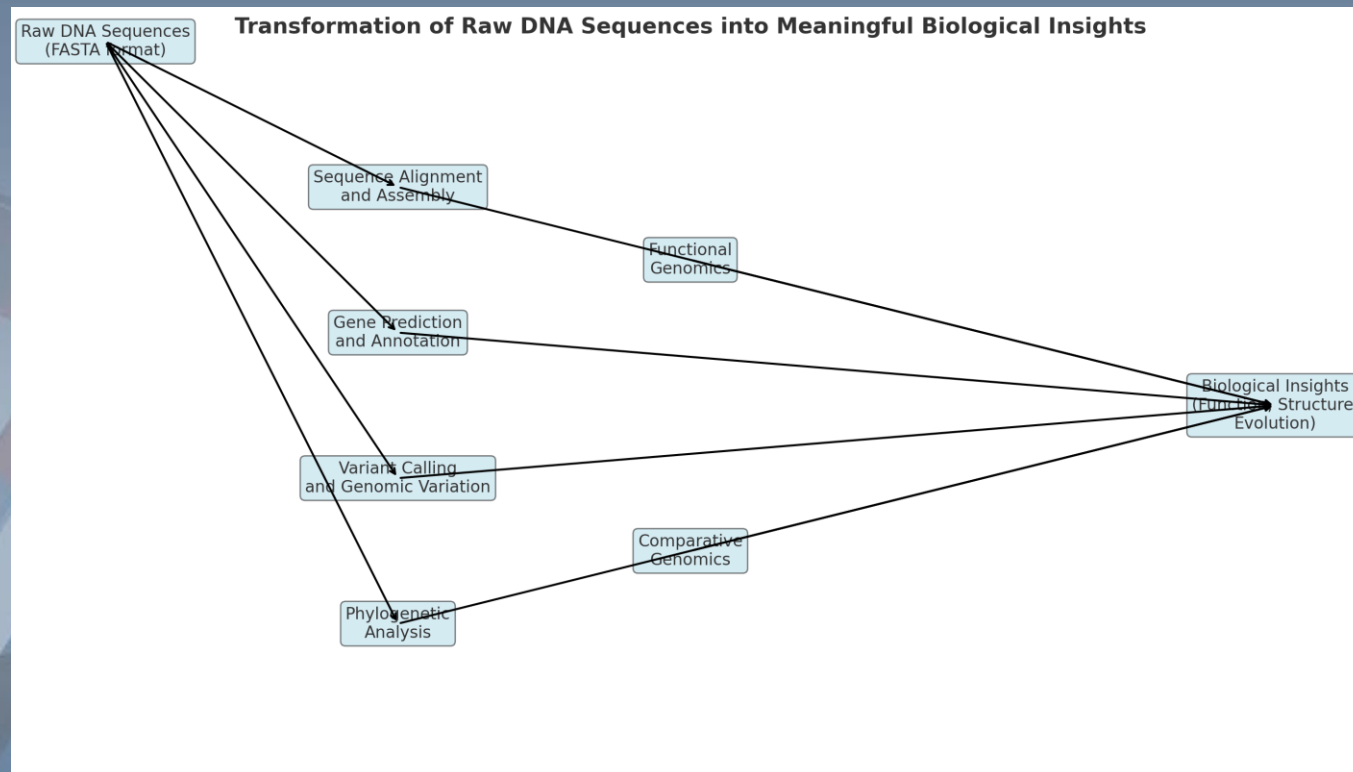
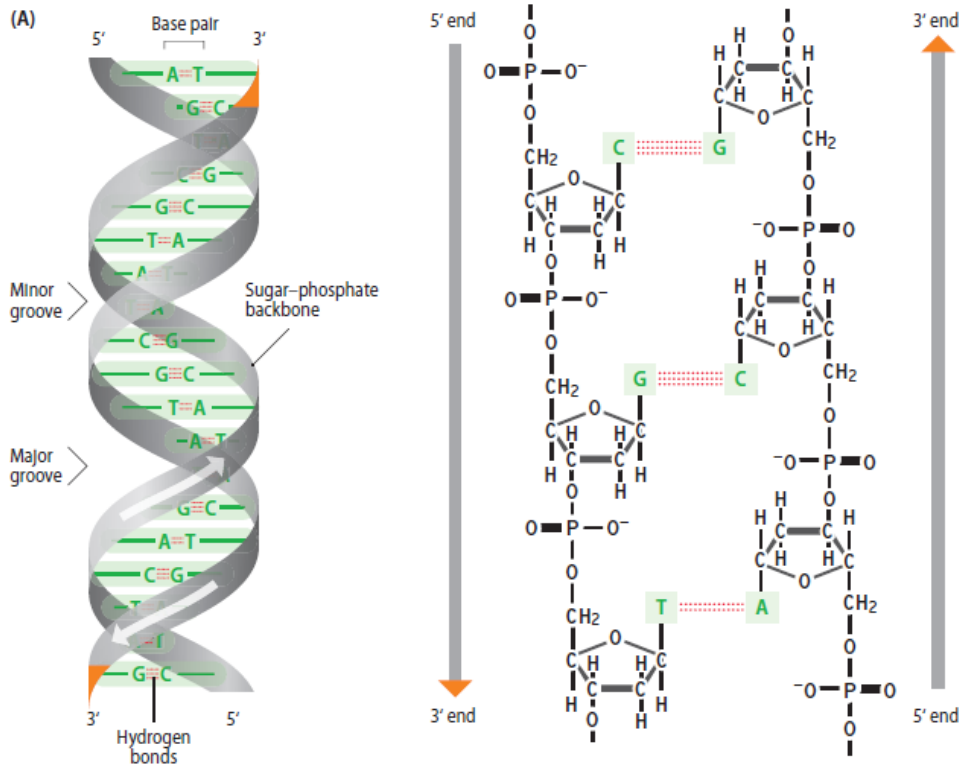


# Analiza sekvenci DNA pomoću bioinformatičkih alata



# Struktura molekule DNA

- povezana s funkcijom nositelja genetske informacije u živim organizmima

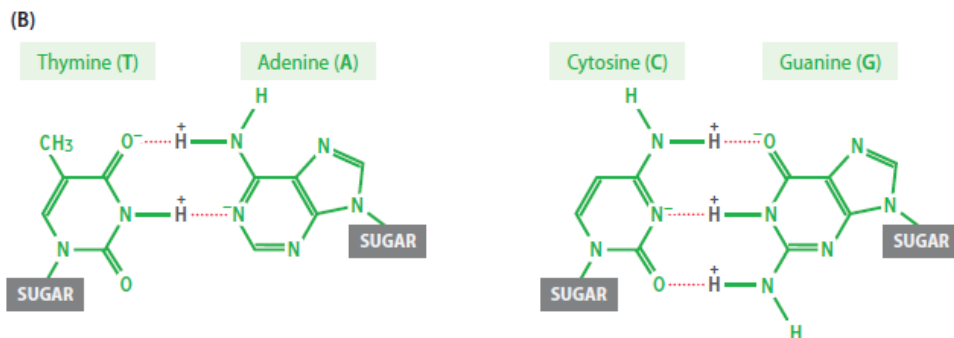


- struktura dvostruke zavojnice
  - dva isprepletena lanca koji se drže zajedno vodikovim vezama između komplementarnih nukleotidnih baza (parovi adenina s timinom i parovi citozina s gvaninom)

- specifično sparivanje baza u DNA (A-T i C-G)
  - ključno je za vjerni prijenos genetskih informacija s jedne generacije na drugu
  - omogućuje točnu replikaciju DNA, osiguravajući da svaka nova stanica dobije identičnu kopiju genetske informacije

- slijed nukleotida duž molekule DNA
  - kodira genetske informacije koje određuju karakteristike organizma
  - geni su specifične sekvence nukleotida koji kodiraju proteine ili funkcionalne RNA molekule, koje su bitne za razne biološke procese

- genetički kod
  - skup pravila prema kojima se informacije kodirane u DNA prevode u proteine
  - redoslijed nukleotida u DNA određuje redoslijed aminokiselina u proteinu, koji zauzvrat određuje strukturu i funkciju proteina



- struktura kromosoma
  - u eukariotskim stanicama DNA je organizirana u kromosome koji se sastoje od dugih niti DNA omotanih oko proteina koji se nazivaju histoni
  - pomaže u zbijanju DNA i reguliranju pristupa genetskim informacijama pohranjenim u molekuli DNA

# Genetički kod

- strukturiran na način da svaki od 64 mogućih tripleta RNA (kodona) sastavljenih od kombinacija 4 nukleotida (adenin, uracil, gvanin, citozin) odgovara jednoj od 20 biološki aktivnih aminokiselina, uz dodatna 3 „stop” kodona koji terminiraju translaciju proteina (Crick, 1968.)

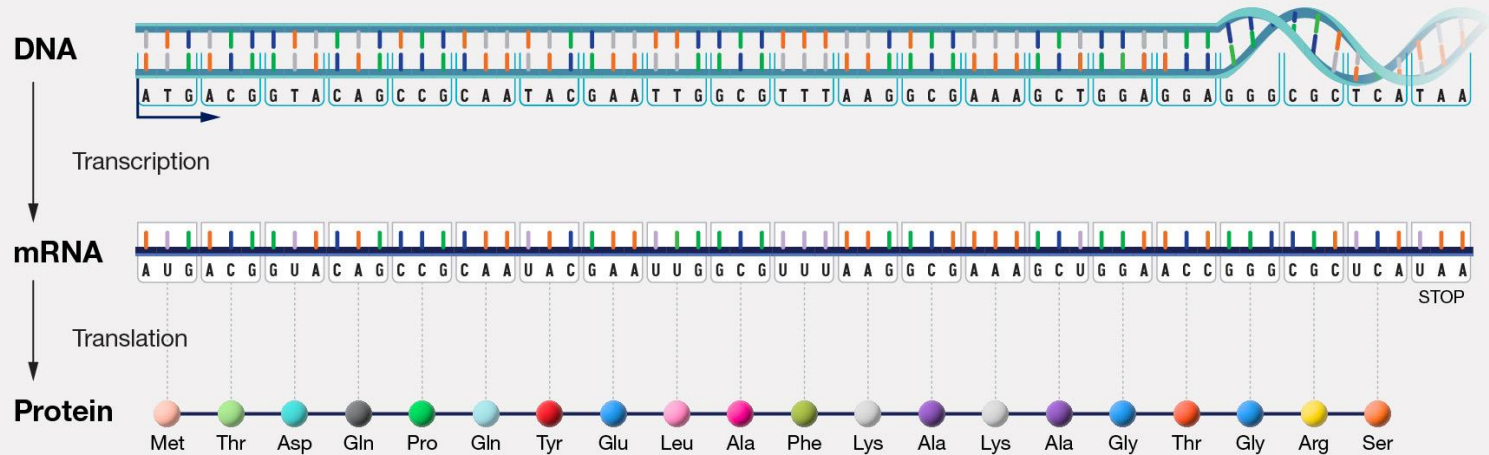
UUU	Phe	UCU	Ser	UAU	Tyr	UGU	Cys
UUC		UCC		UAC		UGC	
UUA	Leu	UCA		UAA	Stop	UGA	Stop
UUG		UCG		UAG		UGG	Trp
CUU	Leu	CCU	Pro	CAU	His	CGU	Arg
CUC		CCC		CAC		CGC	
CUA		CCA		CAA	Gln	CGA	
CUG		CCG		CAG		CGG	
AUU	Ile	ACU	Thr	AAU	Asn	AGU	Ser
AUC		ACC		AAC		AGC	
AUA		ACA		AAA	Lys	AGA	Arg
AUG	Met	ACG		AAG		AGG	
GUU	Val	GCU	Ala	GAU	Asp	GGU	Gly
GUC		GCC		GAC		GGC	
GUA		GCA		GAA	Glu	GGA	
GUG		GCG		GAG		GGG	

# Središnja dogma molekularne biologije

Francis Crick (1958)

- DNA sadrži upute za stvaranje proteina, koje kopira RNA
- RNA zatim koristi upute za stvaranje proteina

DNA -> RNA -> bjelančevine



# Sekvenciranje DNA

- proces određivanja točnog redoslijeda nukleotida unutar molekule DNA



<https://www.novogene.com/eu-en/wp-content/uploads/sites/7/2021/02/novogene-blog-ngsbeginnersguide-20210203.jpg>

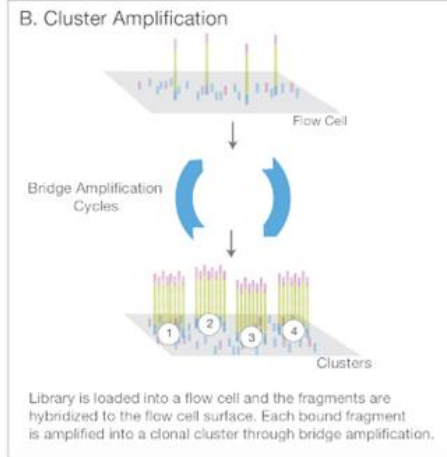
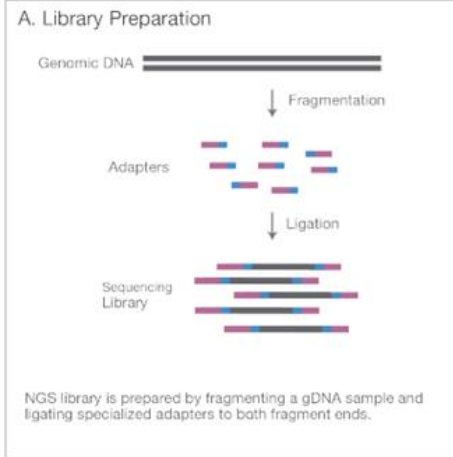


# Sekvenciranje sljedeće generacije

engl. *Next Generation Sequencing*, NGS

## 1. Priprema uzorka

- ekstrakcija DNA iz stanica organizma
- fragmentacija ekstrahirane DNA u manje dijelove, obično duge stotine do tisuće parova baza.
- priprema knjižnice (engl. *library preparation*) - kratke sekvence adaptera dodaju se na krajeve fragmenata DNA kako bi se fragmenti vezali na platformu za sekvenciranje

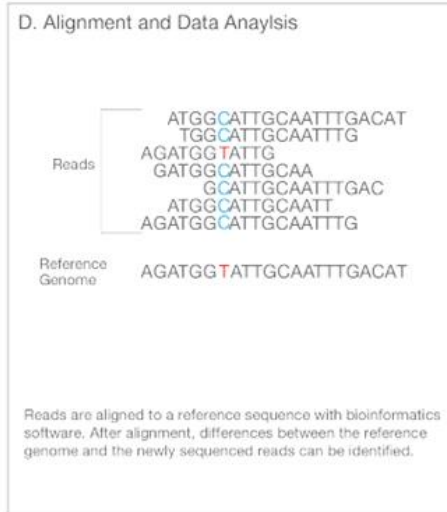
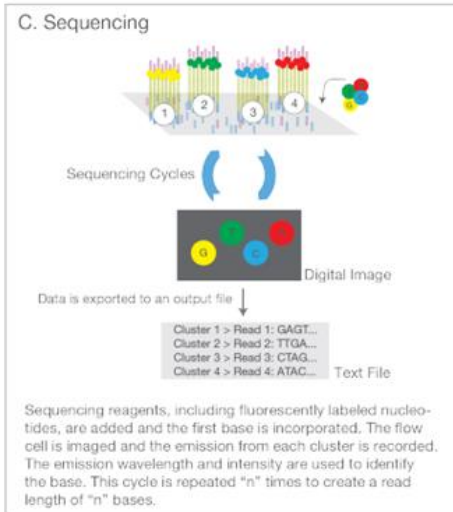


## 2. stvaranje klastera

- fragmenti DNA pričvršćuju se na protočnu stanicu (engl. *flow-cell*) i umnožavaju dajući nakupine (engl. *clusters*) identičnih sekvenci

## 3. sekvenciranje sintezom

- nukleotidi se dodaju jedan po jedan, a svaki dodatak detektira se signalom (npr. fluorescencijom) koji identificira nukleotid ugrađen na svakom položaju



## 4. analiza podataka

- određivanje očitanih baza (engl. *base calling*) - signali otkriveni tijekom sekvenciranja prevode se u niz nukleotida (A, T, C, G)
- sastavljanje očitanih sekvenci (engl. *reads*) - kratke sekvence sastavljaju se u duže sekvence, usklađujući ih s referentnim genomom ako je dostupan ili ih sklapaju *de novo* ako referenca nije dostupna



# FASTA datoteka

- sadrži retke nukleotidnih sekvenci (adenin [A], timin [T], citozin [C] i gvanin [G]) kojima prethodi redak zaglavlja koji počinje simbolom ">"

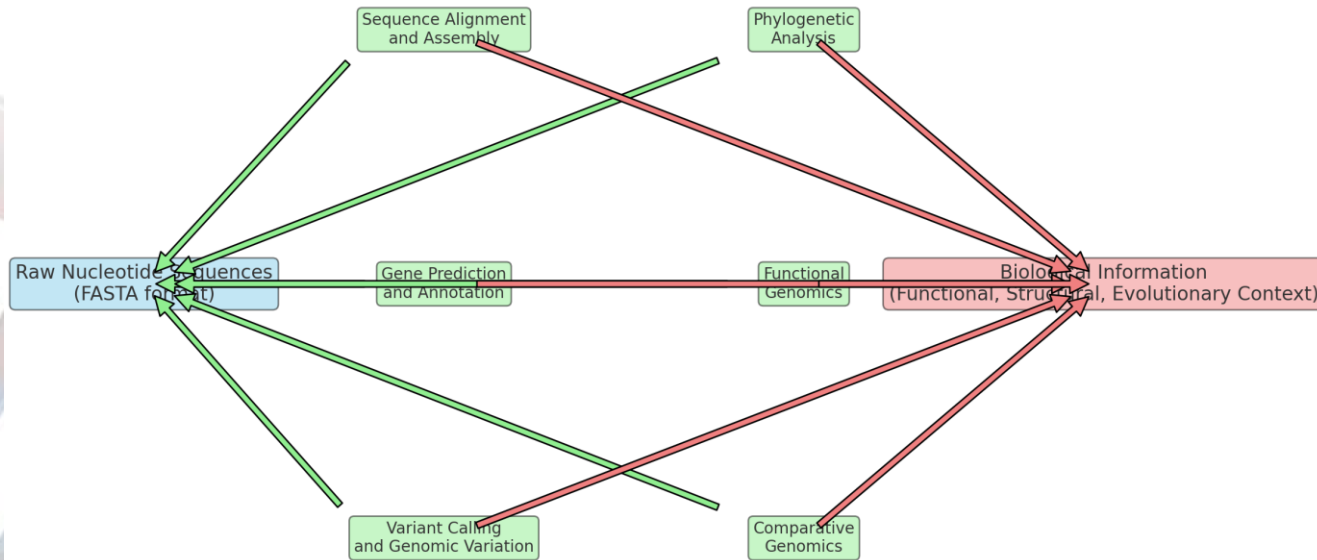
```
>10001f
GTGTGGTAAAAAAGTAGGTAAGAGCAGTTAAAATGCAGAATCCACGTTTCAGTTGACTAAAATTCCTGAGTCAGTCGACTGTAAGTAGAACCCGTGAGTTTTAACATGCATTTCAGTTGAC
>10001r
CTGAGTTAGTCGACTGAATTTTCATAGTAAAATCTCTCGGGTTCAGCTTAGTAAACAGGAAAACTGAGATATTC AATTGAGTTTCATATTAAGCTCTCTGTTTTCGGCTCAGTCGACTGA
>10002f
TTCGGGTTGTTACATTCTCCCTTCTTAAGAGATTTTCGTCGTCGAAATCTCTAAGGATAAAAAATGAGCTAGATATTTTCATCCAATATCTACACCCAACACTCAATCTTCAACCATT
>10002r
ATGATGTTTGAAGCATGCATATAGAAGTCATCATTAGGAATATTTGCATTCCTTTGGATACATAGCATGGGCATTATTGCATAAAAAGTGGTCAATTAATGAGTGGACTGCCTATCTTTT
>10003f
AATAATTTGCTCAATACCGATTTTTCTTAAGAGTTTCAACATTCATTTCTTAAGATTTACTATAAATGGTTACACACTGCTTATACTGATTTTTCACTTCTCTATTGAAGTTCATATTCATT
>10003r
AACTTTTGAAATTTTCATTTAACGACAAATAATTTAAATTTGTTTAGCAAATGAAATTGAAATTGGATATTGGAATCTAATTCAAATGTGAATGTAAGGATAAGTGAAGTGAGAAAACTAA
>10004f
AGGATCGAAAAGTAAATTGGCAACGGCTGTTTCATAAATTAATAAGAAATAATCAACTAAAATTTATAAATTTAGAATTAGAACACAAAAAGCACAGTACTTTTCTTTTCAAGTTTCAGGG
>10004r
ATCTTCCAATGCAAAATCGCGATACAACAAAAGAAAAAAAACAAACCAAGCGGAAGCACCACCCACAGCGATCAACGATGCCTCATGACTGATTGGCGGTAAGCGTATGTGGCTGCG
>10005f
TGGGCCGATAAGTGTGAGGAAGCATTCCAAGGATTGAAACACTTATTGACTAATGCACCATGTTAGTTGTTCCAGAAGGGAATCTGGACTTAGTAATGGATACAAATGCTAGTGGAACTG
>10005r
AAATCTTTCCACACAACACTATCATTGCCAACACTGCATTTTGACATTTTCGGCTAAGTGCATCTGCAACATTATTGTCTTTTCTGGAGTGTGATCAATACAAAAATCATAAGTGGCCATCA
>10006f
TTTAACTAATTTTATCTCAAAATTTTATGTGTATACATTTTATATCATGGTTTATATCTTTATACTATATTATATACACATATATGTATGTATGTATGATTTCAATAATTTGATTTTTA
>10006r
CATACATATATATTATTTTAAATCGAAAAATAAATATTTTATTACAAAAACAATAGTACATACACATACACATTTTGATATATATATATATATATATTATGTGTGTGGTTACTCATTATT
>10007f
TTCGGATAGCCTTTTTTTGGGTGTGAAACTATCCGTTGATGGCTTACAGCGGTCAATGCCGTAGAACTCACTGAGACCTTCGTTTTGACTGGTTATGGCTCCCGTAACTCTAAAAGGGGCA
>10007r
GTTACGGGAGCCATAACCAGTCAAACGAAGGTCTCAGTGAGTTCTACGGCATTGACCGCTGTAAGCCATCAACGGATAGTTTCACACCCAAAAAGGGCTATCCGAACCCCAATCGGGTC
>10008f
TGACGAGGCCGGTAGCGGGCTGTAGTGGCGGCGTGGTTGGTCAAACGTTGATGAAGTGGGTTTCGGTTTTGTTTCGCTGAAGAGTGTTCGCTGATGCTATAACTACTGTTCTTGTGAGGGAG
>10008r
ATTTAAATCGAATTAATAATCAATTTAAACCCAAAATTCGACGGTCTCAGAATCCGTCCAAACAGCATCATCAAGTCCAGAATTTCAATACGAACCTAGCGGAGACAACCGCGTGTAA
>10009f
GGCTACGCCGATGATATAAGGCAATTTTATTATTATCATCACTTTGGTGTGATTTCTCTTTGGAGATTTTTACTGTTGATTGATGATTAGGGATCGTCATTTTATTATGTGGTCCAAAA
>10009r
CAAAGCCACTGCCACACATGTGAGTCAAAAATATGTACATATCAGAGTCTGAATAATAACAGAGTAGCGGAATAGTGGACAAAATAAAAGCGCTATCTATCAGCATAGCAGATAATGTCCA
>10010f
AGGTTACTTTACAAAATATAAACTAAATTTCTCTTTTAAAACAAAATCCAAGTCACCTTAAAATGAGTACTATTATAATAAATTTGAAAAGAATTACTTTAAAATTTAAACATTTAA
>10010r
```



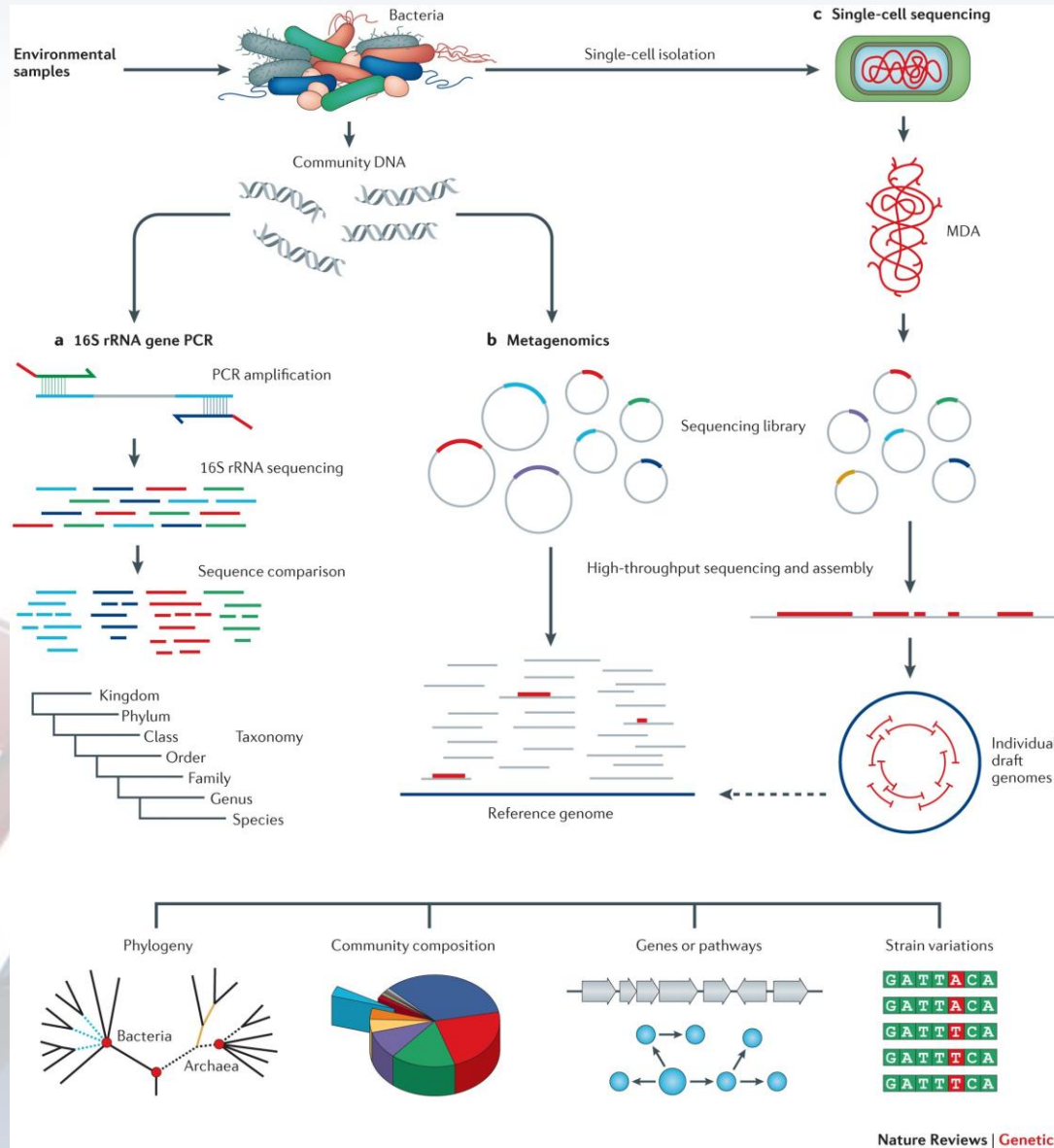
# Bioinformatika

- integrira biologiju, računalnu znanost i informacijsku tehnologiju za analizu i interpretaciju bioloških podataka
- transformira neobrađene nukleotidne sekvence (engl. *raw data*) u složenu kompoziciju bioloških informacija, omogućujući znanstvenicima da dešifriraju funkcionalni, strukturni i evolucijski kontekst genoma
- kroz sofisticirane algoritme i računalne alate, pretvara sekvence u smisleni prikaz uputa za rast, razvoj, preživljavanje i reprodukciju organizma

## Transformation of Raw Nucleotide Sequences into Biological Information



# Uloga bioinformatike u analizi sekvenci DNA i sastavljanju genoma

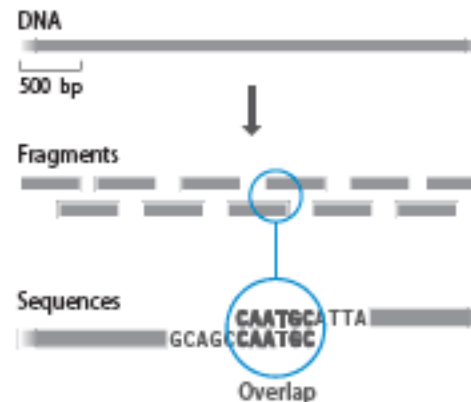


[https://media.springernature.com/full/springer-static/image/art%3A10.1038%2Fnrng3785/MediaObjects/41576\\_2014\\_Article\\_BFnrng3785\\_Fig1\\_HTML.jpg](https://media.springernature.com/full/springer-static/image/art%3A10.1038%2Fnrng3785/MediaObjects/41576_2014_Article_BFnrng3785_Fig1_HTML.jpg)

# Sklapanje genoma (engl. *sequence assembly*)

Kako se mnoštvo kratkih očitanih sekvenci generiranih metodama NGS sastavlja u sekvencu cijelog genoma?

- kontinuirana sekvenca genoma se sastavlja traženjem preklapanja između slijeda nukleotida zbirke pojedinačnih fragmenata
- proces se provodi pomoću specijaliziranog softvera i algoritama koji analiziraju preklapajuća područja sekvenci kako bi ih poravnali i spojili u veće susjedne sekvence, poznate kao kontig sekvence (engl. *contig*) - djelomično sastavljen niz od nekoliko očitanih sekvenci



T.A. Brown - Genomes 4-Garland Science (2018)

- konsenzus sekvenca - reprezentativni slijed nukleotida u kojem je svaki nukleotid onaj koji se najčešće pojavljuje na tom mjestu u različitim sekvencama

```
TAGATTACACAGATTACTGA TTGATGGCGTAA CTA  
TAGATTACACAGATTACTGACTTGATGGCGTAACTA  
TAG TTACACAGATTATTGACTTCCATGGCGTAA CTA  
TAGATTACACAGATTACTGACTTGATGGCGTAA CTA  
TAGATTACACAGATTACTGACTTGATGGGGTAA CTA
```

↓ ↓ ↓ ↓ ↓

```
TAGATTACACAGATTACTGACTTGATGGCGTAA CTA
```

# de novo sklapanje genoma

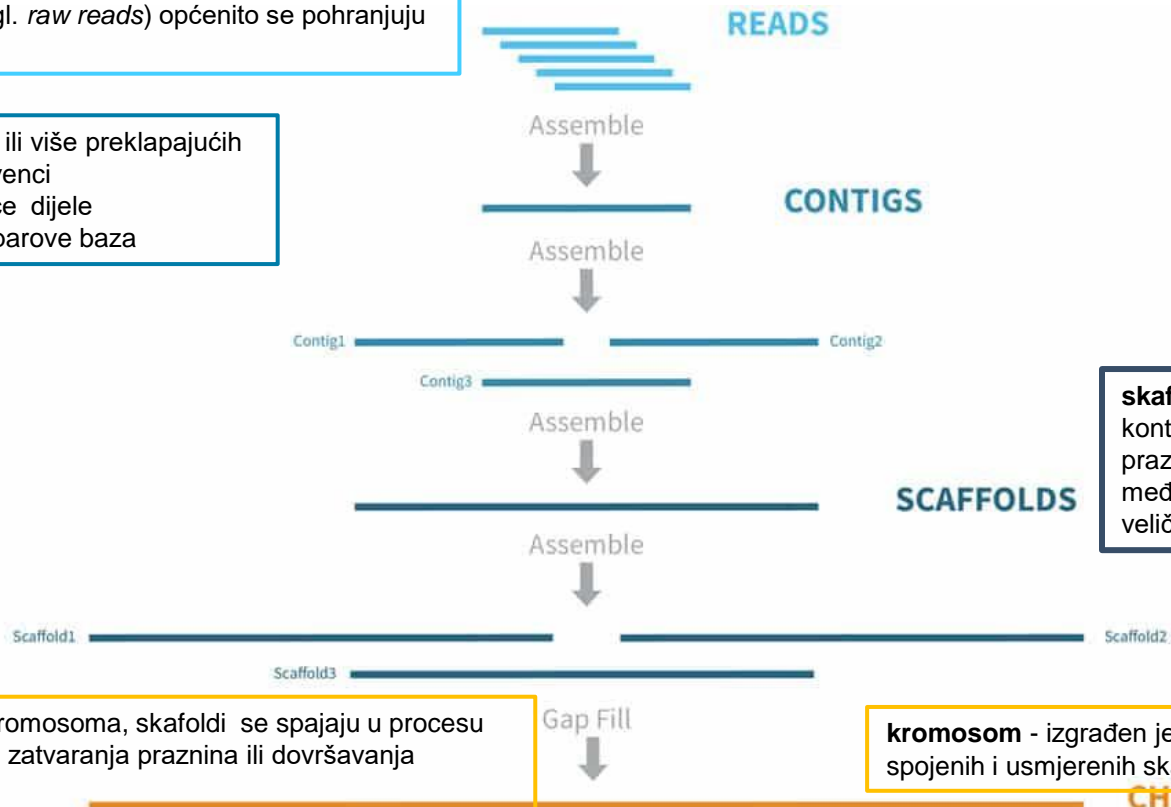
- metoda za konstruiranje genoma iz velikog broja (kratkih ili dugih) fragmenata DNA, bez prethodnog znanja o ispravnom slijedu ili redoslijedu tih fragmenata

## očitanje sekvenca DNA (engl. reads)

- fragменти DNA koja generiraju sekvenatori - obično se kreću u veličini od 35 do 1000 bp ("kratka čitanja") ili "duga čitanja" u veličini od 1000 do 500 000 bp
- neobrađena očitavanja (engl. raw reads) općenito se pohranjuju u FastQ datoteke

## kontigi

- građeni od dvije ili više preklapajućih usmjerenih sekvenci
- očitanje sekvenca dijele podskup ili sve parove baza



**skafoldi** - sastoje se od više kontiga (među kojima često ima praznina) i definiraju njihov međusobni redoslijed, orijentaciju i veličinu procijepa među njima

- u koraku sastavljanja kromosoma, skafoldi se spajaju u procesu popunjavanja praznina, zatvaranja praznina ili dovršavanja genoma

**kromosom** - izgrađen je od dva ili više spojenih i usmjerenih skafolda



# Sastavljanje genoma *de novo*

- metoda za konstruiranje genoma iz velikog broja (kratkih ili dugih) fragmenata DNA, bez prethodnog znanja o ispravnom slijedu ili redoslijedu tih fragmenata
- odnosi se na sekvenciranje novog genoma gdje ne postoji referentna sekvenca dostupna za poravnanje
- pristup referentnog genoma se može koristiti za sekvenciranje *de novo* ako je vrsta čiji se genom sekvencira povezana s drugim vrstama čiji su genomi već sastavljeni
- postojeća sekvenca koristi se kao referentni genom za sastavljanje sekvenci tražeći regije identičnosti ili sličnosti sekvenci
- kratke sekvence novog genoma bit će sastavljenu u kontige prije usporedbe s referencom, kako bi se povećao stupanj točnosti na temelju preklapanja



T.A. Brown - Genomes 4-Garland Science (2018)

