

Master Thesis
Seperating the good from the bad...
Exploring the genomic landscape of
chloroplasts from genomic sequencing data

Simon Pfaff



Julius-Maximilians-Universität Würzburg
Fakultät für Biologie

Master Thesis

Simon Pfaff

August 7, 2018

Seperating the good from the bad... Exploring the genomic landscape of
chloroplasts from genomic sequencing data



Julius-Maximilians-Universität Würzburg

Betreuer: Dr. Markus Ankenbrand

Betreuer: Prof. Dr. Jörg Schulz

Betreuer: Dr. Frank Förster

Lehrstuhl für Bioinformatik

Center for Computational and Theoretical Biology

Contents

1	Zusammenfassung / Abstract	4
1.1	Deutsch	4
1.2	English	4
2	Einleitung	5
2.1	Chloroplasten	5
2.2	Big Data	6
2.3	Open Science	6
2.4	Daten in Daten	9
2.5	Bestehende Programme und ihre Ansätze	9
2.5.1	chloroExtractor	10
2.5.2	fast-plast	12
2.5.3	NOVOPlasty	14
2.5.4	Org.ASM	14
2.5.5	GetOrganelle	14
2.5.6	IOGA	15
2.6	Interesse an Chloroplasten, was tun damit mit diesen Daten?	15
2.7	Aufgaben in der Master Thesis	17
3	Material / Methoden	18
3.1	Evaluation der Programme	18
3.1.1	Testdaten	18
3.1.2	Welche Programme werden weiter verwendet.	20
3.2	Varianz Analyse	22
3.3	GWAS	22
3.4	Struktur Varianz Analyse	23
3.5	Neue Chloroplasten	23
4	Ergebnisse	23
4.1	Automatisierung	23
4.2	Daten: Simulierte Daten	24
4.3	Daten: 1001 Genom Projekt, 11 Testdatensätze	24
4.4	Daten: GO-Preprint	26
4.5	Die besten Programme: fast-plast und chloroExtractor	30
4.6	1001 Genom Projekt	33
4.7	Varianz Analyse	33
4.8	GWAS	34
4.9	Struktur Varianz Analyse	34
4.10	Neue Chloroplasten	36
5	Diskussion	36
5.1	Definition von Success, Einteilung der Erfolge über Genom Länge.	36

5.2	Die Entscheidung für fast-plast und chloroExtractor	38
5.3	Fazit aus der Erfolgschance	38
5.4	Erhöhen der Erfolgsrate	39
5.5	Etablieren einer einfachen scanning Routine	39
5.6	GWAS	39
5.7	Struktur Varianz	40
5.8	Fazit und Zukunftsaussichten	40
6	Abbildungs- und Tabellenverzeichnis	40
7	Anhang	42
7.1	Dockercontainer	42
7.2	Danksagung	44
7.3	Skripte	44
7.4	Tabellen	44

1 Zusammenfassung / Abstract

1.1 Deutsch

In der heutigen Big Data Ära, werden immer mehr neue Daten erzeugt. Doch können auch bereits vorhandene Daten durchaus noch Informationen enthalten welche bisher ungenutzt sind. Dank der Open Science sind viele dieser Daten frei verfügbar. In diesem konkreten Falle werden verschiedene Programme verwendet um in pflanzlichen Genom Daten, Chloroplasten Genome zu suchen. Chloroplasten haben ihre eigene zirkuläre DNA, und wenn beim Sequenzieren diese nicht gefiltert wurde, dann ist sie noch in den Daten vorhanden. So kann man ohne eine neue Sequenzierung machen zu müssen Chloroplasten DNA für Analysen erhalten. Es wurden verschiedene Programme getestet und verglichen, um so viele Chloroplasten wie nur möglich zu bauen und damit Analysen durchzuführen. Es wurde eine vollautomatische Pipeline erstellt mit denen Chloroplasten Genome einfach aus Pflanzen Genomen extrahiert werden können.

1.2 English

In to days Big Data Era every day new data is created. But this is not necessary to get new information. In older data oftentimes hides other data, which may not be used till today. Since open science is growing in the last years many of this data is freely accessible. In this case plant whole genome data is used. This data oftentimes includes chloroplast genome data, when in the sequencing process nobody removed it, by cell nucleus extraction. There were different tools which provided a suit to extract this chloroplast genome data from plant genome data. They were tested, compared and the best were used to calculate a lot of chloroplast genomes. This chloroplast genomes got analysed with different methods. Also there is now a fully automatic pipeline for chloroplast genome extraction from genomic sequencing data.

2 Einleitung

2.1 Chloroplasten

Chloroplasten sind extrem wichtig, ein Großteil der Pflanzen, vor allem grün Pflanzen besitzen diese und nutzen sie um Energie durch Photosynthese zu erzeugen¹. Ihr Genom gilt als hoch konserviert, sowohl in der Gen Orientierung als auch dem Gen Inhalt ², dieser besteht zwischen 100 und 200 Genen³. Die DNA des Chloroplasten hat in etwa eine Konturlänge von 30 bis 60 Mikrometer und besitzt eine Masse von ungefähr 80 - 130 Millionen daltons⁴. Chloroplasten zeigen eine auffällige Genom Struktur, das Genom ist ein Plasmid welches sich in drei Teile aufteilt. Dem Large Single Copy, dem Small Singles Copy, welche beide durch zwei Inverted repeats unterbrochen sind (Fig. 1)⁵. Chloroplasten zeigen eine viel geringere Substitutionsrate als in genomischer DNA, diese ist noch einmal signifikant geringer in den Inverted repeat Regionen ⁶. Dennoch zeigen sich Einzelnucleotid-Polymorphismen (SNPs)⁷. Zudem gibt es in seltenen Fällen eine Genwanderung von Genen auf dem IR zum SSC, wodurch ein IR weg fallen kann, sodass solche Chloroplasten nur noch ein IR besitzen ⁸. Vor allem in gezüchteten Nutzpflanzen finden sich auch Invertierungen des IR ⁹. Durch ihre hohe Konservierung sind Chloroplasten und ihre Gene sehr gut für Barcoding geeignet¹⁰. Mit diesem Barcoding können Pflanzen und ihre Varianten identifiziert werden. Neue Studien zeigen, dass der komplette Chloroplast selbst als ein Art "Ultra-barcode" verwendet werden könnte, da die Variation in Chloroplasten in einer Spezies doch mehr variiert als angenommen ¹¹. Die Gene auf einem Chloroplasten lassen sich in verschiedene Klassen einteilen (Fig. 2). Zunächst

¹Purves Biologie, Sadava D, Hillis D.M, Heller H.C, Berenbaum M.R, (9. Auflage S. 14)

²Raubeson L, Jansen R. (2005). Chloroplast genomes of plants, Plant diversity and evolution: genotypic and phenotypic variation in higher plants. Diversity and Evolution of Plants-Genotypic and Phenotypic Variation in Higher Plants. 3. 10.1079/9780851999043.0045.

³Howe C.J. (2016). Chloroplast Genome. In eLS, John Wiley & Sons, 10.1002/9780470015902.a0002016.pub3

⁴Burgess, Jeremy (1989). An introduction to plant cell development. Cambridge: Cambridge university press. S. 62. ISBN 0-521-31611-1.

⁵Shaw J, Lickey EB, Schilling EE, et al. (2007). Comparison of whole chloroplast genome sequences to choose noncoding regions for phylogenetic studies in angiosperms: The tortoise and the hare III. American Journal of Botany. 10.3732/ajb.94.3.275.

⁶Wolfe KH, Li WH, Sharp PM. (1988). Rates of nucleotide substitution vary greatly among plant mitochondrial, chloroplast and nuclear DNA. Proc Natl Acad Sci USA 10.1073/pnas.84.24.9054.

⁷1001 Genomes Consortium, (2016) 1,135 genomes reveal the global pattern of polymorphism in *Arabidopsis thaliana*. Cell. <https://doi.org/10.1016/j.cell.2016.05.063>

⁸Jansen RK, Wojciechowski MF, Sanniyasi E, et al. Complete plastid genome sequence of the chickpea (*Cicer arietinum*) and the phylogenetic distribution of rps12 and clpP intron losses among legumes (Leguminosae). Molecular phylogenetics and evolution. 10.1016/j.ympev.2008.06.013.

⁹Palmer JD, Jansen RK, Michaels HJ, et al. (1988). Chloroplast DNA Variation and Plant Phylogeny. Annals of the Missouri Botanical Garden, 10.2307/2399279

¹⁰Song Y, Wang S, Ding Y, et al. (2017) Chloroplast genomic resource of Paris for species discrimination. Sci. Rep. 10.1038/s41598-017-02083-7

¹¹Kane N, Sveinsson S, Dempewolf H, et al.(2012), Ultra-barcoding in cacao (*Theobroma* spp.; Malvaceae) using whole chloroplast genomes and nuclear ribosomal DNA. American Journal of Botany, 10.3732/ajb.1100570

gibt es die Gene welche für die Photosynthese wichtig sind, hier unterteilt man nochmal in Photosystem I (psaA, psaB, etc.), Photosystem II (psbA, psbB, etc.), Cytochrome b6f (petA, petB, etc.), ATP Synthase (atpA, atpB, etc.), RuBisCo(rbcL) und NAD(P)H dehydrogenase Gene(ndhA, ndhB, etc.)¹². Die zweite Klasse beinhaltet Gene welche für den Genetischen Apparat, also Transkription, Translation und Replikation nötig sind, sowie RNA Gene. Hierzu zählen Transfer RNA (trnH, trnK, etc.), ribosomale RNA (rrn16, rrn5, etc.), RNA Polymerasen (rpoA, rpoB, etc.) und ribosomale Gene (rps2, rps3, rpl2, rpl16, etc.). Die dritte und letzte Kategorie beschreibt Gene welche konservierte Open Reading Frames (ORFs) haben und ycfs (Hypotheticalchloroplast open reading frames) genannt werden sowie potentiell kodierende Gene wie matK und cemA¹². Vor allem stark konservierte Gene wie rbcL und matK und cemA werden häufig für Barcoding oder zum berechnen von phylogenetischen Bäume verwendet.

2.2 Big Data

Die Big Data Ära zeichnet sich vor allem durch die Flut an Daten aus, welche kaum noch zu bewältigen ist. Im biologischen Sinne zeichnet sich diese Flut an Daten vor allem durch genomische Daten aus. Durch sogenannte Hochdurchsatz Methoden in modernen Sequenzierungs Technologien wie PacBio¹³ oder Illumina¹⁴ sie verwenden, können sehr viele genomische Daten in kurzer Zeit für vergleichsweise wenig Geld erzeugt werden. Um dieser Daten Herr zu werden sind vor allem Programme die diese Daten auswerten wichtig. Die Anforderung an diese Programme sind unter anderem eine hohe Geschwindigkeit und vor allem eine hohe Automatisierbarkeit. Um solche Programme zu entwickeln sind vor allem Kenntnisse in Informatik und in Biologie notwendig. Kenntnisse in Informatik werden benötigt um Programme sinnvoll zu strukturieren und diese anschließend in einer geeigneten Programmiersprache zu erstellen. Zudem müssen diese Programme immer wieder gewartet werden. Auch könnten neue Features eingebaut werden oder dass Programm nochmals Optimiert werden. Kenntnisse in der Biologie werden benötigt um biologische Fragestellungen sowie Sachverhalte zu verstehen. Diese müssen korrekt in das Programm implementiert werden, zudem müssen natürlich alle Ergebnisse korrekt interpretiert werden. Nur mit Kenntnissen in beiden Fächern können Programme welche in der Biologie und Bioinformatik verwendet werden effektiv erstellt und gewartet werden.

2.3 Open Science

Mit der Flut der ansteigenden Daten wächst in den letzten Jahren auch immer mehr die Akzeptanz zur "Open Science". "Open Science" bezeichnet eine Bewegung welche fordert, dass Wissenschaft und alles was dazugehört, Daten, Programme, Ergebnisse öffentlich und für jedermann zugänglich sein soll. Befürworter dieser Bewegung argu-

¹²Ravi V, Khurana JP, Tyagi AK, et al. (2007). An update on chloroplast genome. Plant Systematics and Evolution. 10.1007/s00606-007-0608-0.

¹³<https://www.pacb.com/>

¹⁴<https://www.illumina.com/>

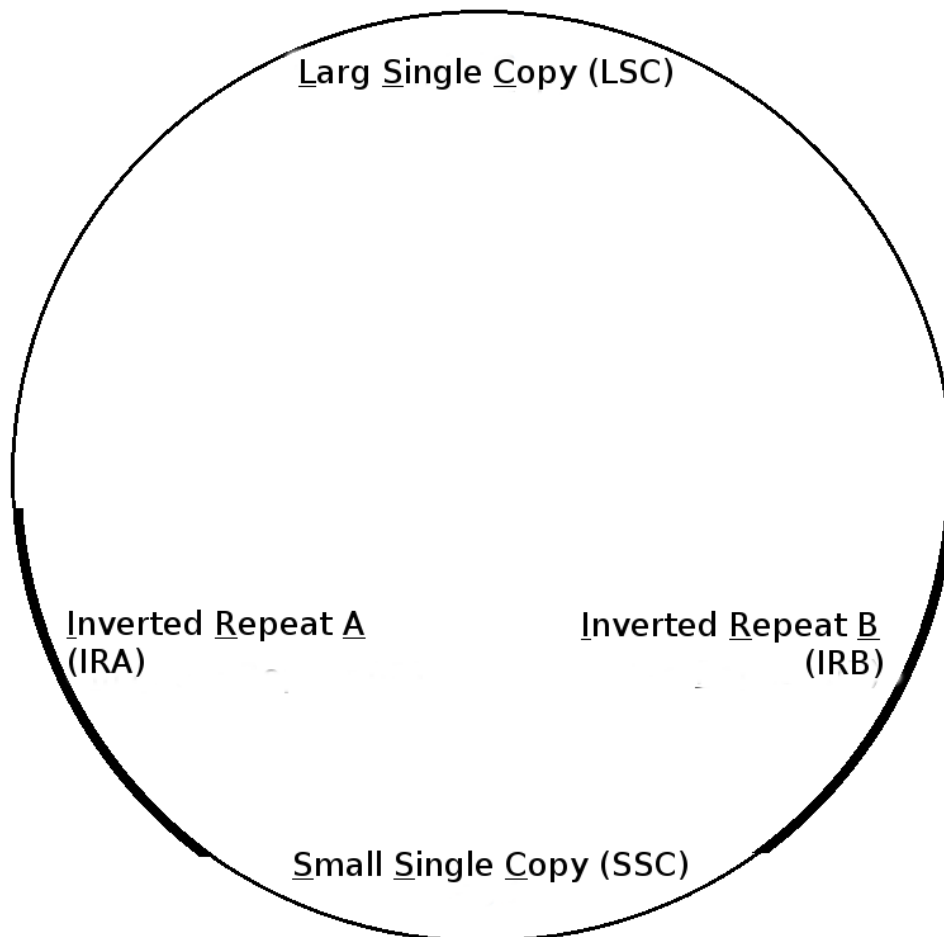


Figure 1: **Chloroplast Genom Einteilung** Der Chloroplast ist unterteilt in Large Single Copy, Small Single Copy und Inverted Repeat, diese unterteilen sich nochmal in IRA und IRB. LSC und SSC werden jeweils von den IRs unterbrochen.

mentieren damit, dass Wissen für jedermann zugänglich sein sollte, und nicht nur für Ausgewählte oder gar gegen Bezahlung. So steigere sich unter anderem die Akzeptanz der Wissenschaft als auch deren Glaubwürdigkeit. Da die Ergebnisse von jedem nachvollziehbar veröffentlicht werden müssen, mit allen Rohdaten und Vorgehensweisen. Dies sei der eigentliche Gedanke hinter der Wissenschaft, sie solle jedem zugänglich sein. Diese Bewegung findet vor allem bei jung Wissenschaftler aber auch bei Älteren immer mehr Anklang. Mittlerweile gibt es mehrere Lizenz Modelle die unter Open Science laufen. Diese Regeln wie die Daten verwendet werden dürfen oder müssen. Dies reicht von Freigeben der Daten und jeglichem Verwendungszweck bis hin zum Zwang, dass alles was mit diesen Daten oder auch Programmen veröffentlicht wird wieder unter der gleichen Open Source Lizenz zu publizieren ist. Alle hier verwendeten Programme und Daten sind unter Open Source Lizenzen veröffentlicht, sonst wäre diese Arbeit gar nicht möglich. Deswegen werden alle Ergebnisse wiederum öffentlich verwendbar sein. Open Science sollte keine Bewegung sein, sondern einfach nur Science.¹⁵

2.4 Daten in Daten

Bei den heutzutage geringen Kosten, Daten vor allem genomische Daten, zu erzeugen ist es nicht verwunderlich dass immer neue Daten generiert werden. Dennoch steckt in bereits erhobenen Daten meist mehr Information als zunächst verwendet. In genomischen Daten zum Beispiel, hier findet sich meistens Daten von Organellen, wie Mitochondrien oder Chloroplasten, welche ihre eigene DNA besitzen. Diese sind dort zu finden da vor einer Sequenzierung häufig keine Kern Extraktion durchgeführt wird, da diese mehr Zeit und Geld kosten würde. Diese Organellen DNA können mit bestimmten Programmen gefiltert werden, hierfür wurde unter anderem der chloroExtractor^{16, 17} programmiert. Dieser kann in genomischen Pflanzen Daten Chloroplasten DNA finden und diese verwenden um einen vollständigen Chloroplasten zu bauen. Hiermit müssen somit keine neuen Sequenzierungen für Chloroplasten mehr durchgeführt werden, wenn man an Chloroplasten forschen möchte.

2.5 Bestehende Programme und ihre Ansätze

Es gibt verschiedene Ansätze um Chloroplasten Genome bzw. ihre DNA aus genomischen Pflanzen Daten zu extrahieren. Die wohl einfachste Möglichkeit ist ein Referenz basiertes Mapping der Daten auf einen Referenz Chloroplasten. Hierzu muss lediglich ein nah verwandter Chloroplast als Referenz benutzt werden. So können die Reads, welche auf diese Referenz passen genommen werden und assembliert werden, mit der gleichen Referenz. Dies funktioniert allerdings nur wenn man eine passende Referenz benutzt, diese sollte von der gleichen Spezies oder zumindest einer Nah verwandten Spezies

¹⁵Watson M. (2015) When will ‘open science’ become simply ‘science’?, Genome Biology <https://doi.org/10.1186/s13059-015-0669-2>

¹⁶Ankenbrand MJ, Pfaff S, Förster F, et al., (2018). chloroExtractor: extraction and assembly of the chloroplast genome from whole genome shotgun data. Journal of Open Source Software, <https://doi.org/10.21105/joss.00464>

¹⁷<https://github.com/chloroExtractorTeam/chloroExtractor>

stammen. Ein anderer Ansatz besteht darin den Chloroplasten de novo zu assemblieren, also ohne Referenz. Um diesen Ansatz zu benutzen müssen aber zunächst die Reads mit Chloroplasten Genom aus den Daten gezogen werden. Hier gibt es wiederum verschiedene Möglichkeiten. Eine Möglichkeit ist es die Reads gegen eine Datenbank von Chloroplasten Genen zu blasten, hierzu muss entweder eine Datenbank von Chloroplasten Genen gestellt werden oder der Benutzer muss eine Pseudo-Referenz einen sogenannten Seed angeben. Ein Seed, was von einigen Basenpaaren bis zu einem kompletten Chloroplasten reichen kann, kann auch eingesetzt werden um durch ein reines Mapping Reads zu finden. Bei kleinen Seeds wird dieser häufig durch gefundene Reads erweitert und eine Liste von Seed erstellt. Auch hier muss aber sichergestellt werden, dass der Seed in den Chloroplasten Daten vorhanden ist. Von diesen Methoden gibt es auch Abwandlungen, wie z.b. das scannen der Daten durch Kmers, hier werden die Daten in verschiedene Kmers zerteilt, durch plotten dieser Kmers können an spezifischen Stellen überrepräsentierte Kmers gefunden werden, diese überrepräsentierten Kmer spiegeln häufig Plastome wieder. Diese sind unter anderem Chloroplasten aber auch Mitochondrien, sie besitzen ihre eigene DNA und kommen im Schnitt häufiger vor als DNA welche im Zellkern zu finden ist. Allerdings gibt es weitaus mehr DNA welche durch häufiges vorkommen überrepräsentiert sind in einem Kmer Plot, hierzu gehören rRNA Gene, Transposons und andere genomische Repeats, welche je nach Art und Spezies variieren kann. Da der chloroExtractor einen Kmer basierten Ansatz benutzt ist ein solches idealisiertes Kmer Diagramm in dessen Logo zu finden (Fig. 3). Abgesehen von den Ansätzen der Programme gibt es zwei verschiedene Arten von Programmen per se, die einen benutzen bereits vorhandene Programme wie Assambler, Mapper oder Kmer-counter. Diese Bauen eine Pipeline um diese Programme, sodass diese in der richtigen Reihenfolge mit den richtigen Parametern mit nur einem Befehl gesteuert werden können. Der Vorteil ist, solche Programme sind einfacher zu warten da sie meist kleiner sind als Programme die dies nicht tun und einfacher zu programmieren. Allerdings sind sie von diesen drittanbieter Programmen abhängig und es können Probleme auftreten wenn diese Änderungen bzw. Updates ausgeben, weswegen meist die kompatiblen Versionen angegeben werden. Ein weiterer Nachteil, der Benutzer muss häufig weitere Programme, sogenannte Abhängigkeiten installieren bevor er das eigentliche Programm nutzen kann. Die andere Möglichkeit ist es die komplette Maschinerie selbst zu programmieren, dies ist sehr aufwendig und bedeutet viel Wartungsarbeit. Vorteil hier ist das keine anderen Abhängigkeiten benötigt werden außer ein System welches das Programm verwenden kann. In dieser Arbeit wurden verschiedene Typen von Programmen verwendet. Es wurden von allen Programmen die jeweils neusten Versionen benutzt, und wenn es zu großen Änderungen wie Bug-fixes kam auf die neuere Version gewechselt, um das bestmögliche Ergebnis für die Daten zu erhalten.

2.5.1 chloroExtractor

Der chloroExtractor (Versionen: 1.0.2, 1.0.3, 1.0.4) ^{16, 17} ist ein Programm welches durch eine Kombination aus Kmer Analyse und Mapping auf bekannte Chloroplasten Gene, Reads von Chloroplasten aus Pflanzen Sequenzierungs Daten extrahiert. Es wurde 2018 vom chloroExtractorTeam veröffentlicht ¹⁶ und besteht hauptsächlich aus Perl und R

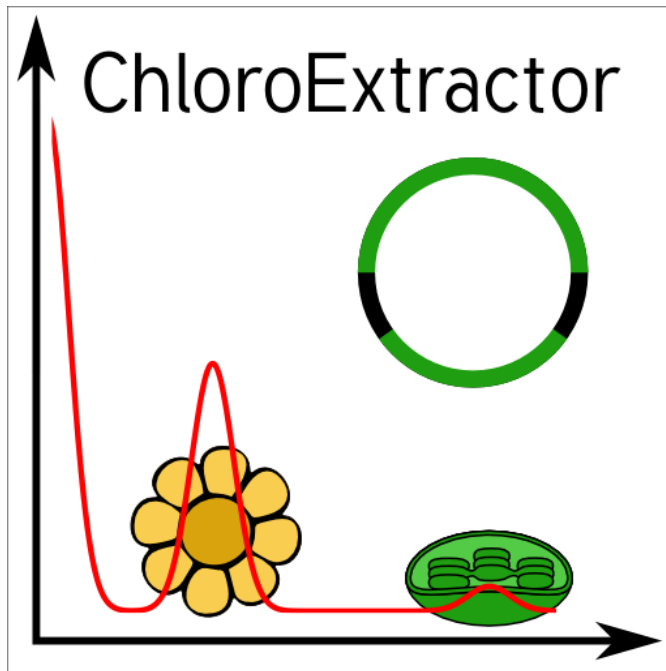


Figure 3: **chloroExtractor Logo** Das Logo des chloroExtractors zeigt die Verteilung der Genomischen Daten in einem Kmer Plot. Der erste Peak zeigt die Kmers Pflanzengenom, der zweite kleinere Peak zeigt die Kmers mit Chloroplasten Genom.

Code. Es verwendet ein Pipeline Programm (PipeWrap.pm¹⁸) um den richtigen Ablauf zu steuern. Dieses Pipeline Tool wird durch eine Konfigurationsdatei gesteuert, sodass ein Benutzer einfach neue Schritte einfügen könnte. Auch können hiermit einfach über eine Datei, Parameter gesteuert werden welche dann in allen verwendeten Programmen gleich sind. Es könnte so auch einzelnen Programmen speziellen Input mitgegeben werden. Auch verfügt der chloroExtractor dank PipeWrap über ein Checkpoint System. Bricht der Ablauf des Programms ab, kann er am genau diesem Punkt wieder gestartet werden ohne das Programm von neu starten zu müssen. Zunächst verwendet der chloroExtractor Bowtie2¹⁹ um die Reads auf eine Referenz aus codierenden Chloroplasten Gene zu mappen. Hierdurch wird der relative Anteil an Chloroplasten geschätzt. Durch ein R Skript wird der Datensatz auf 200-fache Chloroplasten Reads Coverage skaliert. Anschließend wird iterativ durch Jellyfish²⁰ Kmere erzeugt und jene mit zu niedriger Coverage aussortiert. Die Kmere die im richtigen Coverage Bereich liegen werden anschließend mit SPAdes²¹ assembliert, SPAdes arbeitet de novo und benötigt keine Referenz. SPAdes verwendet eine De Bruijn-Graphen Methode um die Reads richtig zusammen zu fügen, diese werden dann durch ein Perl Skript (fcg.pl) zu einem zirkulären Chloroplasten zusammengebaut. Dieses Skript überprüft gleichzeitig mit BLAST+²² ob es sich bei den ausgegebenen Reads wirklich um Chloroplasten handelt (Fig. 4). Falls es dazu kommt das SPAdes den Chloroplasten nicht komplett zusammenbauen kann, dann gibt das fcg.pl Skript die Contigs, welches für den Chloroplasten verwendet werden aus. Hier gibt es verschiedene Fälle. Kann nur die Zirkularität des Chloroplasten nicht aufgelöst werden gibt der chloroExtractor LSC, SSC und IR aus. Sind gar keine Verbindungen der Contigs möglich gibt das fcg.pl Skript jene Contigs aus die einen BLAST+ Treffer besitzen und somit ein teil des Chloroplasten sind. Es wurden drei Verschiedene Versionen des chloroExtractors verwendet, dies brachten unter anderem Bug fixes welche das Programm zum Absturz brachten. Aber auch Verbesserungen am fcg.pl Skript.

2.5.2 fast-plast

Fast-plast (Version: 1.2.8)²³ ist ein weiteres Programm, welches verwendet wird um Chloroplasten DNA zu finden. Es ist in Perl und in C++ programmiert und verwendet auch SPAdes, aber zusätzlich Bowtie1 sowie Bowtie2. Auch hier wird Blast+ verwendet um die richtigen Reads zu finden.

¹⁸<https://github.com/BioInf-Wuerzburg/perl5lib-PipeWrap>

¹⁹Langmead B, Salzberg S. (2012), Fast gapped-read alignment with Bowtie 2. *Nature Methods*, 9: 357-359.

²⁰Marcais G, Kingsford C.(2011), A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics* 10.1093/bioinformatics/btr011

²¹Bankevich A, Nurk S, Antipov D, et al. (2012), SPAdes: A New Genome Assembly Algorithm and Its Applications to Single-Cell Sequencing, *Journal of Computational Biology*, 10.1089/cmb.2012.0021

²²Camacho C, Coulouris G, Avagyan V, et al. (2009), Selecting control genes for RT-QPCR using public microarray data, *BMC Bioinformatics* <https://doi.org/10.1186/1471-2105-10-42>

²³<https://github.com/mrmckain/Fast-Plast>

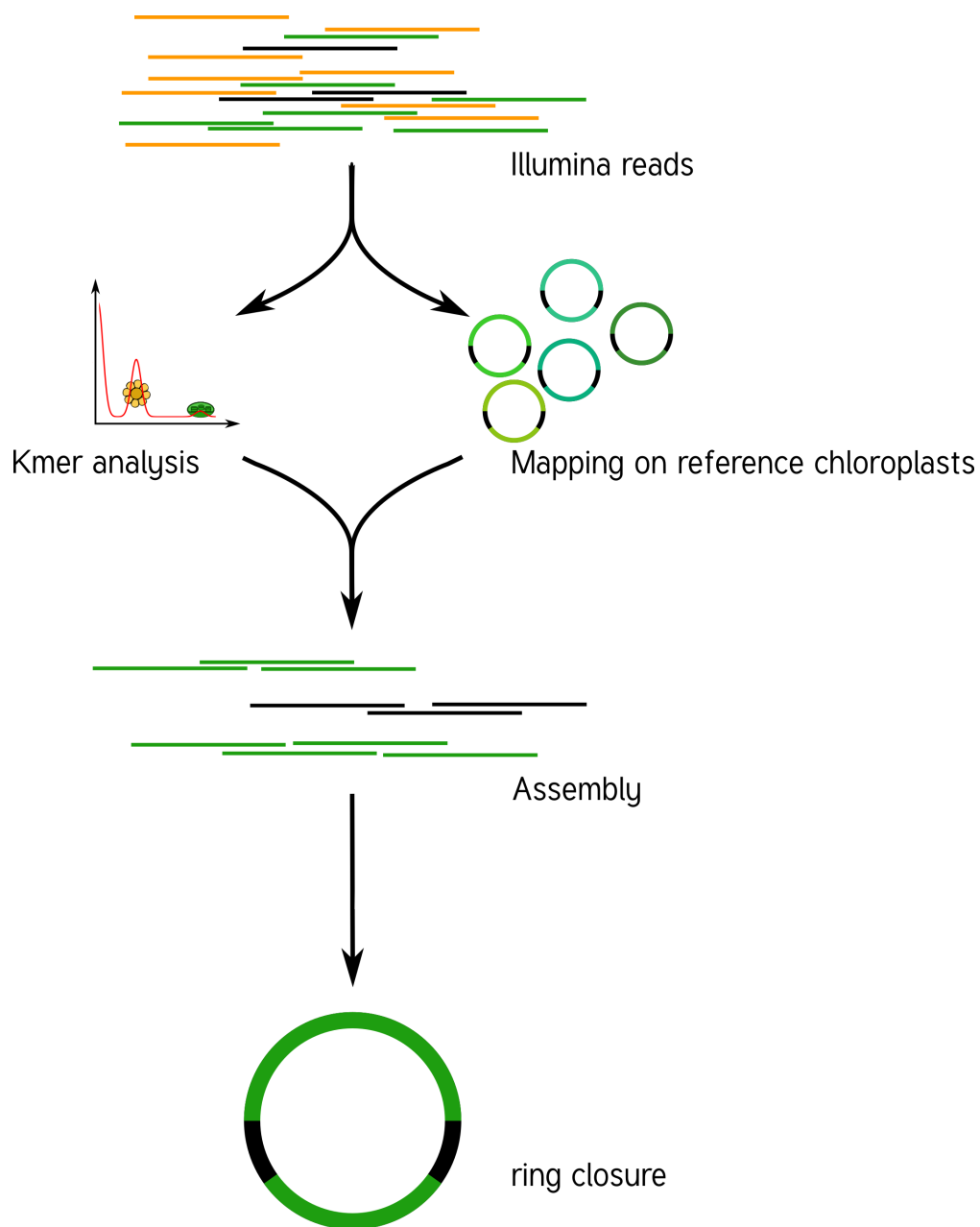


Figure 4: **Ablauf des chloroExtractors** Eine Kombination aus Kmer Analyse und Mapping auf bekannte Chloroplasten rekrutieren Chloroplasten reads um diese anschließend zu Assemblieren um anschließend einen Ringschluss herbeizuführen. (Ankenbrand et al., (2018).)

2.5.3 NOVOPlasty

Im Gegensatz zu den anderen verwendeten Programmen, benutzt NOVOPlasty (Versionen: 2.6.8, 2.6.9, 2.7.0)^{24, 25} keine drittanbieter Programme. Es benötigt somit keine Abhängigkeiten von anderen Programmen und ist komplett in Perl programmiert. NOVOPlasty benutzt sogenannte Seeds um Chloroplasten DNA zu finden, dies können einzelne Chloroplasten Gene sein, aber auch ein kompletter Chloroplast. Die Verschiedenen Verwendeten Versionen versprochen Bug fixes sowie neue Features.

2.5.4 Org.ASM

Org.ASM (Version: 1.0.00-alpha11)²⁶ ist ein Programm hauptsächlich geschrieben in Python. Es versucht überrepräsentierte Sequenzen zu finden und diese zu assemblieren²⁷. Mit Hilfe eines Seeds versucht er diese Sequenzen zu finden. Chloroplasten und andere Organellen wie Mitochondrien sind in Zellen überrepräsentiert, vor allem wenn man eine geringe Coverage über das Pflanzen Genom hat, somit sind diese detektierbar²⁸.

2.5.5 GetOrganelle

GetOrganelle (Versionen: 1.9.82, 1.0.1, 1.0.3)^{29, 30} verwendet zum lokalisieren der Chloroplasten Reads ähnlich wie andere Programme Bowtie2¹⁹ und Blast+, nur muss hier eine Referenz mitgegeben werden. Diese wird nur hierfür verwendet, das assemblieren hingegen geschieht de novo mit SPAdes. Wie auch beim chloroExtractor wird hier der fastg-Graph verwendet um den Chloroplasten zu finden, aber dies muss in fälle des GetOrganelle per Hand, mit Hilfe des Programms Bandage³¹ vollzogen werden. Wie bereits erwähnt nutzt der chloroExtractor ein Perl Skript welchen diesen händischen Schritt automatisiert. (Fig. 5) Vom GetOrganelle wurden drei Versionen verwendet. Zunächst 1.9.82, diese wurde geändert zu 1.0.1 (github commit: b390260 vom 31. März 2018) und 1.0.3. In den Verschiedenen Versionen gab es diverse Bug Fixes, sowie kleine Features. Zudem wurde das Programm GetOrganelle mit 1.0.1 in einer Wissenschaftlichen Arbeit veröffentlicht³⁰.

²⁴<https://github.com/ndierckx/NOVOPlasty>

²⁵Dierckxsens N, Mardulyn P, Smits G. (2016), NOVOPlasty: De novo assembly of organelle genomes from whole genome data. Nucleic Acids Research, 10.1093/nar/gkw955

²⁶<https://pythonhosted.org/ORG.asm/>

²⁷<https://git.metabarcoding.org/org-asm/org-asm/wikis/home>

²⁸<https://pythonhosted.org/ORG.asm/algorithms.html>

²⁹<https://github.com/Kinggerm/GetOrganelle>

³⁰Jin J, Yu W, Yang J, Song Y, et al. (2018), GetOrganelle: a simple and fast pipeline for de novo assembly of a complete circular chloroplast genome using genome skimming data. bioRxiv, <http://doi.org/10.1101/256479>

³¹Wick RR, Schultz .B, Zobel J, et al (2015). Bandage: interactive visualisation of de novo genome assemblies. Bioinformatics, <https://doi.org/10.1093/bioinformatics/btv383>

2.5.6 IOGA

Der Iterative Organellar Genome Assembly, kurz IOGA (Keine Versionsnummer vergeben, github commit: c460ea9 vom 10. Sep. 2016)^{32,33} verwendet BBmap³⁴ für das filtern und trimmen der reads, um anschließend mit SOAPdenovo2³⁵ und SPAdes²¹ die Reads zu assemblieren. Auch dieses Programm benötigt eine Referenz. Der IOGA ist in Python geschrieben.

2.6 Interesse an Chloroplasten, was tun damit mit diesen Daten?

Mit der steigenden Anzahl an frei erhältlichen Chloroplasten Genomen, welche aus NCBI oder CpBase geladen werden können und gegen Ende 2016 erstmals die 1000 Genome überschritten haben³⁶ können immer mehr Versuche mit vielen Chloroplasten durchgeführt werden. So ist immer noch nicht geklärt wie genau die Replikation von Plastid Genomen wie von Chloroplasten wirklich funktioniert. Wie werden Mutationen im Inverted Repeat repariert oder bei der Replikation auf beide IRs übernommen? Da SNPs im IR immer auf beiden gefunden werden. Welche Mutationen treten am häufigsten auf und wie sind diese evtl. an die Struktur des Genoms gekoppelt³⁷? Auch ist immer noch nicht exakt verstanden wie Chloroplasten vererbt werden, es wird zwar angenommen das diese ähnlich wie Mitochondrien maternal vererbt werden doch gibt es bei Pflanzen auch viele Arten die biparental oder uniparental Chloroplasten vererben³⁸. Die in den letzten Jahren stark steigende Anzahl an Chloroplasten Genomen gibt diesen Fragestellungen immer mehr neue Rohdaten die diese Fragen lösen könnten. Auch Probleme die nur mit kleinen Änderungen im Chloroplasten Genom zu tun haben (SNPs) können so auf den Grund gegangen werden, oder auch die Adaption von verschiedenen Chloroplasten Genen in das Pflanzengenom und der daraus folgenden Änderung im Photosynthese Systems³⁹. Auch kann ohne große Änderung an der kodierenden Sequenz, alleine durch Änderung an Transkriptionsfaktoren oder deren Level viel Einfluss auf solche Systeme genommen werden, welche natürlich auch mit dem Chloroplasten zusammenhängen³⁶. Wie bereits erwähnt eignen sich Chloroplasten gut als Barcode Marker. Auch hier kön-

³²<https://github.com/holmrenser/IOGA>

³³Bakker FT, Lei D, et al. (2015), Herbarium genomics: plastome sequence assembly from a range of herbarium specimens using an Iterative Organelle Genome Assembly pipeline, Biol. J. Linnean Soc. <https://doi.org/10.1111/bij.12642>

³⁴<https://jgi.doe.gov/data-and-tools/bbtools/>

³⁵Luo R, Liu B, Xie Y, et al. (2012), SOAPdenovo2: an empirically improved memory-efficient short-read de novo assembler. GigaScience. 10.1186/2047-217X-1-18.

³⁶Tonti-Filippini J, Nevill PG, Dixon K, et al. (2017), What can we do with 1000 plastid genomes?. Plant J, 10.1111/tpj.13491

³⁷Massouh A, Schubert J, Yaneva-Roder L, et al. (2016), Spontaneous Chloroplast Mutants Mostly Occur by Replication Slippage and Show a Biased Pattern in the Plastome of *Oenothera*. The Plant Cell. 10.1105/tpc.15.00879.

³⁸Greiner S, Sobanski J, Bock R. (2015), Why are most organelle genomes transmitted maternally? Bioessays. 10.1002/bies.201400110.

³⁹Wicke S, Schneeweiss GM, dePamphilis CW, et al. (2011) The evolution of the plastid chromosome in land plants: gene content, gene order, gene function. Plant Molecular Biology. 10.1007/s11103-011-9762-4.

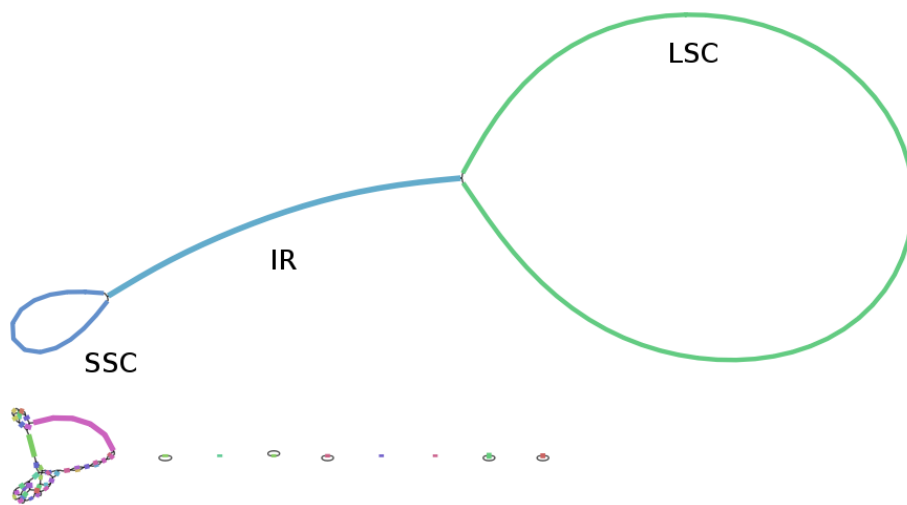


Figure 5: **Bandage - Fastg Visualisierung** Die Visualisierung einer fastg Datei, der eigentlich zirkuläre Chloroplast zeigt sich in einer Form in der SSC (Blau) und LSC (Grün) durch eine Kette welche den IR (Türkis) darstellt verbunden sind. Diese Form wird im fcg.pl Skript des chloroExtractors aufgelöst, wobei beim GetOrganelle diese Struktur per Hand gefunden werden muss.

nen Fortschritte mit mehr Daten erlangt werden. Zudem können mit vielen Chloroplasten Daten sehr gut Phylogenetische Bäume berechnet werden⁴⁰. Dies sind alles Beispiele wie zwischen vielen Spezies mit Hilfe von Chloroplasten Forschung betrieben werden kann. Aber auch innerhalb einer Spezies tauchen Variabilitäten auf, und dies konnte nur mit vielen verschiedenen Chloroplasten der gleichen Spezies herausgefunden werden. So wurden beim 1001 Genom Projekt mehrere Tausend SNPs auf *A.thaliana* Chloroplasten gecalled^{41, 7}. Doch können nicht nur viele Chloroplasten Probleme lösen, schon einzelne neue Chloroplasten können sehr aufschlussreich und informativ sein. So wurde die Idee des chloroExtractors z.B. nur aus dem Grund entworfen einen Chloroplasten aus dem *Dionaea muscipula* (Venusfliegenfalle) Genom zu extrahieren, um diesen separat zu haben, um das Genom leichter zu assemblieren und annotieren. Denn es kann durchaus vorkommen dass bei neuen Genomen, welche de novo assembliert werden müssen, Verunreinigungen durch Chloroplasten auftreten können. Denn ~5 - 20% der kompletten DNA wird von Plastiden-DNA ausgemacht, je nach Spezies und Gewebe³⁶.

2.7 Aufgaben in der Master Thesis

Die Aufgaben dieser Thesis ist grob in drei Teile eingeteilt. Zunächst sollen die verschiedenen Programme, der chloroExtractor^{16, 17}, fast-plast²³, IOGA^{32, 33}, GetOrganelle^{29, 30}, Org.ASM²⁶ und NOVOPlasty^{24, 25} verglichen werden und herausgefunden werden welche das oder die besten Programme sind um damit so viele Chloroplasten Genome zu erzeugen wie möglich. Hier soll vor allem darauf geachtet werden dass die Programme Automatisierbar sind um einen hohen Durchsatz zu haben. Zudem sollen die Programme Ressourcen schonend arbeiten. Der zweite Teil ist das Produzieren von Chloroplasten Genomen, hierzu werden die Pflanzen Genome des 1001 Genom Projektes verwendet. Nach internen Besprechungen und ersten Tests des chloroExtractors, wird angenommen das ca. 10 - 20%⁴² der Datensätze einen kompletten Zirkulären Chloroplasten erbringen könnten. Dies hängt von mehreren Variablen ab. Zunächst wie viel Chloroplasten DNA ist in den Daten vorhanden, dies unterscheidet sich je nachdem welches Gewebe verwendet wurde zum Sequenzieren. Hier haben Wurzeln weniger Chloroplasten als Blüten oder Blätter. Auch hängt es davon ab wie "gut" die Daten sind, generell gilt je größer die Reads desto besser zu assemblieren, auch zu beachten sind insert Size und Anzahl der Reads. Auf den so Produzierten Chloroplasten sollen verschiedene wissenschaftliche Analysen durchgeführt werden, so zum Beispiel eine Varianz Analyse sowie eine Genomweite Assoziationsstudie, kurz GWAS⁴³. Eine GWAS versucht bestimmte Traits, also Eigenschaften mit Genomischen Varianten zu assoziieren, um anschließend eine Aussage darüber treffen zu können ob diese Variante einen Einfluss auf diese Eigen-

⁴⁰Chase MW, Fay MF. (2001) Ancient flowering plants: DNA sequences and angiosperm classification. Genome Biology. <https://doi.org/10.1186/gb-2001-2-4-reviews1012>

⁴¹<http://1001genomes.org/>

⁴²Interne Absprachen mit Förster F. und Ankenbrand M. über geschätzte Ergebnisse

⁴³Korte A, Farlow A. (2013), The advantages and limitations of trait analysis with GWAS: a review. Plant Methods. 10.1186/1746-4811-9-29.

schaft hat oder nicht. Hierzu werden die einzelnen Chromosomen einzeln oder als komplettes Genom angesehen, je nach Ansatz oder Fragestellung. Auch sollte eine Struktur Varianz Analyse durchgeführt werden. Zudem könnten diese Daten benutzt werden um Chloroplasten besser als Genetische Marker zu benutzen. Der dritte Teil ist das Suchen nach bisher noch nicht dokumentierten Chloroplasten Genomen, hierzu sollen Daten verwendet werden welche noch keinen Eintrag in der Chloroplasten Datenbank haben.

3 Material / Methoden

3.1 Evaluation der Programme

Um die oben genannten Programme zu vergleichen habe ich mir verschiedene Ansätze überlegt. Um zunächst zu testen wie genau die Programme funktionieren und ob diese überhaupt funktionieren, wurden sie auf dem Testset SRR5216995 (*Arabidopsis thaliana*: Col-0) mit eine Millionen reads getestet, dieser ist frei zugänglich bei NCBI und dient als Testset beim chloroExtractor¹⁶. Um eine Automatisierung zu erhalten musste für jedes Programm ein Dockercontainer⁴⁴(s. Anhang Tab.9) gebaut werden, falls dieser nicht schon einer vorhanden war, letzteres traf nur für den chloroExtractor zu. Diese Dockercontainer sind auf Dockerhub⁴⁵ frei zur Verfügung⁴⁶. Um das Ziel zu erreichen so viele Chloroplasten wie möglich zu extrahieren, musste eine Automatisierungslösung für alle Programme erstellt werden, damit keine evtl. Manuelle Schritte oder Auswertungen der zeitbestimmende Schritt sind. Um dies zu erreichen mussten zusätzlich einige Bash und Perl Skripte (s. Anhang) geschrieben werden, welche eine volle Automatisierung ermöglichen. Um für alle Programme, welche einen Seed oder eine Referenz benötigen Chancengleichheit herzustellen wurde hier überall die gleiche Datei verwendet. Diese Datei benutzt der chloroExtractor, hier handelt es sich um 4154 Chloroplasten Gene, diese sind aus verschiedenen Arten. Bei den Genen handelt es sich um *ndhB*, *rps12*, *psbD*, *rpl2* und *psbA*. Liefert ein Programm seine eigenen Referenzen mit, wie es der chloroExtractor tut, wurden diese nicht geändert. Da diese als Standardparameter gelten, welche nicht verändert werden sollten.

3.1.1 Testdaten

Es wurden verschiedene Größe von Dateien verwendet. So sind dies alles Illumina short Read Daten, doch unterscheiden sich diese in Readlänge, Insertsize und Anzahl der Reads.

Simulierte Daten Um zu Testen wie gut die verschiedenen Programme mit unterschiedlichen Anteilen von Chloroplasten DNA in Genom Daten zurechtkommen, wurden drei verschiedene Testdatensätze simuliert (Genom : Chloroplast - 1:10, 1:100, 1:1000). Mit diesen sollte auch getestet werden ob die Programme mit viel oder wenig Chloroplasten DNA Anteil zurecht kommen oder einen dieser Fälle bevorzugen. Um diese

⁴⁴<https://www.docker.com/>

⁴⁵<https://hub.docker.com/>

⁴⁶<https://hub.docker.com/u/chloroextractorteam/>

Daten vorzubereiten wurden von *Arabidopsis thaliana* (TAIR10 ⁴⁷) die jeweiligen Chromosomen wie auch die Daten des Chloroplasten von NCBI heruntergeladen. Diese wurden in den jeweiligen Verhältnissen zusammen kopiert. Diese Testdatensätze wurden mit ART^{48, 49}(Version: 2.5.8) erzeugt. ART wird dazu verwendet Short-reads zu erzeugen, ART kann keine zirkulären Daten wie Chloroplasten erzeugen, deswegen wurden diese als lineare Sequenzen verwendet mit der Abfolge LSC-IRB-SSC-IRA, mit einem Overlap zwischen IRA und LSC: (IRA)-LSC-IRB-SSC-IRA. Mitochondrien DNA wurde nicht mit simuliert, da sich zunächst auf die Chloroplasten fokussiert werden sollte und diese Mitochondrien die Komplexität der Daten erhöht. Um die verschiedenen Verhältnisse von Genom und Chloroplasten zu bekommen wurden die Chloroplasten Daten einfach vervielfältigt und anschließend zusammen kopiert. Hiernach wurden sie mit folgenden ART Kommandos als short Reads simuliert. Hiernach wurden die Daten gemischt, da es zu Problemen kommt wenn diese Daten sortiert sind. Für diese Daten wurden 150 Basen paare Reads simuliert, und eine Coverage der Daten welche 100x beträgt. Für die Tests wurden eine Millionen Reads pro Datei benutzt, da diese genug Chloroplasten DNA enthalten sollten.

```
'art_illumina [options] -i <INPUT_SEQ_FILE> -l <READ_LEN> -f <FOLD_COVERAGE>
-o <OUTPUT_FILE_PREFIX> -m <MEAN_FRAG_LEN> -s <STD_DE>'
```

```
'1:10 : ./art_illumina -p -i sequence-arabidopsis-thaliana-kern-chl-1zu10.fa -l 150 -f 100
-o a_thaliana_1_10_sim -m 500 -s 150'
```

```
'1:100 : ./art_illumina -p -i sequence-arabidopsis-thaliana-kern-chl-1zu100.fa -l 150 -f
100 -o a_thaliana_1_100_sim -m 500 -s 150'
```

```
'1:1000 : ./art_illumina -p -i sequence-arabidopsis-thaliana-kern-chl-1zu1000.fa -l 150
-f 100 -o a_thaliana_1_1000_sim -m 500 -s 150'
```

1001 Genom Projekt Um einen ersten Eindruck über die Programme und deren Erfolgs rate zu bekommen wurden parallel zu den Tests mit simulierten Daten, die ersten Tests mit realen Datensätzen vorgenommen. Hierzu wurden Daten aus dem 1001 Genom Projekt⁴¹ verwendet, dies sind alles Daten von *Arabidopsis thaliana*. Es wurden 11 Datensätze (SRR1945435 - SRR1945445) verwendet. Diese sind alle frei verfügbar und wurden von NCBI heruntergeladen. Es wurden jeweils zwei Millionen Reads pro Datei gezogen, mit 150 Basenpaaren pro Read.

GetOrganelle-Paper preprint Um zu weitere Testdaten zu ermitteln und ein Urteil darüber zu fällen welche Programme weiter verwendet werden, wurden 57 Datensätze welche im GetOrganelle Paper ³⁰ verwendet wurden auf allen Programmen getestet.

⁴⁷https://www.ncbi.nlm.nih.gov/assembly/GCF_000001735.3/

⁴⁸Weichun H, Leping L, Jason RM, Gabor TM. (2015), ART: a next-generation sequencing read simulator, Bioinformatics, <https://doi.org/10.1093/bioinformatics/btr708>

⁴⁹<https://www.niehs.nih.gov/research/resources/software/biostatistics/art/index.cfm>

In dieser Arbeit wurden bei 47 Datensätzen von 57, mit dem GetOrganelle erfolgreich zirkuläre Chloroplasten extrahiert. Diese Daten sind auch frei zugänglich und wurden von NCBI heruntergeladen. Gerade hier gab es einige Abweichungen in Dateigrößen. Reads reichten von 75 Basenpaaren bis zu 300 Basenpaare pro Read. Es wurden hier fünf Millionen Reads pro Datei verwendet, da diese im GetOrganelle Paper auch verwendet wurden.

3.1.2 Welche Programme werden weiter verwendet.

Um alle Daten aus dem 1001 Genom Projekt (1135 Datensätze) zu berechnen, mussten aufgrund von Hardware technischen Limitierungen die besten Programme ausgewählt werden. Diese Programme müssen in in Geschwindigkeit sowie in Erfolgs- und Fehlerrate überzeugen. Des weiteren müssen diese Programme gut automatisierbar sein, d.h. am besten mit nur Befehl gestartet werden können, sodass kein weiterer Aufwand anfällt. Dies gilt vor allem auch bei der Wahl der Parameter mit denen das Programm gestartet wird. Diese können nicht für jeden Datensatz angepasst werden, was bedeutet dass die Standardparameter verwendet werden. Dies ist notwendig um einen hohen Durchsatz an Berechnungen zu ermöglichen.

Installation & Automatisierung Alle Programme konnten mit Hilfe von einigen Skripts und dem erstellen eines Dockercontainers so automatisiert werden, dass sie einen hohen Durchsatz erreichen konnten. Das Einzige Programm welches einen Händischen Schritt benötigt ist der GetOrganelle, hier muss die fastg Datei in Bandage geöffnet werden und der zirkuläre Chloroplast selbst heraus gesucht werden. Bei den verschiedenen Skripts handelt es sich vor allem um Start-Skripts. Aber es mussten auch ein paar kleine Skripts verwendet werden um kleine Bugs zu fixen. So kann der IOGA keine Unterordner verwenden da er sonst versucht auf Falsche Dateien zuzugreifen und abstürzt. Dies scheint ein Bug in einem Splitt Befehl zu sein. Beim GetOrganelle mussten zusätzliche Befehle eingebaut werden damit SPAdes keine Fehlermeldungen bringt und abbricht, da er bestimmte Funktionen (hammer.py) nicht ausführen konnte welche für eine Fehler Korrektur verwendet werden, welche GetOrganelle gar nicht nutzt. Org.ASM konnte nur erfolgreich in einem Dockercontainer installiert werden, da dieses Programm sonst verschiedenste Fehlermeldungen brachte. Alle Programme welche Perl verwenden, also chloroExtractor, fast-plast und NOVOPlasty, brachten Fehlermeldungen, da innerhalb des Dockercontainers Globale Variablen nicht vollständig gesetzt waren. Diese Fehler waren aber nicht fatal, und konnten mit dem setzten dieser Variable leicht entfernt werden. Für jedes Programm wurde ein Skript geschrieben welches die Laufzeit überprüft und wenn dieses fertig ist, eine Auswertung startet.

Erfolgsrate Um zunächst zu überprüfen ob ein wirklich ein kompletter Chloroplast zusammengebaut wurden. Wurde bei den ersten Testdatensätzen ein Referenz Mapping auf TAIR10 benutzt. Hierzu wurde mit Bowtie2, später mit minimap2⁵⁰(Version: 2.10-

⁵⁰Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics. 10.1093/bioinformatics/bty191

r761) der Chloroplast auf das TAIR10 Chloroplasten Genom gemapt. Auch wurde mit AliTV⁵¹ eine Visualisierung des Mappings erstellt. Nachdem klar war das es sich bei allen ausgegeben Daten um Chloroplasten handelt, und weil diese Art der Auswertung schlecht Automatisierbar war, wurde ein Bash Skript geschrieben welche die Auswertung übernimmt. Dieses Skript überprüft die Größe des Chloroplasten und in wie vielen Contigs der Chloroplast ausgegeben wurde. Hierzu wurde das SeqFilter⁵²(Version: 2.1.8) Skript verwendet, und anschließend über ein Bash Skript eine Entscheidung getroffen ob es sich um einen kompletten Chloroplasten handelt oder nicht (s. Anhang: Skripte:ev_stat.sh). Hierzu wurden Verschiedene Kategorien eingeführt(s. Tab. 1). Diese Auswertung wurde für dann für alle Testdaten sowie die GetOrganelle Preprint Daten verwendet. Zudem gab es zu den vier Kategorien, in denen Erfolge eingeteilt wurden noch Error, wenn kein Ergebnis vorhanden war wenn das Programm z.B. durch einen Fehler abgebrochen hat. Sowie Cancelled wenn das Programm länger als 14 Tage brauchte wurde es abgebrochen.

Table 1: **Erfolgsraten Einteilung** Das Skript ev_stat.sh scannt die Output Dateien und teilt diese je nach Größe und Anzahl der Contigs in verschiedene Kategorien ein.

Kategorie	Contigs	Basenpaare
Success	1	110 kbp - 180 kbp
Partial	> 1	110 kbp - 180 kbp
Incomp:high	> 1	> 180 kbp
Incom:low	> 1	> 110 kbp

Geschwindigkeit Einer der weniger entscheidenden aber dennoch wichtigen Punkte nach dem gefiltert wurde ist die Geschwindigkeit, oder besser die Laufzeit der Programme. Zunächst wurde hier die Durchschnittszeit genommen, die der Prozess zum rechnen benötigt, anschließend wurde mit dem time Linux Kommando die CPU als auch die Real-zeit gemessen. Die Geschwindigkeit von Programmen mit vielen Abhängigkeiten brauchen im Schnitt länger, da zum benutzen der Dockercontainer Singularity⁵³(Version: 2.4.5-dist) verwendet wurde. Dieses Benötigt Zeit um den Container zu verwenden, zudem wird Zeit in Anspruch genommen wenn viele Daten in den Container gemountet werden müssen.

⁵¹Ankenbrand MJ, Hohlfield S, Förster F, et al. (2017) AliTV—interactive visualization of whole genome comparisons. PeerJ Computer Science, <https://doi.org/10.7717/peerj-cs.116>

⁵²<https://github.com/BioInf-Wuerzburg/SeqFilter>

⁵³<https://singularity.lbl.gov/>

Benötigte Ressourcen Ein weiterer Punkt nachdem aussortiert wurde ist der benötigte RAM Verbrauch. Es wurden verschiedene Größen von Dateien verwendet um in Erfahrung zu bringen wie sich dies auf Ressourcen und Laufzeit auswirkt. Zudem wurde zum Ausführen der Dockercontainer Singularity⁵³ verwendet, welches die benötigte Laufzeit und die benötigten Ressourcen, wie RAM beeinflusst.

3.2 Varianz Analyse

Um mehr über die Chloroplasten und deren Verbreitung, sowie Mutationsrate und somit Varianz zu erfahren wurden zwei verschiedene Varianzanalysen durchgeführt. Zunächst sollte überprüft werden welche Einflüsse die Programme und ihre Strategien den Chloroplasten zu assemblieren, speziell deren Assambler, auf die Varianz der entstehenden Chloroplasten hat. Hierzu wurden die assemblierten Chloroplasten, welche beide verwendeten Programme gemeinsam hatten verwendet. Diese Läufe wurden zunächst zehn mal wiederholt, auch um einen Eindruck über die Reproduzierbarkeit der Ergebnisse zu bekommen. Diese Chloroplasten wurden anschließend mit minimap2⁵⁰ auf das Referenzgenom (TAIR10 Chloroplast⁵⁴) gemapt. Hiernach wurde eine Varianzanalyse mit Samtools⁵⁵ (Version: 0.1.19-96b5f2294a) durchgeführt, hierzu wurde der Befehl 'mpileup/bcftools call'⁵⁶ (bcftools Versionen: 0.1.19-96b5f2294a & 1.8) verwendet. Dieser führt eine Varianzanalyse bzw. ein SNP calling durch. Die zweite Varianzanalyse wurde auf allen Chloroplasten welche aus dem 1001 Genom Projekt gebaut wurden erstellt. Auch diese wurden auf den Referenz Chloroplasten mit minimap2 gemapt und anschließend mit Samtools' 'mpileup' Funktion einem SNP calling unterzogen.

3.3 GWAS

Häufig wird eine GWAS über das komplette Genom berechnet. Doch können auch einzelne Chromosomen oder Organellen bereits signifikante Varianten besitzen. So soll mit dieser GWAS der Einfluss von Chloroplasten Varianten auf Eigenschaften der *A.thaliana* getestet werden. Hierzu wurden die SNP callings aus der Varianzanalyse verwendet. Verschiedene Trait-Tabellen wurden von Arapheno^{57, 58}, einer Trait Datenbank für *A.thaliana*, heruntergeladen und zusammen mit den Varianzanalyse Daten in ein R⁵⁹ Skript gegeben. Dieses R Skript nutzt zunächst vcfr⁶⁰, ein R Paket, um die verschiedenen VCF (Variance Calling File) Daten einzulesen. Anschließend ruft es ein weiteres R

⁵⁴https://www.ncbi.nlm.nih.gov/nuccore/NC_000932.1

⁵⁵Li H, Handsaker B, Wysoker A, et al. (2009) The Sequence alignment/map (SAM) format and SAMtools, Bioinformatics 10.1093/bioinformatics/btp352

⁵⁶Li H, (2011) A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data, Bioinformatics 10.1093/bioinformatics/btr509

⁵⁷<https://arapheno.1001genomes.org/>

⁵⁸Seren Ü, Grimm D, Fitz J, et al. (2017) AraPheno: a public database for Arabidopsis thaliana phenotypes, Nucleic Acids Research, <https://doi.org/10.1093/nar/gkw986>

⁵⁹<https://www.r-project.org/>

⁶⁰<https://cran.r-project.org/web/packages/vcfr/index.html>

Skript auf welches freundlicher weiße von Korte et. al⁴³ zur Verfügung gestellt wurde und eine GWAS Analyse durchführt.

3.4 Struktur Varianz Analyse

Wie bereits erwähnt können Chloroplasten auch verschiedene Strukturelle Änderungen evolvieren. Diese sind durch die Rohdaten, welche meist short Reads sind, nicht aufzudecken. Da diese zu kurz sind um komplette Struktur Varianten zu überspannen.⁷ Hierzu könnten nun die komplett de novo Assemblierten Chloroplasten verwendet werden. Es wurde versucht mit Delly⁶¹ (Version: v0.7.8) und Breakdancer⁶² (Version: 1.3.6) Struktur Varianten in Chloroplasten zu finden.

3.5 Neue Chloroplasten

Um neue Chloroplasten von Spezies zu finden, welche noch nicht in der CpBase^{63, 64} Datenbank sind, wurde eine Liste von Möglichen Daten von NCBI mit CpBase verglichen. Nur 49 Datensätze waren ohne Eintrag in CpBase und hatten somit noch keinen Dokumentierten Chloroplasten für diese Spezies. Auf diese 49 Datensätze wurden sowohl der chloroExtractor als auch der fast-plast angewendet. Um die NCBI liste von Interessanten Daten zu erhalten wurde mit folgendem Befehl gesucht:

```
((((((("green plants"[orgn]) AND "wgs"[Strategy]) AND "illumina"[Platform]) AND "biomol dna"[Properties]) AND "paired"[Layout]) AND "random"[Selection])) AND "public"[Access]
```

Mit einem Skript (s. Anhang, cpbase.sh) wurden alle Spezies Einträge von CpBase geladen welche einen Chloroplasten besitzen. Anschließend wurde mit einem folgendem Perl-Einzeiler die Datensätze herausgegeben welche noch keinen Eintrag in CpBase haben. Zudem musste der Datensatz mindestens zwei Millionen Reads haben und mindestens 200 Basenpaare pro Read aufweisen.

```
perl -F", " -ane 'print if F[6] > 399andF[3]>999999' SraRunInfo_plants.csv | grep -vf species_cpbase.list | sort -u -t, -k29,29 | shuf
```

4 Ergebnisse

4.1 Automatisierung

Um eine Automatisierung aller Programme zu erreichen wurde für jedes Programm ein Dockercontainer gebaut welcher mit Singularity verwendet wird. Zudem wird die komplette Auswertung von einigen Skripts übernommen. Um dies zu Bewerkstelligen wurden mehrere Skripts geschrieben welche sich gegenseitig aufrufen um den kompletten Ablauf sicherzustellen (Fig. 6). Das einzige Skript welches aktiv ausgeführt werden

⁶¹<https://github.com/dellytools/delly>

⁶²<https://github.com/genome/breakdancer>

⁶³http://rocaplab.ocean.washington.edu/old_website/tools/cpbase

⁶⁴http://rocaplab.ocean.washington.edu/tools/cpbase_test/

muss ist das `run_SRRchl.sh`. Dieses Skript setzt Links zu anderen Skripts, zum einen zwei Auswertungs Skripts (`ev_stat.sh` und `percent_stat.sh`) und zu einem Skript namens `cp_skript.sh`. Dieses `cp_skript` übernimmt den kompletten Aufbau der Ordner Struktur, und linkt all die Skripts die jedes Programm braucht, so brauchen IOGA und GetOrganelle eine Referenz, diese wird von diesem Skript in die passenden Ordner kopiert. Auch kopiert und führt dieses `cp_skript.sh` das Skript aus welches die NOVOPlasty Konfigurationsdatei automatisiert für jeden Datensatz schreibt (`make_NP_config.pl`). Für jeden Datensatz wird so ein Ordner erzeugt mit jeweils dem Programm als Unterordner. In jedem Unterordner werden die roh Daten verlinkt, sowie für jedes Programm das passende Evaluierungsskript und Runskript. Als letztes linkt es `sbatch_run_all.sh` und `ev_all.sh` in den jeweiligen Datenordner. Diese werden nun vom `run_SRRchl.sh` Skript ausgeführt. Das `sbatch_run_all.sh` Skript geht nun in jeden Unterordner und startet die jeweiligen Programme über `sbatch` und deren Runskript. Zudem startet es auch die dazugehörigen Evaluierungsskripts, welche auch gleichzeitig als Überwachungsskript dienen. Sobald der Slurm Job fertig ist, startet das Evaluierungsskript des jeweiligen Programms damit, die Finale Output Datei zu überprüfen und diese in eine der vier Erfolgs Kategorien einzuteilen. Zudem schiebt es alle Dateien welche keine Log Dateien oder Finale Output Dateien sind in einen `raw_Programm` Ordner, damit dieser mit dem `clear_skript.sh` gelöscht werden kann, falls diese Daten nicht mehr benötigt werden. Sobald alle Datensätze fertig sind, wird mit dem `ev_stat.sh` Skript eine Datei mit einer Erfolgstabelle mit jedem Programm erstellt. `Percent_stat.sh` kann dann genutzt werden um eine Zusammenfassung über alle Datensätze zu erhalten.

4.2 Daten: Simulierte Daten

Die Simulierten Daten, welche mit ART^{48, 49} erzeugt wurden um das Verhalten der Programme bei verschiedenen Verhältnissen zu testen, konnten von drei Programmen, dem `chloroExtractor`, `fast-plast` und `Org.ASM` bei allen drei Datensätzen geschafft werden. Diese schaffen es einen vollständigen zirkulären Chloroplasten zu bauen. NOVOPlasty baut zwar auch einen kompletten Chloroplasten doch gibt dieser nur die drei verschiedenen Contigs aus (IR, SSC, LSC), und schafft es nicht diese in einen zirkulären Chloroplasten zu vereinen. GetOrganelle wie auch der IOGA schaffen es nicht die simulierten Datensätze zusammen zu bauen da sie mit einem Fehler abbrechen oder wie im Falle des IOGA nach zwei Wochen Laufzeit abgebrochen werden. (s. Tabelle 2)

4.3 Daten: 1001 Genom Projekt, 11 Testdatensätze

Aus den Daten des 1001 Genom Projekts^{41, 7} wurden zunächst elf Testdatensätze verwendet um auch reale Daten auf allen Programmen zu testen. Von den elf Testdatensätzen des 1001 Genom Projekts konnten sechs verschiedene vollständige zirkuläre Chloroplasten zusammengebaut werden. Von diesen sechs bringt der `fast-plast` fünf ein und der `chloroExtractor` einen. Keines der anderen Programme konnte einen weiteren zirkulären Chloroplasten erzeugen (s. Tab.3). Diese elf Datensätze des GetOrganelles wurden per Hand ausgewertet, keiner dieser elf Datensätze konnte einwandfrei mit Bandage zu einem

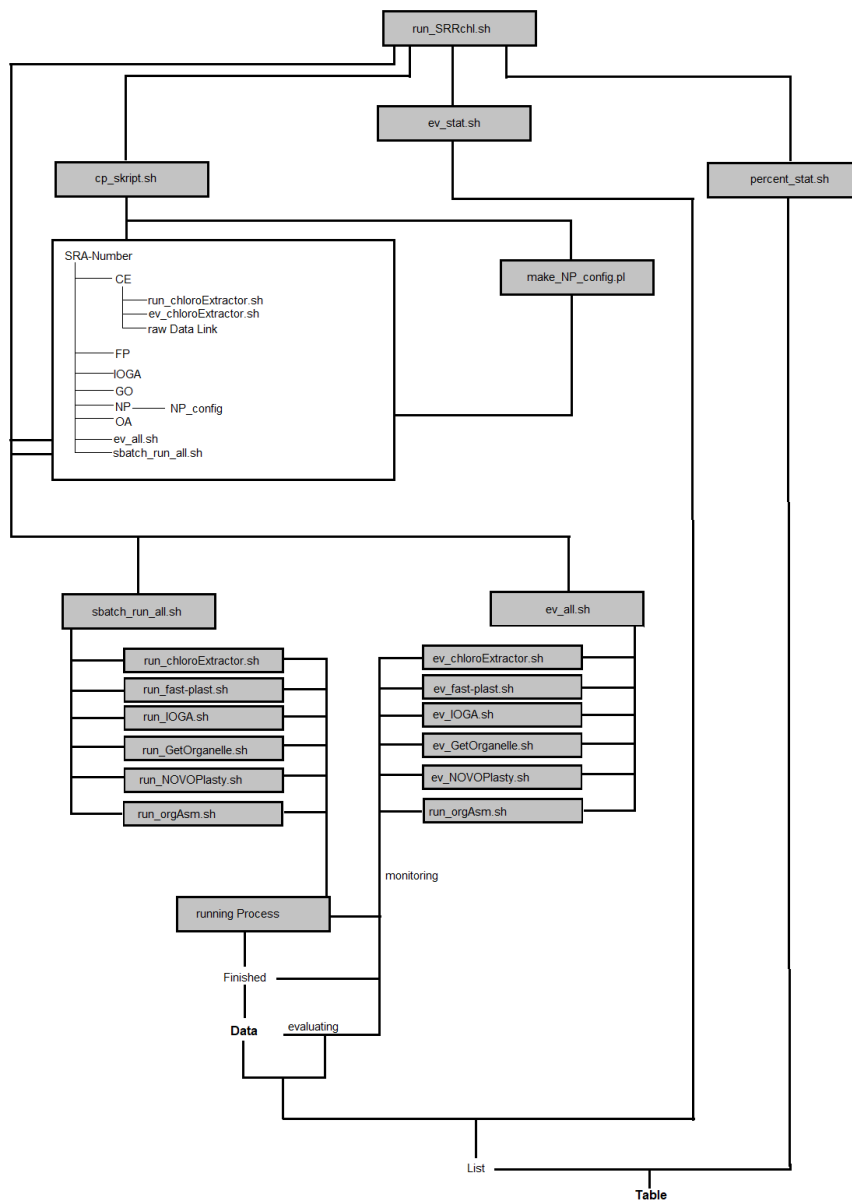


Figure 6: **Automatisierungsskripts** Ablauf der verwendeten Skripte um eine Automatisierung zu erwirken, hier wird nur das run_SRRchl.sh Skript ausgeführt und alle anderen Skripte werden automatisch bis zur Auswertung ausgeführt. So erstellt das cp_sript.sh für jedes Programm einen Unterordner mit den dazugehörigen run und ev Skripts. sbatch_run_all.sh startet alle run Skripts und ev_all.sh startet jedes ev Skript.

Table 2: **Test Datensatz: Simulierte Daten** S steht für Success, E für Error, C für Cancelled die angegebene Zahl steht für die Anzahl der Contigs. Bis auf IOGA und GetOrganelle konnten alle anderen Programme die Simulierten Daten zu einem Chloroplasten zusammenbauen, auch wenn im Falle des NOVOPlasty nicht zirkulär. Die IOGA Läufe mit "C" wurden nach zwei Woche Laufzeit abgebrochen.

Sim(Genome:Chloroplast)	CE	FP	NP	GO	OA	IOGA
1:10	S	S	P-3	E	S	E
1:100	S	S	P-3	E	S	C
1:1000	S	S	P-3	E	S	C

zirkulären Chloroplasten gebaut werden, da immer kein Ringschluss vorhanden war. (vergleiche Fig.7, Fig. 8)

Table 3: **Test Datensatz: 1001 Genom Project** S steht für Success, E für Error, I für Incomplete, C für Cancelled. Sechs verschiedene Chloroplasten konnten zu einem zirkulären Chloroplasten zusammengebaut werden, dabei werden bereits fünf vom fast-plast abgedeckt und einer wird von chloroExtractor beigesteuert.

SRA	CE	FP	NP	GO	OA	IOGA
SRR1945435	I	I	I	I	E	I
SRR1945436	I	S	I	I	I	I
SRR1945437	I	I	I	I	I	I
SRR1945438	P	S	I	I	E	I
SRR1945439	I	S	I	I	I	I
SRR1945440	I	S	E	I	E	I
SRR1945441	I	S	E	I	I	I
SRR1945442	I	I	I	I	C	C
SRR1945443	S	I	I	I	I	I
SRR1945444	I	I	E	I	I	I
SRR1945445	I	I	E	I	E	I

4.4 Daten: GO-Preprint

Um mehr Daten zu testen, wurden alle 57 Datensätze des GetOrganelle Papers ³⁰ benutzt. Da der GetOrganelle diese Daten eigentlich erfolgreich schaffen sollte wurde hier versucht mit dem fcg.pl Skript des chloroExtractors eine Automatisierung der Daten zu

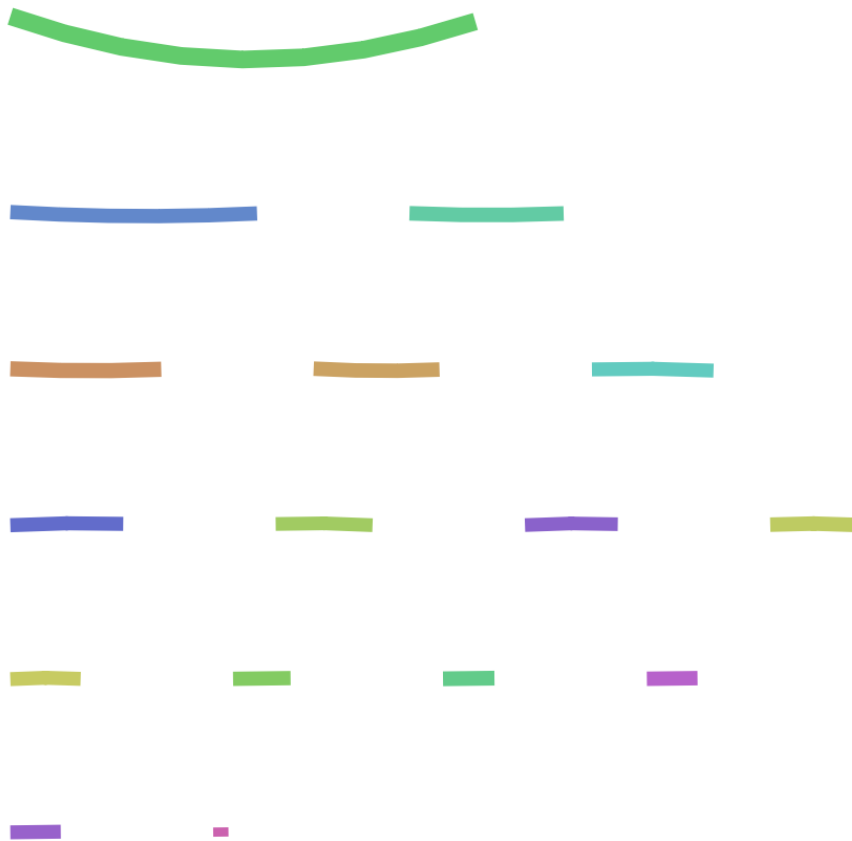


Figure 7: **Bandage: GetOrganelle SRR1945443** Kein Chloroplast erkennbar im fastg Graphen. Keine vernetzung der Contigs per se erkennbar. Dazu im vergleich der graphen des chloroExtractors (Fig. 8)

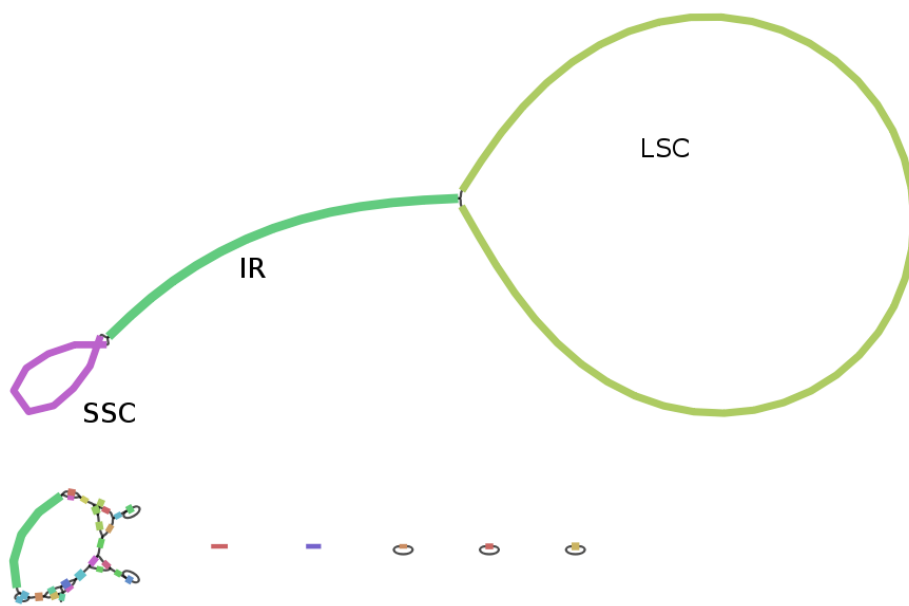


Figure 8: **Bandage: chloroExtractor SRR1945443** Chloroplast klar erkennbar, dieser zeichnet sich aus durch einen Großen Kreis (LSC - Gelb), verbunden über eine Linie (IR - Grün) auf einen Kleinen Kreis (SSC - Violett).

erwirken. Doch versucht der GetOrganelle zunächst die die fastg-Graphen zu verbessern, dies führt dazu dass das fcg.pl Skript nicht mehr funktioniert. So wurden die fastg-Graphen aus SPAdes direkt verwendet, doch ergaben sich hier leider nur zwei Datensätze als zirkuläre Chloroplasten. Auf Nachfrage beim GetOrganelle Team, hieß es dass, wenn es notwendig war alle Parameter angepasst wurden und dass alle Chloroplasten per Hand aus Bandage geholt wurden. Von 57 Datensätzen, welche im GetOrganelle Paper verwendet wurden, konnten insgesamt 40 fertig gestellt werden, diese verteilen sich auf die Programme (s. Tab. 4). So konnten 35 von 40 Erfolgreichen Datensätzen durch den fast-plast und den chloroExtractor erreicht werden. Der fast-plast schafft es 31 Chloroplasten Genome komplett zu bauen, davon sind 17 nur von ihm geschafft worden.(Fig. 9) Der chloroExtractor schafft 14 Chloroplasten, wo von drei nur von ihm geschafft werden. Die zwei Erfolge des GetOrganelles, mit Hilfe des fcg.pl Skripts des chloroExtractors werden auch nur vom GetOrganelle geschafft. Von den sieben des NOVOPlasty ist einer dabei welcher nur von diesem geschafft wird. Von den elf des Org.ASM ist auch einer nur von diesem Geschafft worden. Doch wurden einige auch mehrfach geschafft. So sind vier Stück nur von fast-plast und chloroExtractor geschafft worden. Zwei Stück von fast-plast und Org.ASM, sowie jeweils einer von Org.ASM und NOVOPlasty; chloroExtractor und Org.ASM; fast-plast und NOVOPlasty. Drei Erfolge teilen sich fast-plast, chloroExtractor und Org.ASM. Sowie jeweils ein Erfolg teilen sich fast-plast, Org.ASM, NOVOPlasty, als auch einer von fast-plast, chloroExtractor und NOVOPlasty. Zwei Chloroplasten konnten von allen Programmen bis auf GetOrganelle und IOGA gelöst werden (u.a. Fig. 10: SRR5602602 - Laurus nobilis), wobei letzteres Programm nicht einen Erfolg hat(Fig. 9 & Tab. 4). So können bereits 35 von 40 Chloroplasten alleine durch fast-plast und chloroExtractor geschafft werden.

Table 4: **Test Datensatz: GetOrganelle Preprint** 40 von 57 Datensätze konnten komplett gelöst werden. 31 Datensätze konnten mit dem fast-plast zu einem Chloroplasten gebaut werden, die 14 die der chloroExtractor schafft enthalten weitere welche nicht vom fast-plast geschafft wurden. Somit konnten mit allen Programmen 74% gelöst werden, wenn die drei ohne paired Daten herausgerechnet werden und alleine 64% von diesen 54 von fast-plast und chloroExtractor

Tool	SUCCESS	%	ERROR	PARTIAL	INCOMPI	NO_PAIR	Total
CE	14	~26%	11	17	12	3	
FP	31	~57%	0	18	5	3	
GO	2	~4%	21	26	5	3	
IOGA	0	~0%	22	28	4	3	
NP	7	~13%	19	8	20	3	
OA	11	~20%	36	4	3	3	
Summary	40	~74%	-	-	-	3	57

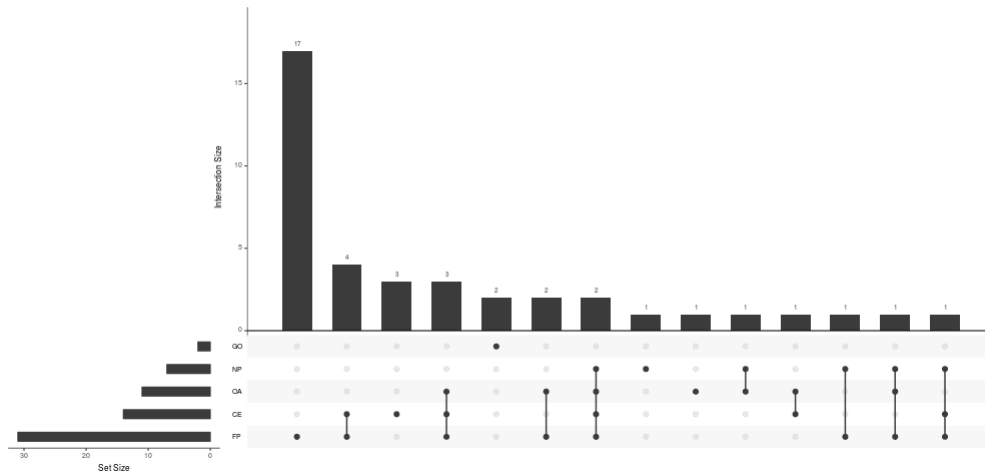


Figure 9: **Upset Diagram GO-Preprint** Hier wird gezeigt wie sich die einzelnen Erfolge auf die 40 Stück aufteilen. So schafft der fast-plast 17 Chloroplasten welches kein anderes Tool schafft. Fast-plast und chloroExtractor haben vier Erfolge gemeinsam, der chloroExtractor schafft drei Chloroplasten welche kein anderes Programm schafft. usw. So werden 35 der geschafften 40 durch fast-plast und chloroExtractor abgedeckt.

4.5 Die besten Programme: fast-plast und chloroExtractor

Da aus zeitlichen und hardwaretechnischen Gründen nicht alle Programme weiterverwendet werden konnten, wurde nach Erfolgsrate, Geschwindigkeit und benötigten Ressourcen (s. Tab 5) gefiltert, am wichtigsten war aber die Automatisierbarkeit der Programme. Bis auf der GetOrganelle konnte für jedes Programm eine Automatisierbarkeit erwirkt werden ohne Daten außen vor zu lassen. Der GetOrganelle benötigt das öffnen der fastg Datei in einem Visualisierungsprogramm für fastg-Graphen, hier wird Bandage empfohlen. Bandage hat allerdings eine schlechte Kommandozeilen Anbindung wodurch auch keine Automatisierbarkeit durch Skripts erfolgen konnte. Es wurde auch versucht mit dem fcg.pl Skript aus dem chloroExtractor, welches genau diesen Schritt im chloroExtractor automatisiert, zu verwenden um auch beim GetOrganelle eine Automatisierbarkeit zu erreichen. Doch führte dies nur bei sehr wenigen Daten zum Erfolg, da der GetOrganelle die von SPAdes erstellte fastg Datei versucht zu verbessern, und die getrimmte Datei nicht mehr vom fcg.pl Skript verwendet werden kann. Dies passiert wohl weil der GetOrganelle beim verbesserten fastg Graphen versucht Namen und Sequenzen anzupassen, womit das fcg.pl Skript nicht zurecht kommt. Es wurde auch versucht die roh fastg Dateien des GetOrganelle zu benutzen dies ergab zwar eine Automatisierbarkeit, doch würden so Teile des GetOrganelles, nämlich das verbessern der fastg Datei unterschlagen. Die Laufzeiten der Programme unterscheiden sich sehr, von 30 Minuten bis über eine Stunde, auch die RAM werte sind sehr unterschiedlich, diese reichten von wenigen 20 Gigabyte bis zu 60 Gigabyte. All diese Werte sind Durchschnittswerte, da ver-

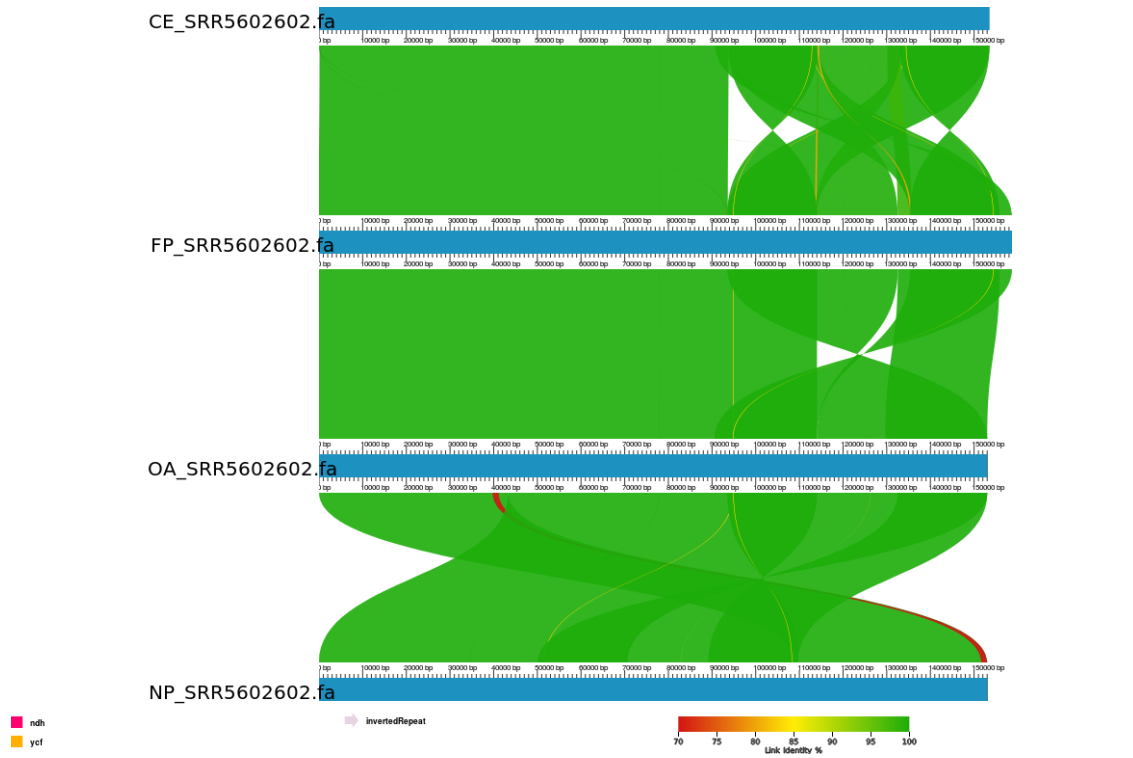


Figure 10: **AliTV SRR5602602 - *Laurus nobilis*** SRR5602602 - *Laurus nobilis*, SRR5602602 - *Laurus nobilis*, Daten des GO-Preprints, geschafft von ChloroExtractor, NOVOPlasty, fast-plast und Org.ASM. Orientierung von LSC sowie von SSC und IRs können nicht perfekt aufgelöst werden und können durchaus Verdreht sein.

schiedene Größen von Dateien als Eingabe verwendet wurden. Da nicht alle Dateien die gleiche Anzahl an Reads hatten, sowie die Größen der einzelnen Reads sich unterschieden. Diese reichten von 75 Basen paare bis zu 300 Basen paare, Anzahl der Reads und somit Größe der Dateien reichten von eine Millionen Reads bis zu fünf Millionen Reads. Die Laufzeiten sind, vor allem bei Programmen mit vielen Abhängigkeiten, erhöht. Da zum nutzen der Dockercontainer Singularity ⁵³ verwendet wurde. Die Programme welche in

Table 5: **Laufzeit und Ressourcenverbrauch** Alle Laufzeiten sind Durchschnittswerte (getrimmtes Mittel), RAM werte zu Peakzeiten. Die Laufzeiten reichen von 30 Minuten (chloroExtractor) bis zu 100 Minuten (IOGA), die RAM Nutzung unterschied sich auch erheblich, diese reichen von 20 GB (chloroExtractor) bis hin zu 60 GB (fast-plast). Aufgrund der Nutzung von verschiedenen großen Datensätzen können nur Durchschnittswerte Angegeben werden.

Tool	Laufzeit	RAM
CE	~ 30 min	~ 20 GB
FP	~ 60 min	~ 60 GB
GO	~ 40 min	~ 50 GB
IOGA	~ 100 min	~ 40 GB
NP	~ 30 min	~ 30 GB
OA	~ 60 min	~ 30 GB

oben genannten Punkte überzeugt haben sind der fast-plast und der chloroExtractor. Der fast-plast benötigt zwar die meisten Ressourcen und ist nicht der schnellste, aber hat mit Abstand die größte Erfolgschance. Zudem ist er voll automatisierbar und erreicht dies mit den vorgegebenen Standard Parametern. Als zweites Programm wird der chloroExtractor verwendet, dieser ist schnell, Ressourcen arm und hat nach dem fast-plast die zweit höchste Erfolgsrate. Mit beiden Programmen konnten 35 von den 40 Erfolgen von 57 Chloroplasten der GetOrganelle-Preprint Daten berechnet werden. Zudem haben diese beiden Programme die wenigsten Probleme bei der Handhabung wie auch bei der Installation zu beginn gemacht. Sie sind durch die gegebenen Parameter einfach zu verwenden und zu Automatisieren. Die von den Programmen geschriebenen Log Dateien sind einfach gehalten um dem Ablauf zu folgen und klar verständlich, der fast-plast gibt sogar drei dieser Dateien aus, da er unterscheidet zwischen Warn- und Fehlermeldungen sowie Standard Meldungen, und eine Datei für den Output der eingebundenen Programme. Der chloroExtractor gibt seine Kompletten Meldungen über ein übergeordnetes Programm aus, welche den Ablauf steuert (PipeWrap). Dieses Programm gibt alles auf STDERROR aus und kann damit einfach mit geloggt werden. In diesem Fall wurde über die slurm Datei, welche von dem verwendeten queueing System ausgegeben wird, mit geloggt. Diese beiden Programme wurden auf allen Daten des 1001 Genom Projekts laufen gelassen, um möglichst viele Chloroplasten zu generieren.

4.6 1001 Genom Projekt

Ziel so viele Chloroplasten wie möglich vollautomatisch aus kompletten Genom Datensätze zu erzeugen, wofür zwei Programme ausgewählt worden sind, wurde zunächst auf Datensätzen des 1001 Genom Projekt versucht. Von den 1135 Datensätzen welche im 1001 Genom Projekt gesammelt wurden, konnten 946 Datensätze erfolgreich von NCBI heruntergeladen werden. Die restlichen 189 konnten nicht richtig heruntergeladen werden aufgrund von Downloadfehlern. Hier handelte es sich um andauerndes Problem, da mehrere male versucht wurde diese Datensätze herunter zu laden. Zudem waren 47 Datensätze keine paired end Datensätze, und konnten deshalb nicht verwendet werden. Von diesen 899 restlichen Datensätzen konnten mit dem fast-plast und dem chloroExtractor 303 komplette zirkuläre Chloroplasten vollautomatisch gebaut werden, dies entspricht etwa 34%. (Tab. 6).

Table 6: **Datensatz: 1001 Genom Projekt** SUCCESS, echte zirkuläre Chloroplasten. Error, Fehler oder Abbrüche im Programm. Partial, keine zirkulären Chloroplasten aber Contigs richtig identifiziert. Incomplete, Nicht richtig identifizierte Chloroplasten.

Tool	SUCCESS	%	ERROR	PARTIAL	INCOMPLETE
CE	136	~15%	54	3	706
FP	266	~30%	29	11	593
Summary	303	~34%	-	-	-

4.7 Varianz Analyse

Um die Varianz Analyse durchzuführen und vor allem zu überprüfen ob die Assambler bzw. die Programme an sich einen Einfluss darauf haben, indem sie z.B. zufällige Seeds verwenden oder zufällige Daten bevorzugen, wurden 89 Datensätze des 1001 Genom Projekts verwendet. Diese 89 Datensätze zeichnen sich dadurch aus, dass sowohl der chloroExtractor als auch der fast-plast diese zu vollständigen Chloroplasten zusammengebaut haben. Diese Datensätze wurden noch zehn weitere Male berechnet. So wurden auf elf mal 89 Datensätzen überprüft welche Einflüsse die Programme auf die Varianz haben. Der chloroExtractor und somit der Assambler SPAdes brachte bei allen elf Durchläufen die exakt gleichen Sequenzen heraus. Dieses Programm arbeitet also 100% Reproduzierbar (Fig. 12). Im Gegensatz dazu der fast-plast, dieser schaffte es nicht wieder bei allen elf Durchläufen alle Chloroplasten wieder korrekt zusammen zubauen, bei bis zu neun verschiedenen Datensätzen konnte kein Erfolgreiches Ergebnis erzielt werden (Fig. 11). Interessanter weiße waren nicht immer die selben Datensätze betroffen, so konnten bei einigen Durchläufen ein Erfolg erreicht werden, bei dem nächsten Durchlauf aber nicht. Ob dies ein Zufalls Effekt des Programms oder der verwendeten Rechner-Infrastruktur ist, konnte nicht überprüft werden. Die zweite Varianz Analyse bzw. SNP calling wurde

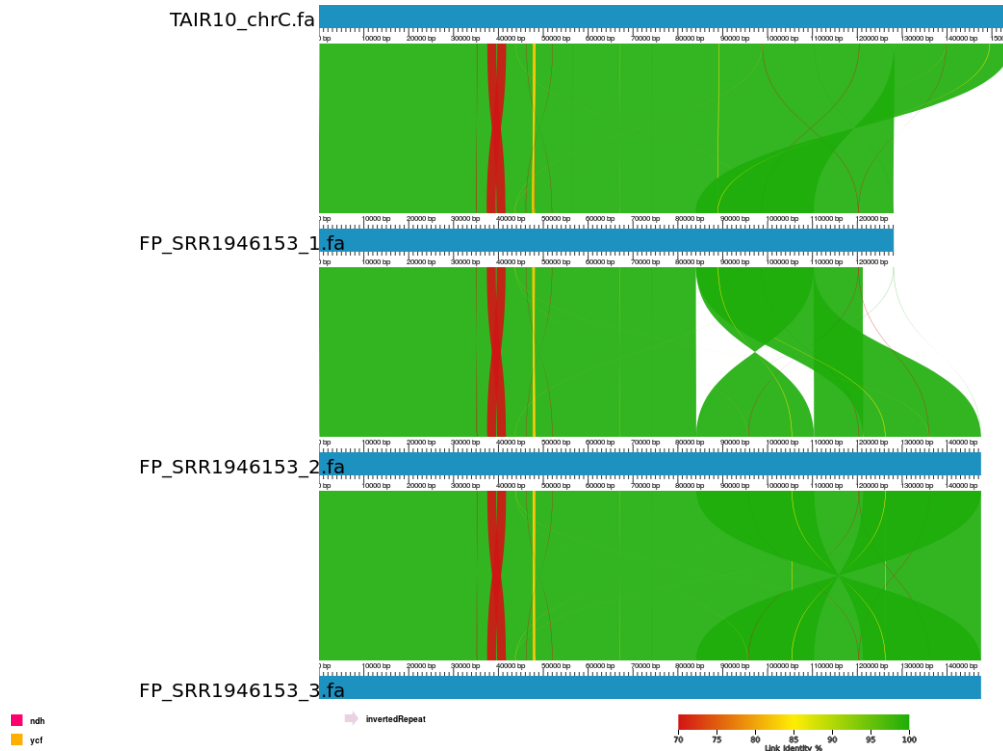


Figure 11: **fast-plast SRR1946153** Drei verschiedene Läufe auf den selben Daten, der fast-plast schafft einen davon nicht (Lauf 1), den anderen aber schon (2 u. 3). Hier fehlt ein Teil des IR, wodurch auch nicht als Erfolg gewertet wird.

auf allen Erfolgreich zusammengebauten Chloroplasten durchgeführt. Das SNP calling ergab dass auf allen 303 Chloroplasten insgesamt 2128 SNPs gefunden wurden. Diese Ergebnisse werden für die GWAS Analyse verwendet.

4.8 GWAS

Die GWAS Analyse, welche mit den 303 kompletten Chloroplasten aus den Daten des 1001 Genom Projekts und den 2128 gefunden SNPs durchgeführt wurde, konnte nur auf zwei verschiedenen Traits berechnet werden. Dies waren die Eigenschaften Flowering Time bei 16°C sowie bei 10°C. Dies sind die beiden Traits am besten untersucht sind und deswegen auch die meisten Daten beinhalten. Für alle anderen Traits konnten keine Berechnungen erstellt werden da die Datenmenge nicht für eine GWAS ausreichend ist.

4.9 Struktur Varianz Analyse

Für die Struktur Varianz Analyse konnten keine Ergebnisse erzielt werden, Grund hierfür war unter anderem fehlende Zeit. Aber auch konnten keine Programme gefunden werden welche mit kompletten Chloroplasten umgehen konnten. Die meisten nutzten direkt

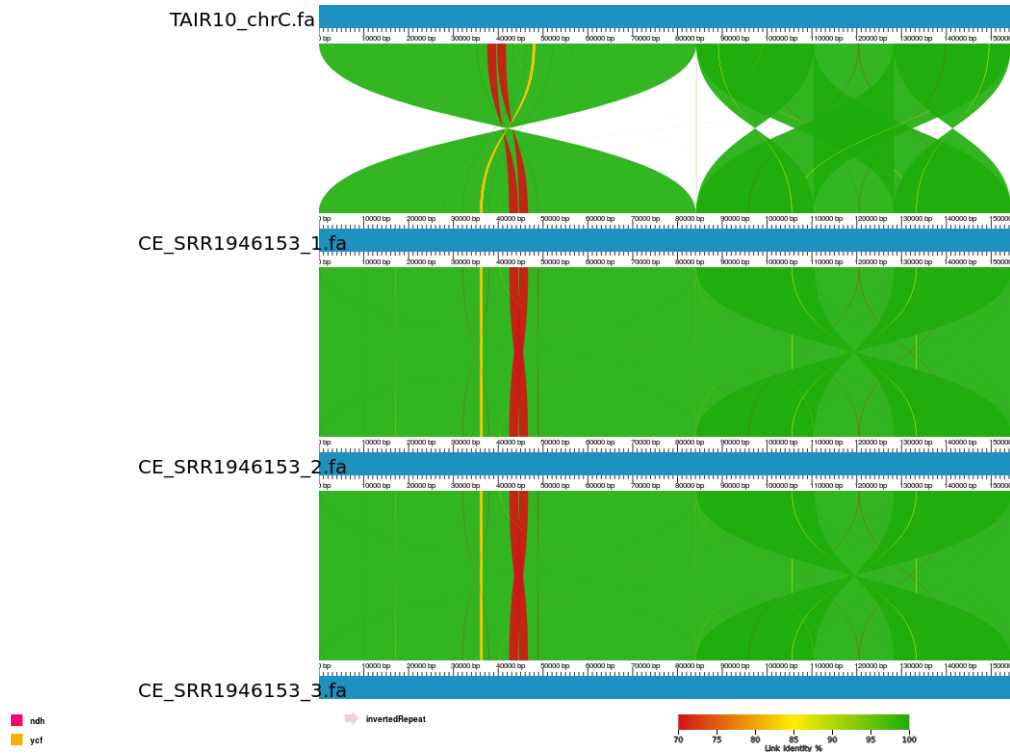


Figure 12: **chloroExtractor SRR1946153** Drei verschiedene Läufe auf den selben Daten, der chloroExtractor bringt das gleiche Ergebnis für alle Durchläufe. Da die Orientierung von LSC, SSC und IR nicht aus short Reads heraus gelesen werden kann, kommt es vor das diese Verdreht sind zur Referenz.

Illumina short reads, wie Delly^{61, 65} oder Breakdancer^{62, 66}.

4.10 Neue Chloroplasten

Aus NCBI wurden 79657 Datensätze heruntergeladen, dies sind alles Pflanzengenome. Diese Liste wurde mit den Einträgen von CpBase (stand 20.06.2018, 2069 Genome von 942 verschiedenen Spezies)verglichen. Es blieben die Übrig, welche keinen Eintrag in CpBase haben. Von diesen 79657 blieben nur 49 Datensätze. Diese wurden auf fast-plast und chloroExtractor benutzt, und es wurden 17 zirkuläre Chloroplasten erfolgreich zusammengebaut (s. Tab. 7). Somit wurden 17 neue Chloroplasten von Spezies welche zuvor noch keinen genomisch bekannten Chloroplasten hatten erfolgreich erstellt (s. Tab. 8).

Table 7: **Neue Chloroplasten** Von den 49 Spezies welche bisher noch keinen Eintrag in CpBase hatte konnten mit Hilfe des fast-plasts und des chloroExtractors 17 neue bisher nicht bekannte Chloroplasten Genome gebaut werden

Tool	SUCCESS	ERROR	PARTIAL	INCOMPL	Total
CE	4	20	16	9	
FP	15	7	22	5	
Summary	17	-	-	-	49

5 Diskussion

5.1 Definition von Success, Einteilung der Erfolge über Genom Länge.

Jegliche Einteilung in die Erfolgs Kategorien: Success, Partial, Incomplete_high und Incomplete_low und Error werden von einem Skript übernommen welches zunächst den SeqFilter benutzt um Informationen über diese Datei zu erhalten. Der SeqFilter zählt die Sequenzen sowie deren Größe. Das Evaluationsskript des jeweiligen Programms teilt aufgrund dieser Daten in die Kategorien ein (s. Tab. 1). Diese Variante ist zwar voll Automatisiert doch nicht Fehlerlos, so können Falsch Positive Sequenzen vorkommen. Diese könnte eine Sequenz aus 150 kbp Adenin sein, und das Skript würde es als einen Success ansehen. Die Daten wurden Stichprobenartig überprüft und dies kam in diesen Stichproben nicht vor, doch ist es nicht auszuschließen. Um sicher zu gehen müsste jeder erstellter Chloroplast auf eine Referenz gemapt werden oder sogar durch Sequenzierung bestätigt werden. Erste Möglichkeit wäre nur Rechenaufwand, könnte aber bei Chloroplasten die noch nicht veröffentlicht wurden oder keine Referenz besitzen schwer werden,

⁶⁵Rausch T, Zichner T, Schlattl A, et al. (2012) Delly: structural variant discovery by integrated paired-end and split-read analysis. Bioinformatics, <https://doi.org/10.1093/bioinformatics/bts378>

⁶⁶Chen K, Wallis JW, McLellan MD, et al. (2009), BreakDancer: an algorithm for high-resolution mapping of genomic structural variation, Nature Methods 10.1038/nmeth.1363

Table 8: **Liste neue Chloroplasten** Liste von 17 Spezies welche mit Hilfe des fast-plast und des chloroExtractors nun ein bekanntes Chloroplasten Genom besitzen.

SRA	Spezies
DRR057122	<i>Momordica charantia</i>
DRR089517	<i>Betula chichibuensis</i>
ERR1462646	<i>Hippophae rhamnoides</i>
ERR2001942	<i>Betula pendula</i>
ERR2003066	<i>Potentilla micrantha</i>
ERR2174632	<i>Solanum pennellii</i>
ERR2187925	<i>Geum urbanum</i>
SRR1503730	<i>Agave tequilana</i>
SRR2847417	<i>Manihot glaziovii</i>
SRR3194007	<i>Artocarpus altilis</i>
SRR3724930	<i>Taraxacum S3</i>
SRR4457832	<i>Pityopsis pinifolia</i>
SRR5046394	<i>Ephedra gerardiana</i>
SRR5464169	<i>Trema orientalis</i>
SRR5590327	<i>Lagenaria siceraria</i>
SRR5799057	<i>Fragaria vesca</i>
SRR5838021	<i>Populus deltoides</i>

zweite Möglichkeit ist sehr Kosten intensiv würde aber letzte Zweifel beseitigen. Eine Verbesserung des Skripts könnte auch eine strengere Beurteilung sein, zumindest wenn man mehr Grundinformationen hat. So könnten bei den Versuchen mit den *A.thaliana* des 1001 Genom Projekts die Grenzen Strenger gewählt werden, da es sich hier immer um die gleiche Spezies handelt. Doch könnten somit die Anzahl der Falsch Negativen erhöht werden, z.B. wenn eine *A.thaliana* Art eine Struktur Variante besitzt mit Verlust eines IR. Die Grenzen wurden bewusst großzügiger Gewählt, da dies den größten Teil der Chloroplasten abdecken dürfte. Gerade bei Chloroplasten welche bisher nicht veröffentlicht oder bekannt sind ist eine Abschätzung schwer, da die Größen von Chloroplasten doch sehr Variieren können. Eine weitere Möglichkeit zu testen ob es sich wirklich um einen Chloroplasten handelt wäre die Verwendung von Benchmarking Universal Single-Copy Orthologs (BUSCO⁶⁷), hierzu werden extrem konservierte orthologe Gene verwendet und überprüft ob diese alle vorhanden sind. Da ein Chloroplast Genom an sich sehr konserviert ist könnte eine Anzahl von Genen genommen werden und diese in einem solchen Modell verwendet werden.

⁶⁷Simão FA, Waterhouse RM, Ioannidis P, et al. (2015), BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. Bioinformatics, 10.1093/bioinformatics/btv351

5.2 Die Entscheidung für fast-plast und chloroExtractor

Es wurde im Ergebnis Teil erklärt warum gerade der fast-plast und der chloroExtractor weiter verwendet wurden. Doch gibt es auch gründe warum sich speziell gegen andere Programme entschieden wurden. So wurde sich gegen den IOGA entschieden, nicht nur weil er langsam ist sondern auch weil er keinerlei Log File während des Prozesses schreibt, erst wenn dieser Komplette beendet ist, so war vor allem zu beginn schwer nachzuvollziehen ob der IOGA nun wirklich noch Arbeitet oder evtl in irgendeinem Loop fest hängt oder sogar aufgehört hat zu arbeiten aber den Prozess nicht beendet. Auch wurde der IOGA zum letzten mal vor zwei Jahren geupdated, es scheint also keine Regelmäßige Wartung oder Verbesserung statt zu finden. Es wurde sich auch gegen den NOVOPlasty entschieden, dieser benötigt zwar keine Abhängigkeiten da er komplett in Perl geschrieben ist, doch hat dies einige Probleme mit sich gebracht. So werden z.B. nicht alle Read header richtig eingelesen wenn der dazugehörige Reguläre Ausdruck (Regular Expression - RegEx) nicht komplett passt, dies kam häufiger vor da nicht alle Header gleich aufgebaut sind und wohl ein paar nicht abgedeckt wurden. Das zweite Problem mit NOVOPlasty ist die Konfigurationsdatei, diese muss exakt dem Beispiel entsprechen und darf nicht ein Zeichen mehr oder weniger enthalten, oder gar Zeilen. Da diese Datei nicht über RegEx eingelesen wird sondern Zeile für Zeile durchgegangen wird. So kam es gerade am Anfang vor das der NOVOPlasty gar nicht funktionierte da ein Leerzeichen in einer nicht verwendeten Option fehlte. Der NOVOPlasty scheint noch Regelmäßig geupdated zu werden, doch Änderte sich bei diesen Updates der Aufbau der Konfigurationsdatei, weswegen jedes mal das Skript zum erstellen dieser Datei umgeschrieben werden musste. Auch warf der NOVOPlasty Fehler in denen gesagt wird dass der Seed nicht lesbar oder inkompatibel sei. Doch wurde bei jedem Versuch als Seed die gleiche Datei verwendet, und dieser Fehler trat nur ab und zu auf. Der Org.ASM brachte zwar erfolgreiche Ergebnisse, im Vergleich würde er auf dem dritten Platz landen, doch gab es einige Probleme bei der Installation. Nur in einem Dockercontainer mit einigen Tricks konnte es geschafft werden dieses Programm erfolgreich zu installieren. Der GetOrganelle konnte zwar mit dem fcg.pl Skript des chloroExtractors automatisiert werden, doch unterschlägt dies dann das eigentliche Endprodukt des GetOrganelles, da das verbesserte bzw. getrimmte fastg nicht vom fcg.pl Skript erkannt wurde und deshalb nur das fastg aus SPAdes selbst verwendet werden kann, dies aber häufig schlechter Ausfällt als das getrimmte oder gar das fastg welches SPAdes im chloroExtractor ausgibt.

5.3 Fazit aus der Erfolgschance

Es wurden in dieser Arbeit 303 Chloroplasten Genome von *Arabidopsis thaliana* und 40 von verschiedenen Spezies (GO-Preprint) sowie 17 neue (s. Tab. 8) erstellt. Nimmt man von den Versuchen die gesamte Zahl, so konnten in etwa 30% der Datensätze zu Chloroplasten Genomen führen. Dies entspricht Tatsächlich mehr als am Anfang der Arbeit angenommen, hier wurden in etwa 10 - 20% geschätzt⁴². Allerdings auch ohne die anderen Programme, abgesehen von chloroExtractor und org.ASM, getestet zu haben. Nimmt man den nur die Erfolgschance von chloroExtractor war diese erste Abschätzung

gar nicht so schlecht.

5.4 Erhöhen der Erfolgsrate

Es gibt mehrere Möglichkeiten wie eine Erfolgsrate erhöht werden könnte. So könnte versucht werden auf die Daten speziell die Start Parameter festzulegen. Dies würde allerdings einiges an Tests benötigen. Auch könnten die Parameter jedes mal geändert werden, dann aber unter dem Verlust einer Automatisierbarkeit. In diesen Versuchen wurden verschiedenen Große Datensätze verwendet, und es lässt sich nicht sagen ob eine Erhöhung dieser einen echten Vorteil bringen würde, hierzu müssten alle Daten noch einmal gestartet werden, dann mit erhöhten oder niedrigeren Readmengen. Theoretisch kann dies einen Zuwachs an Erfolg bringen, wenn das verwendete Programm denn auch alle Daten verwendet, die es bekommt und nicht irgendeinen Cutoff ab einer bestimmten Daten bzw. Read Menge hat. Wenn die kompletten Daten verwendet werden würden hätte dies auch den Vorteil das man sicher gehen kann dass die Daten nicht sortiert wurden, indem man diese einfach nochmal durch mischt. Dies ist wichtig, vor allem bei Programmen welchen einen Cutoff benutzen, denn hier könnte es vorkommen, dass wenn eine Datei sortiert ist Chloroplasten Reads am Ende der Datei liegen und diese somit gar nicht erst benutzt werden. Dennoch ist zu beachten, je mehr Daten natürlich verwendet werden desto länger brauchen die Programme. Zudem kommt eine erhöhte Download Zeit und evtl. die Zeit die gebraucht wird um die Dateien zu mischen bevor diese für die Programme verwendet werden können.

5.5 Etablieren einer einfachen scanning Routine

In dieser Arbeit wurde gezeigt das eine voll Automatische Lösung für das scannen von Chloroplasten in Pflanzen Genom Daten möglich und auch erfolgreich ist. Die hier verwendeten Skripte können frei verwendet und angepasst werden. Doch kann dies alles auch in einem kompletten Dockercontainer benutzt werden. Der chloroExtractorTeam Screening Container⁶⁸, kann verwendet werden um komplett automatisch die Daten von NCBI herunter zu laden, diese zu mischen (mit einem festen Seed) und dann den chloroExtractor und den fast-plast zu verwenden um diese Daten zu verarbeiten. Hierzu muss lediglich der Container gestartet werden und der run.sh Befehl mit der Passenden SRA Nummer gegeben werden. Dieser Container wird gerade verwendet um weitere 12393 Datensätze zu durchsuchen und Chloroplasten zu bauen. Diese Container ist sehr einfach zu benutzen, und alles was dafür gebraucht wird ist Docker⁴⁴ oder ein Programm welches Dockercontainer ausführen kann wie z.B. Singularity⁵³.

5.6 GWAS

Es wurde eine GWAS Studie auf den 303 *A.thaliana* Chloroplasten durchgeführt, doch konnte dies nur auf zwei verschiedenen Eigenschaften berechnet werden. Hierzu gehört die Flowering Time bei 16°C sowie bei 10°C. Die restlichen Arapheno Traits konnten nicht

⁶⁸https://github.com/chloroExtractorTeam/screening_container

berechnet werden. Dies ist vor allem der Fall da zu wenige Daten zur Verfügung stehen, sowohl von unserer Seite aus als auch von der Eigenschaften Seite aus. Eigenschaften wie Größen, Blütenbreite oder Form sind nicht gut genug Katalogisiert um eine geringe Datenmenge, wie sie hier benutzt wurde abzudecken. Dies heißt nicht das keinerlei Assoziation zwischen Chloroplasten Varianz und diesen Eigenschaften besteht. Dies zeigt lediglich dass noch mehr Daten benötigen werden um diese GWAS Studie zu beenden.

5.7 Struktur Varianz

Es wurde angenommen, dass ein kompletter Chloroplast für eine Struktur Varianz Analyse einen großen Vorteil bringt im Vergleich zu Illumina short Reads. Diese sind häufig zu kurz um große Invertierungen oder neu Anordnungen zu überspannen. Leider konnte diese Annahme nicht überprüft werden, da zu wenig Zeit vorhanden war. Als auch Programme, welche eine solche Analyse durchführen nicht erfolgreich benutzt werden konnten

5.8 Fazit und Zukunftsaussichten

Chloroplasten Genome können vielseitig verwendet werden und mit der immer steigenden Anzahl dieser Genome können mehr und mehr Analysen durchgeführt werden. Hier wurde eine Möglichkeit gezeigt wie man solche Chloroplasten Genome aus bereits existierenden Daten bauen kann. Dies ist aber nur möglich da alle Daten dank Open Science verfügbar waren. Es wurden 360 Chloroplasten aus bereits vorhandenen Daten erzeugt, die meisten aus *Arabidopsis thaliana* Daten, die vom 1001 Genom Projekt verwendet wurden. Es wurden aber auch 17 Neue Chloroplasten Genome erzeugt. Zudem wurde ein erster Ausblick auf die verschiedenen Analysen gegeben welche mit Chloroplasten Genomen möglich sind. Schaut man sich die Updates der verschiedenen Programme an, so wird ein Teil davon immer noch geupdated bzw. verbessert. Je mehr Chloroplasten Genome mit der Zeit verfügbar werden, desto mehr Analysen können mit diesen Chloroplasten durchgeführt werden. So zeigte sich hier bei dem versuch einer GWAS, dass es bei zu wenigen Daten zu Problemen kommen kann. Diese Arbeit zeigt vor allem eines, es ist durch Automatisierung möglich Chloroplasten Genome aus Pflanzlichen Sequenzdaten zu bauen. Dank der Automatisierung ist hier lediglich Rechenpower von Nöten. So können in Zukunft noch viel mehr Chloroplasten Genome erzeugt werden.

6 Abbildungs- und Tabellenverzeichnis

List of Figures

1	Chloroplast Genom Einteilung	7
2	Chloroplast Genom: Gen Klassen	8
3	chloroExtractors Logo	11
4	Ablauf des chloroExtractors	13
5	Bandage - Fastg Visualisierung	16

6	Automatisierungsskripts	25
7	Bandage: GetOrganelle SRR1945443	27
8	Bandage: chloroExtractor SRR1945443	28
9	Upset Diagramm GO-Preprint	30
10	AliTV SRR5602602 - Laurus nobilis	31
11	fast-plast SRR1946153	34
12	chloroExtractor SRR1946153	35

List of Tables

1	Erfolgsraten Einteilung	21
2	Test Datensatz: Simulierte Daten	26
3	Test Datensatz: 1001 Genom Project, 11 Datensätze	26
4	Test Datensatz: GetOrganelle Preprint, 11 Datensätze	29
5	Laufzeit und Ressourcenverbrauch	32
6	Datensatz: 1001 Genom Project	33
7	Neue Chloroplasten	36
8	Liste neue Chloroplasten	37
9	Dockercontainer	42
10	Git Links	43

7 Anhang

7.1 Dockercontainer

Table 9: **Dockercontainer** Alle erstellten Dockercontainer stehen zur freien Verfügung.

Programm	Dockerhub link
chloroExtractor	https://hub.docker.com/r/chloroextractorteam/chloroextractor/ Build: b5uvjvdbncyndhjngua85nv
fast-plast	https://hub.docker.com/r/chloroextractorteam/fast-plast_docker/ Build: bgrmngfwpil4sk2kezi9f
NOVOPlasty	https://hub.docker.com/r/chloroextractorteam/novoplasty_docker/ Build: bf9adepndze96bcnteyeabk
IOGA	https://hub.docker.com/r/chloroextractorteam/ioga_docker/ Build: bwnf4xtzhohfcstqjqsqvqvw
GetOrganelle	https://hub.docker.com/r/chloroextractorteam/getorganelle_docker/ Build: bwt3bus2r7utsjprgmjmjuc
Org.ASM	https://hub.docker.com/r/chloroextractorteam/org.asm_docker/ Build: bmcx88c2d79orvuykgivy4q
Screening Container	https://hub.docker.com/r/chloroextractorteam/screening_container/ Build: bdsankaqpukcgfkmkmdq9yud

Table 10: **Git Links** Alle verwendeten Programme stehen zur freien Verfügung.

Programm	git link
chloroExtractor	https://github.com/chloroExtractorTeam/chloroExtractor
fast-plast	https://github.com/mrmckain/Fast-Plast
NOVOPlasty	https://github.com/ndierckx/NOVOPlasty
IOGA	https://github.com/holmrenser/IOGA
GetOrganelle	https://github.com/Kinggerm/GetOrganelle
Org.ASM	https://git.metabarcoding.org/org-asm/org-asm
Screening Container	https://github.com/chloroExtractorTeam/screening_container

7.2 Danksagung

Ich möchte allen Danken, welche mich in meiner Zeit als Student vor allem in den letzten Semestern unterstützt haben. Zunächst Danke ich dem kompletten CCTB, in dem ich in den letzten Semestern mit mehr als nur Kollegen zusammenarbeiten durfte. Hier möchte ich im besonderen meinen Betreuern Markus Ankenbrand und Jörg Schulz danken, welche sich die Zeit nehmen müssen diese Arbeit zu korrigieren und immer mit Rat zur Seite standen. Speziell möchte ich auch meinen ganz besonderen Dank an Frank Förster aussprechen, der uns mehr als nur beratend zur Seite stand und mehrere Nächste geopfert hat um unseren chloroExtractor zu verbessern und zu fixen. Auch möchte ich aber meiner kompletten Familie danken, meinen Eltern Roland und Maria, sowie meinen beiden Brüdern Florian und Christopher, sowie meinem Bruder im Geiste Martin "Löwe" Piekar, sowie natürlich den Damen welche es mit diesen drei Chaoten aushalten, die da wären Vanja, Julia und Charlotte. Natürlich danke ich auch all meinen anderen Freunden welche es mit mir all die Jahre ausgehalten haben, und dies planen noch weiter zu tun, so wie der World of Warcraft Community, der Gilde Maestri delle Arte und dem Raid Invictus, welche mich immer gut unterhalten haben, und das ein oder andere mal evtl. auch zu viel von der Arbeit abgehalten haben. Auch geht mein Dank an Kaffee, vielen dank für Substitution von Schlaf mit Koffein.

7.3 Skripte

Alle verwendeten Skripte können gefunden werden unter:

https://github.com/chloroExtractorTeam/chloroplast_landscape/tree/Documentation_SPfaff Zudem sind diese auf der beiliegenden CD zu finde.

7.4 Tabellen

Alle Ergebnis Tabellen können eingesehen werden unter:

https://github.com/chloroExtractorTeam/chloroplast_landscape/tree/Documentation_SPfaff Zudem sind diese auf der beiliegenden CD zu finde.

Eigenständigkeitserklärung

ERKLÄRUNG gemäß ASPO vom 5.8.2009 § 23 Abs. 10

Hiermit versichere ich, dass ich vorliegende Arbeit selbstständig verfasst, keine anderen als die angegebenen Quellen und Hilfsmittel benutzt und die Arbeit bisher oder gleichzeitig keiner anderen Prüfungsbehörde unter Erlangung eines akademischen Grades vorgelegt habe.

Würzburg, August 7, 2018

Simon Pfaff