

# The 19 Genomes of *Drosophila*: A BAC Library Resource for Genus-Wide and Genome-Scale Comparative Evolutionary Research

Xiang Song,<sup>\*,1</sup> Jose Luis Goicoechea,<sup>\*,1</sup> Jetty S. S. Ammiraju,<sup>\*,1</sup> Meizhong Luo,<sup>\*,1</sup> Ruifeng He,<sup>\*</sup>  
Jinke Lin,<sup>\*</sup> So-Jeong Lee,<sup>\*</sup> Nicholas Sisneros,<sup>\*</sup> Tom Watts,<sup>†</sup> David A. Kudrna,<sup>\*</sup>  
Wolfgang Golser,<sup>\*</sup> Elizabeth Ashley,<sup>\*</sup> Kristi Collura,<sup>\*</sup> Michele Braidotti,<sup>\*</sup>  
Yeisoo Yu,<sup>\*</sup> Luciano M. Matzkin,<sup>†,‡</sup> Bryant F. McAllister,<sup>§</sup>  
Therese Ann Markow<sup>†,‡,2</sup> and Rod A. Wing<sup>\*,2</sup>

<sup>\*</sup>Arizona Genomics Institute and BIO5 Institute, School of Plant Sciences, and <sup>†</sup>Department of Ecology and Evolutionary Biology, University of Arizona, Tucson, Arizona 85721, <sup>‡</sup>Division of Biological Sciences, University of California at San Diego, La Jolla, California 92093 and <sup>§</sup>Department of Biology, University of Iowa, Iowa City, Iowa 52242

Manuscript received January 6, 2011  
Accepted for publication February 5, 2011

## ABSTRACT

The genus *Drosophila* has been the subject of intense comparative phylogenomics characterization to provide insights into genome evolution under diverse biological and ecological contexts and to functionally annotate the *Drosophila melanogaster* genome, a model system for animal and insect genetics. Recent sequencing of 11 additional *Drosophila* species from various divergence points of the genus is a first step in this direction. However, to fully reap the benefits of this resource, the *Drosophila* community is faced with two critical needs: *i.e.*, the expansion of genomic resources from a much broader range of phylogenetic diversity and the development of additional resources to aid in finishing the existing draft genomes. To address these needs, we report the first synthesis of a comprehensive set of bacterial artificial chromosome (BAC) resources for 19 *Drosophila* species from all three subgenera. Ten libraries were derived from the exact source used to generate 10 of the 12 draft genomes, while the rest were generated from a strategically selected set of species on the basis of salient ecological and life history features and their phylogenetic positions. The majority of the new species have at least one sequenced reference genome for immediate comparative benefit. This 19-BAC library set was rigorously characterized and shown to have large insert sizes (125–168 kb), low nonrecombinant clone content (0.3–5.3%), and deep coverage (9.1–42.9×). Further, we demonstrated the utility of this BAC resource for generating physical maps of targeted loci, refining draft sequence assemblies and identifying potential genomic rearrangements across the phylogeny.

**T**HE genus *Drosophila* contains ~2000 species of diverse morphology, ecology, and behavior that are placed in three major lineages: subgenus *Sophophora*, subgenus *Drosophila*, and subgenus *Dorsilopa* (MARKOW and O'GRADY 2006, 2007). The most widely studied species in the genus, *Drosophila melanogaster*, is firmly established as the premier model system for many biological research areas such as neurobiology, medicine, and population biology (RUBIN and LEWIS 2000). Several other species in this genus, such as *D. pseudoobscura* and *D. virilis*, have also been utilized as genetic model systems particularly for evolutionary studies (ORR and

COYNE 1989; ANDERSON *et al.* 1991; POPADIC and ANDERSON 1994; CHARLESWORTH *et al.* 1997; VIEIRA *et al.* 1997; SWEIGART 2010). Recently, the genomes of *D. melanogaster* and 11 other *Drosophila* species, whose most recent common ancestor occurred >45–50 MYA, have been sequenced, assembled, and annotated (ADAMS *et al.* 2000; MYERS *et al.* 2000; CELNIKER *et al.* 2002; RICHARDS *et al.* 2005; DROSOPHILA 12 GENOMES CONSORTIUM 2007; GILBERT 2007). Species were selected for genome sequencing partly on the basis of their relationship with *D. melanogaster*. Nine of the 12 sequenced genomes were sampled from one subgenus, *Sophophora*, to which *D. melanogaster* belongs, and the remaining 3 are from the *Drosophila* subgenus. These sequences have already greatly improved understanding of the evolution and regulation of eukaryotic genes and genomes through comparative analyses (STARK *et al.* 2007). However, to fully reap the benefits from this unique resource, the *Drosophila* community is faced with two critical needs: first, the development of additional genomics resources

Supporting information is available online at <http://www.genetics.org/cgi/content/full/genetics.111.126540/DC1>.

<sup>1</sup>These authors contributed equally to this work.

<sup>2</sup>Corresponding authors: Division of Biological Sciences, University of California at San Diego, La Jolla, CA 92093. E-mail: tmarkow@ucsd.edu; and Arizona Genomics Institute and BIO5 Institute, School of Plant Sciences, University of Arizona, 1657 E. Helen St., Tucson, AZ 85721. E-mail: rwing@ag.arizona.edu

to aid in finishing the 11 existing draft genome sequences and, second, the generation of additional genomic resources that encompass a much broader range of phylogenetic diversity.

Toward this direction, we constructed a comprehensive set of bacterial artificial chromosome (BAC) libraries for 19 different *Drosophila* species representing a broad spectrum of phylogenetic diversity. BAC libraries are powerful tools for comparative genome research (KIM *et al.* 1996; HOSKINS *et al.* 2000; INTERNATIONAL HUMAN GENOME MAPPING CONSORTIUM 2000a,b; LOCKE *et al.* 2000; OSOEGAWA *et al.* 2000, 2001, 2004; EICHLER and DEJONG 2002; GREGORY *et al.* 2002; GIBBS *et al.* 2003; KRZYWINSKE *et al.* 2004; GONZALEZ *et al.* 2005; AMMIRAJU *et al.* 2006; *DROSOPHILA* 12 GENOMES CONSORTIUM 2007; KIM *et al.* 2008; MURAKAMI *et al.* 2008), especially in taxa containing highly repetitive genomes (HAVLAK *et al.* 2004; ELLISON and SHAW 2010; FANG *et al.* 2010). Genome sequences are available for 10 of 19 species for which BAC libraries are constructed, some of which were instrumental in facilitating sequence assemblies (*DROSOPHILA* 12 GENOMES CONSORTIUM 2007), and they remain a high-priority resource for improving and finishing several of the low coverage draft genome assemblies. BAC libraries for species without sequenced genomes present an important resource for positional cloning and large-scale targeted comparative genome analyses.

We selected 19 species within three lineages of the genus *Drosophila* for BAC library construction (Figure 1). These species shared a common ancestor ~40–60 MYA (POWELL 1997) and were selected because of their varied evolutionary distances from *D. melanogaster* and other sequenced species, their diverse ecologies and life history characters, and the fact that they can be reared in the laboratory and used in experimental work in the future. Ten BAC libraries were constructed as a resource for generating BAC end mate-pair sequence to assist in the assembly of whole-genome shotgun sequences and for enabling future genomic research (*DROSOPHILA* 12 GENOMES CONSORTIUM 2007). Beyond those 10 species, we are interested in generating BAC library resources for representative species of lineages not yet targeted for sequencing but that fill in large phylogenetic gaps. The majority of these species have at least one previously sequenced reference genome for immediate comparative benefit. In addition, this new set of species facilitates the “ladder and constellation” approach of modified phylogenetic shadowing proposed by Clark *et al.* ([http://flybase.org/static\\_pages/news/whitepapers/GenomesWP2003.pdf](http://flybase.org/static_pages/news/whitepapers/GenomesWP2003.pdf)) for annotating genome data. In this approach ladder rungs constitute successively increasing divergence points and constellations are clusters of species attaching to these divergence points. This set of 19 BAC libraries documented here will further advance the genus *Drosophila* as an ideal eukaryotic comparative genomics system designed to (1) provide sequencing resources for comparative annotation of the

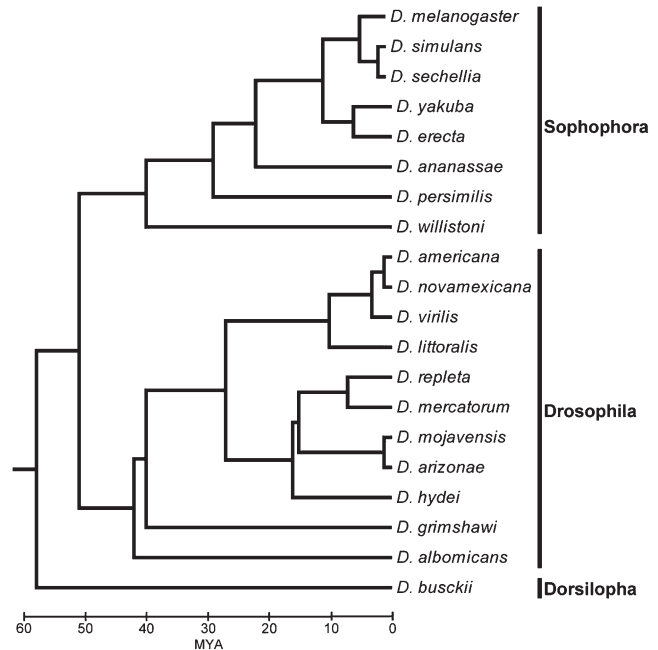


FIGURE 1.—Phylogenetic tree of 19 species and *D. melanogaster* selected for the *Drosophila* BAC resource project. The phylogenetic relationships and approximate divergence times among the *Drosophila* species in our study were determined from a compilation of prior analyses (PITNICK *et al.* 1995; MARKOW and O'GRADY 2006; *DROSOPHILA* 12 GENOMES CONSORTIUM 2007).

*D. melanogaster* genome and (2) provide genomic resources for experimental investigation of gene function throughout the genus *Drosophila*.

## MATERIALS AND METHODS

**Fly culturing and embryo collection:** Fly cultures were expanded on banana/opuntia medium (<http://flyfood.arl.arizona.edu/opuntia.php3>) and healthy sexually mature adult flies were introduced into plexiglass oviposition chambers kept on a 16/8 light/dark cycle at 24°–25° with a relative humidity of 60–80%. Exceptions to this procedure were *D. littoralis*, *D. novamexicana*, *D. americana*, *D. grimshawi*, and *D. persimilis* cultures, which were oviposited at 20°–22°, whereas *D. albomicans* was oviposited at 17°. Medium for *D. sechellia* was supplemented with 0.5% (v/v) hexanoic acid and 0.5% (v/v) octanoic acid to stimulate oviposition. Oviposition medium for *D. grimshawi* was supplemented with 2% (w/v) methylparaben to prevent overgrowth of fungus. *D. busckii* and *D. grimshawi* cultures were grown on Wheeler–Clayton medium (<http://flyfood.arl.arizona.edu/wheeler.php3>). *D. grimshawi* adults were separated by sex until day of placement in the oviposition chamber to enhance embryo production. Adult flies were allowed to oviposit on a given plate for as long as possible without larval hatch. This interval varied between 4 and 48 hr depending on the species. About 1.2–1.5 g wet weight embryos were pooled in batches and stored at –80° at the end of each oviposition session.

**Nuclei preparation and BAC library construction:** Embryos were gently homogenized in PBS buffer (0.76% NaCl, 4 mM NaH<sub>2</sub>PO<sub>4</sub>, 9 mM Na<sub>2</sub>HPO<sub>4</sub>, pH 7.0) using a Dounce Tissue Grinder (Wheaton Science), centrifuged at 4° at 1430 × *g* for 15 min, and resuspended in PBS buffer. The suspension

was then mixed with an equal volume of 1% InCert Agarose (CAMBREX) (in PBS buffer) at 45° and transferred into plug molds. Treatment of plugs to produce unsheared megabase-size DNA was as described (LUO and WING 2003). BAC libraries were constructed as previously described (LUO and WING 2003; AMMIRAJU *et al.* 2006).

**BAC library characterization:** DNA from a random sample of 260–480 BAC clones from each library was isolated, restriction digested with *NotI*, and run on CHEF gels for insert size determination as previously described (LUO and WING 2003; AMMIRAJU *et al.* 2006).

High colony density hybridization filters for each library were prepared using Genetix Q-bots (Genetix) as described previously (AMMIRAJU *et al.* 2006; LUO *et al.* 2006). Nine gene-specific probes were chosen that represented all chromosomes of *D. melanogaster* (supporting information, Table S1 and Table S2). All probe DNA fragments were PCR amplified from the *D. mojavensis* genome and gel purified using a QIAEX II (QIAGEN, Valencia, CA) kit. Table S1 lists the primer sequences used for each probe. Purified DNA fragments were sequenced and similarity searches were conducted to validate their specificity. Probes were prepared by labeling with [<sup>32</sup>P]dCTP using a DecaprimeII random prime labeling kit (Ambion), and hybridizations were carried out as described by AMMIRAJU *et al.* (2006). Positive clones were picked and rearranged onto colony filters, followed by a secondary hybridization with individual probes.

**Fingerprinting and contig assembly:** Positive hybridization clones were fingerprinted using SNaPshot (LUO *et al.* 2003; KIM *et al.* 2008) and assembled into contigs with FPC v 8.5.2 (SODERLUND *et al.* 2000; [www.agcol.arizona.edu](http://www.agcol.arizona.edu)) at a fixed tolerance value 4 and an initial Sulston score  $1e^{-50}$  (AMMIRAJU *et al.* 2006).

**BAC end sequencing and *in silico* analysis:** Fingerprinted BAC clones were end sequenced with a universal T7 primer (5' TAA TAC GAC TCA CTA TAG GG 3') and a custom primer BES\_HR (5' CAC TCA TTA GGC ACC CCA 3') following previously described methods (KIM *et al.* 2008). BAC end sequences (BES) were submitted to GenBank with the following accession numbers: *D. simulans*, EI211963.1–EI212067.1; *D. sechellia*, CZ549016.1–CZ549204.1; *D. yakuba*, EI89369.1–EI189559.1; *D. erecta*, CZ548656.1–CZ548834.1; *D. ananassae*, CZ548467.1–CZ548655.1; *D. persimilis*, EI188778.1–EI189177.1; *D. willistoni*, EI189178.1–EI189368.1; *D. americana*, EI189178.1–EI189368.1; *D. novamexicana*, DU169152.1–DU169329.1; *D. virilis*, CZ549205.1–CZ549371.1; *D. littoralis*, EI211597.1–EI211779.1; *D. repleta*, EI211780.1–EI211962.1; *D. mercatorum*, EI188452.1–EI188610.1; *D. mojavensis*, CZ548835.1–CZ549015.1; *D. arizonae*, EI211417.1–EI211231.1; *D. hydei*, EI188451.1–EI188450.1; *D. grimshawi*, EI188111.1–EI188299.1; *D. albomicans*, EI211043–EI211230.1; and *D. busckii*, EI211418.1–EI211596.1.

All BESs were masked with Repeat Masker (version 3.1.0) against a redundant repeat database with sequences obtained from FlyBase ([www.FlyBase.org](http://www.FlyBase.org)) and Repbase ([www.girinst.org](http://www.girinst.org)). These sequences were used to conduct BLAST analysis against the mitochondrial (NC\_001709, 19,517 bp) and nuclear genome sequences of *D. melanogaster* (Build 5.1) and the freeze 1 genome assemblies from the remaining 11 species (<http://rana.lbl.gov/drosophila/cafl.html> and <http://insects.eugenes.org/species/data/>). To compensate for the lack of whole-genome sequences and to minimize the bias of sequence divergence, the genome sequences of *D. virilis* and *D. mojavensis* were used as pseudoreference sequences for the *D. virilis* and *D. repleta* species group, respectively. BES from *D. albomicans* and *D. busckii* was compared to the *D. grimshawi* sequences.

In addition, similarity searches were conducted with complete gene sequences of each probe against the 12 *Drosophila* whole-

genome sequences (*DROSOPHILA* 12 GENOMES CONSORTIUM 2007). Homologs with a minimum alignment length of 100 bp and 75% of nucleotide identity were retained for further analysis and for a comparison of their presence or absence in FPC-derived contigs.

## RESULTS AND DISCUSSION

***Drosophila* strain selection and genome sizes:** Several criteria were used for careful evaluation of the different *Drosophila* species strains used for BAC resource development in this study. First, all fly lines were inbred for a minimum of eight generations by sib–sib mating to reduce the extent of heterozygosity and subsequently sequenced at six nuclear loci to verify homozygosity (T. A. Markow, unpublished data). Second, to minimize endosymbiont contamination (*Wolbachia* spp. and *Spiroplasma* spp.) at least five adult fly DNA samples from each species were screened with established protocols (MATEOS *et al.* 2006). Finally, species identity was confirmed by both morphological and molecular approaches. When a suitable nuclear or mitochondrial DNA marker was known for a species, that marker was amplified, sequenced, and validated. Additionally, salivary gland chromosomes from third instar larvae were prepared and inspected for inversion polymorphism microscopically. Only homokaryotypic lines were used. All strains (Table 1) are deposited in the University of California at San Diego *Drosophila* Stock Center and are publicly available as a community resource.

Genome size of an organism is the most important factor in determining the depth of a genomic library (reviewed in GREGORY 2005). Previously determined genome sizes (BOSCO *et al.* 2007) were used in this study for estimating the coverage of the BAC libraries for different *Drosophila* species. BOSCO *et al.* (2007) employed two nucleic-acid-binding fluorescent dyes, propidium iodide (PI) and 4',6-diamidino-2-phenylindole (DAPI), in conjunction with flow cytometry to determine genome sizes of 38 species of Drosophilidae, including the 12 sequenced *Drosophila* species (*DROSOPHILA* 12 GENOMES CONSORTIUM 2007).

The genome sizes of 15 of the 19 *Drosophila* species used in this study were based on the PI method and the remaining species (*D. novamexicana*, *D. littoralis*, *D. repleta*, and *D. busckii*) genome sizes were based on the DAPI method alone (for which the PI data were not available) (Table 1). Nine of the *Drosophila* species strains were not the same as the strains analyzed by BOSCO *et al.* (2007). An important finding to consider, as reported by BOSCO *et al.* (2007) and GREGORY and JOHNSTON (2008), is that DAPI may overestimate genome size, which could affect the estimated genome coverage of these four libraries.

Genome sizes of two species, *D. arizonae* and *D. albomicans*, were not known, so the genome sizes of closest relatives *D. mojavensis* and *D. immigrans*, respectively, were applied to estimate the tentative genome coverages of

**TABLE 1**  
**Characteristics of the 19 *Drosophila* BAC library set**

Species	Group <sup>a</sup>	Stock no. <sup>b</sup>	Library name	Enzyme	Genome size (Mb)	Average insert size (kb)	Clone no.	Calculated genome coverage <sup>c</sup>
<i>D. simulans</i>	MEL	DSSC 14021-0251.195	DS_ABa	<i>Hind</i> III	160 <sup>d</sup>	158	18,432	18.2
<i>D. sechellia</i>	MEL	DSSC 14021-0248.25	DS_Ba	<i>Hind</i> III	166 <sup>d</sup>	139	18,432	15.4
<i>D. yakuba</i>	MEL	DSSC 14021-0261.01	DY_Ba	<i>Hind</i> III	188 <sup>d</sup>	148	11,520	9.1
<i>D. erecta</i>	MEL	DSSC 14021-0224.01	DE_TBa	<i>Hind</i> III	145 <sup>d</sup>	149	18,432	18.9
<i>D. ananassae</i>	MEL	DSSC 14024-0371.13	DA_Ba	<i>Bam</i> HI	215 <sup>a</sup>	148	36,864	25.4
<i>D. persimilis</i>	OBS	DSSC 14011-0111.49	DP_Ba	<i>Hind</i> III	183 <sup>d</sup>	151	18,432	15.2
<i>D. willistoni</i>	WIL	DSSC 14030-0811.24	DW_Ba	<i>Hind</i> III	206 <sup>d</sup>	150	18,432	13.4
<i>D. americana</i>	VIR	DSSC 15010-0951.15	DA_ABa	<i>Bst</i> YI	275 <sup>d</sup>	136	11,520	5.7
<i>D. novamexicana</i>	VIR	DSSC 15010-1031.14	DN_Ba	<i>Hind</i> III	244 <sup>e</sup>	155	13,440	8.5
<i>D. virilis</i>	VIR	DSSC 15010-1051.87	DV_VBa	<i>Bst</i> YI	404 <sup>d</sup>	127	55,296	17.4
<i>D. littoralis</i>	VIR	DSSC 15010-1001.11	DL_Ba	<i>Hind</i> III	238 <sup>b</sup>	168	36,864	26
<i>D. repleta</i>	REP	DSSC 15084-1611.10	DR_Ba	<i>Hind</i> III	167 <sup>e</sup>	143	36,864	31.6
<i>D. mercatorum</i>	REP	DSSC 15082-1521.36	DM_Ba	<i>Hind</i> III	128 <sup>d</sup>	125	18,432	18
<i>D. mojavensis</i>	REP	DSSC 15081-1352.22	DM_CBa	<i>Bam</i> HI	152 <sup>d</sup>	143	30,720	28.9
<i>D. arizonae</i>	REP	DSSC 15081-1271.27	DA_CBa	<i>Hind</i> III	152 <sup>f</sup>	133	18,432	16.1
<i>D. hydei</i>	REP	DSSC 15085-1641.58	DH_Ba	<i>Hind</i> III	164 <sup>d</sup>	146	36,864	32.8
<i>D. grimshawi</i>	HAW	DSSC 15287-2541.00	DG_Ba	<i>Hind</i> III	231 <sup>d</sup>	127	18,432	10.1
<i>D. albomicans</i>	IMM	DSSC 15112-1751.08	DA_BBa	<i>Hind</i> III	299 <sup>f</sup>	130	18,432	8
<i>D. busckii</i>	DOR	DSSC 13000-0081.31	DB_Ba	<i>Hind</i> III	194 <sup>e</sup>	166	18,432	15.8

<sup>a</sup> MEL, *melanogaster*; OBS, *obscura*; WIL, *willistoni*; VIR, *virilis*; REP, *repleta*; HAW, Hawaiian; IMM, *immigrans*; DOR, subgenus *Drosophila*.

<sup>b</sup> DSSC: *Drosophila* Species Stock Center.

<sup>c</sup> Calculated genome coverage: by insert size, genome size, and number of clones in the library.

<sup>d</sup> Genome size measured by the PI method (Bosco *et al.* 2007).

<sup>e</sup> Genome size measured by the DAPI method (Bosco *et al.* 2007).

<sup>f</sup> Genome sizes of *D. arizonae* and *D. albomicans* were adopted from the genome size of close relatives, *D. mojavensis* and *D. immigrans*, respectively.

their respective BAC libraries. The genome sizes among the 19 *Drosophila* species varied by ~3.2-fold, with the smallest being *D. mercatorum* and the largest *D. virilis* (Table 1).

**BAC library construction and characterization:** Three different restriction enzymes were used for BAC library construction: *Hind*III, *Bam*HI, and *Bst*YI. Fifteen of the 19 libraries were constructed from DNA partially digested with *Hind*III, followed by size selection and ligation into the *Hind*III site of pIndigoBAC536.SwaI (AMMIRAJU *et al.* 2006) (Table 1). Two libraries each were generated similarly from *Bam*HI (*D. ananassae* and *D. mojavensis*) and *Bst*YI (*D. virilis* and *D. americana*) restriction digests. All libraries, except for the *D. busckii* library (two ligations), were built from single ligations. The number of clones in the 19-BAC library set ranged between 11,520 and 55,296 (Table 1), which were arrayed into 384-well microtiter plates for long-term storage in -80° freezers at the Arizona Genomics Institute's (AGI) BAC/EST Resource Center ([www.genome.arizona.edu](http://www.genome.arizona.edu)).

Insert sizes of individual clones in each library ranged from 10 to 371 kb, with the majority >120 kb (Figure 2). The average insert sizes of these libraries ranged from 125 to 168 kb (Table 1). Percentages of non-insert-

containing clones ranged between 0.3 and 5.3%, which is typical for BAC libraries constructed at AGI (AMMIRAJU *et al.* 2006).

#### Genomic redundancy of the *Drosophila* BAC libraries:

We estimated the genomic depth of the 19 *Drosophila* BAC set by three different, but complementary approaches. First, we estimated the redundancy of each library empirically from the average insert size, total number of clones, and the genome size of the corresponding lineage, which ranged approximately between 5.7- and 32.8-fold (Table 1). To assess the randomness and extent of representational heterogeneity for different genomic regions, we screened the entire set of 19 *Drosophila* BAC libraries with nine gene-specific probes in two successive rounds of hybridizations (MATERIALS AND METHODS, Table S1, and Table S2).

In brief, 4196 putative positive BAC clones were identified in the first round of hybridization, and 3809 (91%) were confirmed by a second hybridization. The number of positive hits per library ranged from 1 to 108 (Table S3). At least one positive hit per each probe was detected for all the libraries with the exceptions of the *D. americana*, *D. repleta*, and *D. hydei* libraries for probe X-CG11387 and *D. ananassae* for probe 3R-CG31247

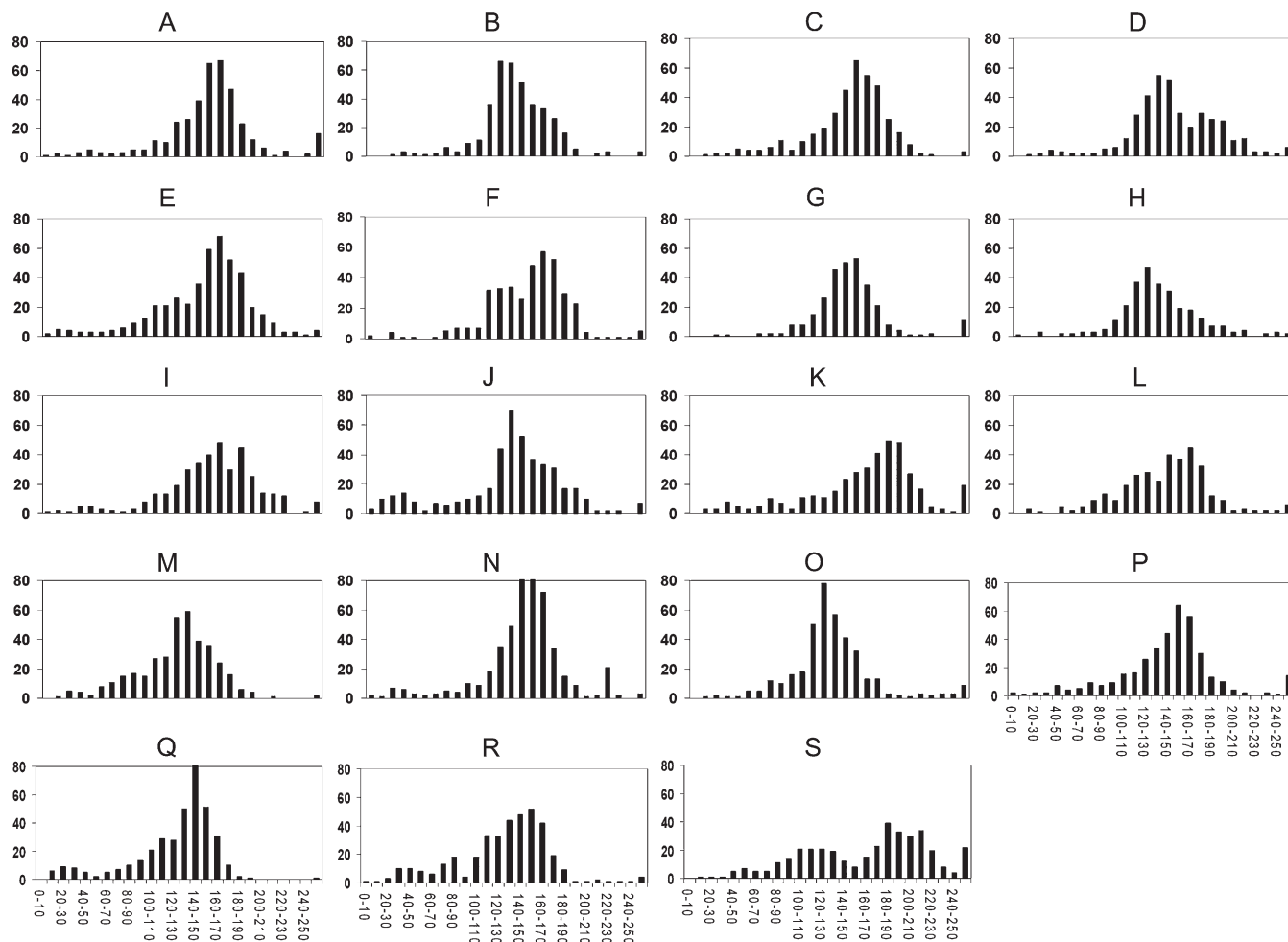


FIGURE 2.—Insert size distribution of 19 *Drosophila* BAC libraries. Histograms A–S depict the insert size distribution in the 19 different libraries. For each histogram, the *x*-axis represents insert size (kilobases) and the *y*-axis represents the number of clones in a particular insert size range. (A) *D. simulans* (DS\_ABa), average insert size 158 kb; (B) *D. sechellia* (DS\_Ba), average insert size 139 kb; (C) *D. yakuba* (DY\_Ba), average insert size 148 kb; (D) *D. erecta* (DE\_TBa), average insert size 149 kb; (E) *D. ananassae* (DA\_Ba), average insert size 148 kb; (F) *D. persimilis* (DP\_Ba), average insert size 151 kb; (G) *D. willistoni* (DW\_Ba), average insert size 150 kb; (H) *D. americana* (DA\_ABa), average insert size 136 kb; (I) *D. novamexicana* (DN\_Ba), average insert size 155 kb; (J) *D. virilis* (DV\_VBa), average insert size 127 kb; (K) *D. littoralis* (DL\_Ba), average insert size 168 kb; (L) *D. repleta* (DR\_Ba), average insert size 143 kb; (M) *D. mercatorum* (DM\_Ba), average insert size 125 kb; (N) *D. mojavensis* (DM\_CBa), average insert size 143 kb; (O) *D. arizonae* (DA\_CBa), average insert size 133 kb; (P) *D. hydei* (DH\_Ba), average insert size 146 kb; (Q) *D. grimshawi* (DG\_Ba), average insert size 127 kb; (R) *D. albomicans* (DA\_BBa), average insert size 130 kb; (S) *D. busckii* (DB\_Ba), average insert size 166 kb.

(Table S3). In these four species no hits were found, even upon three rounds of library screening, with different hybridization stringencies. For *D. ananassae*, the whole-genome draft sequence was available (<http://rana.lbl.gov/drosophila/cafl.html>), and similarity searches revealed the presence of the probe sequence (3R-CG31247; Table S2) in the draft sequence assembly. Therefore, at least in the case of *D. ananassae*, it appears that methodological and/or library coverage issues prevented recovery of this gene via the hybridization-based approach, possibly due to use of heterologous probes, multiple usage of high-density colony filters, or cloning bias (under- and overrepresentation of genomic regions due to usage of a single restriction enzyme during library construction). More data are required to

confirm the absence of the gene X-CG11387 in the other three species (*D. americana*, *D. repleta*, and *D. hydei*).

Hybridization-based genome coverages ranged from  $9.1\times$  (*D. americana*) to  $42.9\times$  (*D. hydei*). In only two species, *D. mercatorum* and *D. willistoni*, the hybridization-based coverage was slightly lower than expected (Table 2). The remaining 17 libraries had either nearly equal or higher coverage than predicted (Table 2 and Table S3). The *D. albomicans* BAC library showed an  $\sim 3.6$ -fold higher than expected coverage on the basis of hybridization (Table 2), which could have resulted from not having accurate genome size estimation for this species (Table 1).

A third and a more rigorous approach using fingerprinted contig (FPC)-based estimations of genomic redundancy of BAC libraries was applied, using a similar

TABLE 2

**A comparison of genomic redundancies of each *Drosophila* BAC library as estimated by empirical, hybridization, and FPC approaches**

Species	Calculated genome coverage <sup>a</sup>	Average hybridization coverage <sup>b</sup>	FPC, general <sup>c</sup>	Ratio of <i>a:b:c</i>
<i>D. simulans</i>	18.2	25.0	17	1:1.4:0.94
<i>D. sechellia</i>	15.4	20.2	14	1:1.3:0.88
<i>D. yakuba</i>	9.1	11.0	9	1:1.2:1.01
<i>D. erecta</i>	18.9	19.7	14	1:1.0:0.75
<i>D. ananassae</i>	25.4	25.3	22	1:1.0:0.87
<i>D. persimilis</i>	15.2	18.3	13	1:1.2:0.86
<i>D. willistoni</i>	13.4	9.6	7	1:0.7:0.52
<i>D. americana</i>	5.7	9.1	8	1:1.6:1.36
<i>D. novamexicana</i>	8.5	14.8	13	1:1.7:1.48
<i>D. virilis</i>	17.4	32.7	19	1:1.9:1.11
<i>D. littoralis</i>	26	25.1	18	1:1.0:0.71
<i>D. repleta</i>	31.6	35.7	14	1:1.1:0.44
<i>D. mercatorum</i>	18	11.7	10	1:0.6:0.54
<i>D. mojavensis</i>	28.9	31.1	17	1:1.1:0.59
<i>D. arizonae</i>	16.1	20.2	10	1:1.3:0.63
<i>D. hydei</i>	32.8	42.9	37	1:1.3:1.12
<i>D. grimshawi</i>	10.1	14.2	9	1:1.4:0.87
<i>D. albomicans</i>	8	28.4	10	1:3.6:1.22
<i>D. busckii</i>	15.8	28.2	9	1:1.8:0.58

<sup>a</sup>Theoretical coverage of each *Drosophila* library from Table 1.

<sup>b</sup>Average hybridization coverage: total number of clones detected by two rounds of hybridization divided by the total number of loci, from Table S3.

<sup>c</sup>FPC-based estimate of genomic redundancy of each *Drosophila* library: total number clones in each FPC assembly divided by the total number of contigs, from Table S4 and Table S5.

strategy to our previous analysis of a set of 11 *Oryza* (cultivated and wild rice) BAC libraries (AMMIRAJU *et al.* 2006). This approach can discriminate the unavoidable cloning bias from those of cross-hybridizations and genetic rearrangements such as duplications. All 3809 hybridization-derived BAC clones were fingerprinted and 3005 (79%) successful fingerprints were assembled into physical contigs (Table S4 and Table S5). Under a scenario of single-copy probes and one contig per probe for each species, the theoretically expected number of contigs is 171 (nine probes for 19 libraries). However, several exceptions were found: (a) as described above, 1 probe, X-CG11387, had no hits in the *D. Americana*, *D. repleta*, and *D. hydei* libraries, and another probe, 3R-CG31247, had no hits in the *D. ananassae* library (Table S3); (b) clones detected from six hybridizations (*D. yakuba*, *D. persimilis*, and *D. willistoni* with probe X-CG11387; *D. mercatorum* with probe 2L-CG4128; and *D. mercatorum* and *D. grimshawi* with probe 4-CG2999) resulted in the presence of singletons (Table S5) (all these instances resulted in less than three positive clones, Table S3). Taking into account the absence of these contigs in these species, 161 contigs are expected.

Our FPC analysis revealed a total of 211 contigs, 50 additional contigs than the expected number of 161 (Table S4). The number of contigs and respective coverage differed among different *Drosophila* libraries for the same probe (Table S5). Five probes (X-CG11387, X-CG32611, 3L-CG10948, 3R-CG31247, and 4-CG2999) essentially behaved as single-copy probes in most *Drosophila* libraries (Table S5). The remaining four detected on average  $\geq 1.4$  contigs per probe (Table S5). To better understand whether these deviations from expectation (50 additional contigs) were due to technical issues (cross-hybridization and assembly artifacts) and/or lineage-specific genetic changes, we gathered data from two additional experiments. First, on the basis of BES mapping information (MATERIALS AND METHODS), we classified 142 contigs as primary (those that map to the expected genomic location) and 69 additional contigs as secondary (27 contigs that cannot be positioned in any genome and 42 contigs that map to nonorthologous locations), a good agreement between the results of FPC analyses and mapping information (Table S2 and Table S6).

Second, nucleotide and protein similarity searches of the probe (or gene) sequences revealed that several secondary sites (17/42 secondary contigs) contained small cross-hybridizing paralogous sequences (Table S6, indicated with \*). It is possible that the 25 remaining secondary sites also contained very small cross-hybridizing sequences that were not easily detected through similarity searches. In addition, sequence analysis of the extended flanking sequences of the primary sites with the secondary sites revealed no evidence of synteny, suggesting cross-hybridization as the main cause for these additional contigs.

To provide a conservative estimate of genome coverage, we considered each identified contig as an independent locus and calculated a weighted FPC coverage that accounts for the presence of several loci (Table S4; AMMIRAJU *et al.* 2006). Estimated FPC coverage for the 19 libraries (Table 2 and Table S4) ranged between 7 $\times$  and 37 $\times$ . Only 2 libraries had coverage  $< 9\times$ : *D. willistoni* (7 $\times$ ) and *D. americana* (8 $\times$ ).

Twelve libraries showed a ratio close to 1:1 between the FPC and empirically estimated coverage (Table 2). The *D. willistoni*, *D. littoralis*, *D. repleta*, *D. mercatorum*, *D. mojavensis*, *D. arizonae*, and *D. busckii* libraries showed ratios  $\leq 0.7:1$  (Table 2 and Table S4). The difference between hybridization-based and contig-based estimates of library coverage is due to the difference in the number of loci used to calculate the coverage. While each probe is considered as a single locus in the hybridization-based approach, each secondary contig is considered as an independent locus in the FPC-based approach (Table 2, Table S3, and Table S4). Together, these results showcase the high quality and deep representational coverage of each of 19 *Drosophila* genomes in their respective libraries.

**Utilization of BAC libraries:** Although a few Drosophila BAC libraries have already been reported in the literature (HOSKINS *et al.* 2000; LOCKE *et al.* 2000; GONZALEZ *et al.* 2005; OSOEGAWA *et al.* 2007; MURAKAMI *et al.* 2008), this is the first synthesis and characterization of a comprehensive set of BAC library resources for the genus, which fills a critical void for the Drosophila research community. Hybridization of nine different probes to the full set of libraries demonstrates the feasibility of isolating homologous regions across the entire genus. Combined with high-throughput sequencing methods (WICKER *et al.* 2006), this set of libraries provides an excellent resource for comparative studies of targeted genomic regions (*e.g.*, LEUNG *et al.* 2010).

First, BAC libraries from species that do not yet have a reference genome sequence themselves provide a source for identifying genome rearrangements in comparisons with the available genome sequences. For example, end sequences of BACs isolated with the X-linked probe CG32611 from *D. novamexicana* map at an unexpected position within contig 12,970 of *D. virilis*, indicating a putative small inversion at the base of the X chromosome that had not been previously identified (VIEIRA *et al.* 1997). Another putative inversion was also revealed in *D. arizonae* by the localization of end sequences of clones hybridizing to CG3139 in the genome sequence of *D. mojavensis*. Targeted analyses inversion breakpoints are also enabled by the availability of these BAC libraries and informed by the reference genome sequences. EVANS *et al.* (2007) used cytological evidence on the position of an inversion in *D. americana* to develop probes for isolating its breakpoints from the respective BAC clones. In addition, the BAC libraries for the nine unsequenced Drosophila species provide robust templates for the whole-genome physical and sequence frameworks. In this direction, the entire *D. persimilis* BAC library was fingerprinted, bidirectionally end sequenced, and assembled into a whole-genome physical map. This map was aligned to the *D. persimilis* and *D. pseudoobscura* draft sequences and is currently under editing (data not shown).

An extremely important application of the BAC resources reported here is in the ability to use functional genomics to test genes underlying the differences between Drosophila species. The tool kit for functional analyses of Drosophila has taken a major leap forward with the recent establishment of the P/ $\Phi$ C31 artificial chromosome manipulation (P[acman]) transgenesis platform (VENKEN *et al.* 2006, 2009; VENKEN and BELLEN 2007). While still reliant on the P-transposable element for transformation, this BAC transgenic system significantly improves upon the size of the DNA to be carried in the vector (>130 kb) and its site-specific integration in the fly genome. An important feature of the P[acman] system is recombineering, which permits cloning/transfer of large DNA fragments from existing Drosophila P1 or BAC clones through a homologous recombination-

mediated gap repair process. Therefore, a combination of the P[acman] system with the 19 Drosophila BAC libraries will provide an unprecedented opportunity to the fly community to access, transfer, and manipulate virtually any genomic region of interest (large genes or even gene clusters) covering the entire phylogenomic range of the genus Drosophila.

Finally, the BAC library set reported here can be used to further improve many of the existing Drosophila draft sequence assemblies (*DROSOPHILA 12 GENOMES CONSORTIUM* 2007) and aid in the characterization of lineage-specific rearrangements. For example, physical mapping of BAC contigs, or individual BAC clones, identified by hybridization probes designed from draft Drosophila genome sequences, has revealed and confirmed chromosomal location of several sequence contigs from the draft assemblies, as well as their relationship to *D. melanogaster* (Table S6). Conserved linkage and physical markers were used to infer the physical organization of the assembled genome assemblies relative to reference chromosome maps (SCHAEFFER *et al.* 2008), and these BAC libraries serve as an appropriate resource to isolate regions at inferred gaps between adjacent contigs (*e.g.*, HOSKINS *et al.* 2000). Using hybridization to recover genome regions containing target genes, combined with end sequencing of positive clones, further reveals the conserved linkage among Drosophila species. For example, scaffolds 20 and 24 map to X[A], 29 to 3L[D], and 30 to 4[F] in *D. sechellia*; 4512 to 4[F] in *D. erecta*, 12,984 to 3R[B] and 12,947 to 4(LR)[F] in *D. ananassae*, 48 to XR[D/A] and 103 to 5[F] in *D. persimilis*; 5 group M to 5[F] in *D. pseudoobscura*; 13,052 to 6[F] in *D. virilis* (*DROSOPHILA 12 GENOMES CONSORTIUM* 2007); 6,498 to 6[F] in *D. mojavensis*; and 14,822 to 6[F] in *D. grimshawi* (Table S6).

These libraries are likely to facilitate a wide array of comparative, evolutionary, and functional genomics studies and play a major role in advancing the Drosophila biology.

This work was supported by National Institutes of Health grant U1HG02525A.

#### LITERATURE CITED

- ADAMS, M. D., S. E. CELNIKER, R. A. HOLT, C. A. EVANS, J. D. GOCAYNE *et al.*, 2000 The genome sequence of *Drosophila melanogaster*. *Science* **287**: 2185–2195.
- AMMIRAJU, J. S., M. LUO, J. L. GOICOEHEA, W. WANG, D. KUDRNA *et al.*, 2006 The *Oryza* bacterial artificial chromosome library resource: construction and analysis of 12 deep-coverage large-insert BAC libraries that represent the 10 genome types of the genus *Oryza*. *Genome Res.* **16**: 140–147.
- ANDERSON, W. W., J. ARNOLD, D. G. BALDWIN, A. T. BECKENBACH, C. J. BROWN *et al.*, 1991 Four decades of inversion polymorphism in *Drosophila pseudoobscura*. *Proc. Natl. Acad. Sci. USA* **88**: 10367–10371.
- BOSCO, G., P. CAMPBELL, J. T. LEIVA-NETO and T. A. MARKOW, 2007 Analysis of *Drosophila* species genome size and satellite DNA content reveals significant differences among strains as well as between species. *Genetics* **177**: 1277–1290.

- CELNIKER, S. E., D. A. WHEELER, B. KRONMILLER, J. W. CARLSON, A. HALPERN *et al.*, 2002 Finishing a whole-genome shotgun: release 3 of the *Drosophila melanogaster* euchromatic genome sequence. *Genome Biol.* **3**: research0079.1-0079.14.
- CHARLESWORTH, B., D. CHARLESWORTH, J. HNILICKA, A. YU and D. S. GUTTMAN, 1997 Lack of degeneration of loci on the neo-Y chromosome of *Drosophila americana*. *Genetics* **145**: 989-1002.
- DROSOPHILA 12 GENOMES CONSORTIUM, 2007 Evolution of genes and genomes on the *Drosophila* phylogeny. *Nature* **450**: 203-218.
- EICHLER, E. E., and P. J. DEJONG, 2002 Biomedical applications and studies of molecular evolution: a proposal for a primate genomic library resource. *Genome Res.* **12**: 673-678.
- ELLISON, C. K., and K. L. SHAW, 2010 Mining non-model genomic libraries for microsatellites: BAC versus EST libraries and the generation of allelic richness. *BMC Genomics* **11**: 428-436.
- EVANS, A. L., P. A. MENA and B. F. MCALLISTER, 2007 Positive selection near an inversion breakpoint on the neo-X chromosome of *Drosophila americana*. *Genetics* **177**: 1303-1319.
- FANG, G., B. P. BLACKMON, D. C. HENRY, M. E. STATON, C. A. SASKI *et al.*, 2010 Genomic tools development for *Aquilegia*: construction of a BAC-based physical map. *BMC Genomics* **11**: 621-628.
- GIBBS, R. A., G. M. WEINSTOCK, M. L. METZKER, D. M. BUZNY, E. J. SODERGRÉN *et al.*, 2003 Genome sequence of the brown Norway rat yields insights into mammalian evolution. *Nature* **428**: 493-521.
- GILBERT, D. G., 2007 DroSpGe: rapid access database for new *Drosophila* species genomes. *Nucleic Acids Res.* **35**: D480-D485.
- GONZALEZ, J., M. NEFEDOV, I. BOSDET, F. CASALS, O. CALVETE *et al.*, 2005 A BAC-based physical map of the *Drosophila buzzatii* genome. *Genome Res.* **15**: 885-892.
- GREGORY, S. G., M. SEKHON, J. SCHEIN, S. ZHAO, K. OSOEGAWA *et al.*, 2002 A physical map of the mouse genome. *Nature* **418**: 743-750.
- GREGORY, T. R., 2005 Synergy between sequence and size in large-scale genomics. *Nat. Rev. Genet.* **6**: 699-708.
- GREGORY, T. R., and J. S. JOHNSTON, 2008 Genome size diversity in the family Drosophilidae. *Heredity* **101**: 228-238.
- HAVLAK, P., R. CHEN, K. J. DURBIN, A. EGAN, Y. REN *et al.*, 2004 The Atlas genome assembly system. *Genome Res.* **14**: 721-732.
- HOSKINS, R. A., C. R. NELSON, B. P. BERMAN, T. R. LAVERTY, R. A. GEORGE *et al.*, 2000 A BAC-based physical map of the major autosomes of *Drosophila melanogaster*. *Science* **287**: 2271-2274.
- INTERNATIONAL HUMAN GENOME MAPPING CONSORTIUM, 2000a A physical map of the human genome. *Nature* **409**: 934-941.
- INTERNATIONAL HUMAN GENOME MAPPING CONSORTIUM, 2000b Initial sequencing and analysis of the human genome. *Nature* **409**: 860-921.
- KIM, H., B. HURWITZ, Y. YU, K. COLLURA, M. GILL *et al.*, 2008 Construction, alignment and analysis of twelve framework physical maps that represent the ten genome types of the genus *Oryza*. *Genome Biol.* **9**: R45.
- KIM, U. J., B. W. BIRREN, T. SLEPAK, V. MANCINA, C. BOYSEN *et al.*, 1996 Construction and characterization of a human bacterial artificial chromosome library. *Genomics* **34**: 213-218.
- KRZYWINSKE, M., J. WALLIS, C. GOSELE, I. BOSDET, R. CHIU *et al.*, 2004 Integrated and sequence-ordered BAC- and YAC-based physical maps for the rat genome. *Genome Res.* **14**: 766-779.
- LEUNG, W., C. D. SHAFFER, T. CORDONNIER, J. WONG, M. S. ITANO *et al.*, 2010 Evolution of a distinct genomic domain in *Drosophila*: comparative analysis of the dot chromosome in *Drosophila melanogaster* and *Drosophila virilis*. *Genetics* **185**: 1519-1534.
- LOCKE, J., L. PODEMSKI, N. AIPPERSBACH, H. KEMP and R. HODGETTS, 2000 A physical map of the polytenized region (101EF-102F) of chromosome 4 in *Drosophila melanogaster*. *Genetics* **155**: 1175-1183.
- LUO, M., and R. A. WING, 2003 An improved method for plant BAC library construction, pp. 3-20 in *Plant Functional Genomics*, edited by E. GROTEWOLD. Humana Press, Totowa, NJ.
- LUO, M., H. KIM, D. KUDRNA, N. B. SISNEROS, S. LEE *et al.*, 2006 Construction of a nurse shark (*Ginglymostoma cirratum*) bacterial artificial chromosome (BAC) library and a preliminary genome survey. *BMC Genomics* **7**: 106.
- LUO, M. C., C. THOMAS, F. M. YOU, J. HSIAO, S. OUYANG *et al.*, 2003 High-throughput fingerprinting of bacterial artificial chromosomes using the snapshot labeling kit and sizing of restriction fragments by capillary electrophoresis. *Genomics* **82**: 378-389.
- MARKOW, T. A., and P. M. O'GRADY, 2006 Evolutionary genetics of reproductive behavior in *Drosophila*: connecting the dots. *Annu. Rev. Genet.* **39**: 263-291.
- MARKOW, T. A., and P. M. O'GRADY, 2007 *Drosophila* biology in the genomic age. *Genetics* **177**: 1269-1276.
- MATEOS, M., S. J. CASTREZANA, B. J. NNAKIVELL, A. M. ESTES, T. A. MARKOW *et al.*, 2006 Heritable endosymbionts of *Drosophila*. *Genetics* **174**: 363-376.
- MURAKAMI, K., A. TOYODA, M. HATTORI, Y. KUROKI, A. FUJIYAMA *et al.*, 2008 BAC library construction and BAC end sequencing of five *Drosophila* species: the comparative map with the *D. melanogaster* genome. *Genes Genet. Syst.* **83**: 245-246.
- MYERS, E. W., G. G. SUTTON, A. L. DELCHER, I. M. DEW, D. P. FASULO *et al.*, 2000 A whole-genome assembly of *Drosophila*. *Science* **287**: 2196-2204.
- OSOEGAWA, K., A. M. TATENO, P. Y. WOON, E. FRENGEN, A. G. MAMMOSER *et al.*, 2000 Bacterial artificial chromosome libraries for mouse sequencing and functional analysis. *Genome Res.* **10**: 116-128.
- OSOEGAWA, K., A. G. MAMMOSER, C. WU, E. FRENGEN, C. ZENG *et al.*, 2001 A bacterial artificial chromosome library for sequencing the complete human genome. *Genome Res.* **11**: 483-496.
- OSOEGAWA, K., B. ZHU, C. L. SHU, T. REN, Q. CAO *et al.*, 2004 BAC resources for the rat genome project. *Genome Res.* **14**: 780-785.
- OSOEGAWA, K., G. M. VESSERE, C. L. SHU, R. A. HOSKINS, J. P. ABAD *et al.*, 2007 BAC clones generated from sheared DNA. *Genomics* **89**: 291-299.
- ORR, H. A., and J. A. COYNE, 1989 The genetics of postzygotic isolation in the *Drosophila virilis* group. *Genetics* **121**: 527-537.
- PITNICK, S., T. A. MARKOW and G. S. SPICER, 1995 Delayed male maturity is a cost of producing large sperm in *Drosophila*. *Proc. Natl. Acad. Sci. USA* **92**: 10614-10618.
- POPADIC, A., and W. W. ANDERSON, 1994 The history of a genetic system. *Proc. Natl. Acad. Sci. USA* **91**: 6819-6823.
- POWELL, J. R., 1997 *Progress and Prospects in Evolutionary Biology: The Drosophila Model*. Oxford University Press, London/New York/Oxford.
- RICHARDS, S., Y. LIU, B. R. BETTENCOURT, P. HRADECKY, S. LETOVSKY *et al.*, 2005 Comparative genome sequencing of *Drosophila pseudoobscura*: chromosomal, geni, and cis-element evolution. *Genome Res.* **15**: 1-18.
- RUBIN, G. M., and E. B. LEWIS, 2000 A brief history of *Drosophila*'s contributions to genome research. *Science* **287**: 2216-2218.
- SCHAEFFER, S. W., A. BHUTKAR, B. F. MCALLISTER, M. MATSUDA, L. M. MATZKIN *et al.*, 2008 Polytenic chromosomal maps of 11 *Drosophila* species: the order of genomic scaffolds inferred from genetic and physical maps. *Genetics* **179**: 1601-1655.
- SODERLUND, C., S. HUMPHRAY, A. DUNHAM and L. FRENCH, 2000 Contigs built with fingerprints, markers, and FPC V4.7. *Genome Res.* **10**: 1772-1787.
- STARK, A., M. F. LIN, P. KHERADPOUR, J. S. PEDERSEN, L. PARTS *et al.*, 2007 Discovery of functional elements in 12 *Drosophila* genomes using evolutionary signatures. *Nature* **450**: 219-232.
- SWEIGART, A. L., 2010 Simple Y-autosomal incompatibilities cause hybrid male sterility in reciprocal crosses between *Drosophila virilis* and *D. americana*. *Genetics* **184**: 779-787.
- VENKEN, K. J. T., and H. J. BELLEN, 2007 Transgenesis upgrades for *Drosophila melanogaster*. *Development* **134**: 3571-3584.
- VENKEN, K. J. T., Y. HE, R. A. HOSKINS and H. J. BELLEN, 2006 P[acman]: a BAC transgenic platform for targeted insertion of large DNA fragments in *D. melanogaster*. *Science* **314**: 1747-1751.
- VENKEN, K. J. T., J. W. CARLSON, K. L. SCHULZE, H. PAN, Y. HE *et al.*, 2009 Versatile P[acman] BAC libraries for transgenesis studies in *Drosophila melanogaster*. *Nat. Methods* **6**: 431-434.
- VIEIRA, J., C. P. VIEIRA, D. L. HARTL and E. R. LOZOVSKAYA, 1997 Discordant rates of chromosome evolution in the *Drosophila virilis* species group. *Genetics* **147**: 223-230.
- WICKER, T., E. SCHLAGENHAUF, A. GRANER, T. J. CLOSE, B. KELLER *et al.*, 2006 454 sequencing put to the test using the complex genome of barley. *BMC Genomics* **7**: 275.



# GENETICS

## Supporting Information

<http://www.genetics.org/cgi/content/full/genetics.111.126540/DC1>

### **The 19 Genomes of *Drosophila*: A BAC Library Resource for Genus-Wide and Genome-Scale Comparative Evolutionary Research**

**Xiang Song, Jose Luis Goicoechea, Jetty S. S. Ammiraju, Meizhong Luo, Ruifeng He, Jinke Lin, So-Jeong Lee, Nicholas Sisneros, Tom Watts, David A. Kudrna, Wolfgang Golser, Elizabeth Ashley, Kristi Collura, Michele Braidotti, Yeisoo Yu, Luciano M. Matzkin, Bryant F. McAllister, Therese Ann Markow and Rod A. Wing\***

Copyright © 2011 by the Genetics Society of America

DOI: 10.1534/genetics.111.126540

**TABLE S1****List of Primer Sequence**

Gene	Primer sequence
Ch2L_CG3139-PA	F: CAAAAATCTGCTCATCAACTTCA R: TAAGGGTTGAGGGTGCATTT
Ch2R_CG14747-PA	F: CAATGCTGCCATTTGAGAAG R: AACACTGCACGAACACGAAG
Ch3R_CG31247-PA	F: CGATATGCCCAAACAATTC R: AGCTGCTGAATCGAGCTTTC
Ch3L_CG32206-PB	F: TTTATCAGCTCCCACTCAG R: ACCAAATCAGGTCACCAGGA
Ch4_CG2999-PA	F: CTCAAGGCCTGATAGCGAAG R: AAAACACAAAGAAAGCGGAAA
Ch2L_CG4128-PA	F: ATAATTTAGCGGGATGAGG R: TTCATTTGCAATGTTGGTC
Ch3L_CG10948_PC	F: TTGCATTATTGTTTCAGTCACTCAG R: TGCCGTAATACATTCTTTGAACA
ChX_CG32611_PB	F: GCGTCAAGTGATCCGAATAG R: CAGCTAGGCTGCTTGGAGAC
ChX_CG11387_PA	F: TTCATACAGACAGCCCACGA R: TTCATACAGACAGCCCACGA

**TABLES S2 and S3**

Tables S2 and S3 are available for download as Excel Files at <http://www.genetics.org/cgi/content/full/10.1534/genetics.111.126540>.

**TABLE S4****Characterization and coverage estimates of *Drosophila* BAC library set via contig analysis**

Species	Group	FPC						Ratio		
		FPC	Clones	Contig	Clones			FPC	Experiment	FPC : Exp
		ID	Total	Total	In Contig	%	Singles	Coverage	coverage	Coverage
<i>D. simulans</i>	MEL	DS_AB	222	12	205	92.3	17	17	18.2	0.9 : 1
<i>D. sechellia</i>	MEL	DS__B	165	12	163	98.8	2	14	15.4	0.9 : 1
<i>D. yakuba</i>	MEL	DY__B	88	9	83	94.3	5	9	9.1	1.0 : 1
<i>D. erecta</i>	MEL	DE_TB	162	11	156	96.3	6	14	18.9	0.8 : 1
<i>D. ananassae</i>	MEL	DA__B	215	9	198	92.1	17	22	25.4	0.9 : 1
<i>D. persimilis</i>	OBS	DP__B	137	10	131	95.6	6	13	15.2	0.9 : 1
<i>D. willistoni</i>	WIL	DW__B	83	11	76	91.6	7	7	13.4	0.5 : 1
<i>D. americana</i>	VIR	DA_AB	75	9	70	93.3	5	8	5.7	1.4 : 1
<i>D. novamexicana</i>	VIR	DN__B	135	10	126	93.3	9	13	8.5	1.5 : 1
<i>D. virilis</i>	VIR	DV_VB	219	10	194	88.6	25	19	17.4	1.1 : 1
<i>D. littoralis</i>	VIR	DL__B	216	11	203	94.0	13	18	26	0.7 : 1
<i>D. repleta</i>	REP	DR__B	185	11	151	81.6	34	14	31.6	0.4 : 1
<i>D. mercatorum</i>	REP	DM__B	107	9	88	82.2	19	10	18	0.5 : 1
<i>D. mojavensis</i>	REP	DM_CB	358	20	339	94.7	19	17	28.9	0.6 : 1
<i>D. arizonae</i>	REP	DA_CB	103	9	92	89.3	11	10	16.1	0.6 : 1
<i>D. hydei</i>	REP	DH__B	379	10	366	96.6	13	37	32.8	1.1 : 1
<i>D. grimshawi</i>	HAW	DG__B	163	14	123	75.5	40	9	10.1	0.9 : 1
<i>D. albomicans</i>	IMM	DA_BB	179	14	143	79.9	36	10	8	1.2 : 1
<i>D. busckii</i>	DOR	DB__B	104	10	91	87.5	13	9	15.8	0.6 : 1

MEL : melanogaster

OBS: obscura

WIL: willistoni

VIR: virilis

REP: repleta

HAW: hawaiian

IMM: immigrans

DOR: subgenus dorsilopha

Ratio FPC : Emp(irical) coverage calculated with data from table 1

**TABLES S5 and S6**

Tables S5 and S6 are available for download as Excel Files at <http://www.genetics.org/cgi/content/full/10.1534/genetics.111.126540>.