

Methoden vor und Nachteile

Rollen spiel in der Prüf.

1. Prüfung immer gut 2. Prüfung schwer
Viel wert auf Sprache gelegt [Sprache lernen]

Kritische Fragen
von der Pipelint
zu fragen

Wie arbeite ich mit Biologen zusammen
Wie werden Proben im Labor bearbeitet
Ansatz schon vor dem Experiment

Wir bewegen uns im Petabyte Bereich

Orga: Merit schreiben, nach 3 Tagen nochmal

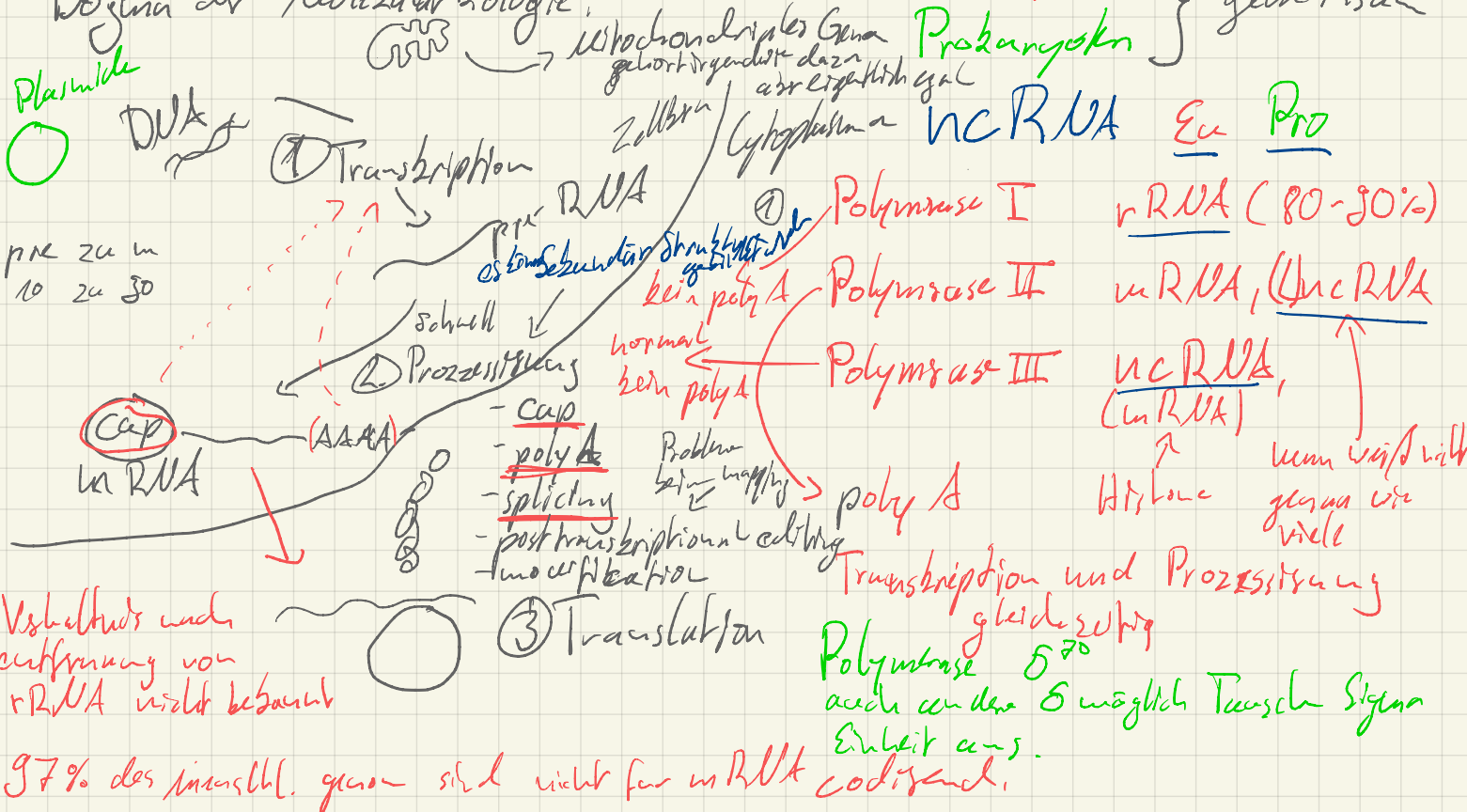
Praktikum: Wann wichtig und es aussuchen sollte, gilt es Lösung

Prüfung: Mündlich.

Genom interessant \rightarrow Genom aus Zelle \rightarrow DNA muss seq werden \rightarrow
DNA zerlegen \rightarrow wieder zusammen bauen \Rightarrow Assemblierung
Bei RNA große Fragmente \rightarrow Mapping

Datenverarbeitung

Dogmen der Molekularbiologie:



Der Großteil von nicht in RNA codierten Bereichen, codieren in RNA

Prokaryoten haben kein capping, aber Phosphat am 5' ende "durch/Seq"

Genom = gesamtele alle Gene + Rest.

total RNA oder TAB fressen alles was nur am 5' hat
 pre-RNA
 pre-RNA
 pre-RNA
 RNA

Transkriptom = Alles was RNA ist [Total RNA] inclusive Ribosomen.

freie Ribo kann aber entfernt gerade werden

was später gemappt wird

Vorbereitung zur Sequenzierung:

Es geht jetzt um Transkriptom

dRNA-Seq (differential) → unterschied zwischen zwei Proben/Koloren
 Bsp. Infiziert vs. nicht infiziert / durch zeitlichen Verlauf
 • total RNA → seq → 90-95% vRNA (schlecht)

• poly A-Seq (nur mRNA aus Pol II, günstig, effektiv)

• Trennung über Sonden (binden an 2' oder RNA) [Ribo zero → 5%-60% vRNA]
 • small RNA → alles kleine microRNA, piRNA selbst herstellen 3' nucleotid für Key-Organismen
 schlechte RNA < 50%

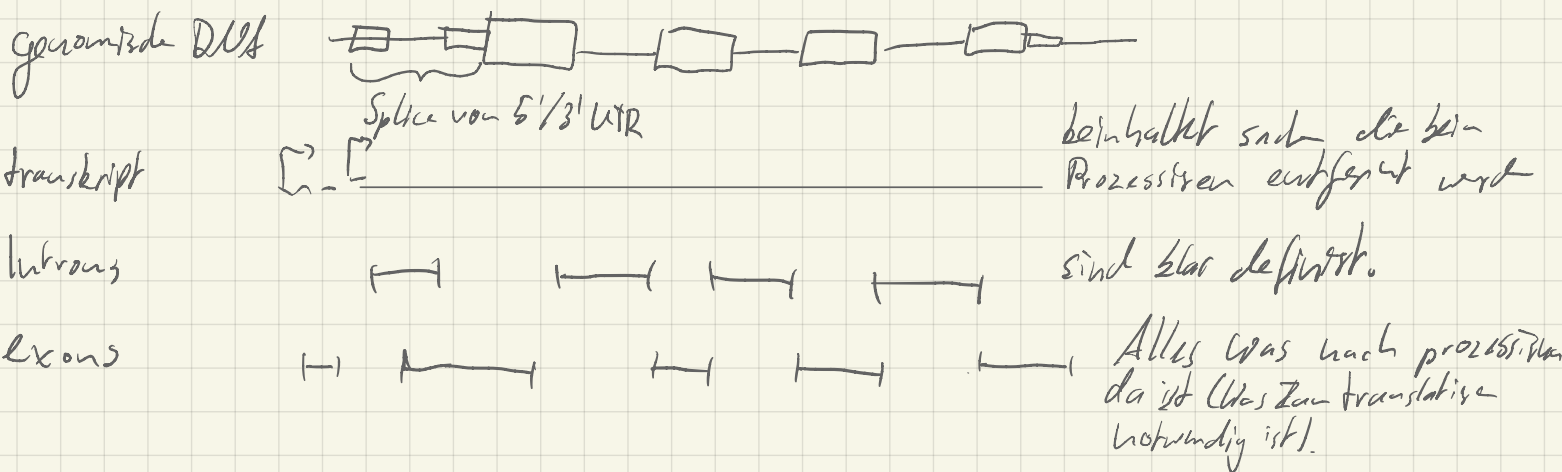
[ChIP-Seq: alles was an target Protein bindet]
 → wenn gar kein Plan.

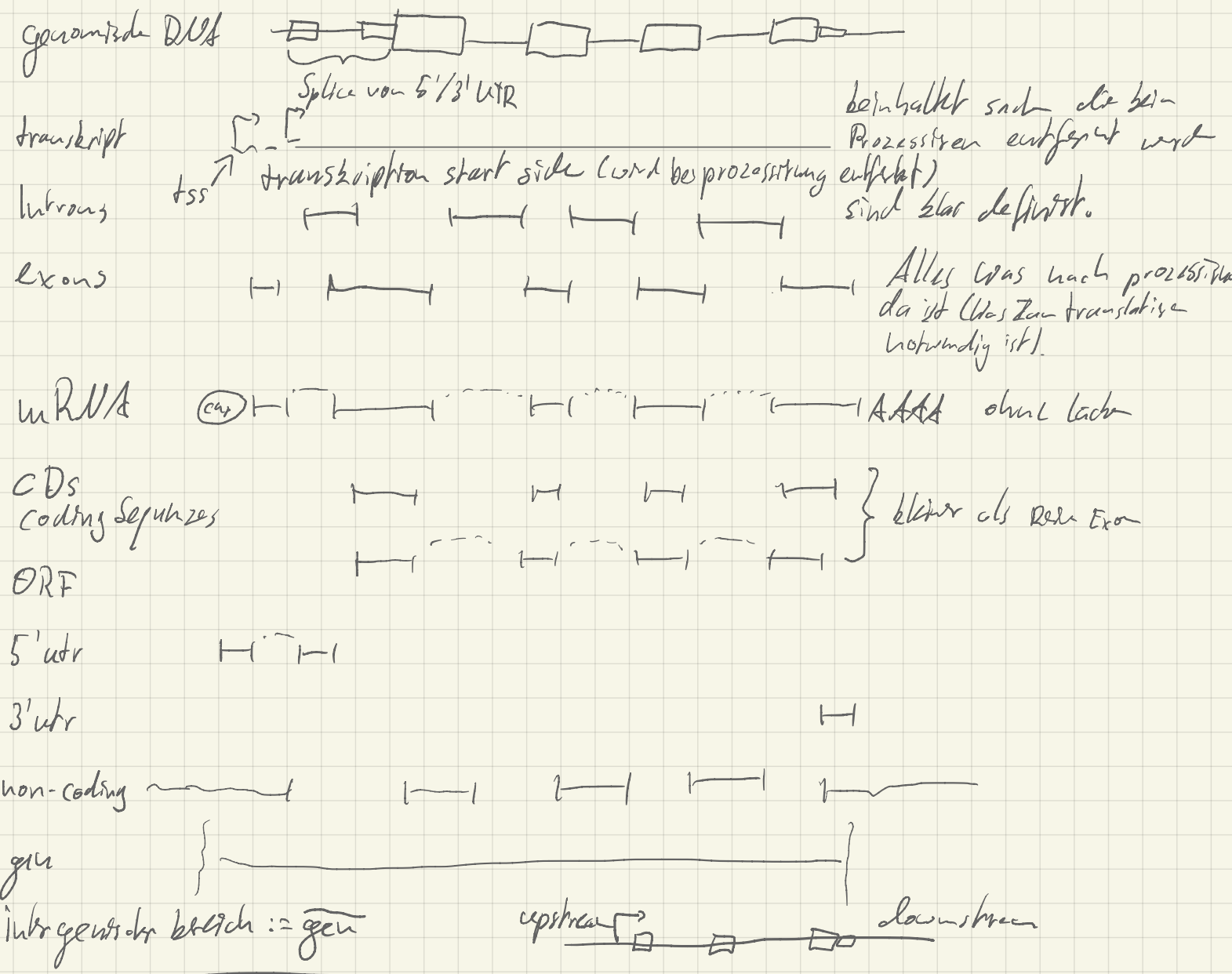
Verstärken warum aber negativ beeinflussen

total RNA, warum? kostengünstiger gemacht als Ribo zero zuerst.
 (Bei Backstein, Archäen, kleine Einzeller)

Man braucht fast immer 1µg pro Methode.
 Wir sollen Protokolle kennen um zu empfehlen

Gen := [führende Gene/Exon sind nicht sicher] alles damit es funktioniert





Ähnlichkeiten:

Mensch Mensch 99,9% } eigentlich alles

Mensch Schimpanse 98,4%

Mensch Maus 90%

Mensch Fliege 60%

Mensch Banane 50%

Mensch Hefe 30%

→ Nur auf beide (Prä) die beide haben

Output bei Sequenzierung: reads $\xrightarrow{\text{zusammenbau}}$ contigs \rightarrow scaffold \rightarrow chromosom

Alles sehr durcheinander gibt auch zu denken.

keine genauen values

→ kein Genom ist fertig Assembled?

Es gibt etwa 100 (nicht Bakterien) Genome im Chromosom Bereich. Bsp.: Bos Taurus (Kuh/Rind)

Die selben reads in verschiedenen Tools ergeben verschiedene Assemblies mit etwa 60% Id (Problem sind repeatregionen)

Genome die man downloaded sind oft nur durch einen Assembler gelaufen (kritisch sehen)

Homologe Gene: = Genen des selben Ursprungs ^{hpa} ^{musculus}

Paraloge Gene: = Kopien von Genen im gleichen Genom, müssen nicht gleiche Funktion haben, ^{homo sapiens} ^{musculus} ^{abgleicher Ursprung}

Orthologe Gene: = haben den selben Vorfahren, in 2 Spezies, macht das selbe

Pseudogen: ^{paraloge Gene die aber nicht abgelesen werden.}
In Datenbanken sind manche falsch annotiert, früher so gedacht und dann das Gegenteil entdeckt und nicht verbessert

Promotor: Region des Gens meist up- (oder down)stream.

Terminator: Bei Bakterien 2 Terminatoren die abhängig oder nicht.
Bei Eukaryoten nicht verstanden. ^{happier}

Theorie: Polymerase liest einfach weiter bis irgendwas abfällt.

Splicing:

Alternatives Splicing

Intron retention: (Intron zurückhaltung)

Enhancer

Silencer

Est: Expressed Sequenz Tag

Est: Einzelne gefundene Exprimierte Sequenzen

~ Viren: 1kilo Basenpaar

Human Genom Projekt: 3,2 und BP

~ Bakterien-genom: 3 Mega Basenpaar

Unterscheidung zwischen De novo und Referenzsequenzierung

Jedes Genom Sequenzierung

Annotation

Gedächtnis Speicherung [Pasha]

Tools entwickeln, die mit Datensatz umgehen können

Ethik → will man wissen, warum man stark Pancreasome keine sinnvolle Datenstruktur

Genregulation verstehen

Funktion jedes Gens verstehen (prot. codieren 50% stuff faster)

Gen Interaktion auf DNA und nach Expression

Chromosomales Interaktion

Konzentration der Gene im Genom

SUPs verstehen (Cicchoh)

gibt nur 1 ref Genom

mittels mehr oder weniger in der Datenbank

und off

Chr. 5: stark stop...

genome feature file
strong gebunden an
faster

Beim Sequenzieren ist jede 1000000te Base falsch

Fehlerrate 10^{-5}

Mensch hat etwa 22000 Gene (Protein-codierende Gene)

(gibt aber mehr, da nicht-codierende auch dabei)

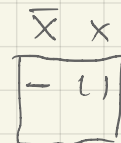
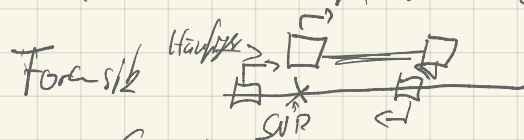
Korrelation zwischen Komplexität der Organismen und Anzahl nicht-codierender Regionen es gibt.

20000 behandeln nicht Isoforme (alternative Formen des Gens)

Von 20000 Genen konnte man 1500 zu Transkripten zuordnen.

Pro Zelltyp hat etwa 6000 Gene aktiv. Viele Gene sind für

extreme cases



SUPs die nur in wenigen Vorkommen

Nutzen: Krankheitsdiagnose, Gentherapie, Frühdiagnose, Medikamente an DNA anpassen, Forensik [Konzentration auf SUPs], Verdauung, Schizophren, Alkohol,

Genom: Protein-codierend 1,2,7% (Exons), Protein-codierende Introns: 25,3%, Lines 20,4%, Sines 13,1% (stark in Viren), Retrotransposons (zu Viren zuordnen) 8,3%, Zugs

Man weiß so gut wie nichts über nicht-proteine

Von prokine Transkriptionsfaktor 12

Transkription 8,8

Nukleinsäure-bindende
proteine 8,5

Rezeptor Promotor 6,5

Transportprot. 6,4

keine Ahnung 23,6%

Sequenzierung

• Maxam-Gilbert

• Sanger / Kette-Abbau

→ vorbereiten.

ncbi full genome Download googlen

ftp.ncbi.nih.gov/genomes/

fua fache nucleotide

faca fache amino acid

Genomnummer m39

Annotationsnummer

normale Weise im Readline

→ ggf. file lesen können

UCSC genome browser

hgdownload.soe.ucsc.edu/download.html

masked → repeats durch definition ersetzen (hard mask)

→ uppercase durch lowercase ersetzen (soft mask)

verschiedene tools oben genannt nutzen RepeatMasker

Centromer und telomere region meist simple repeats

Sines lines prozentual am meisten

k-mer based analysis

animal genome size database

genome size.com

C value gewicht der V nucleotide in pikogrammen

Sequenzmethoden

Illumina: [Illumina Sequencing by Synthesis] Fehler 10^{-3}

Firma Methode Device

Illumina HighSeq / MiSeq / NextSeq (Unterbrechen in Synapsis/Lanes)
PyroSeq

Kommentare Länge 200-300 Nucleotide four five four

Nicht gleich komplettes Gen, 28-stückelte Teile Rada

Adaptonen notwendig ist wichtig, letzter Teil wirklich relevant
Clustering erspart \Rightarrow Flowcell

Oligos auf Flowcells sind fest und exakt komplementär zu
letzten Teil an Adaptor

Nur 5' enden auf Flowcell

Da Licht nicht stark genug \rightarrow Spots (Cluster von gleichem seq)

Oligos in Spots nah genug damit Brücke bauen
Weiter (Waschschritt um nur 1 Bindung zu steuern)
Sequenzierung durch Synthese

Cycle = 1 Waschrunde + Enzymatisches

Base calling = Lichtsignal Übersetzung

2x 150 hilft beim Assemblieren \rightarrow Paired end Sequenzieren

Barcodes in Adaptonen hilft beim aufteilen

PacBio Sequencing Fehler 10^{-4}

Pacific bioscience SART Sequencing

Single Molecule Real Time Sequencing

relativ viele Fehler \rightarrow Circularisierung und 6 mal lesen.

Multiplexing \rightarrow Mehrere Sequenzen gleichzeitig analog Barcoding
Grundlage für Single cell Sequencing

Oxford minion Sequencing

Oxford nanopore technologies Device MinION

5 Gleichzeitig in Pore (200 Pore pro chip) PromethION (8x MinION)

5-Mer chemistry 10x.

Hohe Fehlerrate weil OGCs schlecht, oder zu viele modificationen.

Viele Modifikationen lassen ist aktuelle Forschung

Vorteil hier: RDB und Prokary direkt Sequenzierung, real time, keine Assemblierungsmethoden vorhanden direkt im Feld, günstig

Platform	Illumina MiSeq	Ion Torrent PGM	PacBio RS	Illumina GAIIx	Illumina HiSeq 2000	Oxford Nanopore
Instrument Cost*	\$128 K	\$80 K**	\$695 K	\$256 K	\$654 K	1000\$
Sequence yield per run	1.5-2Gb	20-50 Mb on 314 chip, 100-200 Mb on 316 chip, 1Gb on 318 chip	100 Mb	30Gb	600Gb	22GB DNA 15GB RNA
Sequencing cost per Gb*	\$502	\$1000 (318 chip)	\$2000	\$148	\$41	75\$ (50x)
Run Time	27 hours***	2 hours	2 hours	10 days	11 days	260
Reported Accuracy	Mostly > 99.9%	Mostly Q20	Q20	Mostly > Q30	Mostly > Q30	live (Pore read) ~Q20 DNA 2% Q10 RNA 15%
Observed Raw Error Rate	0.80 %	1.71 %	12.86 %	0.76 %	0.26 %	2% beschränkung Pore und längst 50kb
Read length	up to 150 bases each 2x500	~200 bases	Average 1500 bases**** (C1 chemistry) mit 100kb	up to 150 bases	up to 150 bases each 2x500	=
Paired reads	Yes	Yes	No	Yes	Yes	=
Insert size	up to 700 bases	up to 250 bases - 350	up to 10 kb	up to 700 bases	up to 700 bases	=
Typical DNA requirements	50-1000 ng	100-1000 ng	~1 µg	50-1000 ng	50-1000 ng	100 ng bis zu 1 µg +

* All cost calculations are based on list price quotations obtained from the manufacturer and assume expected sequence yield stated.
 ** System price including PGM, server, OneTouch and OneTouch ES.
 *** Includes two hours of cluster generation.
 **** Mean mapped read length includes adapter and reverse strand sequences. Subread lengths, i.e. the individual stretches of sequence originating from the sequenced fragment, are significantly shorter.

1 Mb Zellen => 1 µg

Coverage: reads auf Genom und schauen wie viel ist abgedeckt (Wiki und Leses)
 60 Faden Illumina ist schon
 Bei Ref nicht schon 10-20 Faden
 Bei langen reads braucht man nicht so hohe Coverage

und NovaSeq (table vervollständigen) iontorrent nachschauen
 1 Sequenzierungsmethode in Prüfung stehen
 Warum Fehler

Kosten berechnung: Für Mensch. Genom

3200 Mb Transkriptom analyse
 2,34% 75 Mb
 6000 Mb 80x (Coverage)

Kosten berechnung: Für Mensch. genau

3200 Mb $\xrightarrow{\text{Transkriptom analyse}}$ 75 Mb $\left. \begin{array}{l} 2,34\% \\ 80 \times (\text{Coverage}) \end{array} \right\}$ unter 30x schlecht!

Abfrage: (Plz 80%)

6000 Mb $\left. \begin{array}{l} 2 \cdot 100 \text{ bp} \end{array} \right\}$

Genomgröße

30 Mio reads

1/4 viel davon

\Rightarrow OUT für 3 TK

Coverage

15 Gb \sim 500 € Flowcell 1x
250 € Libprep 3x

Read output

Illumina

1x 3000 € Flowcell (36 Gb)
6x \sim 200 € Libprep

Fehlt Ignorieren Rechnung in Prüfung

Assemblierungsverfahren:

Wie gut ist das Tool? Welche Probleme hat dieses Tool (immer 10 werden kann das Tool lösen)? Laufzeiten der Methoden (auch Alignments) kennen

GREEDY-Algorithmus (bei 5 Reads 10 Alignments)

- 1.) pairwise Alignments
- 2.) max Score nehmen AIB
- 3.) Merge AIB \rightarrow AB
- 4.) repeat

Kommt klar mit ①, ④ klar

Probleme der Assemblierung

- | | |
|-----------------------------|------------|
| ① Sequenzierfehler | ⑥ Zeit |
| ② Repeats | ⑦ Speicher |
| ③ Orientierung | |
| ④ Lücken (Bug oder Feature) | |

⑤ Chimeras (Zusammenchluss von Sequenzen die nicht miteinander zusammenhängen Barcode).

Algorithmen kommen mit Menge an Reads und Länge nicht klar

Shortest Superstring Problem mit DWA

Gegeben: Set von Strings von 000

Gesucht: kürzester Superstring

Anschauen: Euler Pfad, Hamilton Pfad

111

Assemblierung:

~ das Zusammenführen von Nukleotidsequenzen zu längeren Sequenzfragmenten (durch Sequenz alignment)

Reads \rightarrow contigs \rightarrow (supercontig) \rightarrow scaffold \rightarrow (superscaffold)
 \rightarrow chromosome

Greedy

1. Pairweise Alignment

2. Suche max Score von Alignment ^{diese}

3. Merges max Score-Alignment

4. Wiederhole

Nachteil: feste Zusammenschluss.

Mehrere Contigs, aber wissen über

Nachteile:

- kann mit Repeats nicht umgehen

- sehr hohe Laufzeit, viel

Speicher

- Orientierung

Gute Kriterien für Genom-Assemblierung:

N50-Wert ~ Länge des kürzesten Contigs, das in Summe mit allen längeren Contigs die Hälfte des Genomassembly abdeckt.

\rightarrow Contigs der Länge nach sortieren, dann von lang nach kurz die Contiglänge aufsummieren bis 50% der Basen des Gesamtassembly erreicht sind \rightarrow dann die Länge des letzten addierten Contigs

L50-Wert ~ Wie N50 aber anstatt der Länge des zuletzt addierten Contigs betrachte man Anzahl der Summierte Contigs.

Bsp.: 10, 7, 8, 7, 6, 5, 4, 3, 2 $\Sigma 54$
 $10 + 7 + 8 = 27 \Rightarrow N50 = 8 \quad L50 = 3$

Für Bekannte sind die Werte gut (Dua bin ein Genom (Chromosome))
Viel sagt nichts über Richtigkeit aus.

NGS ~ wie NGS, aber nicht im Bezug auf die im Assembly vorhandene
Nucleotide, sondern bzgl. Genomgröße.

C-Value: $\lambda_{pg} \approx 978 \text{ Mb}$

$$\text{Coverage} = \frac{u \cdot L}{G}$$

$u = \# \text{ Reads}$

$L = \text{durchschnittliche Readlänge}$

$G = \text{Genomlänge}$

so wenn mind. 80

Strategien: Greedy, Hamilton, Euler, Overlapping Consensus (Graph, De Bruijn)
 \hookrightarrow Velvet, Spades, Soap denovo
 \hookrightarrow Newbler \hookrightarrow Velvet, Spades, Soap denovo
 \hookrightarrow SSAsse \hookrightarrow Newbler \hookrightarrow Velvet, Trinity, (Spades), (Soap denovo)

Hamilton:

Shortest Superstring-Problem: ~ SS von zwei Strings X und Y ist ein kurzester String
Z von dem sowohl X als auch Y Substrings sind.

\rightarrow einfach lösbar auf n strings.

- Das SSP ist NP-schwer (NP-vollständig)
- Es gibt nicht immer eine eindeutige Lösung für das SSP
- kann in das Traveling-Sales-Problem umgewandelt werden (NP-vollständig). \hookrightarrow sehr effiziente Heuristiken.

\rightarrow Reduzierung von SSP zum TSP-Hamiltonpfad.

Hamilton-Pfad: ~ ist ein geschlossener Kreis in einem Graphen, der jeden Knoten genau
einmal enthält

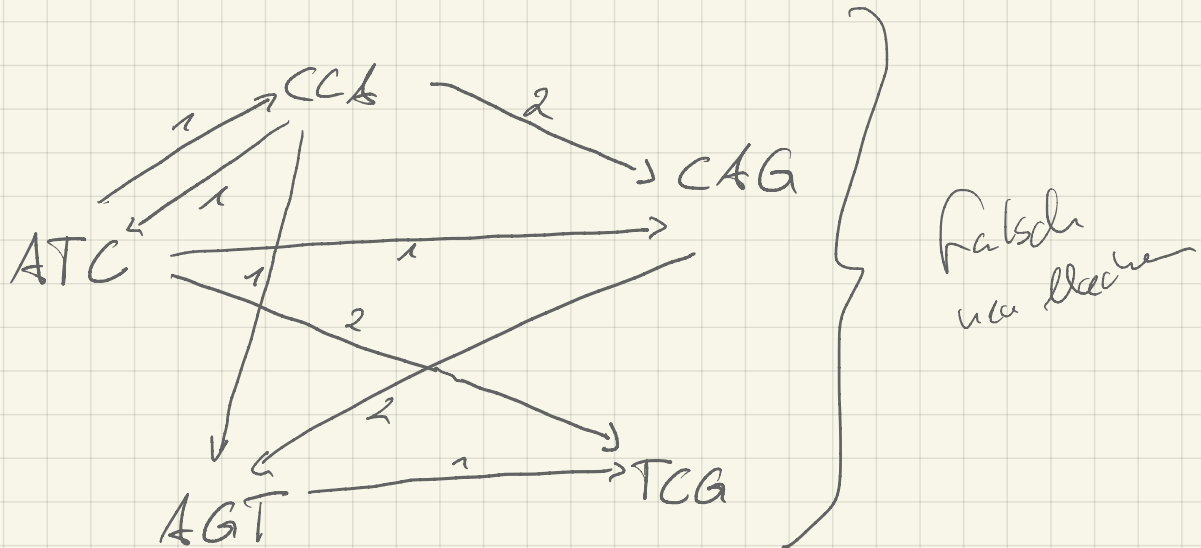
Wir definieren einen Overlap (S_i, S_j) als die Länge des längsten Präfixes von S_j der
auch Suffix von S_i ist. \Rightarrow Reihenfolge ist entscheidend.

Algorithmus - Idee:

- jeder Read ist ein Knoten in unserem Graph
- jedes Kantenpaar zwischen zwei beliebigen Knoten ist Overlap der beiden Knoten
→ Graph ist gerichteter!

Große Probleme: ① Repeats
② Lücke

$S = \{ATC, CCA, CAG, TCG, AGT\}$



1. $\begin{matrix} 2 \\ ATC \\ TCG \end{matrix}$ $\begin{matrix} 2 \\ CCA \\ CAG \end{matrix}$ AGT \Rightarrow ges 4

ATC
CCA 1
CAG 2
TCG 0
TCG 1 \Rightarrow 4

Euler-Algorithmus

Euler-Kreisproblem: Gegeben ein Graph G , finde einen Kreis, der jede Kante von G genau einmal enthält.

Euler-Pfad: Wie Kreis nur nicht ein selbster Endknoten.

Algorithmusidee

① Zerlege alle Reads in ein Fragmentspektrum F der Länge k

② Erstelle einen Graphen G_{DB} (gerichtet)

$$G_{DB}(S) = (V, E, Le)$$

$V =$ alle Prefix / Suffix der Länge $k-1$ aus F

$$E = \{ (a, b) \mid \exists f_i \text{ in } F : \langle a, b \rangle = f_i \} \quad a, b \in V$$

Cherke Knoten sind Teil eines Kreises (Fragments)

$$Le = E \rightarrow \Sigma \quad \text{Labeling}(a, b) = b_n$$

③ Finde Eulerpfad

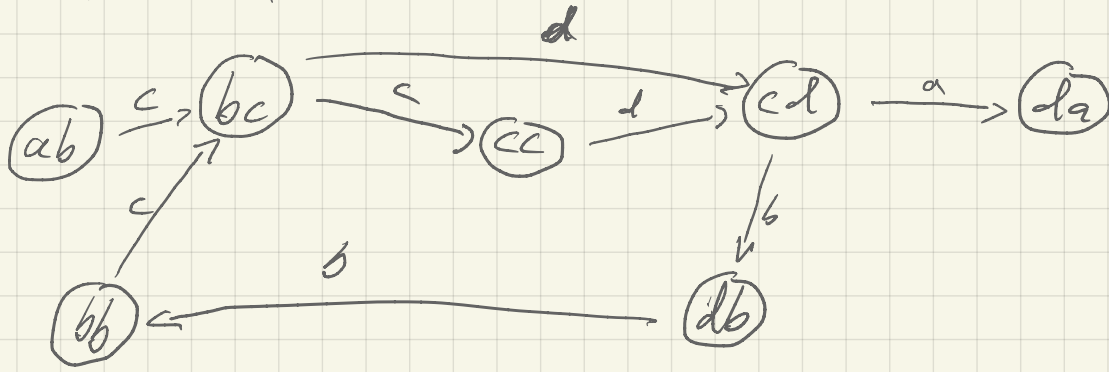
Vorteile

- Keine Berechnung paarweise Alignments oder explizit Overlap zwischen zwei Reads
- Es gibt schon sehr effiziente Algorithmen um Eulerpfade zu finden ($O(n)$ Alg. von Hierholzer) $\hookrightarrow \# \text{ Reads}$
- Ein Eulergraph lässt sich leicht in Knoten disjunkte Kreise zerlegen.

Nachteile:

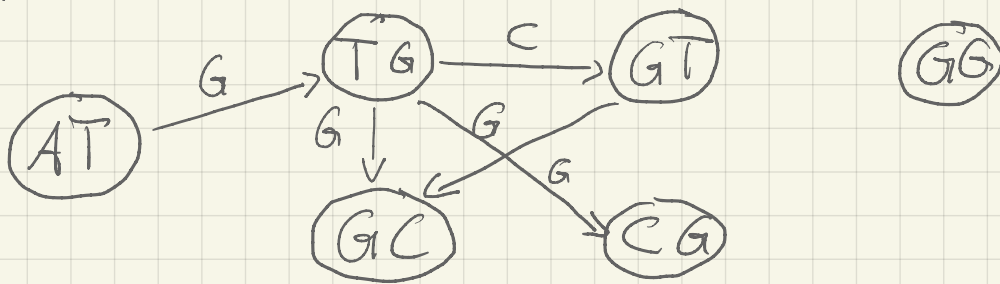
- (erstmal) nicht robust gegen Sequenzfehler
- $\# \text{ Eulerpfade} \sim e^{\# \text{ Repeats}}$
- Zerlegung unserer Reads in k -mer \rightarrow Informationsverlust?

$S = \{abc, bba, bcc, bcd, abd, cda, cdb, dbb\}$



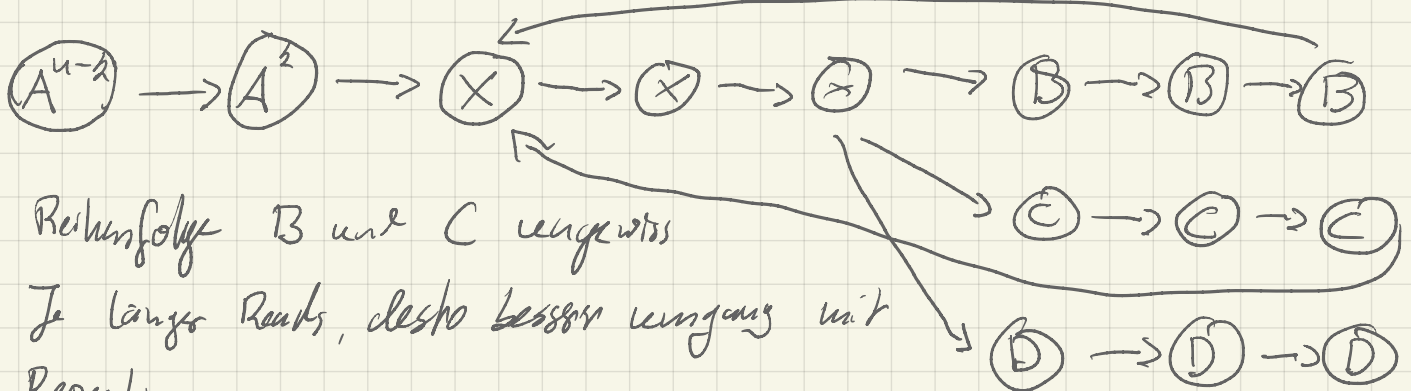
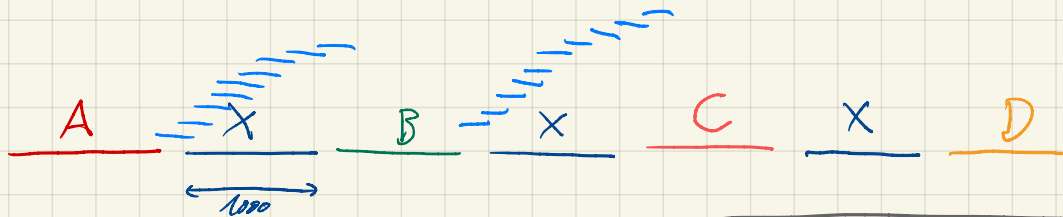
$S = \{ATG, TGC, GTG, GGC, GCA, GCG, CGT, TGG\}$

$k=2$



Paired end
für infos über
Lücken

Euler und Hamilton lernen
↳ hauptsächlich für Illumina



Reihenfolge B und C ungewiss

Je länger Reads, desto besser Umgang mit
Repeats

Euler/Hamilton kann nicht mit Repeats umgehen

Menschliches Genom und ähnliche Eukaryoten hatten keine

rRNAs (doi:10.1101/2013), da Assemblierer nicht mit Repeats klar kommen.

Heute meist 1, aber nicht alle. Assemblierer kommt auch nicht mit Telo-, Leuchtmass
oder

Bakterien repeats:
16 Operon
sonst keine

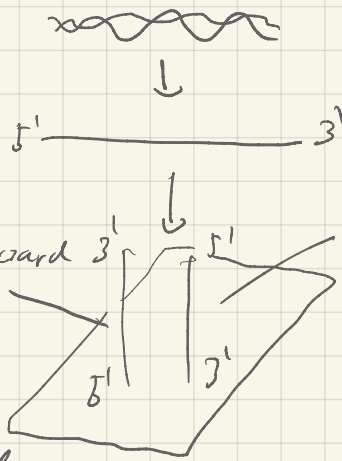
Velvet

Paper:

Velvet Algorithms for

denovo

Jerbino, Birney, 2008



anweisung / reverse
[Kategorie] ← nur für Transkripte

Velvet abgeleitet durch Sanders

Besteht aus K-mer und Euler

Neu Knoten haben Zwillinge Knoten

für 25-50 Bp reads

k-länge wählen, kleine k's klein graph
große k-länge zu viel Rechtleistung
Hardwares. nicht unter 13

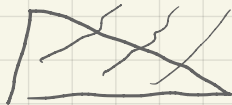
TAGACTG → TAGAC
AGACT
GACTG
können gemischt werden da keine funktionen

TAGAC
AGACT
GACTG

CTG
CTA

ATTG
CAGT

Beispiel aus
Paper



Zwillingknoten nicht rekonstruierbar

Aber k-mer rekonstruierbar

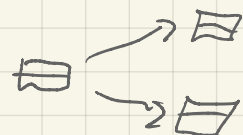
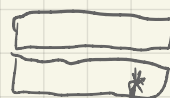
Verbindungen sind auch bei Zwillings äquivalent.

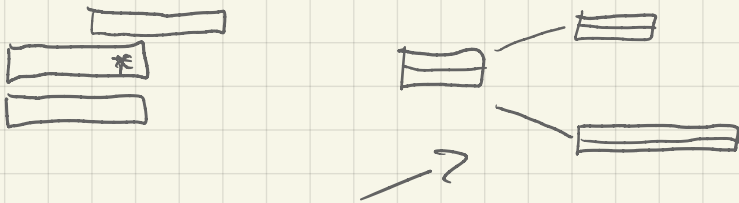
CC

Graphen Aufbau eigentlich unabhängig von Tool geleitet.

Tools unterscheiden sich in Konfliktlösung.

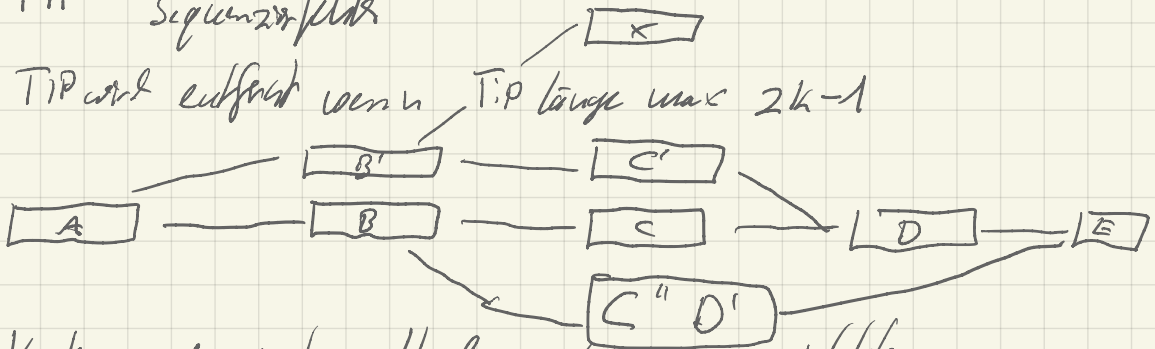
Sequenzfehler





② TIP Sequenzfeld

TIP wird erfüllt wenn TIP Länge max $2k-1$



Kanten werden mit Anzahl der vorkommenden geschildert
man geht meist die stärksten Kanten

⑥ Tourbus (mehr oder weniger Dijkstra)

⇒ Mehr reads liefert auch mehr Fehler, Dijkstra schonint problem
→ Heuristiken → ungenau...

Mapping Verfahren Map, soap, cufflinks

In Review vor einiger Zeit hat zwei tools ergeben.

Bowtie → Bowtie2 → Tophat → Hisat ← State of the Art

SEGEMEHL (Steve Hoffmann, Jena)

Bowtie:

Steven Salzberg

Paper: Ultrafast and memory efficient alignment of short DNA sequences to the human genome, 2009

Entwickelt für Illumina Sequenzierung

2 Arten von Mapping

1. Genom in Speicher → Reads über genom [Soap, Bowtie, Segemehl]

2. Reads in Speicher → Genom über reads [BLAT, Zorro]

Burrow-Index(?) zum schneller suchen in Genom.
Whicker

gegeben: Text T, terminalsymbol \$

gesucht: BWTCTD

1. \$ a c a a c g

2. 0 \$ a c a a c g

1 a a c g \$ a c

a c a a c g \$

3 a c g \$ a c a

4 c a a c g \$ a

5 c g \$ a c a a

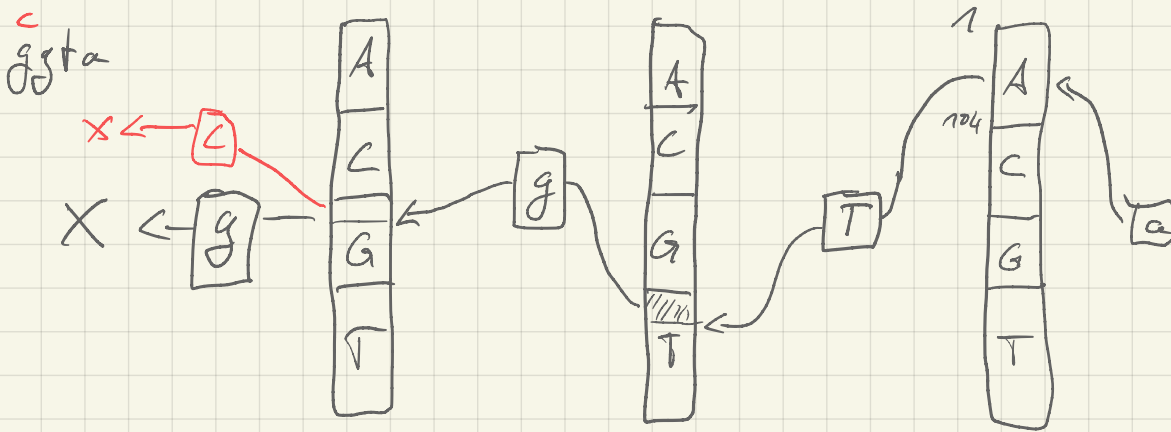
6 g \$ a c a a c

← BWT

3. $\overleftarrow{a a c}$

suchen der Reads

Wenn read sequenzfehler ist schon Problematisches. Wie fixieren?



* fastq

>id

ACTCGATG

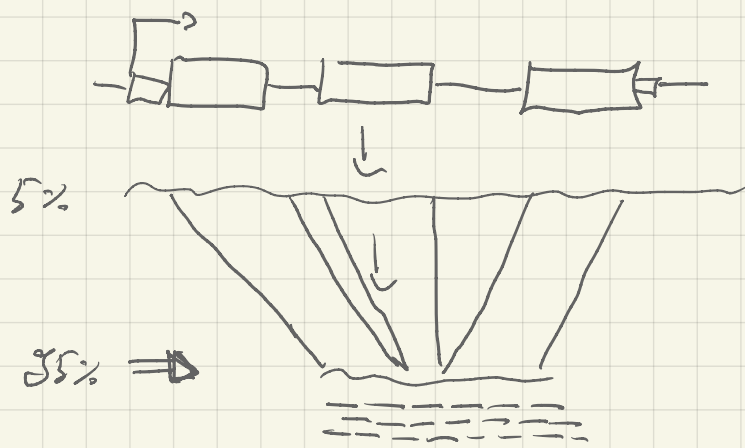
Warten meiste Fehler

Beweis 2 durch reads um

Wenn Beweise eher nicht mehr zugeordnet kann
wird das read verworfen

Leserfehler erhöht sich stark mit Erhöhung der Fehler

"Sequenzierung"?



Baatz in original (kann nicht
mit splitreads umgehen,
→ Baatz2 schon, = splitreads

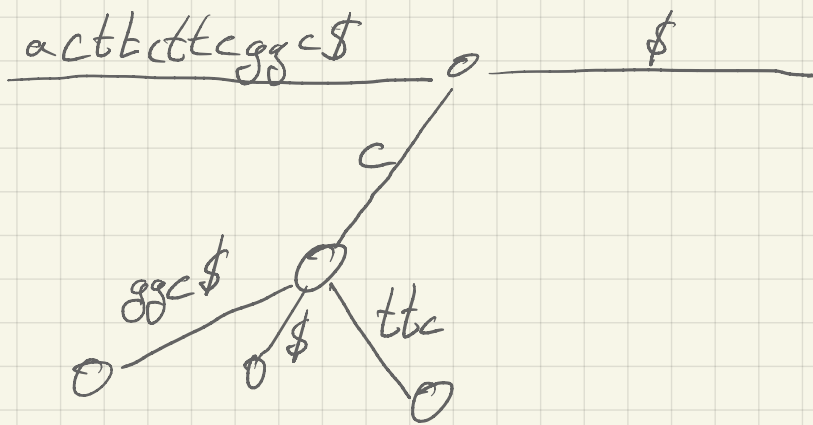
SEGE/MEHL:

- splitreads
- Insertion / Deletions
- ≥ 3 Mutationen

Fast Mapping of Short Sequences
with Mismatches, Insertions and
Deletions Using Index Structures
~~2009~~

genom: 5 := a c t t c t t c g g c
Suffix-Baum terminalsymbol \$

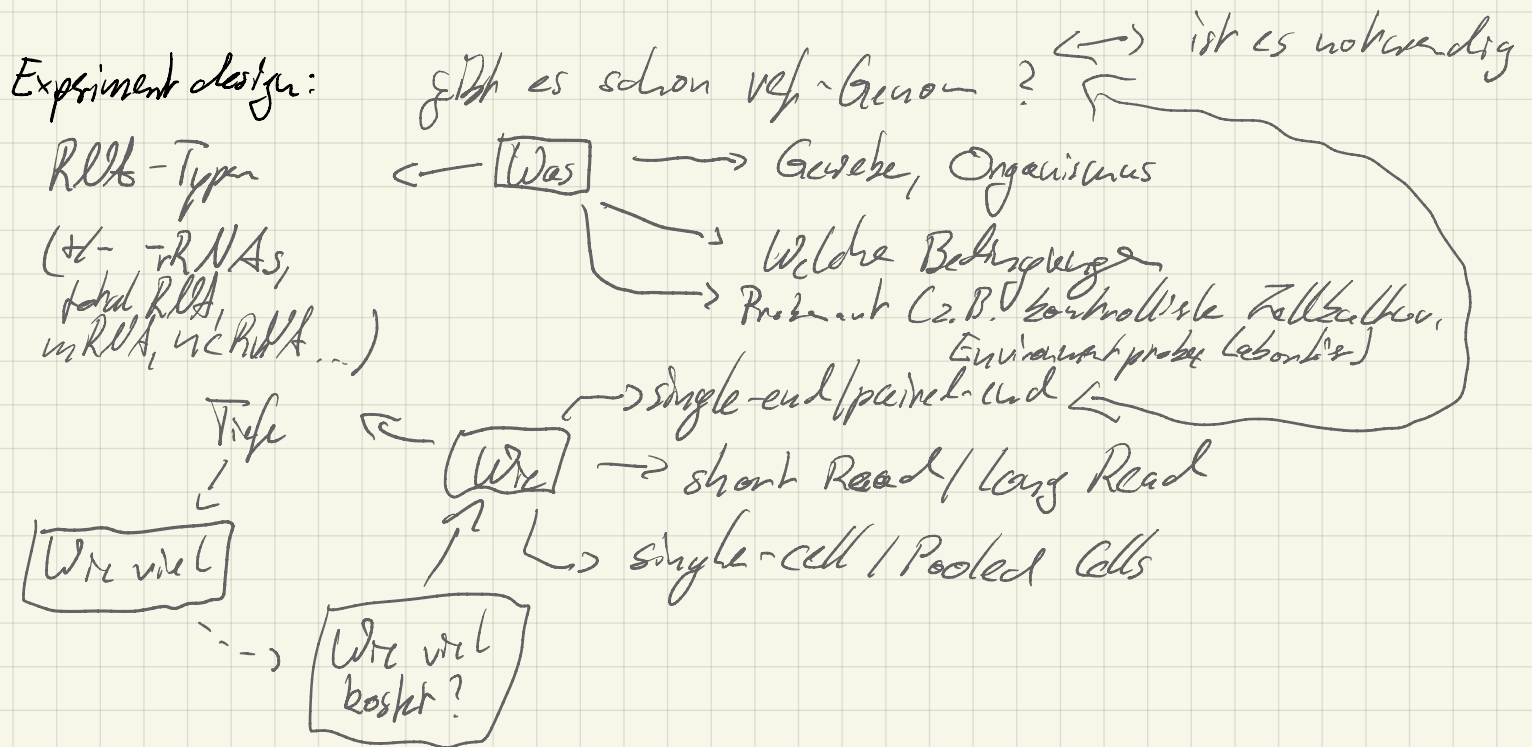
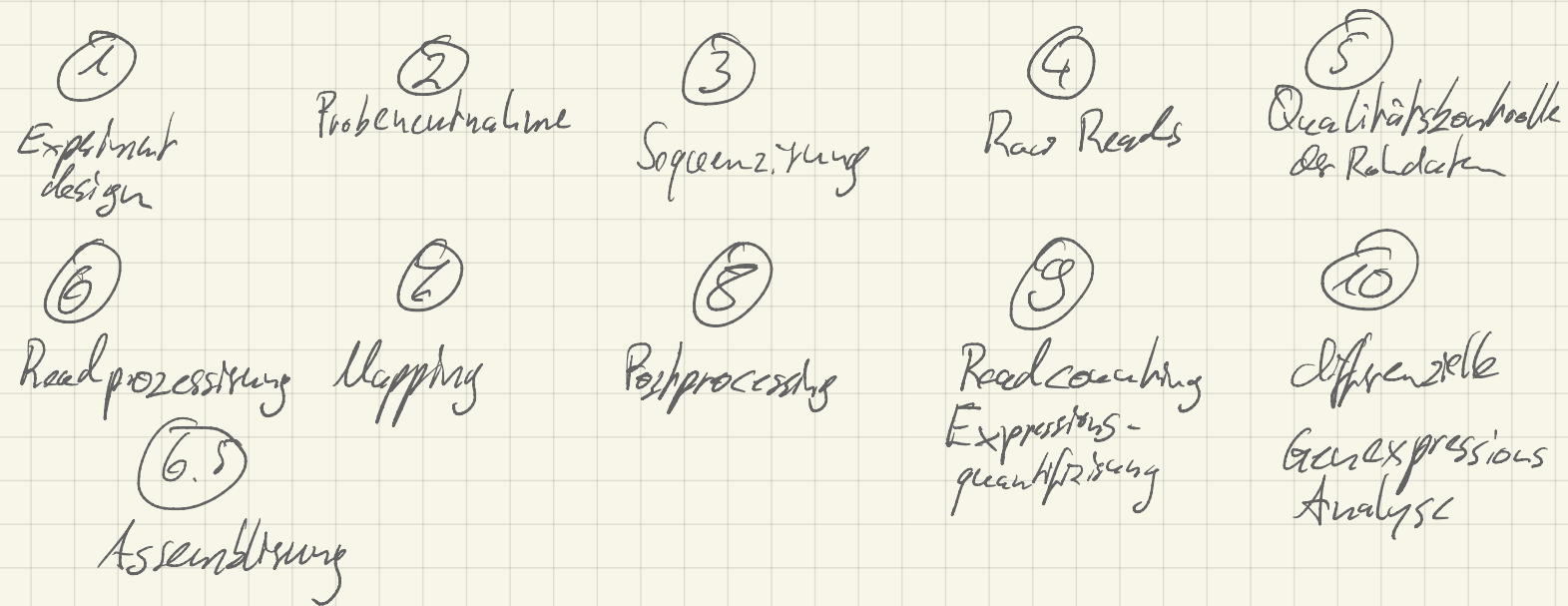
Felder in Figure caption finden,
für Prüfung wichtig



hier ist speziell

Baatz sehr picky vs Segemehl will unappen

Klassische Geneexpressionspipeline



Probencentralsnahme + Library Prep:

1. Probencentralsnahme → Lysieren → RNA isolieren
2. RNA extr. filtern
3. RNA in cDNA umschreiben
4. cDNA spalten (Enzym oder Ultraschall)
5. Adapter-Ligieren

Sequenzierung: siehe vorherige stand

Rohdaten:

Erster Check $\# \text{ Reads} \cdot \varnothing \text{ Readlänge} \approx \text{Cov. / Tiefe}$

Fastq-Format:

- Reiner ASCII-Text

verglichen mit dem
theoretischen Deckungsgrad der
Flowcell

- 1 Read \Rightarrow 4 Zeilen

- 1.: Header beginnt mit @

- 2.: Sequenz des Reads im 5 Buchstaben code

- 3.: Space / Trennsymbol beginnt mit +

- 4.: Sequenzqualität \rightarrow Phred-Score

Phred-Score, Q ist ein Wert, der sich logarithmisch zur Sequenz-
fehlerwahrscheinlichkeit P einer einzelnen Base verhält.

$$Q = -10 \cdot \log_{10}(P)$$
$$P = 10^{-\frac{Q}{10}}$$

 \rightarrow Prüfung! und was macht Phred

$Q = 10 \rightarrow P = 0.1 \approx \text{Accuracy } 90\%$

$Q = 20 \rightarrow P = 0.01 \approx \text{Accuracy } 99\%$

$Q = 30 \quad \dots \quad 99,9\%$

$Q = 60 \quad \dots \quad 99,9999\%$

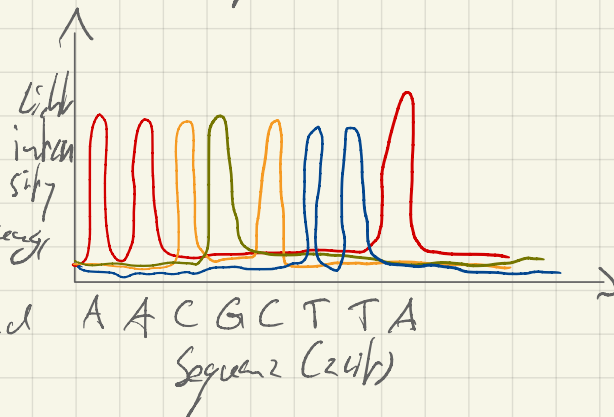
Phred ist eigentlich ein Base-Caller (\rightarrow Software, welche Licht-
signale in Rohdaten umwandelt)

Anhand bestimmter Parameterwerte

z.B. Peak farbe, Peakhöhe, Peakauflösung

Peak überlappung, ... ordnet Q Wert und

Base zu



Qualitätskontrolle:

Beliebte tools: FastQC

- Analysiert eine zufällige Teilmenge an Reads eines Fastq-Files und erstellt aus den Sequenzen und Quality-Scores verschiedene statistische Metriken. → Sollen einem helfen die Gesamtqualität des Fastq-Files abzuschätzen.
- Wurde eigentlich für Whole-Genom Sequenzierungsprojekte entwickelt, deswegen sind ein paar Metriken für RNA-Seq irrelevant und können ignoriert werden.

Read Processing:

- ① Entfernen von Adaptersequenzen (evtl. Barcodes entfernen)
- ② Abschneiden vom 5'- und 3'-Enden der Basen mit schlechter Qualität
- ③ Entfernen ganzer Reads, wenn z.B. \emptyset Qualität zu niedrig, zu kurz, zu viele N Basen

Tools: fastp, prinseq, Cutadapt, trimmomatic

Mapping:

~ Zuordnung der Reads zur ursprünglichen Transkriptionsstelle in der Referenzsequenz.

Tools: Segemehl, Hisat2, STAR 2nd
minimap2 3rd Gen (Illumina, PacBio)

Split-Read Aligner:

Warum Splitten: Introns und Splicing erkennen helfen um Sequenz zu verstehen
⇒ Infos über Expression und evtl. neue (entdeckte) Exons/Isoformen

Wenn Split nicht betrachtet: Wenn zu wenig splits an pos. oder zu langer Split (bei Hsa ≈ 200000 max)

Majorer Treffer: bis zu einem gewissen Punkt alle Positionen unter
später entscheiden wo am wahrscheinlichsten.

- Output: SAM-Format

- Besteht aus mehreren Header zeilen (@)
- Für jedes Read-Alignment eine Zeile aus mind. 11 Spalten

Wenn schlechtes Mapping:

- Probe könnte kontaminiert sein
- Falsches/schlechtes Referenz genom
- Falsche Mapping-Strategie (parameter oder Tools)
- Biologische Gründe (Mutationen etc.)

Postprocessing:

- Konvertierung in BAM (Binär) für Speicher JGV
- Filterung bestimmter Alignments JGV table
- Sortierung (ID, pos, ...), Indizierung → Visualisierung
SAMTools

Read Counting:

~ Ziel des Counting ist es jedem (bekannten) genomischen Feature die korrekte Anzahl an gemappten Reads zuzuordnen

Genomisches Feature: alle genomischen Regionen, die in irgendeiner Art exprimiert werden, aber es gibt auch Ausnahmen.
Ausnahmen: Sequenzierungsmethode bedingt.

Was wollen wir zählen?

- Ist die zugrundeliegende Annotation vollständig (genug)?
- Wie gehen wir mit multimapped Reads um?
- Reads die auf Feature mappen die sich abspalten

Differenzielle Expressionsanalyse:

Begründung der Strategie wichtig!

- Wie sieht das Verhältnis (Foldchange) der Expressionsänderung eines gens in unterschiedlichen Bedingungen aus?

→ Wir vergleichen die Countwerte jeweils eines Gens zwischen den Bedingungen
next Slides.

↳ Countwerte normalisieren!

~~_____~~
Nanopore

Software zu MinION ist MinKNOW

RAW-Signal in HDF5 (.fast5) → "Squiggles" genannt
↳ Array von Integern

Base caller ist Guppy

RNA braucht poly-A
und wird von 3' nach 5'
seq (DNA anders rum)

Anwendung

Projekt mit Gynäkologie (Plazenta labor)

Spülen Plazenta durch (Blut raus)

Schichten von Matrikalen auf embryonaler Seite was rein und
schauen was ankommt.

Trophoblastenschicht (dort treffen sich Mutter und Embryo)

Transkription aus Nabelschnur

Transkription von Trophoblastenschicht

Exosome

↳ was da drin?

→ Embryonale Exosome
kommen zu Mutter

→ Exosome für
Kommunikation?

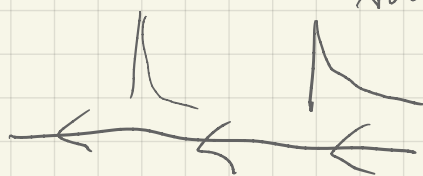
Exosome unterscheiden sich
von Mikrovesen in Größe
und Oberfläche

=> Artefakte finden! mit IGV und UCSC
Artefakt microRNA expression anstatt Protein

Transcriptome genome browser -> google map search

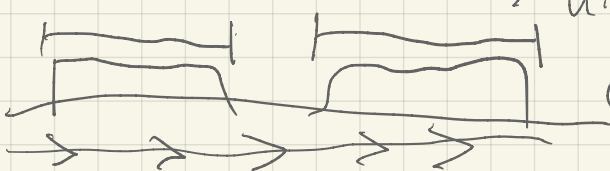
Single Nucleotide -> "fehler" bei reverse Transkript, fehlerhafte
Annotation \nwarrow post Transkript geändert.

Artefakt



-> Splicing / Sashimi plot -> lang drüber

-> nicht das Gen durchlesen



etwa 22 bp lang

2 sind festpol
1 auch möglich

Lesen von Genombildern:

Transkriptionsprofile:

In Prüfung bild lesen können

Prüfungstermin:

17.08.

845

microRNA

Eigener Dicer → Poly A

Intronrest (Lariat) → kein Poly A

Wolke Degradation

DICER teilt rest

Wegen shift auf Hairpin Zwei Micro Nukleotide
Jede non-coding RNA hat eigene Signaturen

Wollen in Proteinen sind auf Sekundärstruktur und Sequenz-
methode zurückzuführen