

# McClintock: An Integrated Pipeline for Detecting Transposable Element Insertions in Whole-Genome Shotgun Sequencing Data

Michael G. Nelson, Raquel S. Linheiro,<sup>1</sup> and Casey M. Bergman<sup>2</sup>

Faculty of Life Sciences, University of Manchester, M13 9PL, United Kingdom

ORCID IDs: 0000-0002-2567-2634 (M.G.N.); 0000-0003-2659-0910 (R.S.L.); 0000-0002-5462-9854 (C.M.B.)

**ABSTRACT** Transposable element (TE) insertions are among the most challenging types of variants to detect in genomic data because of their repetitive nature and complex mechanisms of replication. Nevertheless, the recent availability of large resequencing data sets has spurred the development of many new methods to detect TE insertions in whole-genome shotgun sequences. Here we report an integrated bioinformatics pipeline for the detection of TE insertions in whole-genome shotgun data, called McClintock (<https://github.com/bergmanlab/mcclintock>), which automatically runs and standardizes output for multiple TE detection methods. We demonstrate the utility of McClintock by evaluating six TE detection methods using simulated and real genome data from the model microbial eukaryote, *Saccharomyces cerevisiae*. We find substantial variation among McClintock component methods in their ability to detect nonreference TEs in the yeast genome, but show that nonreference TEs at nearly all biologically realistic locations can be detected in simulated data by combining multiple methods that use split-read and read-pair evidence. In general, our results reveal that split-read methods detect fewer nonreference TE insertions than read-pair methods, but generally have much higher positional accuracy. Analysis of a large sample of real yeast genomes reveals that most McClintock component methods can recover known aspects of TE biology in yeast such as the transpositional activity status of families, target preferences, and target site duplication structure, albeit with varying levels of accuracy. Our work provides a general framework for integrating and analyzing results from multiple TE detection methods, as well as useful guidance for researchers studying TEs in yeast resequencing data.

## KEYWORDS

transposable elements  
bioinformatics  
genomics  
yeast

The widespread availability of genomic data over the last two decades has provided unparalleled opportunities to learn about the abundance, diversity, and functional consequences of transposable

elements (TEs) in modern genomes. However, the computational analysis of TE sequences in both reference and resequenced genomes remains a challenging area of bioinformatics research because of the repetitive nature of these sequences. Development of bioinformatics tools for the detection and annotation of TEs in reference genomes is now a relatively mature field (Bergman and Quesneville 2007; Saha *et al.* 2008; Lerat 2010), although many open questions remain about choosing the best tools for specific biological applications (Hoen *et al.* 2015). In contrast, detection of reference and nonreference TE insertions in whole-genome shotgun (WGS) resequencing data are an active research area (reviewed in Ewing 2015), with a large number of methods published in recent years (Sackton *et al.* 2009; Ewing and Kazazian 2010, 2011; Hormozdiari *et al.* 2010; Quinlan *et al.* 2010; Fiston-Lavier *et al.* 2011, 2015; Kofler *et al.* 2012, 2016; Lee *et al.* 2012; Linheiro and Bergman 2012; Nellaker *et al.* 2012; Platzer *et al.* 2012; Chen *et al.* 2013, 2017; Cridland *et al.* 2013; Robb *et al.* 2013; Gilly *et al.* 2014; Nakagome *et al.* 2014; Thung *et al.* 2014; Wu *et al.* 2014; Zhuang *et al.*

Copyright © 2017 Nelson *et al.*

doi: <https://doi.org/10.1534/g3.117.043893>

Manuscript received February 14, 2017; accepted for publication June 15, 2017; published Early Online June 21, 2017.

This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Supplemental material is available online at [www.g3journal.org/lookup/suppl/doi:10.1534/g3.117.043893/-/DC1](http://www.g3journal.org/lookup/suppl/doi:10.1534/g3.117.043893/-/DC1).

<sup>1</sup>Present address: Centro de Biotecnologia e Química Fina, Universidade Católica Portuguesa, Porto, 4202-401, Portugal.

<sup>2</sup>Corresponding author: Department of Genetics and Institute of Bioinformatics, University of Georgia, 120 E. Green St., Athens, GA 30602. E-mail: [cbergman@uga.edu](mailto:cbergman@uga.edu)

2014; Hawkey *et al.* 2015; Hénaff *et al.* 2015; Jiang *et al.* 2015; Rahman *et al.* 2015; Quadrana *et al.* 2016).

Because of the wide array of available methods, it remains unclear which method for detecting TEs in resequenced genomes is best suited for particular genomic problems, leading to substantial investigator effort in terms of installation and testing, or the application of suboptimal bioinformatic approaches. Most papers reporting new methods to detect reference or nonreference TEs in WGS data provide some measure of their own performance relative to using simulations, benchmark genomic data, or PCR-based validation. However, only a handful of papers have reported new methods that include performance evaluation relative to other methods (Gilly *et al.* 2014; Ewing 2015; Hawkey *et al.* 2015; Hénaff *et al.* 2015; Jiang *et al.* 2015; Rahman *et al.* 2015; Chen *et al.* 2017), and these are often limited in scope to only a single organism or TE family. In addition to being incomplete, comparative analysis of bioinformatic systems in papers that report new methods can fall victim to the “self-assessment trap” (Norel *et al.* 2011). Moreover, there is no common format for the annotation of nonreference TE insertions (Bergman 2012; Rishishwar *et al.* 2016), making direct comparison of predictions from different methods more challenging. Recently, Rishishwar *et al.* (2016) performed an independent comparative evaluation of seven WGS-based TE detection methods using human genomic data, which revealed many method-specific predictions and recommended combining the results of multiple systems followed by manual curation (see also Ewing 2015). Rishishwar *et al.* (2016) also highlighted the challenges users face when installing and running multiple TE detection methods, and provide helpful advice for users and developers.

As a step toward a fully automated framework for running and evaluating multiple methods to detect TEs in WGS resequencing data, we have developed an integrated pipeline called McClintock (<https://github.com/bergmanlab/mcclintock>) that generates standardized output for multiple WGS-based TE detection methods. The primary goal of the McClintock pipeline is to lower the barrier to installation, use, and evaluation of multiple WGS-based TE detection methods. Several key features of the McClintock pipeline are that it automates formatting of key input files and standardizes output of multiple TE detection methods to allow easy comparisons of results from different systems, as recommended by Rishishwar *et al.* (2016). In the initial version of McClintock, we incorporate six complementary TE detection methods that make predictions based on split-read- or read-pair-based evidence in Illumina WGS data. Here we describe the McClintock system and its component methods, and perform a comparative evaluation using simulated and real yeast genome data. Our analysis supports previous conclusions that no single TE detection method provides a comprehensive detection of nonreference TEs (Ewing 2015; Rishishwar *et al.* 2016), but provides a framework for further testing, development, and integration to achieve this ultimate aim, as well as useful guidance for yeast researchers to select appropriate TE detection tools.

## MATERIALS AND METHODS

### Analysis of simulated WGS data sets with single artificial TE insertions

To investigate the performance of McClintock component methods on data containing known, nonreference TE insertions, we created simulated *Saccharomyces cerevisiae* genomes, each containing a single synthetic nonreference TE insertion from one of the four active TE families in an otherwise unmodified *S. cerevisiae* reference genome. Since active *S. cerevisiae* TEs (*Ty1*, *Ty2*, *Ty3*, and *Ty4*) are known to target tRNA genes (Ji *et al.* 1993; Chalker and Sandmeyer 1990, 1992; Devine and Boeke 1996; Kim *et al.* 1998; Baller *et al.* 2012; Mularoni *et al.* 2012; Qi

*et al.* 2012), each synthetic insertion was placed upstream of a different annotated tRNA gene in the reference genome, taking the orientation of the tRNA gene into consideration. The annotation for 299 tRNAs was extracted from the SGD genome annotation for *sacCer2* (SGD version R61.1.1). *Ty1*, *Ty2*, and *Ty4* have been shown to insert predominantly within the first 200 bp upstream of tRNA genes, and *Ty3* appears to target more specifically the region of RNA polymerase III transcription initiation, 16 or 17 nucleotides from the 5' ends of tRNA genes (Ji *et al.* 1993; Chalker and Sandmeyer 1990, 1992; Devine and Boeke 1996; Kim *et al.* 1998; Baller *et al.* 2012; Mularoni *et al.* 2012; Qi *et al.* 2012). All active *S. cerevisiae* TEs produce 5-bp target site duplications (TSDs) on insertion (Gafner and Philippsen 1980; Rinckel and Garfinkel 1996; Chalker and Sandmeyer 1990; Kim *et al.* 1998; Zou *et al.* 1996). To mimic these insertion preferences in our simulations, *Ty1*, *Ty2*, *Ty3*, and *Ty4* were alternately selected for insertion; a 5-bp TSD was created (either 200–195 bp upstream of a tRNA gene for *Ty1*, *Ty2*, and *Ty4*; or 17–12 bp upstream of tRNA genes for *Ty3*); and the corresponding full-length *Ty* canonical sequence was inserted in the reference genome. 299 insertions were produced with the TE sequence inserted on the positive strand of the genome, and 299 were produced with the TE sequence reverse complemented to test the effects of TE orientation on method performance.

We simulated resequencing of single-insertion synthetic genomes using Wgsim (<https://github.com/lh3/wgsim>) (Li *et al.* 2009) with a 1% base error rate ( $-e$  0.01). Read lengths were chosen to be 101 bases each with an insert size of 300 bp (42-bp SE) and 100× coverage to mimic the properties of a large sample of WGS data sets collected by Strobe *et al.* (2015), which we use in our analysis of real yeast genomes (see below). To generate an average read depth of 100× across the length of *sacCer2* reference genomes with additional single TE insertions, *in silico* WGS samples were created with 6,024,220 read pairs for *Ty1* insertions, 6,024,237 read pairs for *Ty2* insertions, 6,023,936 read pairs for *Ty3* insertions, and 6,024,369 read pairs for a *Ty4* insertion.

McClintock (version e945d20da22dc1186b97960b44b86bc21c96ac27) was run on each of these simulated data sets using reference TE annotations and canonical TE sequences from Carr *et al.* (2012), plus a manually produced hierarchy file based on the reference TE annotation in Carr *et al.* (2012). We used the standard, unmodified reference genome sequence option of McClintock for these single synthetic insertion simulations. The mean of the number of nonreference and reference TEs predicted per sample was calculated across all 299 simulated samples for each strand. The proportion of correct predictions of nonreference TEs was calculated at four thresholds of accuracy: (i) requiring the exact TSD to be annotated correctly, (ii) requiring a prediction to be within a 100-bp window either side of the TSD, (iii) within a 300-bp window either side of the TSD (the insert size of the simulated sequencing data set), or (iv) within a 500-bp window either side of the TSD. BEDtools window (Quinlan and Hall 2010) was used to calculate correct predictions within the given windows. A prediction was classified as exactly correct only if the same TE family was predicted to occur at the exact coordinates of the TSD of the synthetic TE insertion location. For nonexact overlaps, BEDtools window allows a permissive definition of a true positive, where a correct TE prediction is counted when any part of a predicted insertion falls within the given threshold distance if the correct TE family is predicted. The orientation of a predicted insertion was not taken into account for determining a correct prediction because some methods do not predict orientation.

To visualize the accuracy of nonreference TE predictions, the results files for the 299-positive strand and 299-negative strand single-insertion samples were converted into two bigWig files (one for each strand) using BEDtools and wigToBigWig (Kent *et al.* 2010). This was performed for

each TE family and each component method of McClintock. SeqPlots (Stempor 2014) was then used to produce plots of the genome coverage of predictions for each TE family, centered around the simulated insertion locations for that family. Visualization of predicted insertions negative strand simulations were reverse complemented and depicted on the same plot as positive strand simulations in different colors. Plots were centered on the 5-bp TSD and extended  $\pm 10$  bp for split-read methods, and  $\pm 500$  bp for read-pair methods, respectively. Results for TEMP were partitioned based on whether or not split-read support was available for a prediction. Prior to visualization, we attempted to filter out any obvious false-positive predictions using the fact that each synthetic insertion location should only be predicted in one simulated sample. Thus, any locations where a predicted nonreference insertion was observed across multiple simulated samples indicated a potential false positive. This filtering was necessary to prevent a nonreference insertion that was predicted by RelocaTE in the same location in 149 single synthetic insertion samples from dominating the visualization for this component method. False-positive filtering prior to visualization only affected five other potential insertions for PoPoolationTE, and thus this filtering procedure does not substantially alter positional accuracy results. To further investigate the accuracy of TSDs predicted by split-read methods, the length of the predicted TSD was plotted for each active yeast TE family. To be consistent with analysis of real yeast genomes (see below) and to mitigate effects of false-positive predictions found at the same site in multiple samples, TSD lengths predicted in simulated data were only plotted for unique insertion sites rather than all insertions.

To investigate the concordance of nonreference TE predictions made by different McClintock component methods, we first determined whether or not each method had made a “correct” prediction in each of the simulated samples with a synthetic TE insertion. Predictions for ngs\_te\_mapper, RelocaTE, TEMP (both split read and read pair), and PoPoolationTE were classified as correct if they overlapped with the true location of the TSD. Predictions for RetroSeq and TE-locate were classified as correct if they occurred within a 100- or 500-bp window of the correct location of the TSD, respectively. The orientation of a prediction was not taken into account when classifying a prediction as “correct” or not, because not all methods predict orientation. The overlap of these correct predictions was then plotted as Venn diagrams using jvenn (Bardou *et al.* 2014), comparing split-read methods, read-pair methods, and finally the total set of correct predictions from all split-read vs. all read-pair methods.

### Analysis of real WGS data sets

To assess the relative performance characteristics of the component methods on real data, McClintock was run on a large sample of *S. cerevisiae* data sets from Strobe *et al.* (2015) that includes 93 *S. cerevisiae* strains from different geographical locations and clinical origins. The Strobe *et al.* (2015) samples were sequenced on an Illumina HiSeq 2000 with paired-end reads of 101 bases each, an average insert size of 300 bases, and a median coverage of  $>117\times$ . We used these general library characteristics in our single synthetic insertion simulations (see above) to allow more direct comparison with analysis of these real yeast genomes. The raw fastq files for the 93 sequenced strains were obtained from the EBI Sequence Read Archive (SRA072302).

McClintock (version 354acec977e37c354f6f05046940b0dabf09b331) was run on each of these samples using reference TE annotations and canonical TE sequences from Carr *et al.* (2012), and a manually produced hierarchy file based on the annotation in Carr *et al.* (2012). The McClintock version used for the analysis of real yeast data differs slightly from that used for simulated data in terms of three small improvements that were required to handle variation in sample names (for ngs\_te\_mapper)

and differences in read lengths of paired-end fragments (for PoPoolationTE) which were encountered when analyzing real yeast genome data. We used the standard, unmodified reference genome sequence option of McClintock for these analyses. The average number of non-reference and reference TEs predicted per strain was plotted as box plots for each method. In addition, the total numbers of nonreference and reference TE insertions per TE family were summarized across all strains for each McClintock component method, both genome wide and in tRNA gene regions.

To biologically validate results of different component methods of McClintock, we took advantage of the fact that *Ty* elements are known to insert in close proximity to tRNA genes in *S. cerevisiae* (Ji *et al.* 1993; Chalker and Sandmeyer 1990, 1992; Devine and Boeke 1996; Kim *et al.* 1998; Baller *et al.* 2012; Mularoni *et al.* 2012; Qi *et al.* 2012). A prediction was counted as within a tRNA gene region if any part of the annotation was within 1000 bp upstream or 500 bp downstream of the transcription start site of one of the 299 annotated tRNA genes, taking tRNA gene orientation into account. To visualize the patterns of non-reference TE predictions around tRNA genes, all results for all 93 samples were converted to a single genome coverage bigWig file for each TE and each component method. SeqPlots (Stempor 2014) was used to produce plots of the genome coverage averaged across the 299 tRNA genes. Plots were centered on the start of the tRNA gene and extended 1000 bp upstream and 500 bp downstream, taking into account the orientation of each tRNA gene. Results for TEMP were subset into two groups based on whether split-read support for a prediction was available or not.

The lengths of TSDs for nonreference TE insertions predicted by the split-read methods were plotted by TE family. To prevent any non-reference TE insertions present at the same location in multiple samples from biasing the results, only unique insertion sites were plotted. If a method called an insertion at nearly the same location but with a longer or shorter TSD in different samples, these were classed as unique sites.

### Data availability

The McClintock pipeline is available under the FreeBSD license at <https://github.com/bergmanlab/mcclintock>. Supplemental Material, File S1 contains a combined supplement text including: detailed descriptions of McClintock components; an overview of the McClintock execution process; details of postprocessing of component method output; methods, results, and discussion for analysis of McClintock applied to simulated resequencing data created for unmodified *S. cerevisiae* reference genomes; and Figures S1–S4 and Tables S1 and S2. Supporting data sets of McClintock predictions for real yeast genomes in SRA072302 are available in File S2. Code used to generate simulated yeast genomes and apply McClintock to simulated and real yeast genome data are provided in File S3.

## RESULTS AND DISCUSSION

### McClintock component methods and their dependencies

We initiated our design of McClintock with a literature search for candidate bioinformatic systems that can detect TE insertions from NGS data in 2014, which yielded 33 potential systems. Our main project objective was to develop a system that automatically detects nonreference TE insertions in raw WGS data for any species. Thus, we excluded systems that required any wet-laboratory enrichment from further consideration. Systems that did not make their code available were also rejected. This left a list of 12 candidate software systems. After preliminary testing of these 12 methods, six were rejected from further testing because of difficulties during installation [Tangram (Wu *et al.*

■ **Table 1** An overview of the features of the component TE detection methods in the McClintock pipeline

Method	ngs_te_mapper	RelocaTE	TEMP	RetroSeq	PoPoolationTE	TE-locate
Split read	✓	✓	✓			
Read pair			✓	✓	✓	✓
Nonreference TEs	✓	✓	✓	✓	✓	✓
Reference TEs	✓	✓	✓ <sup>a</sup>		✓	✓
Orientation	✓	✓ <sup>b</sup>	✓			✓ <sup>c</sup>
TSD	✓	✓	✓ <sup>d</sup>			
Detects TE families not in reference genome	✓	✓	✓	✓ <sup>e</sup>	✓	

Split read and read pair refer to what type of evidence is used to make TE-insertion predictions (see main text for details).

<sup>a</sup>TEMP reports whether a reference TE is absent from the resequenced sample rather than providing direct evidence for the presence of a reference TE.

<sup>b</sup>RelocaTE output provides information about the orientation of nonreference TEs, but not for reference TEs. McClintock annotates the orientation of reference TEs in RelocaTE output using the original reference TE annotation.

<sup>c</sup>TE-locate provides information about the orientation of nonreference TEs where possible, but not for reference TEs. McClintock annotates the orientation of reference TEs in TE-locate output using the original reference TE annotation.

<sup>d</sup>TEMP only makes TSD predictions for insertions with split-read support.

<sup>e</sup>RetroSeq can detect TE families not present in the reference genome when using Exonerate to generate a reference TE annotation, but not when using a user-supplied reference TE annotation which is the default option in McClintock.

2014) and VariationHunter (Hormozdiari *et al.* 2010)], reliance on data for a specific organism [TEA (Lee *et al.* 2012) and VirusSeq (Chen *et al.* 2013)], inability to detect nonreference insertions [T-lex (Fiston-Lavier *et al.* 2011)], or the inability to distinguish general structural variations from TE insertions [HYDRA (Quinlan *et al.* 2010)]. Six remaining methods [ngs\_te\_mapper (Linheiro and Bergman 2012), TE-locate (Platzer *et al.* 2012), PoPoolationTE (Kofler *et al.* 2012), RetroSeq (Keane *et al.* 2013), RelocaTE (Robb *et al.* 2013), and TEMP (Zhuang *et al.* 2014)] had publicly available code that could be installed reproducibly and met project objectives were selected for incorporation into the initial McClintock pipeline. Since the original selection of methods for inclusion in McClintock, a number of additional methods that meet the initial project requirements [“pecnv teclust” (Cridland *et al.* 2013), TIF (Nakagome *et al.* 2014), TE-Tracker (Gilly *et al.* 2014), Mobster (Thung *et al.* 2014), ITIS (Jiang *et al.* 2015), Jitterbug (Hénaff *et al.* 2015), TIDAL (Rahman *et al.* 2015), ISmapper (Hawkey *et al.* 2015), MELT (Sudmant *et al.* 2015), SPLITREADER (Quadrona *et al.* 2016), and TEPID (Stuart *et al.* 2016)] and new versions of some methods [PoPoolationTE2 (Kofler *et al.* 2016) and RelocaTE2 (Chen *et al.* 2017)] have been released. These methods have not yet been incorporated into McClintock, but the flexible architecture of our system permits their inclusion in the future.

A summary of the main features of the six component methods included in McClintock is shown in Table 1. A more detailed overview of the component methods, their original use case, software/data dependencies, and limitations is provided in the “Description of McClintock Component Methods” section of File S1. While none of the McClintock component methods were originally designed for detecting TEs in yeast, using the yeast system as a test bed does not favor any particular component method and realistically models application of component systems to a new species. The six component systems each have many dependencies on other pieces of software, which must all be correctly installed before the component system will function correctly. These software dependencies are listed in Table 2. Several of these component dependencies require end-user licenses, and thus it was not possible to fully automate installation of all component methods. McClintock therefore assumes component dependencies are installed system wide, but automates installation of the component methods themselves. A passive check is performed during installation of McClintock that reports whether component dependencies are available, though installation is not halted if they are missing. Because of the large number of component dependencies and subsequent development of components

themselves, we developed McClintock to use specific versions of components and their dependencies. Table 2 also lists the version of each dependency that was used with McClintock to obtain the results presented here.

McClintock component methods also have a variety of data dependencies that are required as inputs, which are listed in Table 3. The component methods incorporated into McClintock together require a total of 13 different data dependencies to run. However, since many of these data dependencies can be automatically generated or are format alterations that can be automatically achieved with simple preprocessing steps, the number of data dependencies can be reduced to three required inputs for McClintock: a fasta file of the reference genome, a fasta file of the canonical TE sequences, and fastq files of NGS reads (paired or single ended).

### The McClintock pipeline

An overview of the data flow and processing steps performed by the McClintock pipeline is shown in Figure 1. A detailed description of how the McClintock pipeline is executed can be found in the “Overview of the McClintock process” section of File S1. In the following sections, we describe the options for running the McClintock pipeline, then describe how component methods are parsed in the context of the McClintock pipeline to create standardized output for downstream analysis.

**Reference TE annotation options:** Several McClintock component systems rely on information about TEs in the reference genome as part of their workflow, which can be either supplied by the user or automatically generated by McClintock. If a preexisting annotation of the TE sequences in the reference genome is available, a one-based GFF file of this data can be used as input for the McClintock pipeline. If such a reference TE annotation is provided, then the user must also create and supply a TE “hierarchy” file as another input. The hierarchy file contains two tab-delimited columns, the first listing the name of each instance in the reference TE annotation and the second listing the canonical TE family to which that instance belongs. If no reference TE annotation is provided, then a reference TE annotation and hierarchy file is created automatically by running RepeatMasker and postprocessing RepeatMasker output files.

**Reference genome sequence options:** McClintock provides options to automatically create various different modified reference genomes. These options were implemented because some component methods



■ **Table 2 Software dependencies required to install and run each component TE detection method in the McClintock pipeline**

Software	ngs_te_mapper	RelocaTE	TEMP	RetroSeq	PoPoolationTE	TE-locate	Version Used in this Study
Linux	✓	✓	✓	✓	✓	✓	CentOS 6
Perl		✓	✓	✓	✓	✓	5.18.1
R (R Core Team 2013)	✓						3.0.2
BioPerl (Stajich et al. 2002)		✓	✓				1.006001
RepeatMasker (Smit et al. 2013)					✓		4.0.2
BEDTools (Quinlan and Hall 2010)			✓	✓			2.17.0
SAMTools (Li et al. 2009)		✓	✓ <sup>a</sup>	✓	✓		0.1.19-44428cd
BCFTools (Li et al. 2009)				✓			0.1.19-44428cd
twoBitToFa (Kuhn et al. 2013)			✓				294
BLAT (Kent 2002)		✓					35 × 1
Exonerate (Slater and Birney 2005)				✓			2.2.0
Bowtie (Langmead et al. 2009)		✓					1.0.0
BWA (Li 2013)	✓		✓	✓	✓	✓	0.7.4-r385 <sup>b</sup>

<sup>a</sup>Only compatible with SAMTools 0.1.19 or earlier (Rishishwar et al. 2016).

<sup>b</sup>This specific version of BWA is needed to ensure compatibility between PoPoolationTE, which uses BWA-ALN, and other component methods which use BWA-MEM.

(RetroSeq and TE-locate) require an instance of a TE to exist in the reference genome for nonreference instances of that family to be detected in a resequenced sample. This is important because, in some cases, like the *Drosophila melanogaster* P-element (Kaminker et al. 2002), the reference genome does not include any copies of a TE family that occurs in natural populations. This situation may also occur when a TE family has been introduced experimentally into a strain lacking that TE to study its transposition. To allow for these cases, McClintock has an option to generate modified reference genomes that include additional “chromosomes” comprised of canonical TE sequences or TE sequences extracted from the reference genome. An annotation of TEs in the additional “chromosomes” is then appended to the reference TE annotation file. PoPoolationTE requires a modified reference genome with canonical TE sequences and reference TE sequences added as additional “chromosomes.” Thus these reference genome modifications are always made specifically for PoPoolationTE, regardless of whether user-supplied options to modify the reference genome are provided globally for other component methods.

**Run options:** McClintock offers additional options to customize the way the pipeline is run. It is possible to specify which component methods are executed, allowing tailored output and shorter run times. McClintock and its component methods produce short-read alignment files and other intermediate files that can be very large, and thus an option is provided to remove unwanted intermediate files. BAM files output by McClintock may be useful for other purposes, so an option is provided to eliminate all intermediate files other than BAM files. The location of all output files can be changed to any absolute path that the user requests. Within the specified location, all output files will be produced in a directory named after the reference genome sequence with results for each sample stored in subdirectories named after the fastq files for that sample, allowing multiple samples for the same reference genome to reuse common index files.

**Postprocessing and standardization of output format:** The component methods within McClintock produce their output in different file formats and annotation frameworks (see Bergman 2012 for discussion). Therefore, McClintock performs a number of postprocessing steps to standardize outputs from different methods into a common annotation framework. Details of the native annotation framework for component methods and the postprocessing steps made by McClintock can be found in the “Postprocessing and Standardization of Component

Method Output” section of File S1. Before performing these steps, the original (unedited) results for each method are saved in the output directory for that sample. If TE predictions are made by any component method in the additional “chromosomes” added in modified reference genomes (see above), these results are removed from the standard results files and retained in a subdirectory within the results directory called “non-ref\_chromosome\_results.”

The output file format chosen to standardize results for all component methods is a zero-based BED6 format because it allows easy integration with the BEDTools and UCSC genome browser. The BED format provides a fourth column to contain a name for the annotated feature. All records in these BED files contain the name of the TE family predicted at that location and whether the prediction is of a nonreference or reference TE. The name column also reports the sample ID from the fastq input file and the name of the component method that made the prediction. The type of evidence used for the prediction is also listed, either “sr” representing a prediction made from split-read evidence, “rp” representing a prediction made from read-pair evidence, or “nonab” for TEMP reference TE predictions that rely on no evidence for the absence of the TE in the sample. In addition, filtering and redundancy removal was performed within the result file for each component method. No redundancy filtering is performed by McClintock across component methods, allowing users to more directly compare output from different methods. To facilitate viewing of results on the UCSC genome browser, a header is included in each BED file. This header is read by the UCSC browser and lists the sample name and McClintock component system that produced the results as the track name and description, allowing multiple result files for the same sample to be merged and visualized simultaneously.

### Application of McClintock to simulated *S. cerevisiae* genomes with single synthetic TE insertions

To test McClintock and its component methods, we used simulated WGS data sets based on the genome of the model eukaryote, *S. cerevisiae*. We chose *S. cerevisiae* for testing McClintock because its reference genome is relatively small and has been completely determined (Goffeau et al. 1996), it has large samples of publicly available resequenced genomes (Liti et al. 2009; Almeida et al. 2015; Strobe et al. 2015), and the genome biology of its TEs is relatively simple and well characterized (Kim et al. 1998; Carr et al. 2012). Briefly, the 12-Mb *S. cerevisiae* reference genome contains 483 annotated TEs from six long terminal repeat (LTR) retrotransposon families (*Ty1*, *Ty2*, *Ty3*,

■ **Table 3** Data dependencies required to successfully run each component of the McClintock pipeline

	ngs_te_mapper	RelocaTE	TEMP	RetroSeq	PoPoolationTE	TE-locate
Reference genome (fasta)	✓	✓	✓	✓	✓	✓
Canonical TE sequences (fasta)	✓	✓ <sup>a</sup>	✓	✓ <sup>b</sup>	✓	✓
Annotation of reference TEs (GFF)					✓	✓
Annotation of reference TEs (BED)			✓	✓ <sup>c</sup>		
Annotation of reference TEs (custom format)		✓				
Unaligned reads (single-end fastq)	✓	✓				
Unaligned reads (paired-end fastq)					✓	
Aligned reads (BAM)			✓	✓		
Aligned reads (lexically sorted SAM)						✓
TE hierarchy (custom format)			✓		✓	

<sup>a</sup>Must include an entry in the format “TSD=...” for each TE in the file on the same line as the header, where “...” is the TSD sequence if known, or a string of periods with equal to the TSD length if the TSD sequence is unknown. If neither length nor the sequence of the TSD is known, “TSD=UNK” can be supplied.

<sup>b</sup>Must be formatted as one fasta file per TE family and a file of files listing their locations.

<sup>c</sup>Must be one BED file for each entry in the reference TE annotation and a file of files listing their locations.

*Ty3\_1p*, *Ty4*, and *Ty5*) (Carr *et al.* 2012), a type of TE that can be processed effectively by all six McClintock component methods. *Ty1* and *Ty2* share a nearly identical LTR sequence but differ in their internal regions (Kim *et al.* 1998), while *Ty3* and *Ty3\_1p* have 82% nucleotide identity over their entire length (Fingerman *et al.* 2003). Most TEs in *S. cerevisiae* are solo LTRs or otherwise truncated copies, with only ~50 full-length elements from four active families in the reference genome (*Ty1*, *Ty2*, *Ty3*, and *Ty4*) (Kim *et al.* 1998; Carr *et al.* 2012). *Ty1* and *Ty2* have the most full-length copies in the *S. cerevisiae* reference genome, with very few full-length copies being observed for *Ty3* and *Ty4* (Kim *et al.* 1998; Carr *et al.* 2012). The active TE families in *S. cerevisiae* are known to target tRNA genes (Ji *et al.* 1993; Chalker and Sandmeyer 1990, 1992; Devine and Boeke 1996; Kim *et al.* 1998; Baller *et al.* 2012; Mularoni *et al.* 2012; Qi *et al.* 2012) and create a 5-bp TSD on insertion (Gafner and Philippsen 1980; Rinckel and Garfinkel 1996; Chalker and Sandmeyer 1990; Kim *et al.* 1998; Zou *et al.* 1996).

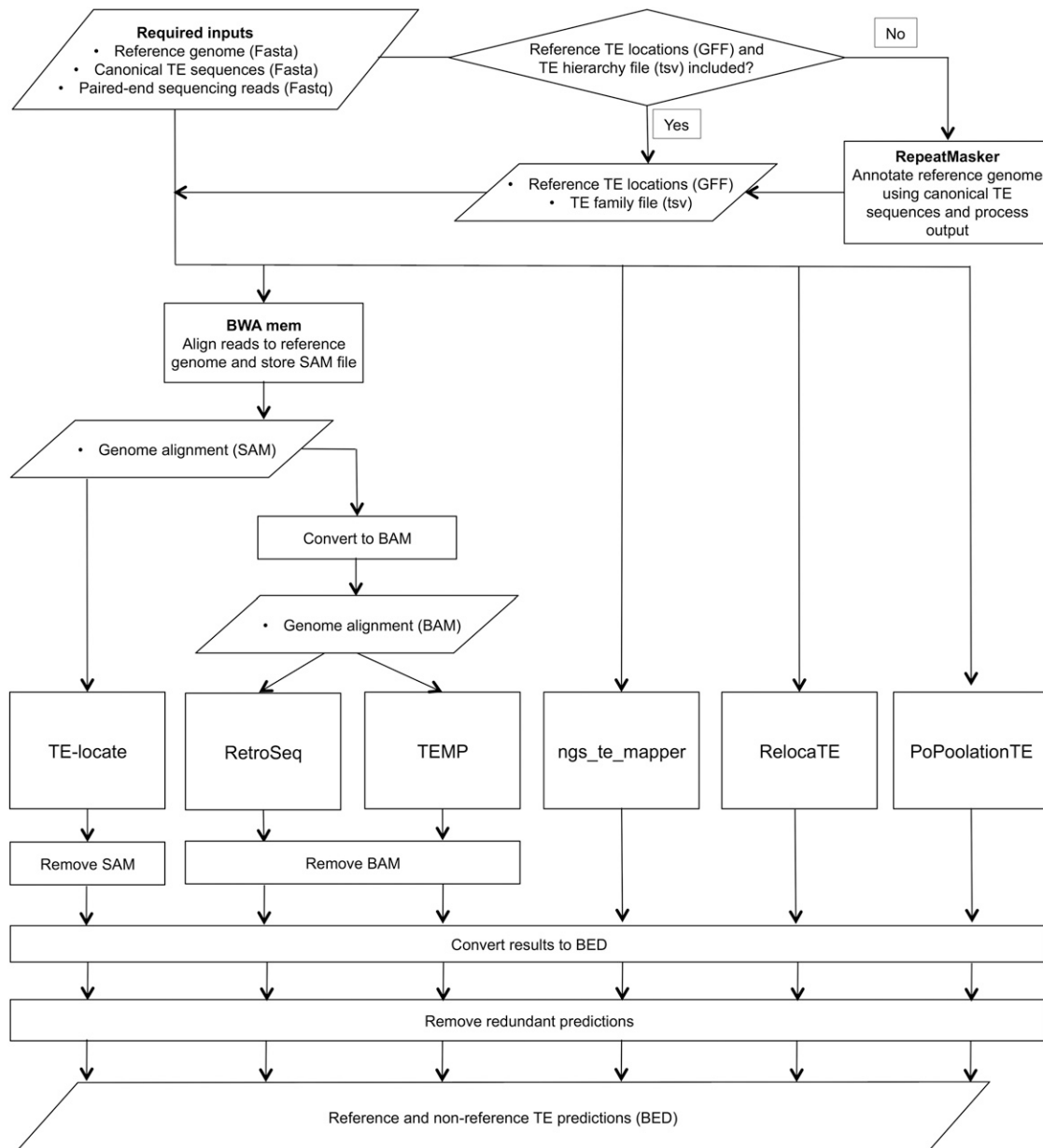
We first performed control analyses by simulating WGS resequencing of unmodified *S. cerevisiae* reference genome samples and applying McClintock to these data sets (see “Simulating Resequencing of the *S. cerevisiae* Reference Genome” in File S1). While not the major focus of this study, these reference genome simulations allowed us to evaluate how often McClintock component methods detected reference TEs and, more importantly, how often component methods detected false-positive nonreference TEs (in the absence of any true, nonreference TE insertions). An example of reference TE predictions for all six component methods is shown in Figure S1A in File S1. In general, analysis of unmodified simulated reference genomes showed that McClintock component methods cannot detect all reference TEs (Table S1 in File S1), but also typically have low false-positive rates for predicting nonreference TE insertions when they are truly absent (Table S2 in File S1). Additionally, these simulations showed that McClintock had better performance at 100× vs. 10× coverage, and that neither the choice of reference TE annotation nor reference genome options substantially affected the detection of reference or nonreference TEs for most McClintock component methods.

Next, we simulated WGS samples for reference genomes that include a single synthetic TE insertion from one of the four active TE families (placed at biologically realistic locations upstream of tRNA genes) to evaluate the ability of McClintock component methods to detect true positive nonreference TE insertions. To do this, WGS reads were simulated for 598 samples, each with a different synthetic TE insertion placed upstream of one of the 299 tRNA genes in the yeast genome. 299 samples were created for single synthetic insertions in the positive orientation upstream of tRNA genes, and 299 samples for single

synthetic insertions in the negative orientation. Genomes with synthetic insertions were created by alternately selecting one of the four active TE families and creating a 5-bp sequence 12–17 bp upstream of a tRNA start site for *Ty3* or 195–200 bp upstream of a tRNA start site for *Ty1*, *Ty2*, and *Ty4*. This 5-bp sequence formed the basis of a synthetic TSD and became the location into which a full-length *Ty* canonical sequence was inserted in the *sacCer2* reference genome. All single-insertion samples were simulated at 100× coverage since the ability of component methods to detect reference TEs improved with increasing coverage and to better match properties of the real yeast genomes analyzed below. An illustration of nonreference TE predictions for all six component methods in a genomic segment containing a synthetic TE insertion is shown in Figure S1B in File S1. In the following sections, we detail the analysis of these single synthetic insertion simulated samples in terms of overall numbers of reference and nonreference TE predictions and positional accuracy of nonreference TE predictions.

**Numbers of reference and nonreference TE predictions:** Table 4 shows the mean number of reference and nonreference TE insertions predicted across all 299 simulated single-insertion samples on the positive and negative strands, respectively. The proportion of correct predictions of nonreference TEs was calculated at four thresholds of accuracy: (i) requiring the exact TSD to be annotated correctly, (ii) requiring a prediction to be within a 100-bp window either side of the TSD, (iii) within a 300-bp window either side of the TSD (the insert size of the simulated sequencing libraries), or (iv) within a 500-bp window either side of the TSD. If all single TE insertion samples were predicted correctly for a method, it would lead to an average value of exactly one nonreference TE predicted per sample. Comparing row one of Table 4 (single-insertion simulation) with row nine of Table S1 in File S1 (unmodified reference simulation), we can infer that the inclusion of single synthetic insertions into the yeast genome does not substantially alter the ability of any McClintock component method to predict reference TEs. As expected, comparing row two of Table 4 (single-insertion simulation) with row nine of Table S2 in File S1 (unmodified reference simulation), we see gains in the numbers of nonreference TE insertions predicted for all methods; demonstrating that McClintock components can detect true positives above false-positive baselines in our simulation framework.

For *ngs\_te\_mapper*, the average number of nonreference predictions shows this method systematically underpredicts nonreference TE insertions. However, the average number of predictions made overall per sample is only slightly higher than the average number of exact



**Figure 1** Overview of the McClintock pipeline. In the flowchart, important processes are shown as boxes, decision points as diamonds, and data at important steps as parallelograms. Note that the last three steps of the pipeline are applied independently to each method. Final results from each component method are output independently by McClintock, allowing the user to easily merge output or assess for overlap among methods.

predictions. Consistent with unmodified reference genome simulations (see row nine of Table S2 in File S1), this result indicates that only a small number of nonreference predictions made by *ngs\_te\_mapper* are false positives. Moreover, whenever *ngs\_te\_mapper* makes a prediction of a nonreference TE (that is within 500 bp of the true insertion site), the prediction was always at the exact TSD, suggesting high accuracy in terms of position and TSD structure for this method (see below). We also observed that *ngs\_te\_mapper* detected fewer insertions when the synthetic insertion is on the negative strand relative to the tRNA gene, suggesting there can be strand bias in the detection of nonreference TEs. This bias could be due to yeast genome organization, our simulation framework, the *ngs\_te\_mapper* algorithm, or a combination of these factors.

RelocaTE produced, on average, slightly more than one nonreference TE prediction per sample. At face value, this result suggests that RelocaTE may detect essentially every synthetic insertion, but also makes occasional false-positive predictions. In fact, the average excess number of predictions made by RelocaTE in single-insertion simulated genomes is very close to the false-positive rates observed in simulations of unmodified reference genomes (see row nine of Table S2 in File S1). However, only ~50% of the total RelocaTE predictions are made within 500 bp of the true insertion. Thus, it appears that the inclusion of single synthetic insertions increases the rate of false-positive nonreference TE predictions by RelocaTE relative to unmodified reference genomes. Nevertheless, RelocaTE produces more correct predictions within 100 bp of the true insertion site than *ngs\_te\_mapper*, the other purely

■ **Table 4 Average numbers of predictions and correct predictions, by method, for simulated yeast WGS samples with a single synthetic TE insertion upstream of tRNA genes**

Insertion strand	ngs_te_mapper		RelocaTE		TEMP		RetroSeq		PoPoolationTE		TE-locate	
	+	-	+	-	+	-	+	-	+	-	+	-
Reference TEs mean	41.26	41.26	130.42	130.42	482.98	482.98	N.A.	N.A.	163.50	163.50	271.32	271.32
Nonreference TEs mean	0.42	0.32	1.12	1.11	0.90	0.90	0.87	0.86	1.18	1.14	0.98	0.92
Exact	0.40	0.29	0.30	0.24	0.36	0.36	0.00	0.00	0.00	0.00	0.00	0.00
Within 100 bp	0.40	0.29	0.63	0.61	0.90	0.90	0.68	0.66	0.16	0.16	0.07	0.06
Within 300 bp	0.40	0.29	0.63	0.61	0.90	0.90	0.69	0.67	0.16	0.16	0.70	0.54
Within 500 bp	0.40	0.29	0.63	0.61	0.90	0.90	0.69	0.67	0.16	0.16	0.82	0.78

Simulated WGS samples had 100x coverage, and McClintock was run using the reference TE annotation from Carr *et al.* (2012) and the unmodified reference genome option. The first two rows show the mean number of reference and nonreference predictions per sample, averaged across all simulated samples for that strand. Rows three to six show the average number of nonreference predictions of the correct TE family across samples that fell within the given distance of the known synthetic TE insertion site. For each method, the first column corresponds to insertions on the positive strand and the second column corresponds to insertions on the negative strand. For a prediction to be considered “exact,” the location of the TSD had to be predicted correctly. Numbers for TEMP combine predictions with split-read and read-pair support.

split-read method, despite producing fewer exact predictions than ngs\_te\_mapper. Thus many of the nonexact RelocaTE predictions within 100 bp of the true location are likely to be accurately positioned, but simply not have the correct TSD structure (see below). Like ngs\_te\_mapper, RelocaTE also appears to have a slightly higher true-positive rate for positive strand insertions, with the difference in the number of correct predictions on the positive strand being greater in the exact prediction category.

The average total number of nonreference TE predictions for TEMP is nearly one (0.90), confirming results from unmodified reference genome simulations (see row nine in Table S2 in File S1) that TEMP makes very few false-positive nonreference predictions. Moreover, the total number of nonreference TE predictions for TEMP is the same as the average number that are accurate within 100 bp of the true insertion site. These results suggest TEMP is correctly predicting most simulated insertions, but not to base pair accuracy (see below). Some positional inaccuracy is expected for TEMP since not all predictions for this method are supported by split-read evidence. For TEMP, there appears to be no difference in detection ability for TE insertions on the positive or negative strand.

RetroSeq predicted nearly as high an average number of nonreference TE predictions per sample as TEMP, but the proportion predicted correctly was lower than TEMP for all length thresholds. The fact that not all RetroSeq predictions are within 500 bp of the true insertion suggests that RetroSeq can produce some false-positive predictions of nonreference TE insertions when the sample is not identical to the reference genome, unlike what was observed for simulations of unmodified reference genomes (see row nine of Table S2 in File S1). Because RetroSeq does not use split-read information, no predictions from this method were exact, however most predictions were generally within 100 bp of the true location. For RetroSeq there is a only slight reduction in ability to detect nonreference TE insertions on the negative strand compared with the positive strand at all length thresholds.

PoPoolationTE produces an average of slightly more than one nonreference TE prediction per sample, but this method shows the lowest proportion of true-positive predictions at the most permissive length thresholds, suggesting most predictions are false positives. This result supports those obtained from unmodified reference genomes that PoPoolationTE makes approximately one false-positive prediction per genome in the absence of any synthetic nonreference TE insertions (see row nine of Table S2 in File S1). Because PoPoolationTE does not use split-read information and the span predicted by this method is often large (see Figure S1B in File S1), no predictions made by PoPoolationTE were exact. For PoPoolationTE there appears to be no difference in

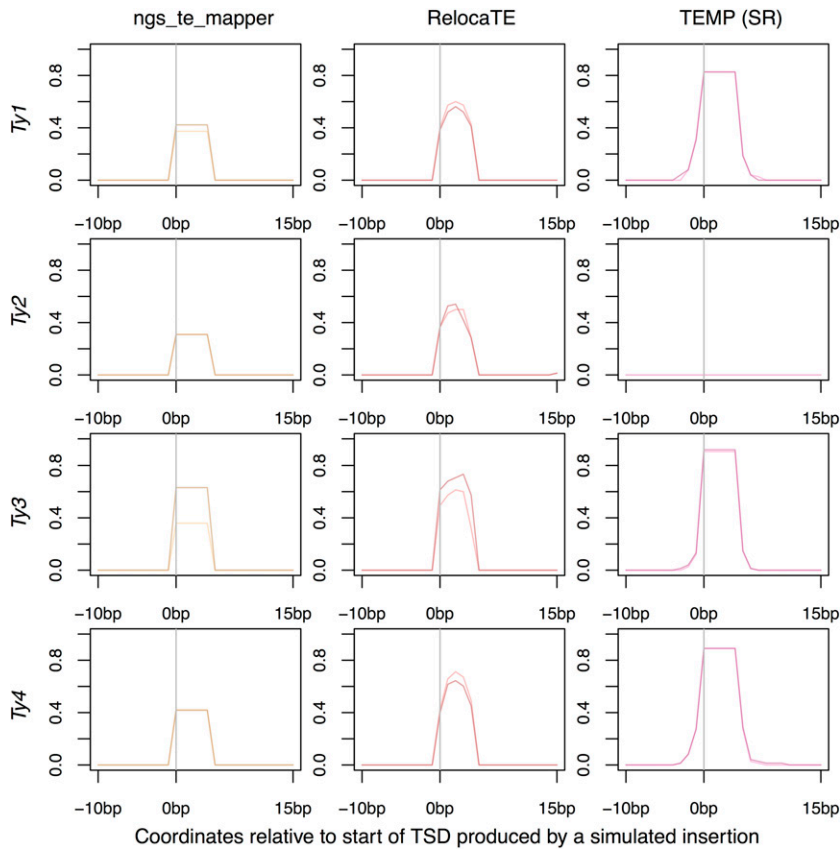
ability to detect nonreference TE insertions correctly on the positive or negative strand.

TE-locate produced an average of nearly one nonreference TE prediction per sample. However, these include some false-positive predictions or at least predictions that are >500 bp from the actual insertion location. The proportion of correct nonreference TE insertions predicted by TE-locate drops steadily from 500 to 100 bp, with TE-locate predicting the lowest number of correct insertions for any method at the 100-bp scale. As with the other read-pair methods, no predictions could be considered exact because TE-locate does not predict a TSD. These numbers indicate that, though the ability of TE-locate to detect the presence of a TE in the general vicinity of its true location is good, the annotation will not be as positionally accurate as other read-pair methods like TEMP or RetroSeq. For TE-locate there appears to be a reduction in detection ability at all thresholds for TE insertions on the negative strand compared with the positive strand.

**Positional accuracy of nonreference TE predictions:** To visualize more clearly the positional accuracy of McClintock component methods, predicted nonreference insertions were plotted around the known location of synthetic insertions (Figure 2 and Figure 3). Plots were produced for each TE family and method to determine if the family of the synthetic TE insertion affected results for a particular method. Table 4 showed that for split-read methods, there was no increase in the accuracy at thresholds of ~100 bp and many predictions were exactly correct. For read-pair methods, it appeared predictions could be several hundred base pairs from the correct location. As such, split-read (Figure 2) and read-pair (Figure 3) results were plotted on different spatial scales. Since TEMP could use both split-read and read-pair evidence, results for this method were partitioned into two categories for visualization. For a small number of cases, RelocaTE (one location) and PoPoolationTE (five locations) predicted nonreference TE insertions at the same genomic location in multiple samples. These predictions must include false positives based on the fact that each synthetic genome had only a single insertions at different genomic locations. Inclusion of these high-frequency, false-positive predictions dominated the visualization of results for these two methods, and thus predictions for these six cases were filtered prior generating Figure 2 and Figure 3 (see *Materials and Methods* for details).

Figure 2 shows that when ngs\_te\_mapper makes a prediction, it produces the TSD at the correct location, apparently with no TSDs called too long or too short. Direct analysis of TSD length distributions supports this conclusion: for simulated data, ngs\_te\_mapper always predicts the correct TSD length for nonreference insertions (Figure





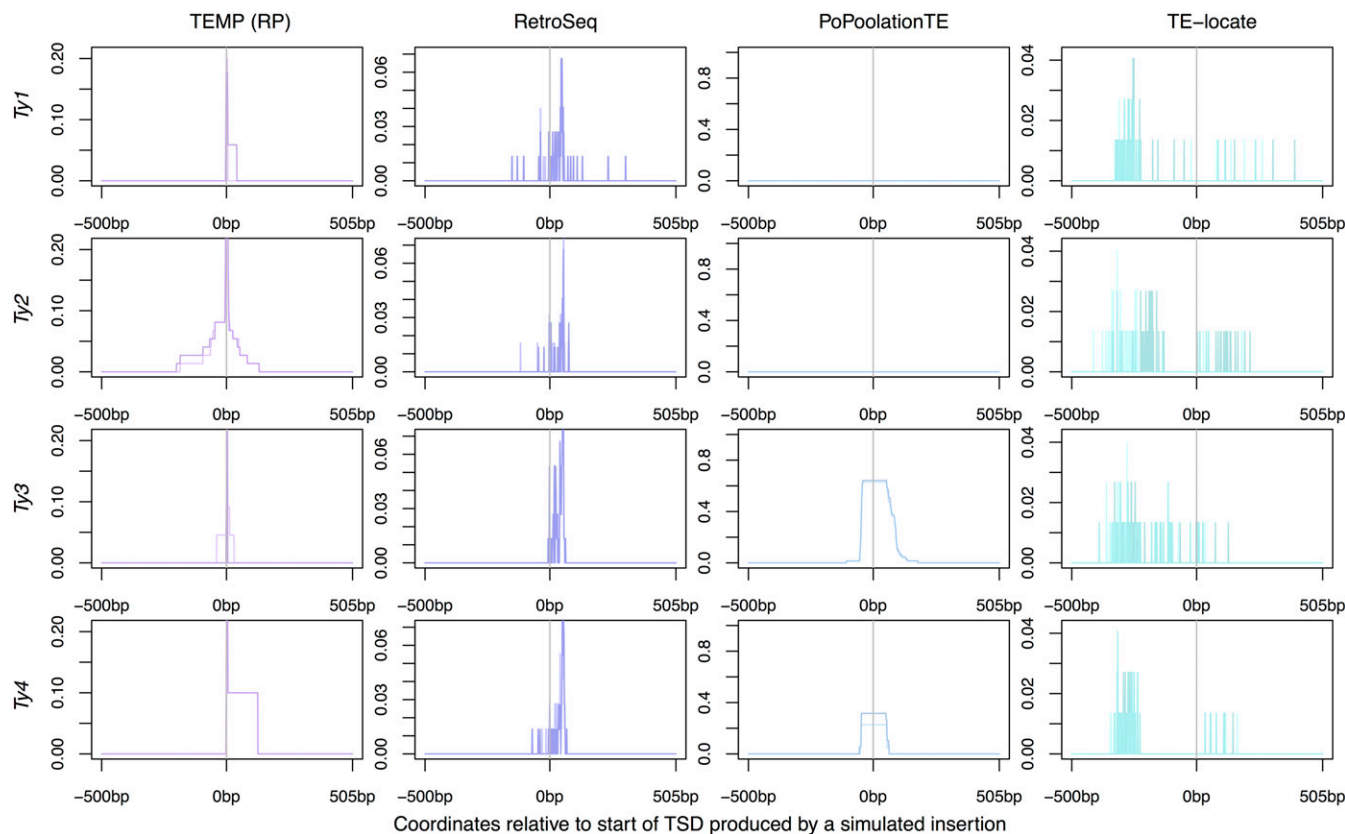
**Figure 2** Positional accuracy of nonreference TE insertions made by methods using split-read evidence on single-insertion synthetic genomes. Data for TEMP are for predictions that do have split-read evidence and may or may not have read-pair evidence. The location of the synthetic TSD is from position 0–5 bp on each plot. The darker line for each method indicates predictions averaged across simulated genomes with insertions on the positive strand; the lighter line indicates predictions averaged across simulated genomes with insertions on the negative strand. A value of one would indicate a perfect prediction in all samples since there is one synthetic insertion per genome.

S2 in File S1). For *Ty1*, *Ty2*, and *Ty4*, *ngs\_te\_mapper* detected insertions on the positive or negative strand with similar accuracy. Thus, the main difference in detection rates on the positive and negative strands for *ngs\_te\_mapper* observed in Table 4 appears to be for *Ty3* insertions, where many fewer insertions were detected correctly on the negative strand. For *RelocaTE*, the predicted TSDs of nonreference insertions are in approximately the correct locations but with coordinate ranges that are frequently too short (see also Figure S2 in File S1). As with *ngs\_te\_mapper*, *RelocaTE* shows the biggest difference in ability to detect *Ty3* insertions on the negative strand relative to the positive strand. TEMP split-read predictions for *Ty1*, *Ty3*, and *Ty4* are often predicted correctly but with the TSD often annotated to be longer than its true length (see also Figure S2 in File S1). Surprisingly, TEMP made no predictions for nonreference *Ty2* insertions using split-read evidence, perhaps because of the ambiguous signal arising from the similarity of *Ty1* and *Ty2* LTR sequences. For TEMP, there is no difference in detection ability for insertions on the positive or negative strand for any family.

Results of the positional accuracy for read-pair methods are shown in Figure 3. For *Ty1*, *Ty3*, and *Ty4* there were very few insertions (only three per family) that TEMP did not have split-read supporting evidence for, and thus few insertions for these families are plotted in Figure 3. In contrast, all *Ty2* predictions made by TEMP in the single-insertion simulations had read-pair evidence. For all families, when only read-pair evidence is used, TEMP generally predicts an insertion at the correct site, but with some slight inaccuracy on either side. The majority of RetroSeq predictions appear to be clustered close to the true insertion locations, but there appears to be a slight bias for RetroSeq to predict insertions 3' of where the true TE is located on reference genome coordinates. This bias is potentially introduced by the breakpoint determination step of RetroSeq, which always scans in the 5' to 3'

direction (see section “Description of McClintock Component Methods” in File S1). PoPoolationTE produced the highest number false-positive predictions (Table 4). When these false-positive nonreference predictions are filtered from the results, all predictions for *Ty1* and *Ty2* in the windows around simulated insertions are eliminated. The effect of removing false positives is probably most pronounced for *Ty1* because it is the most common TE family in *S. cerevisiae*, and thus would be the most likely family to have a reference insertion with sequence similarity to the synthetic insertion in the vicinity of tRNA genes. PoPoolationTE makes no predictions for *Ty2*, even including false positives. For *Ty3* and *Ty4*, PoPoolationTE has the capability of producing relatively accurate predictions, albeit with low resolution (nearly 100 bp around the true insertion site). For TE-locate, many predictions are made within 500 bp of the true insertion, but they are clearly spread further from the true insertion location than other methods. TE-locate also appears to have a slight bias to predict insertions 5' of the true insertion location on reference genome coordinates.

**Overlap between methods:** To understand the concordance of predictions made by the McClintock components, we investigated the overlap among methods for predictions that were made correctly at the sites of synthetic insertions. As shown in Figure 2 and Figure 3, different methods have different positional accuracy, and thus we used different windows to classify if a method made a “correct” prediction for a known insertion or not. Predictions for *ngs\_te\_mapper*, *RelocaTE*, TEMP (both split read and read pair), and PoPoolationTE were classified as correct if they had any overlap with the true location of the TSD; while predictions for RetroSeq and TE-locate were classified as correct if they occurred within a 100- or 500-bp window, respectively, of the correct location of the TSD. Neither the orientation nor the TE family was taken



**Figure 3** Positional accuracy of nonreference TE insertions made by methods using read-pair evidence on single-insertion synthetic genomes. Data for TEMP are for predictions that do not have split-read evidence but do have read-pair evidence. Note that the y-axes of plots are scaled differently for each method. The location of the synthetic TSD is from position 0–5 bp on each plot. The darker line for each method indicates predictions averaged across simulated genomes with insertions on the positive strand; the lighter line indicates predictions averaged across simulated genomes with insertions on the negative strand. A value of one would indicate a perfect prediction in all samples since there is one synthetic insertion per genome.

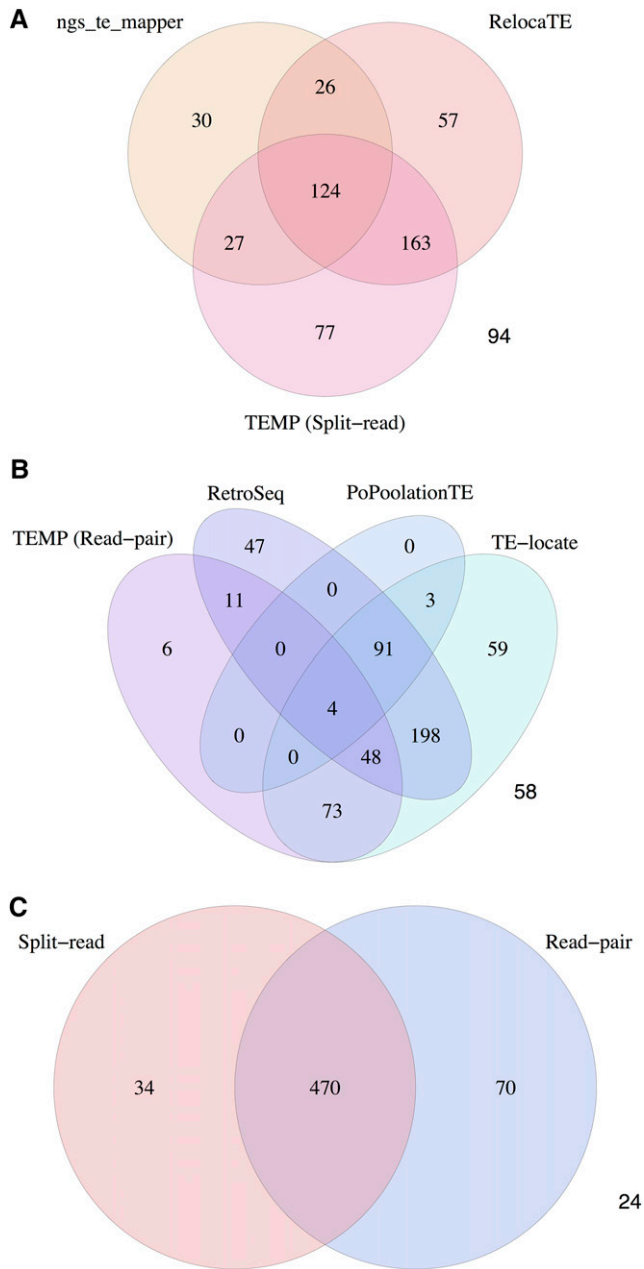
into account when classifying a prediction as correct or not. The overlap of correctly detected insertions are shown in Figure 4, A and B, for split-read and read-pair insertions, respectively. The overlap of correct predictions made by all split-read methods vs. all read-pair methods is shown in Figure 4C.

Figure 4A shows that the majority of split-read predictions are supported by at least two methods ( $n = 340$ , 57%) but that each method made many correct TE predictions that were not made by any other method. RelocaTE and TEMP made a greater number of correct overlapping predictions with each other than either of these methods did with `ngs_te_mapper`. Figure 4A also shows that 16% ( $n = 94$ ) of synthetic insertions were not predicted by any split-read method at the threshold of positional accuracy used here. Figure 4B shows that the vast majority of synthetic TE insertions ( $n = 428$ , 72%) are predicted by at least two of the read-pair methods, but that only 24% ( $n = 143$ ) of insertions are supported by three or more methods. RetroSeq and TE-locate make the highest number of unique correct predictions. ~10% ( $n = 58$ ) of synthetic insertion samples were not predicted by any read-pair method at the threshold of positional accuracy used here. Finally, Figure 4C shows that, while the overwhelming majority of insertions are predicted by at least one split-read and one read-pair method ( $n = 470$ , 79%), there are many insertions that are only predicted using one type of evidence or the other ( $n = 104$ , 17%) given the thresholds of positional accuracy used here. Nevertheless, use of all six methods recovers nearly 96% of synthetic insertions, demonstrating the

utility of integrating multiple TE-identification methods enabled by McClintock.

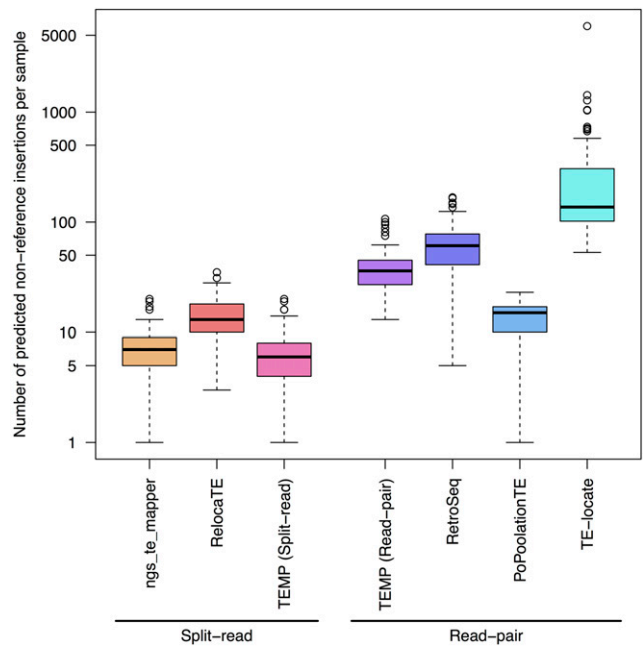
### Application of McClintock to 93 yeast genomes

The previous sections presented results on the accuracy of McClintock component methods on simulated resequencing data. Simulations are useful for testing methods under controlled settings, but do not capture all aspects of how methods perform when applied to real data. Since much is known about the expected insertion preferences of TEs in *S. cerevisiae* (Gafner and Philippsen 1980; Ji *et al.* 1993; Chalker and Sandmeyer 1990, 1992; Devine and Boeke 1996; Rinckel and Garfinkel 1996; Zou *et al.* 1996; Kim *et al.* 1998; Baller *et al.* 2012; Mularoni *et al.* 2012; Qi *et al.* 2012), analysis of real WGS data sets can be used as an alternative approach to evaluate if McClintock component methods can recapitulate the known genome biology of yeast TEs. To do this, we analyzed 93 high-coverage *S. cerevisiae* WGS data sets from Strope *et al.* (2015) using McClintock to generate TE predictions for all six component methods. Figure 3 and Figure S3 in File S1 show how many of the nonreference and reference TEs per strain, respectively, are detected by the different McClintock component methods across all 93 samples. In general, split-read methods predict between 5 and 20 nonreference TE insertions per strain, whereas read-pair methods predict ~40–100 nonreference TE insertions per strain (Figure 3). Numbers of reference TEs predicted per strain in real data (Figure 5) are generally lower than in simulated genomes (Table 4 and Table S1 in File S1). The exceptions to this pattern



**Figure 4** Concordance of correctly predicted nonreference insertions among McClintock component methods. (A) The concordance of nonreference predictions by methods that use split-read evidence that overlap with the true location of a synthetic insertion. (B) The concordance of nonreference predictions made by methods that use read-pair evidence that either overlap (TEMP, PoPoolationTE), or are within 100 bp (RetroSeq), or 500 bp (TE-locate) of the true location of a synthetic insertion. (C) The concordance of correctly predicted synthetic nonreference TEs with split-read or read-pair evidence. Predictions for TEMP were partitioned based on whether they had split-read evidence (split-read) or not (read-pair). Counts in all diagrams total 598, the number of simulated samples with single synthetic insertions.

are TEMP and PoPoolationTE, which show similar or higher numbers of reference TE predictions per strain in real data relative to simulations. We note that for a few strains in the Strope *et al.* (2015) data set, TE-locate predicted several hundred nonreference insertions; these



**Figure 5** Numbers of nonreference TE insertions per strain predicted by McClintock component methods in real yeast genomes. Predictions for TEMP were partitioned based on whether they had split-read evidence (split-read) or not (read-pair). Data are from 93 yeast strains taken from Strope *et al.* (2015). Methods are classified based on whether they use split-read or read-pair evidence to make a nonreference TE prediction. The box plot is shown on a  $\log_{10}$  scale. The thick line indicates the median, the colored box is the interquartile range, the whiskers mark the most extreme data point which is no more than 1.5 times the interquartile range from the box, and the  $\circ$ 's are outliers. Note that for TE-locate, several outlier samples generated hundreds of predicted nonreference TE insertions.

strains did not appear to be outliers in terms of their nonreference TE content based on other methods (results not shown).

We evaluated the quality of nonreference TE predictions made by McClintock component methods on the Strope *et al.* (2015) data set using three aspects of the known biology of TEs in *S. cerevisiae*: (i) activity of families, (ii) tRNA targeting, and (iii) TSD length. Our expectations based on prior knowledge of yeast TE biology are that methods that make high quality nonreference TE predictions should (i) show few nonreference predictions for inactive TE families (*Ty3\_1p* and *Ty5*), (ii) show a high proportion of nonreference predictions in the vicinity of tRNA genes, and (iii) show characteristic 5-bp TSDs for nonreference predictions made by split-read methods.

**Prediction of active and inactive families:** Table 5 shows numbers of nonreference TE predictions made by McClintock component methods across all strains in the Strope *et al.* (2015) data set. As expected, all methods predicted multiple nonreference insertions for TE families that are known to be active in this species. Additionally, ngs\_te\_mapper and TEMP make no nonreference TE predictions for both inactive families in *S. cerevisiae*, supporting simulation results above that show these methods have low false-positive rates. RelocaTE makes nonreference TE predictions for *Ty3\_1p* but not *Ty5*, PoPoolationTE makes nonreference TE predictions for *Ty5* but not *Ty3\_1p*, and both RetroSeq and TE-locate predict nonreference insertions for *Ty3\_1p* and *Ty5*. RelocaTE is the only split-read method that predicts nonreference insertions for an inactive family, suggesting that split-read methods

■ **Table 5** Number and location of nonreference TEs predicted by McClintock component methods in 93 yeast genomes

	Carr	ngs_te_mapper	RelocaTE	TEMP	RetroSeq	PoPoolationTE	TE-locate
Ty1	218/313 (70%)	93/101 (92%)	15/18 (83%)	827/1093 (76%)	1854/2835 (65%)	139/194 (72%)	2082/16388 (13%)
Ty2	30/46 (65%)	58/77 (75%)	303/425 (71%)	1343/1853 (72%)	839/1169 (72%)	27/36 (75%)	1110/8132 (14%)
Ty3	43/45 (96%)	378/387 (98%)	670/678 (99%)	991/1008 (98%)	1299/1445 (90%)	1006/1013 (99%)	1748/3813 (46%)
Ty3_1p	12/15 (80%)	0/0 (N.A.)	23/23 (100%)	0/0 (N.A.)	12/16 (75%)	0/0 (N.A.)	83/86 (97%)
Ty4	29/49 (59%)	95/118(81%)	143/190(75%)	259/310(84%)	238/292(82%)	15/20(75%)	324/1083(30%)
Ty5	0/15 (0%)	0/0 (N.A.)	0/0 (N.A.)	0/0 (N.A.)	3/74 (4%)	0/12 (0%)	0/887 (0%)

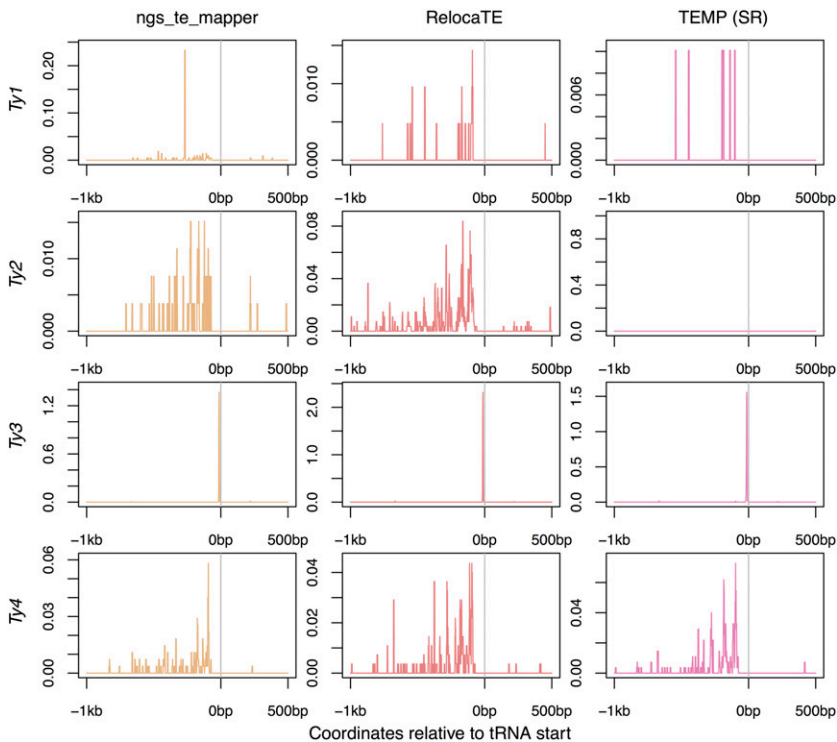
Each cell shows the number of nonreference TEs predicted in tRNA regions followed by the total number of nonreference TEs predicted genome wide. Data are for numbers of insertions, not numbers of nonredundant insertion sites, so TE insertion alleles present in more than one sample are counted independently. A prediction is counted in a tRNA region if any portion of the annotation is within 1000 bp upstream and 500 bp downstream of the tRNA start site, taking into account the orientation of the tRNA gene. The first column applies the same analysis to the reference TE annotations from Carr *et al.* (2012). N.A. indicates that no nonreference TE insertions were predicted by a method for that TE family.

generally have a higher ability to discriminate active from inactive TE families. Compared to the total numbers predicted for other active TE families, the three pure read-pair methods predicted fewer nonreference insertions for both *Ty3\_1p* and *Ty5*, suggesting false-positive rates for these methods are not so high as to overwhelm true signal. The one exception is for TE-locate, which predicted relatively high numbers of *Ty5* insertions, which is likely related to the outlier samples noted above where TE-locate predicts hundreds of presumably false-positive nonreference insertions.

**Predicted insertions in tRNA regions:** Active TE families in *S. cerevisiae* are known to target tRNA genes (Ji *et al.* 1993; Chalker and Sandmeyer 1990, 1992; Devine and Boeke 1996; Kim *et al.* 1998; Baller *et al.* 2012; Mularoni *et al.* 2012; Qi *et al.* 2012). The highest density of *Ty1* and *Ty2* insertions are in the 200 bp upstream of the tRNA transcription start site (Ji *et al.* 1993; Devine and Boeke 1996; Kim *et al.* 1998; Baller *et al.* 2012; Mularoni *et al.* 2012). *Ty3* targets a specific location just upstream of tRNA gene transcription start sites (Chalker and Sandmeyer 1990, 1992; Kim *et al.* 1998; Qi *et al.* 2012).

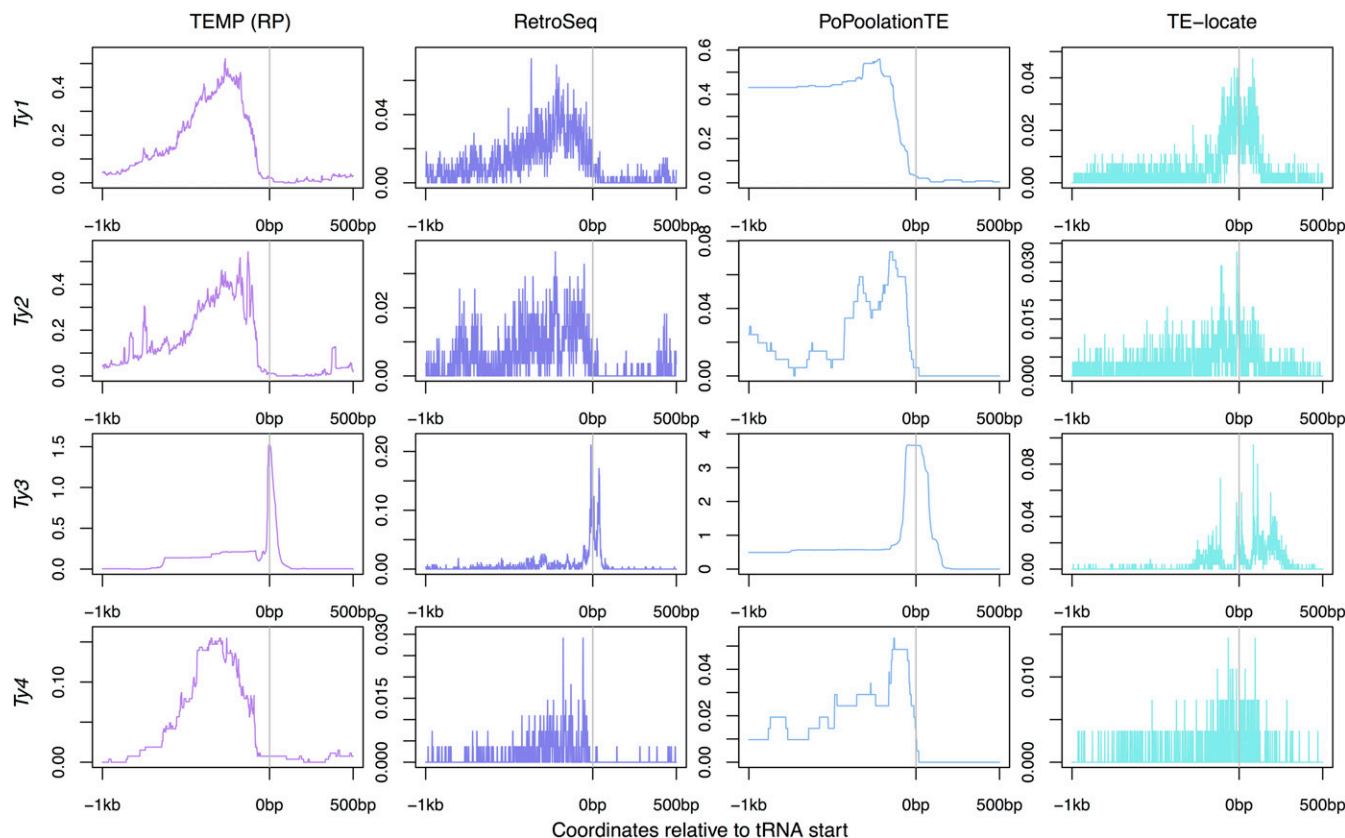
Patterns of *Ty4* insertions have not been experimentally determined, although the locations of insertions in the reference genome suggest a similar pattern to that of *Ty1* and *Ty2* (Kim *et al.* 1998).

To evaluate if nonreference TE insertions predicted by McClintock component methods show expected hallmarks of tRNA targeting, we plotted locations of nonreference TE insertions identified in the Strobe *et al.* (2015) strains using split-read evidence and read-pair evidence in Figure 6 and Figure 7, respectively. The expected profiles of insertion into tRNA gene regions is observed for all *Ty* families for ngs\_te\_mapper, RelocaTE, TEMP, and RetroSeq, albeit with the different levels of resolution that are characteristic of each method. Consistent with simulation data (Figure 2), TEMP appears to have difficulty predicting *Ty2* using split-read data in real yeast genomes, and this effect also appears to impact prediction of *Ty1* insertions using split-read data in real data (Figure 6). PoPoolationTE can predict meaningful profiles of insertion for *Ty3* and *Ty4* (Figure 7), as expected based on simulation data (Figure 3). However, in contrast to simulation data where only putative false positives are predicted (Figure 3), PoPoolationTE also predicts nonreference insertions for *Ty1* and *Ty2* in real data (Figure 7). Since



**Figure 6** Locations of nonreference TE predictions relative to tRNA genes made by methods using split-read evidence on real yeast genomes. Data for TEMP are for predictions that do have split-read evidence and may or may not have read-pair evidence. The transcription start site of each tRNA gene is aligned at position zero, taking into account tRNA gene orientation, for a window extending 1 kb upstream and 500 bp downstream. The frequency of a prediction at each base is counted across all 93 strains in the Strobe *et al.* (2015) data set, then averaged across the 299 tRNA genes, and plotted as a line for each method and TE family. These plots show all predictions and therefore include allelic predictions present in more than one strain. Also, any given strain may have more than one insertion at the same relative location in a tRNA gene, and thus the scale for these plots can go above one.





**Figure 7** Locations of nonreference TE predictions relative to tRNA genes made by methods using read-pair evidence on real yeast genomes. Data for TEMP are for predictions that do not have split-read evidence but do have read-pair evidence. The transcription start site of each tRNA gene is aligned at zero on the plots, taking into account orientation, for a window extending 1 kb upstream and 500 bp downstream. The frequency of a prediction at each base is counted across all 93 strains in the Strope *et al.* (2015) data set, then averaged across the 299 tRNA genes, and plotted as a line for each method and TE family. These plots show all predictions and therefore include allelic predictions present in more than one strain. Also, any given strain may have more than one insertion at the same relative location in a tRNA gene, and thus the scale for these plots can go above one.

PoPoolationTE predicts reference and nonreference insertions in the same way, and since many *Ty1* and *Ty2* insertions exist in the reference genome upstream regions of tRNA genes, it is possible that these *Ty1* and *Ty2* insertions predicted in the Strope *et al.* (2015) data set are actually reference insertions that are mislabeled by PoPoolationTE as nonreference insertions. Finally, nonreference insertions predicted by TE-locate are only weakly enriched in tRNA regions for all families, and the positional profiles produced by TE-locate are shifted relative to expectations and predictions made by other methods.

To quantify the proportion of nonreference TEs that were predicted in tRNA regions, we counted predictions 1000 bp upstream and 500 bp downstream of a tRNA gene, taking into account the orientation of the tRNA gene but not the orientation of the TE insertion. The expected percentage of TEs located in these regions if they were inserted randomly in the genome would be 0.037% [(299 tRNA genes × 1500 bp window)/12,162,995 bp genome]. Previous analyses of tRNA targeting of TEs in the *S. cerevisiae* reference genome (Kim *et al.* 1998) assessed whether TEs were within 750 bp of a tRNA gene or other RNA polymerase III gene (excluding other intervening TE sequences). Here we use extended regions for tRNA targeting based on the inaccuracy in nonreference predictions observed for some methods in the simulations above. For comparison with previous results, we first applied our definition of tRNA targeting to the reference TE annotation from Carr *et al.* (2012) (Table 5). Estimated proportions of *Ty* elements in tRNA

regions for the Carr *et al.* (2012) reference annotation are lower than those reported by Kim *et al.* (1998), however, they still show highly biased targeting toward tRNA regions.

Nonreference TE predictions of all four active *Ty* elements show the expected enrichment in tRNA regions for each McClintock component method (Table 5). For all methods, *Ty3* is the active TE family most strongly associated with tRNA regions, consistent with experimental data and observations based on the reference genome (Chalker and Sandmeyer 1990, 1992; Kim *et al.* 1998; Qi *et al.* 2012). Split-read methods predict a higher proportion of nonreference TEs in tRNA regions relative to expectations based on TEs in the reference genome. For read-pair methods, at least one TE family showed a lower proportion of nonreference TEs in tRNA regions relative to reference TEs. We interpret this observation to be due to the lower positional accuracy of read-pair methods. TE-locate consistently predicted the lowest number of TEs in tRNA regions for active *Ty* families, though predicted insertions for this method still showed an enrichment in tRNA regions relative to random expectation. We interpret the low tRNA enrichment for TE-locate to be a consequence of the low positional accuracy of read-pair methods combined with the presence of outlier samples for this method which have very high numbers of nonreference predictions.

As discussed above, nonreference predictions were made by RelocaTE, RetroSeq, and TE-locate for the inactive *Ty3\_1p* family. Despite most likely being false positives, these predictions were predominantly

in tRNA regions, suggesting they could either be nonreference *Ty3* insertions that are miscalled as nonreference *Ty3\_1p*, or reference *Ty3\_1p* insertions called as nonreference *Ty3\_1p* insertions. Nonreference predictions were also made by RetroSeq, PoPoolationTE, and TE-locate for the inactive *Ty5* family. The majority of these predictions are made outside of the tRNA regions, as is expected based on the known location of *Ty5* insertions in the reference genome prior knowledge about *Ty5* target preferences (Zou *et al.* 1996; Kim *et al.* 1998; Baller *et al.* 2011). These nonreference TE predictions may be false positives (possibly caused by mapping inconsistencies in heterochromatic regions where *Ty5* elements typically insert) or real nonreference “insertions” that arose by recombination events rather than transposition events (Zou *et al.* 1995).

**Prediction of TSDs by split-read methods:** Finally, we evaluated the performance of split-read methods to predict the known TSD lengths of active yeast *Ty* families in real WGS data. All available experimental and genomic data indicates that active yeast *Ty* families create 5 bp TSDs on insertion (Gafner and Philippsen 1980; Rinckel and Garfinkel 1996; Chalker and Sandmeyer 1990; Kim *et al.* 1998; Zou *et al.* 1996). TSD length distributions of unique insertion sites for McClintock predictions in real yeast genomes are shown for *Ty1*, *Ty2*, *Ty3*, and *Ty4* in Figure S4 in File S1. As observed in simulated data (Figure 2 and Figure S2 in File S1), *ngs\_te\_mapper* predictions had the highest proportion of correct TSD lengths predicted per family. However, in contrast to simulated data, *ngs\_te\_mapper* can infrequently make incorrect TSD-length predictions in real data. Confirming simulation results, RelocaTE generally underpredicts the length of TSDs, and TEMP consistently overpredicts the lengths of TSDs for all families in real data. For all split-read methods, the modal value of the TSD-length distribution reflects the true TSD length for all families. Thus, the modal TSD length provided by each of the split-read methods yields biologically meaningful inferences about TSD structure.

## CONCLUSIONS AND FUTURE DIRECTIONS

Here we describe McClintock, an integrated pipeline for detecting TE insertions in WGS resequencing data. McClintock offers many advantages relative to running multiple TE detection methods in isolation. Specific versions of compatible software dependencies required to run each method are fully documented, allowing users to easily set up their environment. The number of input files required to run all methods is reduced and complex processing of input files to create the correct custom formats and file relationships is automated. In addition, the pipeline is structured to allow parallel computations for multiple samples, so population data sets can be analyzed more quickly. Finally, results from individual methods are standardized to facilitate comparisons across methods and easy visualization in the UCSC genome browser. Overall, McClintock greatly lowers the barriers to running multiple TE detection methods, allowing users to gain more insight into how various methods work for their samples. McClintock does not currently include all published TE detection methods, although additional methods can be easily incorporated into the pipeline due to the flexible architecture and open-source nature of the project.

In addition, we have applied McClintock to simulated and real yeast WGS samples to evaluate the performance of McClintock component methods. Simulations on the unmodified *S. cerevisiae* reference genomes reveal that sequencing coverage influences detection of reference TEs, but that recovery of reference TE insertions and false-positive rates for nonreference TE insertions are generally low even at high sequencing coverage. Simulations on *S. cerevisiae* reference genomes

including a single nonreference insertion showed that pure split-read methods may detect fewer TE insertions than read-pair methods, but they have much higher positional accuracy. Single-insertion simulations also revealed that the TE family affects the ability of methods to detect nonreference TE insertions. We find substantial difference in the ability of McClintock component methods to detect subsets of nonreference insertions in the yeast genome, but that by combining multiple methods that use split-read and read-pair data, nonreference TEs at nearly all biologically realistic locations can be detected in simulated data. Finally, application of McClintock to a large sample of real yeast genomes reveals that most but not all McClintock component methods can recover known aspects of TE biology in yeast such as family activity status, tRNA gene targeting, and TSD structure. Together, our results suggest that even in the context of a simplified model eukaryotic genome like *S. cerevisiae*, current TE detection methods using short-read data do not provide comprehensive recovery of all TE insertions in WGS resequencing samples. Further performance studies in other genomic contexts, including newer methods not currently included in McClintock, are needed to generalize the results presented here, and to provide a road map for developing more advanced systems for the detection of TEs in unassembled short-read genomic data.

## ACKNOWLEDGMENTS

We thank Matthew Ronshaugen, Douda Bensasson, and members of the Bergman Laboratory for helpful comments throughout the project; Robert Kofler, Thomas Keane, Alexander Platzter, Sofia Robb, and Jiali Zhuang for code fixes and insights into McClintock component methods; and Nick Gresham for high-performance cluster computing assistance. This work was supported by Wellcome Trust Award 096602/B/11/Z (M.G.N.) and Human Frontier Science Program grant RGY0093/2012 (C.M.B.).

## LITERATURE CITED

- Almeida, P., R. Barbosa, P. Zalar, Y. Imanishi, K. Shimizu *et al.*, 2015 A population genomics insight into the Mediterranean origins of wine yeast domestication. *Mol. Ecol.* 24: 5412–5427.
- Baller, J. A., J. Gao, and D. F. Voytas, 2011 Access to DNA establishes a secondary target site bias for the yeast retrotransposon Ty5. *Proc. Natl. Acad. Sci. USA* 108: 20351–20356.
- Baller, J. A., J. Gao, R. Stamenova, M. J. Curcio, and D. F. Voytas, 2012 A nucleosomal surface defines an integration hotspot for the *Saccharomyces cerevisiae* Ty1 retrotransposon. *Genome Res.* 22: 704–713.
- Bardou, P., J. Mariette, F. Escudié, C. Djemiel, and C. Klopp, 2014 Jvarkit: an interactive Venn diagram viewer. *BMC Bioinformatics* 15: 293.
- Bergman, C. M., 2012 A proposal for the reference-based annotation of de novo transposable element insertions. *Mob. Genet. Elements* 2: 51–54.
- Bergman, C. M., and H. Quesneville, 2007 Discovering and detecting transposable elements in genome sequences. *Brief. Bioinform.* 8: 382–392.
- Carr, M., D. Bensasson, and C. M. Bergman, 2012 Evolutionary genomics of transposable elements in *Saccharomyces cerevisiae*. *PLoS One* 7: e50978.
- Chalker, D. L., and S. B. Sandmeyer, 1990 Transfer RNA genes are genomic targets for de Novo transposition of the yeast retrotransposon Ty3. *Genetics* 126: 837–850.
- Chalker, D. L., and S. B. Sandmeyer, 1992 Ty3 integrates within the region of RNA polymerase III transcription initiation. *Genes Dev.* 6: 117–128.
- Chen, J., T. Wrightsman, S. R. Wessler, and J. E. Stajich, 2017 RelocaTE2: A High Resolution Transposable Element Polymorphism Mapping Tool for Population Resequencing. *PeerJ*. 5: e2942.
- Chen, Y., H. Yao, E. J. Thompson, N. M. Tannir, J. N. Weinstein *et al.*, 2013 VirusSeq: software to identify viruses and their integration sites using next-generation sequencing of human cancer tissue. *Bioinformatics* 29: 266–267.

- Cridland, J. M., S. J. Macdonald, A. D. Long, and K. R. Thornton, 2013 Abundance and distribution of transposable elements in two *Drosophila* QTL mapping resources. *Mol. Biol. Evol.* 30: 2311–2327.
- Devine, S. E., and J. D. Boeke, 1996 Integration of the yeast retrotransposon Ty1 is targeted to regions upstream of genes transcribed by RNA polymerase III. *Genes Dev.* 10: 620–633.
- Ewing, A. D., 2015 Transposable element detection from whole genome sequence data. *Mob. DNA* 6: 24.
- Ewing, A. D., and H. H. Kazazian, 2010 High-throughput sequencing reveals extensive variation in human-specific L1 content in individual human genomes. *Genome Res.* 20: 1262–1270.
- Ewing, A. D., and H. H. Kazazian, Jr., 2011 Whole-genome resequencing allows detection of many rare LINE-1 insertion alleles in humans. *Genome Res.* 21: 985–990.
- Fingerman, E. G., P. G. Dombrowski, C. A. Francis, and P. D. Sniegowski, 2003 Distribution and sequence analysis of a novel Ty3-like element in natural *Saccharomyces paradoxus* isolates. *Yeast* 20: 761–770.
- Fiston-Lavier, A.-S., M. Carrigan, D. A. Petrov, and J. González, 2011 T-lex: a program for fast and accurate assessment of transposable element presence using next-generation sequencing data. *Nucleic Acids Res.* 39: e36.
- Fiston-Lavier, A.-S., M. G. Barrón, D. A. Petrov, and J. González, 2015 T-lex2: genotyping, frequency estimation and re-annotation of transposable elements using single or pooled next-generation sequencing data. *Nucleic Acids Res.* 43: e22.
- Gafner, J., and P. Philippsen, 1980 The yeast transposon Ty1 generates duplications of target DNA on insertion. *Nature* 286: 414–418.
- Gilly, A., M. Etcheverry, M.-A. Madoui, J. Guy, L. Quadraña *et al.*, 2014 TE-Tracker: systematic identification of transposition events through whole-genome resequencing. *BMC Bioinformatics* 15: 377.
- Goffeau, A., B. G. Barrell, H. Bussey, R. W. Davis, B. Dujon *et al.*, 1996 Life with 6000 genes. *Science* 274: 563–567.
- Hawkey, J., M. Hamidian, R. R. Wick, D. J. Edwards, H. Billman-Jacobe *et al.*, 2015 ISMapper: identifying transposase insertion sites in bacterial genomes from short read sequence data. *BMC Genomics* 16: 667.
- Hénaff, E., L. Zapata, J. M. Casacuberta, and S. Ossowski, 2015 Jitterbug: somatic and germline transposon insertion detection at single-nucleotide resolution. *BMC Genomics* 16: 768.
- Hoen, D. R., G. Hickey, G. Bourque, J. Casacuberta, R. Cordaux *et al.*, 2015 A call for benchmarking transposable element annotation methods. *Mob. DNA* 6: 13.
- Hormozdiari, F., I. Hajirasouliha, P. Dao, F. Hach, D. Yorukoglu *et al.*, 2010 Next-generation VariationHunter: combinatorial algorithms for transposon insertion discovery. *Bioinformatics* 26: i350–i357.
- Ji, H., D. P. Moore, M. A. Blomberg, L. T. Braiterman, D. F. Voytas *et al.*, 1993 Hotspots for unselected Ty1 transposition events on yeast chromosome III are near tRNA genes and LTR sequences. *Cell* 73: 1007–1018.
- Jiang, C., C. Chen, Z. Huang, R. Liu, and J. Verdier, 2015 ITIS, a bioinformatics tool for accurate identification of transposon insertion sites using next-generation sequencing data. *BMC Bioinformatics* 16: 72.
- Kaminker, J. S., C. M. Bergman, B. Kronmiller, J. Carlson, R. Svirskas *et al.*, 2002 The transposable elements of the *Drosophila melanogaster* euchromatin: a genomics perspective. *Genome Biol.* 3: research0084.
- Keane, T. M., K. Wong, and D. J. Adams, 2013 RetroSeq: transposable element discovery from next-generation sequencing data. *Bioinformatics* 29: 389–390.
- Kent, W. J., 2002 BLAT—the BLAST-like alignment tool. *Genome Res.* 12: 656–664.
- Kent, W. J., A. S. Zweig, G. Barber, A. S. Hinrichs, and D. Karolchik, 2010 BigWig and BigBed: enabling browsing of large distributed datasets. *Bioinformatics* 26: 2204–2207.
- Kim, J. M., S. Vanguri, J. D. Boeke, A. Gabriel, and D. F. Voytas, 1998 Transposable elements and genome organization: a comprehensive survey of retrotransposons revealed by the complete *Saccharomyces cerevisiae* genome sequence. *Genome Res.* 8: 464–478.
- Kofler, R., A. J. Betancourt, and C. Schlötterer, 2012 Sequencing of pooled DNA samples (pool-seq) uncovers complex dynamics of transposable element insertions in *Drosophila melanogaster*. *PLoS Genet.* 8: e1002487.
- Kofler, R., D. Gómez-Sánchez, and C. Schlötterer, 2016 PoPoolationTE2: comparative population genomics of transposable elements using pool-seq. *Mol. Biol. Evol.* 33: 2759–2764.
- Kuhn, R. M., D. Haussler, and W. J. Kent, 2013 The UCSC genome browser and associated tools. *Brief. Bioinform.* 14: 144–161.
- Langmead, B., C. Trapnell, M. Pop, and S. L. Salzberg, 2009 Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol.* 10: R25.
- Lee, E., R. Iskow, L. Yang, O. Gokcumen, P. Haseley *et al.*, 2012 Landscape of somatic retrotransposition in human cancers. *Science* 337: 967–971.
- Lerat, E., 2010 Identifying repeats and transposable elements in sequenced genomes: how to find your way through the dense forest of programs. *Heredity (Edinb)* 104: 520–533.
- Li, H., 2013 Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. arXiv: 1303.3997.
- Li, H., B. Handsaker, A. Wysoker, T. Fennell, J. Ruan *et al.*, 2009 The sequence alignment/map format and SAMtools. *Bioinformatics* 25: 2078–2079.
- Linhaeiro, R. S., and C. M. Bergman, 2012 Whole genome resequencing reveals natural target site preferences of transposable elements in *Drosophila melanogaster*. *PLoS One* 7: e30008.
- Liti, G., D. M. Carter, A. M. Moses, J. Warringer, L. Parts *et al.*, 2009 Population genomics of domestic and wild yeasts. *Nature* 458: 337–341.
- Mularoni, L., Y. Zhou, T. Bowen, S. Gangadharan, S. J. Wheelan *et al.*, 2012 Retrotransposon Ty1 integration targets specifically positioned asymmetric nucleosomal DNA segments in tRNA hotspots. *Genome Res.* 22: 693–703.
- Nakagome, M., E. Solovieva, A. Takahashi, H. Yasue, H. Hirochika *et al.*, 2014 Transposon Insertion Finder (TIF): a novel program for detection of de novo transpositions of transposable elements. *BMC Bioinformatics* 15: 71.
- Nellaker, C., T. M. Keane, B. Yalcin, K. Wong, A. Agam *et al.*, 2012 The genomic landscape shaped by selection on transposable elements across 18 mouse strains. *Genome Biol.* 13: R45.
- Norel, R., J. J. Rice, and G. Stolovitzky, 2011 The self-assessment trap: can we all be better than average? *Mol. Syst. Biol.* 7: 537.
- Platzer, A., V. Nizhynska, and Q. Long, 2012 TE-Locate: a tool to locate and group transposable element occurrences using paired-end next-generation sequencing data. *Biology (Basel)* 1: 395–410.
- Qi, X., K. Daily, K. Nguyen, H. Wang, D. Mayhew *et al.*, 2012 Retrotransposon profiling of RNA polymerase III initiation sites. *Genome Res.* 22: 681–692.
- Quadraña, L., A. B. Silveira, G. F. Mayhew, C. LeBlanc, R. A. Martienssen *et al.*, 2016 The Arabidopsis thaliana mobilome and its impact at the species level. *Elife* 5: e15716.
- Quinlan, A. R., and I. M. Hall, 2010 BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* 26: 841–842.
- Quinlan, A. R., R. A. Clark, S. Sokolova, M. L. Leibowitz, Y. Zhang *et al.*, 2010 Genome-wide mapping and assembly of structural variant breakpoints in the mouse genome. *Genome Res.* 20: 623–635.
- Rahman, R., G. W. Chirn, A. Kanodia, Y. A. Sytnikova, B. Brembs *et al.*, 2015 Unique transposon landscapes are pervasive across *Drosophila melanogaster* genomes. *Nucleic Acids Res.* 43: 10655–10672.
- R Core Team, 2013 *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Rinckel, L. A., and D. J. Garfinkel, 1996 Influences of histone stoichiometry on the target site preference of retrotransposons Ty1 and Ty2 in *Saccharomyces cerevisiae*. *Genetics* 142: 761–776.
- Rishishwar, L., L. Mariño-Ramírez, and I. K. Jordan, 2016 Benchmarking computational tools for polymorphic transposable element detection. *Brief. Bioinform.* bbw072.
- Robb, S. M. C., L. Lu, E. Valencia, J. M. Burnette, Y. Okumoto *et al.*, 2013 The use of RelocaTE and unassembled short reads to produce high-resolution snapshots of transposable element generated diversity in rice. *G3 (Bethesda)* 3: 949–957.
- Sackton, T. B., R. J. Kulathinal, C. M. Bergman, A. R. Quinlan, E. B. Dopman *et al.*, 2009 Population genomic inferences from sparse high-throughput sequencing of two populations of *Drosophila melanogaster*. *Genome Biol. Evol.* 1: 449–465.

- Saha, S., S. Bridges, Z. V. Magbanua, and D. G. Peterson, 2008 Computational approaches and tools used in identification of dispersed repetitive DNA sequences. *Trop. Plant Biol.* 1: 85–96.
- Slater, G. S. C., and E. Birney, 2005 Automated generation of heuristics for biological sequence comparison. *BMC Bioinformatics* 6: 31.
- Smit, A. F. A., R. Hubley, and P. Green, 2013 RepeatMasker. Available at: <http://www.repeatmasker.org>.
- Stajich, J. E., D. Block, K. Boulez, S. E. Brenner, S. A. Chervitz *et al.*, 2002 The Bioperl toolkit: perl modules for the life sciences. *Genome Res.* 12: 1611–1618.
- Stempor, P., 2014 seqplots: an interactive tool for visualizing NGS signals and sequence motif densities along genomic features using average plots and heatmaps. R package version 1.6.0. Available at: <http://github.com/przemol/seqplots>.
- Strope, P. K., D. A. Skelly, S. G. Kozmin, G. Mahadevan, E. A. Stone *et al.*, 2015 The 100-genomes strains, an *S. cerevisiae* resource that illuminates its natural phenotypic and genotypic variation and emergence as an opportunistic pathogen. *Genome Res.* 25: 762–774.
- Stuart, T., S. R. Eichten, J. Cahn, Y. V. Karpievitch, J. O. Borevitz *et al.*, 2016 Population scale mapping of transposable element diversity reveals links to gene regulation and epigenomic variation. *Elife* 5: e20777.
- Sudmant, P. H., T. Rausch, E. J. Gardner, R. E. Handsaker, A. Abyzov *et al.*, 2015 An integrated map of structural variation in 2,504 human genomes. *Nature* 526: 75–81.
- Thung, D. T., J. de Ligt, L. E. Vissers, M. Steehouwer, M. Kroon *et al.*, 2014 Mobster: accurate detection of mobile element insertions in next generation sequencing data. *Genome Biol.* 15: 488.
- Wu, J., W.-P. Lee, A. Ward, J. A. Walker, M. K. Konkel *et al.*, 2014 Tangram: a comprehensive toolbox for mobile element insertion detection. *BMC Genomics* 15: 795.
- Zhuang, J., J. Wang, W. Theurkauf, and Z. Weng, 2014 TEMP: a computational method for analyzing transposable element polymorphism in populations. *Nucleic Acids Res.* 42: 6826–6838.
- Zou, S., D. A. Wright, and D. F. Voytas, 1995 The *Saccharomyces* Ty5 retrotransposon family is associated with origins of DNA replication at the telomeres and the silent mating locus HMR. *Proc. Natl. Acad. Sci. USA* 92: 920–924.
- Zou, S., N. Ke, J. M. Kim, and D. F. Voytas, 1996 The *Saccharomyces* retrotransposon Ty5 integrates preferentially into regions of silent chromatin at the telomeres and mating loci. *Genes Dev.* 10: 634–645.

Communicating editor: J. Ross-Ibarra