base of a cluster is read off from sequential images

Flowcell imaging

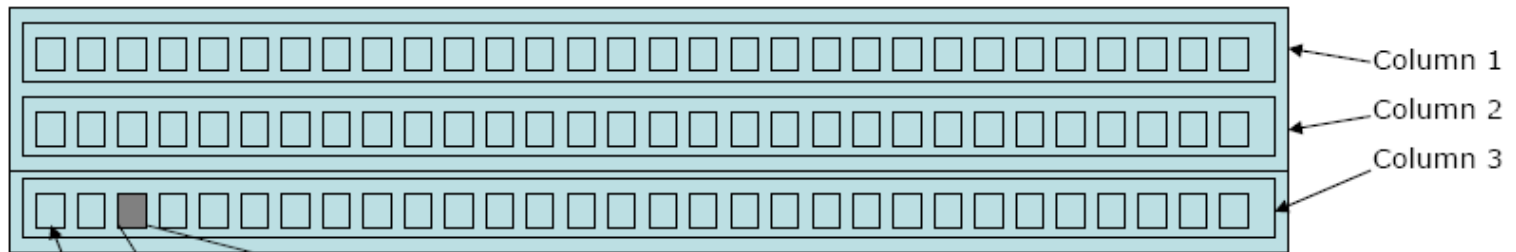




A **flow cell** contains eight lanes



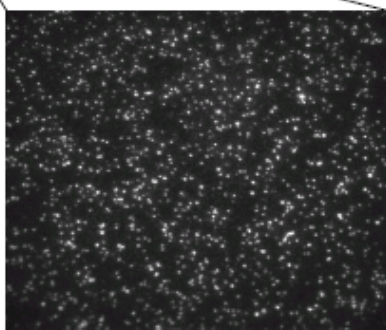
Each **lane/channel** contains **three columns** of tiles



Each **column** contains **100 tiles**

Tile

20K-30K
Clusters



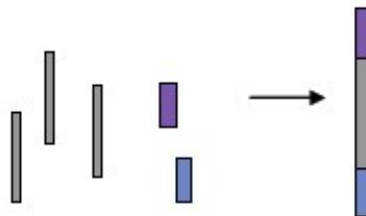
350 X 350 μm

Each tile is imaged four times per cycle – one image per base.

345,600 images for a 36-cycle run

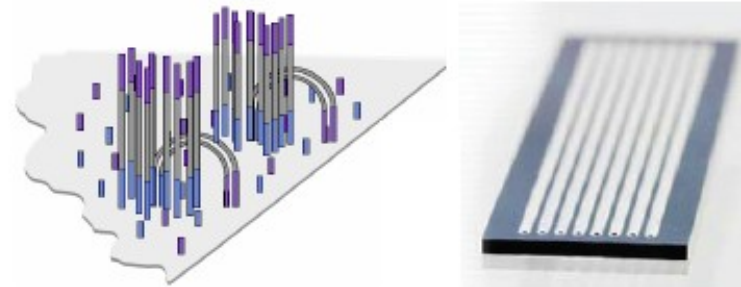
Illumina Sequencing pipeline

1. Sample Prep (1-5 days)



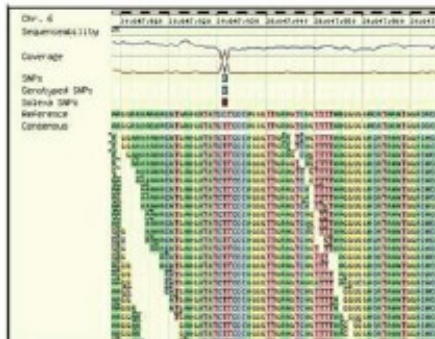
Ligate adapters

2. Cluster generation on flow cell (1.5 day)



Clonal Single molecular Array

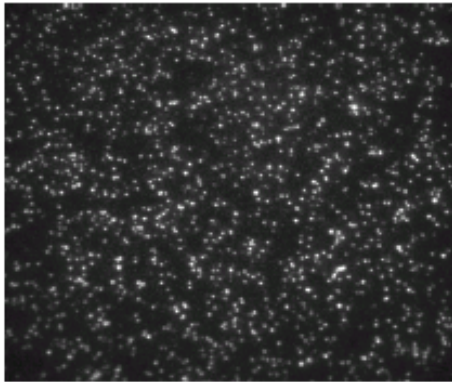
4. Data Analysis (days-months)



3. Sequencing and imaging (2-3 days)



Data Analysis Pipeline



tiff image files
(345,600)

Firecrest

1	T	130	543	140.0	347.7	739.1	24046.0	202.2	209.7	297.0	2104.4
1	T	180	421	231.0	341.9	497.7	21423.8	229.3	380.8	14319.2	20217.9
1	T	240	426	216.4	356.0	501.6	21362.3	345.5	319.7	467.9	19749.5
1	T	241	509	187.7	382.7	597.4	20747.7	1489.2	1034.1	161.0	482.7
1	T	224	285	178.5	372.1	486.5	20302.6	8297.1	12746.0	159.4	286.8
1	T	155	544	170.2	339.5	530.3	18408.9	307.6	418.8	364.9	17172.9
1	T	301	307	355.8	472.1	782.0	20449.1	1891.2	12332.1	191.9	743.0
1	T	175	406	210.4	323.8	522.3	16249.2	544.4	208.7	535.9	20587.5
1	T	240	522	287.9	533.0	456.0	15096.7	4285.6	10442.1	3394.7	2486.9
1	T	194	522	220.2	455.9	486.6	18895.6	189.5	152.8	12299.4	14131.7
1	T	237	422	147.6	457.7	521.0	16025.2	712.0	990.0	416.4	10774.0
1	T	160	526	170.4	400.7	481.9	14486.9	1249.7	4305.8	241.3	524.1
1	T	164	549	205.7	385.0	480.4	13465.5	2410.3	9408.2	76.7	243.0
1	T	179	381	207.2	372.3	562.1	10442.2	240.7	282.3	314.4	16462.8
1	T	224	423	216.3	460.4	474.4	18360.9	1331.1	10764.6	159.2	446.3
1	T	139	583	241.0	358.9	542.7	18183.9	226.9	302.0	13425.1	15107.5
1	T	220	428	225.1	486.8	553.2	15716.8	3338.0	10291.0	311.3	594.4
1	T	300	307	194.0	329.0	460.3	24628.4	294.7	590.4	403.0	16946.9
1	T	334	512	249.8	599.6	430.9	24101.4	4787.9	11274.9	602.5	177.3
1	T	150	327	216.7	349.4	536.6	17715.4	2413.2	9446.9	377.4	523.2
1	T	243	541	182.5	375.9	470.2	22603.1	4711.0	11481.7	139.5	604.9
1	T	241	408	206.4	341.2	497.0	17248.9	4030.2	9318.9	112.1	34.4
1	T	174	509	226.3	328.4	457.9	17172.1	179.5	306.5	387.3	14274.9
1	T	371	582	230.4	546.4	426.1	21245.9	4630.4	10982.2	146.3	216.1
1	T	271	608	176.8	391.5	487.5	21381.2	1832.2	11093.9	191.9	409.8
1	T	195	303	236.4	389.5	465.4	14629.3	4094.2	8305.9	289.5	3794.0
1	T	301	392	181.8	378.0	553.4	22549.7	8013.1	13222.2	899.6	1211.8
1	T	249	549	197.7	525.1	543.4	14512.2	1640.8	10451.3	171.3	504.9
1	T	140	517	108.7	388.0	510.1	14448.1	1755.0	8400.2	155.7	381.8

intensity files

Bustard

1	T	130	543	TTTGAACAGCATATTATAGCGACG
1	T	180	421	TGTTTTTTTTTTTTTTTGGACAGG
1	T	240	426	TTTGATCTGTPTTCTGCTGCGAGG
1	T	241	509	TCTGCTGCTGCTGCTGCTGCTGCT
1	T	214	595	TACAAAATCCCTGCCCATATGGACT
1	T	135	544	TTATCTGCATCCGATGCAATTTTAG
1	T	301	507	TCCTGCTTATTTGCTCTTTTJATTT
1	T	175	604	TTGGATCCGGGTAAAGGGAGAGAT
1	T	242	522	TACTAATATACAGATATGTTGAAA
1	T	196	522	TGTGACGGAGGGACGGCTGACAT
1	T	237	612	TTGCTGCTGCTGCTGCTGCTGCT
1	T	160	528	TCTGATTTTTCACAGTAAAGAAAAC
1	T	164	543	TCTGAGAAACCTGCTGCTGCTGCT
1	T	179	581	TCTGAAATCTTGCATGCTGCTGCT
1	T	224	623	TATTAGAGGCTGAGCGCTGCTGCTG
1	T	129	583	TTATGGATGGGAGCGAGGGAGGCT
1	T	220	418	TCCAAAATGTTTAAATATAGAGGCA
1	T	340	507	TTATTTGAGATTAATGTTTCCAAAT
1	T	334	512	TTATTTGTTTCCACTAATGGGAGTC
1	T	155	517	TCCCAAAAGAAAAAGAGAGAGGAG
1	T	343	541	TATTTCCATGCTGCTAATGATAGAT
1	T	241	608	TATTAGCCAGGTGAGTGGTGTACACC
1	T	174	520	TTTTTTAGTAGAGTGGGATTTACACC
1	T	371	592	TATTCCTATAGAAACAGCCATAGGG
1	T	271	508	TCTCTGGAAATATAGCTTACCGAG
1	T	195	603	TACTGAGTGGGGCCCTGGTACTCTG
1	T	501	710	AAAAAAAAAAAAAAAAAAAAAAAAAAAA

Sequence files

Additional
Data Analysis

Alignment to Genome

Eland

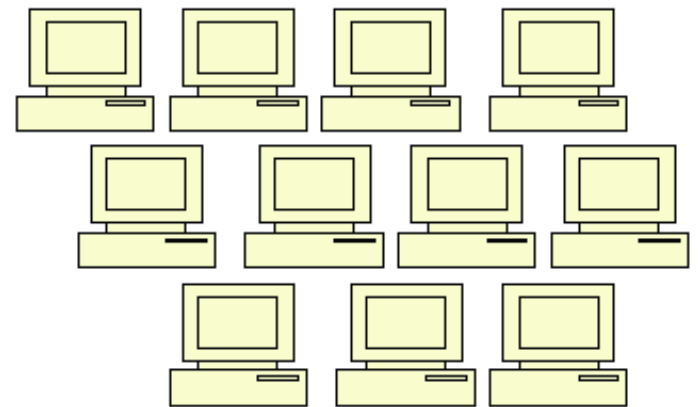


Data Analysis Pipeline



→
rsync

Computer Cluster



Imaging
Base calling
Sequence alignment

(10-12hrs)

All this generates a lot of Data!

1.5 TB data/run

- 1 Gig of Space
 - 125,000 pages of text
 - 11 CDs of Music
 - 4000 (1024x768) JPEG images
 - 40,000 pages of PDF
- 1 TB of space
 - 220 Million pages of text
 - 300 hours of video
 - 4,000,000 JPEG images
 - 1,000 copies of the Encyclopedia Britannica
 - 1/10 of the printed Library of Congress

Illumina sequencers

Illumina MiSeq

4 millions reads/run
150 bp/read



Illumina GAIx

300 millions reads/run
150 bp/read

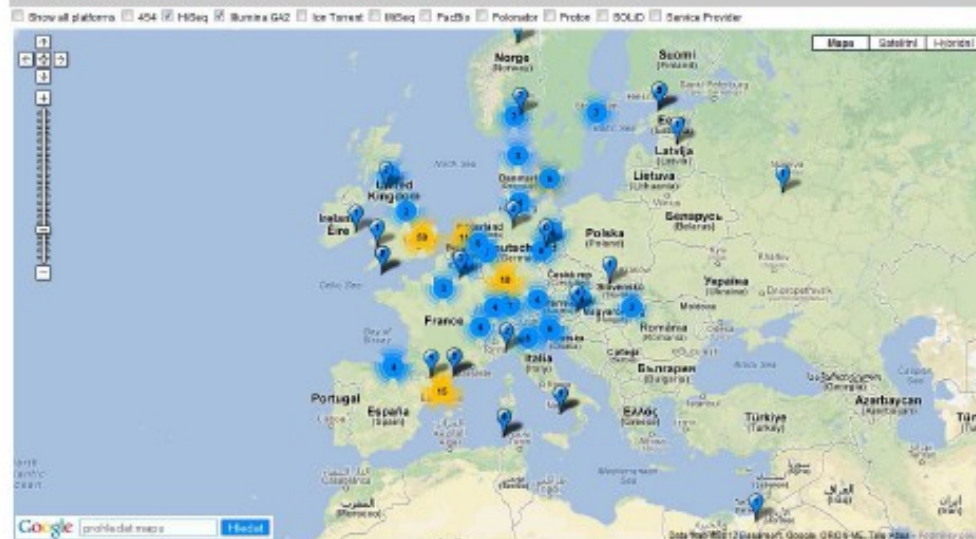


Illumina HighSeq

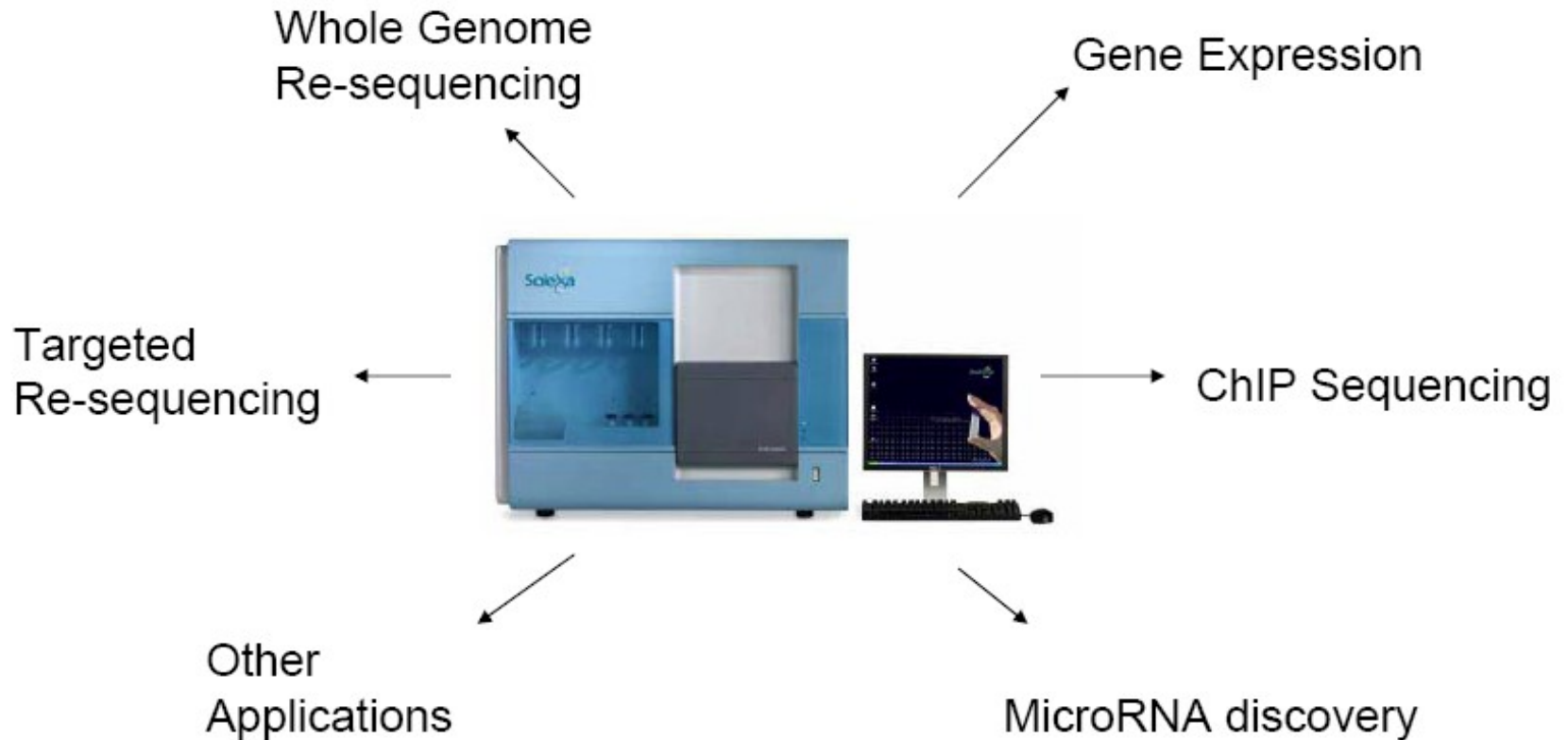
1500 – 3000 millions reads/run
100 bp/read



Next Generation Genomics: World Map of High-throughput Sequencers



Applications of the Technology



SOLiD

(sequencing by Oligonucleotide Ligation and



... a další (každého půlroku nová technologie - bouřlivý rozvoj !!!)

Přehled současných metod NGS

Platform	Year	Sequencing Method	Amplification	Detection	Features
454	2005	Pyro-sequencing	Emulsion PCR	Light	First NGS
Illumina	2007	Synthesis	Bridge PCR	Light	90% of Market
SOLiD	2008	Ligation	Emulsion PCR	Light	Lowest Error Rate
Ion Torrent	2010	Synthesis	Emulsion PCR	Hydrogen Ion	Semiconductor Chip
Pacific Biosciences	2010	Synthesis	None = Single Molecule	Light	Anchored Polymerases
Oxford Nanopore	2012	Nanopore	None = Single Molecule	Electrical Conductivity	"Run Until" Sequencing

Výkonnost jednotlivých metod

Instrument	Run time	Millions of Reads/run	Bases / read	Yield MB/run
3730xl (capillary)	2 hrs	0.000096	650	0.06
PacBio RS	2 hrs	0.01	860 – 1,500	5-10
454 GS Jr. Titanium	10 hrs	0.1	400	50
Ion Torrent – 314 chip	2.5 hrs	0.25	200	50
454 FLX Titanium	10 hrs	1	400	400
454 FLX+	20 hrs	1	650	650
Ion Torrent – 316 chip	3 hrs	1.6	200	320
Illumina MiSeq	26 hrs	4	150+150	1200
Ion Torrent – 318 chip	4.5 hrs	4	200	800
Illumina GAIIx	14 days	300	150+150	96,000
SOLiD – 5500xl	8 days	>1,410 ^d	75+35	155,100
Illumina HiSeq 1000	8.5 days	≤1500	100+100	≤300,000
Illumina HiSeq 2000	11.5 days	≤3000	100+100	≤600,000

Chybovost jednotlivých metod

Platform	Primary Errors	Single-pass Error Rate (%)	Final Error Rate (%)
3730xl (capillary)	Substitution	0.1-1	0.1-1
454	Indel	1	1
Illumina	Substitution	~0.1 (85% of reads)	~0.1 (85% of reads)
SOLiD	A-T bias	~5	≤0.1
Ion Torrent	Indel	~1	~1
PacBio RS	CG deletions	~15	≤15
Oxford Nanopore	Deletions	≥4	4

Traditional Sequencing vs. Next Generation Sequencing: Data Throughput

1 x Illumina GAI



200+ of 3730xl



Vs.

Days vs. Years

The Sequencing Landscape is Changing

Is Sanger sequencing dead?

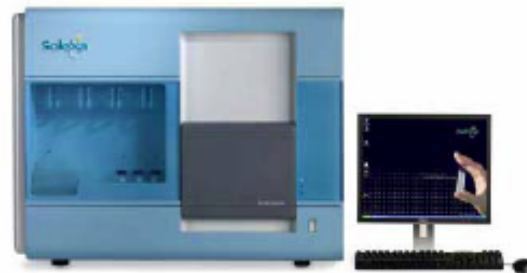
Future of sequencing centers

ABI 3730XL



- Routine sequencing
- Verify SNPs from next gen
- 1X scaffold for novel genomes

Next Gen
short read instrument
(Solexa)



“When quantity matters
but length doesn’t”

- Expression tags
- Chip Seq
- Re-sequencing

Next Gen
long read instrument
(454)



“When length matters”

- Novel genomes
- Metagenomics

Využití

1. Celogenomové sekvenování de novo
2. Celogenomové resekvenování
3. Sekvenování amplikonů (PCR produktů)
+ to samé i s RNA (resp. cDNA)



1. Celogenomové sekvenování de novo

Problém: **KRÁTKÝ READ LENGTH**

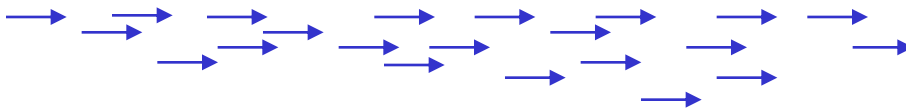
- **400bp** 454 FLX Roche
- **35-75bp** Solexa, Solid
- vs **800-1000bp** Sanger



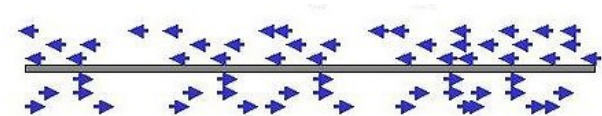
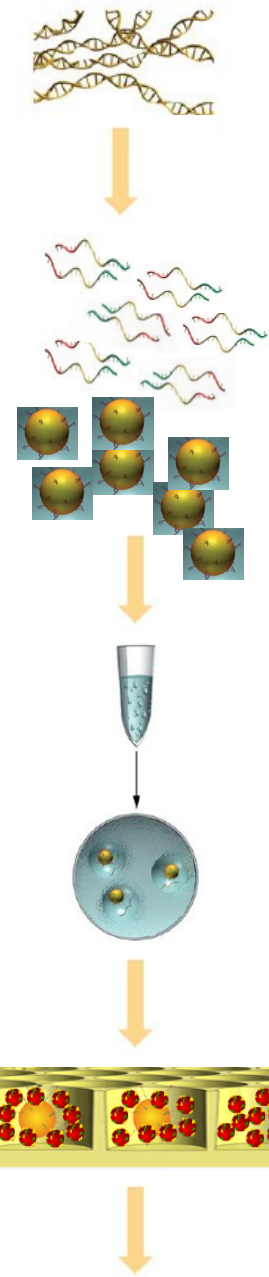
→ Uspořádání (assembly) už není problém z hlediska výpočetní kapacity

!!!! **REPETITIVNÍ OBLASTI** delší než read length !!!!

GTAAAAAAAAAAAAAAAAAAAAAAC



Zvláště komplexní eukaryotické genomy - úseky souvislých oblastí přerušovaných mezerami



1. Celogenomové sekvenování de novo

- získání kompletní uspořádané sekvence celých velkých eukaryotních genomů pomocí next-generation sequencing de novo je problém (ale to je nakonec i u Sangera)
- viry, prokaryota, malá eukaryota, mitochondrie/plastidy/plasmidy

x Ale čá



Genetic Detection and Characterization of Lujo Virus, a New Hemorrhagic Fever–Associated Arenavirus from Southern Africa

Thomas Briese^{1,3*}, Janusz T. Paweska^{2,3}, Laura K. McMullan³, Stephen K. Hutchison⁴, Craig Street¹, Gustavo Palacios¹, Marina L. Khristova⁵, Jacqueline Weyer², Robert Swanepoel², Michael Egholm⁴, Stuart T. Nichol³, W. Ian Lipkin^{1*}

¹Center for Infection and Immunity, Mailman School of Public Health, Columbia University, New York, New York, United States of America, ²Special Pathogens Unit, National Institute for Communicable Diseases of the National Health Laboratory Service, Sandringham, South Africa, ³Special Pathogens Branch, Division of Viral and Rickettsial Diseases, Centers for Disease Control and Prevention, Atlanta, Georgia, United States of America, ⁴454 Life Sciences, Branford, Connecticut, United States of America, ⁵Biotechnology Core Facility Branch, Centers for Disease Control and Prevention, Atlanta, Georgia, United States of America

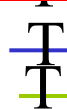
Abstract

Lujo virus (LUJV), a new member of the family *Arenaviridae* and the first hemorrhagic fever–associated arenavirus from the Old World discovered in three decades, was isolated in South Africa during an outbreak of human disease characterized by nosocomial transmission and an unprecedented high case fatality rate of 80% (4/5 cases). Unbiased pyrosequencing of RNA

node of the Old World arenaviruses. The virus G1 glycoprotein sequence was highly diverse and almost equidistant from that of other Old World and New World arenaviruses, consistent with a potential distinctive receptor tropism. LUJV is a novel, genetically distinct, highly pathogenic arenavirus.



u



Cílené sekvenování

= sekvenování jen určité části genomu či vybrané skupiny genů

Restriction enzyme genome reduction (RAD-Seq)

- sekvenování náhodných oblastí genomu vybraných na základě délky po restričním štěpení genomové DNA. Lze kombinovat vzorky DNA z více jedinců. Identifikace polymorfních markerů.

RAD-Seq

- Štěpení genomové DNA pomocí jednoho či více restričních enzymů.
- Výběr restričních fragmentů jen určité velikosti
- Sekvenování kusů vybraných fragmentů (stačí konce fragmentů).

2. Celogenomové resekvenování

- podobné problémy jako u de novo, ale méně (větší strukturální přestavby..)

KOMPARATIVNÍ GENOMIKA

- viry, prokaryota, malá eukaryota
- mitochondrie/plastidy/plasmidy

ANCIENT (mt) DNA

- různé směsné, degradované vzorky, např. fosilie

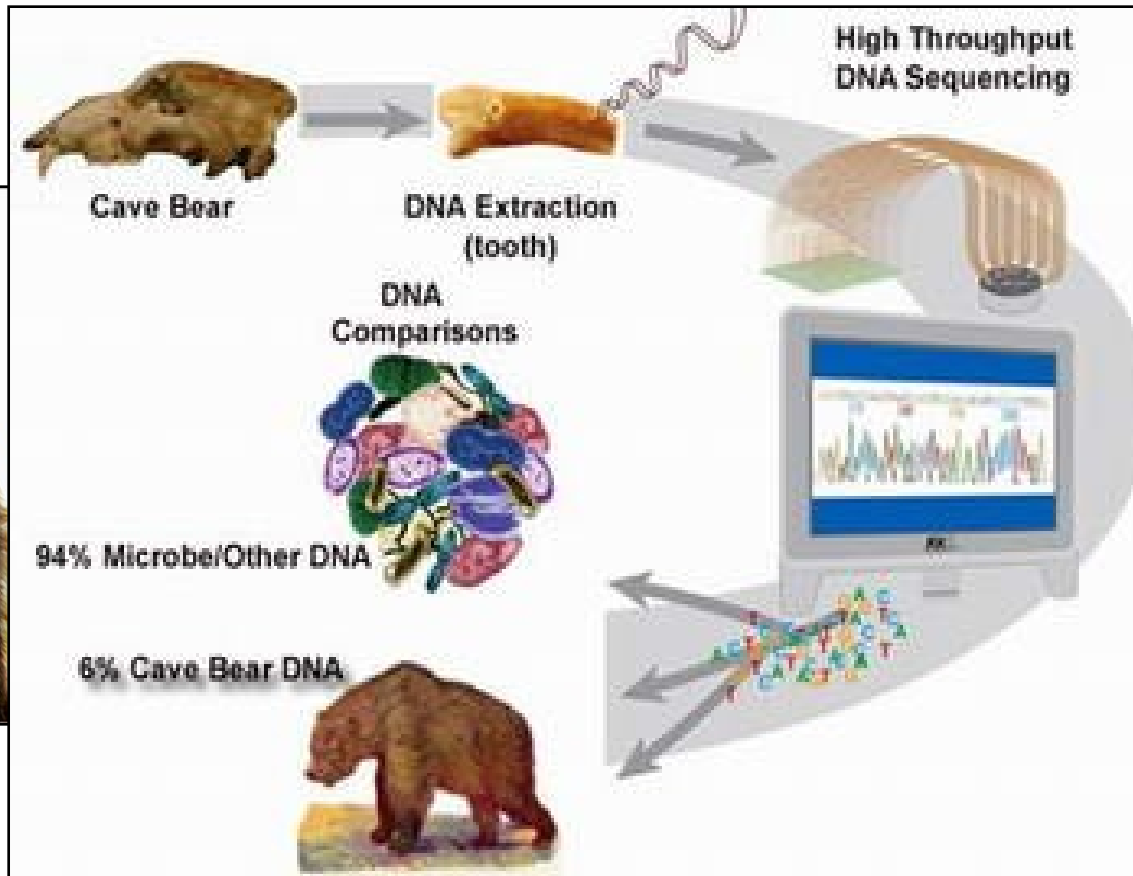
Cell

A Complete Neandertal Mitochondrial Genome Sequence Determined by High-Throughput Sequencing

Richard E. Green,^{1,*} Anna-Sapfo Malaspinas,² Johannes Krause,¹ Adrian W. Briggs,¹ Philip L.F. Johnson,³ Caroline Uhler,⁴ Matthias Meyer,¹ Jeffrey M. Good,¹ Tomislav Maricic,¹ Udo Stenzel,¹ Kay Prüfer,¹ Michael Siebauer,¹ Hernán A. Burbano,¹ Michael Ronan,⁵ Jonathan M. Rothberg,⁶ Michael Egholm,⁵ Pavao Rudan,⁷ Dejana Brajković,⁸ Željko Kučan,⁷ Ivan Gušić,⁷ Märten Wikström,⁹ Liisa Laakkonen,¹⁰ Janet Kelso,¹ Montgomery Slatkin,² and Svante Pääbo¹

Ancient Genomes Resurrected

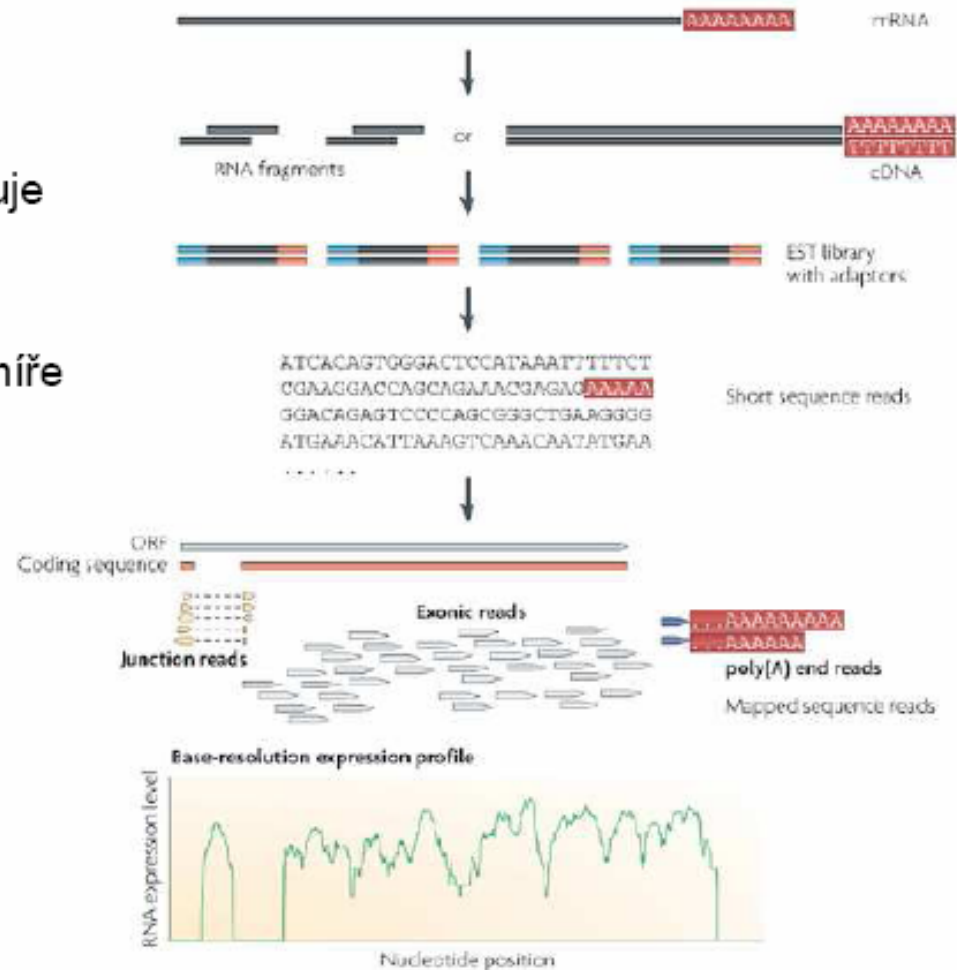
- Degraded state of the sample → mitDNA sequencing
- Nuclear genomes of ancient remains: cave bear, mommoth, Neanderthal (10^6 bp)



Problems: contamination modern humans and coisolation bacterial DNA

Sekvenování transkriptomu (RNA-Seq)

- Odpadá nutnost klonování cDNA.
- Hluboké sekvenování umožňuje identifikovat i dosud neznámé transkripty.
- Možnost získat informaci i o míře transkripce jednotlivých genů (přesnější než microarrays).
- RNA lze normalizovat – vyrovnání početnosti jednotlivých transkriptů.



3. Sekvenování ampliconů (PCR produktů)

SMĚSNÉ VZORKY

1. Metagenomika/metatranskriptomika

- Celé společenstvo půdních, vodních mikroorganismů, střevní mikroflóra
- PCR genu 16S (18S) rRNA
- lze i kvantifikovat

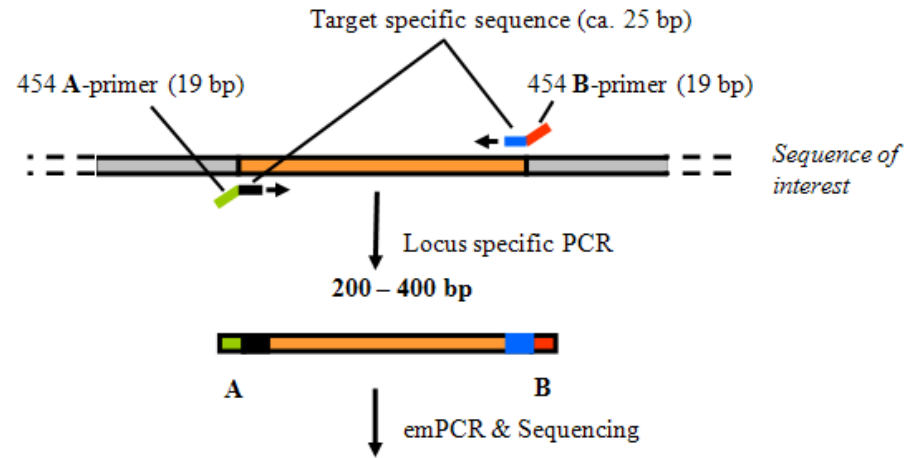
2. Složení potravy (COI barcoding)

4. Studie u kandidátních genů

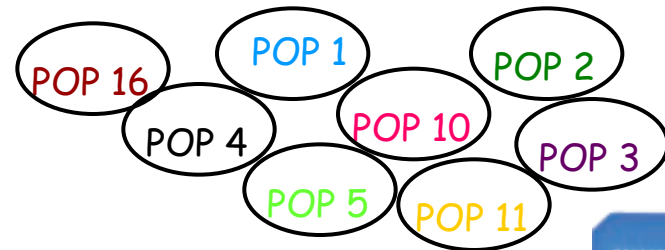
20x
NEMOCNÉ MYŠI

20x
ZDRAVÉ MYŠI

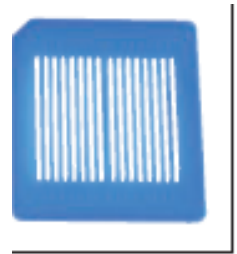
1. PCR např. imunitního genu/genů
2. Sekvenování
3. Které varianty jsou asociovány s chorobou??



3. Populační genetika



1. PCR genu/genů
2. Sekvenování
3. Zjištění sekvencí variant a frekvencí variant v každé populaci (záleží na pokrytí)



Metagenomics

- Characterizing the biodiversity found on Earth
- The growing number of sequenced genomes enables us to interpret partial sequences obtained by direct sampling of specific environmental niches.
- Examples: ocean, acid mine site, soil, coral reefs, human microbiome which may vary according to the health status of the individual

THE METAGENOMICS PROCESS



Extract all DNA from
microbial community in
sampled environment

DETERMINE WHAT THE GENES ARE (Sequence-based metagenomics)

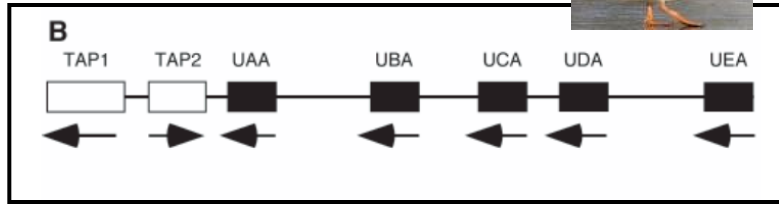
- Identify genes and metabolic pathways
- Compare to other communities
- and more...

DETERMINE WHAT THE GENES DO (Function-based metagenomics)

- Screen to identify functions of interest, such as vitamin or antibiotic production
- Find the genes that code for functions of interest
- and more...

3. Sekvenování amplikonů (PCR produktů)

5. Genové duplikace

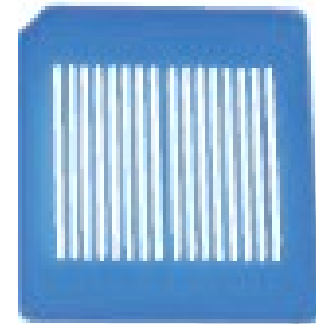
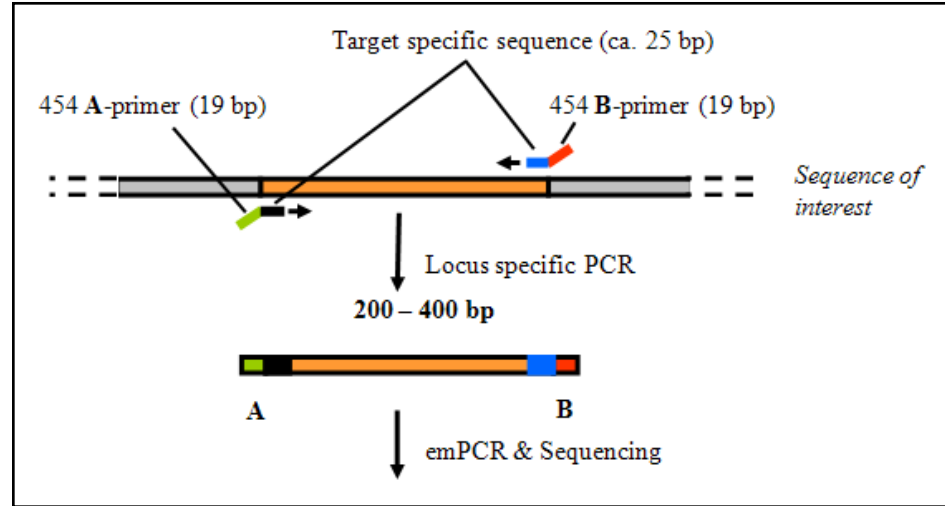


A-adaptor MID Target specific

Označí jedince

Amplifikuje všechny kopie MHC genů

Potřeba k emPCR, sekvenování..

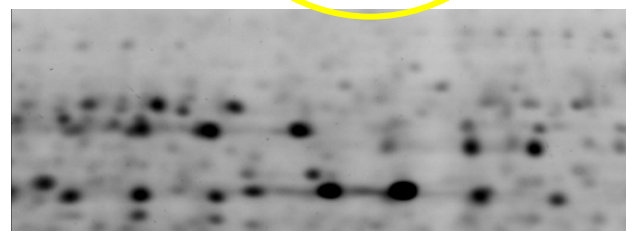
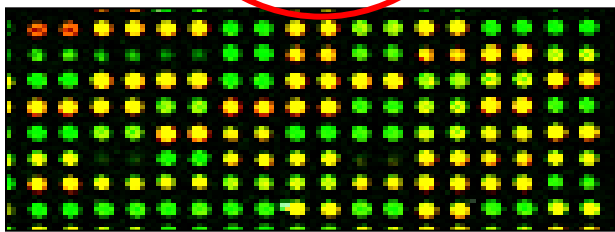
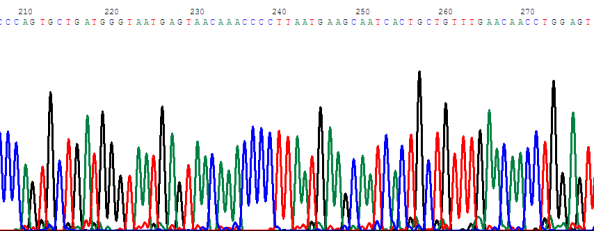
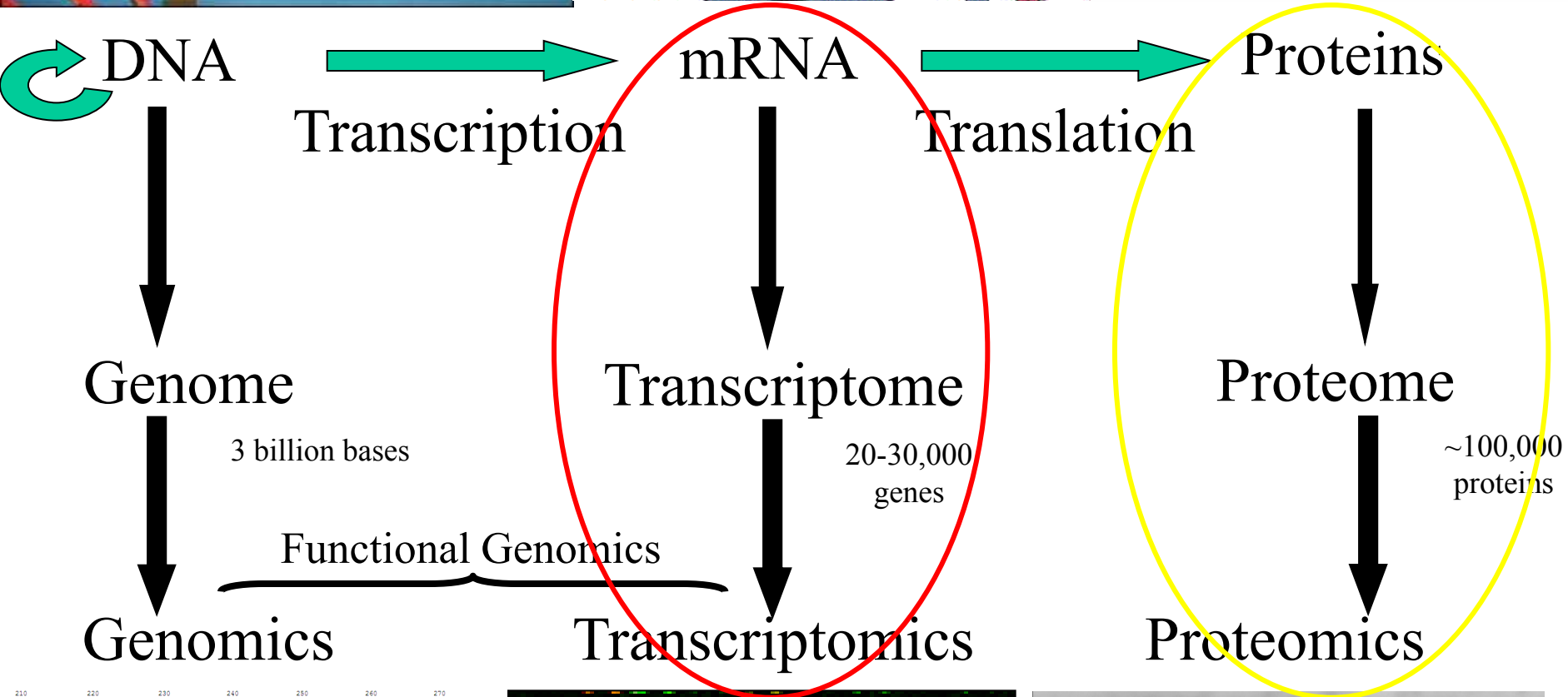
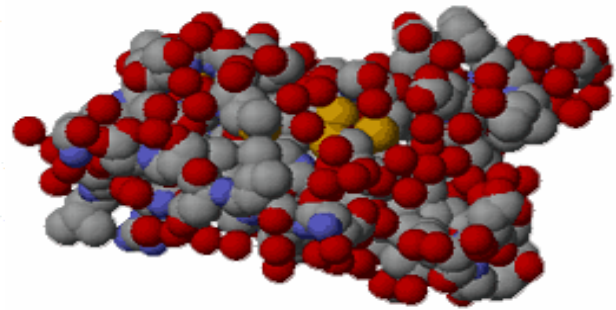
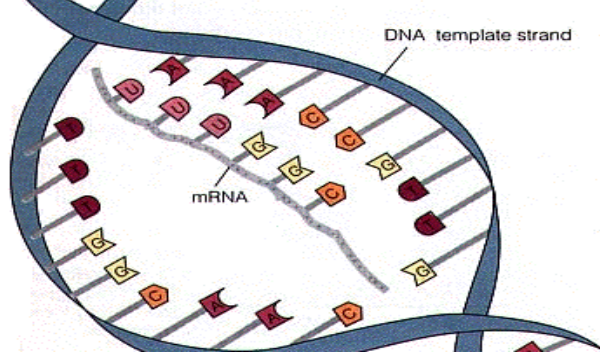
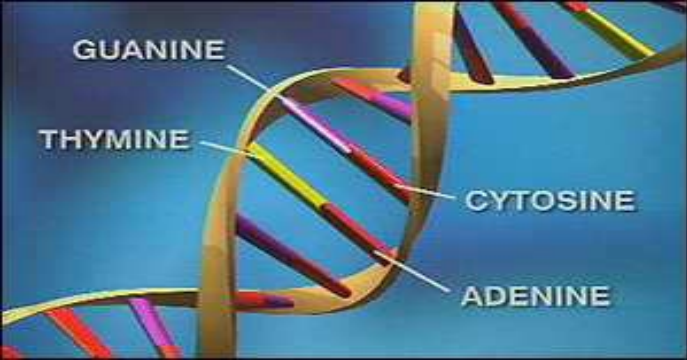


192 jedinců

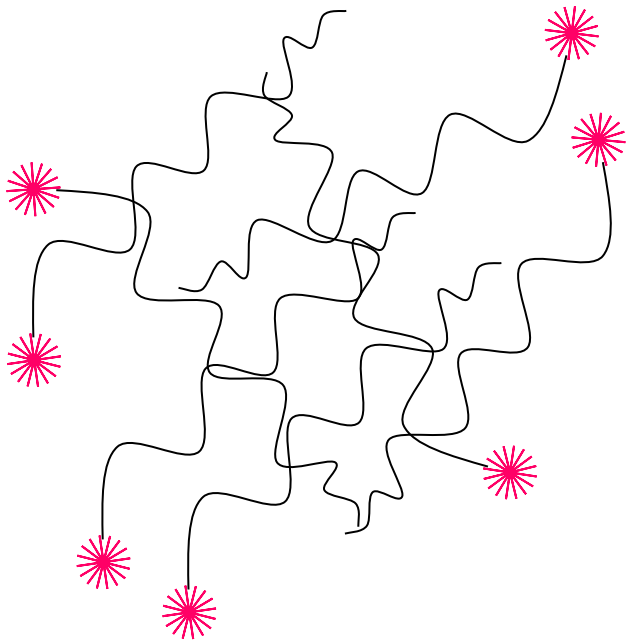
Budoucnost genetických metod v ekologickém výzkumu

2. Analysis of expression by microarrays („transcriptomics“)

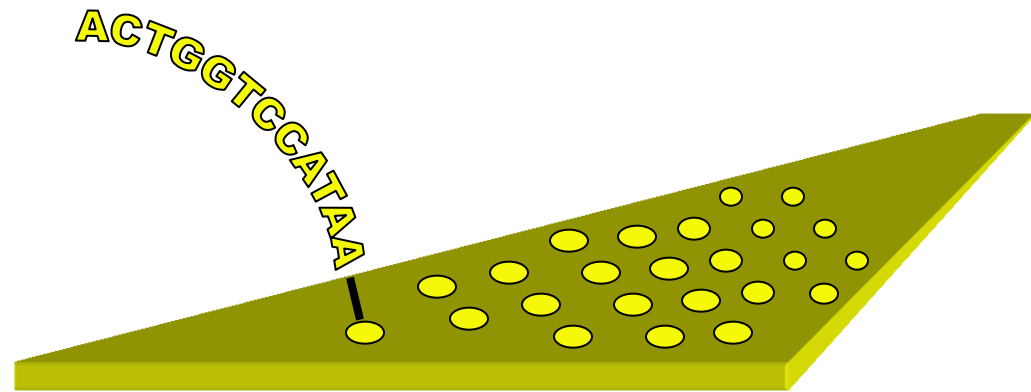
Ranz JM, Machado CA: Uncovering evolutionary patterns of gene expression using microarrays. TREE, 21(1): 29-37



Microarray analysis of transcriptome (~ specific DNA hybridization)



Target (i.e. mix of transcripts in a form of cDNA)

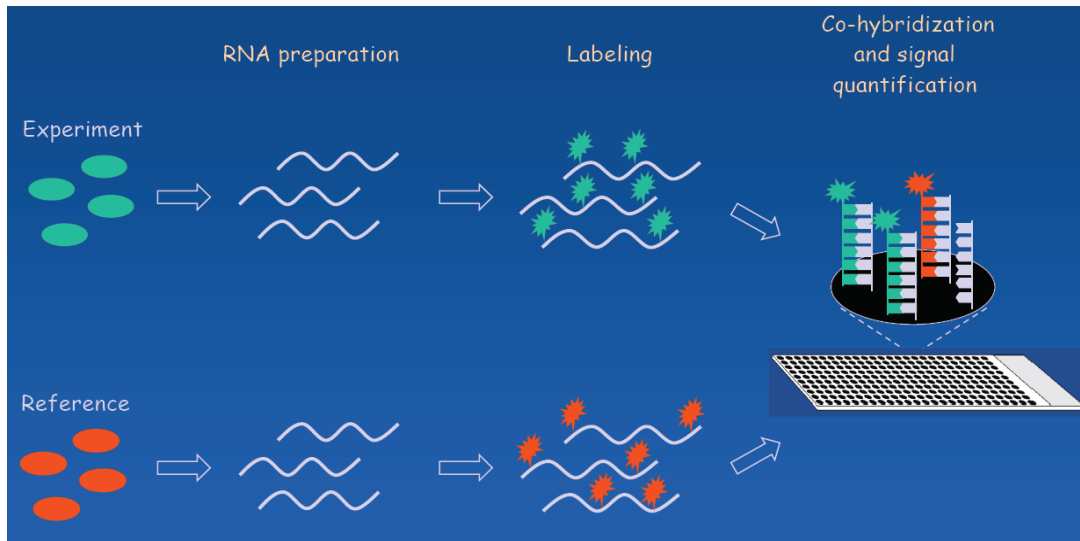


Probe (i.e. synthesized oligonucleotides complementary to particular genes)

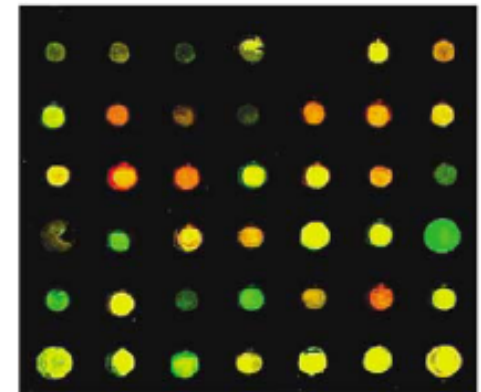
How to get a transcription profile

- vždy srovnání kontroly a „treatment“

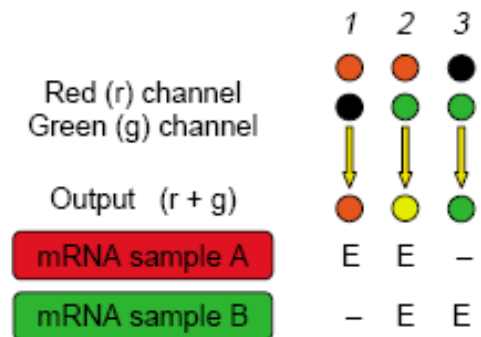
(a)



(b)



(c)



TRENDS in Ecology & Evolution

Case study: Joop Ouborg et al.

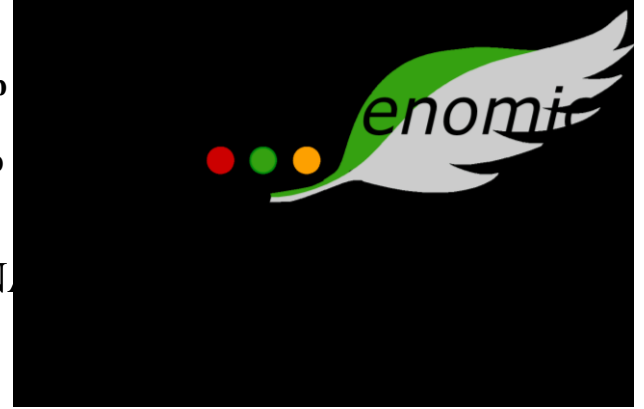
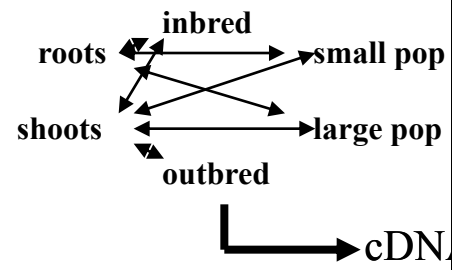
Transcriptional profiling of inbreeding depression and genetic erosion in *Scabiosa columbaria*: the balance between genetic drift and selection in the genetic erosion process.





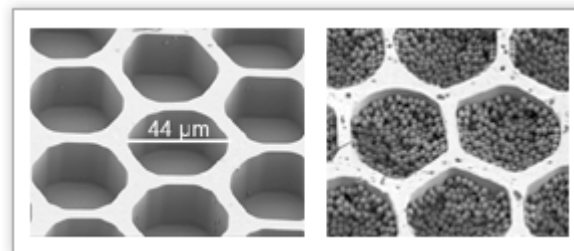
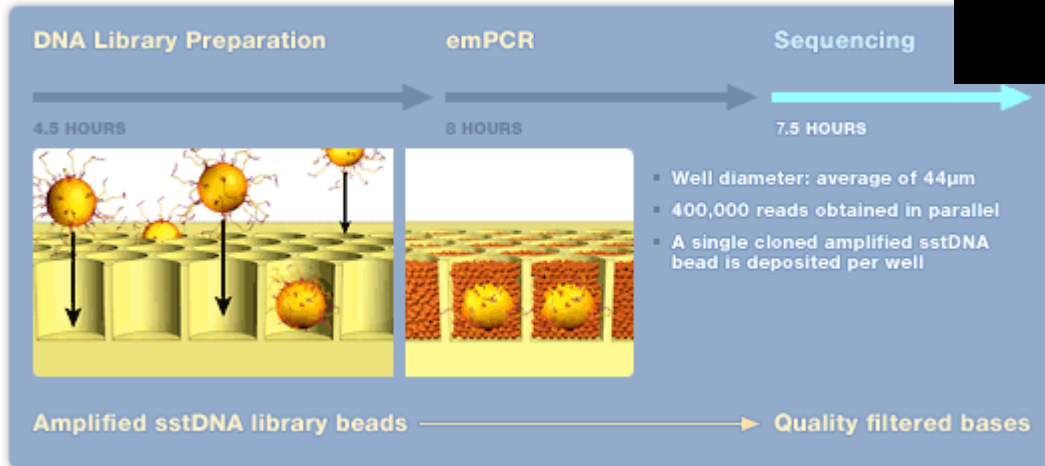
Example:

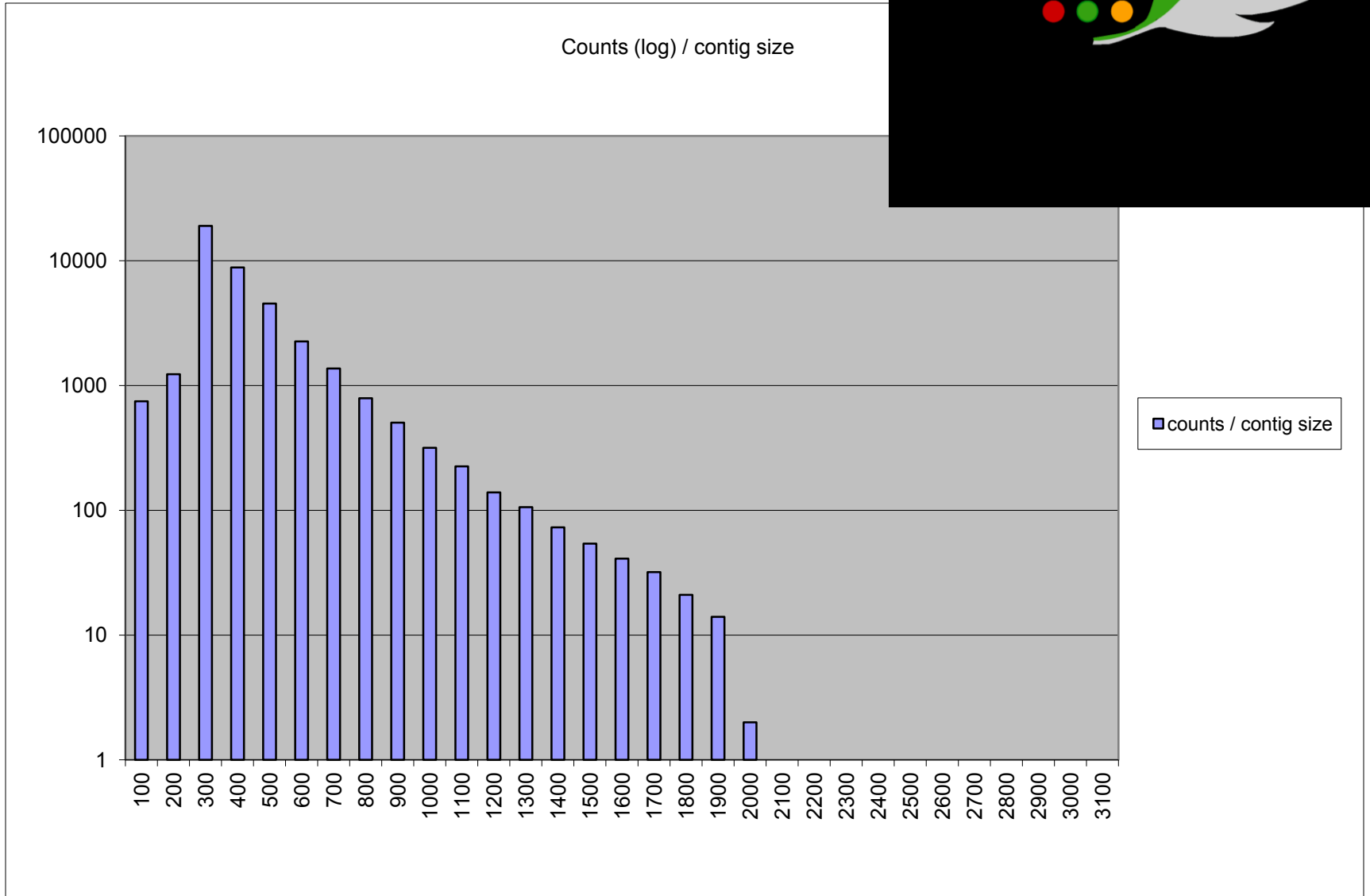
Scabiosa columbaria



cDNA library preparation – 454 sequencing

FIGURE 9





Total number of reads: **528557** Number of contigs: **40302**



In the next phase:

Annotation of these 40.000⁺ ESTs („e
sequence tags“)

Automated programs available, like **BLAST2GO** (<http://www.blast2go.de/>):

just feed a file with the ESTs into the program, and turn it on.....

1 week later you will have the results, being:

- Homology with known sequences
- Known function

The sequences may also be searched for:

EST-associated SSR markers: MISA (<http://pgrc.ipk-gatersleben.de/misa/>)

SNP markers: SNP-mining software like PolyBayes

(<http://genome.wustl.edu/tools/software/polybayes.cgi>)

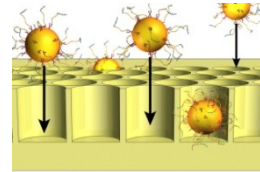
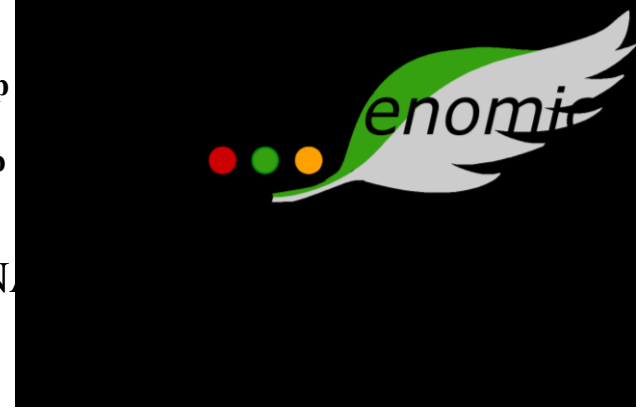
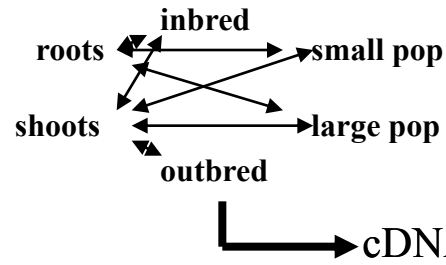
Again by using search software, freeware

ALMOST HALF OF GENES (ESTs) ARE UNKNOWN !!!



Example:

Scabiosa columbaria

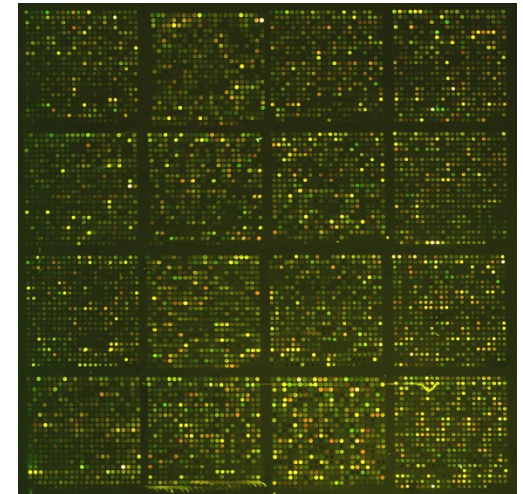
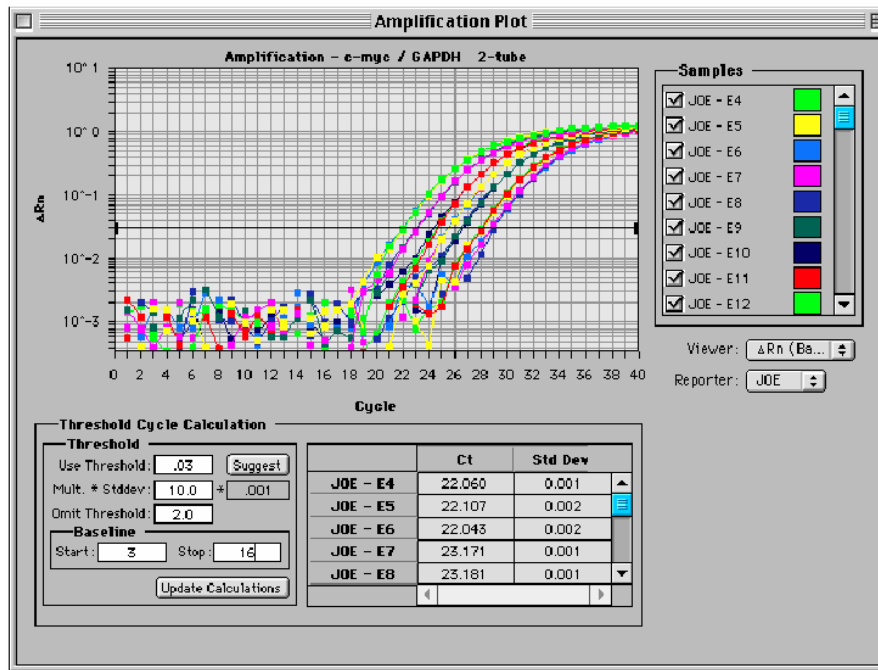


530.000 sequences in one run, leading to ~ 40.000 ESTs

Two methods of detecting transcrip

1. Design of quantitative RealTime-PCR method
EST sequences

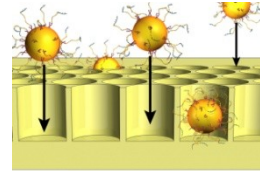
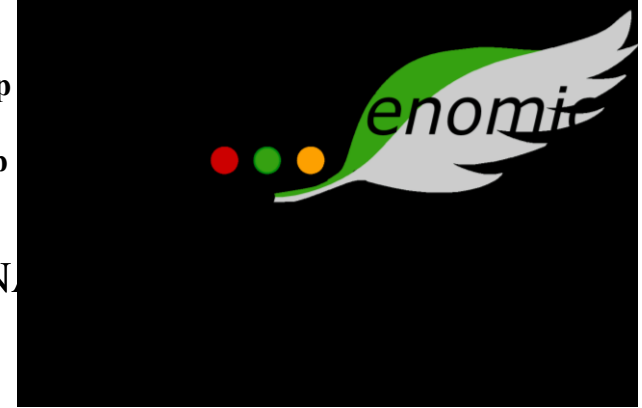
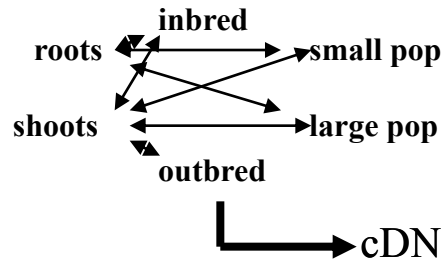
2. Design of a *Scabiosa* specific microarray



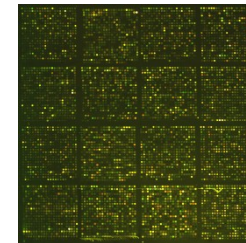


Example:

Scabiosa columbaria



530.000 sequences in one run, leading to ~ 40.000 ESTs



↓
15k – 30k
60-mer
microarrays

Experiment: transcriptional profiling of inbreeding depression



Expected pay-off:

- Ecogenomic approach to conservation genetics: effects of genetic erosion on functional genetic variation
- How does genetic erosion affect evolutionary potential?
- What is the **balance between genetic drift and natural selection** in effects of habitat fragmentation?
- Are there general **inbreeding depression genes**, or is inbreeding depression a random phenomenon?
- **Which genes are involved in inbreeding depression in different life history stages**, and can this explain the non-correlation of IBD between these stages?
- What are the **footprints of selection** in the genomes of individuals from small and large populations?
- What is the **selective value of variation in gene expression**?



Costs/requirements (2008/2009):

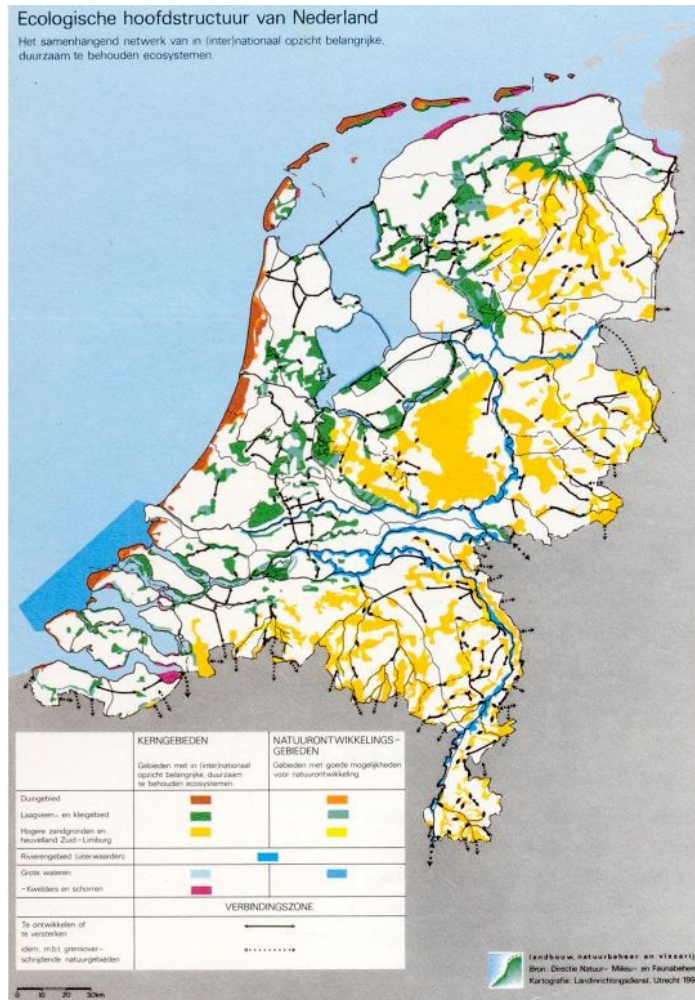
Costs are diminishing continuously

454 FLX-cDNA sequencing : 1 month, 15.000 €
(used to be
200.000 € with Sanger technology)

microarray production: 100 € per array
microarray screening: 150 € per array

cheaper options (like SOLEXA technology) are
becoming available, at much lower costs

Relative costs of conservation genomics:



Projected costs (but this is almost certain a severe underestimation):

20 billion Euro

That is:

20.000.000.000 Euro

That is equivalent to
40.000.000 microarray
runs.....

We live in exciting times !!!