

# BIOINFORMATIK

Föreläsare:  
Mats Nilsson  
mats.nilsson.6088@med.lu.se

= SÄTT ATT SAMLA INFO OCH GÖRA DEN PRAKTISKT ANVÄNDBAR FÖR ÄNDAMÅLET.

## DEFINITIONER

GÖRA EN FÖRESTÄLLNING OM MOLEKYLERS FYSISKA KEMI OCH TILLÄMPA INFORMATIONSTEKNIK FÖR ATT FÖRSTÅ OCH ORGANISERA INFORMATIONEN ASSOCIERAD MED DESSA MOLEKYLER. PÅ STOR SKALA.

*ENLIGT OXFORD ORDBOK.*

RESEARCH, UTVECKLING, TILLÄMPNING AV DATORISERADE VERKTYG OCH INSTÄLLNINGAR FÖR EXPANSION AV BIOLOGISK-, MEDICINSK-, BETEENDEMÄSSIG- ELLER HÄLSOMÄSSIG DATA. DENNA DATA INSAMLAS, LAGRAS, ORGANISERAS, ARKIVERAS, ANALYSERAS OCH VISUALISERAS.

*ENLIGT COMMITTEE NATIONAL INSTITUTE OF MENTAL HEALTH*

TVÄRVETENSKAPLIG DISCIPLIN DÄR ALGORITMER FÖR ANALYS AV BIOLOGISKA ( FFA. MOLEKYLÄRBIOLOGISKA) DATA UTVECKLAS. MAN ANVÄNDER OCKSÅ TERMEN SYNONYMT MED ENGELSKANS COMPUTATIONAL BIOLOGY.

*ENLIGT WIKIPEDIA*

(MAN VILL ÖKA FÖRSTÅElsen FÖR BIOLOGISKA PROCESSER. TAR FRAM ALGORITMER OSV. METODER FÖR ATT LOKALISERA GENER INOM GENSEKVENSER. SKAPA STORA DATABASER MED DENNA TYP AV INFORMATION.)

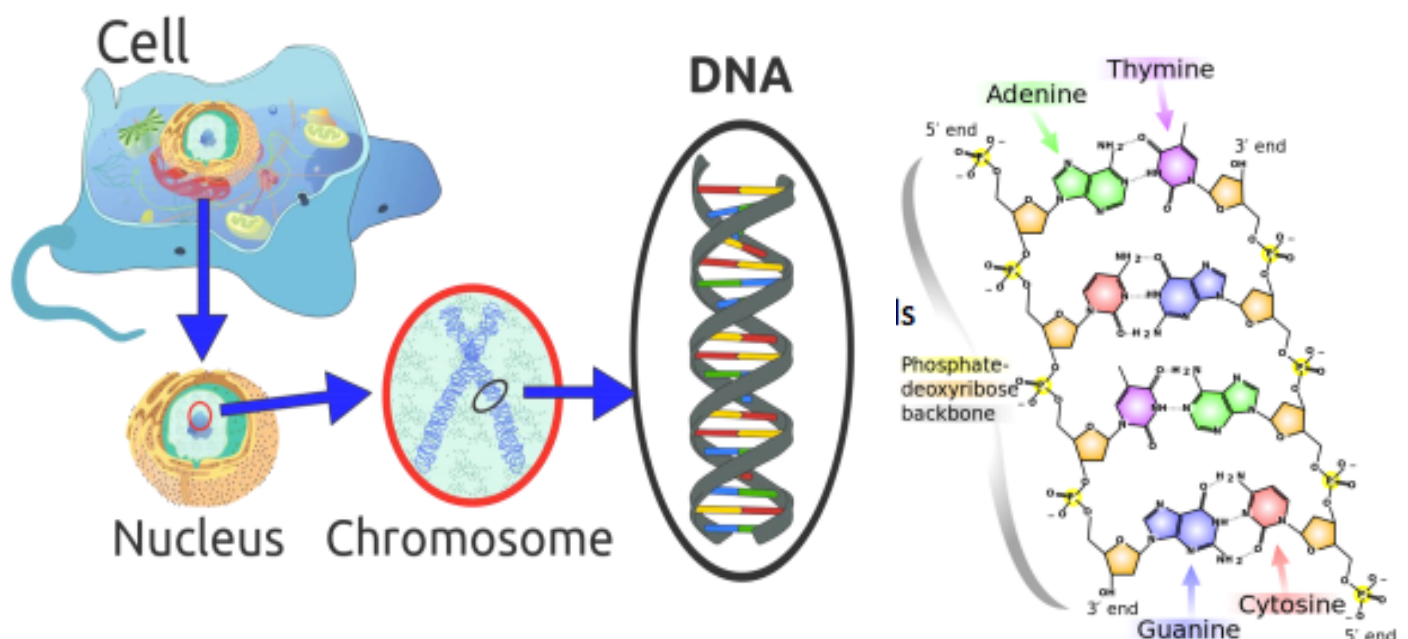
PRIMÄRA MÅLET ÄR ATT ÖKA FÖRSTÅElsen FÖR BIOLOGISKA PROCESSER.

- UTVECKLING AV NYA ALGORITMER, STATISTISKA MÄTMETODER OCH DATORPROGRAM FÖR UTVECKLINGEN AV STORA DATABASER.
- IMPLEMENTERING AV UTVECKLADE ALGORITMER, PROGRAM I DATAUTVÄRDERING OCH UTVÄRDERING AV RESULTATEN.
- SKAPA OCH FÖRBÄTTRA PUBLIKA TILLGÄNGLIGA DATABASER.

## GENETISK KOD – DNA

C-G, T-A.

STRÄNGAR BLIR KOPIOR AV VARANDRA. OM NÅGONTING BLIR FEL KAN ENZYM GÅ IN OCH RÄTTA TILL FELET.



## OM DNA

SJÄLVREPLIKERANDE MATERIAL SOM FINNS I NÄSTAN ALLA LEVANDE ORGANISMER SOM HUVUDSAKLIG KONSTITUENT AV KROMOSOMER. BÄRARE AV GENETISK INFORMATION.

BESTÅR AV TVÅ BIOPOLYMERSTRÄNGAR TVINNANDE KRING VARANDRA I EN DUBBELHELIX. BASER PARAS A-T OCH C-G. VÄTEBINDNINGAR SAMMANLÄNKAR KVÄVEBASERNA FRÅN DE TVÅ POLYNUKLEOTIDERNÄ → DUBBELSTRÄNGAT DNA.

### DET CENTRALA DOGMAT



BESTÅR AV CODON. VILKA BESTÅR AV TRE BASER. INNEHÅLLER G,A,U OCH C. VARJE CODON MOTSVARAR I SIN TUR EN AMINOSYRA VILKEN DEN "OMVANDLAS TILL" I TRANSLATIONEN. FINNS ÄVEN START- OCH STOPPSEKVENSER.

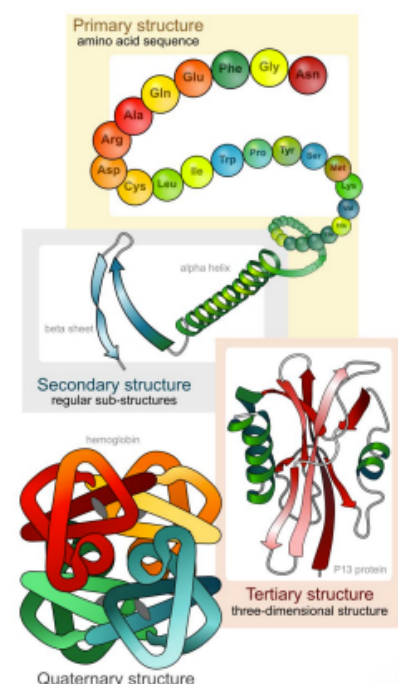
AMINOSYROR SOM TILLSAMMANS BLIR PROTEIN. FINNS TRE CODON SOM ENDAST KODAR FÖR START OCH STOPPSIGNALER FÖR DE GENETISKA MEDDELANDENA.

TITTAR PÅ ENSKILDA AMINOSYRORS SIDOKEDJOR. ADIABATISKA – KVÄVE OCH SYRE.

FINNS SYROR, BASER OCH POLÄRA. KODAS FÖR ATT VETA VAD DE KAN ERSÄTTAS MED. EN LADDAD ERSÄTTTS MED EN ANNAN MOTSVARANDE LADDAD.

### NIVÅER PÅ PROTEINSTRUKTURER

NIVÅ	BESKRIVNING	STABILISERAS AV
PRIMÄR	AMINOSYRASEKVENSER	PEPTIDBINDNINGAR
SEKUNDÄR	ALFA-HELIXAR OCH BETA-SHEETS I POLYPEPTID	VÄTEBINDNINGAR MELLAN GRUPPER LÄNGS PEPTIDSTRÄNGEN
TETRIÄR	TREDIMENSIONELLA FORMEN AV EN POLYPEPTID	INTERAKTION MELLAN OLIKA R-GRUPPER OCH PEPTIDLÄNGDEN.
KVARTENÄR	FORMATIONEN AV KOMBINERADE POLYPEPTIDER	INTERAKTION MELLAN R-GRUPPER OCH PEPTIDLÄNGDEN PÅ OLIKA POLYPEPTIDER.



INFORMATION SORTERAS I GRUPPER BASERAT PÅ DESS LIKHETER.

## HOMOLOG

INDIKERAR GENER ELLER PROTEINER SOM ÄR EVOLUTIONÄRT RELATERADE. INNEBÄR GENER MED GEMENSAMT EVOLUTIONÄRT URSPRUNG. ANTINGEN ORTOLOGA ELLER PARALOGA.

*”Man kan aldrig vara säker på att två nästan identiska organ, proteiner eller nukleinsyrasekvenser verkligen har sitt ursprung hos en gemensam anfäder, så termen används allmänt om organ etc. som är så lika varandra att det skulle vara mycket osannolikt att likheterna bara beror på slumpen. Märk att homologgi är en absolut egenskap – antingen är proteinerna, organen etc. homologa eller så är de det inte. Det är därför inte korrekt att tala om olika grader av homologgi, däremot kan man tala om olika grader av likhet. Två proteiner som är homologa kan alltså visa t.ex. 70 % sekvensidentitet, eller 96 % sekvensidentitet. Organ etc. som är homologa kan ha utvecklats i olika riktningar och se mycket olika ut, men förblir homologa, t.ex. våra armar och fåglarnas vingar som båda har sitt ursprung i framben”*

## OTROLOG

EN GEN SOM ÄR HOMOLOG TILL EN GEN HOS EN ORGANISM AV EN ANNAN ART.

DEN GEMENSAMMA STAMFORMEN FÖR ORTOLOGA GENER EXISTERADE INNAN DE UTVECKLINGSLINJER SOM LEDDE FRAM TILL RESPEKTIVE ART SKILJDES ÅR. EXEMPELVIS HEMOGLOBIN-B-GENEN HOS MÄNNISKA OCH MUS.

## PARALOG

EN GEN SOM ÄR EN HOMOLOG TILL EN GEN HOS SAMMA ORGANISM. HAR UPPKOMMIT GENOM ATT URSPRUNGLIG GEN HAR DUPLICERATS SÅ ATT DE TVÅ KOPIORNA UTVECKLATS OBEROENDE AV VARANDRA OCH FÅTT OLIKA FUNKTIONER. EXEMPELVIS HEMOGLOBINER OCH MYOGLOBINER HOS MÄNNISKAN.

## KLASSIFIKATION

### 1. MÖNSTERIGENKÄNNING

IGENKÄNNING AV SÄRSKILD SEKVENS ELLER STRUKTUR SOM KAN ASSOCIERAS MED SÄRSKILD KARAKTÄRISTIK.

### 2. PREDIKTION

#### ○ HOMOLOGIBASERAD PREDIKTION

FÖR OKÄNT PROTEIN IDENTIFIERAS DE MED HJÄLP AV ATT IDENTIFIERA HOMOLOGER MED KÄND FUNKTION ELLER STRUKTUR. MED HJÄLP AV INFORMATIONEN OM DEN KÄNDA HOMOLOGEN KAN MAN FÖRUTSPÅ STRUKTUREN/FUNKTIONEN HOS DET OKÄNDA PROTEINET.

#### ○ AB INITIO (DE NOVO) PREDIKTION

”FROM THE BEGINNING”. PREDIKTION GÖRS GENOM ANVÄNDNING AV DATORISERADE MODELLER UTAN DIREKT JÄMFÖRELSE MED EXISTERANDE DATA

#### ○ SEQUENCE ALIGNMENT

IDENTIFIKATION GENOM ATT IDENTIFIERA LIKANDE REGIONER.

LÄGGER UPP BREDVID VARANDRA. HITTAT SEKVENSER SOM DELVIS LIKNAR VARANDRA.

##### ▪ GLOBAL ALIGNMENT

ARRANGERAR VARJE REST I VARJE SEKVENS. BRA FÖR SEKVENSER MED I GROVA DRAG SAMMA STORLEK OCH LIKHET.

##### ▪ LOCAL ALIGNMENT

ANVÄNDBAR FÖR OLIKA SEKVENSER INNEHÅLLANDE LIKA REGIONER ELLER LIKA SEKVENSTEMAN INOM DE STÖRRE SEKVENSKONTEXTEN.

#### ○ PARWISE SEQUENCE ALIGNMENT METHODS

##### ▪ DOT-MATRIX METHOD

ENKELT VIS ATT EVALUERA LIKHETER MELLAN TVÅ SEKVENSER. I EN GRAF ÄR ENA SEKVENSEN PÅ Y-AXELN OCH DEN ANDRA PÅ X-AXELN. MARKERING FINNS DÄR DET FINNS EN MATCHNING. OM LIKHET FINNS BILDAS EN LINJE.

##### ▪ DYNAMISK PROGRAMMERING

LIKHET GER SCORE. STRAFFPOÄNC FÖR OLIKHET.

VANLIGTVIS ANVÄNDER PROTEINALIGNMENTS EN SUBSTITUTIONSMATRIS FÖR ATT TILDELA AMINOSYRAN SCORES FÖR DESS MATCHES OCH MISMATCHES. KLYFT-STRAFF FÖR ATT MATCHA EN AMINOSYRA I EN SEKVENS MED EN KLYFTA EN ANNAN.

- **ORDMETOD**  
LETAR EFTER GANSKA EXAKT ÖVERENSSTÄMMELE I STOR DATABAS.
- **MULTIPEL SEKVENSALIGNMENT**  
UTÖKA PARVISA METODEN. ALLTSÅ SAMMA SAK MEN I MYCKET STÖRRE SKALA.  
CLUSTAL PROGRAM ELLER PSI-BLAST (POSITION SPECIFIC ITERATED BASIC LOCAL ALIGNMENT SEARCH TOOL).

## ANVÄNDNINGSSOMRÅDEN

- GENOMIK
- GENEXPRESSIONSANALYS
- PROTEOMIK
- BIOLOGISKA NÄTVERK

### I. GENOMIK

#### - **SHOTGUN – GENOME ASSEMBLY**

MAN HAR ETT ANTAL EXAKT LIKADANA GENOM, SÖNDERDELAR DEM I SLUMPVIS STORA DELAR. SEKVENSBESTÄMMER SMÅDELARNA. HITTAR BITAR SOM ÖVERLAPPAR → TALAR OM VAR DE ANDRA SKA LIGGA.

ÄR SVÅR METOD TY STORA DELAR ÄR IDENTISKA, UPPREPADE ALLTSÅ.

KALLAS FÖR "REPEATS" OCH KAN VARA TUSENTALS NUKLEOTIDER LÅNGA OCH FINNAS PÅ TUDENRALS OLIKA PLATSER. SÄRSKILT I STORA GENOM, SÅSOM I PLANTOR OCH MÄNNISKOR.

#### - **GENPREDIKTION**

MAN LISTAR UT VILKA DELAR I ETT DNA SOM KODAR FÖR VILKA GENER.

MAN SORTERAR UT DE PROTEINER SOM ENDAST FINNS KVAR AV HISTORISKA SKÄL. DVS FILTRERAR BORT ICKE-KODANDE BITAR. DETEKTERA FUNKTIONELLA DELAR, MÖNSTER OCH ORF (OPEN READING FRAME) MED HJÄLP AV

##### ○ **EMPIRISKA METODER**

LETAR EFTER SEKVENSER SOM LIKNAR REDAN KÄNDA. KAN ANVÄNDA LOKAL ALIGNMENT

MATCHNINGAR. DESSA KAN VARA TOTALA ELLER PARTIELLA. SÖKS UPP MHA ALGORITMER FÖR LOCAL ALIGNMENT. T.EX BLAST. UTFALL AV METODEN BEGRÄNSAS AV UPPLÖSNING OCH INNEHÅLLSMÄNGD I SEKVENS DATABASEN.

##### ○ **AB INITIO**

BASERAD PÅ GENINNEHÅLL OCH SIGNALDETEKTION. STARTDEKVENSS – mRNA

HEL SEKVENS – OPEN READING FRAME. SOM EJ HAR VERKLIG FUNKTION. LIGGER MELLAN VIKTIGA, KODANDE DELAR. ALLT FINNS I OCF. "GRÅ ZONER".

DET DÄREMELLAN – INTRON (FÖRSVINNEN, TAS BORT MED SPLICING)

DET AV INTRESSE – EXON.

SÄTTTS PÅ EN POLY-A-SVANS. MER KOMPLEXT I EUKARYOT ÄN PROKARYOT.

CpG-ÖAR – GENER SOM FINNS MEN EJ KODAR TY STÄNGTS AV. BLAND ANNAT PCA ATT DET FINNS SÅ MYCKET INTRONER I EUKARYOTER ÄR DET SVÅRARE ATT DETEKTERA PERIODICITET I EUKARYOTER.

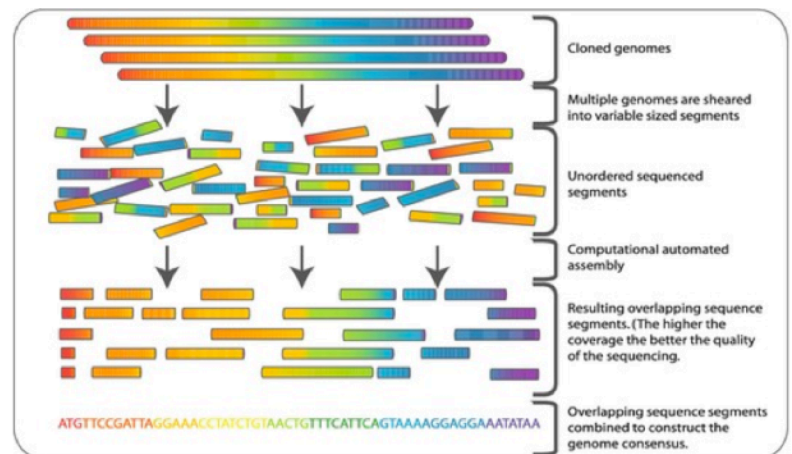
#### - **PHYLOGENETISKT TRÄD**

MOLEKYLÄR SEKVENSERING OCH MORFOLOGISKA DATAMATRISER.

PHYLOGENETISKA TRÄDET REPRESENTERAR EN HYPOTES GÄLLANDE DET EVOLUTIONÄRA FÖRHÅLLET MELLAN OLIKA BIOLOGISKA ARTER. SE HUR ORGANISMER ÄR BESLÄKTADE. TITTA PÅ SEKVENS.

NUKLEOTIDSEKVENSKODANDE GENER ÄR GRUNDEN FÖR KLASSIFIKATIONEN. MULTIPEL SEKVENSALIGNMENT SPELAR VIKTIG ROLL. METODER:

- MULTIPEL SEKVENS
- AVSTÅND MELLAN
- ANTAL MUTATIONER MELLAN.



### SANNOLIKHET FÖR EN MUTATION I FÖRHÅLLANDE TILL EN ANNAN.

- **DISTANSMATRISMETOD**  
GENETISKA AVSTÅNDET – ANDEL SOM INTE MATCHAR  
MELLANRUM IGNORERAS ELLER RÄKNAS SOM MISSMATCH.
- **MAXIMAL PARSIMONY METHOD**  
SNÅLASTE. TITTA PÅ VAD SOM MEST EKONOMISKT FRÅN EN SEKvens TILL EN ANNAN. HUR MÅNGA MUTATIONER SOM SKETT.  
IDENTIFIERAR DEN POTENTIELLA PHYLOGENETISKA TRÄD SOM KRÄVER LÄGST TOTALT NUMMER AV EVOLUTIONÄRA HÄNDELSER SOM FÖRKLARING TILL OBSERVERAD SEKvensDATA.  
ANVÄNDBAR DÅ ALLA EVENT EJ HAR SAMMA SANNOLIKHET.
- **MAX. LIKELYHOOD**  
HUR TROLIGT ATT EN SAK HÄNT I FÖRHÅLLANDE TILL EN ANNAN.  
FLER MUTATIONER = LÄGRE SANNOLIKHET

### **GENOMSPÄNNANDE ASSOCIATIONSSTUDIER , GWAS**

TITTA PÅ ANTAL INDIVIDER. TITTA PÅ ENKLASTE NUKLEOTIDERN I DNA SOM BYTTS UT OCH GIVIT SKILLNAD MELLAN INDIVIDERN. SE VAD SOM GER UPPHOV TILL VISST SÄRDRAG. BERÄKNAS STATISTISKT. OM DEN ENA ELLER ANDRA HAR SÅCELIEN BENÄGENHET ATT UTVECKLA SÄRSKILD SJUKDOM ELLER DYLIKT. FÅR UT T-VÄRDEN. SKA VARA SIGNIFIKANT SKILJT FRÅN  $\pm 1$ . ALLEL ÄR DÅ ASSOCIERAD MED SJUKDOM.

C-DNA = KOMPLEMENTÄRT DNA.

FOKUS LIGGER PÅ ATT FINNA ASSOCIATIONER MELLAN ENKEL POLYNUKLEOTID POLYMORFOS (SNPs) OCH DRAG SÅSOM STORA SJUKDOMAR. TITTA PÅ MÅNGA VANLIGA GENETISKA VARIANTER HOS OLIKA INDIVIDER FÖR ATT SE OM VARIATIONEN KAN ASSOCIERAS MED ETT VISST MÖNSTER. ÄR EJ EN SÖKANDE METOD TY DEN UNDERSÖKER HELA GENOMET.

ODDS RATIO BERÄKNAS. ODDS FÖR SJUKDOMEN OM MAN HAR SEN SPECIFIKA ALLELENS AKT ODDS FÖR SJUKDOM OM MAN INTE HAR ALLELEN. P-VÄRDE FÖR SIGNIFIKANSEN HOS ODDSET RÄKNAS SEDAN UT MED TJI-TVÅ-TEST.

SVAGHET I DESSA STUDIER ÄR ATT DE ENDAST FOKUSERAR PÅ VANLIGA GENETISKA VARIATIONER DÅ DESS ANTAGANDE ÄR ATT VANLIGA GENETISKA VARIATIONER FÖRKLARAR ÄRFTLIGHETEN I VANLIGA SJUKDOMAR.

### 2. GENEXPRESSIONSANALYS

MÄTNING AV EXPRESSION AV TUSENTALS GENER SIMULTANT. DETTA FÖR ATT SKAPA EN GLOBAL BILD AV DEN CELLULÄRA FUNKTIONEN. SEKVENSEN TALAR OM VAD CELLEN MÖJLIGEN KAN GÖRA. EXPRESSIONSPROFILER TALAR OM VAD DEN FAKTISKT GÖR VID VISS PUNKT I TIDEN.

#### TEKNIKER

##### ○ DNA MICRO ARRAY

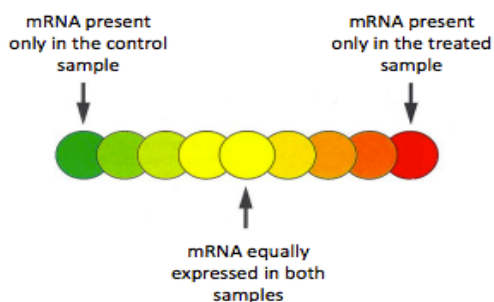
MÄTER DEN RELATIVA AKTIVITETEN AV TIDIGARE IDENTIFIERADE TARGET-GENER.

ETT DNA-CHIP ÄR EN SAMLING AV MIKROSKOPISKA DNA-SPOTS LÄNKADE TILL EN FAST YTA. VARJE DNA-SPOT INNEHÅLLER EN SPECIFIK DNA-SEKvens, KALLAD PROB. (KORT GENSEKTION) PROBEN ANVÄNDS FÖR ATT HYBRIDISERAS MED ETT cDNASAMPTEL UNDER STRIKTA FÖRHÅLLANDEN. PROB-TARGETHYBRIDISERINGEN DETEKTERAS OCH KVANTIFIERAS GENOM DETEKTION AV FLUOROFORER, SILVER OCH CHEMILUMESCENCE-MÄRKTA TARGETS. DETTA FÖR ATT AVGÖRA DEN RELATIVA BINDNINGEN AV NUKLEISYRASEKVENSER I TARGET.

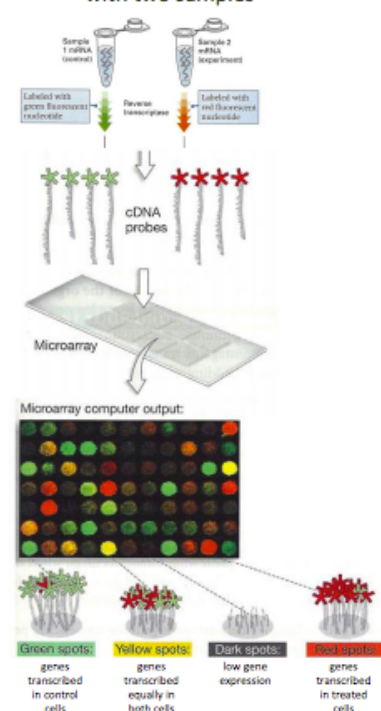
KAN ANVÄNDS FÖR ATT MÄTA FÖRÄNDRINGAR I EXPRESSIONSNIVÅ, DETEKTERA ENKLA NUKLEOTIDPOLYMORFOSER OCH FÖR GENOTYP.

##### DATAANALYS AV MIKROARRAY

ÄR VIKTIG. VARJE MIKROARRAY GER TIOUSENTALS DATAPUNKTER.



#### Comparing gene expression with two samples



##### ○ SERIAL ANALYSIS OF GENE EXPRESSION (SAGE)

SKAPAR SNAPSHOT AV mRNA POPULATIONEN I SAMPLET I FORM AV SMÅ TAGGAR SOM KORRESponderAR TILL FRAGMENT AV TRANskRIPTEN.

- RNA-SEQ (RNA-SEQUENCING)

ANVÄNDER FÖRMÅGAN AV NEXT-GENERATION SEQUENCING FÖR ATT VISA SNAPSHOT AV RNA-NÄRVARO OCH KVANTITET VID SAMMA TIDPUNKT.

### 3. PROTEOMICS

STRUKTUR OCH PREDIKTION

VARJE PROTEIN BESTÅR AV EN OVECKAD POLYPEPTID ELLER RANDOM COIL. DET VECKLAS SEDAN IN TILL EN KARAKTÄRISTISK OCH FUNKTIONELL TREDIMENSIONELL STRUKTUR.

KORREKT TREDIMENSIONELL STRUKTUR ÄR ESSENTIELL FÖR FUNKTIONEN. DOCK KAN VISSA DELAR HOS FUNKTIONELLA PROTEIN FÖRBLI OVECKADE.

SOMLIGA NEURODEGENERATIVA SJUKDOMAR BEROR PÅ EN ACKUMULERING AV MISSVECKNING. DVS ATT DE FÖRBLIR OVECKADE. SEKUNDÄR STRUKTURPREDIKTION. LOKALISERA SEKUNDÄRA STRUKTURER BASERAT PÅ KUNSKAP OM DESS AMINOSYRASEKVENSN.

TETRIÄR STRUKTURPREDIKTION

- *KOMPARATIV PROTEINMODELLERING.*

ANVÄNDER TIDIGARE LÖSTA STRUKTURER SOM STARTPUNKT.

HOMOLOGISK MODELLERING: BASERAT PÅ FÖRUTSATSEN ATT TVÅ HOMOLOGA PROTEINER DELAR LIKA STRUKTURER.

ANVÄNDS FÖR ATT AVGÖRA VILKEN DEL AV EN PROTEIN SOM ÄR VIKTIG I STRUKTURFORMATIONEN OCH INTERAGERAR

MED ANDRA PROTEIN. INFORMATIONEN ANVÄNDS FÖR ATT FÖRUTSE STRUKTUREN AV EN PROTEIN FRÅN DESS HOMOLOG.

FOLD RECOGNITION. METOD FÖR MODELLERING AV PROTEINER MED SAMMA VIKNING MED KÄNDA STRUKTURER, MEN SOM SAKNAR KÄNDA HOMOLOGER.

- *DE NOVO PHYSICS-BASED MODELING*

REFERERAR TILL ALGORITMPROCESS GENOM VILKEN PROTEINERS TETRIÄRA STRUKTUR KAN FÖRUTSPÅS FRÅN AMINOSYRASEKVENSEN. ( DEN PRIMÄRA STRUKTUREN )

### 4. BIOLOGISKA NÄTVERK

REPRESENTERAR OCH ANALYSERAR BIOLOGISKA KOMPLEXA SYSTEM.

NODER OCH FÖRBINDELSLINJER. STÅR FÖR INTERAKTION. LINJER – EDGES.

GRAD – ANTALLINJER SOM FÖRBINDER NODER.

BETWEENNESS – MÄTER HUR CENTRAL EN NOD ÄR I ETT NÄTVERK. NODER MED HÖG BETWEENNESS FUNGERAR SOM BROAR MELLAN OLIKA PUNKTER I NÄTVERKET.

INTERAKTION: FYSISKA INTERAKTIONER MELLAN MOLEKYLER I EN CELL.

TYPEN AV INTERAKTIONER:

- PROTEIN-PROTEIN NÄTVERK
- PROTEIN-DNA NÄTVERK
- METABOLISKA NÄTVERK –
- SIGNALNÄTVERK

PROTEOM: HELA SETTET AV PROTEINER UTTRYCKTA AV ETT GENOM, EN CELL, VÄVNAD ELLER EN ORGANISM VID ETT SÄRSKILT TILLFÄLLE. ”UTTRYCKTA PROTEIN AV EN GIVEN CELL ELLER ORGANISM VID EN GIVEN PUNKT I TIDEN UNDER DEFINIERADE FÖRHÅLLANDEN”