

otb: Creating a HiC/HiFi Pipeline to Assemble the *Prosapia bicincta* Genome

David C Molik^{*1}, Amanda Stahlke^{2,3}, Sharu Paul Sharma⁴, Taylor J. Simmonds⁵, Renee Corpuz⁵, Angela Kauwe⁵, Jeremy Schrader⁵, Charles Mason⁵, Sheina Sim⁵ and Scott Geib⁵

¹Arthropod-borne Animal Diseases Research Unit, Center for Grain and Animal Health Research, United States Department of Agriculture, Agricultural Research Service, 1515 College Ave, Manhattan, KS 66502 USA

²Colorado Mesa University, Department of Biological Sciences, Wubben Hall and Science Center, 1100 North Avenue, Grand Junction, CO 81501-3122 USA

³RiversEdge West DNC, Tamarisk Beetle Monitoring Program, 125 N 8th St, Grand Junction, CO 81501 USA

⁴Genome Informatics Facility, Iowa State University, 2200 Osborn Drive, Ames, IA 50011-4009 USA

⁵Tropical Pest Genetics and Molecular Biology Research Unit, Daniel K. Inouye U.S. Pacific Basin Agricultural Research Center, United States Department of Agriculture, Agricultural Research Service, 64 Nowelo St, Hilo HI 96720 USA

*Corresponding author: david.molik@usda.gov

1 Abstract

2 Two-lined spittlebug (Hemiptera: Cercopidae) is an insect herbivore that is widely distributed throughout the eastern portion of the United States. This insect
3 is a significant pest of turf grasses and pasturelands, where feeding causes wilting and phytotoxemia, resulting in plant mortality. This pest is endemic to
4 tropical and subtropical regions of the world and has a long-term impact on the composition of the plant community in grassland landscapes, altering plant
5 communities in rangelands and posing a substantial threat to cattle production, citing a potential threat to livestock production. Using phased genomes
6 to better understand the physiology, ecology, and evolution of this insect pest can yield novel targets to exploit for sustainable pest management. The
7 introduction of a new HiC/HiFi phased genomics assembly pipeline called Only The Best (otb) has solved several data management problems that were
8 prevalent in the creation and storage of genome assemblies.

9 **Keywords:** genome assembly; non-model organism; haplotype phasing; next-generation sequencing; assembly error correction

1 Introduction

2 The USDA-ARS AgPest 100 Initiative (Ag100pest) aims to gener-
3 ate high-quality genome assemblies of existing and/or emerging
4 pest insect species that threaten agricultural production (Childers
5 *et al.* 2021). High-quality genome assemblies can root and inform
6 both basic and applied research. The time cost in production of
7 multiple genomes assemblies can cause inefficiencies in projects
8 such as the ag100pest; a HiFi/HiC assembly pipeline is required
9 for efficiencies in these projects. In order to induce implementa-
10 tion of such a pipeline, it is required to utilize the pipeline in a
11 test scenario. As the Ag100Pest often works with true bugs, and
12 true bug genomics are relatively unexplored Jiang *et al.* (2021), an
13 agriculturally significant true bug is the ideal test case for such a
14 pipeline, Two-line spittlebugs (*Prosapia bicincta*) are such test case.
15 Two-lined spittlebug (Hemiptera: Cercopidae) is an insect herbi-
16 vore that is widely distributed throughout the eastern portion of
17 the United States (Potter *et al.* 1991; Braman and Abraham 1995). In
18 2016, the two-lined spittlebug was first detected on Hawai'i Island
19 (Thorne *et al.* 2017). The immature lifestages of this species are sig-
20 nificant pests of turfgrass and pasturelands, where feeding causes
21 wilting and phytotoxemia resulting in plant mortality (Byers and
22 Wells 1966; Fagan and Kuitert 1969; Joseph and Jespersen 2021).
23 Since its establishment in Hawai'i, two-lined spittlebugs have had
24 significant cascading effects on plant communities in rangelands,
25

altering the composition of plant communities of grass-dominated
landscapes and posing a significant threat to cattle production
(Bremer *et al.* 2021). Uncovering the behavioral and metabolic
strategies insects employ to exploit plants is an important step
in determining their pest status. Like other cercopids, the two-
lined spittlebug eats a nutritionally impoverished diet in xylem
sap (Mattson Jr 1980). The processing of this diluted diet produces
the characteristic spittle masses at the base of the plant. Spittlebugs
have some metabolic innovations to contend with these diluted
diets. Like other hemipterans that feed on phloem and xylem, spit-
tlebugs harbor endosymbiotic bacteria that reside in specialized
structures called bacteriomes. Spittlebugs harbor two symbionts
in independent organs. *Sulcia muelleri* is ubiquitous in Auchen-
orrhyncha, while the other symbiont can be *Zinderia insectola* or
a *Sodalis*-like microorganism (Koga *et al.* 2013; Koga and Moran
2014). These symbionts help provide complementary sets of es-
sential amino acids through complex and intertwined metabolic
pathways (Ankrah *et al.* 2020).

Biological control of two-lined spittlebugs in grassland ecosys-
tems is inherently challenging. Adults are long-lived and highly
fecund (Peck 1998), few commercial grass cultivars exhibit resis-
tance and / or tolerance to this insect (Braman *et al.* 2014; Joseph
and Jespersen 2021), and nymphs feed in protected areas at the
base of the plant, which facilitates their escape from natural en-
emies (Nachappa *et al.* 2006). In the case study we use here, a
genome assembly of two-lined spittlebug is a critical first step to-
wards understanding the physiology, ecology, and evolution of

1 this herbivorous pest and may yield novel targets to exploit for
2 sustainable pest management.

3 Phased genomes, which determine both chromosomes of a
4 diploid, are a valuable resource for analyzing genetic variation
5 within populations (Snyder *et al.* 2015). This is particularly impor-
6 tant in agricultural research, where understanding genetic vari-
7 ation is crucial to understanding the evolution and spread of
8 resistance genes that can impact outbreaks of insect pests (Leftwich
9 *et al.* 2015). The creation of a large number of phased arthropod
10 genomes, a number of which could be enabled by more hands-off
11 and accessible bioinformatic pipelines, has several applications
12 (Tewhey *et al.* 2011). Having phased genomes to better understand
13 the genetic variation of two-lined spittlebugs in an invasive popu-
14 lation could lead to valuable insights into how the insect adapts to
15 insecticides or adapts to new environments.

16 To create phased genomes of the two-lined spittlebug and other
17 insects, we developed a new HiC / HiFi phased genomic assembly
18 pipeline called otb, or Only The Best [genome assembly tools]. Our
19 pipeline reduces the time spent organizing data, installing and
20 calibrating bioinformatic tools, and, therefore, performing analysis.
21 By implementing this pipeline, we reduced the amount of time re-
22 quired to produce a usable genome. The careful implementation of
23 data management and standardization also significantly reduced
24 team effort in genome assembly creation. otb is a software tool
25 that utilizes the nextflow programming language (Tommaso *et al.*
26 2017) and is accessed through a bash script. To ensure a consistent
27 compute environment across users, otb is implemented within a
28 singularity container management software, which enables users
29 to share containers with other users within the same environment.
30 The use of nextflow provides the benefit of parallel task execution
31 and efficient management of compute resources, while singularity
32 ensures a consistent and reproducible compute environment. This
33 also eliminates the need for software duplication across a high-
34 performance computing cluster (HPC). The development of otb
35 was primarily for the United States Department of Agriculture,
36 Agricultural Research Services' ag100pest and Beenome projects,
37 where large numbers of reference genomes needed to be created.
38 However, otb can be used for any project that requires the automa-
39 tion of HiFi genomes until human involvement is necessary.

40 Materials and Methods

41 Sampling

42 Two-lined spittlebug were collected from the University of
43 Hawai'i's Kona research station (79-7381 Hawai'i Belt Road, Holu-
44 aloa, HI) were Shannon Wilson reared them. The samples were
45 flash-frozen on site on August 13th, 2020.

46 DNA Extraction, PacBio Library Preparation, and Sequenc- 47 ing

48 High molecular weight DNA (HMW DNA) was extracted for the
49 preparation of the PacBio HiFi library from a single adult male *P.*
50 *bicincta*. The sample was cryoground using a Spex GenoGrinder
51 2010, and DNA was extracted from the ground tissue using the
52 Qiagen MagAttract HMW DNA kit (Cat# 67563) following the
53 kit protocol. The concentration of the extracted HMW DNA was
54 quantified using the Qubit 1x dsDNA HS kit (Q33230), and DNA
55 purity was assessed using UV-vis spectroscopy. The size distribu-
56 tion of the extracted HMW DNA was evaluated using an Agilent
57 Femto Pulse instrument with the Genomic DNA 165kbp kit (Cat #
58 FP-1002-0275).

59 Before preparation of the PacBio HiFi library, the extracted
60 HMW DNA was sheared using the Diagenode Megaruptor 2 with

61 the 20 kbp shearing program, to target a sheared DNA size of ap-
62 proximately 10-15kbp. Sheared DNA was used to prepare PacBio
63 HiFi libraries, using PacBio's Express Template Prep Kit 2.0 (PN:
64 102-088-900) following the kit protocol, with the optional nuclease
65 digestion step after library preparation. PacBio libraries were size
66 selected with 40% diluted AMPure PB beads (PN: 102-182-500) to
67 remove library molecules shorter than 3kbp following PacBio's
68 protocol. The final libraries selected for size were quantified using
69 the Qubit 1x dsDNA HS kit, and the library size was checked using
70 the Agilent Femto Pulse with the Genomic DNA 165kbp kit. The
71 PacBio libraries were sequenced on a PacBio Sequel IIe using a
72 30 hour movie time with 2 hours of pre-extension. Sequencing
73 reaction was prepared using the Sequel II Binding Kit 2.2.

74 HiC Library Preparation

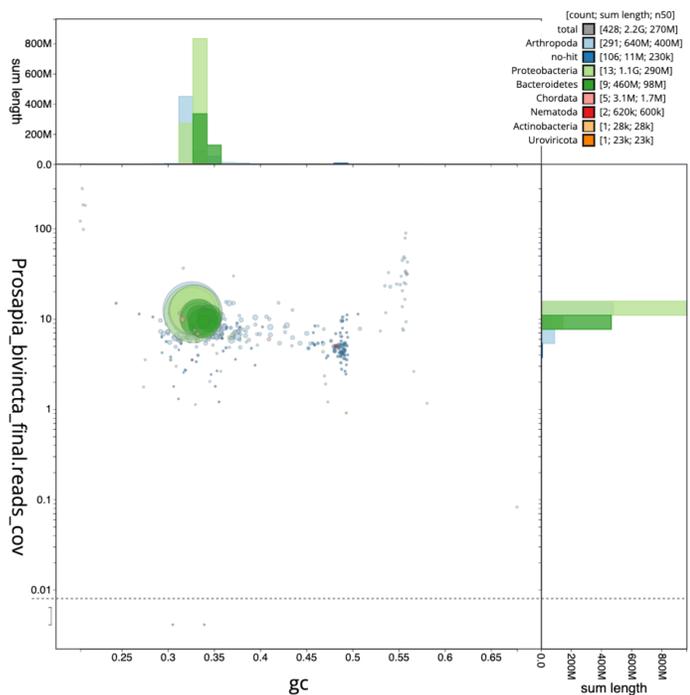
75 To prepare a HiC library, a single adult male *P. bicincta* was cryo-
76 ground and then fixed in freshly prepared TC fixation buffer, fol-
77 lowing the low-input crosslinking protocol for the Arima HiC 2.0
78 kit. Proximity ligation was performed on cross-linked samples
79 using the Arima HiC 2.0 kit following the manufacturer's protocol.
80 Prior to preparation of the Illumina library, the proximity-ligated
81 DNA was sheared using a Diagenode Biorupter and then size-
82 selected to enrich the sheared DNA in the 200-600bp range. The
83 size distribution of the selected size sheared DNA was checked
84 using an Agilent TapeStation with the High Sensitivity D5000
85 ScreenTape (Cat # 5067-5588) before proceeding to library prepara-
86 tion. The Illumina HiC library was prepared using the Swift Accel
87 NGS 2S Plus DNA Library kit following the protocol outlined in
88 the Arima HiC 2.0 kit. Library amplification was performed us-
89 ing the KAPA Library Amplification Kit with Primer Mix (Cat#
90 KK2620), with 8 cycles of PCR. The final libraries were quantified
91 using the Qubit 1x dsDNA HS kit, and the library size distribution
92 was checked using an Agilent TapeStation. The HiC library was
93 sequenced on an Illumina NovaSeq 6000.

94 otb Genome Assembly

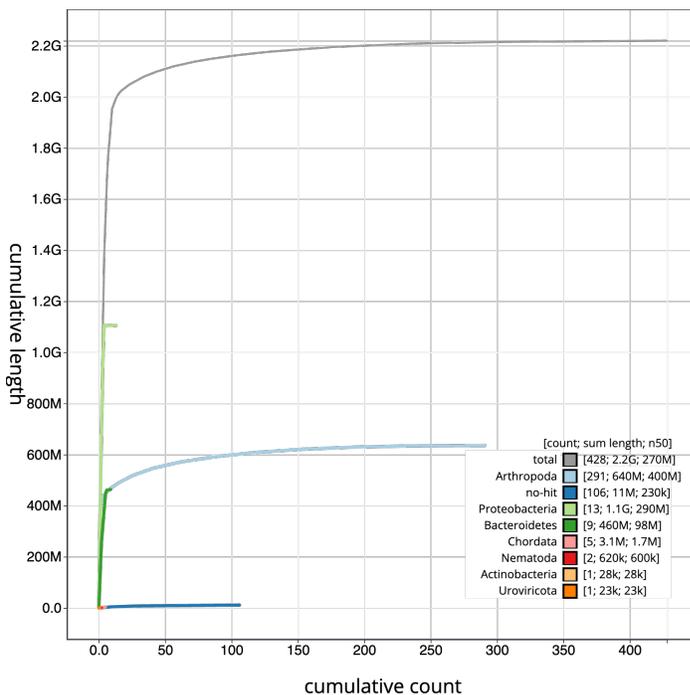
95 Genome assembly of Two-lined spittle-bug raw HiFi and HiC reads
96 was completed with otb. otb was set to use the "merfin" polishing
97 option, additionally, at the k-mer creation step, KMC3 was used.
98 Juicebox was then utilized for HiC correction and assembly was
99 modified accordingly. shhquis.jl was utilized post assembly to re-
100 arrange chromosomes for hic mapping, reducing the manual time
101 required in the hic rearrangement step. Manual hic rearrangement,
102 which was completed with Yabs exported HiC data was exported
103 to Juicebox and final genome assembly report which was com-
104 pleted with Blobtools. otb was written for this project in Nextflow,
105 in addition to a functionalized bash script, which pre-downloads
106 and checks software containers of the constituent assembly tools.
107 As a nextflow pipeline was written so that if errors occur in the
108 pipeline, otb can be rerun from that point. otb also was written so
109 that software versions of all tools are exported into a final reports
110 folder (see Supplement for software versions). otb was also written
111 with configs for slurm, local, and sge cluters, and comes with an
112 optional slurm template script to run the pipeline.

113 Results and Discussion

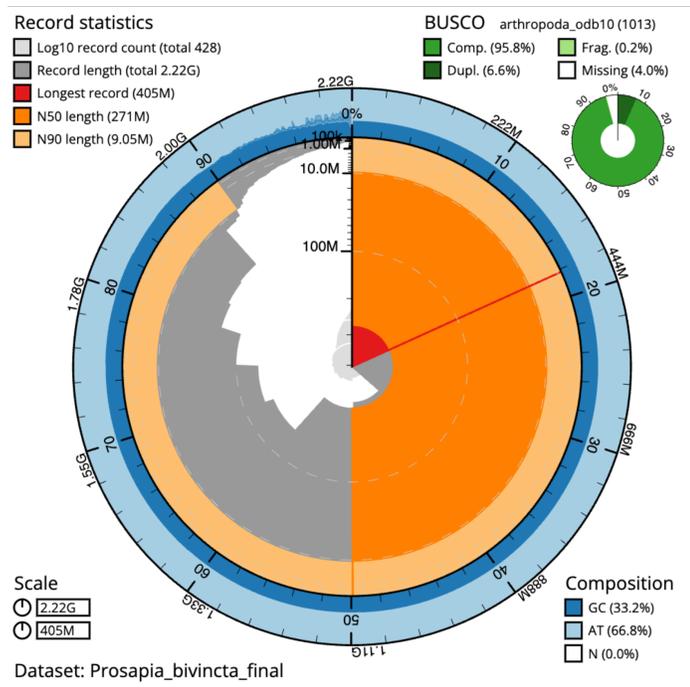
114 The goal of otb is to deliver a genome assembly as close to a
115 polished genome as possible (i.e. reduce manual task time). otb
116 takes several steps in order to reach this point (see Fig. 2, Table
117 1). Starting with a setup, first otb will check its environment and
118 any set environmental variables, as well as called flags and any
119 modes it should be running int, then otb will check if all required



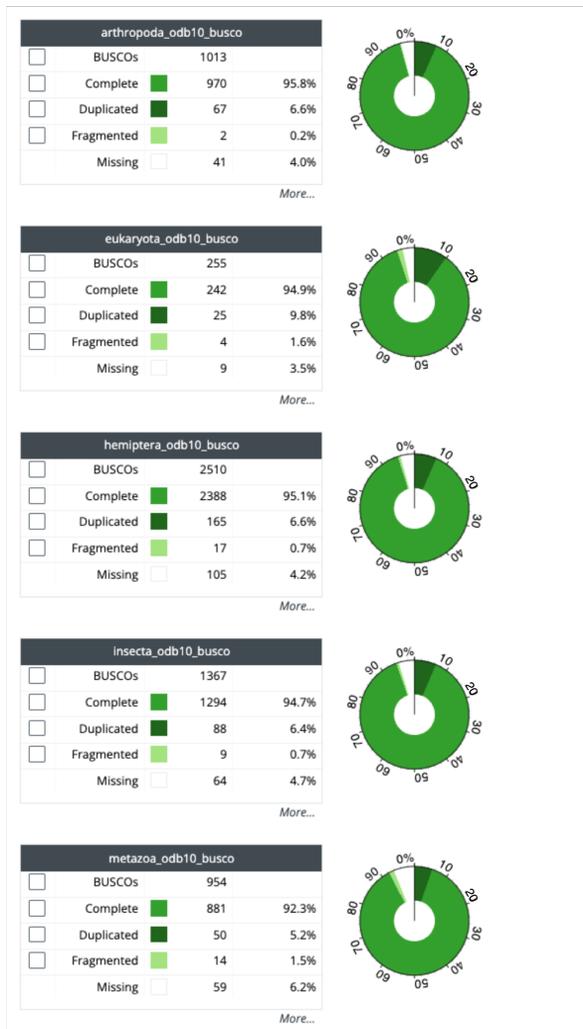
A



B



C



D

Figure 1 Blobtools reports: A: blobplot, showing length and GC content B: snail plot showing record statistics and Busco C: cumulative plots showing length and assignment of reads D: Busco plots showing complete buscos of several taxonomic classifications

Table 1 Software Tools Utilized by otb

Software	References	Latest Title
BamTools	(Barnett <i>et al.</i> 2011)	BamTools: a C++ API and toolkit for analyzing and managing BAM files
BBTools	(Bushnell <i>et al.</i> 2017)	BBMerge – Accurate paired shotgun read merging via overlap
BCFTools	(Li 2011; Danecek <i>et al.</i> 2021)	Twelve years of SAMtools and BCFtools
bwa	(Li and Durbin 2009)	Fast and accurate short read alignment with Burrows–Wheeler transform
fcs-adaptor	(CGR 2022a,b)	Foreign Contamination Screen (FCS) tool for GenBank submissions
GFastats	(Formenti <i>et al.</i> 2022a)	Gfastats: conversion, evaluation and manipulation of genome sequences using assembly graphs
BUSCO	(Manni <i>et al.</i> 2021b,a; Seppey <i>et al.</i> 2019) (Waterhouse <i>et al.</i> 2018, 2017; Simão <i>et al.</i> 2015)	BUSCO Update: Novel and streamlined workflows along with a wider and deeper phylogenetic coverage for the scoring of eukaryotic, prokaryotic, and viral genomes
DeepVariant	(Poplin <i>et al.</i> 2018)	A universal SNP and small-indel variant caller using deep neural networks
GenomeScope2	(Ranallo-Benavidez <i>et al.</i> 2020; Vurture <i>et al.</i> 2017)	GenomeScope 2.0 and Smudgeplot for reference-free profiling of polyploid genomes
hicstuff	(Matthey-Doret <i>et al.</i> 2020, 2021)	koszullab/hicstuff: Use miniconda layer for docker and improve P(s) normalization.
HiFiAdapterFilt	(Sim <i>et al.</i> 2022)	HiFiAdapterFilt, a memory efficient read processing pipeline, prevents occurrence of adapter sequence in PacBio HiFi reads and their negative impacts on genome assembly
hifiasm	(Cheng <i>et al.</i> 2021)	Haplotype-resolved de novo assembly using phased assembly graphs with hifiasm
Jellyfish	(Marçais and Kingsford 2011)	A fast, lock-free approach for efficient parallel counting of occurrences of k-mers
KMC 3	(Kokot <i>et al.</i> 2017)	KMC 3: counting and manipulating k-mer statistics
Merfin	(Formenti <i>et al.</i> 2022b)	Merfin: improved variant filtering, assembly evaluation and polishing via k-mer validation
RagTag	(Alonge <i>et al.</i> 2021, 2019)	Automated assembly scaffolding elevates a new tomato system for high-throughput genome editing
SAMTools	(Li 2011; Danecek <i>et al.</i> 2021)	Twelve years of SAMtools and BCFtools
Shhqis.jl	(Molik 2022)	molikd/Shhqis.jl: Inital Release
VCFTools	(Danecek <i>et al.</i> 2011)	The variant call format and VCFtools
Yahs	(Kokot <i>et al.</i> 2017; Deorowicz <i>et al.</i> 2015, 2013)	KMC 3: counting and manipulating k-mer statistics

containers are available and working, if not, otb will download the required singularity containers and proceed to calling the nextflow run, the main body of the analysis. The sequencing data will be filtered and then assembled with HiFiASM and hicstuff (Cheng *et al.* 2021; Matthey-Doret *et al.* 2020). If the user requests it, Busco can be optionally run at this point. Shhqis.jl, an in-house script is run after the initial assembly to cluster and orient the contigs. Optionally, otb can undergo genome assembly polishing, and it provides this in three ways: Merfin, DeepVariant (Formenti *et al.* 2022b; Poplin *et al.* 2018), or a ‘Simple’. In “Simple” the error corrected reads are used to scaffold, both Merfin and DeepVariant do this as well. Busco can optionally be run at this point. The user can then optionally run KMC/Yahs, and Busco can optionally run at that point as well (Kokot *et al.* 2017; Manni *et al.* 2021a). The result is that the maximum number of assembly steps are carried out, saving the user from having to perform each step individually. By including Hi-C contact map rearranging in the pipeline, and steps to reduce the number of contigs in the draft assembly, less work is needed in the scaffolding step (Molik 2022).

otb was tested with Two-lined spittle bug data. The resultant genome was fairly typical of a true bug. The genome assembly had a N50 of 270.86 Megabases, a Total scaffold length of 2.22 Gigabases, a GC content of 33.22%, had an average scaffold length 5.19 Megabases, had a Total scaffold length 2.22 Gigabases, had 33.22% GC Content, and had a N50 270.86 Megabases. Blobtools hemiptera showed a 95.1% Buscos found (see Fig. 1. The genome assembly of the two-lined spittle bug will provide a basis for further work into its interaction with its obligate hosts, and no doubt into the control of the pest.

Nextflow and Singularity are two software infrastructure tools commonly used in scientific computing workflows. Nextflow is an open-source software solution that enables users to define pipelines and processes for data analysis in a concise and reproducible manner, supporting various languages and technologies.

This versatile platform is built on a reactive programming model, allowing dynamic, parallel processing, and data flow management. This makes Nextflow an ideal tool for managing large-scale, multi-step data analysis processes in a scalable and efficient manner.

Singularity is a container technology specifically designed for high-performance and scientific computing environments. This technology provides a way to bundle applications, dependencies, and the environment into a single executable package, ensuring reproducibility, compatibility with existing Linux container technologies, and security. Singularity offers several key features, such as its ability to run containers with root-level access on any system while ensuring that the host system remains isolated and protected. These features make Singularity an attractive solution for many scientific and engineering workflows.

Together, Nextflow and Singularity provide a powerful and versatile solution for managing complex scientific computing workflows. They enable users to share and run applications, tools, and workflows in a consistent and reproducible manner, across different HPC environments and operating systems. As such, they have become essential tools for managing the large-scale and complex data analysis processes required in modern scientific research.

The introduction of otb, a new HiC/HiFi phased genomics assembly pipeline, has solved several data management problems that were prevalent in the creation and storage of genome assemblies. otb is written in nextflow and utilizes singularity to ensure uniformity of the computing environment. Offers parallel task execution and resource management, while also reducing the time spent organizing data, installing tools, and performing analysis. With the implementation of otb, genome creation can be automated, especially in regards to projects such as the Ag100pest.

Data availability

Code used in the creation of this genome is in the public domain per United States 17 U.S.C. § 105. The code is freely available for

1 use and modification:

	github	DOI
2	otb molikd/otb	10.5281/zenodo.6689816
	Shhquis.jl molikd/Shhquis.jl	10.5281/zenodo.6315237

3 Data published for this article are in the public domain per
4 United States 17 U.S.C. § 105. The data are freely available for use
5 and modification:

	Identifier
6	BioProject PRJNA987615
	BioSample SAMN35984262

7 Acknowledgments

8 The McDonnell Genome Institute (MGI) provided several docker
9 containers used by otb, especially through the works of contribu-
10 tors to these containers: John Garza, Thomas B. Mooney, and
11 Alex Paul. Special thanks to Paul Stodghill of the USDA ARS in
12 Ithaca for providing another container used by otb. The authors
13 also thank Coffee (*Coffea arabica*) for always being there for them
14 and Tina Turner for providing the music that made up the sound
15 track during the writing of this manuscript.

16 Funding

17 David Molik was supported through this project by the USDA
18 Agricultural Research Service (ARS) HQ Research Associate pro-
19 gram in Big Data. This work was carried out in part by the Tropical
20 Pest Genetics and Molecular Biology Research Unit, ARS Project
21 number 2040-22430-028-000D. This research used resources pro-
22 vided by the SCINet project of the USDA Agricultural Research
23 Service, ARS project number 0500-00093-001-00-D.

24 Conflicts of interest

25 none declared.

26 Footnotes

27 The U.S. Department of Agriculture is an equal opportunity lender,
28 provider, and employer.

29 Mention of trade names or commercial products in this report is
30 solely for the purpose of providing specific information and does
31 not imply recommendation or endorsement by the U.S. Depart-
32 ment of Agriculture.

33 Literature cited

34 Alonge M, Lebeigle L, Kirsche M, Aganezov S, Wang X, Lippman
35 ZB, Schatz MC, Soyk S. 2021. Automated assembly scaffold-
36 ing elevates a new tomato system for high-throughput genome
37 editing. *bioRxiv*. doi: 10.1101/2021.11.18.469135, pre-print: not
38 peer-reviewed.
39 Alonge M, Soyk S, Ramakrishnan S, Wang X, Goodwin S, Sedlazeck
40 FJ, Lippman ZB, Schatz MC. 2019. RaGOO: fast and accurate
41 reference-guided scaffolding of draft genomes. *Genome Biology*.
42 20. doi: 10.1186/s13059-019-1829-6.
43 Ankrah NY, Wilkes RA, Zhang FQ, Zhu D, Kaweesi T, Aristilde
44 L, Douglas AE. 2020. Syntrophic splitting of central carbon
45 metabolism in host cells bearing functionally different symbiotic
46 bacteria. *The ISME journal*. 14:1982–1993.

Barnett DW, Garrison EK, Quinlan AR, Strömberg MP, Marth GT. 47
2011. Bamtools: a c++ api and toolkit for analyzing and manag- 48
ing bam files. *Bioinformatics*. 27:1691–1692. doi: 10.1093/bioin- 49
formatics/btr174. 50

Braman S, Hanna W, Schwartz B, Nair S. 2014. Evaluation of chi- 51
nese centipedegrasses and other turfgrass taxa for potential re- 52
sistance to twolined spittlebug, *prosapia bicincta* (say). *Journal* 53
of entomological science. 49:121–129. 54

Braman SK, Abraham CM. 1995. Twolined spittlebug. *Handbook* 55
of turfgrass insect pests. Entomological Society of America, Lan- 56
ham, MD. pp. 88–90. 57

Bremer LL, Nathan N, Trauernicht C, Pascua P, Krueger N, Jokiel J, 58
Barton J, Daily GC. 2021. Maintaining the many societal benefits 59
of rangelands: The case of hawai'i. *Land*. 10:764. 60

Bushnell B, Rood J, Singer E. 2017. Bbmerge—accurate paired 61
shotgun read merging via overlap. *PloS one*. 12:e0185056. doi: 62
10.1371/journal.pone.0185056. 63

Byers R, Wells HD. 1966. Phytotoxemia of coastal bermuda- 64
grass caused by the two-lined spittlebug, *prosapia bicincta* (ho- 65
moptera: Cercopidae). *Annals of the Entomological Society of* 66
America. 59:1067–1071. 67

CGR NCGR. 2022a. Fcs. <https://github.com/ncbi/fcs>. 68

CGR NCGR. 2022b. Foreign contamination screen (fcs) tool for 69
genbank submissions. Technical report. NCBI Insights. Bethesda, 70
MD. Accessed March 2023. [https://ncbiinsights.ncbi.nlm.nih.gov/](https://ncbiinsights.ncbi.nlm.nih.gov/2022/07/28/fcs-beta-tool/) 71
[2022/07/28/fcs-beta-tool/](https://ncbiinsights.ncbi.nlm.nih.gov/2022/07/28/fcs-beta-tool/). 72

Cheng H, Concepcion GT, Feng X, Zhang H, Li H. 2021. Haplotype- 73
resolved de novo assembly using phased assembly graphs with 74
hifiasm. *Nature Methods*. 18:170–175. doi: 10.1038/s41592-020- 75
01056-5. 76

Childers AK, Geib SM, Sim SB, Poelchau MF, Coates BS, Simmonds 77
TJ, Scully ED, Smith TP, Childers CP, Corpuz RL *et al.* 2021. The 78
usda-ars ag100pest initiative: high-quality genome assemblies 79
for agricultural pest arthropod research. *Insects*. 12:626. 80

Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, 81
Handsaker RE, Lunter G, Marth GT, Sherry ST *et al.* 2011. The 82
variant call format and VCFtools. *Bioinformatics*. 27:2156–2158. 83
doi: 10.1093/bioinformatics/btr330. 84

Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, 85
Whitwham A, Keane T, McCarthy SA, Davies RM *et al.* 2021. 86
Twelve years of SAMtools and BCFtools. *GigaScience*. 10. doi: 87
10.1093/gigascience/giab008. 88

Deorowicz S, Debudaj-Grabysz A, Grabowski S. 2013. Disk- 89
based k-mer counting on a PC. *BMC Bioinformatics*. 14. doi: 90
10.1186/1471-2105-14-160. 91

Deorowicz S, Kokot M, Grabowski S, Debudaj-Grabysz A. 2015. 92
KMC 2: fast and resource-frugal k-mer counting. *Bioinformatics*. 93
31:1569–1576. doi: 10.1093/bioinformatics/btv022. 94

Fagan EB, Kuitert L. 1969. Biology of the two-lined spittlebug, 95
prosapia bicincta, on florida pastures (homoptera: Cercopidae). 96
Florida Entomologist. pp. 199–206. 97

Formenti G, Abueg L, Brajuka A, Brajuka N, Gallardo-Alba C, 98
Giani A, Fedrigo O, Jarvis ED. 2022a. Gfastats: conversion, eval- 99
uation and manipulation of genome sequences using assembly 100
graphs. *Bioinformatics*. 38:4214–4216. 101

Formenti G, Rhie A, Walenz BP, Thibaud-Nissen F, Shafin K, Koren 102
S, Myers EW, Jarvis ED, Phillippy AM. 2022b. Merfin: improved 103
variant filtering, assembly evaluation and polishing via k-mer 104
validation. *Nature Methods*. 19:696–704. doi: 10.1038/s41592- 105
022-01445-y. 106

Jiang T, Yin Z, Cai R, Yu H, Lu Q, Zhao S, Tian Y, Yan Y, Guo J, Chen 107
X. 2021. Chromosomal-Level Genome Assembly of a True Bug, 108

- 1 Aspongopus chinensis Dallas, 1851 (Hemiptera: Dinidoridae).
2 Genome Biology and Evolution. 13:evab232.
- 3 Joseph SV, Jespersen D. 2021. Influence of relative humidity on the
4 expression of twolined spittlebug (hemiptera: Cercopidae) feed-
5 ing injury in turfgrass genotypes. Arthropod-Plant Interactions.
6 15:197–207.
- 7 Koga R, Bennett GM, Cryan JR, Moran NA. 2013. Evolutionary
8 replacement of obligate symbionts in an ancient and diverse
9 insect lineage. Environmental microbiology. 15:2073–2081.
- 10 Koga R, Moran NA. 2014. Swapping symbionts in spittlebugs:
11 evolutionary replacement of a reduced genome symbiont. The
12 ISME journal. 8:1237–1246.
- 13 Kokot M, Długosz M, Deorowicz S. 2017. KMC 3: counting and
14 manipulating k-mer statistics. Bioinformatics. 33:2759–2761. doi:
15 10.1093/bioinformatics/btx304.
- 16 Leftwich PT, Bolton M, Chapman T. 2015. Evolutionary biology
17 and genetic techniques for insect control. Evolutionary Applica-
18 tions. 9:212–230. doi: 10.1111/eva.12280.
- 19 Li H. 2011. A statistical framework for SNP calling, mutation dis-
20 covery, association mapping and population genetical parameter
21 estimation from sequencing data. Bioinformatics. 27:2987–2993.
22 doi: 10.1093/bioinformatics/btr509.
- 23 Li H, Durbin R. 2009. Fast and accurate short read alignment with
24 Burrows–Wheeler transform. Bioinformatics. 25:1754–1760.
- 25 Manni M, Berkeley MR, Seppey M, Simão FA, Zdobnov EM. 2021a.
26 BUSCO update: Novel and streamlined workflows along with
27 broader and deeper phylogenetic coverage for scoring of eu-
28 karyotic, prokaryotic, and viral genomes. Molecular Biology and
29 Evolution. 38:4647–4654. doi: 10.1093/molbev/msab199.
- 30 Manni M, Berkeley MR, Seppey M, Zdobnov EM. 2021b. BUSCO:
31 Assessing genomic data quality and beyond. Current Protocols.
32 1. doi: 10.1002/cpz1.323.
- 33 Marçais G, Kingsford C. 2011. A fast, lock-free approach for effi-
34 cient parallel counting of occurrences of k-mers. Bioinformatics.
35 27:764–770. doi: 10.1093/bioinformatics/btr011.
- 36 Matthey-Doret C, Baudry L, Bignaud A, Cournac A, Remi-
37 Montagne, Guiglielmoni N, Foutel-Rodier T, Scolari VF. 2020.
38 hicstuff: Simple library/pipeline to generate and handle hi-c
39 data. doi: 10.5281/zenodo.4066351.
- 40 Matthey-Doret C, Baudry L, Mortaza S, Moreau P, Koszul R,
41 Cournac A. 2021. Normalization of chromosome contact maps:
42 Matrix balancing and visualization, In: , Springer US. pp. 1–15.
43 doi: 10.1007/978-1-0716-1390-0_1.
- 44 Mattson Jr WJ. 1980. Herbivory in relation to plant nitrogen content.
45 Annual review of ecology and systematics. 11:119–161.
- 46 Molik D. 2022. molikd/shhquis.jl: Inital release. doi: 10.5281/ZEN-
47 ODO.6315238.
- 48 Nachappa P, Guillebeau L, Braman S, All J. 2006. Susceptibility
49 of twolined spittlebug (hemiptera: Cercopidae) life stages to
50 entomophagous arthropods in turfgrass. Journal of economic
51 entomology. 99:1711–1716.
- 52 Peck DC. 1998. Natural history of the spittlebug prosapia nr.
53 bicincta (homoptera: Cercopidae) in association with dairy pas-
54 tures of costa rica. Annals of the Entomological Society of Amer-
55 ica. 91:435–444.
- 56 Poplin R, Chang PC, Alexander D, Schwartz S, Colthurst T, Ku A,
57 Newburger D, Dijamco J, Nguyen N, Afshar PT *et al.* 2018. A uni-
58 versal SNP and small-indel variant caller using deep neural net-
59 works. Nature Biotechnology. 36:983–987. doi: 10.1038/nbt.4235.
- 60 Potter DA, Braman SK *et al.* 1991. Ecology and management of
61 turfgrass insects. Annual Review of Entomology. 36:383–406.
- 62 Ranallo-Benavidez TR, Jaron KS, Schatz MC. 2020. GenomeScope
2.0 and smudgeplot for reference-free profiling of polyploid
63 genomes. Nature Communications. 11. doi: 10.1038/s41467-020-
64 14998-3.
- 65 Robinson JT, Turner D, Durand NC, Thorvaldsdóttir H, Mesirov JP,
66 Aiden EL. 2018. Juicebox.js provides a cloud-based visualization
67 system for hi-c data. Cell Systems. 6:256–258.e1.
- 68 Seppey M, Manni M, Zdobnov EM. 2019. BUSCO: Assessing
69 genome assembly and annotation completeness, In: , Springer
70 New York. pp. 227–245. doi: 10.1007/978-1-4939-9173-0_14.
- 71 Sim SB, Corpuz RL, Simmonds TJ, Geib SM. 2022. Hifiadapterfilt,
72 a memory efficient read processing pipeline, prevents occur-
73 rence of adapter sequence in pacbio hifi reads and their nega-
74 tive impacts on genome assembly. BMC Genomics. 23:157. doi:
75 10.1186/s12864-022-08375-.
- 76 Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdob-
77 nov EM. 2015. BUSCO: assessing genome assembly and anno-
78 tation completeness with single-copy orthologs. Bioinformatics.
79 31:3210–3212. doi: 10.1093/bioinformatics/btv351.
- 80 Snyder MW, Adey A, Kitzman JO, Shendure J. 2015. Haplotype-
81 resolved genome sequencing: experimental methods and
82 applications. Nature Reviews Genetics. 16:344–358. doi:
83 10.1038/nrg3903.
- 84 Tewhey R, Bansal V, Torkamani A, Topol EJ, Schork NJ. 2011. The
85 importance of phase information for human genomics. Nature
86 Reviews Genetics. 12:215–223. doi: 10.1038/nrg2950.
- 87 Thorne M, Fukumoto G, Curtiss R, Hamasaki R. 2017. New spittle-
88 bug on pasture grasses in hawai’i two-lined spittlebug, prosapia
89 bicincta. Research communication. College of Tropical Agricul-
90 ture, University of Hawai’i at Manoa: Honolulu, HI, USA. Hon-
91 olulu, HI.
- 92 Tommaso PD, Chatzou M, Floden EW, Barja PP, Palumbo E,
93 Notredame C. 2017. Nextflow enables reproducible compu-
94 tational workflows. Nature Biotechnology. 35:316–319. doi:
95 10.1038/nbt.3820.
- 96 Vurture GW, Sedlazeck FJ, Nattestad M, Underwood CJ, Fang H,
97 Gurtowski J, Schatz MC. 2017. GenomeScope: fast reference-free
98 genome profiling from short reads. Bioinformatics. 33:2202–2204.
99 doi: 10.1093/bioinformatics/btx153.
- 100 Waterhouse RM, Seppey M, Simão FA, Manni M, Ioannidis P,
101 Klioutchnikov G, Kriventseva EV, Zdobnov EM. 2017. BUSCO
102 applications from quality assessments to gene prediction and
103 phylogenomics. Molecular Biology and Evolution. 35:543–548.
104 doi: 10.1093/molbev/msx319.
- 105 Waterhouse RM, Seppey M, Simão FA, Zdobnov EM. 2018. Using
106 BUSCO to assess insect genomic resources, In: , Springer New
107 York. pp. 59–74. doi: 10.1007/978-1-4939-8775-7_6.
- 108

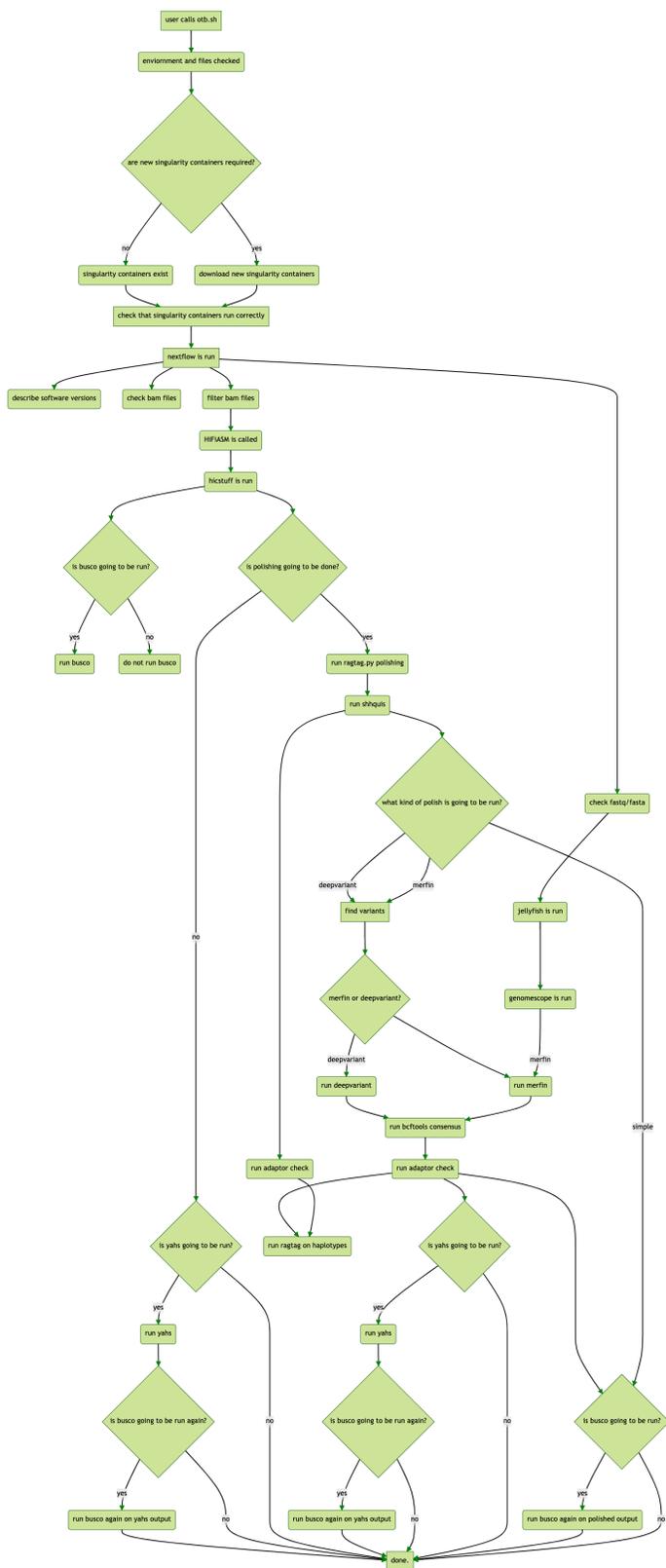


Figure 2 otb flowchart. Workflow diagram of otb showing the process of otb running, otb.sh the entry point for otb run software and container checks, followed by the assessment of the type of hifiasm assembly to be created, otb also allows for multiple types of a sequence based polishing run, including a "simple" or reuse of error correct reads remapped using ragtag.py, a "deep variant" which uses deepvariant, and "merfin". Busco and sequence stats are run at multiple points in the pipeline. Yahs is run to produce HiC maps which can be utilized in JuiceBox [Robinson et al. \(2018\)](#).