

# Jak vypadají genomická data (a co se s nimi dělá)

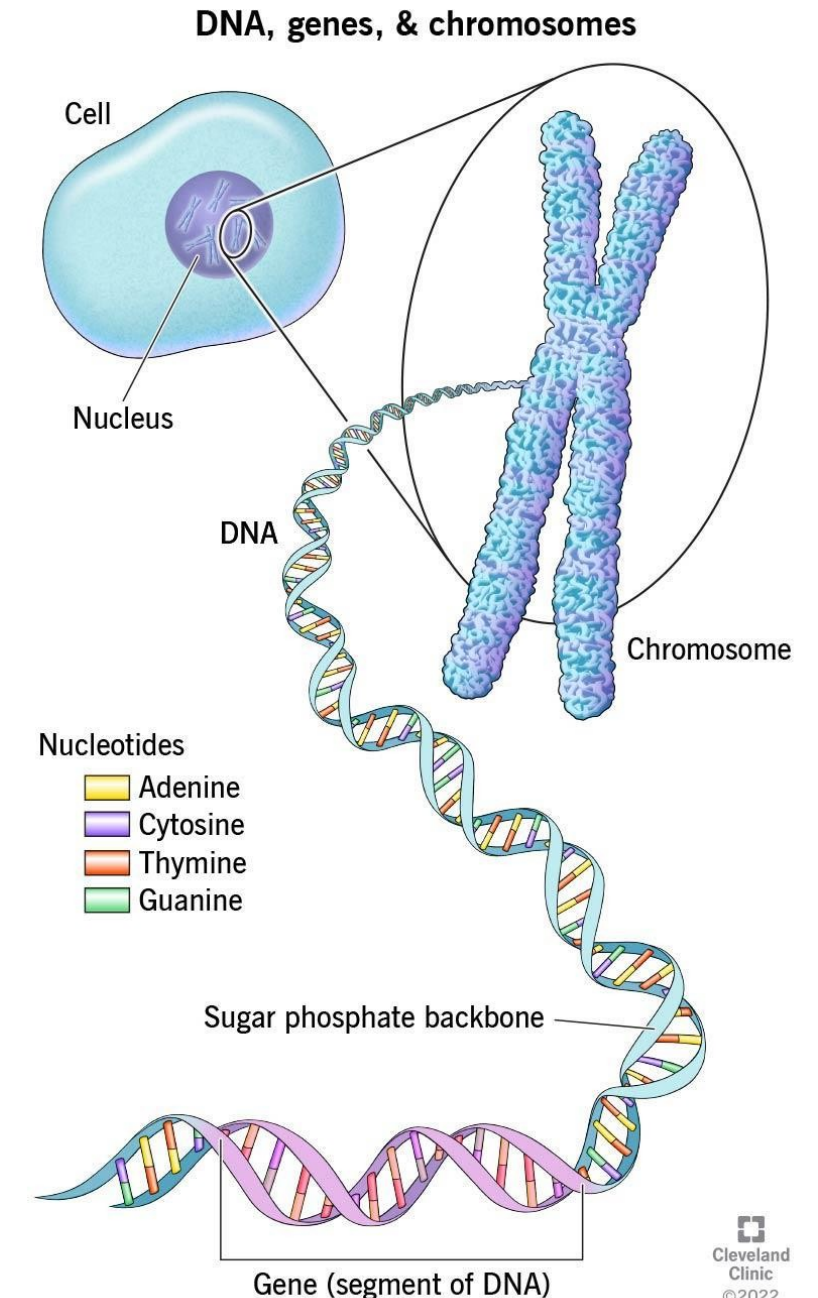
# Co je to genom?

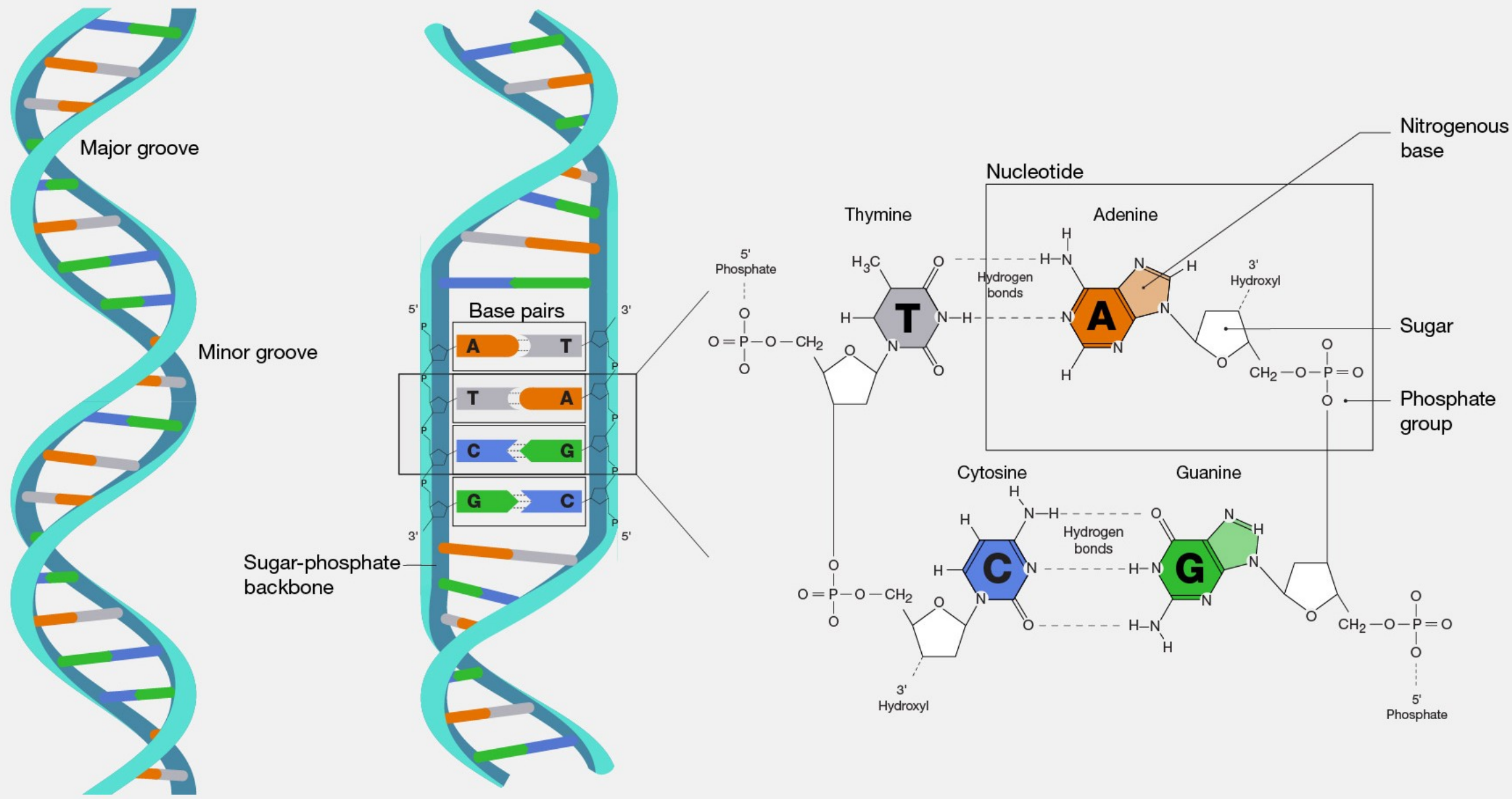
**Genom** – celková genetická informace nebo genetický materiál, který je obsažen v jádře buňky organismu. Genom zahrnuje celkovou DNA a obsahuje informace o všech vlastnostech jedince.

**Gen** – základní informační a funkční jednotka dědičné informace v živých organismech. Geny obsahují informace pro vytváření proteinů a regulačních molekul, které jsou nezbytné pro fungování organismu. Každý gen má specifickou funkci a přispívá k určitým vlastnostem a charakteristikám jedince.

\* Lidský genom zahrnuje cca 20 000 protein-kódujících genů

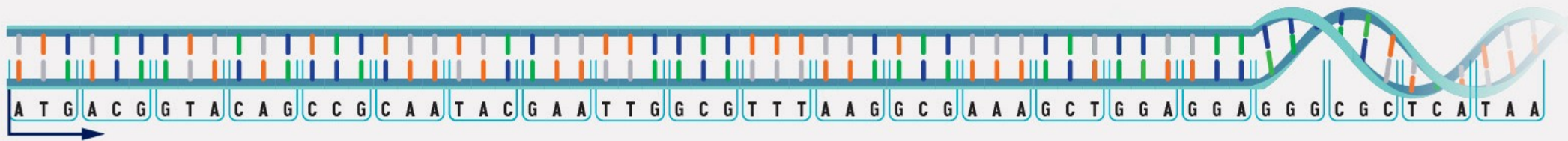
\*\* To představuje jen asi 1-2 % lidského genomu





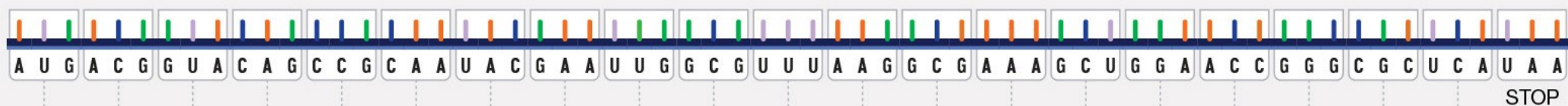
\* V lidském genomu se nachází asi 3,2 miliardy párů bází

DNA



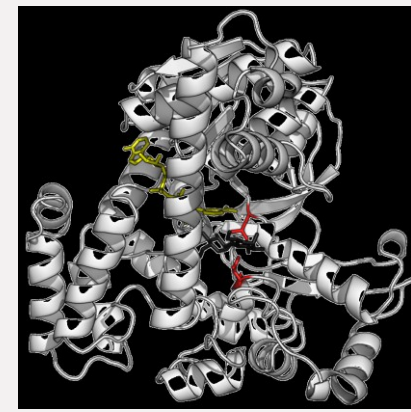
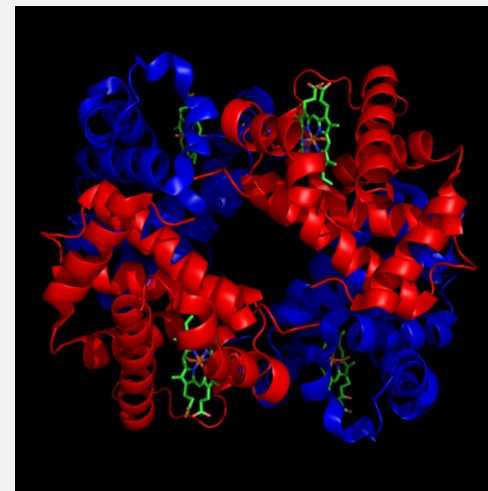
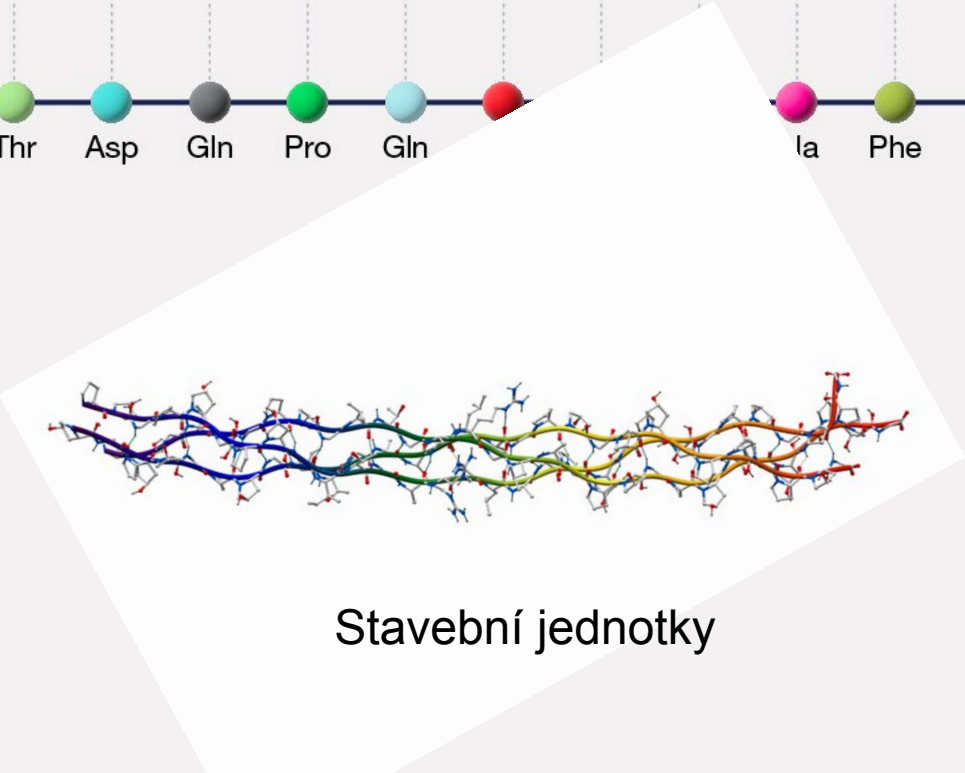
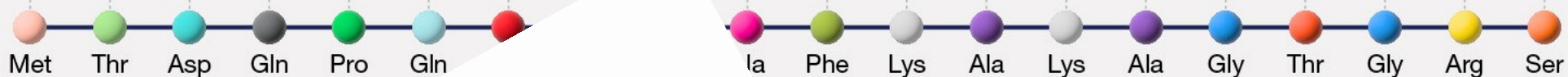
Transcription

mRNA



Translation

Protein

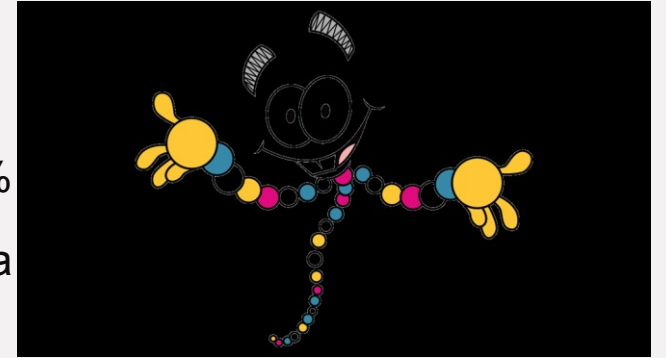


# Proč nás zajímá primární struktura

## DNA?

### 1 Výjimečnost

Genom každého člověka je **unikátní** (99,9 % genetického materiálu sdílíme, 0,1 % unikátní). Studium DNA a její primární struktury je klíčové pro **forezní vědu** a **identifikaci osob**.



### 2 Principy dědičnosti

Studiem primární struktury můžeme lépe porozumět tomu, jak se genetická informace ukládá, replikuje a předává z jednoho pokolení na druhé

### 3 Medicína

Genetické choroby jsou spojeny s konkrétními změnami v primární struktuře DNA. Studium těchto změn je nezbytné pro **diagnostiku** genetických onemocnění a vývoj **léčebných** postupů.

### 4 Genové inženýrství

Manipulace s primární strukturou DNA vede k vytvoření nových genetických konstrukcí, transgenních organismů a terapeutických postupů, jako je genová terapie.

### 5 Genová exprese

### 6 Studium evoluce/biodiverzity

### 7 Studium mikrobiomu

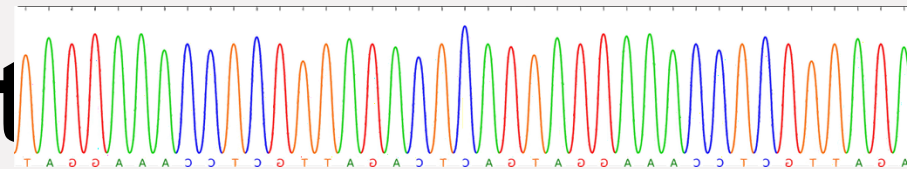
### 8 ...

# Co jsou tedy genomická data?

Genomická data jsou soubory informací získaných z analýzy genetického materiálu (DNA nebo RNA) organismu. Tato data obsahují podrobné informace o **struktuře, sekvenci a funkci** genů a dalších genetických elementů v genomu daného organismu. Genomická data mohou zahrnovat následující informace:

- 1.Sekvence DNA:** Pořadí nukleotidů (A, C, G, T) v molekule DNA. Sekvence DNA umožňují identifikovat geny, regulační sekvence a další důležité úseky DNA.
- 2.Geny:** Genomická data obsahují informace o umístění a struktuře genů v genomu. Identifikujeme, jaký protein nebo RNA každý gen kóduje, a jakým způsobem může ovlivnit funkci organismu.
- 3.Genetické varianty:** Změny v sekvenci DNA, jako jsou jednonukleotidové polymorfismy (SNP), delece, inserce a další mutace. Genomická data obsahují informace o těchto variantách, které mohou mít význam pro dědičnost, vývoj chorob a individuální rozdíly mezi lidmi.
- 4.Genová exprese:** Informace aktuální aktivaci genů a jak se projevují prostřednictvím produkce proteinů nebo RNA molekul. To umožňuje studovat, jaké geny jsou zapojeny do různých biologických procesů a jak mohou být ovlivněny různými podmínkami.
- 5.Struktura chromozomů:** Data mohou obsahovat informace o uspořádání chromozomů v jádře buňky, což je důležité pro studium chromozomálních abnormalit a genetických onemocnění.
- 6.Mikrobiom:** Analýza genetického materiálu mikroorganismů v rámci lidského těla – zastoupení a funkce/vliv těchto mikroorganismů.

# The Human Genome Project



Human Genome Project měl za cíl zmapovat a rozluštit celý lidský genom. Projekt byl zahájen v roce **1990** a oficiálně dokončen v roce **2003**. Jednalo se o společný úsilí mezinárodního vědeckého komunity. Klíčové cíle projektu zahrnovaly:

- 1. Zmapování lidského genomu:** Pořadí nukleotidů v lidské DNA v každém z 23 lidských chromozomů.
- 2. Identifikace genů:** Identifikace a anotace všech genů v lidském genomu. To zahrnovalo určení, kde se jednotlivé geny nacházejí, jaké mají funkce a jaké proteiny kódují.
- 3. Studium genetických variant:** Projekt zkoumal genetické varianty, jako jsou SNP (jednonukleotidové polymorfismy).
- 4. Aplikace v medicíně:** Využití genomických znalostí pro zlepšení diagnostiky, prevence a léčby genetických chorob.

Dokončení projektu v roce 2003 bylo historickým milníkem v oblasti biologie a genetiky. Výsledky projektu poskytly základ pro rozvoj genomiky a personalizované medicíny.

Důležitým aspektem projektu bylo, že data z něj byla **veřejně dostupná** pro vědeckou komunitu, což umožnilo mnoha dalším výzkumným projektům a studiím využívat tuto cennou informační základnu pro další pokroky v oblasti genetiky a biomedicíny.

# The Human Genome Project

Na sekvenování lidského genomu se podíleli vědci z **20 různých univerzit a výzkumných center** ze Spojených států, Velké Británie, Francie, Německa, Japonska a Číny.

Sekvence lidského genomu **nepochází od jediného člověka**, ale několika lidí, jejichž identita byla záměrně anonymizována, aby bylo chráněno jejich soukromí.

Původně předpokládané náklady činily **3 miliardy dolarů**, přičemž předpokládaná doba trvání projektu byla 15 let. Tato přibližná částka se blíží přesnému číslu.

Projekt ve výsledku **nevygeneroval** kompletní lidský genom. V dubnu 2003 konsorcium oznámilo, že vytvořilo v podstatě kompletní sekvenci lidského genomu. Konkrétně představovala 92 % lidského genomu a obsahovala méně než 400 mezer.

Dne 31. března 2022 konsorcium Telomere-to-Telomere (T2T) oznámilo, že doplnilo zbývající mezery a vytvořilo **první skutečně kompletní sekvenci** lidského genomu.

01 + info

## WHAT IS A GENOME?

A genome is all of the base sequences in an organism's DNA (the letter). In other words, the genome is the whole of the genetic information of an organism.

03

## WHAT WERE THE OUTCOMES OF THE HUMAN GENOME PROJECT?

**Identification:** The number, position, scale, and sequence of human genes have all been determined in the human genome project.

**Development:** This has enabled the development of unique gene probes for detecting genetic disorder sufferers and carriers.

**Medicine:** The discovery of different proteins which resulted in better treatments (rational drug design)

**Bloodline:** Comparisons with other genomes also revealed information about human evolution and history.

05

## WHAT ARE THE BENEFITS OF THE HUMAN GENOME PROJECT?

This has been an advantage in making the information available to researchers to aid with the following:  
Identifying genes and genetic diseases (more research)

Development of gene probes to help detect sufferers and carriers of genetic diseases (e.g. sickle cell anemia)  
Additionally, thanks to the human genome project it has been possible to determine what type of molecules healthy individuals make and what genes make up those molecules.

Along with how to copy these genes to distribute those healthy molecules to sick people to help them.

02 + info

## WHAT IS THE HUMAN GENOME PROJECT AND WHAT'S ITS AIM?

The human genome project is an international cooperation project created to sequence the complete human genome

04 + info

## WHO CREATED THE HUMAN GENOME PROJECT?

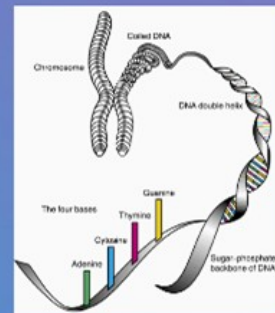
An international team of scientists set out the human genome project in which they were able to find the entire base sequence of human genes and found out that they were made up of over 3 billion letters. They then made all their findings public in 2003 which was an advantage to a lot of scientists and researchers.

06

## NOTE THAT:

The human genome is made up of:

- 46 chromosomes
- ~21,000 genes
- ~3 billion base pairs.



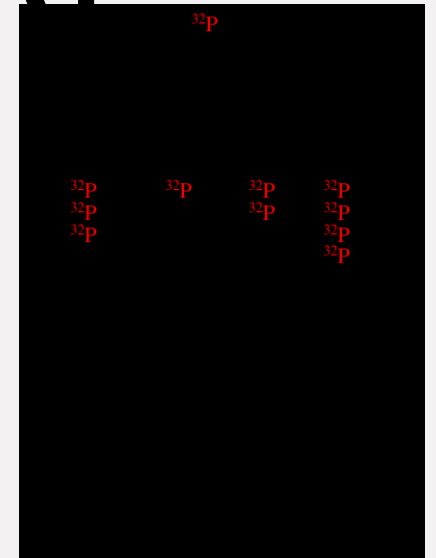


# Jak přečíst primární sekvenci DNA?

Určení sekvence = metoda SEKVENOVÁNÍ

## MAXAM-GILBERTOVA METODA

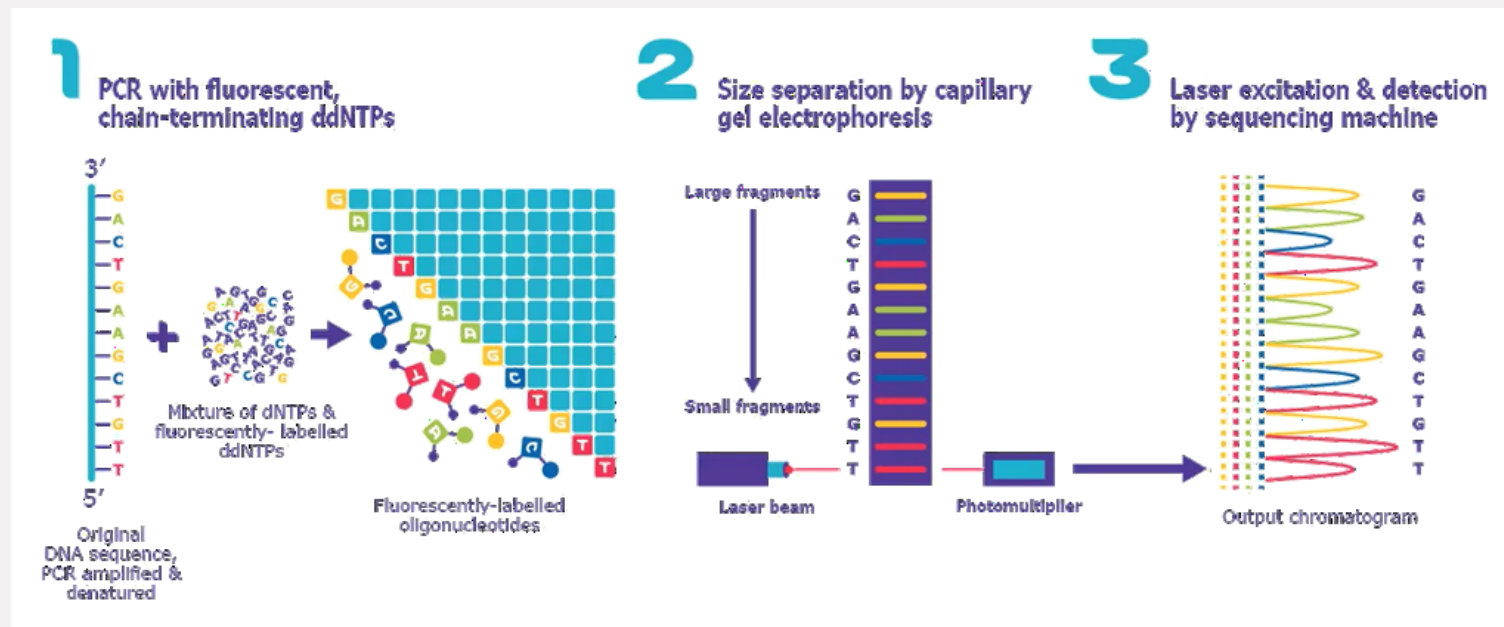
Vysoce toxické chemikálie (radioaktivní značení 5' konce DNA pomocí  $^{32}\text{P}$ ), štěpení DNA činidly, gelová elektroforéza, autoradiografie



## SANGEROVA METODA

Založená na selektivním začleňování dideoxynukleotidů (ddNTP) – přerušení syntézy DNA.

Vzorek DNA je rozdělen do 4 oddělených reakcí, které obsahují všechny standardní deoxynukleotidy. Ke každé reakci je přidán pouze jeden ze čtyř dideoxynukleotidů (ddATP, ddGTP, ddCTP, nebo ddTTP). Dochází k syntéze komplementárního řetězce do začlenění značených ddNTP. Fragmenty jsou separovány pomocí GE a analyzován fluorescenční signál (každý ddNTP jiná barva).



# Sekvenování nové generace (NGS)

- Miniaturizace a paralelizace sekvenačních technologií
- Analýza mnoha molekul/fragmentů najednou
- Rychlejší než tradiční metody

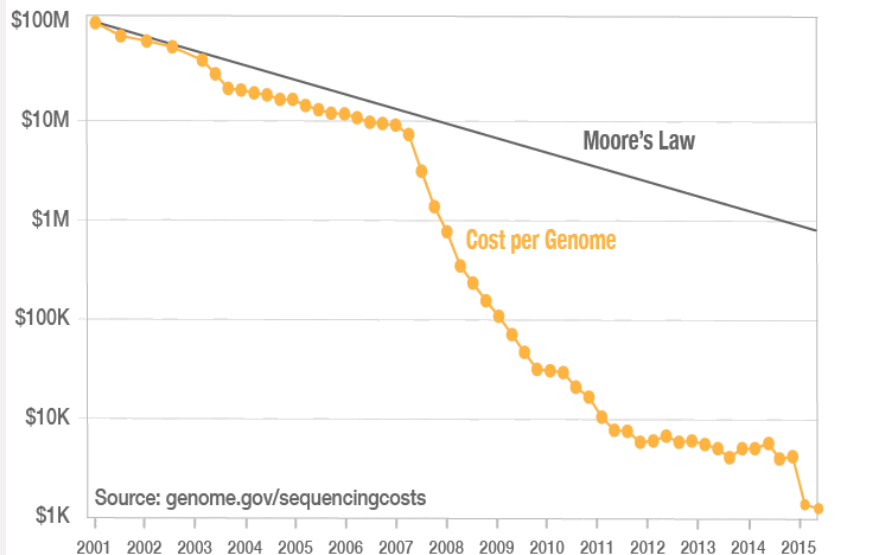
- Snížené náklady na sekvenování (automatizace)
- Vyšší přesnost sekvenování
- Fragmentace DNA – 50 až 500 bazí
- **Několik technologií (každá má výhody i**



**omezení a je vhodná pro různé typy**

Metoda	Sangerova	Illumina
Kapacita (bp / hod)	76 000	1 800 000 000
Cena (€ / Gbp)	1 250 000	50

**Náklady na sekvenování lidského genomu.** Odklon křivky nákladů na sekvenování od Moorova zákona se shoduje s nástupem sekvenování nové generace (NGS). Moorův zákon pochází z odvětví počítačového hardwaru, který zahrnuje zdvojnásobení "výpočetního výkonu" každé dva roky. Má se za to, že technologie, které se řídí tímto zákonem, jsou považovány za úspěšné. Představuje tedy užitečný vztah pro porovnávání technologického pokroku.



Sample/genome



Fragmentation



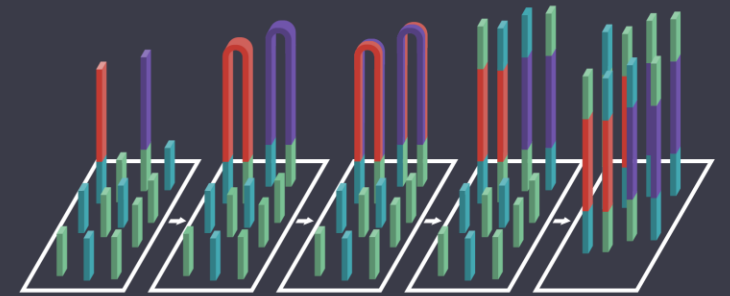
End repair & adapter ligation



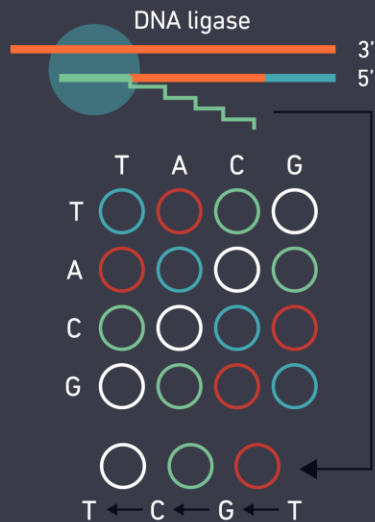
Clonal amplification by emulsion PCR



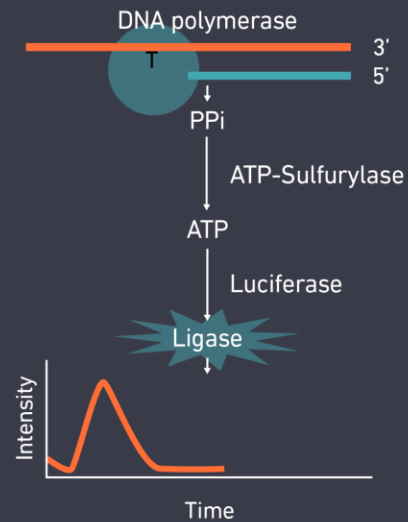
Clonal amplification by bridge PCR



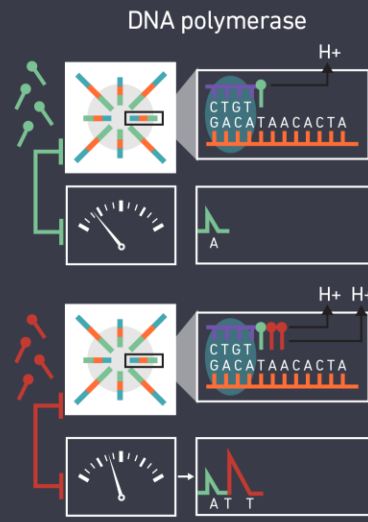
Sequencing by ligation



Pyrosequencing

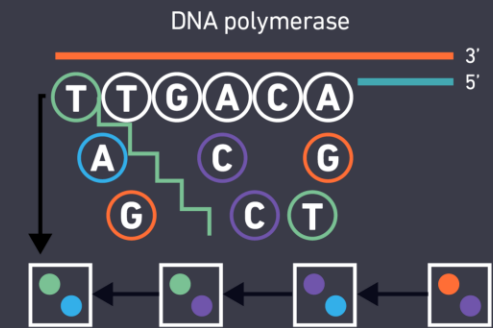


Proton detection sequencing



Cluster generation

Reversible terminator sequencing



# Illumina

- Dříve Solexa
- Momentálně nejrozšířenější technologie masivně paralelního sekvenování
- Princip sekvenování syntézou za pomoci reverzibilních terminátorů
- Klonální amplifikace (můstková PCR)
- 99% přesnost



+ vysoká přesnost

nejnižší cena za jednu osekvenovanou bázi

množství publikací využívající Illumina technologii

množství komerčně dostupných kitů pro různé aplikace

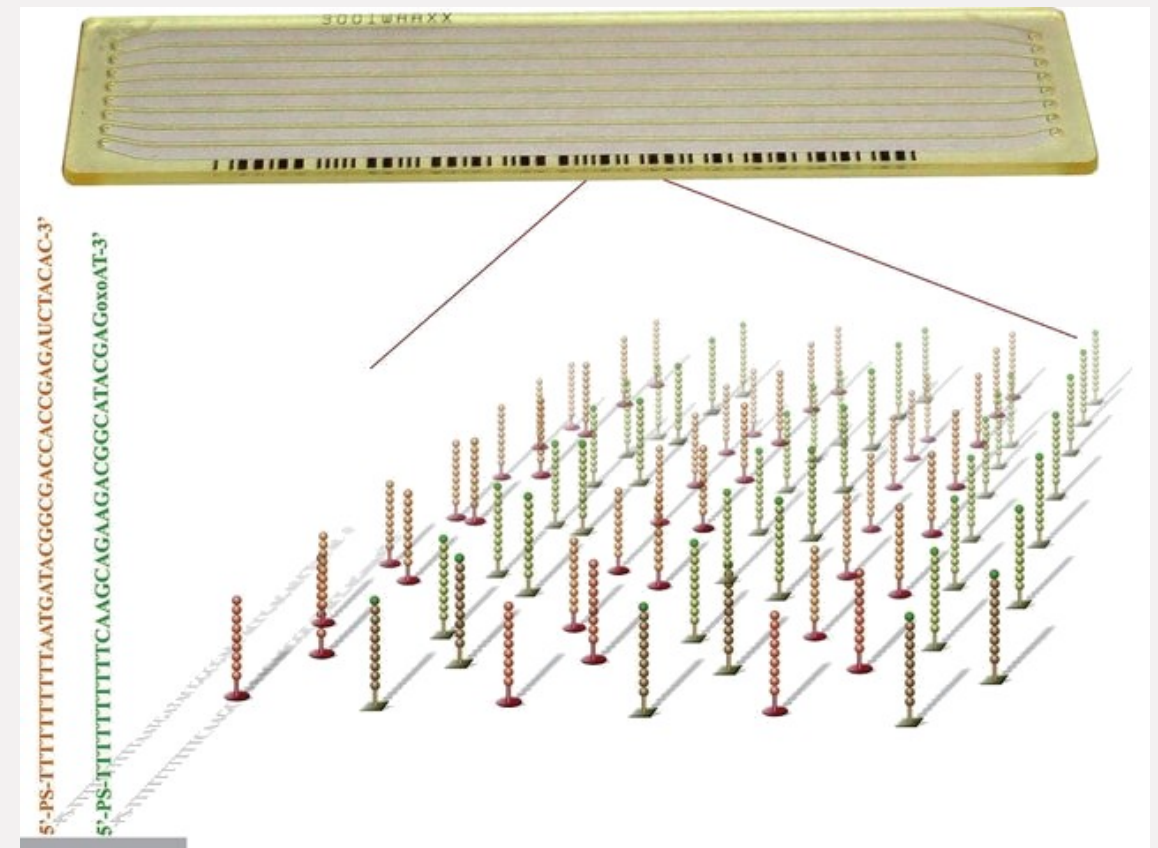
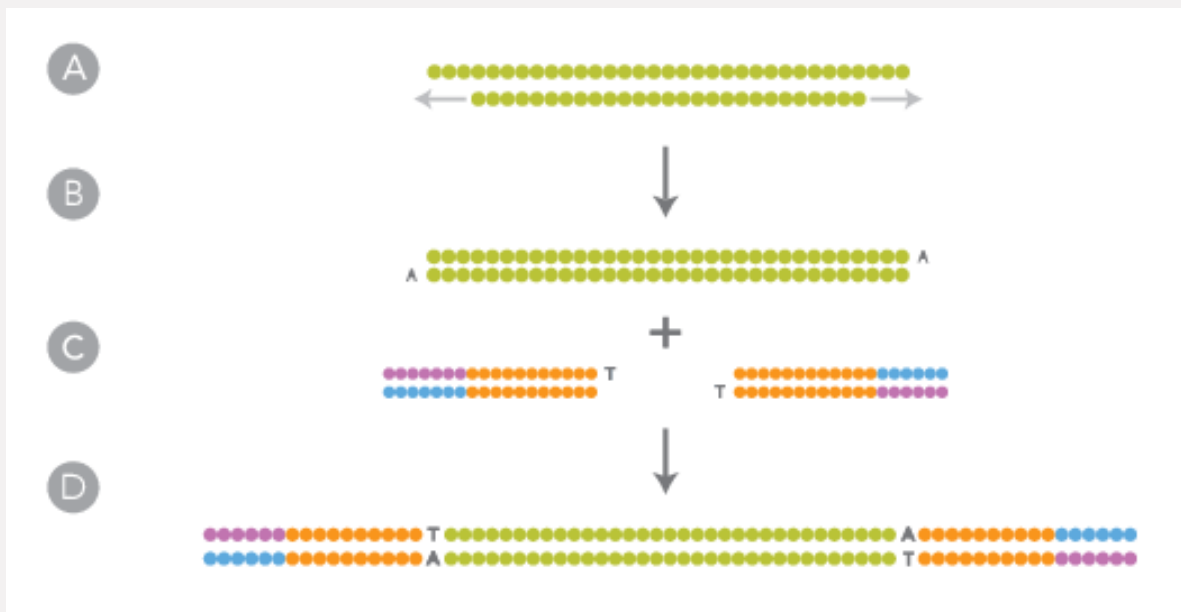
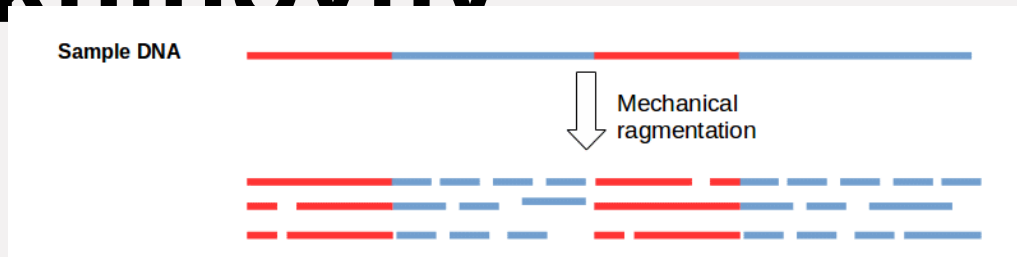
- vysoké pořizovací náklady

krátká délka čtení – max 300 bazí – pouze některé přístroje

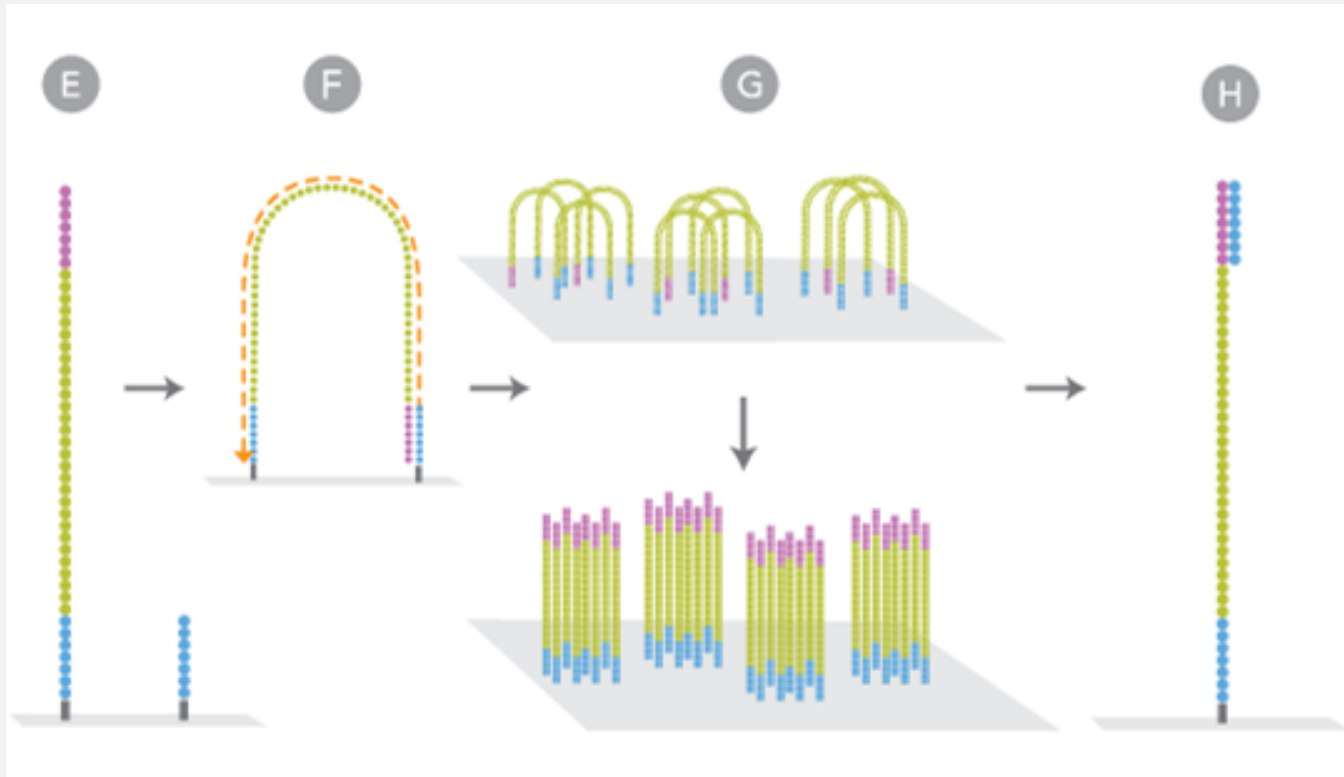
komplikovaná de novo assembly z Illumina NGS dat

dlouhá doba sekvenování 12h až 4 dny

# Illumina – příprava sekvenační knihovny

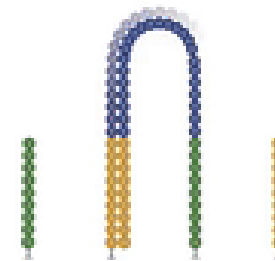


# Illumina – amplifikace

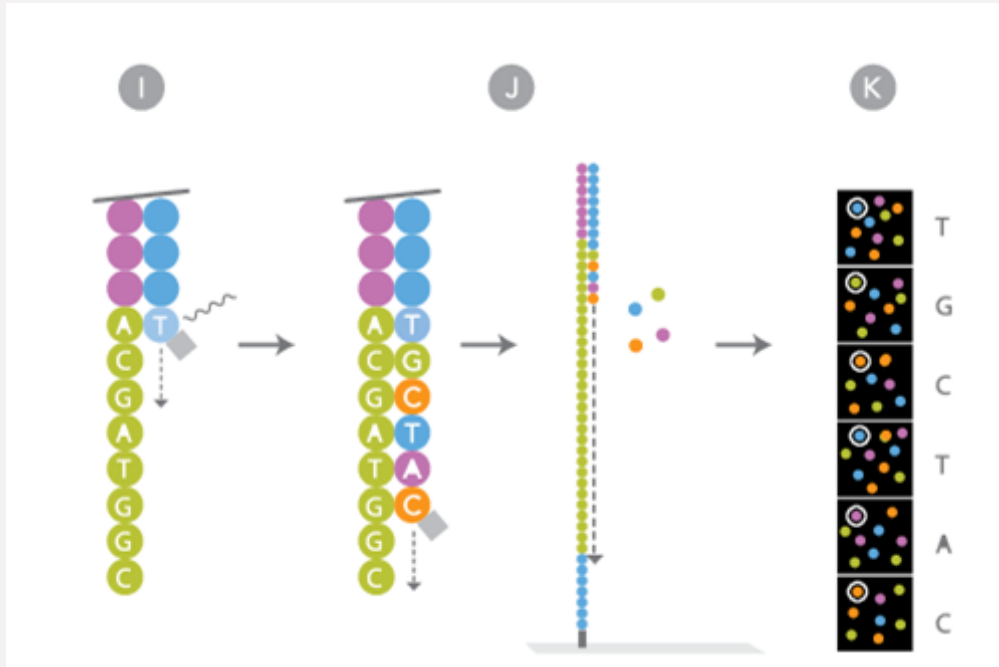


Fragmenty jsou amplifikovány pomocí můstkové PCR, dochází k tvorbě klastrů (jeden fragment = jeden klastr, obsahuje tisíce až miliony kopií)

eration



# Illumina – sekvenace



# Bioinformatická analýza NGS dat

## 1 Primární

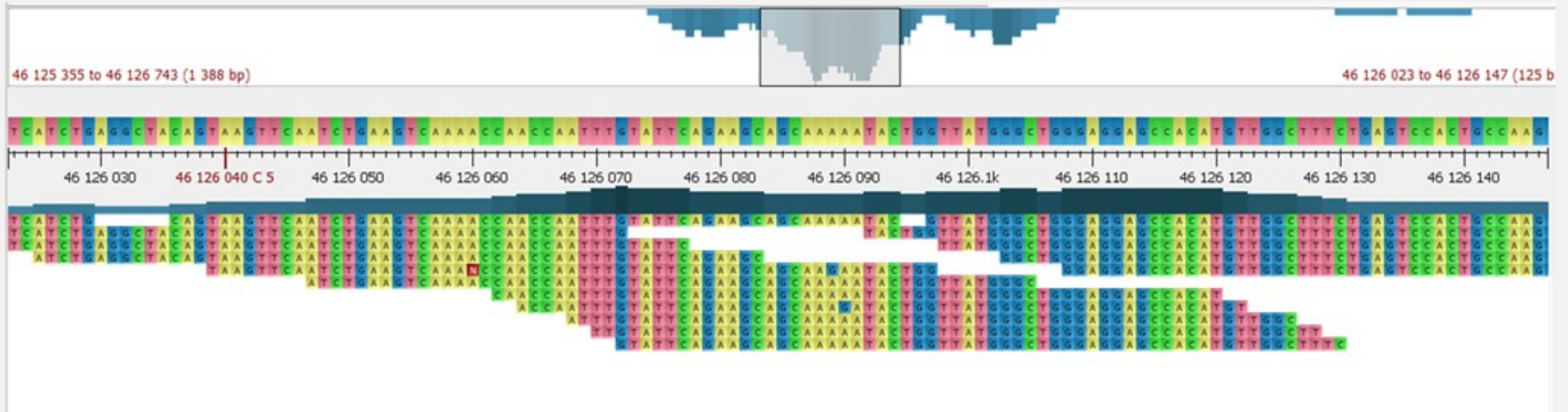
Konverzi surových přístrojových dat na sekvenční data sestávající z pořadí nukleotidových bází. Často prováděna přímo na sekvenátoru.

## 2 Sekundární

Sestavení genomu/ části genomu a detekce variací.

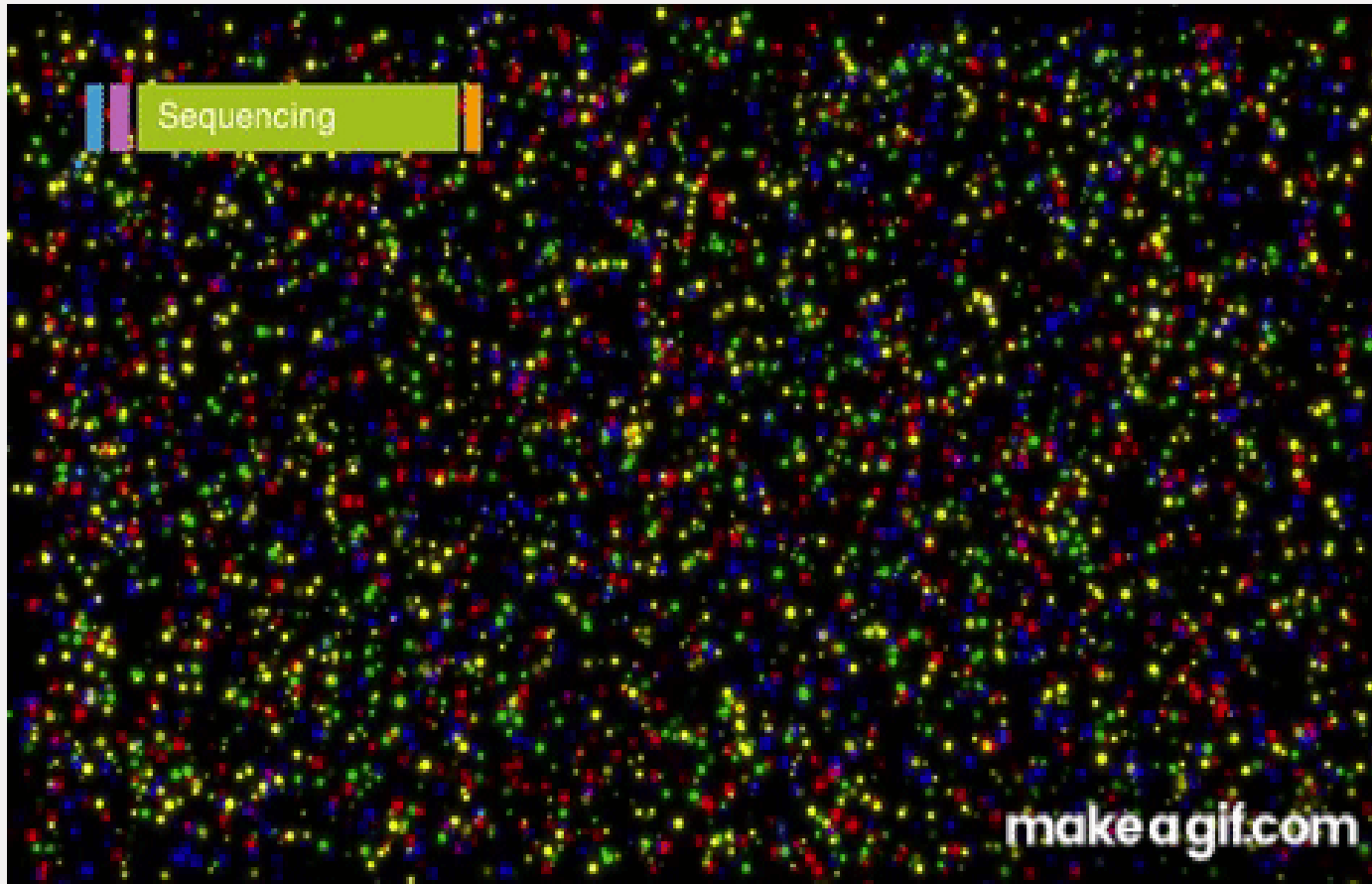
## 3 Terciární

Anotace získaných dat, biologický kontext.





# Primární analýza dat



4,8 TB

\* Každý z miliard klastrů vytvoří 150 bp čtení

Surová obrazová data se převedou do tabulky intenzity, která zaznamenává umístění každého klastru a hodnoty intenzity barev (jedna barva pro každou ze čtyř bází). Tyto číselné hodnoty se převedou na tzv. base call a použijí se k přiřazení hodnoty kvality pro dané pořadí bází.

Lane	Tile	X	Y	Cycle 1 - A C G T				Cycle 2 - A C G T			
5	12	924	1580	493.1	388.9	3826.7	2359.4	185.6	122.3	360.4	307.8
5	12	773	395	85.5	113.0	2327.5	1158.0	156.3	166.9	113.5	909.6
5	12	105	786	1243.8	741.1	45.8	67.4	318.4	692.6	48.3	41.7
5	12	598	690					3.6	505.7	1919.1	959.3
5	12	1107	1207					8.6	230.5	815.1	512.1
5	12	1074	466					38.4	41.8	64.9	1102.9
5	12	887	356	743.1	488.4	42.2	305.0	230.3	603.6	-63.1	-20.1
5	12	642	1789	63.2	54.3	861.7	595.7	81.5	86.0	54.9	385.4
5	12	599	314	845.5	533.2	45.2	581.0	269.9	569.9	13.0	78.4
5	12	839	1103	372.0	812.6	16.7	70.5	58.4	89.4	35.4	1394.9
5	12	347	1792	343.8	766.9	108.4	638.5	73.2	43.9	121.6	1882.2
5	12	807	1114	63.9	63.8	828.3	1369.0	1074.4	714.3	-39.9	29.4

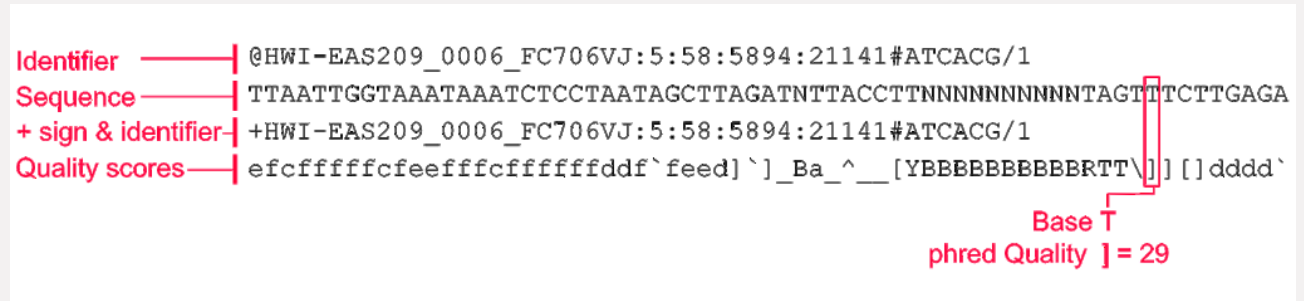
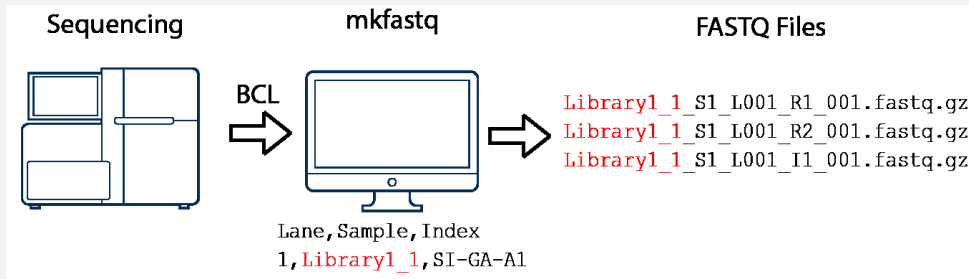
250GB

Lane	Tile	X	Y	Cycle 1 - A C G T				Cycle 2 - A C G T			
1	7										
1	7										
1	7										
1	7										
1	7	214	595								
1	7	155	544								
1	7	301	507								
1	7	175	606								
1	7	242	522								
1	7	196	522								
1	7	237	612								
1	7	160	528								
1	7	164	543								

60GB

```
s_1_0007_seq.txt - WordPad
File Edit View Insert Format Help
| 7 | GAACAAGCATAT
1 7 | TTTTTTTTTTTT
1 7 | GATCATGTTTTC
1 7 | CCTGCCTCAGCC
1 7 | TACAAAATCCCTGCC
1 7 | TTATCTGCATCCGGT
1 7 | TCCCTGCTTATTGAC
1 7 | TTGGAATCGGGGTTA
1 7 | TAACAAATATACAGG
1 7 | TGTCACAGGAGGGAA
1 7 | TTGCTGCAAGCTCAG
1 7 | TCTGATTTTACACA
1 7 | TCTCAGAGAAACGTG
```

# Primární analýza dat



## Pre-procesování dat

**Filtrování:** Čtení jsou z dat filtrována na základě kvality base call a délky čtení. Báze s nízkou intenzitou mohou vést k detekci falešně pozitivních variant, proto je třeba je odstranit. Čtení, která jsou příliš krátká, se pravděpodobně zarovnávají k více oblastem v genomu a způsobují špatné mapovací metriky.

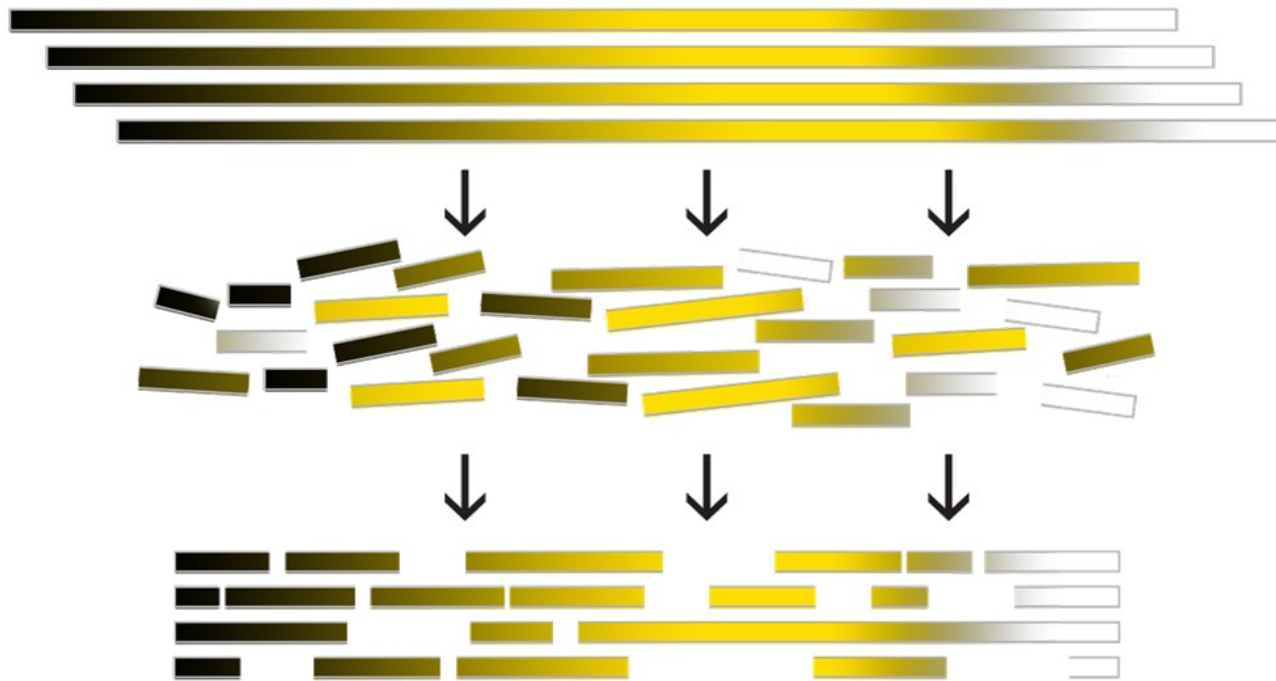
**Demultiplexování:** Multiplexování v NGS znamená sekvenování více vzorků současně na stejném přístroji. Demultiplexování znamená rozdělení sekvenačních čtení do samostatných souborů podle indexu „čárového kódu“ použitého pro každý vzorek.

**Ořezávání:** Adaptorové sekvenční ligované na konce knihoven během procesu přípravy knihovny je třeba ze sekvenačních čtení odstranit, protože mohou narušovat mapování a sestavování



# Sekundární analýza dat

*“Imagine a book cut by scissors into 10 million small pieces. Assuming that 1 million pieces are lost and the remaining 9 million are splashed with ink... try to recover the original text!” [P.Pevzner, UCSD]*



ATGTTCCGATTAGGAAACCTATACTGCATTTCAGTAAACG

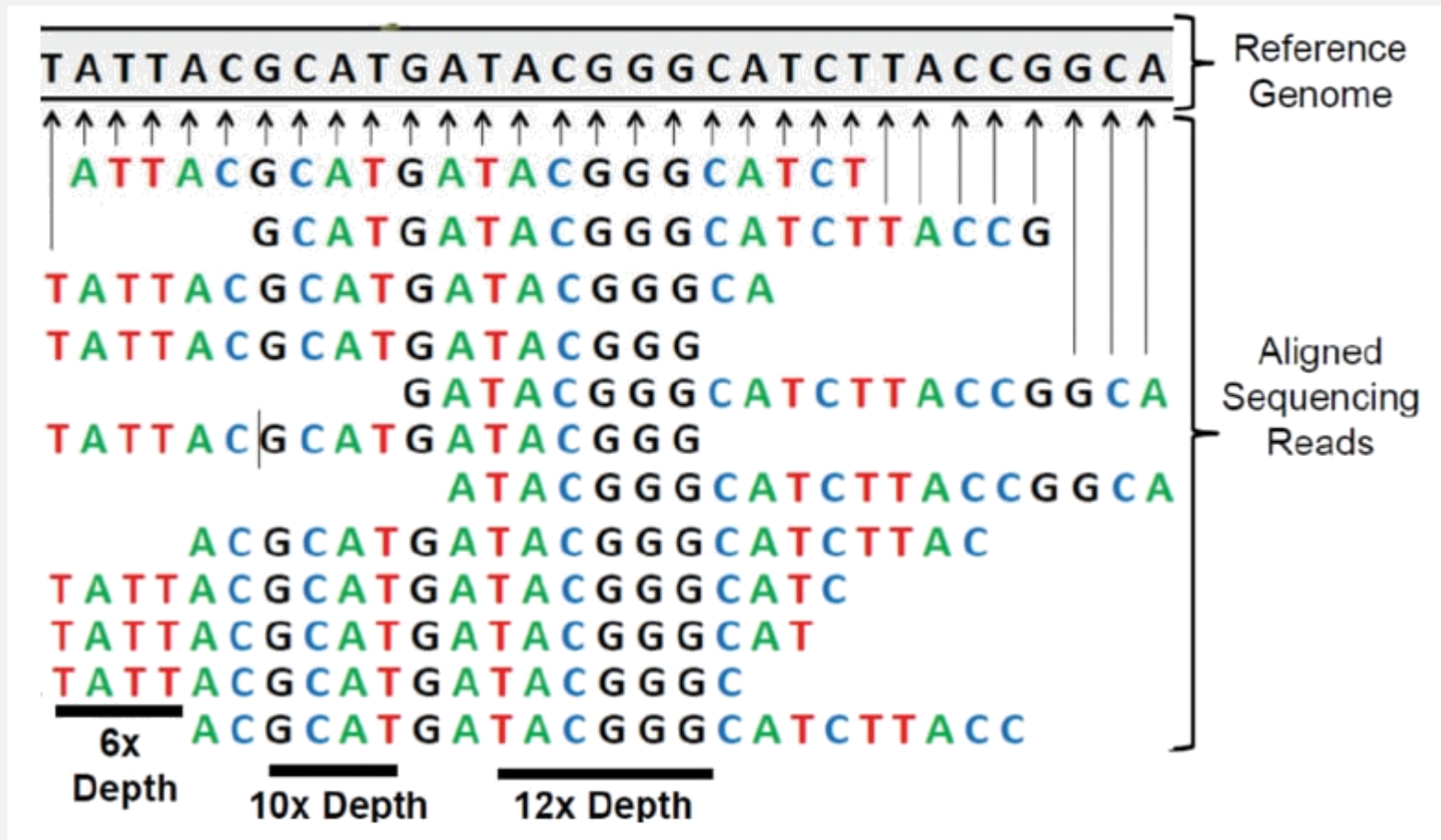
## De novo sestavení

Čtení se zarovnávají navzájem na základě jejich sekvenční podobnosti, aby se vytvořila dlouhá konsenzuální sekvence nazývaná kontig.



## Sestavení na základě referenčního genomu

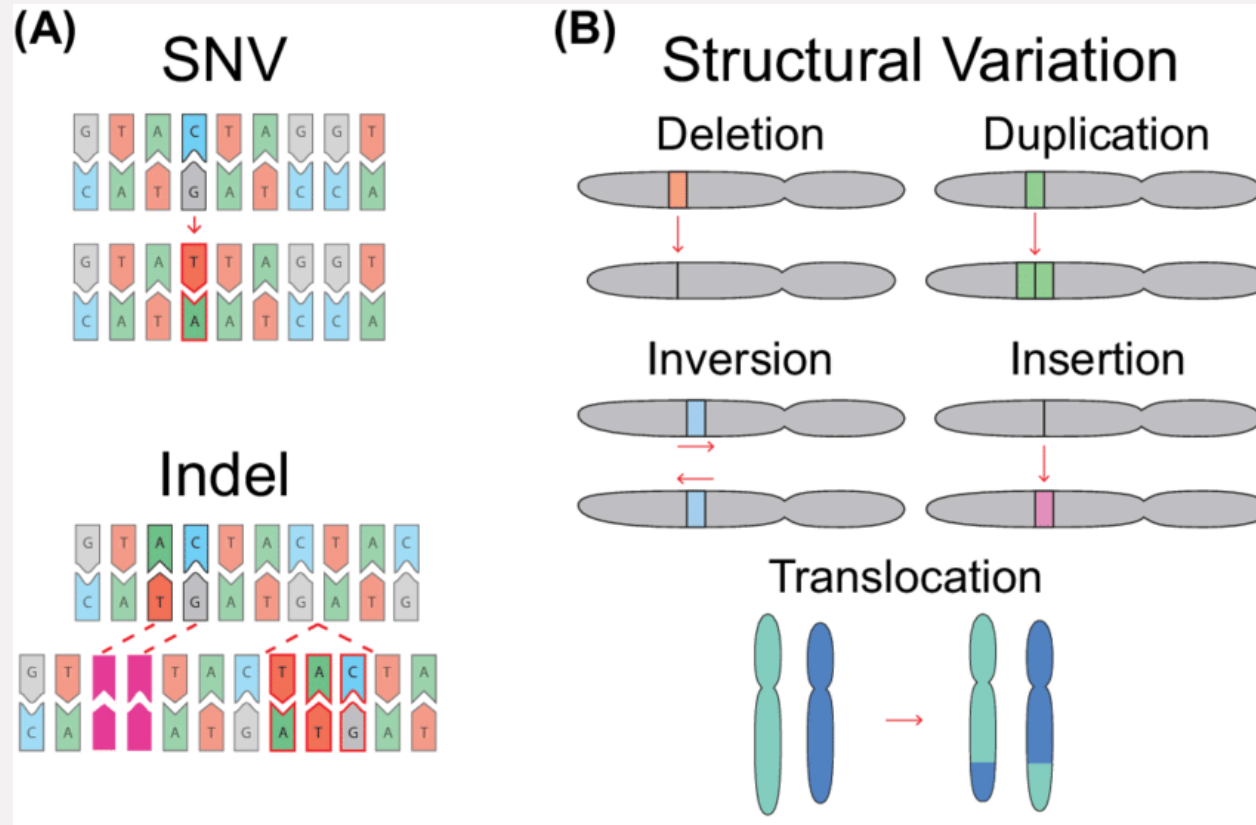
# Sekundární analýza dat



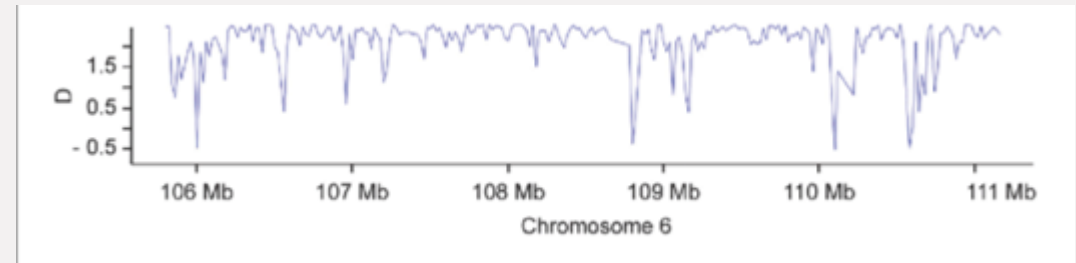
Cílem zarovnání sekvencí je najít místo, odkud čtení pochází, a určit, kolik čtení se k tomuto místu zarovnálo.

# Sekundární analýza dat

## Identifikace variací



Referenční genom je standard, dokážeme identifikovat odchylky.



Strukturální variace na úrovni chromozomu

# Sekundární analýza dat

## Používané formáty souborů

FILE TYPE	DESCRIPTION	WHERE IT IS USED
FASTQ	Text-based file format containing raw sequence reads and the associated quality score of each base	Storage of raw sequence data and input into sequence alignment
BED	Browser Extensible Data file is a tab-delimited text file that is used to store genomic regions as coordinates	In variant calling pipelines to direct the analysis to a genomic region
SAM	Sequence Alignment Map file, used to store text-based information for reads aligned to a reference sequence	Store information on read alignment, e.g. position and quality
BAM	Binary Alignment Map file is a compressed binary version of a SAM file. Can be opened in genome browsers to view read alignment	Used for input into variant calling pipelines
VCF	The Variant Call Format is a text file which stores sequence variants, each variant occupies a single row	Generated by variant calling pipelines. Used as input into variant annotation

# Terciální analýza dat

## Biologický kontext získaných dat

**Anotace** variant je proces předpovídání biologického vlivu nebo funkce genetických variant. Využívají se anotační nástroje, které pracují s VCF soubory, výstupem je zpráva o anotovaných variantách a jejich biologickém účinku.

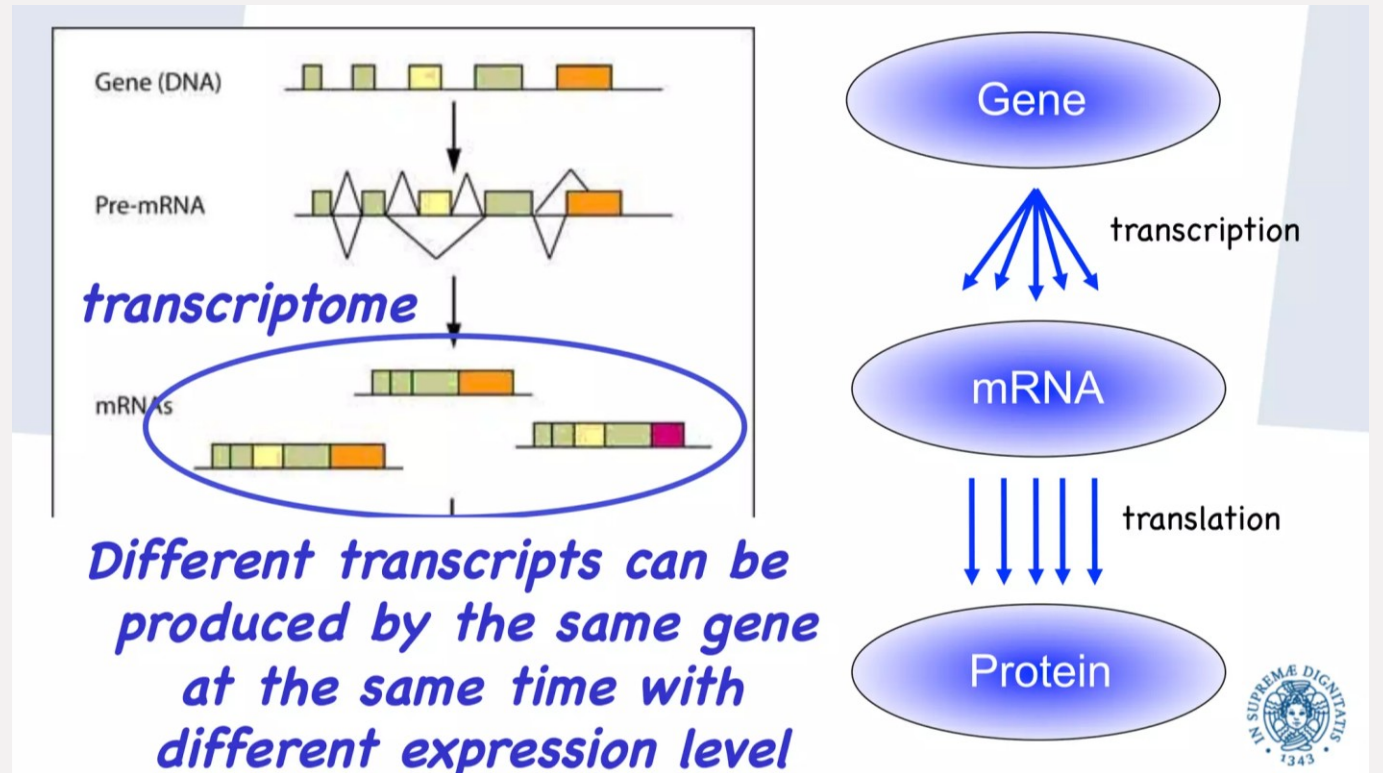
**Interpretaci** variant obvykle provádí kvalifikovaná osoba, například klinický genetik. Jejich práce zahrnuje shromáždění všech dostupných informací o pacientovi, včetně rodinné anamnézy onemocnění, a porovnání genotypu pacienta s klinickým fenotypem.

Součástí je také **statistická** analýza a **vizualizace** dat.

# Další aplikace NGS

## Genová exprese

Proces, kterým je v genu uložená informace převedena v reálně existující buněčnou strukturu (RNA) – množství transkriptu je úměrné míře genové exprese genů. Jedná se o reakci organismu na vnější i vnitřní vlivy.



U sekvenace RNA se provádí zpětný přepis do DNA, tzv. reverzní transkripce.

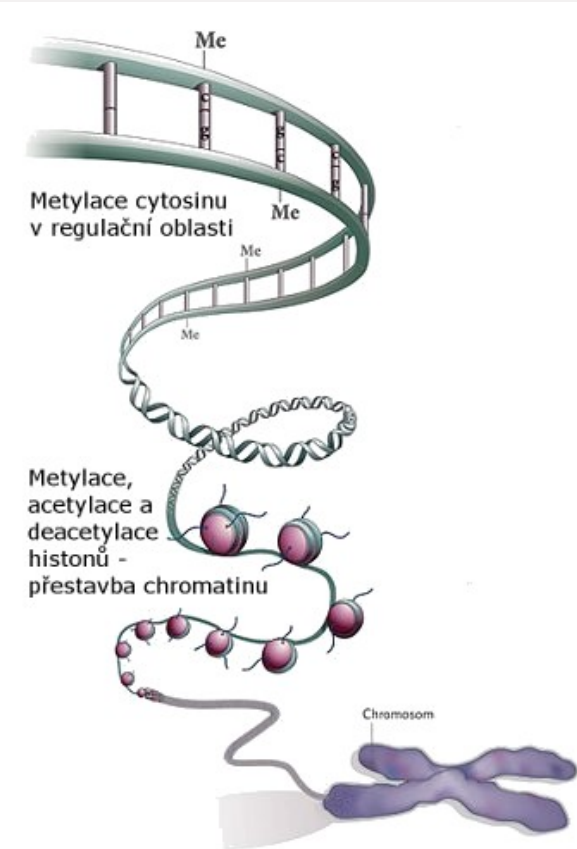


# Další aplikace NGS

## Epigenetika

Změny v genové expresi, které nejsou způsobeny změnou nukleotidové sekvence DNA.

Methylace DNA, acetylace histonů, **mikroRNA**, ...



Journal of Cellular and Molecular Medicine

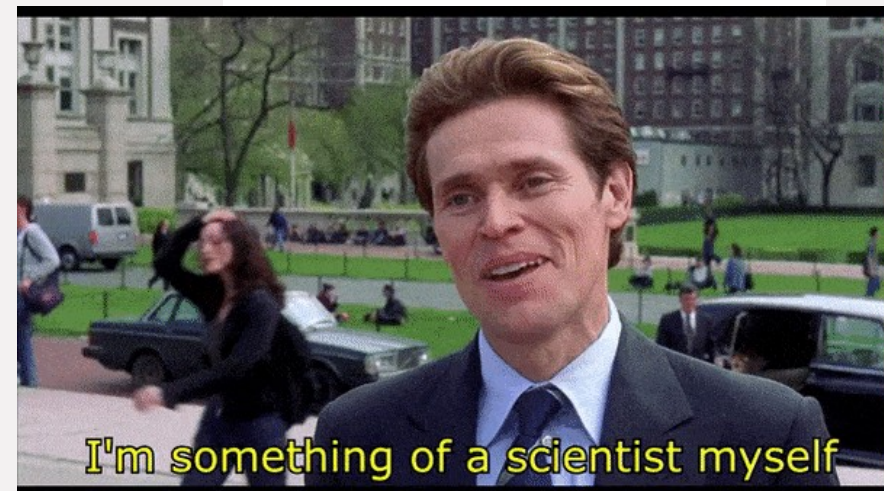
Open Access

Short Communication | Open Access |

### Genome-wide identification of urinary cell-free microRNAs for non-invasive detection of bladder cancer

Jaroslav Juracek, Barbora Peltanova, Jan Dolezel, Michal Fedorko, Dalibor Pacik, Lenka Radova, Petra Vesela, Marek Svoboda, Ondrej Slaby , Michal Stanik

First published: 24 January 2018 | <https://doi.org/10.1111/jcmm.13487> | Citations: 28



# Shrnutí na závěr

- Genomická data jsou unikátní a vysoce citlivá data, která do určité míry předpovídají délku života a dispozice k onemocněním jedince
- Znalost primární struktury (pořadí nukleotidů) umožňuje získat informace o umístění, struktuře i funkci genů a dalších oblastí DNA
- Zlatá metoda analýzy primární struktury NK je Sekvenování nové generace (NGS)
- Součástí procesu je jak práce v laboratoři (wet lab), tak bioinformatické zpracování dat
- NGS generuje velké množství dat (TB / genom)

DĚKUJI ZA  
POZORNOST!

