

Genexpressionsanalyse Projekt 2 (Emmelie E. Tessema)

Aufgabe 0

Zusammenfassend beziehen sich die Daten auf die Genexpression in *Saccharomyces cerevisiae* unter verschiedenen Bedingungen, und die Studie vergleicht die Ergebnisse der RNA-Seq- und Mikroarray-Analysen, wobei verschiedene analytische Schritte und Methoden im Prozess bewertet werden.

Überschrift	Beschreibung
Art der Daten	<ul style="list-style-type: none">Daten: transkriptomische Daten → RNA<i>Generierung der Daten</i>: RNA-Seq-Analyse und Mikroarray-Analyse<i>Transkriptom-Analyse</i>: Enthält den vollständigen Satz von RNA-Transkripten, die zu einem bestimmten Zeitpunkt vom Genom produziert werden.
Organismus	<ul style="list-style-type: none"><i>Untersuchter Organismus</i>: <i>Saccharomyces cerevisiae</i>, Stamm CEN.PK 113-7DEs wird die <i>Saccharomyces cerevisiae</i> als Bäckerhefe und Modellorganismus in der biologischen Forschung verwendet
Vergleich	<ul style="list-style-type: none">Studie vergleicht Ergebnisse zwischen RNA-Seq-Analyse und Mikroarray-AnalyseFür die Arbeit wuchs der <i>Saccharomyces cerevisiae</i> Stamm unter zwei verschiedenen Bedingungen (Batch und Chemostat)<ul style="list-style-type: none">Chemostat: kontinuierliche Zufuhr von einem Nährmedium und Entfernen von verbrauchtem MediumBatch: diskontinuierliche Kultur. Medium bleibt gleich und wird weder ergänzt noch aufgefüllt oder ausgewechseltEs werden die analytische Schritte bei RNA-Seq-Datenanalyse mit <u>Illumina-Plattform</u> bewertet und ein Vergleich basierend auf <u>Affymetrix-Mikroarrays</u> durchgeführt
Analytische Schritte	<ul style="list-style-type: none">Studie zur <i>Bewertung des Einflusses genetischer Variationen</i> auf die <i>Schätzung der Genexpression</i>Verwendung von drei verschiedenen Read-Mapping-Alignern (Gsnap, Stampy und TopHat) auf dem S288c-GenomUntersuchung der Fähigkeiten von fünf verschiedenen statistischen Methoden zur Erkennung differentieller Genexpression
Ergebnisse	<ul style="list-style-type: none">Hohe Reproduzierbarkeit zwischen biologischen ReplikatenHohe Konsistenz zwischen den beiden Plattformen für die Genexpressionsanalyse (Korrelation $\geq 0,91$)Gute Übereinstimmung bei der Identifizierung differentieller Genexpression durch verschiedene statistische Methoden
Fazit	<ul style="list-style-type: none">Die Studie vergleicht RNA-Seq und Mikroarrays für die Genexpressionsanalyse.Es wird untersucht, wie verschiedene Schritte die Analyse von RNA-Seq-Daten beeinflussen.

Aufgabe 1

```
fastq -dump--gzip --split-files SRR4535xx
```

- **fastq-dump**: SRA-Daten (Sequence Read Archive) in das FASTQ-Format zu konvertieren
- **--gzip**: komprimiert die Ausgabedateien im GZIP-Format
- **--split-files**: teilt die Reads in separate Dateien auf
 - i.d.R. Forward- & Reverse-Reads

Befehl: *fastqc *.fastq.gz*

Aufgabe 2

Probe SRR453570_2

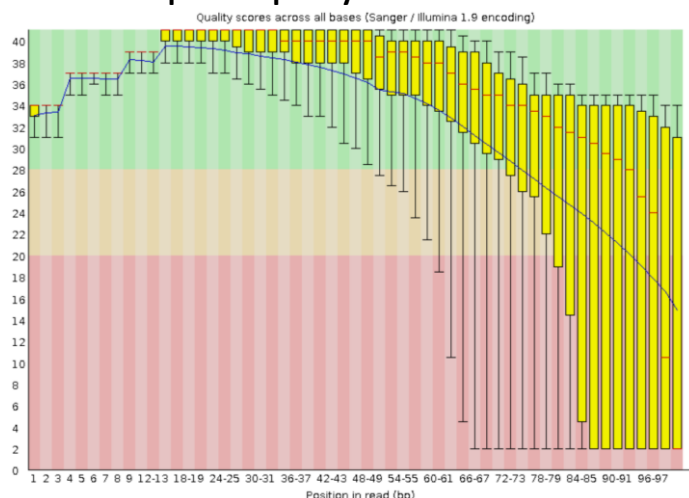
Basic Statistics

Measure	Value
Filename	SRR453570_2.fastq.gz
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	6745975
Total Bases	681.3 Mbp
Sequences flagged as poor quality	0
Sequence length	101
%GC	42

- **Fastqc**: Ruft das Programm FastQc auf
- ***.fastq.gz**: Wendet das Programm auch alle Files an

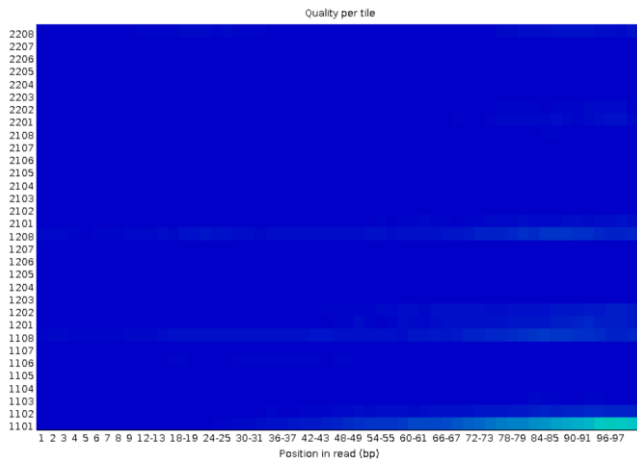
- Hier sind die allgemeinen Daten über die Sequenzierte RNA zusammengefasst
- Die Quality-Score-Kodierung basiert auf Sanger/Illumina 1.9
- Insgesamt wurden 681.3 Mbp sequenziert
- Die Read-Länge beträgt zudem 101 bp

Per base sequence quality



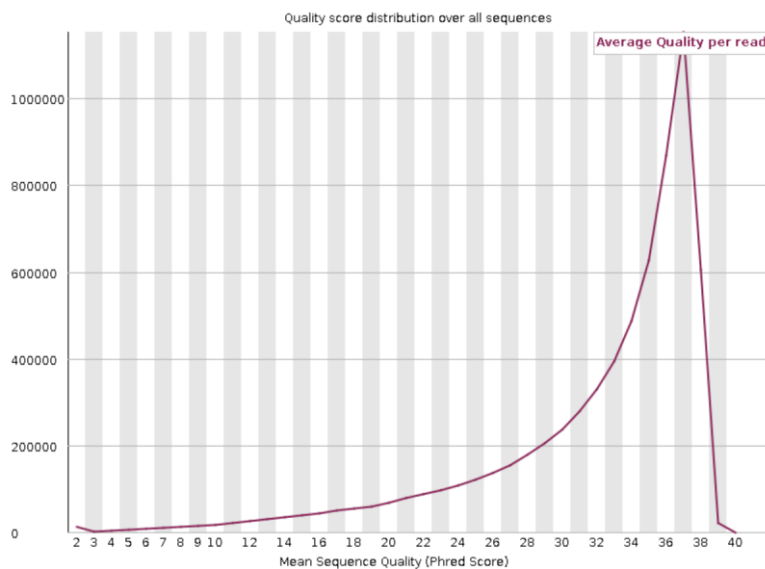
- Mit zunehmender Read-Länge ist eine deutliche Abnahme der Qualität der sequenzierten Basen zu erkennen
- Es erfolgt im Bereich der ersten Basenpaare (ca. 1bp-12bp) ein Anstieg des Quality Scores
- Im Bereich von ca. 13-19bp sind die Basen, mit den höchsten Quality Scores
- Danach folgt eine deutliche Abnahme der Qualität, was man anhand der abnehmenden Quality Scores erkennt

Per title sequence quality



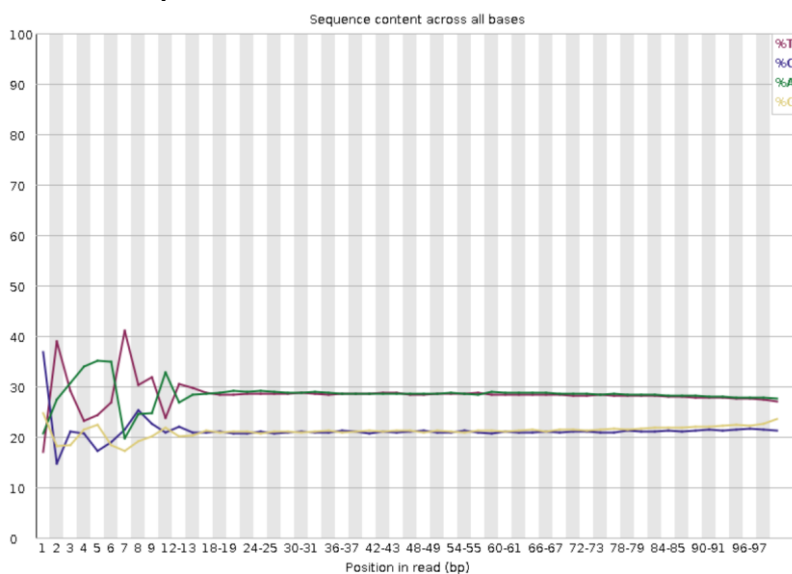
- Es ist erkennbar, dass nur minimale Fehler in der Sequenzierung von Illumina selbst vorlagen
- Die Signale sind alle überwiegend eindeutig
- Mit zunehmender Readlänge, ist das Signal nur leicht qualitativ schlechter

Per sequence quality scores



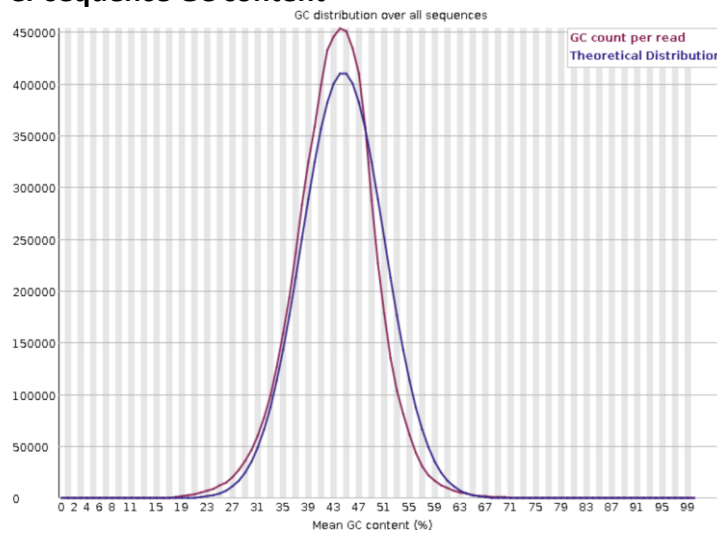
- Der größte Teil der gemessenen Basen haben durchschnittlich gesehen einen hohen Quality Score
- Es sind keine weiteren Peaks im unteren Bereich erkennbar
- Der Anstieg beginnt jedoch bereits im niedrigeren Bereich des Phred Scores (bei ca. 12-14)
- Ein erkennbarer Anteil der Reads besitzt einen niedrigen Phred Score

Per base sequence content



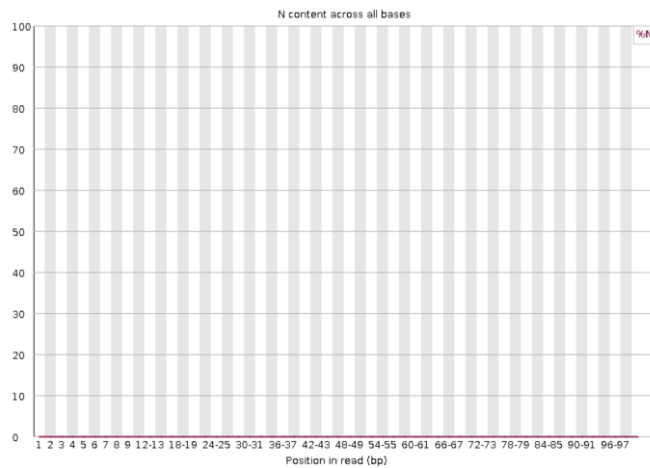
- Hier ist anfangs eine sehr schlechte Qualität (Bereich 1 bis 12bp) und danach eine relativ gute Qualität der einzelnen Basen pro Position über alle Reads zu erkennen
- In der ersten Phase sind hohe Schwankungen jedoch zu erwarten, da der Sequenzierer sich da noch in einer Art „Findungsphase“ befindet

Per sequence GC content



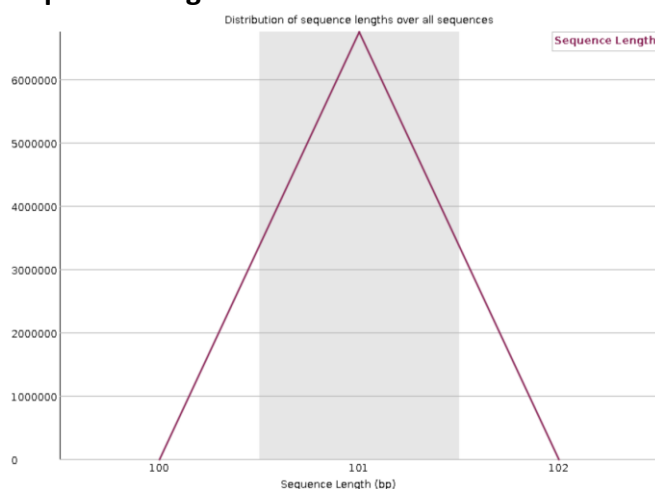
- An sich ist eine Überlappung des Peaks bezüglich des relativen durchschnittlichen GC-Gehalts im Bereich des mittleren GC-Gehalts erkennbar
- Die Daten sind zudem auch normalverteilt, was zu einer guten Bewertung der Qualität der Daten beiträgt, da der GC-Gehalt ausgewogen vorliegt
- Bei den Gemessenen Reads ist der GC-Gehalt aber bei deutlich mehr Reads aufgetreten als in dem theoretischen Vergleichsmessung
- Positiv zu bewerten ist ebenfalls, dass es keine weiteren Peaks gibt

Per base N content



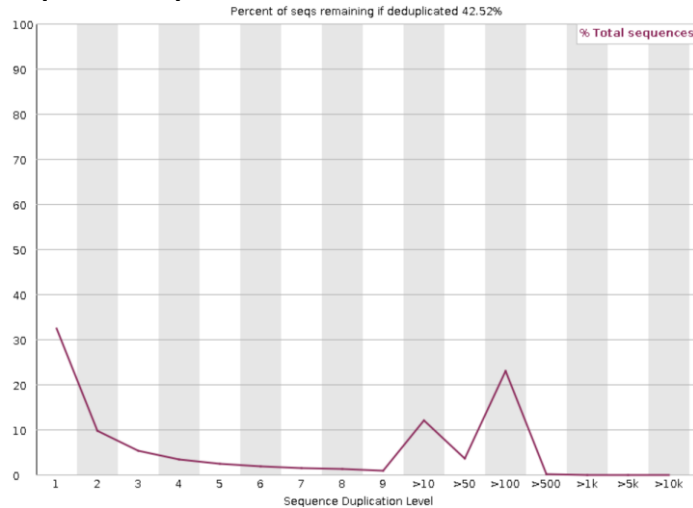
- Keine uneindeutige-Determinierung der gemessenen Basen graphisch erkennbar
- Daraus lässt sich schließen, dass die gemessenen Basen eindeutig erkannt worden sind

Sequence Length Distribution



- Die Länge der Sequenzen sind überwiegend gleich, da ein eindeutiger Peak erkennbar ist
- Der Peak liegt ungefähr bei ca. 101bp
 - Im Intervall von 100 und 102 bp
- Es ist erkennbar, dass Fragmente nur von gleicher Länge generiert worden sind
- Dies trägt zu einer positiven Bewertung der Datenqualität bei

Sequence Duplication Levels



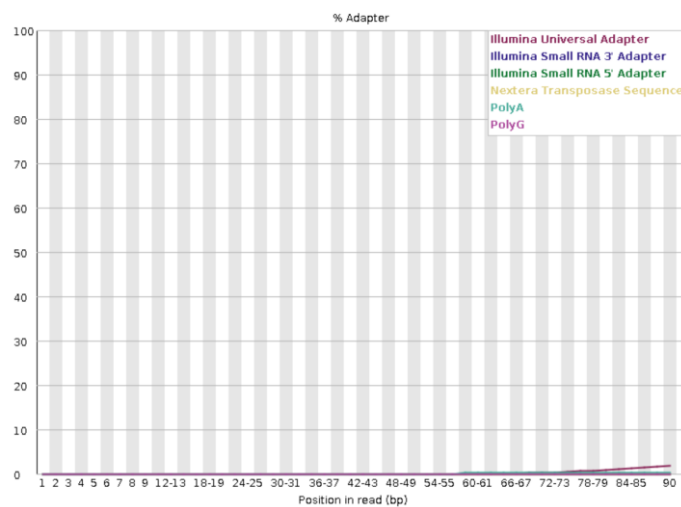
- Ein Teil der Reads, haben eine erkennbare hohe Duplikationsrate
- Die Peaks liegen im Bereich von 9 bis >500 bezüglich des Duplikationslevels
- Es sind 2 große Peaks erkennbar
 - Der größere befindet sich auf einem Duplikationslevel von ca. >100
 - Der kleinere Peak befindet sich bei ca. >10
- Da es sich um Transkriptionsdaten handelt, waren jedoch auch Duplikate zu erwarten

Overrepresented sequences

Sequence	Count	Percentage	Possible Source
GATCGGAAGAGCGTCGTGTAGGGAAAGAGTGTAGATCTCGGTGGTCGCCG	11924	0.1767572515462924	Illumina Single End PCR Primer 1 (100% over 50bp)

- Da nur eine Sequenz gelistet ist, spricht dies für eine relativ gute Qualität der gemessenen Daten

Adapter Content



- Der Illumina Universal Adapter zeigt einen leichten Anstieg im Bereich von ca. ab 78-79 bp
- Es scheinen somit viele Sequenzen am Ende einen leichten Anteil der Adaptersequenz zu enthalten

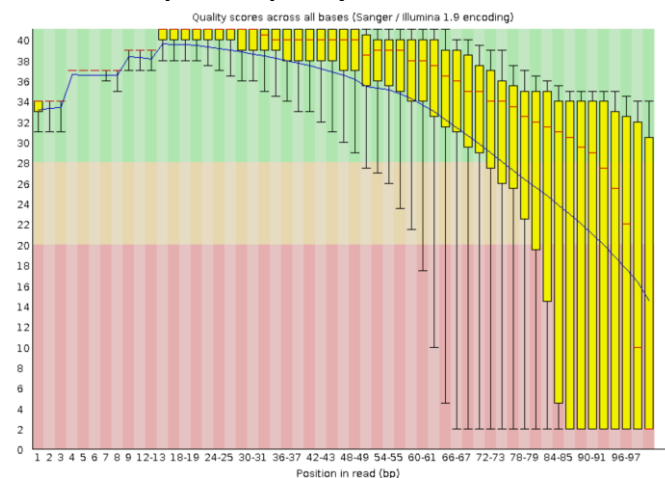
Probe SRR453571_2

Basic Statistics

Measure	Value
Filename	SRR453571_2.fastq.gz
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	6163396
Total Bases	622.5 Mbp
Sequences flagged as poor quality	0
Sequence length	101
%GC	41

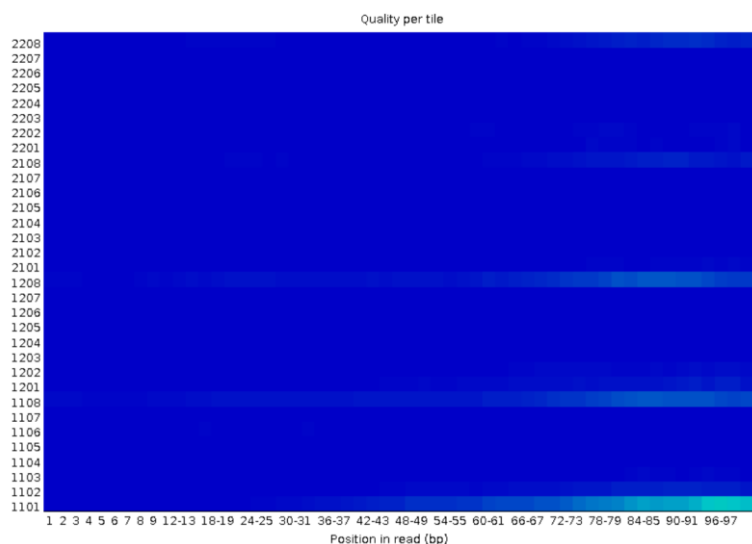
- Hier sind die allgemeinen Daten über die Sequenzierte RNA zusammengefasst
- Die Quality-Score-Kodierung basiert auf Sanger/Illumina 1.9
- Insgesamt wurden 622.5 Mbp sequenziert
- Die Read-Länge beträgt zudem 101 bp

Per base sequence quality



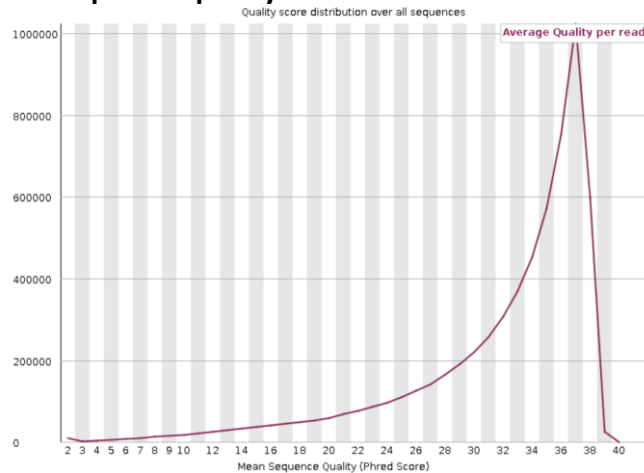
- Mit zunehmender Read-Länge ist auch hier eine deutliche Abnahme der Qualität der sequenzierten Basen zu erkennen
- Es erfolgt im Bereich der ersten Basenpaare (ca. 1bp-12bp) ein Anstieg des Quality Scores
- Im Bereich von ca. 13-19bp sind die Basen, mit den höchsten Quality Scores
- Danach folgt eine starke Abnahme der Qualität, was man anhand der kleiner werdenden Quality Scores erkennt

Per title sequence quality



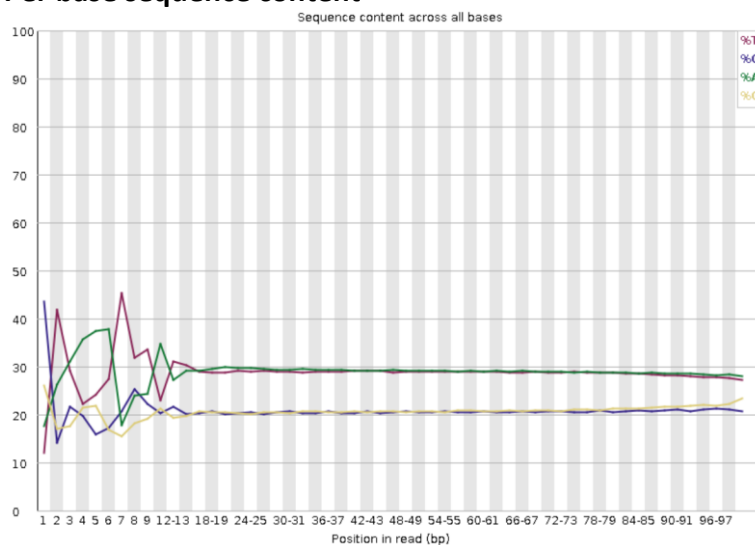
- Es ist erkennbar, dass nur minimale Fehler in der Sequenzierung von Illumina selbst vorlagen
- Die Signale sind alle überwiegend eindeutig
- Mit zunehmender Readlänge, ist das Signal nur leicht qualitativ schlechter

Per sequence quality scores



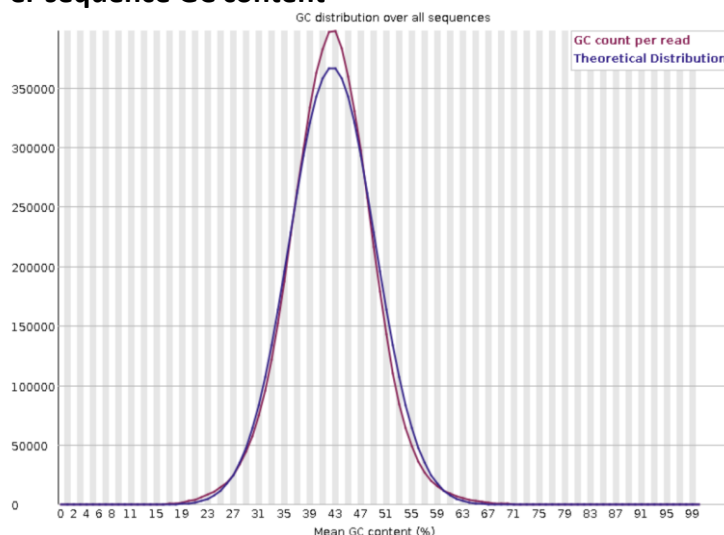
- Der größte Teil der gemessenen Basen haben durchschnittlich gesehen einen hohen Quality Score
- Es sind keine weiteren Peaks im unteren Bereich erkennbar
- Der Anstieg beginnt jedoch bereits im niedrigeren Bereich des Phred Scores (bei ca. 12-14)
- Ein erkennbarer Anteil der Reads besitzt einen niedrigen Phred Score

Per base sequence content



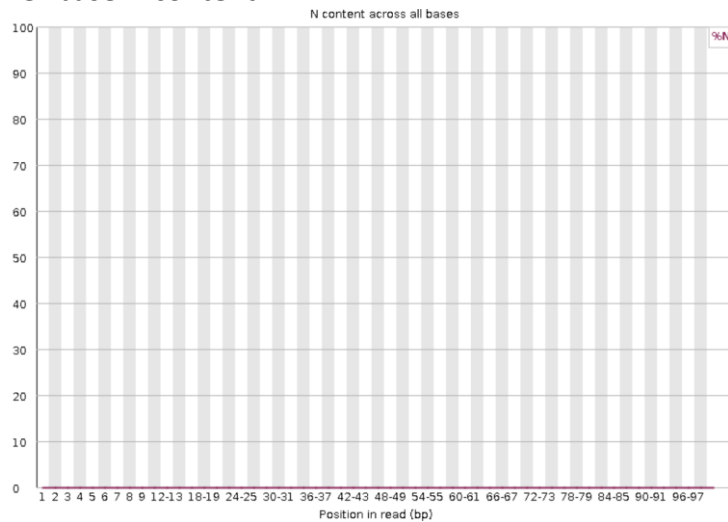
- Hier ist anfangs eine sehr schlechte Qualität (Bereich 1 bis 12bp) und danach eine relativ gute Qualität der einzelnen Basen pro Position über alle Reads zu erkennen
- In der ersten Phase sind hohe Schwankungen jedoch zu erwarten, da der Sequenzierer sich da noch in einer Art „Findungsphase“ befindet

Per sequence GC content



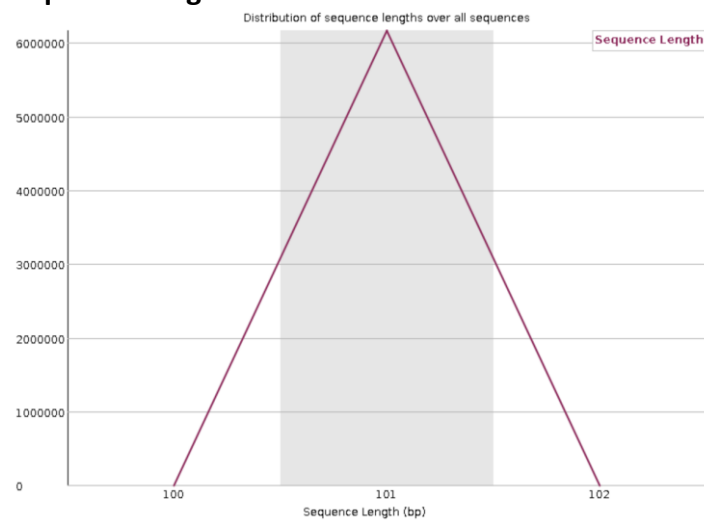
- An sich ist eine Überlappung des Peaks bezüglich des relativen durchschnittlichen GC-Gehalts im Bereich des mittleren GC-Gehalts erkennbar
- Die Daten sind zudem auch normalverteilt, was zu einer guten Bewertung der Qualität der Daten beiträgt, da der GC-Gehalt ausgewogen vorliegt
- Bei den Gemessenen Reads ist der GC-Gehalt etwas mehr Reads, als in dem theoretischen Vergleichsmessung
- Positiv zu bewerten ist ebenfalls, dass es keine weiteren Peaks gibt

Per base N content



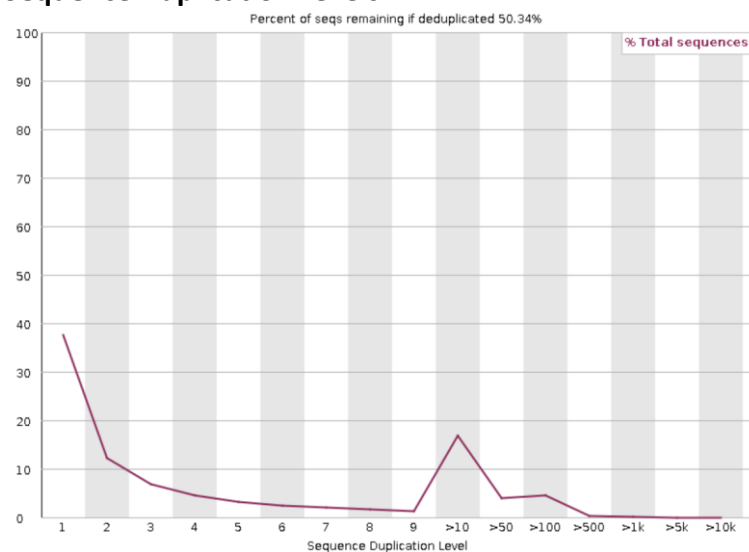
- Keine uneindeutige-Determinierung der gemessenen Basen graphisch erkennbar
- Daraus lässt sich schließen, dass die gemessenen Basen eindeutig erkannt worden sind

Sequence Length Distribution



- Die Länge der Sequenzen sind überwiegend gleich, da ein eindeutiger Peak erkennbar ist
- Der Peak liegt ungefähr bei ca. 101bp
 - Im Intervall von 100 und 102 bp
- Es ist erkennbar, dass Fragmente nur von gleicher Länge generiert worden sind
- Dies trägt zu einer positiven Bewertung der Datenqualität bei

Sequence Duplication Levels

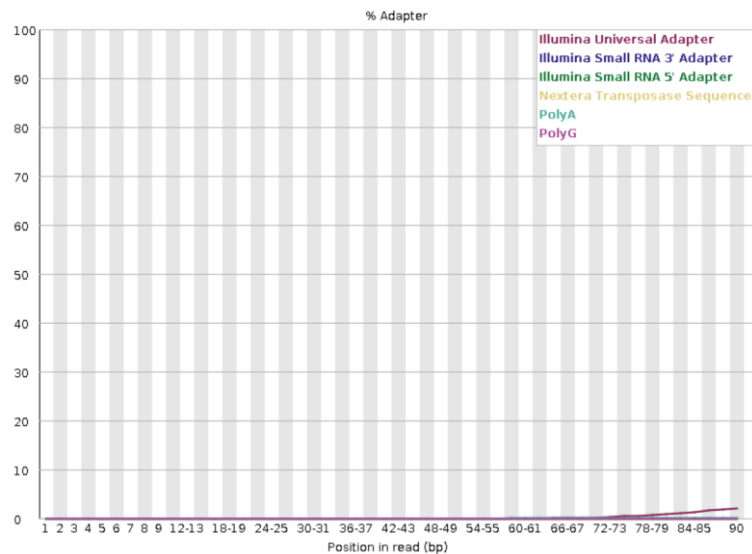


- Ein Teil der Reads, haben eine erkennbare hohe Duplikationsrate
- Es gibt einen großen Peak bei ca. >10 Duplikationslevel und ein etwas kleinerer Peak bei ca. >100 Duplikationslevel
- Die Peaks liegen im Bereich von 9 bis >500 bezüglich des Duplikation Levels
- Da es sich um Transkriptionsdaten handelt, waren jedoch auch Duplikate zu erwarten

Overrepresented sequences

No overrepresented sequences

Adapter Content



- Da keine überrepräsentierte Reads erkannt worden sind, spricht dies zusätzlich von einer guten Qualität der Daten

- Der Illumina Universal Adapter zeigt einen leichten Anstieg im Bereich von ca. ab 78-79 bp
- Es scheinen somit viele Sequenzen am Ende einen leichten Anteil der Adaptersequenz zu enthalten

Aufgabe 3

Cutadapt

- Das Programm ist in der Lage die Adaptersequenzen aus den Sequenzen zu entfernen

Trim-Galore

- Dies ist eine "Wrapper-Software", welche verschiedene Werkzeuge für das Trimmen von den Reads kombiniert
- Einschließlich Cutadapt
- Es erkennt die Adaptersequenzen und führt das Trimmen basierend auf den Qualitätsschwellenwerten durch
- Zudem ist das Tool in der Lage Reads von zu niedriger Qualität und nicht ausreichende Länge zu entfernen

Durchgeführter Befehl:

```
(cutadaptenv) → ~ trim_galore -j 4 --paired
```

- **Trim_galore**: ruft das Program Trim Galor auf
- **-j 4**: gibt an, dass der Befehl 4 Prozesse parallel ausführen soll
 - dies kann die Geschwindigkeit der Datenverarbeitung verbessern
- **--paired**: diese Option gibt an, dass es sich um gepaarte (Paired-End) Reads handelt
 - Sie enthalten somit sowohl den Vorwärts- als auch den Rückwärtsstrang

Warum ist es wichtig, dass die gepaarten Reads zusammen beschnitten werden?

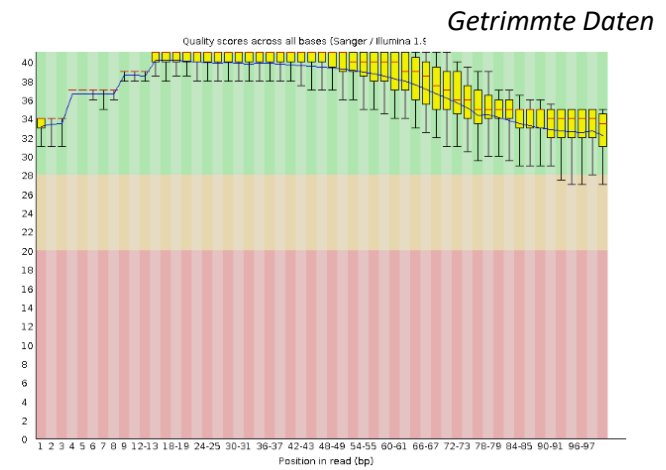
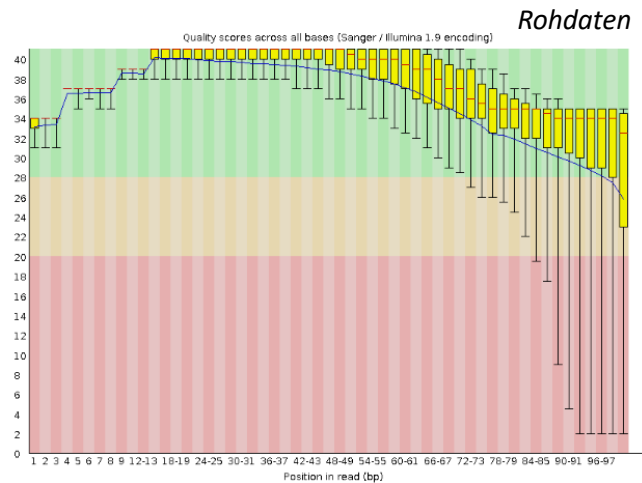
- es ist wichtig, da sie aus derselben DNA-Fragmentsequenz stammen und daher miteinander verbunden sind
- es dient dem Sicherstellen von konsistenten beschnittenen Reads
 - sonst können Unterschiede in der Datenqualität auftreten
- zudem kann man durch das gleiche Trimming der Reads ein genaueres Mapping der Reads auf das Referenzgenom ausführen
- Zudem können durch die Verbindung von Informationen aus beiden Read-Paaren, Fehler in der Sequenzierung besser erkannt werden

Vergleich zwischen den Rohdaten und den getrimmten Daten

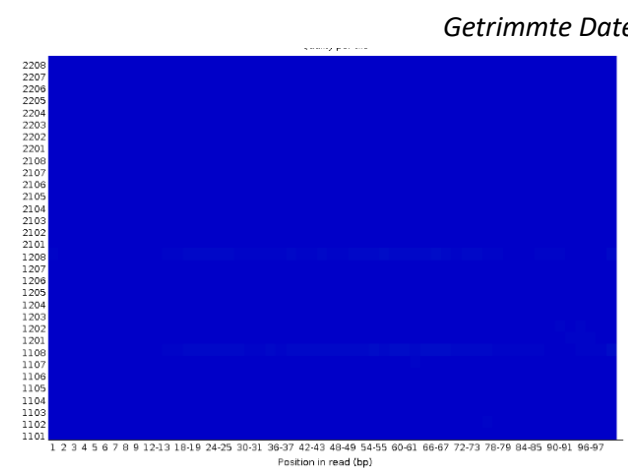
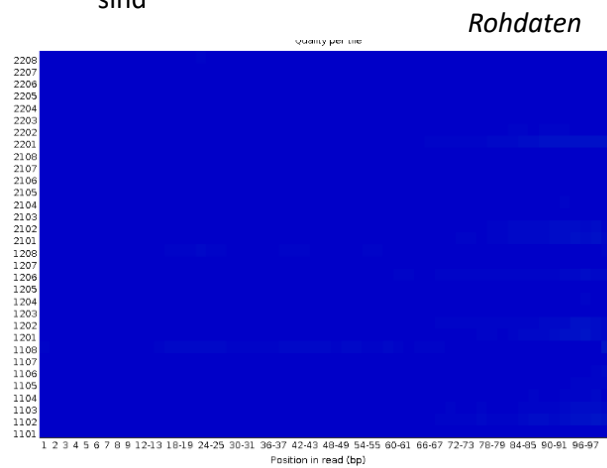
→ SRR453566_1

Rohdaten		Getrimmte Daten	
Measure	Value	Measure	Value
Filename	SRR453566_1.fastq.gz	Filename	SRR453566_1_val_1.fq.gz
File type	Conventional base calls	File type	Conventional base calls
Encoding	Sanger / Illumina 1.9	Encoding	Sanger / Illumina 1.9
Total Sequences	5725730	Total Sequences	5650497
Total Bases	578.2 Mbp	Total Bases	540.4 Mbp
Sequences flagged as poor quality	0	Sequences flagged as poor quality	0
Sequence length	101	Sequence length	20-101
%GC	41	%GC	41

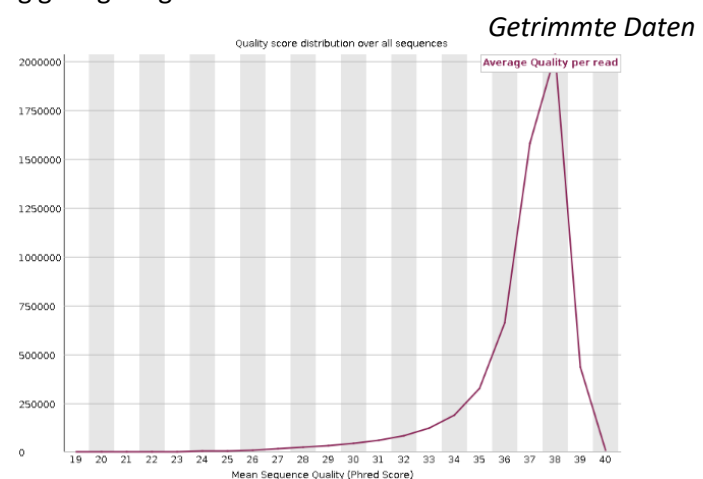
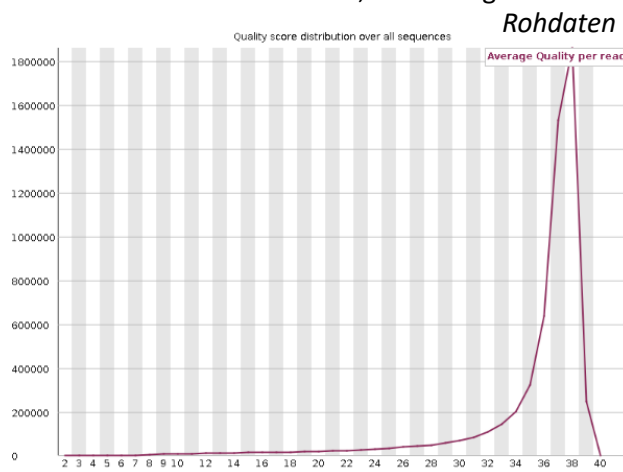
- In den Basic Statistics ist eine deutliche Abnahme der Gesamtbasen-Anzahl, der Gesamtzahl der Sequenzen sowie der Sequenzlänge erkennbar



- Es ist deutlich erkennbar, dass Reads mit einem niedrigen Qualitätsscore entfernt worden sind, sodass nur noch Daten mit einem guten bis mäßigen Qualitätsscore erhalten worden sind



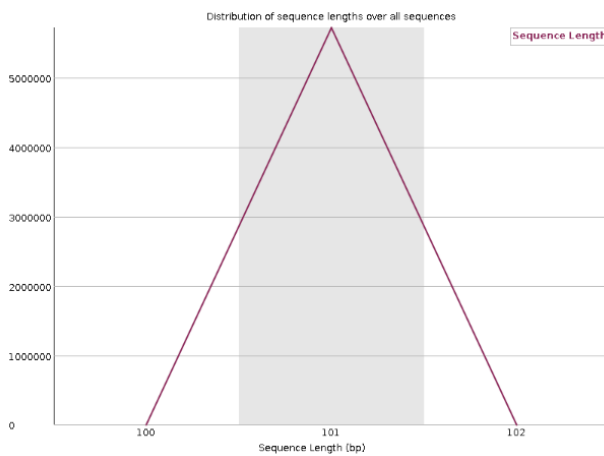
- Bei genauerem Hinschauen ist erkennbar, dass auch hier das Reads oder Teile eines Rads entfernt worden sind, dessen Signal nicht eindeutig genug ausgefallen sind



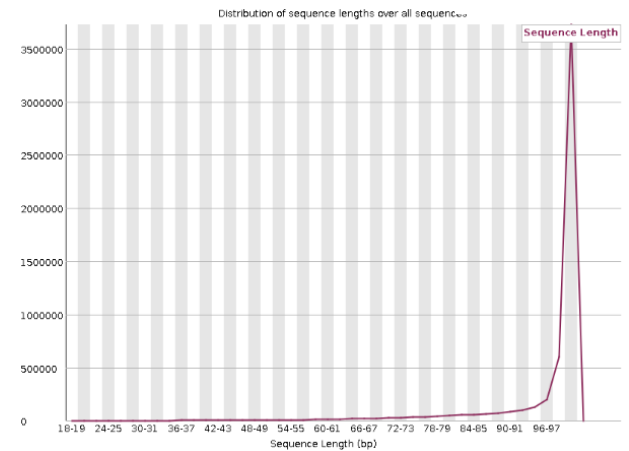
- Da alle Reads mit zu niedrigen Qualitätsscore entfernt worden sind (was man am abgeschnittenen Index erkennt), ist verhältnismäßig erkennbar, dass der Anteil an Reads mit hoher Qualität höher als vorher ist

- Beim Durchschnittlichen GC-Gehalt, beim Basengehalt pro Sequenz, beim Duplikationslevel der Sequenzen und beim N-Gehalt pro Sequenz, waren keine graphischen Unterschiede erkennbar, weshalb sie hier auch nicht mehr aufgeführt werden

Rohdaten



Getrimmte Daten



- Durch das Trimmen sind unterschiedlich lange Read-Längen entstanden, wodurch sich hier kein symmetrischer Verlauf wie bei den Rohdaten ergibt
- Jedoch bleibt der allergrößte Teil auf einer Sequenzlänge von 96-97 bis 102 erhalten

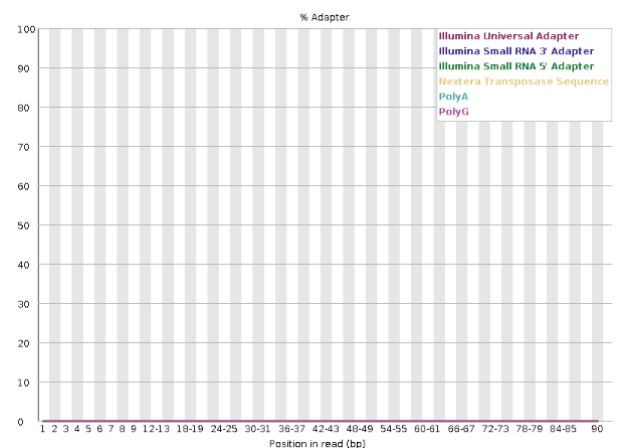
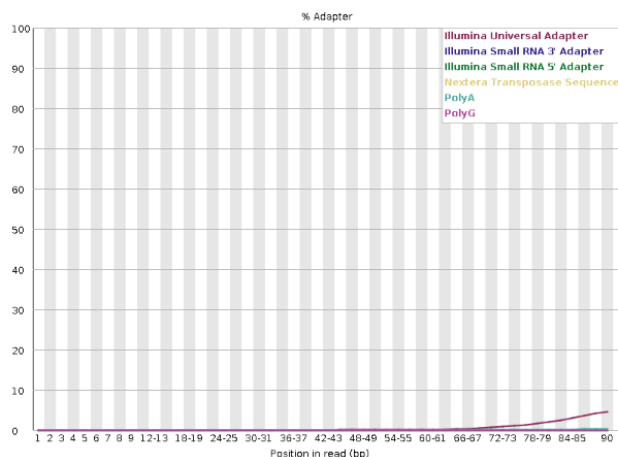
Snalfeenan

Sequence	Count	Percentage	Possible Source
GATCGGAAGAGCACACGTCTGAACTCCAGTCACCGATGATCTCGTATGC	13818	0.24119195281649677	TruSeq Adapter, Index 2 (100% over 50bp)

Getrimmte Daten

No overrepresented sequences

- Durch das Trimming wurden auch die Adaptersequenzen entfernt, die vorher zu den überpräsentierten Sequenzen gehörten



- Auch hier ist zu erkennen, dass die Adaptersequenzen aus den Reads entfernt worden sind
- Im Gegensatz zu den Rohdaten (wo die Illumina Universal Adaptersequenz noch deutlich in den Reads vorhanden war), erkennt man das anhand der Nulllinie entlang der x-Achse

→ SRR453566_2

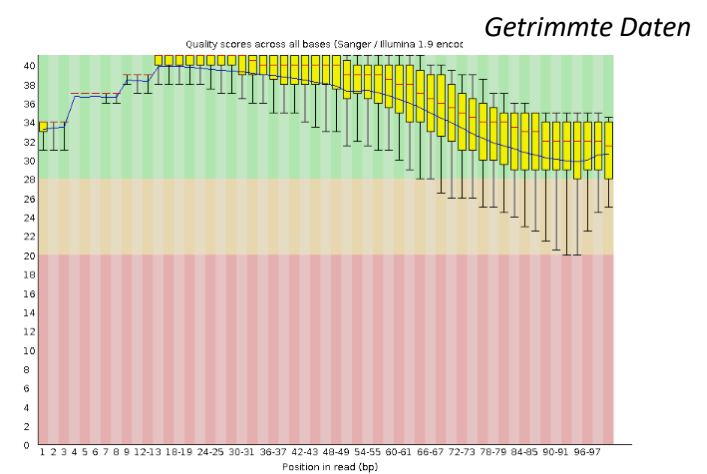
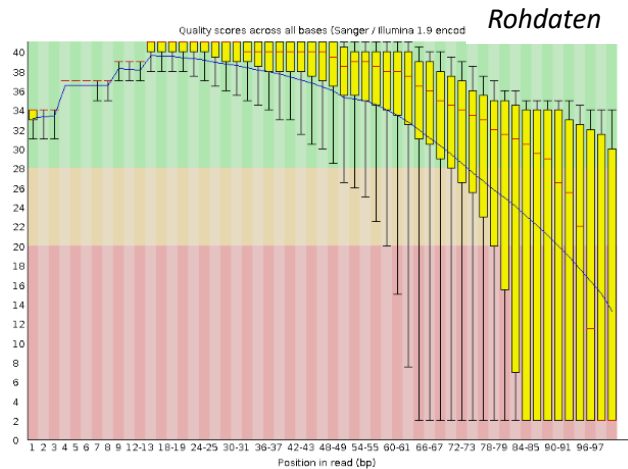
Rohdaten

Measure	Value
Filename	SRR453566_2.fastq.gz
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	5725730
Total Bases	578.2 Mbp
Sequences flagged as poor quality	0
Sequence length	101
%GC	42

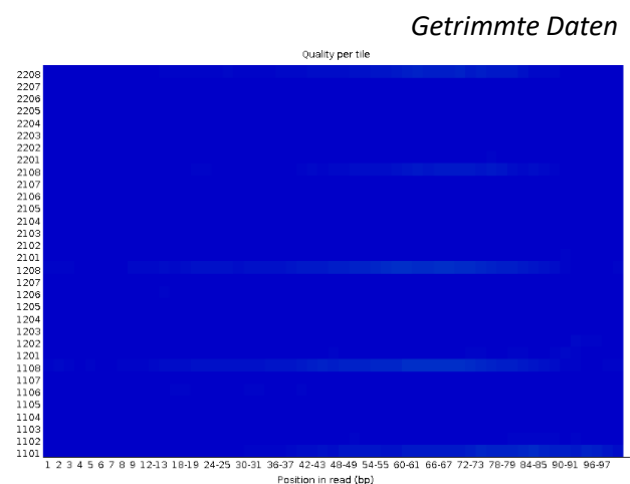
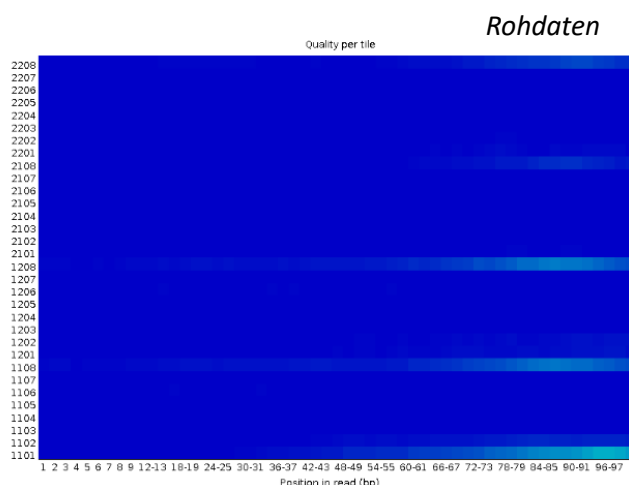
Getrimmte Daten

Measure	Value
Filename	SRR453566_2_val_2.fq.gz
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	5650497
Total Bases	496.7 Mbp
Sequences flagged as poor quality	0
Sequence length	20-101
%GC	41

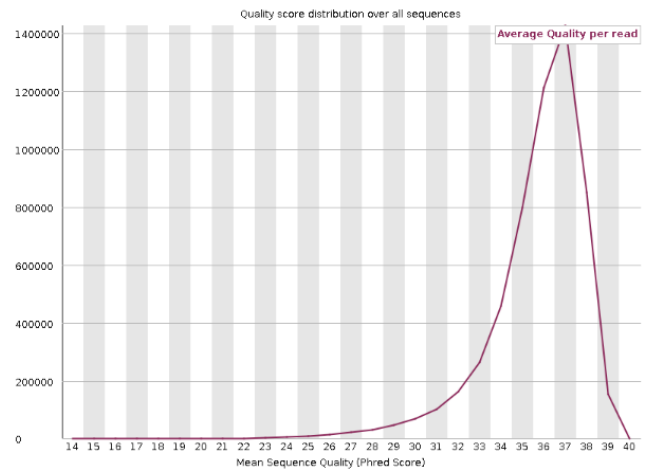
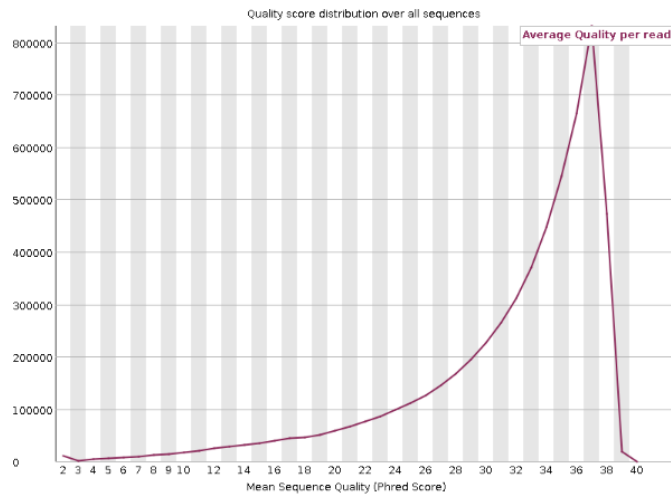
- In den Basic Statistics ist eine deutliche Abnahme der Gesamtbasen-Anzahl, der Gesamtzahl der Sequenzen sowie der Sequenzlänge erkennbar



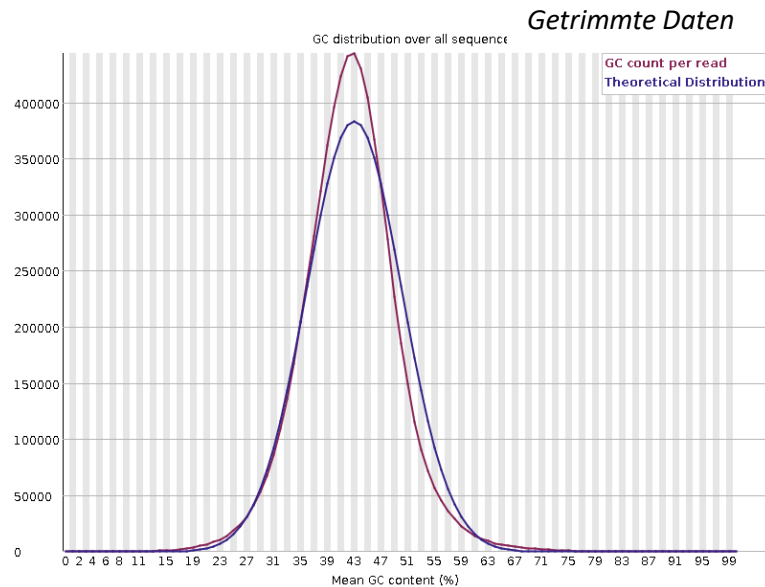
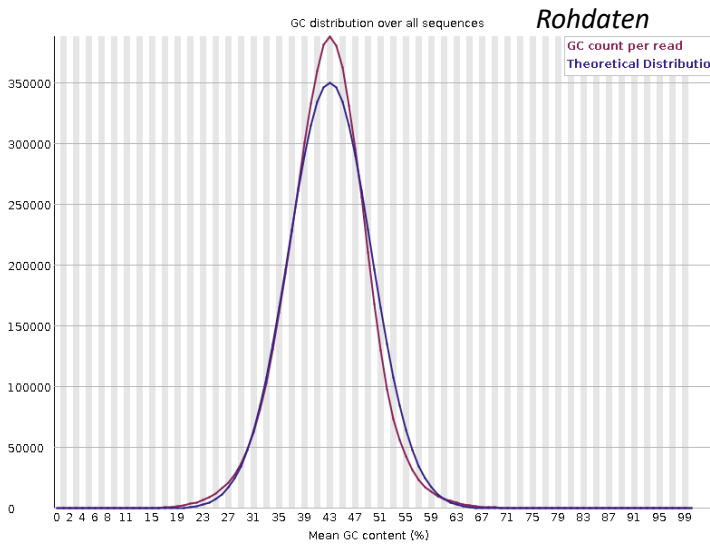
- Es ist deutlich erkennbar, dass Reads mit einem niedrigen Qualitätsscore entfernt worden sind, sodass nur noch Daten mit einem guten bis mäßigen Qualitätsscore erhalten worden sind
- Der niedrigste Qualitätsscore beträgt hier 20



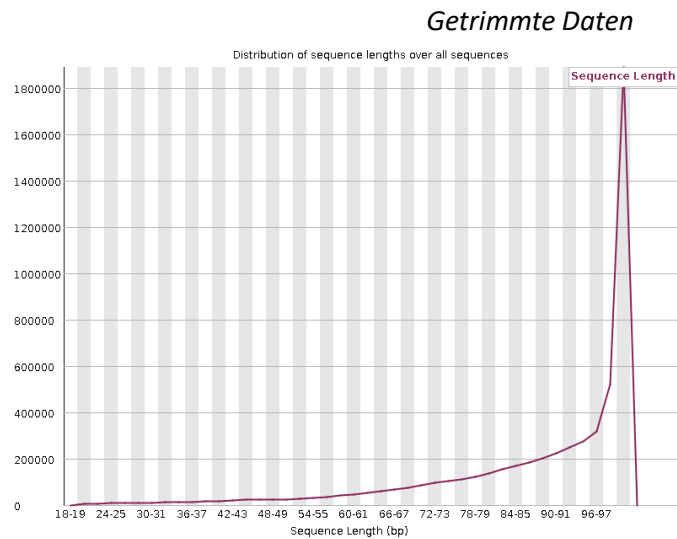
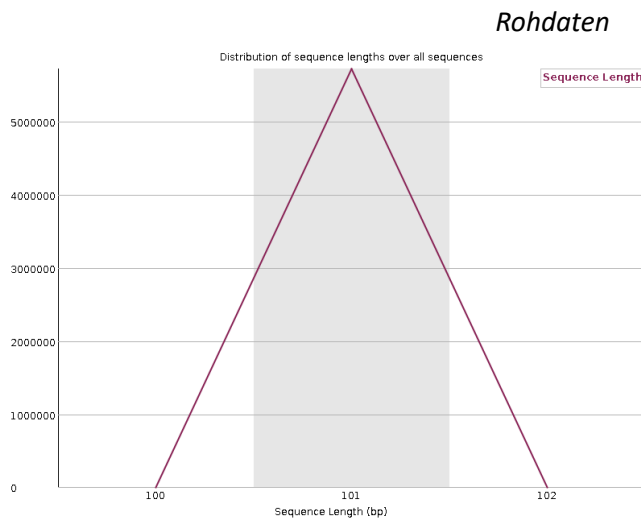
- Es ist deutlich erkennbar, dass auch hier das Reads oder Teile eines Rads entfernt worden sind, dessen Signal nicht eindeutig genug ausgefallen sind



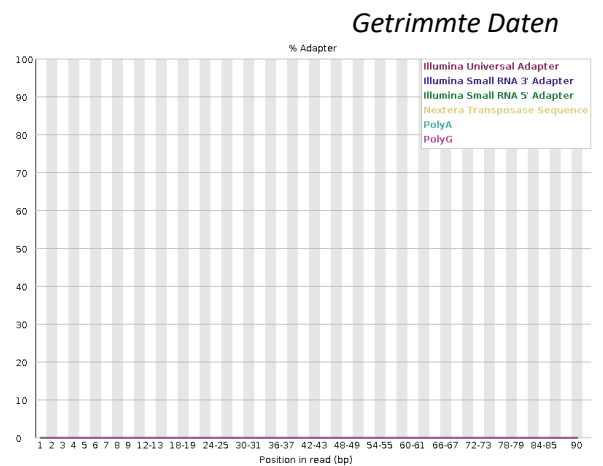
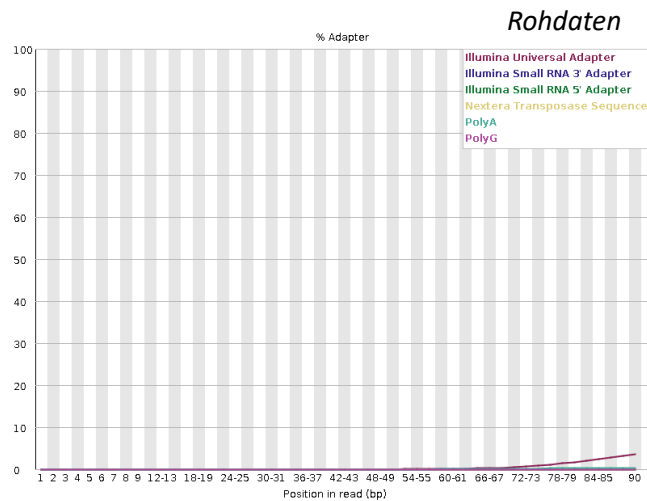
- Da alle Reads mit zu niedrigen Qualitätsscore entfernt worden sind (was man am abgeschnittenen Index erkennt), ist verhältnismäßig erkennbar, dass der Anteil an Reads mit hoher Qualität höher als vorher ist



- Da ein großer Anteil der Reads entfernt worden sind, ist nach dem Trimming auch der relative Gehalt bzw. der Durchschnittliche GC-Gehalt gestiegen, was man an dem höher liegendem Peak erkennt



- Durch das Trimmen sind unterschiedlich lange Read-Längen entstanden, wodurch sich hier kein symmetrischer Verlauf wie bei den Rohdaten ergibt
- Jedoch bleibt der allergrößte Teil auf einer Sequenzlänge von 96-97 bis 102 erhalten
- Beim Basengehalt pro Sequenz, bei der Angabe der Überpräsentierten Sequenzen, beim Duplikationslevel der Sequenzen und beim N-Gehalt pro Sequenz, waren keine graphischen Unterschiede erkennbar, weshalb sie hier auch nicht mehr aufgeführt werden



- Auch hier ist zu erkennen, dass die Adaptersequenzen aus den Reads entfernt worden sind
- Im Gegensatz zu den Rohdaten (wo die Illumina Universal Adaptersequenz noch deutlich in den Reads vorhanden war), erkennt man das anhand der Nulllinie entlang der x-Achse

Aufgabe 4

Was ist ein Index?

- Bowtie-built erstellt einen Bowtie-Index aus einer Reihe von DNA-Sequenzen
- Es wird ein Satz mit 6 Dateien rausgegeben mit folgenden Suffixen:
 - .1.ebwt, .2.ebwt, .3.ebwt, .4.ebwt, .rev.1.ebwt, .rev.2.ebwt.
- Diese Dateien bilden zusammen den Index
- Sie sind alles, was benötigt wird, um die Reads am Referenzgenom auszurichten
 - Sie bildet eine vorgefertigte Datenstruktur, wo effizient nach bestimmten Mustersequenzen gesucht werden, kann
- Die Originalsequenzdateien werden von Bowtie nicht mehr verwendet sobald er Index erstellt worden ist

Quelle: <https://bowtie-bio.sourceforge.net/manual.shtml#the-bowtie-build-indexer>

Durchgeführter Befehl:

```
GenExProjekt2 bowtie2-build ./index/Saccharomyces_cerevisiae.R64-1-1.dna.toplevel.fa.gz ./index/yeast
```

- **Bowtie2-build:** erstellt einen Index für eine Referenzsequenz, sodass das Bowtie2-Programm schneller nach Übereinstimmungen zwischen kurzen DNA-Sequenzen und dieser Referenzsequenz suchen kann.
- **./index/Saccharomyces_cerevisiae.R64-1-1.dna.toplevel.fa.gz:** Anwendung auf die im Verzeichnis './index/' liegende Datei und im gz-komprimierten FASTA-Format vorliegend
- **./index/yeast:** Dies ist der Basisname für den Index, der erstellt wird. Bei diesem Befehl ist es './index/yeast'

Durchgeführter Befehl:

```
bowtie2 -p 6 -x ./index/yeast -1 SRR453566_1_val_1.fq.gz -2 SRR453566_2_val_2.fq.gz -S SRR453566_mapped.sam
```

- **Bowtie2:** ruft das Bowtie2-Programm auf, um Alignments durchzuführen
- **-p 6:** Anzahl an Threads oder Prozessorkerne, die Bowtie verwenden soll (hier: 6 Kerne)
- **-x ./index/yeast:** Pfad zum Bowtie2-Index wird hier angegeben
- **-1 SRR4535xx_1_cutadapt_val_1.fq.gz:** Die erste jeweilige Eingabedatei.
- **-2 SRR4535xx_2_val_2.fq.gz:** Die jeweilige zweite Eingabedatei
- **-S SRR4535xx_mapped.sam:** Gibt an, dass die Ausgabe in einer SAM-Datei gespeichert werden soll

Was sind SAM-Dateien?

- = "Sequence Alignment Map"
- Ein spezielles Dateiformat, das Informationen über das Alignment von DNA- oder RNA-Sequenzen gegen eine Referenzsequenz enthält
- besteht aus einem Kopf [Header] + dem Alignment-Abschnitt

Woraus bestehen sie?

Spaltenname	Beschreibung
Header	<ul style="list-style-type: none"> • beginnt mit "@" + Identifier • indiziert den Typen + Subtypen der Header Lin <pre>@SQ SN:chr14 LN:107349540 @PG ID:bwa PN:bwa VN:0.7.7-r441 CL:bwa mem ref/seq.fa r1.fastq r2</pre> <ul style="list-style-type: none"> • @SQ: gibt den Referenzort zurück • SN: hier Chr 14 • LN: gibt die Länge des Referenzort zurück • @PG: beschreibt das Programm, womit der SAM-File generiert worden ist • PN: gibt den Namen des Programms zurück • CL: gibt eine Kopie, des Befehls, welches aufgeführt wurde zurück
QNAME	<ul style="list-style-type: none"> • = eindeutiger Identifier für den Read • Ausnahme: Gepaarte Enden haben denselben QNAME, da sie von derselben DNA-Sequenz stammen • Unterscheidung von gepaarten Reads erfolgt anhand ihrer Ausrichtung (= bestimmt durch den FLAG-Wert)
FLAG	<ul style="list-style-type: none"> • eine dezimale (Basis-10) Zahl <ul style="list-style-type: none"> ◦ verwendet, um eine binäre (Basis-2) Zahl darzustellen ◦ Ziffern repräsentieren verschiedene Wahrheits-/Falschaussagen bezüglich der Ausrichtung des Gelesenen <ul style="list-style-type: none"> ▪ 0=false ▪ 1=true • Bitfeld, das Informationen über das Mapping des Reads enthält <ul style="list-style-type: none"> ◦ über versch. Aussagen

Decimal	Binary	Exp.	Meaning
1	1	2 ⁰	This is a paired read
2	10	2 ¹	This read is part of a pair that aligned properly*
4	100	2 ²	This read was not aligned
8	1000	2 ³	This read is part of a pair and its mate was not aligned
16	10000	2 ⁴	This read aligned in the reverse direction**
32	100000	2 ⁵	This read is part of a pair and its mate aligned in the reverse direction*
64	1000000	2 ⁶	This read is the first in the pair (read 1)
128	10000000	2 ⁷	This read is the second in pair (read 2)
256	100000000	2 ⁸	The given alignment is a secondary alignment***
512	1000000000	2 ⁹	Read failed quality check (such as Illumina quality filtering)
1024	10000000000	2 ¹⁰	Read was flagged as a duplicate (such as a PCR duplicate)
2048	100000000000	2 ¹¹	Supplementary alignment (Exact meaning varies by aligner)

* Proper alignment indicates both reads in a pair are oriented towards one another (one forward, one reverse), are both on the same contig, and are within the expected distance from one another.

** Direction is relative to the reference sequence used for alignment

*** The read had multiple potential alignments; this was one of them, but not the first choice from among them

Beispiel:

- FLAG value → 99
 - 64+32+2+1
- Paired Reads kreieren immer ungerade zahlen

RNAME	<ul style="list-style-type: none"> • Name bzw. Kennung des referenzierten DNA-Segments • erscheint auch im Header
POS	<ul style="list-style-type: none"> • Startposition der gelesenen DNA-Sequenz auf dem Referenzgenom (leftmost) • 1 → die erste Position im Referenzgenom
MAPQ	<ul style="list-style-type: none"> • Mapping-Qualität des Reads • Score der angibt, wie wahrscheinlich es ist, dass die Sequenz richtig oder falsch gemappt worden ist • 255 → keine Wahrscheinlichkeit gegeben, und wird als Platzhalter Wert benutzt • Formel: $-10\log_{10}(e)$ <ul style="list-style-type: none"> ◦ e= <i>Wahrscheinlichkeit, dass die mapping Position falsch ist</i> ◦ insg. gerundet auf einen Integer • Abhängig von <ul style="list-style-type: none"> ◦ wie einzigartig die alignte Region im Genom ist ◦ Länge des Alignments ◦ Anzahl Mismatches & Gaps
CIGAR	<ul style="list-style-type: none"> • Beschreibung der Alignment-Operationen • = Concise Idiosyncratic Gapped Alignment Report String.

	<ul style="list-style-type: none"> Sequenz von Zahlen & Buchstaben, die Kontinuitäten oder Diskontinuitäten in der Ausrichtung anzeigen <ul style="list-style-type: none"> die durch eingefügte/gelöschte Basen (oder andere Gründe für Diskontinuität) verursacht werden <p>Beispiel:</p> <ul style="list-style-type: none"> 5M2D5M Die Buchstaben in der CIGAR-Sequenz haben folgende Bedeutungen: <ul style="list-style-type: none"> M: Match (Übereinstimmung) <ul style="list-style-type: none"> CAVE: es zählt auch als Match, wenn die Base nicht übereinstimmen aber zueinander gemapp sind! D: Deletion (Löschung) I: Insertion (Einfügung) S: Soft clip (Teilweise Abschneiden) H: Hard clip (Komplettes Abschneiden) N: Skipped region (Übersprungener Bereich) Die Zahlen in der CIGAR-Sequenz geben die Anzahl der aufeinanderfolgenden Vorkommen des jeweiligen Buchstabens an.
RNEXT	<ul style="list-style-type: none"> Name bzw. Kennung des nächsten referenzierten Segments entspricht Feld 3 (Referenzname)+ folgt denselben Regeln <ul style="list-style-type: none"> Ausnahme: beschreibt das gepaarte End-Mate des Reads (wenn vorhanden) Wert ist "=", wenn es identisch mit dem Referenznamenwert ist (Platz sparen)
PNEXT	<ul style="list-style-type: none"> entspricht Feld 4 + die gleichen Regeln
TLEN	<ul style="list-style-type: none"> gibt die Länge der Referenzsequenz an, auf die sich das Read abbildet. <ul style="list-style-type: none"> wird manchmal mit der Read-Länge verwechselt, ist aber oft gleichwertig Ein Read mit mehreren Einfügungen kann eine kleinere Referenzlänge haben als die Read-Länge Ein Read mit mehreren Deletionen kann eine längere Referenzlänge haben als die Read-Länge <i>Merke:</i> Wenn RNA oder cDNA auf genomische DNA abgebildet wird, kann die Referenzlänge aufgrund eines Introns bei einem kurzen Read mehrere Zehntausend Basen betragen
SEQ	<ul style="list-style-type: none"> die Read Sequenz
QUAL	<ul style="list-style-type: none"> Qualitätswerte der Basen in SEQ Phred-scaled vom FASTQ-File generiert

Quelle: <https://www.zymoresearch.de/blogs/blog/what-are-sam-and-bam-files>

Aufgabe 5: Von Alignments zu Genen

Tool: Samtools

- = eine Reihe von Programmen für die Interaktion und Nachbearbeitung kurzen DNA-Read-Alignments im SAM (oder BAM, CRAM) Format
- Es konvertiert zwischen den Formaten, führt Sortierung, Zusammenführung und Indizierungen und kann Lesevorgänge in beliebigen Regionen schnell abrufen

Durchgeführter Befehl:

```
(cutadaptenv) → GenExProjekt2 samtools view -bS SRR453566_mapped.sam > SRR453566_mapped.bam  
samtools sort SRR453566_mapped.bam > SRR453566_mapped_sorted.bam  
samtools index SRR453566_mapped_sorted.bam
```

- `samtools view -bS SRR4535XX_mapped.sam > SRR4535XX_mapped.bam`
 - Der Befehl konvertiert die übergebene SAM-Datei in das BAM Format
 - `-b`: gibt an, dass die Ausgabe im BAM-Format sein soll
 - `-s`: gibt an, dass die Eingabedatei im SAM-Format vorliegt
- `samtools sort SRR4535XX_mapped.bam > SRR4535XX_mapped_sorted.bam`
 - Der Befehl sortiert die BAM-Datei und speichert die sortierte Ausgabe in der Datei `SRR4535XX_mapped_sorted.bam`
 - Durch das Sortieren wird die BAM-Datei nach genomischer Position sortiert
- `samtools index SRR4535XX_mapped_sorted.bam`
 - `index`: Der Befehl erstellt einen Index für sortierte BAM-Datei
 - der Index erleichtert einen schnellen Zugriff auf bestimmte Bereiche des Genoms in der BAM Datei

Ausgeführter Befehl:

```
GenExProjekt2 sed -i 's/^chr//g' yeast_genes.bed  
GenExProjekt2
```

- mit diesem Befehl wird "chr" am Anfang jeder Zeile in der `yeast_genes.bed` entfernt und speichert die Änderungen direkt in der Datei

Tool: Bedtools

- ein Toolkit, welches mehrere Operationen ermöglicht, um sie auf genomischen Intervallen (BED-Dateien) auszuführen
- Hauptfunktionen: genomische Intervalle aus mehreren Dateien überschneiden, zusammenzuführen, zu zählen, zu ergänzen und zu mischen

Quelle: <https://bedtools.readthedocs.io/en/latest/>

Ausgeführter Befehl:

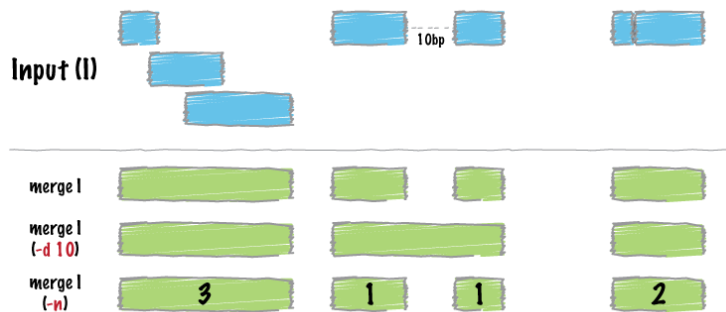
```
(cutadaptenv) → GenExProjekt2 bedtools multicov -bams SRR453566_mapped_sorted.bam SRR453567_mapped_sorted.bam SRR453568_mapped_sorted.bam SRR453569_mapped_sorted.bam SRR453570_mapped_sorted.bam SRR453571_mapped_sorted.bam -bed yeast_genes.bed >Bedtools_out_all.bed
```

- `bedtools multicov`: Anzahl der Reads bestimmen, die sich in bestimmten genomischen Bereichen befinden
 - Bereiche zwischen einer .bed-Datei und BAM-Dateien (die Alignments der Reads) analysieren
- `-bams SRR453566_mapped_sorted.bam SRR453567_mapped_sorted.bam SRR453568_mapped_sorted.bam SRR453569_mapped_sorted.bam`

SRR453570_mapped_sorted.bam SRR453571_mapped_sorted.bam: gibt die BAM-Dateien an, die für die Überlappungsberechnung verwendet werden sollen

- **-bed yeast_genes.bed**: Gibt die .bed-Datei an, die die genomischen Regionen definiert, über die die Überlappungen berechnet werden sollen
- **> Bedtools_out_all.bed**: Leitet die Ausgabe des Befehls in eine Datei namens Bedtools_out_all.bed

Graphische Darstellung



Quelle: https://bioweb.pasteur.fr/docs/modules/bedtools/2.19.1/_images/intersect-glyph.png