



HAL
open science

Efficient k-mer based curation of raw sequence data: application in *Drosophila suzukii*

Mathieu Gautier

► **To cite this version:**

Mathieu Gautier. Efficient k-mer based curation of raw sequence data: application in *Drosophila suzukii*. Peer Community Journal, 2023, 3, pp.e79. 10.24072/pcjournal.309 . hal-04667623

HAL Id: hal-04667623

<https://hal.inrae.fr/hal-04667623v1>

Submitted on 5 Aug 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License



Peer Community Journal

Section: Genomics

RESEARCH ARTICLE

Published
2023-09-01

Cite as
Mathieu Gautier (2023)
*Efficient k-mer based curation of
raw sequence data: application
in Drosophila suzukii*, Peer
Community Journal, 3: e79.

Correspondence
mathieu.gautier@inrae.fr

Peer-review
Peer reviewed and
recommended by
PCI Genomics,
[https://doi.org/10.24072/pci.
genomics.100244](https://doi.org/10.24072/pci.genomics.100244)



This article is licensed
under the Creative Commons
Attribution 4.0 License.

Efficient *k-mer* based curation of raw sequence data: application in *Drosophila suzukii*

Mathieu Gautier ¹

Volume 3 (2023), article e79

<https://doi.org/10.24072/pcjournal.309>

Abstract

Several studies have highlighted the presence of contaminated entries in public sequence repositories, calling for special attention to the associated metadata. Here, we propose and evaluate a fast and efficient *k-mer*-based approach to assess the degree of mislabeling or contamination. We applied it to high-throughput whole-genome raw sequence data for 236 Ind-Seq and 22 Pool-Seq samples of the invasive species *Drosophila suzukii*. We first used Clark software to build a dictionary of species-discriminating *k-mers* from the curated assemblies of 29 target drosophilid species (including *D. melanogaster*, *D. simulans*, *D. subpulchrella*, or *D. biarmipes*) and 12 common drosophila pathogens and commensals (including Wolbachia). Counting the number of *k-mers* composing each query sample sequence that matched a discriminating *k-mer* from the dictionary provided a simple criterion for assignment to target species and evaluation of the entire sample. Analyses of a wide range of samples, representative of both target and other drosophilid species, demonstrated very good performance of the proposed approach, both in terms of run time and accuracy of sequence assignment. Of the 236 *D. suzukii* individuals, five were reassigned to *D. simulans* and eleven to *D. subpulchrella*. Another four showed moderate to substantial microbial contamination. Similarly, among the 22 Pool-Seq samples analyzed, two from the native range were found to be contaminated with 1 and 7 *D. subpulchrella* individuals, respectively (out of 50), and one from Europe was found to be contaminated with 5 to 6 *D. immigrans* individuals (out of 100). Overall, the present analysis allowed the definition of a large curated dataset consisting of > 60 population samples representative of the worldwide genetic diversity, which may be valuable for further population genetics studies on *D. suzukii*. More generally, while we advocate careful sample identification and verification prior to sequencing, the proposed framework is simple and computationally efficient enough to be included as a routine post-hoc quality check prior to any data analysis and prior to data submission to public repositories.

¹CBGP, INRAE, CIRAD, IRD, Montpellier SupAgro, Université de Montpellier, Montpellier, France

Contents

1	Introduction	2
2	Material and Methods	3
2.1	Construction of the Clark and Clark-I target dictionaries of species-discriminating <i>k</i> -mers	3
2.2	Query short-read sequencing data	5
2.3	Assignment of query sequences and contamination estimation	7
3	Results	7
3.1	Clark and Clark-I run times	7
3.2	Proportion of assigned sequences	8
3.3	Assignment accuracy for samples representative of target and other species	10
3.4	Scanning 236 Ind-Seq and 22 Pool-Seq <i>D. suzukii</i> WGS data	11
4	Discussion	14
	Acknowledgements	16
	Fundings	16
	Conflict of interest disclosure	16
	Data, script, code, and supplementary information availability	16
	References	17

1. Introduction

With the democratization of sequencing technologies, the availability of genomic sequence in public repositories is increasing at an unprecedented rate. This is enabling the construction of large and highly informative combined datasets for an increasing number of model and non-model species, which in turn is refining the power and resolution of population genomics inference (e.g. Kapun et al., 2021). However, this increased availability of data comes at the cost of increased heterogeneity in the resulting combined dataset. For example, data sets may combine different sequencing library preparation protocols or technologies that are rapidly evolving with variable sequence quality or coverage. Similarly, for a given species, publicly available data may refer to original studies based on different sampling strategies consisting of either sequencing individuals (aka Ind-Seq) or pools of individuals (aka Pool-Seq) representative of some populations, the latter approach being quite popular due to its cost-effectiveness (Schlötterer et al., 2014). Nevertheless, such technical characteristics can be taken into account in downstream analyses if an appropriate statistical framework is used.

More problematically, several recent studies have highlighted the high level of contamination in public repositories, which requires special attention when relying on the associated metadata description files (Cornet and Baurain, 2022; Francois et al., 2020; Steinegger and Salzberg, 2020). For example, working with wild-caught samples of species that are difficult to distinguish from other closely related species sharing the same habitat may lead to taxonomic errors or biological contamination of the sample. Such potential problems have already been reported in population genetic studies of *Drosophila melanogaster*, where sample contamination with *D. simulans* individuals was not uncommon (Kapun et al., 2021; Machado et al., 2021). In addition to biological sources, contamination may be of experimental (e.g., sample contamination or mislabeling) and/or computational (e.g., during data processing) origin (Cornet and Baurain, 2022). It should also be noted that these contamination problems are obviously not specific to publicly available data and may be even more pronounced in newly generated data that have not yet been analyzed.

In recent years, several software packages have been developed to assess the level of contamination in genomic datasets, which has been greatly facilitated by the active field of metagenomics. As recently reviewed by Cornet and Baurain (2022), the available approaches can be classified into either database-free or reference-based methods. Database-free methods roughly

consist of partitioning sequences based on their DNA composition (e.g. GC content or frequencies in short DNA sequences of a few nt), but they are not well suited for the analysis of large amounts of samples as they require a case-by-case inspection of the results (Cornet and Baurain, 2022). Reference-based methods consist of aligning sequences to a set of tagged sequences representative of all or part (e.g., genes) of the genomes of candidate species. In practice, this may allow either negative and/or positive filtering (i.e., removal of contaminating sequences or identification of sequences from some species of interest) of the sequencing data (Cornet and Baurain, 2022). To accomplish this task, approaches based on the exact matching of *k*-mers (i.e., *k* nt long DNA words) constituting the query sequences to a dictionary of labeled *k*-mers (built from target species genomes) have proven highly efficient and are now very popular for sequence taxonomic classification in the metagenomic field (Ounit and Lonardi, 2016; Ounit et al., 2015; Wood et al., 2019; Wood and Salzberg, 2014).

Taking advantage of the high quality assemblies available for several dozen drosophilid genomes (Kim et al., 2021), the aim of this study was to rely on a *k*-mer-based approach to assess the level of contamination in public sequence data for the spotted wing *Drosophila D. suzukii*. *D. suzukii* originates from Asia and has recently invaded the entire European and American continents to become a major invasive insect pest causing dramatic losses in fruit production (Asplen et al., 2015; Cini et al., 2012). This species has thus become of great scientific interest, particularly to population geneticists, and several recent studies have provided informative samples for characterizing the structuring of its genetic diversity at global and whole-genome scales. Here, we focused on two recently published and publicly available Pool-Seq and Ind-Seq datasets, consisting of whole-genome sequences (WGS) for i) 22 pools of individual DNA (with $n=50$ to $n=100$ individuals per pool) representative of populations sampled both in the Asian native range ($n=6$) and in the European ($n=8$) and American invaded ranges ($n=8$) (Olazcuaga et al., 2020); and ii) 236 individuals collected mainly in North America but also at several sites in Europe, Brazil and Asia (Lewald et al., 2021). A combined analysis of these two datasets using standard descriptive approaches revealed anomalous behavior of some samples (not shown), thus motivating a systematic screening of all samples for putative contamination or (taxonomic) misidentification problems. Indeed, as highlighted by Piper et al. (2022), rapid morphological identification of *D. suzukii* on wild-caught specimens can be tricky. For example, the distinctive spots observed on the wing extremities are present only in (non-juvenile) males, and this feature is shared with two of its sister species *D. biarmipes* and *D. subpulchrella*, whose distributions overlap all or part of that of *D. suzukii* in its native Asian range (Ometto et al., 2013; Takamori et al., 2006). In addition, both *D. suzukii* and *D. subpulchrella* females possess a large and serrated ovipositor that allows them to penetrate under the skin of ripening fruits and lay eggs (Atallah et al., 2014), making the distinction between these two species even more difficult.

To assess contamination in publicly available *D. suzukii* raw sequencing data, we developed and evaluated a fast and efficient approach based on *k*-mer-based methods implemented in the software Clark (Ounit et al., 2015). We first build dictionaries of species-discriminating *k*-mers from the curated assemblies of 29 target drosophila species and 12 common drosophila pathogens and commensals. WGS data for individual samples representative of both the target and other drosophilid species were then analyzed to evaluate the performance of the proposed approaches, both in terms of run time and accuracy of sequence assignment. Finally, we analyzed publicly available WGS data for the aforementioned 236 Ind-Seq (Lewald et al., 2021) and 22 Pool-Seq (Olazcuaga et al., 2020) samples of the invasive species *Drosophila suzukii*, allowing us to identify unambiguously contaminated samples.

2. Material and Methods

2.1. Construction of the Clark and Clark-I target dictionaries of species-discriminating *k*-mers

Of the 136 reference genome assemblies available for species belonging to the genus *Drosophila* in the NCBI repository (<https://www.ncbi.nlm.nih.gov/datasets/genomes/> accessed in February 2022), 29 were retained based on assembly quality criteria such as contiguity (evaluated with contig N50) and completeness (using BUSCO scores, Manni et al., 2021); but also and mostly

based on phylogenetic criteria (Figure 1). Our goal was to obtain a good representation of species closely related to *D. suzukii*, focusing on those belonging to the two subgenera *Sophophora* and *Drosophila* that are not unambiguously resolved (see Discussion). For subgroups or groups represented by multiple species (among those with good quality assemblies available), only one target species was selected, favoring the most cosmopolitan or temperate species (Jezovit et al., 2017), except for the species most closely related to (and likely to be confounded with) *D. suzukii* (e.g., *D. subpulchrella* and *D. biarmipes*). To further improve the representation of *D. suzukii* in the *k*-mer dictionary, the draft assembly of Ometto et al. (2013) was also downloaded from the ENA repository (<https://www.ebi.ac.uk/ena/browser/home>). Although this assembly was of lower quality than the reference (Paris et al., 2020), it was obtained from a different isofemale line and was based on short read sequences from a pool of females and males. Similarly, for *D. subpulchrella* (the sister species of *D. suzukii*), the assembly from (Kim et al., 2021) was considered in addition to the latest NCBI reference assembly (Durkin et al., 2021), because it is based on male individuals and therefore contain Y-linked contigs. The high quality *D. simulans de novo* assembly from Chang et al. (2022) was also included for similar reasons.

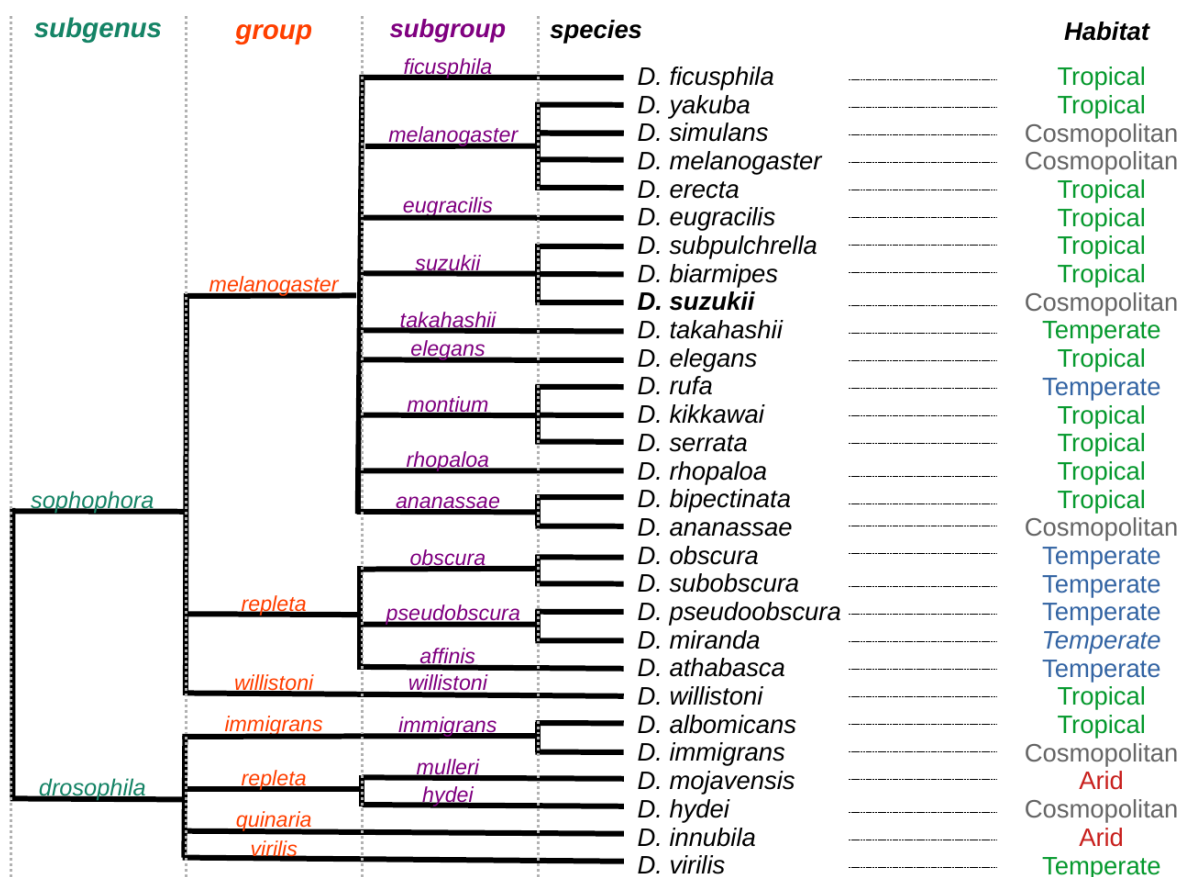


Figure 1 – Relationship between the 29 target drosophilid species (adapted from <https://www.ncbi.nlm.nih.gov/taxonomy>). Species habitat was defined according to Jezovit et al. (2017), except for *D. miranda*. *D. suzukii* is highlighted in bold.

The resulting 32 assemblies, described in Table 1, were further screened for non-*Drosophila* contaminating sequences using the program Kraken2 v2.1.2 (Wood et al., 2019) by querying a database constructed from the NCBI non-redundant nucleotides (nt) released in February 2020. A contig or scaffold sequence was considered contaminating if it was assigned to a taxonomic identifier unrelated to any drosophilid species. Note that contigs assigned to *Wolbachia* endosymbionts were also flagged as contaminating, as we chose to consider *Wolbachia* specifically here (see below). Of the 90,071 sequences (i.e., contigs or scaffolds) from all 32 assemblies (5.96 Gb in total), 16,123 sequences (17.9%) were found to be contaminating. As detailed in Table S1, these contaminating sequences were mostly short, ranging from 110 bp to 1,478,327

bp (median size of 1,522 bp), totaling only 102.7 Mb (i.e. 1.72% of all sequences). It should be noted that Wolbachia-related sequences represented only 6,173,139 bp of the contaminating sequences (6.01%), with the major contributor being the *D. ananassae* assembly (6,078,940 bp), which may be explained by the widespread lateral gene transfer from Wolbachia described in this species (Klasson et al., 2014). The other Wolbachia contaminating sequences belonged to the assemblies for *D. suzukii* (83,189 bp) of Ometto et al. (2013) and *D. willistoni* (11,010 bp). Finally, out of the 32 assemblies, only three (for *D. albomicans*, *D. innubila* and *D. melanogaster* species) were found to be free of any contaminant, the most contaminated assemblies being those for *D. immigrans* and *D. willistoni* with 10.6% and 8.1% of their total length contaminated, respectively (Table 1). As expected, the completeness of the assemblies (except for the draft assembly for *D. suzukii* from Ometto et al., 2013 mentioned above) remained quite good after filtering out contaminating sequences, with more than 98% (resp. 95%) of the 3,285 BUSCO genes of the *diptera_odb10* dataset (Manni et al., 2021) identified in 26 (31) of the assemblies (Table 1). Finally, in addition to the drosophilid species and following Kapun et al. (2021), 13 genome assemblies representing twelve different common drosophilid commensals and pathogens were included in the construction of the *k-mer* dictionaries (Table 1). Note that the two reference assemblies for the Wolbachia endosymbiont of *D. melanogaster* and *D. simulans* were used to represent Wolbachia.

From the 45 reference assemblies representing the 29 drosophilid, commensal, and pathogen species, two different dictionaries of species-discriminating *k-mers* (i.e., *k-mers* that occur exclusively in the genome of a species represented by one or more assemblies) were then constructed using versions 1.2.6.1 of Clark (default $k=31$) and Clark-I (default $k=27$), respectively (Ounit et al., 2015). Clark-I is a variant of Clark designed for use when the amount of RAM is limited, with minimal impact on assignment accuracy. Both Clark and Clark-I were run in single-threaded mode on a computer cluster grid. Building the *k-mer* dictionary (on a single thread of a cluster node equipped with a processor Intel® Xeon® CPU E5-2683 v4 @2.10GHz) took 2h46min with a peak RAM usage of 128G using Clark and 55s with a peak RAM usage of 2.65G using Clark-I. The resulting database consisted of 3,714,249,662 31-mers and 50,311,519 27-mers, respectively, and required 47.8 Gb and 1.97 Gb of RAM to load when computing the query sequence classification with Clark and Clark-I, respectively.

2.2. Query short-read sequencing data

A total of 301 short read WGS data sets were downloaded from the public SRA repository (<https://www.ncbi.nlm.nih.gov/sra>). These include 43 samples used for the empirical evaluation of *k-mer*-based assignment accuracy, derived from the sequencing of laboratory strains representative of different drosophilid species (including data on 12 of the 29 target species available for the strains used to generate the corresponding assemblies) and the Wolbachia endosymbiont of *D. melanogaster* (Table S2). As detailed in Table S2, all of these data were obtained from paired-end (PE) sequencing (2×150 nt) on an Illumina HiSeq 4000 instrument, except for ten samples sequenced on an Illumina i) GAIIX in PE125 mode ($n=1$); ii) NextSeq550 in PE150 mode ($n=4$); iii) HiSeq 2000 in PE100 ($n=2$) and PE150 ($n=1$) modes; iv) MiSeq in PE300 mode ($n=1$); or v) HiSeq Ten X in PE150 ($n=1$). The second type of data corresponded to WGS data for 236 *D. suzukii* individuals (Ind-Seq data) representative of 40 population samples (4-10 ind. per sample, mean=5.9) published by Lewald et al. (2021). These samples were mainly collected in the continental USA ($n=31$). The other regions represented are Brazil ($n=1$); Europe ($n=2$; Ireland and Italy) for two of them; China ($n=2$); South Korea ($n=2$), but also Japan ($n=1$) and Hawaii ($n=1$), via two laboratory strains. These were all sequenced on an Illumina HiSeq4000 in PE150 ($n=201$) or PE100 ($n=35$) mode (Table S3). The last type of data corresponded to WGS data from 22 pools of *D. suzukii* individuals (Pool-Seq data) representing 22 worldwide populations representative of the Asian native range ($n=6$) and the European ($n=8$) and American ($n=8$) invaded ranges, published by Olazuaga et al. (2020). These were all sequenced on an Illumina HiSeq2500 in PE125 mode (Table S3).

Raw PE reads were filtered with *fastp* 0.23.1 (Chen et al., 2018) with the default options to remove contaminating adaptor sequences and trimmed for poor quality bases (i.e. with a phred

Table 1 – Description of the reference genome assemblies for the 29 drosophilid species (n=32 assemblies) and 12 common commensals and pathogens (n=13 assemblies) used to build the target *k*-mer dictionaries. All genome assemblies were downloaded from the NCBI repository (<https://www.ncbi.nlm.nih.gov>), except for the two additional assemblies for *D. simulans* and *D. suzukii*, which were downloaded from the Dryad (<https://datadryad.org>) and ENA (<https://www.ebi.ac.uk/ena>) repositories, respectively (with accession ID in italics in the third column). The size and N50 of all assemblies are given in the fourth and fifth columns. For drosophilid species, these correspond to the assemblies after filtering out the identified contaminant contigs (or scaffolds), the percentage of the original assembly retained is given in parentheses in the fourth column. Similarly, the BUSCO scores in parentheses correspond to the percentage of complete genes identified among the 3,285 genes of the *diptera_odb10* dataset (Manni et al., 2021).

ID	Species	Reference	Size in Mb (% init.)	N50 in Mb (BUSCO)
Dalbo	<i>Drosophila albomicans</i>	GCA_009650485.1	165.85 (100.0)	33.43 (96.6)
Danan	<i>Drosophila ananassae</i>	GCA_017639315.1	207.74 (97.16)	26.43 (99.1)
Datha	<i>Drosophila athabasca</i>	GCA_008121215.1	191.06 (99.17)	52.10 (98.3)
Dbiar	<i>Drosophila biarmipes</i>	GCA_018148935.1	183.51 (99.03)	23.38 (98.9)
Dbipe	<i>Drosophila bipectinata</i>	GCA_018153845.1	189.91 (98.71)	15.79 (99.1)
Deleg	<i>Drosophila elegans</i>	GCA_018152505.1	177.57 (99.51)	21.93 (99.1)
Derec	<i>Drosophila erecta</i>	GCA_003286155.1	146.49 (99.97)	22.15 (99.2)
Deugr	<i>Drosophila eugracilis</i>	GCA_018153835.1	158.76 (96.33)	2.299 (98.5)
Dficu	<i>Drosophila ficusphila</i>	GCA_018152265.1	158.79 (94.61)	14.22 (98.7)
Dhyde	<i>Drosophila hydei</i>	GCA_003285905.1	151.30 (98.41)	5.150 (98.9)
Dimmi	<i>Drosophila immigrans</i>	GCA_018153375.1	163.77 (89.36)	11.45 (98.9)
Dinnu	<i>Drosophila innubila</i>	GCA_004354385.2	166.28 (100.0)	29.57 (98.8)
Dkikk	<i>Drosophila kikkawai</i>	GCA_018152535.1	185.80 (98.41)	21.81 (98.8)
Dmela	<i>Drosophila melanogaster</i>	GCA_000001215.4	143.73 (100.0)	25.29 (98.6)
Dmira	<i>Drosophila miranda</i>	GCA_003369915.1	286.71 (99.87)	35.26 (98.9)
Dmoja	<i>Drosophila mojavensis</i>	GCA_018153725.1	162.96 (99.87)	24.88 (99.0)
Dobsc	<i>Drosophila obscura</i>	GCA_018151105.1	179.77 (99.97)	3.93 (98.4)
Dpseu	<i>Drosophila pseudoobscura</i>	GCA_009870125.1	163.10 (99.89)	32.42 (98.7)
Drhop	<i>Drosophila rhopaloa</i>	GCA_018152115.1	193.38 (99.93)	15.81 (98.5)
Drufa	<i>Drosophila rufa</i>	GCA_018153105.1	196.67 (94.35)	24.72 (98.7)
Dserr	<i>Drosophila serrata</i>	GCA_002093755.1	193.27 (97.60)	1.010 (97.3)
Dsimu	<i>Drosophila simulans</i>	GCA_016746395.2	154.00 (99.76)	21.50 (99.0)
		<i>dryad.280gb5mr6</i>	131.51 (99.89)	23.40 (99.0)
Dsubo	<i>Drosophila subobscura</i>	GCA_008121235.1	126.19 (99.96)	24.18 (98.7)
Dsubp	<i>Drosophila subpulchrella</i>	GCA_014743375.2	263.87 (99.52)	11.59 (98.9)
		GCA_018150325.1	265.10 (98.87)	1.467 (96.8)
Dsuzu	<i>Drosophila suzukii</i>	GCA_013340165.1	266.69 (99.51)	2.610 (97.4)
		<i>CAKG01000000</i>	162.25 (94.55)	0.005 (83.8)
Dtaka	<i>Drosophila takahashii</i>	GCA_018152695.1	164.65 (99.47)	12.38 (97.8)
Dviri	<i>Drosophila virilis</i>	GCA_003285735.1	189.28 (99.91)	8.697 (99.0)
Dwill	<i>Drosophila willistoni</i>	GCA_000005925.1	220.00 (92.94)	4.707 (98.9)
Dyaku	<i>Drosophila yakuba</i>	GCA_016746365.2	147.66 (99.84)	25.18 (99.0)
Apomo	<i>Acetobacter pomorum</i>	NZ_AEUP00000000.1	3.332	0.076
Cinte	<i>Commensalibacter intestine</i>	NZ_AGFR00000000.1	2.454	0.476
Efaec	<i>Enterococcus faecalis</i>	NC_004668.1	2.870	2.807
Gmorb	<i>Gluconobacter morbifer</i>	NZ_AGQV00000000.1	2.887	0.423
Lbrev	<i>Lactobacillus brevis</i>	NC_008497.1	2.552	2.553
Lplan	<i>Lactobacillus plantarum</i>	NC_004567.2	3.231	3.231
Palca	<i>Providencia alcalifaciens</i>	NZ_AKKM01000049.1	3.990	3.99
Pburh	<i>Providencia burhodogranariea</i>	NZ_AKKL00000000.1	4.579	2.508
Pento	<i>Pseudomonas entomophila</i>	NC_008027.1	5.889	5.889
Prett	<i>Providencia rettgeri</i>	NZ_AJSB00000000.1	4.454	4.309
Scere	<i>Saccharomyces cerevisiae</i>	GCF_000146045.2_R64	12.16	0.924
		NC_002978.6	1.268	1.268
Wolb	<i>Wolbachia pipientis</i>	NC_012416	1.446	1.446

quality score <15). In addition, the `--merge` and `--include_unmerged` options were used to merge the detected overlapping PE reads into a single sequence. Finally, the `--stdout` option was enabled to generate an interleaved fastq output, which was converted to fasta format (losing quality and pairing information) with a simple `awk` one-liner for assignment analysis. As shown in Figure S1 and detailed in Tables S2 and S3, the quantity (and quality) of sequencing data was highly variable between samples, with the percentage of non-overlapping sequences ranging from 5.79 to 94.4 (median 35.0) as a consequence of different insert sizes; and the estimated percentage of duplicate reads ranging from 0.69 to 24.8 (median 4.44) (Figure S1B). Note that the sequencing data were not de-duplicated here, although this may be possible using the latest version of `fastp` (Chen et al., 2018).

2.3. Assignment of query sequences and contamination estimation

For each sample, the sequences contained in the filtered `fasta` files were matched to the target dictionaries of species-discriminating *k*-mers using Clark and Clark-I (Ounit et al., 2015). Briefly, analyzing a sequence consists of first decomposing it into its constituent *k*-mers (i.e., a *L* nt long sequence can be decomposed into $L - k + 1$ *k*-mers of length *k* nt) of length $k = 31$ and $k = 27$ for Clark and Clark-I, respectively. Each *k*-mer is then searched in the corresponding target dictionary and, if found, assigned to the underlying target species. Counting the number of *k*-mers assigned to the different species then provides a simple decision criterion for sequence classification. More specifically, for a given sequence, let t_1 and t_2 be the target species with the highest and second highest counts ($k_q(t_1)$ and $k_q(t_2) \leq k_q(t_1)$) of matching *k*-mers. If no species-discriminating *k*-mer was found in the sequence (i.e., $k_q(t) = 0$ for all target species *t*), the sequence is unassigned. If $k_q(t_1) > 0$, the sequence is assigned to species t_1 with a 'confidence score' defined as $c_q(t_1) = \frac{k_q(t_1)}{k_q(t_1) + k_q(t_2)}$, noting that $c_q(t_1) = 1$ if all the matching *k*-mers are assigned to t_1 (i.e., $k_q(t) = 0$ for all $t \neq t_1$). At the sample level, the origin and level of contamination can then be further assessed by counting the number of sequences assigned to the different target species. In practice, Clark was run with option `-s 2` to load only half of the species-discriminating *k*-mers in the target dictionary, following the manual recommendation indicating that this value 'represents a good trade-off between speed, accuracy and RAM usage'. Both Clark and Clark-I were run with the options `-n 1` (i.e., on a single thread) and `-m 0` (to compute the confidence score). The resulting `csv` files were parsed with a custom `awk` script to count for each sample i) the total number of sequences with no matching *k*-mer; ii) the total number of sequences with at least one matching *k*-mer; and iii) the proportion of sequences assigned to each target species. Four different criteria were considered for assigning sequences to their inferred species *t*, taking into account both the minimum number $nk_{\min}(t)$ of matching *k*-mers and the confidence score $c_q(t)$: i) $nk_{\min}(t) \geq 1$ and $c_q(t) > 0.9$; ii) $nk_{\min}(t) \geq 1$ and $c_q(t) > 0.95$; iii) $nk_{\min}(t) \geq 5$ and $c_q(t) > 0.9$; and the most stringent iv) $nk_{\min}(t) \geq 5$ and $c_q(t) > 0.95$. All subsequent analyses were performed using the R software (R Core Team, 2017).

3. Results

3.1. Clark and Clark-I run times

Publicly available short-read WGS data for 301 different samples derived from i) laboratory strains representing different drosophilid species ($n=43$); ii) 236 (putative) *D. suzukii* individuals representing 40 different populations; and iii) 22 pools of *D. suzukii* individuals representing 22 different populations were assigned to two different species-discriminating *k*-mer dictionaries built from the curated assemblies available for 29 drosophilid species (Figure 1) and 12 common drosophila commensals and pathogens (Table 1), using the *k*-mer-based approaches implemented in Clark and Clark-I (Ounit et al., 2015). Although this step is not required for assignment, the raw PE reads were filtered to limit the potential impact of varying sequence quality on the assessment of assignment efficiency and accuracy, particularly with respect to the observed proportion of unassigned sequences per query sample. After filtering, the total number of sequences per sample ranged from 1.61×10^6 to 367×10^6 (median of 18.5×10^6) for a total

Table 2 – Mean Clark and Clark-I run times (minimum-maximum) across the analyses of the 305 short-read sequencing datasets. Each analysis was run on a single thread of a cluster node equipped with a processor Intel® Xeon® CPU E5-2683 v4 @2.10GHz

Running time in min mean (min-max)	Clark	Clark-I
Loading of the <i>k</i> -mer dictionary	2.23 (1.13-5.04)	0.075 (0.032-0.137)
Assignment per 10 ⁶ sequences	1.05 (0.444-2.51)	0.619 (0.228-1.50)

length ranging from 0.248 Gb (i.e. $\sim 0.9X$ of the *D. suzukii* genome) to 36.9 Gb (i.e. $\sim 137X$ of the *D. suzukii* genome). The sequence length was representative of typical short read datasets, with a sample mean length (after merging overlapping reads) ranging from 92.7 bp to 287 bp (Figure S1C).

Tables S2 and S3 show the total Clark and Clark-I run times t_r for each sample, together with the time t_l required to load the corresponding *k*-mer target dictionary and the time t_a required to assign all sequences ($t_r = t_l + t_a$). As summarized in the Table 2, t_l was a few seconds for Clark-I and a few minutes for Clark, the Clark-I target dictionary containing about 75 times less *k*-mer than Clark's (see M&M). In addition, Clark-I required much less RAM than Clark (1.97 Gb vs 47.8 Gb), allowing it to run on a standard laptop. Note that Clark and Clark-I were run sequentially on each sample on a computer grid, but the samples were analyzed in parallel. Therefore, the run times between samples may be somewhat dependent on the characteristics of the underlying node, which explains the observed variation in dictionary loading times.

Given the size of the data sets, most of the analysis time was spent on sequence assignment which was almost linearly related to the number of sequences (Figure S2) as sequence length was similar across samples (Figure S1C). On average, the analysis of 1 million sequences (i.e., $\sim 0.56X$ of the *D. suzukii* genome with 150 nt reads) took 0.619 and 1.05 minutes with Clark-I and Clark, respectively (Table 1), making both approaches highly computationally efficient.

3.2. Proportion of assigned sequences

The percentage of sequences with no matching *k*-mer (i.e., not assignable) was similar between Clark (ranging from 2.29% to 85.5% with a median value of 20.1%) and Clark-I (ranging from 4.07% to 86.1% with a median value of 15.7%) (Figures 2A and 2B). Surprisingly, this percentage tended to be slightly lower for the *D. suzukii* sample (Ind-Seq or Pool-Seq) when analyzed with Clark-I, which may be related to the smaller *k*-mer size ($k=27$ for Clark-I and $k=31$ for Clark) leading to lower specificity. However, the proportion of sequences with no matching *k*-mer remained higher for Clark-I analyses for samples representative of the other species either represented or not represented in the target dictionaries (Figure 2B). As expected, and regardless of the program used, the highest percentages were observed for samples belonging to species not represented in the target dictionaries (up to 85.5% and 86.1% of sequences with no matching *k*-mer for the *D. repleta* sample analyzed with Clark and Clark-I, respectively), although the distribution was very wide and almost bimodal due to some samples being represented by closely related target species (see below). The sample representing target species had the lowest number of sequences with no matching *k*-mer, most of them (including *D. suzukii*) corresponding to short-read sequence data obtained from the same strains used to generate the reference assembly, with the notable exception of *D. melanogaster*, *D. simulans*, and the Wolbachia sample (see below), which were also outliers in the distributions of Figure 2B (see Table S4). Their values was actually similar to wild-caught *D. suzukii* samples (see below). The *D. simulans* sample was obtained from Madagascar individuals (Palmieri et al., 2014) thus distantly related to the two reference assembly strains, which may explain the observed pattern (see Discussion). Likewise, the analyzed *D. melanogaster* sample corresponded to a pool of 162 isogenic strains from the DGRP panel and may thus display higher genetic diversity (Zhu et al., 2012).

Consistent with a lower specificity of Clark-I (suggested by the unexpectedly slightly lower proportion of sequences with no matching *k*-mer in *D. suzukii* individuals), the percentages of

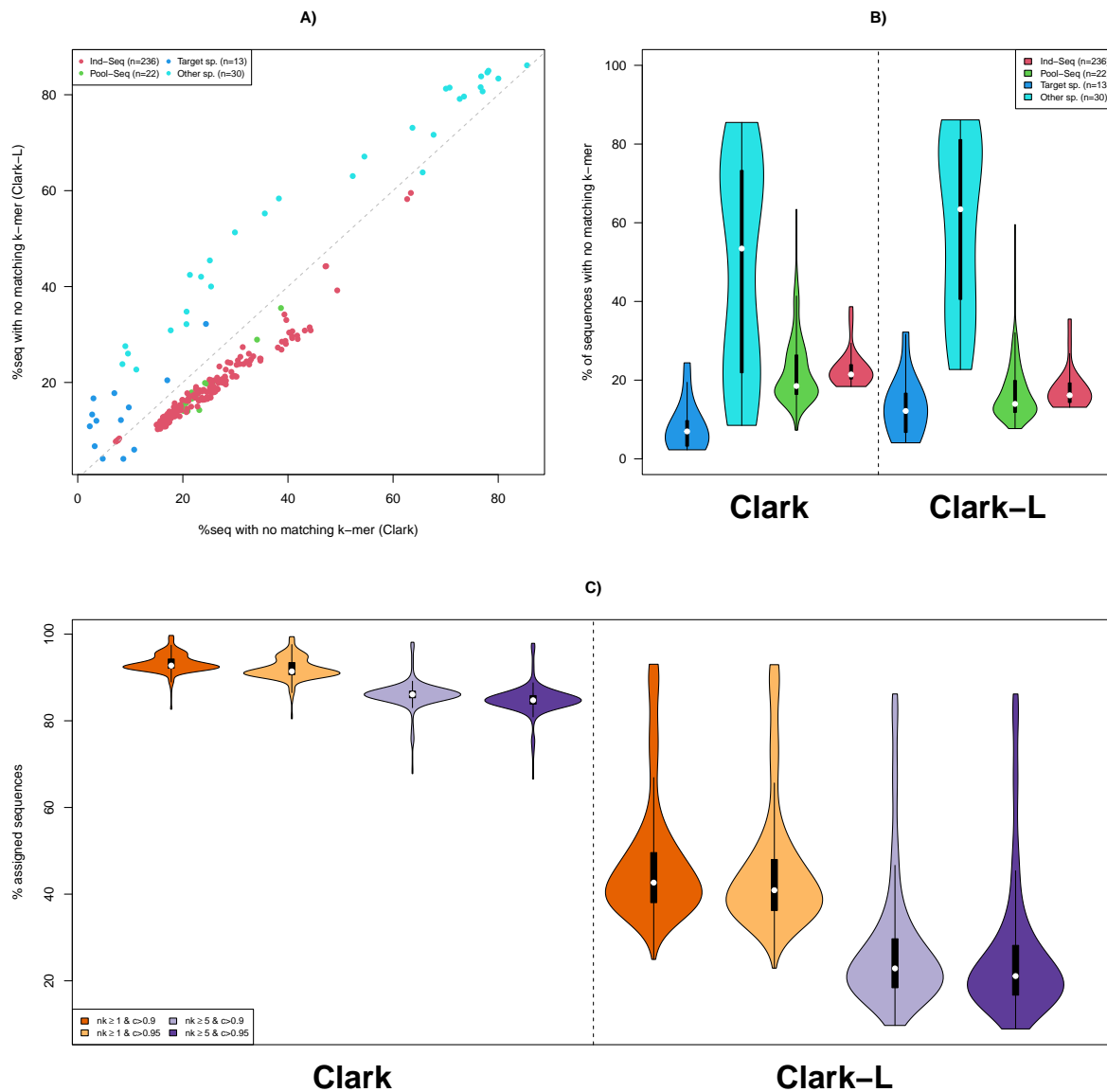


Figure 2 – Sequence assignment rate for the 301 samples analyzed with Clark and Clark-L. A) Percentage of sequences with no matching k -mer in the target dictionaries. Samples are colored according to their origin, i.e. i) dark blue if from species represented in the target dictionary ('Target sp. '); ii) light blue if from drosophilid species not represented in the target dictionary ('Other sp. '); iii) green for *D. sukuzii* individuals from Lewald et al. (2021) ('Ind-Seq'); and iv) red for pools of *D. sukuzii* individuals from Olazcuaga et al. (2020) ('Pool-Seq'). B) Violin plots showing the distribution of the percentage of sequences with no matching k -mer in the corresponding target dictionary with Clark (left panel) and Clark-L (right panel) analyses. For each analysis, four distributions are shown for the different sample origins (same color code as in A). C) Distribution of the percentage of assigned sequences (among those with at least one species-discriminating k -mer from the target dictionary) for four filtering criteria on i) the number nk of matching k -mers ($nk \geq 1$ or $nk \geq 5$); and ii) the assignment confidence score c as defined in the main text ($c > 0.9$ or $c > 0.95$).

assigned sequences among assignable sequences (i.e., containing at least one k -mer matching the dictionary of target species discriminating k -mers) were much lower with Clark-L than with Clark (Figure 2C). The percentages of assigned sequences always decrease with the stringency of the filtering criteria on the number nk of matching k -mers ($nk \geq 1$ or $nk \geq 5$) and the assignment confidence score c (as defined above in the M&M section), with the threshold on nk having the strongest effect. At the most stringent criterion ($nk \geq 5$ and $c > 0.95$), which was chosen for

this latter Wolbachia sample was actually obtained from sequencing a *D. melanogaster* strain, and the observed level of contamination was in close agreement with the 5% of reads mapping to the Wolbachia wMel genome by the original authors (Newton and Sheehan, 2015). Similar results were obtained when scanning these 13 samples with Clark-I (Figure S3 and Table S5), with some notable differences. Indeed, the percentage of sequences assigned to their species of origin was also above 99% (including the *D. subpulchrella* one) or close to it (with 98.0% for *D. yakuba*) for 8 of the 9 samples that showed similarly high assignment rates with Clark. However, it was substantially lower for the *D. sukuzii* sample (92.1%), with 7.22% of its sequences assigned to the *D. subpulchrella* sister species. Similarly, only 86.2% of the *D. melanogaster* sample sequences were assigned to *D. melanogaster*, with 6.86%, 2.62%, 1.65%, and 1.40% assigned to *D. sukuzii*, *D. simulans*, *D. virilis*, and Wolbachia, respectively. Conversely, the percentage of correctly assigned sequences was higher with Clark-I than with Clark for the *D. biarmipes* (96.0%); *D. simulans* (98.1%) and Wolbachia (40.0% with 55.8% assigned to *D. melanogaster*) samples, the latter apparently being overestimated.

Of the 30 samples from non-target species, 16 had more than 96% of their reads assigned to a single target species by Clark (Figure 3). As expected, the corresponding species was generally the most closely related (Kim et al., 2021). More precisely, samples from i) *D. paulistorum* and *D. insularis* (*D. willistoni* subgroup) and *D. sucinea* and *D. nebulosa* (*bocainensis* subgroup from the *willistoni* group) had 99.7%, 99.7%, 98.1%, and 97.9% of their sequences assigned to *D. willistoni*, respectively; ii) *D. parabiptinata*, *D. malerkotliana pallens*, *D. malerkotliana malerkotliana*, *D. pseudoananassae*, and *D. pseudoananassae nigrens*, all of which belong to the *ananassae* subgroup, had 99.2%, 99.1%, 99.0%, 96.5%, and 96.0% of their sequences assigned to *D. biptinata* (*ananassae* subgroup), respectively; iii) *D. ambigua* and *D. tristis* (*obscura* subgroup) had 98.7% and 97.3% of their sequences assigned to *D. obscura*, respectively; iv) *D. americana* and *D. littoralis* (*virilis* group) had 99.2% and 98.6% of their sequences assigned to *D. virilis*, respectively; and finally v) *D. carrolli*, *D. fuyamai*, and *D. kurseongensis* (*rhopaloo* subgroup) had 98.2%, 98.0%, and 97.7% of their sequences assigned to *D. rhopaloo*, respectively. As shown in Figure S4A, these 16 samples also had percentages of sequences with no matching *k-mer* in the range of those observed for samples from target species (Figure 2), i.e. <40% except for *D. sucinea* and *D. nebulosa*. For the other samples from the most distantly related species, both the highest observed assignment rate (to a target species) and the percentage of sequences with no matching *k-mer* clearly suggested that the target repository was not representative. At the extreme, the most represented target species capture less than 30% of the assigned sequences for the samples from *D. repleta*, *D. pruinosa*, *D. ohnishii*, and *D. bocqueti* (Figures 3 and S4A). Such species may therefore be considered unassignable with the current version of the *k-mer* dictionary. Despite a higher proportion of sequences with no matching species-discriminating *k-mer*, very similar results were obtained with Clark-I (Figure S3 and S4B).

3.4. Scanning 236 Ind-Seq and 22 Pool-Seq *D. sukuzii* WGS data

As summarized in Figure 4 (see Table S4 for details), sequences from the 236 Ind-Seq (Lewald et al., 2021) and 22 Pool-Seq (Olazcuaga et al., 2020) *D. sukuzii* were generally assigned to *D. sukuzii* by Clark. More precisely, 215 of the 236 Ind-Seq and 17 of the 22 Pool-Seq showed > 95% of their (assigned) sequences assigned to *D. sukuzii*, with a median proportion of 97.5% over the 258 samples. It should be noted that these 215 individuals and 17 pools, which can be unambiguously considered as fully *D. sukuzii*, all had a non-negligible fraction of their sequences assigned to *D. subpulchrella* with a median of 1.94% (ranging from 1.50% to 3.12%) and 2.20% (ranging from 1.96% to 2.64%), respectively. These proportions were higher than the one observed for the *D. sukuzii* reference sample (0.433%) and may be related to the incomplete representation of genetic diversity within *D. sukuzii* by the *k-mer* dictionary (see Discussion). Conversely, the results allowed 16 clearly mislabeled *D. sukuzii* individuals to be identified as *D. simulans* (n=5) or *D. subpulchrella* (n=11). These consist of i) the 5 individuals (with US-Ca2 ID prefix, Table S2) sampled simultaneously in Watsonville (California, USA) with 92.1% to 96.9% of their sequences assigned to *D. simulans* (96.9% to 98.4% if Wolbachia is also included); ii) the 5 individuals (with Ko-Nam ID prefix, Table S2) sampled in Namwon (South Korea) with 97.9% to

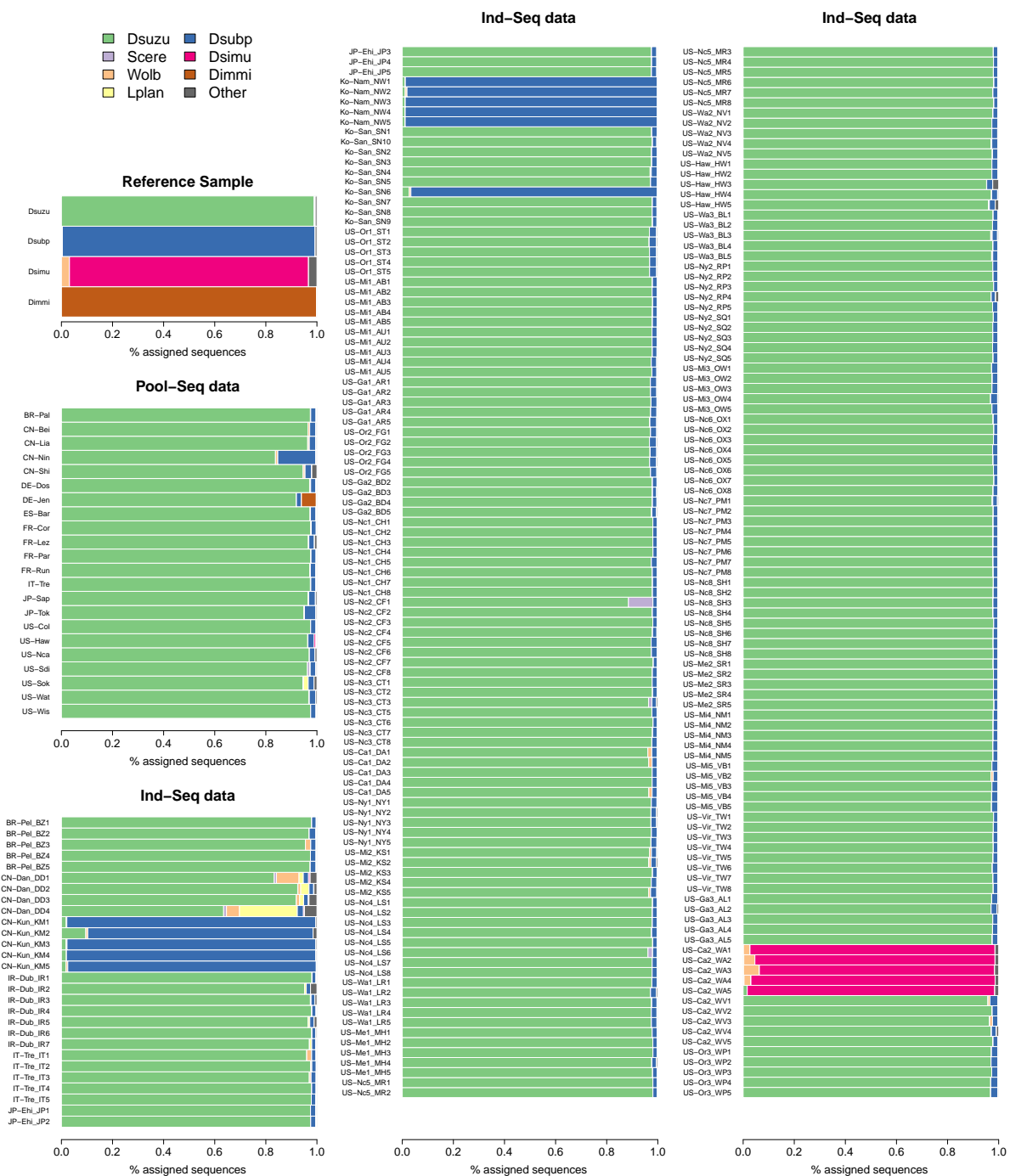


Figure 4 – Barplots summarizing assignment results obtained with Clark using the most stringent sequence assignment criterion (i.e., $nk \geq 5$ and $c > 0.95$, see the main text) for the *D. sukukii* Ind-Seq (n=236) and Pool-Seq (n=22) samples. For each sample, the proportions of sequences assigned to the 7 target species that contribute at least 5% of the sequences of one of any of the 258 samples are shown using the color code indicated in the top-left legend. The proportions of sequences assigned to the 34 other target species are shown in gray.

98.7% of their sequences assigned to *D. subpulchrella*; iii) one of the 10 individuals (with Ko-San ID prefix, Table S2) sampled in Sancheong (South Korea) with 96.4% of its sequences assigned to *D. subpulchrella* (the other 9 individuals showing only 1.71% to 2.09% of their sequences assigned to *D. subpulchrella*); and iv) four of the five individuals (with CN-Kun ID prefix, Table S2) sampled in Kunming (Yunnan, China) with 97.3% to 97.6% of their sequences assigned to

D. subpulchrella. The last CN-Kun individual had a unique pattern with 88.1% of its sequences assigned to *D. subpulchrella* and 9.58% assigned to *D. suzukii*, which may be consistent with a recent hybrid origin (see Discussion). For the 10 individuals that can be unambiguously considered as fully *D. subpulchrella* (i.e. with >95% of their sequences assigned to *D. subpulchrella*), an assignment pattern opposite to that of the *D. suzukii* individuals was observed, as all of them had a non-negligible fraction of their sequences assigned to *D. suzukii* with a median value of 1.61% (ranging from 1.14% to 2.78%).

Among the 22 Pool-Seq samples, two to three pools were found to be likely contaminated with non-*D. suzukii* individuals. These are i) the DE-Jen pool of 100 individuals sampled in Jena (Germany), which contains 5.79% of sequences assigned to *D. immigrans*; ii) the CN-Nin pool of 50 individuals sampled in Ningbo (Zhejiang, China), which contains 15.0% of sequences assigned to *D. subpulchrella* (and 83.8% to *D. suzukii*); and iii) the JP-Tok pool of 50 individuals sampled in Tokyo (Japan), with 4.47% of sequences assigned to *D. subpulchrella* (and 94.9% to *D. suzukii*). Assuming an equal contribution of pool individuals to the Pool-Seq sequences, the DE-Jen pool may actually contain up to 6 *D. immigrans* individuals (and 94 *D. suzukii* individuals). Furthermore, to estimate the number of *D. subpulchrella* individuals in contaminated pools while accounting for *D. suzukii* and *D. subpulchrella* cross-assignment of sequences, let $\alpha = \frac{p_{\text{sub}}}{p_{\text{sub}} + p_{\text{suz}}}$ be the relative proportion of sequences assigned to *D. subpulchrella*. Based on the median proportions observed in the Ind-Seq samples, the following rough estimates were obtained: $\hat{\alpha}_{\text{suz}} = \frac{0.0194}{0.977 + 0.0194} = 0.0195$ for *D. suzukii* individuals and $\hat{\alpha}_{\text{sub}} = \frac{0.0151}{0.0151 + 0.0978} = 0.985$ for *D. subpulchrella* individuals. The number of *D. subpulchrella* individuals n_{sub} in a contaminated pool of n individuals can then simply be derived from these estimates using their observed relative proportion α_o as $n_{\text{sub}} = n \frac{\alpha_o - \hat{\alpha}_{\text{suz}}}{\hat{\alpha}_{\text{sub}} - \hat{\alpha}_{\text{suz}}}$. This leads to an estimated number of *D. subpulchrella* individuals of $\hat{n}_{\text{sub}}^{\text{CN-Nin}} = 6.85$ and $\hat{n}_{\text{sub}}^{\text{JP-Tok}} = 1.32$, i.e. probably 7 and 1 *D. subpulchrella* individuals within the CN-Nin and JP-Tok pools, respectively.

Overall, very low levels of Wolbachia contamination were detected within the Ind-Seq and Pool-Seq samples, with median proportions of assigned sequences of $3.80 \times 10^{-4}\%$ and 0.145%, respectively. However, 14 samples (Ind-Seq only) had more than 1% of their sequences assigned to Wolbachia. They consisted of i) the five US-Ca2 individuals mentioned above, which are actually *D. simulans*, with proportions ranging from 1.08% to 6.17%; ii) the four individuals with the CN-Dan ID prefix (Table S2), sampled in Dandeong (China), with proportions ranging from 1.07% to 8.82%; iii) three of the five individuals with the US-Ca1 ID prefix (Table S2) sampled in Davis (California, USA) with proportions ranging from 1.29% to 1.55%; iv) one of the five individuals with the BR-Pel ID prefix sampled in Pelotas (Brazil) with a proportion of 2.07%; and v) one of the five individuals with the IT-Tre ID prefix sampled in Trento (Italy) with a proportion of 1.92%. Finally, a few Ind-Seq and Pool-Seq samples showed non-negligible to substantial contamination with five of the 11 other microbial species represented in the k -mer target dictionary. For example, more than 1% of the sequences were assigned to the *L. plantarum* bacterial gut symbiont for five samples corresponding to i) the four CN-Dan individuals (see above), with proportions ranging from 1.55% to 22.6%; and ii) the US-Sok pool of 50 individuals sampled in Dayton (Oregon, USA) with a proportion of 1.87%. Similarly, > 1% of the sequences were assigned to *S. cerevisiae* yeast for five samples corresponding to i) one of the four CN-Dan individuals with a proportion of 1.12%; ii) three individuals (with ID prefixes US-Nc2, US-Nc3, and US-Nc4, Table S2) sampled in different locations in North Carolina (USA) with proportions ranging from 1.33% to 9.58%; and iii) the US-Sdi pool of 50 individuals sampled in San-Diego (California, USA) with a proportion of 1.02%. At the margin, three other microbial species were also found to be represented by more than 1% of the sequences in at least one sample. These are i) the *A. pomorum* gut bacteria in two Chinese (CN-Dan) individuals (with proportions of 1.04% and 1.46%) and in the CN-Shi pools of 50 individuals sampled in Shiping (China) with 1.56%; ii) the *L. brevis* intestinal bacteria also found in two Chinese (CN-Dan) individuals with proportions of 1.12% and 4.19%; and iii) the *E. faecalis* pathogens in an Irish individual with proportion of 1.65%.

As expected from the assignment of *D. suzukii* and *D. subpulchrella* reference samples, Figure S5 (see Table S5 for details) suggested a worse performance of Clark-I. The proportions of *D.*

suzukii sequences appeared to be substantially underestimated, with a higher effect of cross-assignment with *D. subpulchrella*. In addition, Clark-I did not allow to detect the presence of the microbial target species as detected by Clark.

4. Discussion

The primary objective of this study was to propose and evaluate a computationally fast and accurate method for assessing contamination levels in publicly available WGS data for the *D. suzukii* species, which has been increasingly studied over the past decade. The availability of high quality genome assemblies for a wide range of drosophilid species (Kim et al., 2021) made it possible to rely on a *k-mer*-based approach consisting of constructing and querying dictionaries of species-discriminating *k-mers*. Such an approach has already proven to be quite valuable and benefits from the availability of optimized software, such as Kraken (Wood et al., 2019; Wood and Salzberg, 2014) or Clark (Ounit and Lonardi, 2016; Ounit et al., 2015), which were primarily developed for metagenomics applications but have also been proposed for contaminant detection (Cornet and Baurain, 2022). As in the latter case, our primary goal here was to classify sequences at the level of predefined (target) species, and Clark thus seemed particularly attractive due to its computationally efficient, tractable, and flexible way of both constructing and querying user-defined *k-mer* dictionaries. Although Kraken may be able to further assign higher-level taxonomic labels by considering phylogenetic relationships among target species, this feature was not critical for our purpose. In fact, it may have made it more difficult in practice, since the phylogeny among Drosophilidae species is far from being fully and unambiguously resolved. In particular, Finet et al. (2021) recently provided evidence for a paraphyletic status of the subgenus *Sophophora*, to which most of the target species belong (Figure 1). However, as illustrated by the assignment of sequences from species closely related to one of the represented groups or subgroups (e.g., *ananassae* or *obscura*) but not included in the construction of the *k-mer* dictionary, species-level assignment provided consistent results about their origin. Yet, assignment of samples to species belonging to groups or subgroups less well represented by the target species should be interpreted with caution, especially when the observed proportion of non-matching *k-mers* is high (Figure S4). In such cases, analysis with a newly built *k-mer* dictionary including more closely related species may be valuable. Indeed, our main focus was on the evaluation of *D. suzukii* samples. We therefore chose to deliberately overrepresent the *suzukii* subgroup in the *k-mer* dictionary construction by including the high quality genome assemblies available for *D. suzukii*, *D. subpulchrella*, and *D. biarmipes*. The latter two species were in fact the most likely confounders in field-collected samples from the Asian range of *D. suzukii* (see Introduction). Interestingly, the inclusion of these closely related species seemed to have only a limited effect on the number of discriminating *k-mers* in the resulting dictionary, with the percentage of sequences with no assigned *k-mer* for their corresponding reference samples being in the range of that observed for reference samples from other target species (Figure 2).

Searching the resulting *k-mer* dictionary of target species sequences with Clark (Ounit et al., 2015) was highly efficient in terms of both run time and memory requirements. This makes analyses of common short-read sequencing data tractable on standard workstations or computer grids, and even on a standard laptop when using the lighter Clark version (Ounit et al., 2015), although at some moderate cost in assignment accuracy. More specifically, it took only a few minutes and about 50 Gb of RAM to load the Clark dictionary (<1 min and <2 Gb of RAM for the Clark-I dictionary), and the mapping took about one minute per million of typical 150 nt short reads. Such assignment analyses could thus be performed routinely and may be worth including as a standard part of the quality control of sequencing data, at least for the *D. suzukii* sample. Note that here we have chosen to screen sequences after filtering raw PE reads with *fastp* (Chen et al., 2018), primarily to limit the potential impact of varying sequence quality across samples on the assessment of assignment accuracy (e.g., proportion of sequences assigned). Although this is not required in practice when trying to assign samples or assess their contamination levels, it seems to be a reasonable strategy when combined with other quality control procedures. Finally, for contamination assessment at the whole-sample level, *k-mer*-based approaches represent an

attractive and efficient mapping-free alternative to competitive mapping methods that consist of mapping sequencing reads to hologenomes constructed from target species assemblies (e.g. Kapun et al., 2021). It also allows for easy interrogation of a wider range of target species, providing good quality genome assemblies are available. For sequence filtering purposes, however, such approaches must be used with caution because they rely on species-discriminating *k*-mers and thus may leave a substantial fraction of sequences unassigned. More advanced (and computationally expensive) methods may then be valuable, such as the one implemented in Clark-s (Ounit and Lonardi, 2016), which allows some mismatches in *k*-mer matching to improve the sensitivity of sequence assignment or even Kraken (Wood et al., 2019; Wood and Salzberg, 2014), which was used here to identify contaminating contigs in the assemblies of the target species. Indeed, this program can rely on *k*-mers shared by several species for sequence assignment, and not only species discriminating *k*-mers, since all the *k*-mers of the target dictionary (possibly built from very large databases such as the NCBI nt) are mapped to the nodes of a phylogenetic tree (species discriminating *k*-mers to terminal nodes and shared *k*-mers to internal nodes).

Overall, the results obtained from the analysis of WGS data for reference samples belonging to different target species and single or pools of *D. sukuzii* individuals demonstrated the high accuracy of the *k*-mer-based approach. It also allowed the unambiguous identification of 16 mislabeled *D. sukuzii* individuals among the 236 (i.e. 6.78%) from the Lewald et al. (2021) study. Five corresponded to *D. simulans* individuals collected at the same site in Watsonville (California, USA). It should be noted that Lewald et al. (2021) discarded these samples from their analysis because they displayed too low mapping rates like the Dandong (China) sample, which was found here to be substantially contaminated with microbial (and Wolbachia) sequences. The eleven other non-*D. sukuzii* individuals from three different locations in Asia could all be assigned to *D. subpulchrella* individuals. These were also identified as *D. subpulchrella* by Lewald et al. (2021) (and discarded from their analysis) using a phylogenetic analysis of the mitochondrial COX2 gene. Two of the 22 Pool-Seq samples of (Olazuaga et al., 2020) collected in the Asian native area were also, and unexpectedly, found to be contaminated with *D. subpulchrella* individuals, namely CN-Nin with 7 *D. subpulchrella* individuals and to a lesser extent JP-Tok with 1 *D. subpulchrella* individual (both out of 50 individuals in total). More surprisingly, but confirming a gene-based analysis by D. Obbard (pers. comm.), the DE-Jen pool collected in Jena (Germany) was found to be contaminated with 5 to 6 *D. immigrans* individuals (out of 100). These observations may indicate that great care should be taken when analyzing sequencing data from wild-caught samples, and that more attention should probably be paid to species identification prior to sequencing. High-throughput metabarcoding and non-destructive approaches, such as those recently proposed by Piper et al. (2022), may represent valuable alternatives to sometimes difficult morphological identification by allowing rapid and efficient diagnosis of *D. sukuzii* samples at any life stage. Such efforts may be even more critical for Pool-Seq experiments, since filtering out contaminated sequences (e.g., using competitive mapping) is far more challenging than discarding mislabeled Ind-Seq samples, especially when the sample is contaminated by individuals from very closely related species (such as *D. subpulchrella* for *D. sukuzii*).

Although two different *D. sukuzii* genome assemblies were used to build the species discriminating *k*-mer dictionary, all (pure) *D. sukuzii* Ind-Seq and Pool-Seq samples showed a small but non-negligible fraction of their sequences (from 1.14% to 2.78%) assigned to *D. subpulchrella* by the most stringent criterion. Because i) the *D. sukuzii* reference genome assemblies were derived from isofemale lines established from individuals sampled in the North American (Chiu et al., 2013) and European (Ometto et al., 2013) invaded areas; and ii) *D. subpulchrella* has not been yet described (to our knowledge) outside the Asian native range of *D. sukuzii*; it is highly unlikely that this pattern is the result of pervasive gene flow between the two species, but rather can be explained by the close phylogenetic relationship between the two species. Indeed, some *D. subpulchrella*-discriminating *k*-mers may actually map to orthologous regions not represented in the *D. sukuzii* reference assemblies and/or capture shared genetic variation between the two species due to incomplete lineage sorting (ILS). Including more reference assemblies (e.g., from different strains) for each target species may be considered as a valuable strategy to improve both the sensitivity (by 'positive filtering' of the discriminating *k*-mers that capture intraspecific

genetic variation) and specificity (by ‘negative filtering’ of the incompletely sorted k -mers). The optimal number of representative assemblies is thus likely to both depend on the relatedness of the selected target species and for each target species on their genetic diversity. Alternatively, the misassigned short read sequences found in the analyzed samples can be included in the construction of the k -mer dictionary, assuming that the considered samples are not contaminated and are ‘pure’ representatives of the corresponding target species. Such refined target dictionaries may even further allow providing (rough) estimates of the genome-wide level of interspecific gene flow, or at least the identification of highly admixed individuals. Hence, in the sample of identified *D. subpulchrella* individuals, if about 2% of the short-read sequences were assigned to *D. suzukii* (in a similar but reversed pattern as observed for *D. suzukii* individuals), one (presumably) *D. subpulchrella* individual had nearly 10% of its sequences assigned to *D. suzukii*. The status of this sample may be of special interest for further study as it could represent a previously unreported case supporting some recent (i.e., only a few generations back) admixture events between *D. suzukii* and *D. subpulchrella*. As discussed by Lalyer et al. (2021), if no such recent events have been reported to date, several studies suggest that hybridization has occurred between these two sister species (Conner et al., 2017).

Overall, the present analysis allowed the definition of a large curated dataset consisting of > 60 population samples representative of global genetic diversity, which may be valuable for further *D. suzukii* population genetics studies. Although constructed with the analysis of *D. suzukii* samples in mind, the k -mer dictionary developed here may be directly relevant to the analysis of the level of contamination of samples from other target species such as *D. simulans* or *D. melanogaster*. Likewise, the current dictionary also allows for the rapid identification of Wolbachia-infected samples, which may be of interest for a first rapid screening of drosophilids samples since the set of Wolbachia-discriminating k -mers was built by combining *D. simulans* and *D. melanogaster* Wolbachia assemblies. More generally, while we advocate careful sample identification and verification prior to sequencing, the proposed framework is straightforward and computationally efficient. It thus could be considered as a routine post-hoc quality check approach to be applied prior to any data analysis and prior to data submission to public repositories.

Acknowledgements

I am grateful to the genotoul bioinformatics platform Toulouse Occitanie (Bioinfo Genotoul, <https://doi.org/10.15454/1.5572369328961167E12>) for providing computing resources. I also wish to thank Darren Obbard, Alan Bergland, Joaquin Nunez and Nicolas Rode for helpful discussion and feedback. Preprint version 2 of this article (Gautier, 2023a, <https://doi.org/10.1101/2023.04.18.537389>) has been peer-reviewed and recommended by Peer Community In Genomics (Galtier, 2023, <https://doi.org/10.24072/pci.genomics.100244>).

Fundings

The author declares he has received no specific funding for this study.

Conflict of interest disclosure

The author declares that he complies with the PCI rule of having no financial conflicts of interest in relation to the content of the article. The author is a PCI Evol Biol recommender.

Data, script, code, and supplementary information availability

The Clark and Clark-I k -mer databases and the (cleaned) assemblies used to build them have been made publicly available from the Data INRAE repository (Gautier, 2023b, <https://doi.org/10.57745/HYTIBH>). The compressed archive also contains scripts used to run Clark and Clark-I analyses and parse the results. All sequencing data analyzed in this study are publicly available under the accession IDs reported in Tables 1, S2 and S3. Supplementary

Tables S1 to S5 and Supplementary Figures S1 to S5 are available online with the latest preprint version of the manuscript (Gautier, 2023a, <https://doi.org/10.1101/2023.04.18.537389>).

References

- Asplen MK, Anfora G, Biondi A, Choi DS, Chu D, Daane KM, Gibert P, Gutierrez AP, Hoelmer KA, Hutchison WD, Isaacs R, Jiang ZL, Kárpáti Z, Kimura MT, Pascual M, Philips CR, Plantamp C, Ponti L, Véték G, Vogt H, et al. (2015). *Invasion biology of spotted wing Drosophila (Drosophila suzukii): a global perspective and future priorities*. *Journal of Pest Science* **88**, 469–494. <https://doi.org/10.1007/s10340-015-0681-z>.
- Atallah J, Teixeira L, Salazar R, Zaragoza G, Kopp A (2014). *The making of a pest: the evolution of a fruit-penetrating ovipositor in Drosophila suzukii and related species*. *Proceedings of the Royal Society B: Biological Sciences* **281**, 20132840. <https://doi.org/10.1098/rspb.2013.2840>.
- Chang CH, Gregory LE, Gordon KE, Meiklejohn CD, Larracuenta AM (2022). *Unique structure and positive selection promote the rapid divergence of Drosophila Y chromosomes*. *eLife* **11**, e75795. <https://doi.org/10.7554/eLife.75795>.
- Chen S, Zhou Y, Chen Y, Gu J (2018). *fastp: an ultra-fast all-in-one FASTQ preprocessor*. *Bioinformatics* **34**, i884–i890. <https://doi.org/10.1093/bioinformatics/bty560>.
- Chiu JC, Jiang X, Zhao L, Hamm CA, Cridland JM, Saelao P, Hamby KA, Lee EK, Kwok RS, Zhang G, Zalom FG, Walton VM, Begun DJ (2013). *Genome of Drosophila suzukii, the spotted wing drosophila*. *G3* **3**, 2257–71. <https://doi.org/10.1534/g3.113.008185>.
- Cini A, Ioriatti C, Anfora G (2012). *A review of the invasion of Drosophila suzukii in Europe and a draft research agenda for integrated pest management*. *Bulletin of Insectology* **65**, 149–160.
- Conner WR, Blaxter ML, Anfora G, Ometto L, Rota-Stabelli O, Turelli M (2017). *Genome comparisons indicate recent transfer of wRi-like Wolbachia between sister species Drosophila suzukii and D. subpulchrella*. *Ecology and Evolution* **7**, 9391 – 9404. <https://doi.org/10.1002/ece3.3449>.
- Cornet L, Baurain D (2022). *Contamination detection in genomic data: more is not enough*. *Genome Biology* **23**, 60. <https://doi.org/10.1186/s13059-022-02619-9>.
- Durkin SM, Chakraborty M, Abrieux A, Lewald KM, Gadau A, Svetec N, Peng J, Kopyto M, Langer CB, Chiu JC, Emerson JJ, Zhao L (2021). *Behavioral and Genomic Sensory Adaptations Underlying the Pest Activity of Drosophila suzukii*. *Molecular Biology and Evolution* **38**, 2532–2546. <https://doi.org/10.1093/molbev/msab048>.
- Finet C, Kassner VA, Carvalho AB, Chung H, Day JP, Day S, Delaney EK, De Ré FC, Dufour HD, Dupim E, Izumitani HF, Gautério TB, Justen J, Katoh T, Kopp A, Koshikawa S, Longdon B, Loreto EL, Nunes MDS, Raja KKB, et al. (2021). *DrosoPhyla: Resources for Drosophilid Phylogeny and Systematics*. *Genome Biology and Evolution* **13**. evab179. <https://doi.org/10.1093/gbe/evab179>.
- Francois CM, Durand F, Figuet E, Galtier N (2020). *Prevalence and Implications of Contamination in Public Genomic Resources: A Case Study of 43 Reference Arthropod Assemblies*. *G3* **10**, 721–730. <https://doi.org/10.1534/g3.119.400758>.
- Galtier N (2023). *Decontaminating reads, not contigs*. *Peer Community in Genomics*, 100244. <https://doi.org/10.24072/pci.genomics.100244>.
- Gautier M (2023a). *Efficient k-mer based curation of raw sequence data: application in Drosophila suzukii*. *bioRxiv*, 2023.04.18.537389, ver. 2, peer-reviewed and recommended by Peer Community in Genomics. <https://doi.org/10.1101/2023.04.18.537389>.
- Gautier M (2023b). *kmer dictionaries and associated scripts for kmer contaminant detection in Drosophila suzukii sequencing data using Clark program*. Version V1. Data INRAe, Recherche Data Gouv. <https://doi.org/10.57745/HYTIBH>.
- Jezovit JA, Levine JD, Schneider J (2017). *Phylogeny, environment and sexual communication across the Drosophila genus*. *Journal of Experimental Biology* **220**, 42–52. <https://doi.org/10.1242/jeb.143008>.

- Kapun M, Nunez JCB, Bogaerts-Márquez M, Murga-Moreno J, Paris M, Outten J, Coronado-Zamora M, Tern C, Rota-Stabelli O, Guerreiro MPG, Casillas S, Orengo DJ, Puerma E, Kankare M, Ometto L, Loeschcke V, Onder BS, Abbott JK, Schaeffer SW, Rajpurohit S, et al. (2021). *Drosophila Evolution over Space and Time (DEST): A New Population Genomics Resource*. *Molecular Biology and Evolution* **38**, 5782–5805. <https://doi.org/10.1093/molbev/msab259>.
- Kim BY, Wang JR, Miller DE, Barmina O, Delaney E, Thompson A, Comeault AA, Peede D, D'Agostino ERR, Pelaez J, Aguilar JM, Haji D, Matsunaga T, Armstrong EE, Zych M, Ogawa Y, Stamenković-Radak M, Jelić M, Veselinović MS, Tanasković M, et al. (2021). *Highly contiguous assemblies of 101 drosophilid genomes*. *eLife* **10**. <https://doi.org/10.7554/eLife.66405>.
- Klasson L, Kumar N, Bromley R, Sieber K, Flowers M, Ott SH, Tallon LJ, Andersson SGE, Dunning Hotopp JC (2014). *Extensive duplication of the Wolbachia DNA in chromosome four of Drosophila ananassae*. *BMC Genomics* **15**, 1097. <https://doi.org/10.1186/1471-2164-15-1097>.
- Lalyer CR, Sigsgaard L, Giese B (2021). *Ecological vulnerability analysis for suppression of Drosophila suzukii by gene drives*. *Global Ecology and Conservation* **32**, e01883. <https://doi.org/10.1016/j.gecco.2021.e01883>.
- Lewald KM, Abrieux A, Wilson DA, Lee Y, Conner WR, Andrezza F, Beers EH, Burrack HJ, Daane KM, Diepenbrock L, Drummond FA, Fanning PD, Gaffney MT, Hesler SP, Ioriatti C, Isaacs R, Little BA, Loeb GM, Miller B, Nava DE, et al. (2021). *Population genomics of Drosophila suzukii reveal longitudinal population structure and signals of migrations in and out of the continental United States*. *G3* **11**. <https://doi.org/10.1093/g3journal/jkab343>.
- Machado HE, Bergland AO, Taylor R, Tilk S, Behrman E, Dyer K, Fabian DK, Flatt T, González J, Karasov TL, Kim B, Kozeretska I, Lazzaro BP, Merritt TJ, Pool JE, O'Brien K, Rajpurohit S, Roy PR, Schaeffer SW, Serga S, et al. (2021). *Broad geographic sampling reveals the shared basis and environmental correlates of seasonal adaptation in Drosophila*. *eLife* **10**, e67577. <https://doi.org/10.7554/eLife.67577>.
- Manni M, Berkeley MR, Seppey M, Simão FA, Zdobnov EM (2021). *BUSCO Update: Novel and Streamlined Workflows along with Broader and Deeper Phylogenetic Coverage for Scoring of Eukaryotic, Prokaryotic, and Viral Genomes*. *Molecular Biology and Evolution* **38**, 4647–4654. <https://doi.org/10.1093/molbev/msab199>.
- Newton ILG, Sheehan KB (2015). *Passage of Wolbachia pipientis through Mutant Drosophila melanogaster Induces Phenotypic and Genomic Changes*. *Applied and Environmental Microbiology* **81**, 1032–1037. <https://doi.org/10.1128/AEM.02987-14>.
- Olazcuaga L, Loiseau A, Parrinello H, Paris M, Fraimout A, Guedot C, Diepenbrock LM, Kenis M, Zhang J, Chen X, Borowiec N, Facon B, Vogt H, Price DK, Vogel H, Prud'homme B, Estoup A, Gautier M (2020). *A Whole-Genome Scan for Association with Invasion Success in the Fruit Fly Drosophila suzukii Using Contrasts of Allele Frequencies Corrected for Population Structure*. *Molecular Biology and Evolution* **37** (8), 2369–2385. <https://doi.org/10.1093/molbev/msaa098>.
- Ometto L, Cestaro A, Ramasamy S, Grassi A, Revadi S, Siozios S, Moretto M, Fontana P, Varotto C, Pisani D, Dekker T, Wrobel N, Viola R, Pertot I, Cavalieri D, Blaxter M, Anfora G, Rota-Stabelli O (2013). *Linking genomics and ecology to investigate the complex evolution of an invasive Drosophila pest*. *Genome Biology and Evolution* **5**, 745–757. <https://doi.org/10.1093/gbe/evt034>.
- Ounit R, Lonardi S (2016). *Higher classification sensitivity of short metagenomic reads with CLARK-S*. *Bioinformatics* **32**, 3823–3825. <https://doi.org/10.1093/bioinformatics/btw542>.
- Ounit R, Wanamaker S, Close TJ, Lonardi S (2015). *CLARK: fast and accurate classification of metagenomic and genomic sequences using discriminative k-mers*. *BMC Genomics* **16**, 236. <https://doi.org/10.1186/s12864-015-1419-2>.
- Palmieri N, Nolte V, Chen J, Schlötterer C (2014). *Genome assembly and annotation of a Drosophila simulans strain from Madagascar*. *Molecular Ecology Resources* **15**, 372–381. <https://doi.org/10.1111/1755-0998.12297>.
- Paris M, Boyer R, Jaenichen R, Wolf J, Karageorgi M, Green J, Cagnon M, Parinello H, Estoup A, Gautier M, Gompel N, Prud'homme B (2020). *Near-chromosome level genome assembly of the*

- fruit pest *Drosophila suzukii* using long-read sequencing. *Scientific reports* **10** (1), 11227. <https://doi.org/10.1038/s41598-020-67373-z>.
- Piper AM, Cunningham JP, Cogan NOI, Blacket MJ (2022). DNA Metabarcoding Enables High-Throughput Detection of Spotted Wing *Drosophila* (*Drosophila suzukii*) Within Unsorted Trap Catches. *Frontiers in Ecology and Evolution* **10**. <https://doi.org/10.3389/fevo.2022.822648>.
- R Core Team (2017). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing. Vienna, Austria. url: <https://www.R-project.org/>.
- Schlötterer C, Tobler R, Kofler R, Nolte V (2014). Sequencing pools of individuals - mining genome-wide polymorphism data without big funding. *Nature Reviews Genetics* **15** (11), 749–763. <https://doi.org/10.1038/nrg3803>.
- Steinegger M, Salzberg SL (2020). Terminating contamination: large-scale search identifies more than 2,000,000 contaminated entries in GenBank. *Genome Biology* **21**, 115. <https://doi.org/10.1186/s13059-020-02023-1>.
- Takamori H, Watabe Ha, Fuyama Y, Zhang Yp, Aotsuka T (2006). *Drosophila subpulchrella*, a new species of the *Drosophila suzukii* species subgroup from Japan and China (Diptera: Drosophilidae). *Entomological Science* **9**, 121–128. <https://doi.org/10.1111/j.1479-8298.2006.00159.x>.
- Wood DE, Lu J, Langmead B (2019). Improved metagenomic analysis with Kraken 2. *Genome Biology* **20**, 257. <https://doi.org/10.1186/s13059-019-1891-0>.
- Wood DE, Salzberg SL (2014). Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome Biology* **15**, R46. <https://doi.org/10.1186/gb-2014-15-3-r46>.
- Zhu Y, Bergland AO, González J, Petrov DA (2012). Empirical Validation of Pooled Whole Genome Population Re-Sequencing in *Drosophila melanogaster*. *PLoS One* **7**, 1–7. <https://doi.org/10.1371/journal.pone.0041901>.