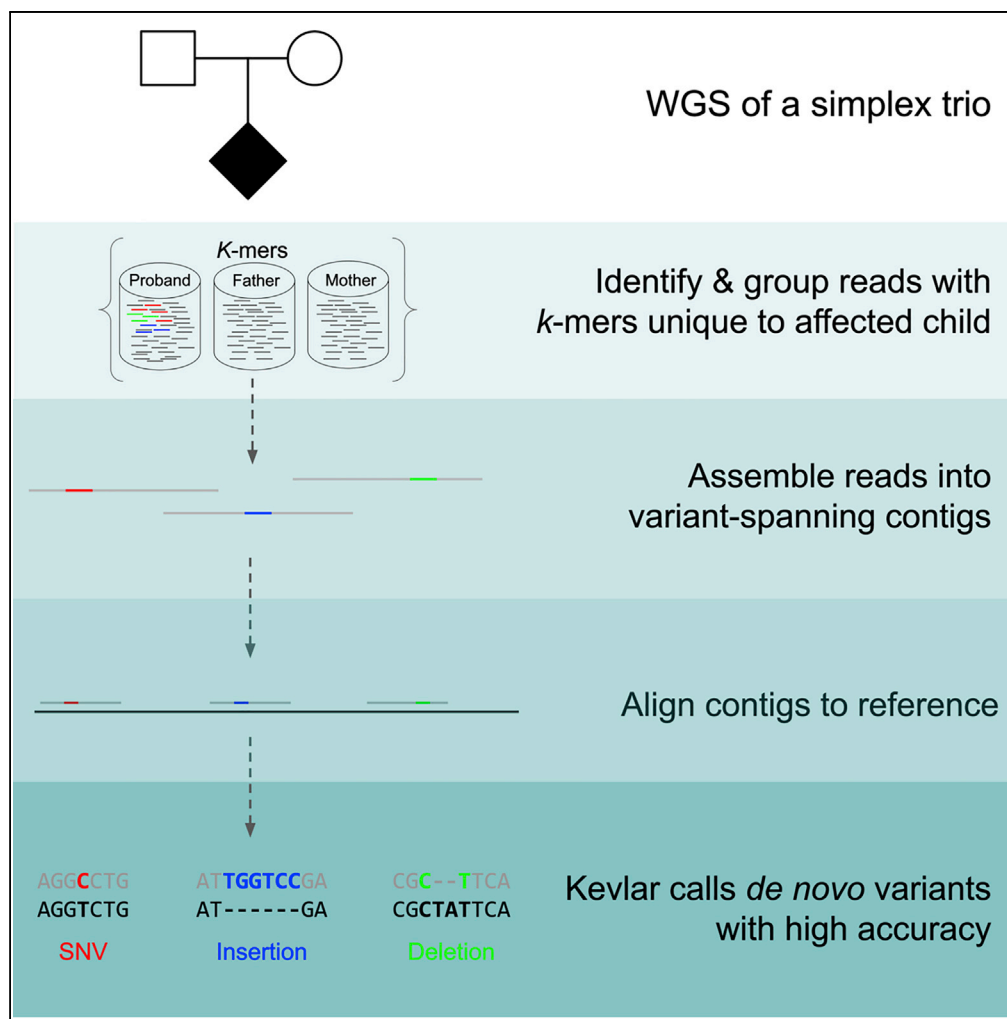


Special Issue: RECOMB-Seq 2019

Article

# Kevlar: A Mapping-Free Framework for Accurate Discovery of *De Novo* Variants



Daniel S. Standage, C. Titus Brown, Fereydoon Hormozdiari

daniel.standage@nbacc.dhs.gov (D.S.S.)  
 ctbrown@ucdavis.edu (C.T.B.)  
 fhormozd@ucdavis.edu (F.H.)

**HIGHLIGHTS**

Method for discovery of *de novo* variants without mapping reads to a reference genome

Novel probabilistic score for ranking variant predictions as confidently *de novo*

Predicts *de novo* SNVs, indels, and structural variants with high accuracy

Higher accuracy than competing methods for predicting long (>100 bp) variants

Standage et al., iScience 18, 28–36  
 August 30, 2019 © 2019 The Authors.  
<https://doi.org/10.1016/j.isci.2019.07.032>



Special Issue: RECOMB-Seq 2019

## Article

# Kevlar: A Mapping-Free Framework for Accurate Discovery of *De Novo* Variants

Daniel S. Standage,<sup>1,5,\*</sup> C. Titus Brown,<sup>1,2,\*</sup> and Fereydzoun Hormozdiari<sup>2,3,4,6,\*</sup>

## SUMMARY

***De novo* genetic variants are an important source of causative variation in complex genetic disorders. Many methods for variant discovery rely on mapping reads to a reference genome, detecting numerous inherited variants irrelevant to the phenotype of interest. To distinguish between inherited and *de novo* variation, sequencing of families (parents and siblings) is commonly pursued. However, standard mapping-based approaches tend to have a high false-discovery rate for *de novo* variant prediction. Kevlar is a mapping-free method for *de novo* variant discovery, based on direct comparison of sequences between related individuals. Kevlar identifies high-abundance *k*-mers unique to the individual of interest. Reads containing these *k*-mers are partitioned into disjoint sets by shared *k*-mer content for variant calling, and preliminary variant predictions are sorted using a probabilistic score. We evaluated Kevlar on simulated and real datasets, demonstrating its ability to detect both *de novo* single-nucleotide variants and indels with high accuracy.**

## INTRODUCTION

It is speculated that genetic variation is a major contributing factor in complex genetic disorders. The genetic heritability of many disorders is estimated to be relatively high. For example, the heritability of autism spectrum disorder is over 0.6, and the heritability of schizophrenia is over 0.5 (Cardno et al., 1999; Hallmayer et al., 2011). Only a fraction of this heritability is explained by known genetic variants, however, a phenomenon termed *missing heritability* (Manolio et al., 2009). One hypothesis is that *de novo* mutations, in particular indels and structural variants (SVs), are a large source of causative variation (and consequently missing heritability) in developmental disorders (Eichler et al., 2010; Manolio et al., 2009; Veltman and Brunner, 2012). However, the complexity of *de novo* variant discovery, especially *de novo* indel and SV discovery, has resulted in incomplete accounting of their contribution to these disorders. The discovery of genetic variants in general, and *de novo* variants in particular, remains a topic of intense research interest. In addition to illuminating the role of genetic variation in the etiology of complex disorders, improved discovery and cataloging of *de novo* variants across many samples or cohorts will shed additional light on important unresolved questions in human genomics, including rates, biases, and mechanisms of new mutation.

Whole genome sequencing of simplex families (presenting an isolated case of a genetic disorder) is a proven successful approach for discovery of novel genetic variants resulting from *de novo* mutation in the germline (Fromer et al., 2014; Iossifov et al., 2014; O’Roak et al., 2012; Veltman and Brunner, 2012; Zaidi et al., 2013). A “trio” composed of an individual affected by the disorder (the proband), the mother, and the father (alternatively, a “quad” or “quartet” composed of the proband, mother, father, and a sibling) provides a rich information source for discriminating between shared and unique variation. Following standard variant calling protocols, mapping-based methods for *de novo* variant prediction begin by aligning reads to the reference genome. Variants are then predicted for each individual based on artifacts observed in the read alignments, such as mismatches, gaps, abrupt shifts in coverage, and discordant read pair distances or orientations (Hormozdiari et al., 2009; Layer et al., 2014; Medvedev et al., 2010; Rausch et al., 2012; Sindi et al., 2012; Soylev et al., 2017; Ye et al., 2009). This initial process typically results in millions of variant predictions, which *de novo* variant discovery algorithms must examine to discern between inherited variation, true *de novo* variation, and spurious variant calls.

<sup>1</sup>Population Health and Reproduction, University of California, Davis, USA

<sup>2</sup>Genome Center, University of California, Davis, USA

<sup>3</sup>MIND Institute, University of California, Davis, USA

<sup>4</sup>Biochemistry and Molecular Medicine, University of California, Davis, 1 Shields Avenue, Davis, CA 95616, USA

<sup>5</sup>Present address: National Biodefense Analysis and Countermeasures Center, Fort Detrick, MD 21702, USA

<sup>6</sup>Lead Contact

\*Correspondence: [daniel.standage@nbacc.dhs.gov](mailto:daniel.standage@nbacc.dhs.gov) (D.S.S.), [ctbrown@ucdavis.edu](mailto:ctbrown@ucdavis.edu) (C.T.B.), [fhormozd@ucdavis.edu](mailto:fhormozd@ucdavis.edu) (F.H.)  
<https://doi.org/10.1016/j.isci.2019.07.032>



Although reference-based variant discovery methods have proved valuable in the study of complex genetic disorders, we note some of their limitations. Despite consistent improvements in read alignment algorithms, finding the correct mapping for each read is still complicated by sequencing errors, repetitive DNA content, and misassemblies in the reference. Reads that do not map to the reference genome because they span mutation breakpoints or contain novel sequence are ignored completely by mapping-based variant predictors. Also, few methods are able to predict multiple variant types simultaneously using a single strategy, instead focusing exclusively on single-nucleotide variants (SNVs), short indels, or SVs separately. Finally, most variant calls determined by analysis of read alignments are not unique to the individual of interest (child, or *proband*) but instead reflect divergence in ancestry between the family and the reference genome donors. Estimates of human germline mutation rates give an expectation of approximately 80 novel mutations per generation (Campbell and Eichler, 2013), and distinguishing true *de novo* variation events from millions of inherited or false variants is a substantial challenge.

More generally, accurate and comprehensive *de novo* variant discovery is complicated by several computational and biological factors, and remains an elusive goal. Any algorithm must be confident not only in the existence of the variant in the proband but also in its *non-existence* in both parents. And although SNVs are the most common variant type, larger variants that are less frequent, nevertheless, affect more nucleotides overall and are hypothesized to have an even greater impact in genetic disorders. Accurate discovery of these larger *de novo* variants is particularly challenging due to the inherent complexity of indel and SV prediction. In a reference-mapping context, calling indels with confidence requires accurate mapping of each read spanning the indel, with all gaps arranged consistently. This is possible only for short indels and tends to be prone to error and misalignment. Thus prediction of indels with length >10 bp has proved to be very challenging and accompanied by high false-positive and false-negative rates. Furthermore, the prediction of SVs via read mapping is only possible through indirect signatures such as alterations in read depth or read-pair signatures. These signatures can be quite noisy and result in high rate of false-negative and false-positive prediction. As a result, some basic properties of *de novo* SVs, including their rate of occurrence, remain unknown. It is important to note that there also exists no method for predicting more complex types of *de novo* SVs, such as inversion-duplication.

Many of the challenges with *de novo* variant prediction can be mitigated by an approach that compares sequence content between related individuals directly, rather than indirectly via a reference genome. Such an approach neither requires any read alignments nor is it sensitive to off-target shared or inherited variation. What a mapping-free approach *does* require is a signature of variation that is not defined in terms of artifacts observed in read alignments.

One of the first tools to explore a mapping-free strategy for predicting and genotyping variants was Cortex, which introduced the concept of a “colored de Bruijn graph” to compare sequence content from two or more samples and predict variants between samples (Iqbal et al., 2012). Cortex was used successfully for predicting variants in the 1000 Genomes Project. The DiscoSnp method (Uricaru et al., 2014) implemented a very efficient strategy for scanning a de Bruijn graph for “bubbles” reflective of isolated SNVs. More recently, DiscoSnp++ has improved on this strategy and is capable of predicting isolated SNVs, proximal SNVs, and indels without the use of a reference genome (Peterlongo et al., 2017). At the core of both methods is the analysis of *k*-mers, or sequences of a fixed length *k*.

Increased attention is being given to these kinds of *k*-mer-based methods that avoid read alignments altogether. Indeed, mapping-free strategies for a variety of genomic and transcriptomic applications have become increasingly prominent, in large part due to their efficiency and robustness to the shortcomings of reference genomes. (It is important to note that these and other developments have greatly benefited from the availability of software libraries for rapid exact and approximate *k*-mers; these libraries include Jellyfish, Marçais and Kingsford, 2011; khmer, Crusoe et al., 2015; ntHash, Mohamadi et al., 2016; DSK, Rizk et al., 2013; and KMC, Deorowicz et al., 2013). In the realm of transcriptome analysis, tools such as Kallisto (Bray et al., 2016) and Sailfish (Patro et al., 2014) are capable of accurate RNA-sequencing quantification at a fraction of the time and computational cost of previous mapping-based strategies. A recent study has also introduced a novel mapping-free method for performing genome-wide association studies from whole-genome sequence data (Rahman et al., 2018) using *k*-mer counts. The tool HAWK (Rahman et al., 2018) performs rapid and accurate discovery of variant-phenotype associations by directly comparing *k*-mer frequencies between arbitrary numbers of case and control samples. HAWK counts all *k*-mers in the

sequenced samples and finds  $k$ -mers that are significantly associated with the phenotype or trait of interest (“significant  $k$ -mers”), and then performs a local assembly of these significant  $k$ -mers to predict the corresponding significant variants associated with the traits. This approach provides an efficient method for discovery of significant associations between all types of variants (i.e., SNVs, indels, and SVs) and the phenotype or trait of interest (Rahman et al., 2018).

Developments in variant prediction frameworks continue to spur improvements in a variety of contexts. Scalpel (Narzisi et al., 2014) implements a hybrid method for *de novo* indel discovery from whole-exome sequencing of quads. Read mapping is used only to localize reads to the reference genome. In subsequent steps, Scalpel performs localized *de novo* assembly of reads at loci of interest and aligns assembled contigs back to the loci to annotate any *de novo* variants present (Narzisi et al., 2014). More recently, NovoBreak (Chong et al., 2017) introduced a method that utilizes  $k$ -mer counts to predict somatic variants, including SVs, by comparison of paired tumor and normal whole-genome sequence samples. COBASI (Gómez-Romero et al., 2018) performs rapid and accurate *de novo* SNV discovery on whole-genome sequencing of trios by computing perfect matches to unique strings in the reference genome and then identifying abrupt shifts in the coverage of the resulting alignments. Finally, mapping-free approaches such as LAVA (Shajii et al., 2016), VarGeno (Sun and Medvedev, 2018), MALVA (Bernardini et al., 2019), and Nebula (Khorsand and Hormozdiari, 2019) were recently developed for fast and accurate genotyping of common variation using whole-genome sequencing data.

The present study introduces a new mapping-free strategy grounded on a  $k$ -mer-based formulation of the *de novo* variant discovery problem—see Figure 1A. Intuitively, a novel germline mutation should result in new sequence content in the proband compared with the parental genomes. Even in the simplest case, a single-nucleotide substitution, most of the  $k$ -mers spanning the mutation should be unique, given a sufficiently large value of  $k$ . Incidentally, this is also true for other classes of variants, such as indels and various types of structural variation. And with sufficiently deep sampling of the proband genome, the expectation is that these novel  $k$ -mers are present in the read data at levels that can be readily distinguished from sequencing errors. Thus, it should be possible to detect both SNVs and larger variants (indels, SVs) simultaneously using a single mapping-free model.

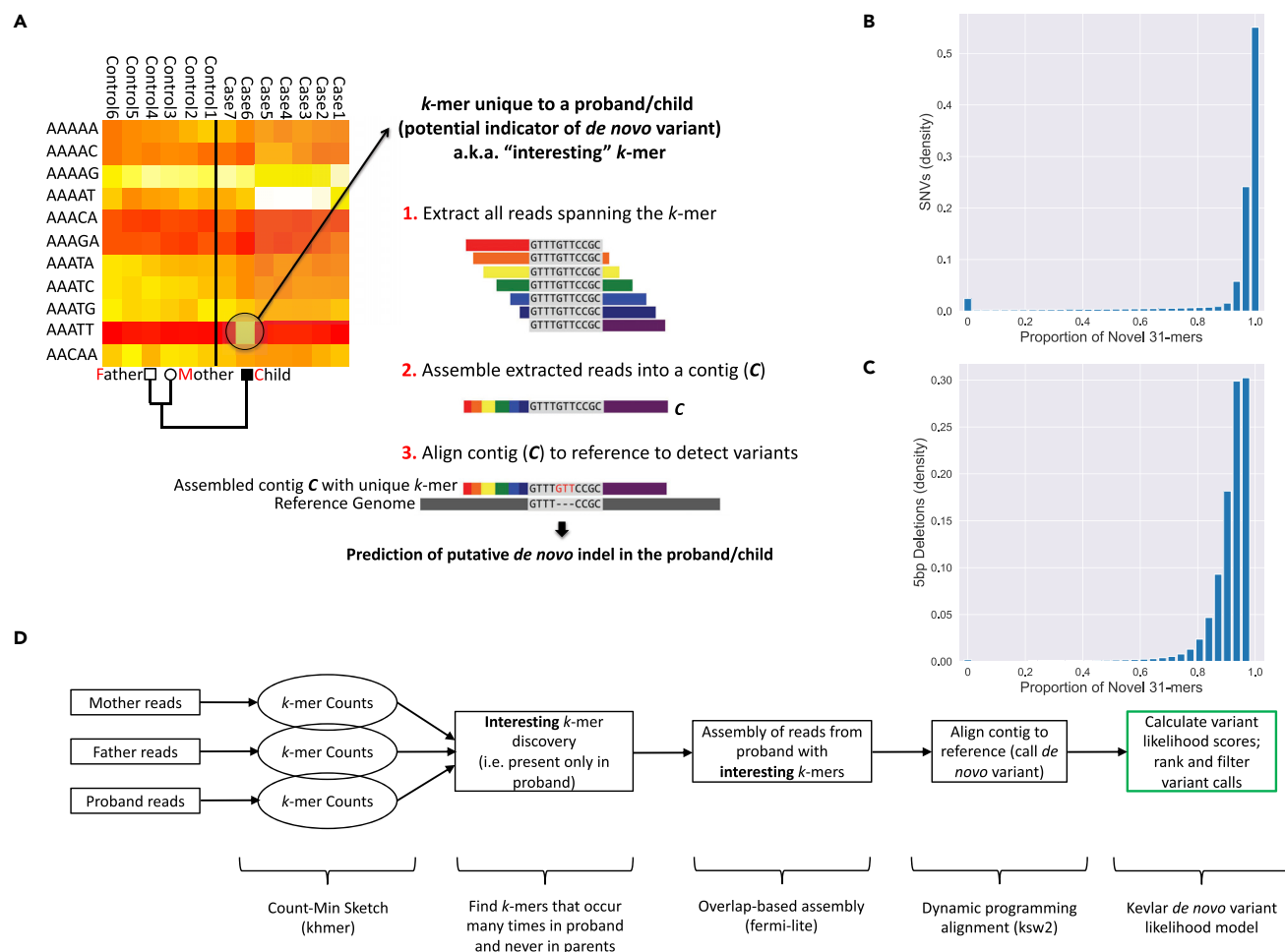
Building on this intuition, we developed Kevlar, a new method based on a mapping-free formulation of the *de novo* variant discovery problem. Kevlar examines  $k$ -mer abundances to identify “interesting”  $k$ -mers, which we define as having significantly high abundance in the proband or child reads, whereas being effectively absent in the reads from both parents. These interesting  $k$ -mers are an indicator of the potential existence of a *de novo* variant in the proband and are conceptually similar to the “significant”  $k$ -mers used by HAWK (Rahman et al., 2018). We next group the reads containing interesting  $k$ -mers into disjoint sets, each reflecting a putative variant, based on the  $k$ -mers shared between pairs of reads. Kevlar then uses standard algorithms to assemble each set of reads into contigs and align the assembled contigs to a reference genome to make preliminary variant calls. Finally, Kevlar employs a probabilistic model to score predicted variants to distinguish between miscalled inherited variants and true *de novo* mutations.

We demonstrate the utility of this new method on simulated and real data. We show that Kevlar achieves similar predictive performance to best-in-class tools for SNV and short indel discovery, while at the same time predicting larger events with high accuracy. We also demonstrate Kevlar’s ability to accurately predict large-scale SV events, defining breakpoints with nucleotide-level precision.

Kevlar is available as an open source software project and can be invoked via a Python API, a command-line interface, or a standard Snakemake workflow (Köster and Rahmann, 2012). The stable and actively developed source code is available at <https://github.com/kevlar-dev/kevlar>, and documentation is available at <https://kevlar.readthedocs.io>.

## RESULTS

We present a novel framework for discovery of *de novo* variants based on direct comparisons of sequence content between related individuals, requiring no mapping of short reads to a reference genome. This framework utilizes a single strategy that accurately predicts SNVs, insertions and deletions (indels), and structural variation events simultaneously.



**Figure 1. Overview of Kevlar**

(A) Visual summary of the mapping-free approach for *de novo* variant discovery.

(B) The likelihood that novel mutation results in unique mutation-spanning k-mers, determined by simulating single-nucleotide substitutions genome-wide and measuring the proportion of SNV-spanning k-mers that are not observed elsewhere in the genome. The trend observed for  $k = 31$  holds for a wide range of  $k$  values (approximately 20–60).

(C) The same as (B) except for 5-bp deletion mutations.

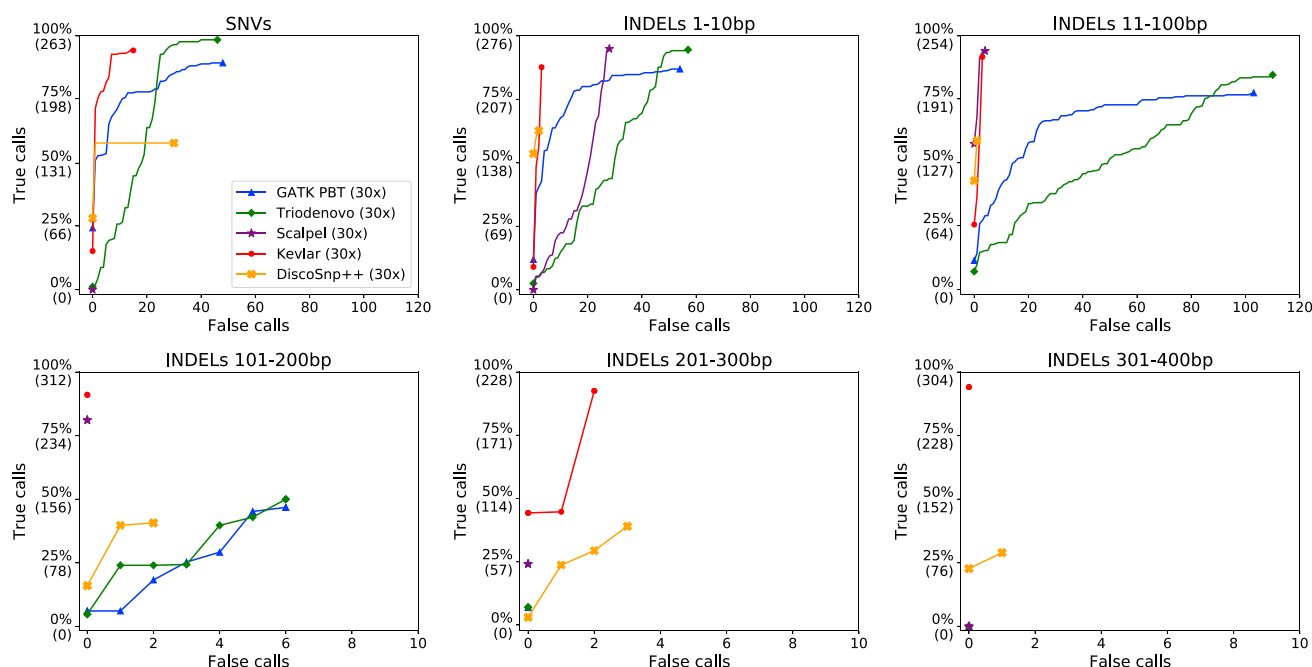
(D) An overview of the Kevlar workflow.

## Overview of Kevlar

Our variant discovery strategy is fundamentally a search for novel DNA content in the sample of interest. It is based on the observation that k-mers (short subsequences of fixed length  $k$ ) spanning a *de novo* mutation will be novel with high probability (Figures 1B and 1C). Often the subject is a child affected by a disorder or other trait of interest (referred to as *proband*), with related individuals being the two parents.

Figure 1D summarizes the Kevlar workflow. In brief, DNA sequence reads from the case and control samples are processed independently. For each sample, the reads are split into k-mers and the abundance of each k-mer is stored for subsequent lookup. A second pass over reads from the case sample then identifies all k-mers that are unique to the proband—that is, k-mers that are abundant in the proband but effectively absent in both parents. Reads containing any novel k-mers are retained for subsequent processing.

After applying filters for contamination and erroneous k-mer abundances, the reads containing novel k-mers are partitioned such that any two reads sharing at least one novel k-mer are grouped together.



**Figure 2. Accuracy of Five De Novo Variant Prediction Algorithms**

Receiver operating characteristic (ROC) curves compare variant prediction performance on a simulated dataset. Average sequencing depth was approximately 30x. Each of the six panes shows prediction accuracy for a different variant type: SNVs, insertions or deletions (indels) 1–10 bp in length, 11- to 100-bp indels, 101- to 200-bp indels; 201- to 300-bp indels; and 301- to 400-bp indels. Note that the scale of the x axis for long indels is an order of magnitude smaller than the x axis scale for SNVs and short (< 100 bp) indels.

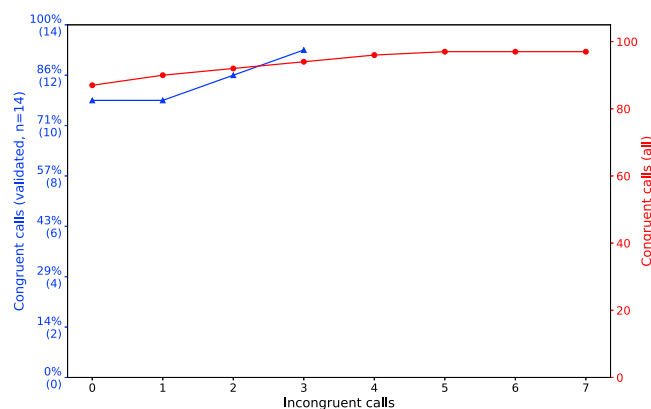
The reads in each partition are then analyzed independently: they are assembled into a contig, the contig is aligned to the reference genome, and the alignment is used to assess the presence of a variant and make a variant call. Finally, Kevlar employs a likelihood-based score to rank and filter the variant calls.

Each step of the Kevlar workflow is discussed in detail in the [Transparent Methods](#).

### Performance on Simulated Data

We simulated whole-genome shotgun sequencing of a mock family for a fine-grained assessment of Kevlar's accuracy in identifying different variant types at different levels of sequencing depth. Our simulation not only realistically modeled the inheritance of parental variants but also included hundreds of "de novo" (unique to the proband) SNVs and indels ranging in size from <10 to 400 bp. The sequencing was simulated at 10x, 20x, 30x, and 50x coverage with low error rate. We compared Kevlar's accuracy on this dataset with two widely used mapping-based *de novo* variant callers (GATK PhaseByTransmission, [Francioli et al., 2016](#); and TrioDenovo, [Wei et al., 2015](#)) as well as two mapping-free or hybrid variant callers (Scalpel, [Narzisi et al., 2014](#); and DiscoSnp++, [Peterlongo et al., 2017](#)).

The accuracy of all variant callers evaluated is poor at low (10x) coverage (see [Figure S1](#)). GATK PhaseByTransmission makes very few variant predictions at 10x coverage. The remaining variant callers report numerous predictions, but in general suffer from both low sensitivity (failing to predict many true variants) and poor specificity (predicting many false variants). TrioDenovo shows the best prediction performance for SNVs and short (1–100 bp) indels at 10x coverage. At 20x coverage ([Figure S2](#)), all five algorithms show marked improvement in SNV detection, in particular TrioDenovo, which achieves  $\geq 90\%$  sensitivity. Scalpel exhibits both improved sensitivity and improved specificity at 20x coverage and approaches or surpasses TrioDenovo's performance for indels of most lengths. Kevlar's ability to accurately detect indels >100 bp becomes evident at 20x coverage.



**Figure 3. Performance of Kevlar on SSC Trio 14153**

ROC curves showing congruence between *de novo* variant calls made by Kevlar on the SSC 14153 trio and corresponding calls from the denovo-db variant database. The red curve shows Kevlar's performance compared with all denovo-db calls, and the blue curve shows Kevlar's performance compared with denovo-db calls with experimental validation.

At higher levels of coverage (30x and 50x), Kevlar consistently achieves top performance across all variant types (see Figures 2 and S3). Notably, Kevlar recovers  $\geq 90\%$  of true variants while making very few false predictions across all variant types at high coverage. TrioDenovo shows marginally better sensitivity than Kevlar for predicting SNVs at 30x and 50x (as does GATK PhaseByTransmission at 50x), but at the expense of numerous false predictions. Kevlar also rivals Scalpel's impressive short indel prediction performance and exceeds it for predicting long (>100 bp) indels.

### Performance on the SSC 14153 Autism Trio

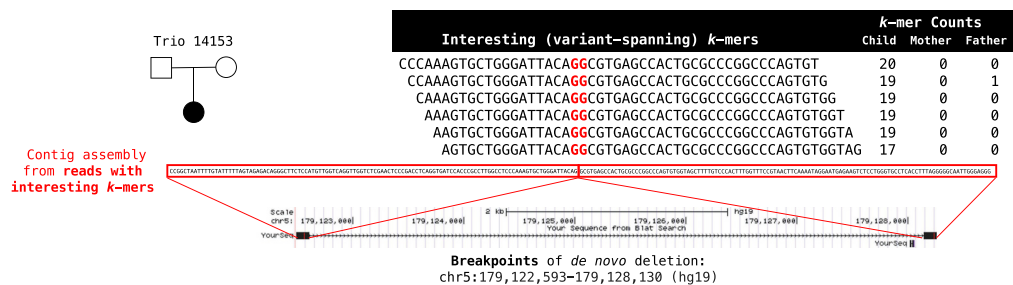
To assess Kevlar's performance on real data, we applied Kevlar to predict *de novo* variants in the proband of an autism trio from the Simons Simplex Collection (family 14153). As a reference for comparison, we obtained a potential "truth set" from the denovo-db database (<http://denovo-db.gs.washington.edu/denovo-db/>). This truth set includes 196 *de novo* variant predictions and represents the union of predictions made for this trio by several recent studies (Turner et al., 2016, 2017; Werling et al., 2018). Note that the expected number of *de novo* variants per generation is estimated to be around 100 (Campbell and Eichler, 2013; Turner et al., 2016), or about half of the number of predictions in the truth set. Annotations in the denovo-db database indicate that experimental validation has confirmed 14 of the 196 calls.

In total, Kevlar predicts 219 *de novo* variants for trio 14153, including 150 SNVs, 68 indels/SVs, and a single 2-bp multinucleotide variant. We note that Kevlar assigned many of these predicted variants a low likelihood of the variant being a true *de novo* event. Figure 3 shows the congruence between the 100 top-ranked Kevlar calls and the denovo-db calls for this trio.

Of the 14 denovo-db calls with experimental validation, 13 (92.9%) were predicted accurately by Kevlar and assigned a high likelihood score, indicative of a confident *de novo* variant call. Overall, the 100 Kevlar variant calls ranked highest by the likelihood score include only four calls not present in denovo-db (probable false calls). On the other hand, only five Kevlar variant calls present in denovo-db (probable true variants) are not among the 100 highest ranked Kevlar calls. Of the 196 denovo-db calls, 95 are absent from the Kevlar predictions. The majority of these calls (75/95,  $\approx 80\%$ ) occur in regions of repetitive DNA and have shown to be unreliable in experimental validation (Tychele Turner, personal communication).

Finally, a recent study verified the presence of a *de novo* deletion of approximately 6 kbp in the proband of this trio (Turner et al., 2016), removing the 5' UTR of the gene *CANX*. Kevlar also predicted this *de novo* deletion successfully and identified the precise (and previously undetermined) breakpoints at chr5:179,122,593 and chr15:179,128,130 (GRCh37). Inspection of the variant reveals that both the deletions' breakpoints occur in *Alu* repeat elements abundant throughout the genome (Figure 4). As a result, only seven of the *k*-mers spanning the variant are unique signatures of mutation not already present elsewhere





**Figure 4. An Experimentally Validated 6-kbp De Novo Deletion as Predicted by Kevlar**

The interesting k-mers, their abundances in each sample, the variant-spanning contig assembly, and the breakpoints are depicted.

in the genome. We observe with interest that both breakpoints occur inside a 20-bp identical repeat, indicating this *de novo* deletion is the result of non-allelic homologous recombination.

## DISCUSSION

*De novo* variants are a major contributing factor in many disorders (e.g., intellectual disability, autism, and epilepsy). Accurate discovery of these variants has been challenging as prediction methods need to be confident not only in the existence of the event in the proband or child but also in the absence of the variant in the parents. Current approaches depend on correct alignments of sequence reads to a reference genome. Any complications in computing read alignments due to repeats, gaps in the reference, or variant complexity can result in false predictions or failure to discover a true *de novo* variant.

The method proposed in this study compares *k*-mers between related individuals to find the *k*-mers indicating a *de novo* variant in the sample of interest. We acknowledge recently proposed methods Novo-Break (Chong et al., 2017) and HAWK (Rahman et al., 2018), which are conceptually similar and likewise capable of accurately predicting *de novo* variants. Kevlar, HAWK (Rahman et al., 2018), and other related methods do not depend on mapping reads to a reference genome, but instead rely on direct comparison of sequence content between related individuals. This strategy enables Kevlar to accurately predict several classes of *de novo* mutations (substitutions, insertions, deletions, SVs) simultaneously with a single simple mathematical model. As long as the *de novo* mutation creates a *k*-mer not already present in the reference genome, the proposed algorithm should be able to accurately discover the event. We have also developed a *k*-mer-based likelihood model for scoring and ranking variant calls according to their probability of being true *de novo* events. This likelihood score is effective in discerning *de novo* variants from inherited mutations and false variant calls. We have demonstrated the effectiveness of our discovery method and scoring model using both simulated and real data. Kevlar is competitive with best-in-class tools for discovery of a variety of variant types, and substantially outperforms available methods for discovery of larger *de novo* variants. Kevlar not only predicts indels and SVs with high sensitivity and specificity but also reports the exact breakpoints of these variants with single base pair precision.

*De novo* variants are, by definition, expected to be unique for each individual. Aggregating multiple simplex trios will not increase the rate of recall. However, multiple trios could potentially be aggregated to identify any systematic errors resulting in the same *k*-mers being marked as “interesting” in multiple samples. Identifying and removing these *k*-mers and any corresponding variant calls could improve precision.

Development of completely reference-free methods is tremendously valuable in scenarios where the availability, quality, or completeness of a reference genome is insufficient. Kevlar’s preliminary steps—identifying variant-spanning reads, binning reads into groups corresponding to distinct putative variants, and assembling each read group into a variant-spanning contig—are performed without the use of a reference genome. We note, however, that subsequent steps in the Kevlar workflow to annotate, filter, and score the preliminary variant calls still depend on a reference genome. One promising approach to developing a completely reference-free *de novo* variant discovery method would be to annotate variants by aligning variant-spanning contigs directly to an assembly or variation graph.



### Limitations of the Study

Misclassification of heterozygous inherited variants as *de novo* is one of the main sources of false prediction. These errors are enriched at loci with low coverage in the donor parent. This is due to the difficulty of distinguishing true variation from sequencing error. It is possible that utilizing a probabilistic approach for selecting “interesting” *k*-mers, as proposed in HAWK (Rahman et al., 2018), can reduce the false *de novo* prediction rate.

Kevlar will successfully annotate *k*-mers that span the breakpoints of large insertions. It will also assemble the reads containing these *k*-mers into breakpoint-spanning contigs. However, unless the inserted sequence is entirely novel, Kevlar is unlikely to assemble a single contig that spans the entire variant and is thus capable of annotating its precise coordinates.

Even using a probabilistic *k*-mer counting strategy, Kevlar’s memory requirements can be quite demanding. Applying error correction to the input reads will substantially reduce Kevlar’s memory requirements, but this typically leads to a small reduction in sensitivity for discovering SNVs.

Finally, in scoring and ranking of the predicted *de novo* variants Kevlar assumes independence between *k*-mers in likelihood calculation. While this assumption simplifies the likelihood calculation, a more sophisticated formulation that does not have this limitation may yield improvements in scoring and ranking the final variant calls.

### METHODS

All methods can be found in the accompanying [Transparent Methods supplemental file](#).

### DATA AND CODE AVAILABILITY

The Kevlar software is hosted as an open source software project at <https://github.com/kevlar-dev/kevlar> and is freely available under the MIT license. User documentation is available at <https://kevlar.readthedocs.io>. Reads from the simulated dataset are available in FASTQ format from DOI <https://doi.org/10.1706/ODF.IO/4CHPB>. Reads from the 14153 trio are available in BAM format from the Simons Simplex Collection at <https://www.sfari.org/2015/12/11/whole-genome-analysis-of-the-simons-simplex-collection-ssc-2/#chapter-how-to-access-the-data>.

### SUPPLEMENTAL INFORMATION

Supplemental Information can be found online at <https://doi.org/10.1016/j.isci.2019.07.032>.

### ACKNOWLEDGMENTS

We would like to acknowledge Dr. Tamer Mansour, Luiz Irber Jr., Camille Scott, and Lisa Johnson for helpful discussions on method development and implementation and Dr. Tychele Turner for helpful discussions on the method evaluation. We also thank reviewers and several colleagues for comments on earlier versions of the manuscript, which have improved the final paper.

This work is funded in part by the Gordon and Betty Moore Foundation’s Data-Driven Discovery Initiative through Grant GBMF4551 and NIH R01 HG007513, both to C.T.B., and by the Sloan Research Fellowship number FG-2017-9159 to F.H..

### AUTHOR CONTRIBUTIONS

D.S.S., C.T.B., and F.H. conceived the study. D.S.S. implemented the method and performed the experiments. D.S.S. and F.H. designed the experiments and wrote the manuscript. D.S.S., C.T.B., and F.H. edited and approved the final manuscript.

### DECLARATION OF INTERESTS

The authors declare no competing interests.

Received: February 11, 2019

Revised: June 24, 2019

Accepted: July 19, 2019

Published: August 30, 2019

## REFERENCES

- Bernardini, G., Bonizzoni, P., Denti, L., Previtali, M., and Schönhuth, A. (2019). Malva: genotyping by mapping-free allele detection of known variants. *bioRxiv*, 575126.
- Bray, N.L., Pimentel, H., Melsted, P., and Pachter, L. (2016). Near-optimal probabilistic RNA-seq quantification. *Nat. Biotechnol.* 34, 525.
- Campbell, C.D., and Eichler, E.E. (2013). Properties and rates of germline mutations in humans. *Trends Genet.* 29, 575–584.
- Cardno, A.G., Marshall, E.J., Coid, B., Macdonald, A.M., Ribchester, T.R., Davies, N.J., Venturi, P., Jones, L.A., Lewis, S.W., Sham, P.C., et al. (1999). Heritability estimates for psychotic disorders: the Maudsley twin psychosis series. *Arch. Gen. Psychiatry* 56, 162–168.
- Chong, Z., Ruan, J., Gao, M., Zhou, W., Chen, T., Fan, X., Ding, L., Lee, A.Y., Boutros, P., Chen, J., et al. (2017). novobreak: local assembly for breakpoint detection in cancer genomes. *Nat. Methods* 14, 65.
- Crusoe, M.R., Alameldin, H.F., Awad, S., Boucher, E., Caldwell, A., Cartwright, R., Charbonneau, A., Constantinides, B., Edverson, G., Fay, S., et al. (2015). The Khmer software package: enabling efficient nucleotide sequence analysis. *F1000Res.* 4, 900.
- Deorowicz, S., Debudaj-Grabysz, A., and Grabowski, S. (2013). Disk-based k-mer counting on a pc. *BMC Bioinformatics* 14, 160.
- Eichler, E.E., Flint, J., Gibson, G., Kong, A., Leal, S.M., Moore, J.H., and Nadeau, J.H. (2010). Missing heritability and strategies for finding the underlying causes of complex disease. *Nat. Rev. Genet.* 11, 446.
- Francioli, L.C., Cretu-Stancu, M., Garimella, K.V., Fromer, M., Kloosterman, W.P., Genome of the Netherlands consortium, Samocha, K.E., Neale, B.M., Daly, M.J., Banks, E., DePristo, M.A., and de Bakker, P.I. (2016). A framework for the detection of *de novo* mutations in family-based sequencing data. *Eur. J. Hum. Genet.* 25, 227–233.
- Fromer, M., Pocklington, A.J., Kavanagh, D.H., Williams, H.J., Dwyer, S., Gormley, P., Georgieva, L., Rees, E., Palta, P., Ruderfer, D.M., et al. (2014). *De novo* mutations in schizophrenia implicate synaptic networks. *Nature* 506, 179.
- Gómez-Romero, L., Palacios-Flores, K., Reyes, J., García, D., Boege, M., Dávila, G., Flores, M., Schatz, M.C., and Palacios, R. (2018). Precise detection of *de novo* single nucleotide variants in human genomes. *Proc. Natl. Acad. Sci. U S A* 115, 5516–5521.
- Hallmayer, J., Cleveland, S., Torres, A., Phillips, J., Cohen, B., Torigoe, T., Miller, J., Fedele, A., Collins, J., Smith, K., et al. (2011). Genetic heritability and shared environmental factors among twin pairs with autism. *Arch. Gen. Psychiatry* 68, 1095–1102.
- Hormozdiari, F., Alkan, C., Eichler, E.E., and Sahinalp, S.C. (2009). Combinatorial algorithms for structural variation detection in high-throughput sequenced genomes. *Genome Res.* 19, 1270–1278.
- Iossifov, I., O’Roak, B.J., Sanders, S.J., Ronemus, M., Krumm, N., Levy, D., Stessman, H.A., Witherspoon, K.T., Vives, L., Patterson, K.E., et al. (2014). The contribution of *de novo* coding mutations to autism spectrum disorder. *Nature* 515, 216.
- Iqbal, Z., Caccamo, M., Turner, I., Flicek, P., and McVean, G. (2012). *De novo* assembly and genotyping of variants using colored de bruijn graphs. *Nat. Genet.* 44, 226.
- Khorsand, P., and Hormozdiari, F. (2019). Nebula: Ultra-efficient mapping-free structural variant genotyper. *bioRxiv*, 566620.
- Köster, J., and Rahmann, S. (2012). Snakemake: a scalable bioinformatics workflow engine. *Bioinformatics* 28, 2520–2522.
- Layer, R.M., Chiang, C., Quinlan, A.R., and Hall, I.M. (2014). Lumpy: a probabilistic framework for structural variant discovery. *Genome Biol.* 15, R84.
- Manolio, T.A., Collins, F.S., Cox, N.J., Goldstein, D.B., Hindorf, L.A., Hunter, D.J., McCarthy, M.I., Ramos, E.M., Cardon, L.R., Chakravarti, A., et al. (2009). Finding the missing heritability of complex diseases. *Nature* 461, 747.
- Marçais, G., and Kingsford, C. (2011). A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics* 27, 764–770.
- Medvedev, P., Fiume, M., Dzamba, M., Smith, T., and Brudno, M. (2010). Detecting copy number variation with mated short reads. *Genome Res.* 20, 1613–1622.
- Mohamadi, H., Chu, J., Vandervalk, B.P., and Birol, I. (2016). ntHash: recursive nucleotide hashing. *Bioinformatics* 32, 3492–3494.
- Narzisi, G., O’Rawe, J.A., Iossifov, I., Fang, H., Lee, Y.-h., Wang, Z., Wu, Y., Lyon, G.J., Wigler, M., and Schatz, M.C. (2014). Accurate *de novo* and transmitted indel detection in exome-capture data using microassembly. *Nat. Methods* 11, 1033.
- O’Roak, B.J., Vives, L., Girirajan, S., Karakoc, E., Krumm, N., Coe, B.P., Levy, R., Ko, A., Lee, C., Smith, J.D., et al. (2012). Sporadic autism exomes reveal a highly interconnected protein network of *de novo* mutations. *Nature* 485, 246.
- Patro, R., Mount, S.M., and Kingsford, C. (2014). Sailfish enables alignment-free isoform quantification from rna-seq reads using lightweight algorithms. *Nat. Biotechnol.* 32, 462.
- Peterlongo, P., Riou, C., Drezen, E., and Lemaître, C. (2017). Discosnp++: *de novo* detection of small variants from raw unassembled read set(s). *bioRxiv*, 209965.
- Rahman, A., Hallgrímsson, I., Eisen, M., and Pachter, L. (2018). Association mapping from sequencing reads using k-mers. *Elife* 7, e32920.
- Rausch, T., Zichner, T., Schlattl, A., Stütz, A.M., Benes, V., and Korb, J.O. (2012). DELLY: structural variant discovery by integrated paired-end and split-read analysis. *Bioinformatics* 28, i333–i339.
- Rizk, G., Lavenier, D., and Chikhi, R. (2013). DSK: k-mer counting with very low memory usage. *Bioinformatics* 29, 652–653.
- Shajii, A., Yorukoglu, D., William Yu, Y., and Berger, B. (2016). Fast genotyping of known snps through approximate k-mer matching. *Bioinformatics* 32, i538–i544.
- Sindi, S.S., Önal, S., Peng, L.C., Wu, H.-T., and Raphael, B.J. (2012). An integrative probabilistic model for identification of structural variation in sequencing data. *Genome Biol.* 13, R22.
- Soylev, A., Kockan, C., Hormozdiari, F., and Alkan, C. (2017). Toolkit for automated and rapid discovery of structural variants. *Methods* 129, 3–7.
- Sun, C., and Medvedev, P. (2018). Toward fast and accurate SNP genotyping from whole genome sequencing data for bedside diagnostics. *bioRxiv*, 239871.
- Turner, T.N., Coe, B.P., Dickel, D.E., Hoekzema, K., Nelson, B.J., Zody, M.C., Kronenberg, Z.N., Hormozdiari, F., Raja, A., Pennacchio, L.A., et al. (2017). Genomic patterns of *de novo* mutation in simplex autism. *Cell* 171, 710–722.
- Turner, T.N., Hormozdiari, F., Duyzend, M.H., McClymont, S.A., Hook, P.W., Iossifov, I., Raja, A., Baker, C., Hoekzema, K., Stessman, H.A., et al. (2016). Genome sequencing of autism-affected families reveals disruption of putative noncoding regulatory dna. *Am. J. Hum. Genet.* 98, 58–74.
- Uricaru, R., Rizk, G., Lacroix, V., Quillery, E., Plantard, O., Chikhi, R., Lemaître, C., and Peterlongo, P. (2014). Reference-free detection of isolated snps. *Nucleic Acids Res.* 43, e11.
- Veltman, J.A., and Brunner, H.G. (2012). *De novo* mutations in human genetic disease. *Nat. Rev. Genet.* 13, 565.
- Wei, Q., Zhan, X., Zhong, X., Liu, Y., Han, Y., Chen, W., and Li, B. (2015). A Bayesian framework for *de novo* mutation calling in parents-offspring trios. *Bioinformatics* 31, 1375–1381.
- Werling, D.M., Brand, H., An, J.-Y., Stone, M.R., Zhu, L., Glessner, J.T., Collins, R.L., Dong, S., Layer, R.M., Markenscoff-Papadimitriou, E., et al. (2018). An analytical framework for whole-genome sequence association studies and its implications for autism spectrum disorder. *Nat. Genet.* 50, 727–736.
- Ye, K., Schulz, M.H., Long, Q., Apweiler, R., and Ning, Z. (2009). Pindel: a pattern growth approach to detect break points of large deletions and medium sized insertions from paired-end short reads. *Bioinformatics* 25, 2865–2871.
- Zaidi, S., Choi, M., Wakimoto, H., Ma, L., Jiang, J., Overton, J.D., Romano-Adesman, A., Bjornson, R.D., Breitbart, R.E., Brown, K.K., et al. (2013). *De novo* mutations in histone-modifying genes in congenital heart disease. *Nature* 498, 220.

**ISCI, Volume 18**

**Supplemental Information**

**Kevlar: A Mapping-Free  
Framework for Accurate  
Discovery of *De Novo* Variants**

**Daniel S. Standage, C. Titus Brown, and Fereydoun Hormozdiari**

## Supplemental Information

# Kevlar: a mapping-free framework for accurate discovery of *de novo* variants

Daniel S. Standage, C. Titus Brown, Fereydoon Hormozdiari

## Transparent Methods

### *Assessing diagnostic utility of novel k-mers*

We expect that a *de novo* mutation will result in numerous novel *k*-mers, given a sufficiently large value of *k*. We also expect that these novel *k*-mers will be present in high abundance, given sufficiently deep sampling of the proband genome. Intuitively, we can use these novel *k*-mers to identify reads that span the *de novo* variant—see Figure 1a.

We assessed this intuition by traversing the human reference genome (GRCh38) base by base, simulating variants (SNVs and 5 bp deletions) at each position. For each simulated mutation, we determined the fraction of *k*-mers spanning the mutation that exist nowhere else in the genome, and thus act as a diagnostic signature of that particular variant. We then aggregated over the entire genome the probability that *k*-mer spanning a mutation (in this case 31-mers) will be novel—see Figure 1b and 1c.

Based on the results of this experiment, we formulate the *de novo* variant discovery problem as a search for putatively novel *k*-mers that are abundant in the proband and effectively absent in each parent. For sake of simplicity, we are using the term *proband* to refer generally to the subject or focal individual, and *parent* to refer generally to control individuals.

Here, *abundant* and *effectively absent* are defined in terms of a simple threshold model. Let *X* be the *absence threshold*, and *Y* be the *presence threshold*, and  $A = \{A_p, A_m, A_f\}$  be the abundances of a *k*-mer in the proband, mother, and father. We designate this *k*-mer as “*interesting*” (putatively novel) if and only if  $A_p \geq Y$ ,  $A_m \leq X$ , and  $A_f \leq X$ . Based on our experience, the values  $Y = 5$  and  $X = 1$  produce desirable results for 30x sequencing coverage.

### *Kevlar workflow*

The steps of the Kevlar workflow, summarized at a high level in Figure 1d, are described in detail in the subsequent sections.

### *Step 0: Compute k-mer counts*

Preliminary to identifying novel *k*-mers, the abundance of each *k*-mer in each sample must be counted. Storing exact counts of every *k*-mer requires a substantial amount of space (dozens of gigabytes or more per sample), so Kevlar exploits several strategies to reduce the space required for keeping *k*-mer counts in memory.

First, Kevlar stores approximate *k*-mer counts in a Count-Min sketch, a probabilistic data structure similar to a Bloom filter that operates in a fixed amount of memory, exchanging accuracy for space efficiency (Zhang et al., 2014). A separate Count-Min sketch is used for each sample. The accuracy of each Count-Min sketch depends on its size and the number of distinct elements (*k*-mers in this case) being tracked. The Count-Min sketch exhibits a one-sided error, meaning that *k*-mer counts are sometimes overestimated but never underestimated. The extent of inaccuracy in the *k*-mer counts is summarized by the *false discovery rate (FDR)* statistic computed from the occupancy of the Count-Min sketch.

Second, Kevlar uses a masked counting strategy in which *k*-mers present in the reference genome and a contaminant database (composed of bacterial, viral, vector, and adapter sequences) are ignored. This substantially reduces the number of *k*-mers to be stored in the Count-Min sketch, and as a consequence the desired level of accuracy can be maintained using a smaller amount of space.

Third, *k*-mer counts are recomputed with exact precision in subsequent steps of the Kevlar workflow, which means any *k*-mer retained erroneously due to an inflated count can be compensated for at a later stage. As a consequence, it is possible to reduce the size of the Count-Min sketch even further, resulting in a FDR of 0.5 or greater.

Kevlar’s *k*-mer counting operations are invoked with the `kevlar count` command, and rely on bulk sequence loading procedures and an implementation of the Count-Min sketch data structure from the `khmer` library (Crusoe

et al., 2015; Standage et al., 2017; Zhang et al., 2014). Note that several alternative  $k$ -mer counting libraries and tools (Marçais and Kingsford, 2011; Rizk et al., 2013) have been developed and utilized to solve a variety of different biological problems (Bray et al., 2016; Patro et al., 2014; Rahman et al., 2018; Sun and Medvedev, 2018).

### *Step 1: Identifying novel $k$ -mers and reads*

To identify sequences spanning *de novo* variants, Kevlar scans each read sequenced from the proband. The per-sample abundances of each  $k$ -mer are queried from the Count-Min sketches computed in previous steps. If a  $k$ -mer is present in high abundance in the proband and absent from the parents (that is, it satisfies user-specified abundance thresholds), it is designated as “interesting” or putatively novel. This operation is similar to the selection of “*novo*”  $k$ -mers by NovoBreak (Chong et al., 2017) and “significant”  $k$ -mers by HAWK (Rahman et al., 2018). Any read containing one or more interesting  $k$ -mers is retained for subsequent processing. This step is implemented in the `kevlar novel` command.

### *Step 2: Contamination, reference, and abundance filters*

Reads containing putative novel  $k$ -mers are filtered prior to subsequent analysis. This filtering step serves two purposes.

First, Kevlar re-computes the abundance of each interesting  $k$ -mer in the proband sample. The relatively small volume of these reads allows Kevlar to re-compute  $k$ -mer counts with perfect accuracy in a small amount of memory and time. Any  $k$ -mer whose corrected count no longer satisfies the required abundance threshold is discarded. Note that since only proband reads are retained, only the proband  $k$ -mer abundances can be recomputed. This filtering step will not recover a  $k$ -mer that is erroneously discarded in the previous step due to an erroneously inflated  $k$ -mer abundance in one of the control (parent and sibling) samples.

Second, if for any reason  $k$ -mers from the reference genome and contaminants are not ignored in the initial  $k$ -mer counting step, this filtering step provides another opportunity to discard these  $k$ -mers.

After these filters are applied, any read that no longer contains any novel  $k$ -mers is discarded, and the remainder of the reads are retained for subsequent analysis.

The `kevlar filter` command is used to execute these contamination, reference, and abundance filters.

### *Step 3: Partitioning reads using shared novel $k$ -mers*

Interesting reads spanning the same variant are expected to share numerous interesting  $k$ -mers. These shared novel  $k$ -mers provide a mechanism for grouping the reads into disjoint sets reflecting distinct variants.

To be precise, we define a *read graph*  $G$  as follows: every read containing one or more novel  $k$ -mers is represented by a node in  $G$ , and a pair of nodes is connected by an edge if they have one or more novel  $k$ -mers in common. With this formulation, if two reads share a novel  $k$ -mer they are part of the same connected component in  $G$ . Overall  $G$  is sparse, but typically each connected component of the graph is highly connected. In subsequent steps, each component or partition  $p \in G$  is analyzed independently.

The `kevlar partition` command implements this partitioning strategy.

### *Step 4: Contig assembly and reference target selection*

For each connected component  $p \in G$ , we assemble the corresponding reads using the overlap-based algorithm implemented in the *fermi-lite* library (Li, 2017a). Briefly, *fermi-lite* performs error correction, trims reads at unique  $l$ -mers, constructs an FM-index of the trimmed reads, and constructs a transitively reduced overlap graph. The optimal path in the final graph is output as a contig  $C_p$  suitable for variant calling.

Next, we select a target reference sequence (or set of candidate targets) for the contig  $C_p$ . Briefly, Kevlar decomposes the contig into overlapping subsequences of length  $l$  (*seeds*;  $l = 51$  by default), and uses BWA MEM (Li, 2013) to identify locations of exact matches for each seed sequence in the reference genome. The genomic interval that spans all seed exact matches, plus  $\Delta$  nucleotides in each direction ( $\Delta = 50$  by default), is then selected as the target reference sequence for  $C_p$ . If any adjacent seed matches are separated by more than  $D$  nucleotides ( $D = 10,000$  by default), then the seed matches are split at that point and multiple reference targets are selected. The set of reference target sequences corresponding to contig  $C_p$  is denoted  $T_{C_p}$ .

Read assembly is invoked with the `kevlar assemble` command, and reference target selection is invoked with the `kevlar localize` command.

### Step 5: Contig alignment and variant annotation

The contig  $C_p$  is aligned to each reference target sequence  $t \in T_{C_p}$  using the ksw2 library (Li, 2017b)—specifically its implementation of Green’s formulation of dynamic programming global alignment and extension (`ksw2_extz`). If there are multiple candidate targets, only the highest scoring alignment is retained. When a contig aligns to multiple locations with the same optimal score, all optimal alignments are retained for variant calling.

Prior to variant calling, kevlar right-aligns any gaps at the right end of the alignment to minimize the number of alignment blocks/operations. Next, Kevlar inspects the alignment path (represented as a CIGAR string) of each alignment and tests for matches against expected patterns. Alignments matching the pattern

$\text{^}(\backslash\text{d+}[\text{DI}])\text{?}\backslash\text{d+M}(\backslash\text{d+}[\text{DI}])\text{?}\text{\$}$  are classified as SNV events, and the “match” block of the alignment is scanned for mismatches between the contig and the reference target. Any mismatch is reported as a single nucleotide variant. Alignments matching the pattern  $\text{^}(\backslash\text{d+}[\text{DI}])\text{?}\backslash\text{d+M}\backslash\text{d+}[\text{ID}]\backslash\text{d+M}(\backslash\text{d+}[\text{DI}])\text{?}\text{\$}$  are classified as indel events. In addition to reporting the internal gap of this alignment as an indel variant, the flanking “match” blocks are also scanned for mismatches between the contig and the target to be reported as putative SNVs. Any alignment not matching the two patterns described above is designated as an uninterpretable “no-call” and listed in the output along with the corresponding contig sequence.

In some cases, there is a possibility that kevlar will report two or more calls in close proximity. While the probability of two *de novo* variants occurring in close proximity is effectively nil, it is common for an inherited variant to occur proximal to a *de novo* variant. Occasionally one of these inherited variants will not be spanned by any interesting  $k$ -mers, in which case it can immediately be designated as a “passenger” variant call. However, in cases where an inherited variant is spanned by one or more interesting  $k$ -mers, we rely on subsequent examination of  $k$ -mer abundances to distinguish novel variants from inherited variants.

The `kevlar call` command computes the contig alignments and makes preliminary variant calls.

### Step 6: Likelihood scoring model for ranking and filtering variant calls

Given the filters already discussed, false *interesting*  $k$ -mer designations are rare throughout the genome overall. Redundancy from a high depth of sequencing coverage prevents sequencing errors from driving the reported abundance of  $k$ -mers present in the parents to 0. If a  $k$ -mer is present in either parent, it is disqualified from the *interesting* or novel designation.

We observed false *interesting*  $k$ -mer designations are enriched around inherited mutations. It is very common for variants present in one parent to be absent from the other parent. If by chance the depth of sequencing coverage is low at such a locus in the donor parent, there may not be enough redundancy to compensate for sequencing errors. As a result, some  $k$ -mers that are truly present in the donor parent will have a reported abundance of 0. Being truly absent from the other parent, these  $k$ -mers are erroneously designated as unique to the proband.

A related complication occurs when a novel variant is proximal to an inherited variant. Both variants are reflected in the alignment of the associated contig (assembled from proband-derived *interesting* reads) to the reference genome. In both of these cases, distinguishing novel variants from inherited variants benefits from examination of the abundances of all  $k$ -mers containing each variant, as well as the corresponding reference  $k$ -mers.

We utilize a likelihood based model to score and rank the predicted *de novo* variants. We consider the abundance of the interesting  $k$ -mers to calculate the likelihood of the event observed being *de novo*, inherited, or a false call. Using these likelihood probabilities, we calculate a score for each variant being truly a *de novo* variant based on ratio of likelihoods.

First, for each variant we define a set of alternate  $k$ -mers  $\mathbf{A}$  as the  $k$ -mers indicating existence of the variant (alternate genotype). We consider only  $k$ -mers that are unique to this variant (that is, they don’t appear in any other location in the reference genome). We assume that there are a total of  $n$  alternate  $k$ -mers.

Let the random variables  $v_c$ ,  $v_f$ , and  $v_m$  indicate the genotype (i.e.  $\{0/0, 0/1, 1/1\}$ ) of the putative variant in the proband/child, father, and mother respectively. The random variable  $\mathbf{A}_c = \{A_{c_1}, A_{c_2}, \dots, A_{c_n}\}$  denotes the counts of the alternate allele  $k$ -mers in the proband,  $\mathbf{A}_m = \{A_{m_1}, A_{m_2}, \dots, A_{m_n}\}$  the alternate allele  $k$ -mer counts in the mother, and  $\mathbf{A}_f = \{A_{f_1}, A_{f_2}, \dots, A_{f_n}\}$  the alternate allele  $k$ -mer counts in the father. The likelihood that a putative variant is *de novo* can be calculated as follows.

$$\begin{aligned}
L(\text{dn} = 1) &= P(\mathbf{A}_c, \mathbf{A}_m, \mathbf{A}_f \mid \text{dn} = 1) \\
&= P(\mathbf{A}_c, \mathbf{A}_m, \mathbf{A}_f \mid v_c = 0/1, v_m = 0/0, v_f = 0/0) \\
&= P(\mathbf{A}_c \mid v_c = 0/1)P(\mathbf{A}_m \mid v_m = 0/0)P(\mathbf{A}_f \mid v_f = 0/0)
\end{aligned}$$

We note that there are dependencies between  $k$ -mer counts within a sample. However, to simplify the calculation of likelihoods, we assume independence between the  $k$ -mer counts and provide an approximation of the likelihoods. For calculating the probability of an observed  $k$ -mer count conditioned on a 1/1 genotype, we assume a normal distribution where parameters are learned empirically for each sample using only exonic  $k$ -mers that occur only once in the reference genome. For the genotype 0/0 we use binomial distribution to calculate the likelihood of the observed  $k$ -mer abundance assuming the  $k$ -mer is generated by, e.g., sequencing error. Similarly, we calculate the likelihood that a putative variant is a false positive prediction by conditioning on the variant's non-existence (genotype 0/0) in all three samples, i.e.  $L(fp = 1) = P(\mathbf{A}_c, \mathbf{A}_m, \mathbf{A}_f \mid v_c = 0/0, v_m = 0/0, v_f = 0/0)$ .

Finally, we calculate the likelihood of observed  $k$ -mer counts under the inheritance assumption. As there are several different valid scenarios to represent variant inheritance the likelihood calculation requires additional steps as explained below (again assuming independence of  $k$ -mer abundances as an approximation).

$$\begin{aligned}
L(\text{ih} = 1) &= P(\mathbf{A}_c, \mathbf{A}_m, \mathbf{A}_f \mid \text{ih} = 1) \\
&\approx \prod_{i=1}^n P(A_{c_i}, A_{m_i}, A_{f_i} \mid \text{ih} = 1) \\
&= \prod_{i=1}^n \frac{P(\text{ih} = 1 \mid A_{c_i}, A_{m_i}, A_{f_i})P(A_{c_i}, A_{m_i}, A_{f_i})}{P(\text{ih} = 1)} \\
&= \prod_{i=1}^n \frac{P(A_{c_i}, A_{m_i}, A_{f_i})}{P(\text{ih} = 1)} \times P(\text{ih} = 1 \mid A_{c_i}, A_{m_i}, A_{f_i})
\end{aligned}$$

We calculate the  $P(\text{ih} = 1 \mid A_{c_i}, A_{m_i}, A_{f_i})$  as summation of probability of possible trio-genotype combinations representing inheritance scenarios (e.g.,  $(v_c = 1/0, v_f = 1/0, v_m = 0/0)$  or  $(v_c = 1/0, v_f = 0/0, v_m = 1/0)$ ). Furthermore, we assume a constant prior value for  $P(\text{ih} = 1)$  based on all possible valid inheritance scenarios.

Finally, we utilize a heuristic score motivated from the likelihood ratio test to score and rank any predicted variant as being a *de novo* variant. Note that, as numerical calculation of the likelihoods is numerically prone to error we consider the logarithm of the score. Thus, we formally define the score assigned to each variant for being *de novo* as  $S_L = \log L(\text{dn} = 1) - \max\{\log(L(\text{ih} = 1)), \log(L(\text{fp} = 1))\}$ . The `kevlar simlike` command computes likelihoods for preliminary variant calls, sorts the calls, and filters out low scoring and otherwise problematic calls.

### Data simulations

We simulated whole-genome shotgun sequencing for a hypothetical trio (father, mother, and proband) to evaluate the accuracy of our *de novo* variant discovery algorithm. Using the human reference genome (GRCh38) and a catalog of common variants (dbSNP), we constructed two independent diploid genomes representing the two parents. We randomly selected SNPs and indels from dbSNP and assigned the variants to each parental haplotype at a rate of 1 for every 1000 bp.

We then constructed the diploid proband genome through recombination of the parental diploid genomes and simulated germline mutation. SNVs and indels ranging from <10 bp to 400 bp in length were simulated as heterozygous events unique to the proband, representing *de novo* variation.

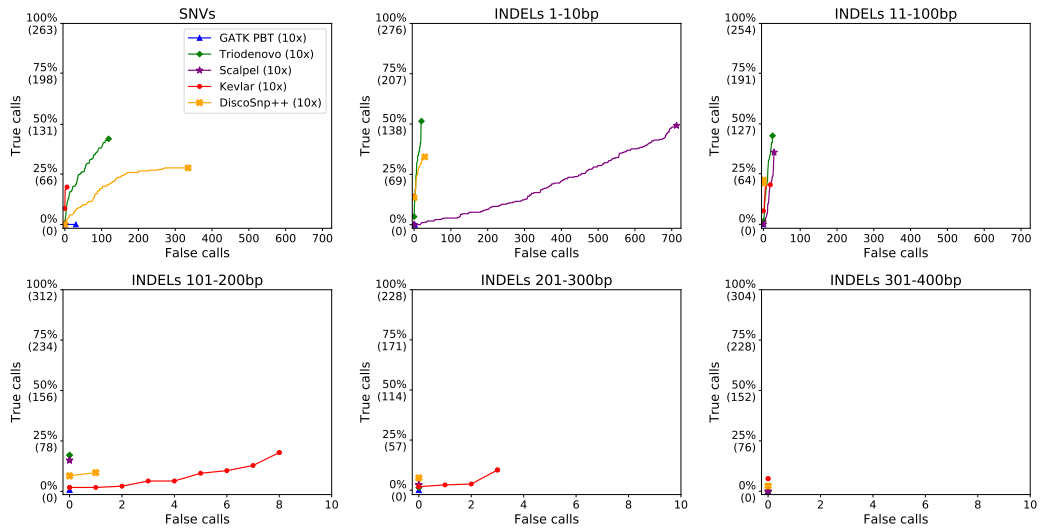
Finally, we used `wgsim` (Li, 2011) to simulate whole-genome shotgun sequencing of each individual. This produced sequences resembling Illumina 2x150bp paired-end reads with low sequencing error rate. The sequencing was repeated at four different average depths of sequencing coverage: 10x, 20x, 30x, and 50x.



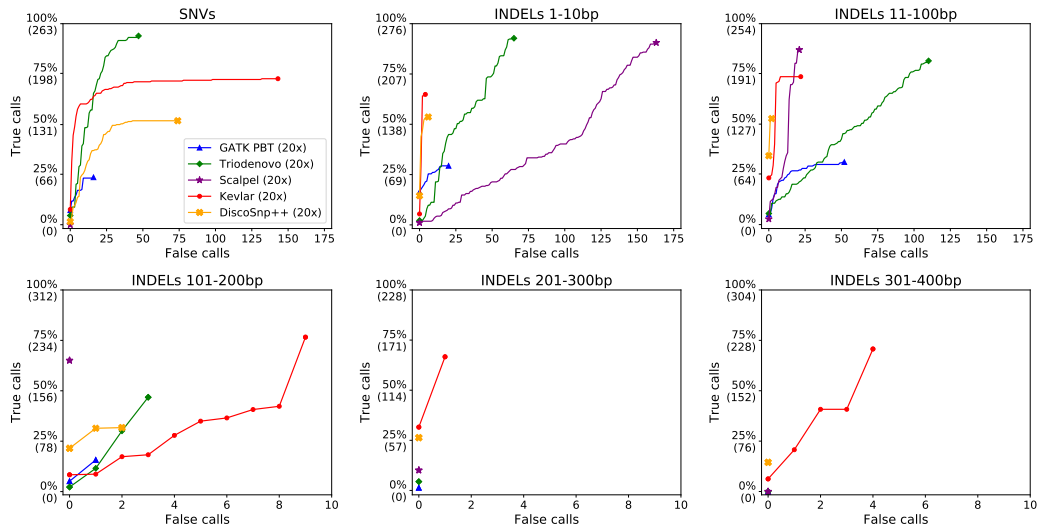
## References

- Bray, N. L., Pimentel, H., Melsted, P., and Pachter, L. (2016). Near-optimal probabilistic RNA-seq quantification. *Nature Biotechnology*, 34(5):525.
- Chong, Z., Ruan, J., Gao, M., Zhou, W., Chen, T., Fan, X., Ding, L., Lee, A. Y., Boutros, P., Chen, J., et al. (2017). novobreak: local assembly for breakpoint detection in cancer genomes. *Nature methods*, 14(1):65.
- Crusoe, M. R., Alameldin, H. F., Awad, S., Boucher, E., Caldwell, A., Cartwright, R., Charbonneau, A., Constantinides, B., Edverson, G., Fay, S., et al. (2015). The khmer software package: enabling efficient nucleotide sequence analysis. *F1000Research*, 4.
- Li, H. (2011). wgsim: Read simulator for next generation sequencing. <https://github.com/lh3/wgsim>.
- Li, H. (2013). Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. arXiv:1303.3997.
- Li, H. (2017a). fermi-lite: Standalone C library for assembling illumina short reads in small regions. <https://github.com/lh3/fermi-lite>.
- Li, H. (2017b). KSW2: Global alignment and alignment extension. <https://github.com/lh3/ksw2>.
- Marçais, G. and Kingsford, C. (2011). A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics*, 27(6):764–770.
- Patro, R., Mount, S. M., and Kingsford, C. (2014). Sailfish enables alignment-free isoform quantification from rna-seq reads using lightweight algorithms. *Nature Biotechnology*, 32(5):462.
- Rahman, A., Hallgrímsdóttir, I., Eisen, M., and Pachter, L. (2018). Association mapping from sequencing reads using k-mers. *eLife*, 7:e32920.
- Rizk, G., Lavenier, D., and Chikhi, R. (2013). DSK: k-mer counting with very low memory usage. *Bioinformatics*, 29(5):652–653.
- Standage, D. S., Aliyari, A., Cohen, L. J., Crusoe, M. R., Head, T., Irber, L., Joslin, S. E., Kingsley, N. B., Murray, K. D., Neches, R., Scott, C., Shean, R., Steinbiss, S., Sydney, C., and Brown, C. T. (2017). khmer release v2.1: software for biological sequence analysis. *The Journal of Open Source Software*, 2(15):272.
- Sun, C. and Medvedev, P. (2018). Toward fast and accurate SNP genotyping from whole genome sequencing data for bedside diagnostics. *bioRxiv*, page 239871.
- Zhang, Q., Pell, J., Canino-Koning, R., Howe, A. C., and Brown, C. T. (2014). These are not the k-mers you are looking for: Efficient online k-mer counting using a probabilistic data structure. *PLoS ONE*, 9(7):1–13.

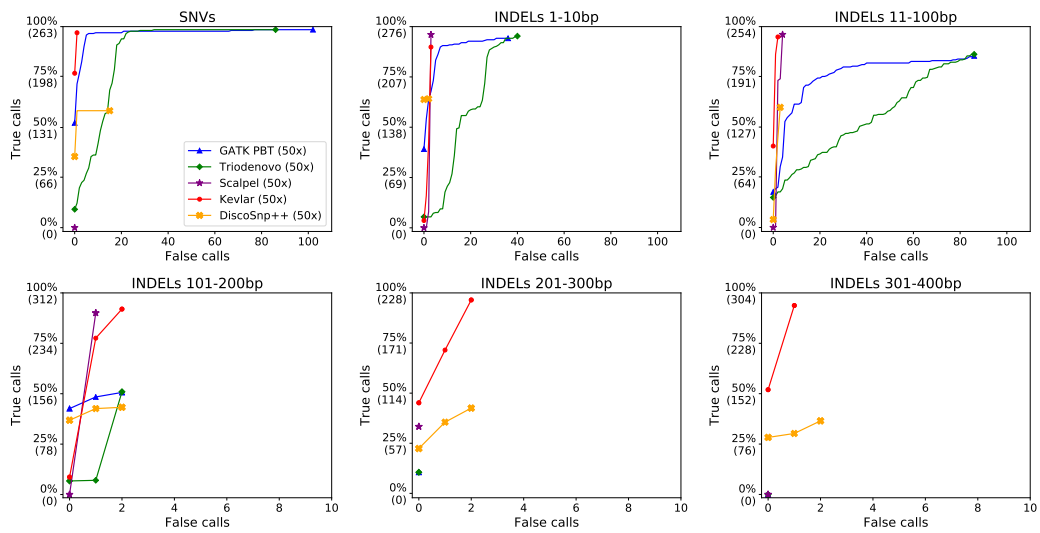
# Supplementary Figures



**Figure S1** Accuracy of five *de novo* variant prediction algorithms at 10x coverage, Related to Figure 2.



**Figure S2** Accuracy of five *de novo* variant prediction algorithms at 20x coverage, Related to Figure 2.



**Figure S3** Accuracy of five *de novo* variant prediction algorithms at 50x coverage, Related to Figure 2.