

METHOD

Open Access

# MINTIE: identifying novel structural and splice variants in transcriptomes using RNA-seq data



Marek Cmero<sup>1,2,3</sup>, Breon Schmidt<sup>1,2,4</sup>, Ian J. Majewski<sup>5,6</sup>, Paul G. Ekert<sup>1,2,7,8</sup>, Alicia Oshlack<sup>1,2,3,4\*</sup> and Nadia M. Davidson<sup>1,2,4\*</sup>

\* Correspondence: [Alicia.Oshlack@petermac.org](mailto:Alicia.Oshlack@petermac.org); [nadia.davidson@petermac.org](mailto:nadia.davidson@petermac.org)

Alicia Oshlack and Nadia M. Davidson are joint corresponding authors and supervised this project equally.

<sup>1</sup>Peter MacCallum Cancer Centre, Melbourne, VIC, Australia

Full list of author information is available at the end of the article

## Abstract

Calling fusion genes from RNA-seq data is well established, but other transcriptional variants are difficult to detect using existing approaches. To identify all types of variants in transcriptomes we developed MINTIE, an integrated pipeline for RNA-seq data. We take a reference-free approach, combining de novo assembly of transcripts with differential expression analysis to identify up-regulated novel variants in a case sample. We compare MINTIE with eight other approaches, detecting > 85% of variants while no other method is able to achieve this. We posit that MINTIE will be able to identify new disease variants across a range of disease types.

## Introduction

Rearrangements of the genome can disrupt or modify gene function and have been implicated as the causal event in disease. In cancer, somatic genomic rearrangements are common and can alter the genomic landscape to drive oncogenesis and cancer progression [1–3]. Whole genome sequencing (WGS) has successfully been used to detect structural variants (SVs) and profile their frequency in cancers [4] and rare diseases [5, 6]. However, clinically relevant variants can occur alongside benign events, making prioritisation of important events difficult. This is especially true in cancers with genomic instability.

Transcriptome profiling has been used to interpret the functional impact of genomic variants through alterations in gene expression, transcript sequence, or both [7]. Some rearrangements that alter gene structure may be more effectively detected from RNA sequencing (RNA-seq). In particular, numerous computational approaches have been developed to reliably call fusion genes [8, 9]. Similarly, several methods now exist to detect novel gene splicing [10–12], and the utilisation of RNA-seq to detect and interpret splicing variants has been shown to improve diagnostic yields in rare Mendelian disorders [13, 14].



© The Author(s). 2021 **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

While there are many methods to detect fusion genes and novel splice variants (NSVs) from RNA-seq, transcribed structural variants (TSVs) involving a single gene, such as deletions, inversions, internal tandem duplications (ITDs) and partial tandem duplications (PTDs) are difficult to detect, and only a handful of tools have been developed [15–18]. However, these variant types can be clinically important, for example, IKZF1 deletions in acute lymphoblastic leukaemia [19], FLT3 internal tandem duplications (ITDs) and MLL partial tandem duplications (PTDs) in acute myeloid leukaemia [15–17, 20]. The prevalence of TSVs in disease is not yet known, but small SVs (< 100kbp) that could give rise to TSVs are frequently seen in WGS [4, 21]. In addition, fusion-finding tools are likely to miss non-canonical fusions, which we define as a fusion transcript that includes non-reference sequence (i.e. outside of the expected transcriptome) at the fusion boundary or a fusion product created between a gene and an intergenic region.

Fusion finding tools typically use discordant or split read detection approaches, where reads are first mapped to either a reference genome [22, 23] or transcriptome [24, 25], then discordant reads are detected and a set of filtering steps are performed. While mapping information from genome alignments can be used to detect non-canonical fusions and other TSVs, most fusion finding tools do not consider these variants. However, some exceptions exist; Arriba [26] is a method that utilises chimeric reads from STAR's alignments to detect fusions, including non-canonical fusions, as well as PTDs. Similarly, CICERO [18] utilises chimeric read information and combines this with local assembly to detect canonical and non-canonical fusions, as well as ITDs. SQUID [16], which is a method based on the idea of rearranging the reference genome using concordant and discordant reads, can detect non-canonical fusions and some TSVs. This approach has also been extended to handle heterogeneous (i.e. multi-allelic) contexts [27]. These approaches rely heavily on the reference genome and generally have limitations in detecting small variants (< 100 bp) such as INDELs and ITDs. These methods are biased towards detecting certain classes of variants.

An alternative to detecting variants from read alignment is to use de novo assembly. This approach attempts to reconstruct a sample's transcriptome, and then aligns the assembly back to the reference genome to identify variants. Assembly can reconstruct transcripts containing variants of all types and sizes, which avoids potential biases caused by alignment to the reference genome. KisSplice [11] is an example of a reference free method focused on splice variants. It calls splice variants based on bubble-detection within the assembly De Bruijn graph. De novo assembly approaches typically generate many transcripts, many of which will be misassemblies, leading to a potentially large false positive rate. Additionally, these methods are generally slower to run due to the computational requirements of de novo assembly. One strategy to reduce false positives and improve speed is to only consider a specific list of target genes, a strategy employed by TAP [15]. Another approach, used by DE-kupl [28], is to perform differential expression of k-mers against a set of control samples. DE-kupl performs k-mer counting, filtering, DE and then extends the resulting DE k-mers into contigs. However, DE-kupl reports a very large number of variants (in the order of 100,000) [28], which makes interpretation and prioritisation difficult, particularly as the tool does not provide variant type information. Another limitation of the DE-kupl approach is that it

compares two groups, each with multiple samples, and is therefore not designed for detecting variants in a single sample.

Reference independent assembly combined with differential expression against controls, as introduced by DE-kupl, is a powerful strategy for the unbiased detection of transcribed variants. Here we propose an alternative approach, called MINTIE, which can detect any kind of anomalous insertion/deletion ( $\geq 7$  bp by default) or splicing (flanked by  $\geq 20$  bp by default) in any gene. MINTIE is an RNA-seq analysis pipeline that combines the advantages of full de novo assembly with differential expression to identify unique variants in a case sample versus a set of controls. MINTIE's pipeline includes further steps to filter, annotate and prioritise variants, which reduces transcriptional noise and aids in variant discovery.

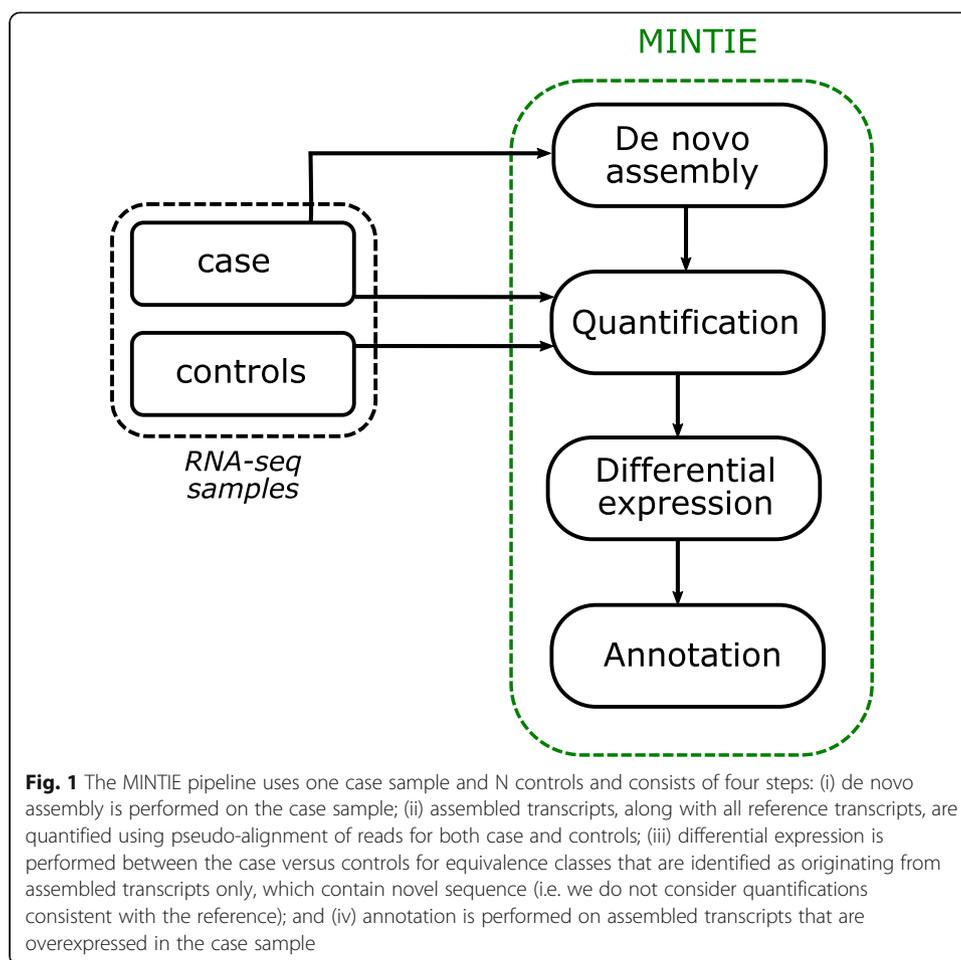
We tested MINTIE on a simulation of 1500 variants including fusions, TSVs and NSVs and showed that it is highly sensitive across many different variant types. We compared the performance to eight other methods and found that MINTIE was the only one able to consistently detect all classes of variants we simulated at high recall rates ( $> 85\%$  of total variant transcripts). In addition, we ran MINTIE on RNA-seq samples from acute myeloid leukaemia (AML) and normal blood, where we confirmed its sensitivity in detecting FLT3-ITDs, KMT2A-PTDs and fusions on real data. Normal (non-cancer) samples were used to assess the background rate of variants. In these samples, MINTIE reported a median of 122 genes per sample containing transcriptional variants, demonstrating a low background rate.

We used MINTIE to discover several interesting alterations in known driver or tumour suppressor genes in a cohort of paediatric acute lymphoblastic leukaemia (ALLs) from the Royal Children's Hospital, Melbourne. We found a recurrent non-canonical fusion involving RB1 with downstream intergenic regions of the genome. We also identified novel splicing of ETV6, and PTDs of IKZF1 and PAX5. Finally, we demonstrate that MINTIE is able to detect novel splice variants identified in a prior study in a rare disease cohort. In addition, we discovered a previously missed complex rearrangement in the disease-causing DMD gene in a patient with muscular dystrophy from this cohort.

## Results

### Algorithm overview

MINTIE is an approach for detecting novel transcribed variants that utilises de novo assembled transcripts, which are then prioritised for novelty using differential expression. The MINTIE pipeline (Fig. 1) is divided into four main steps: transcriptome assembly of the case sample, pseudo-alignment of the reads from the case and controls to an index composed of the assembled and reference transcripts, differential expression to identify upregulated novel features, and annotation of novel transcripts. One case sample is required as input. The pipeline is always run in a 'single case versus N controls' fashion (although multiple cases versus the same control set can be run in parallel, each case is run in a 1 vs. N controls fashion). Ideally, controls should be samples of the same tissue type but do not have to be normals or matched normals. While we recommend running MINTIE with controls, in some cases, this may not be feasible,



and thus, the method can also be run without controls, with the caveat that more background variants are likely to be identified.

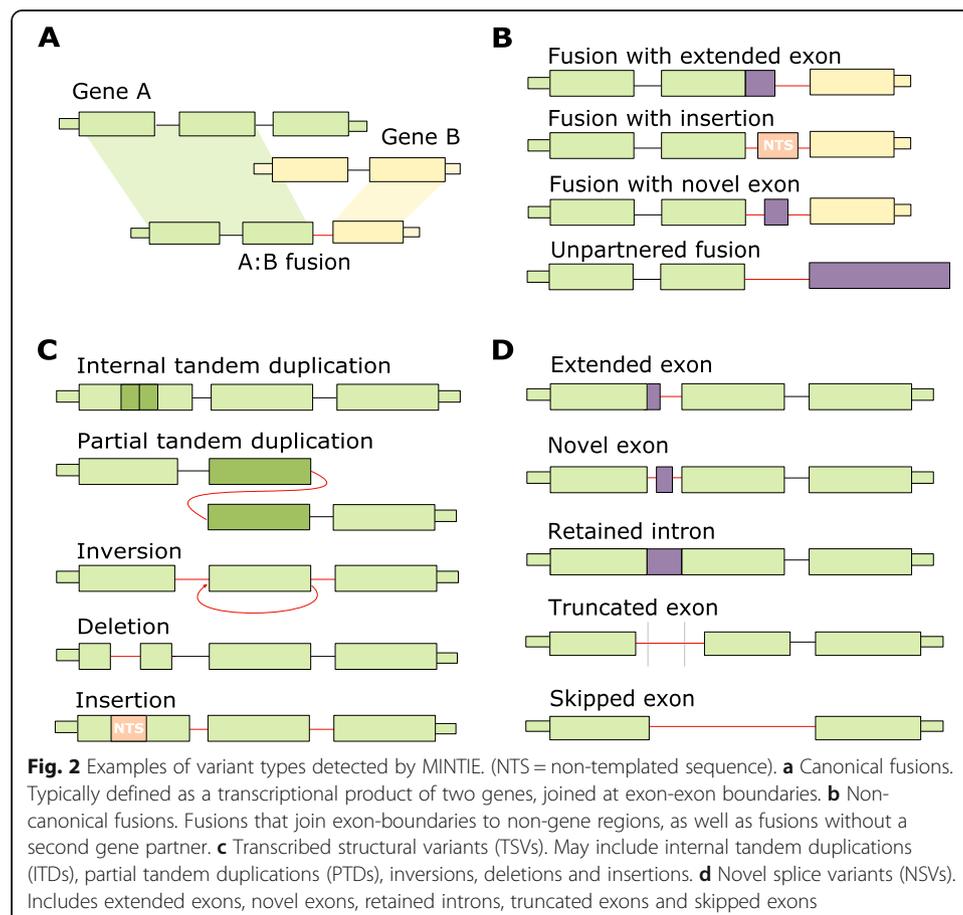
Conceptually, the idea of MINTIE is to perform whole-transcriptome de novo assembly, remove features consistent with reference transcripts, and then perform differential expression on non-reference features versus a set of controls. Significantly over-expressed assembled transcripts are aligned to the genome and variants are identified. The novelty of this approach is that we use differential expression testing to select assembled transcripts that contain highly expressed novel sequence. Therefore, we avoid using alignment to a reference genome to define novel sequence. Instead, we identify novel sequence by testing gene expression on equivalence classes that are unique to the assembly.

In more detail, once a de novo assembly of a case sample is obtained, by default using SOAPdenovo-Trans [29], the assembled transcripts are merged with a standard transcriptome reference (we use CHES v2.2 [28] by default), and a Salmon [30] index is created for this file. Salmon is then run on the case and N controls. Salmon equivalence class (EC) counts are matched across all samples. ECs represent the set of transcripts that a given read is equally compatible with, and have been used previously as a basis for differential expression [31], differential isoform usage [32], as well as fusion detection [33]. MINTIE retains EC counts that are compatible with de novo assembled

transcripts only (no reference transcripts), and all other EC counts are discarded. This step identifies novel ECs where there is no compatible match to an existing reference transcript, as ECs containing reference transcripts are assumed to be consistent with the reference, and therefore uninteresting. The expression of these novel transcript ECs are then compared by performing differential expression testing between 1 case and N controls with edgeR [34]. This step is required to remove common but unannotated transcripts, as well as to remove erroneous assemblies where there is no discernable expression. Significant ECs (false discovery rate, FDR < 0.05 and log<sub>2</sub> fold change, log<sub>2</sub>FC > 2 by default) are then retained, after which the transcripts corresponding to these ECs are extracted and aligned to the genome. MINTIE then performs annotation, filtering, and estimates variant allele frequency for each given variant. See the “Methods” section for more detail.

### MINTIE detects more types of variants than other tools

In order to test the utility of MINTIE across a wide range of variant types, we simulated 1500 variants by inserting, deleting, combining or duplicating sequences from 100 randomly selected hg38 UCSC RefSeq [30] transcripts. We then generated reads from this modified reference using ART-Illumina [31] v2.5.8. Fifteen different classes of variant were simulated (Fig. 2): 500 fusions (100 each of canonical fusions, unpartnered

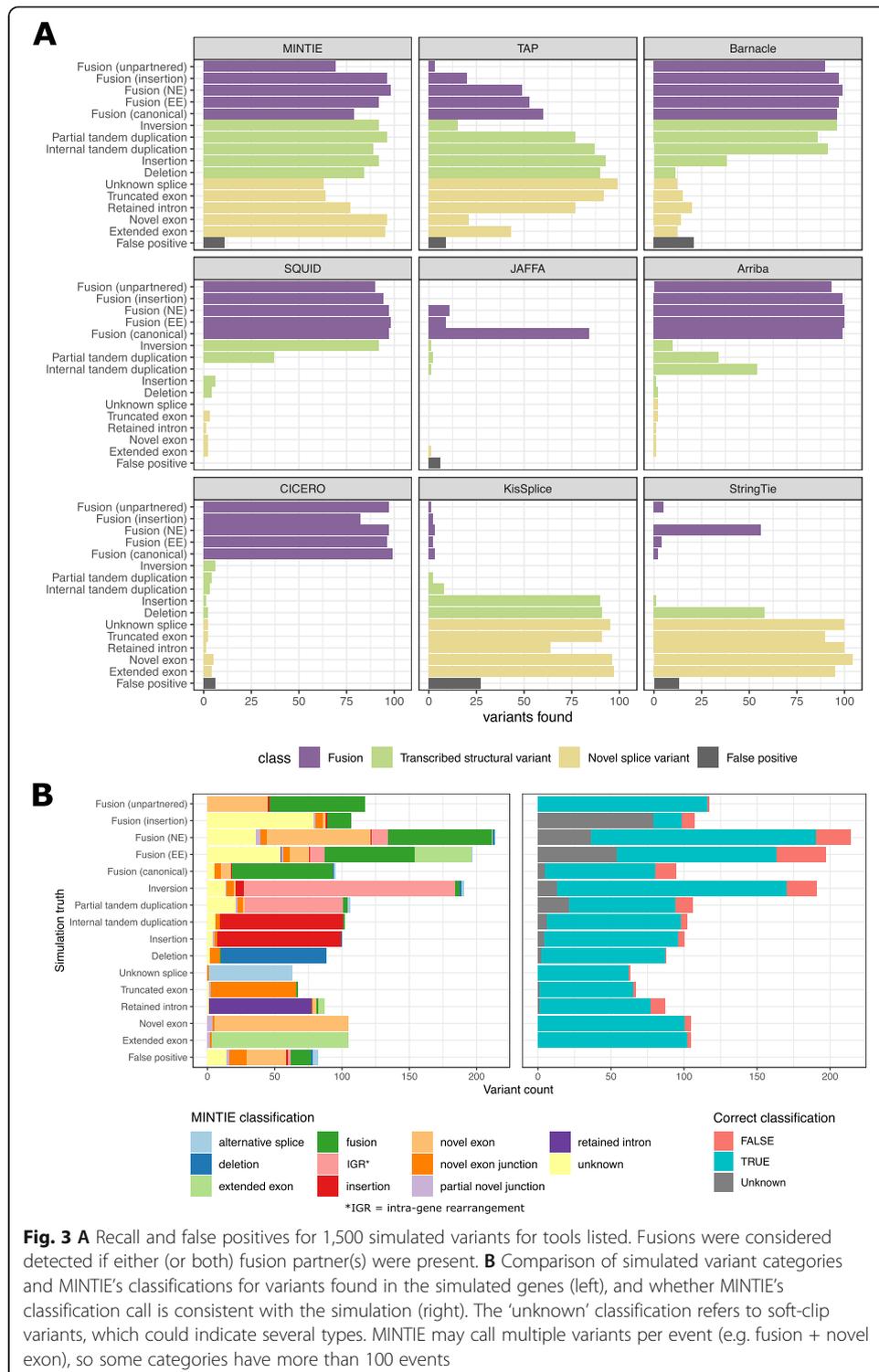


fusions, and fusions with extended exons, novel exons and insertions at their fusion boundaries), 500 transcribed structural variants (TSVs) (100 each of deletions, insertions, ITDs, PTDs and inversions), 500 novel splice variants (NSVs) (100 each of extended, novel and truncated exons, retained introns and unannotated splice variants) and 100 unmodified background genes (see the “Methods” section). Variants were simulated at 50x coverage in a heterozygous fashion (i.e. both modified and unmodified transcripts were simulated for the variant sample, resulting in 100x total gene coverage: 50x variant coverage and 50x wild-type coverage). We also generated a control sample by simulating reads over wild-type (unmodified) transcripts at 100x for all transcripts used in the variant sample.

We ran MINTIE and eight other variant detection methods (TAP [15], Barnacle [17], SQUID [16], JAFFA [24], Arriba [26], CICERO [18], KissSplice [11] and StringTie [12] on the reads generated from the simulated variants. Although DE-kupl is conceptually similar to MINTIE, we were not able to test its performance because it cannot be run on a single case sample. Due to the large number of tools, each with their own output formats that had to be individually processed, we utilised a liberal approach to counting variant calls as true positives (calling a variant in a gene at either end of a fusion was counted as a hit and variant classifications were not considered; this is detailed in the “Methods” section). MINTIE was run with control samples generated from the same transcripts without the variants (see the “Methods” section).

MINTIE detected 86.8% of fusions, 90.6% of TSVs and 79% of NSVs (Fig. 3A). MINTIE had its lowest recall on unknown splice variants (skipped exons) (63%) and truncated exons (64%). Manual inspection revealed that missed cases were mostly due to the variant not being assembled or insufficient read coverage. Assembly issues could potentially be improved with another assembler. Insufficient read coverage particularly affected some deletion variants that were near the end or start of the variant transcript and thus less likely to have adequate read support. To assess the impact of read coverage, we downsampled the simulated reads to three coverage values: 20x, 10x and 5x. Only a modest drop in performance was observed (Additional file 1: Fig. S1A). Even at 5x, MINTIE found most simulated variants (769 of 1500). We also found that MINTIE retained high detection rates across all sizes of variants, even below the default insertion/deletion threshold size of  $\geq 7$  bp and down to 1 bp, apart from single base-pair deletions (Additional file 1: Fig. S2).

While MINTIE had lower sensitivity when compared to specialised tools for novel splice variant detection (StringTie and KissSplice), our method was able to find the most total variants (86.2%), while no other tool came close to this. CICERO has been reported to have high sensitivity for ITD detection [18]; however, it only identified 3 ITDs in our simulated data set. CICERO ranks ITDs by whether the gene is known for recurrent ITDs. We did not simulate ITDs in such genes and hypothesise that this is the reason they were missed. These ITDs were reported (93%) in CICERO’s unfiltered output (Additional file 1: Fig. S3); however, the unfiltered output also included many false positives (596 in total). Given that only 100 unmodified (wild-type) transcripts were simulated, we expected minimal false positives to be identified. No false positive calls were identified for SQUID and Arriba. JAFFA reported six false positive fusion genes (across five calls, all at low confidence). CICERO, TAP, MINTIE, StringTie and Barnacle called 6, 9, 11, 13 and 21 false positives



respectively, while KisSplice had the highest number at 27. StringTie, TAP and Barnacle reported calls in background genes (those simulated without any variants) with 1, 2 and 3 hits respectively. MINTIE, KisSplice and JAFFA did not report any false positive calls in background genes. All other reported false positives were outside of simulated gene regions. Manual inspection of MINTIE's false positives

revealed that these calls were largely due to poor alignment to homologous sequences (for example, to a pseudogene of the variant gene).

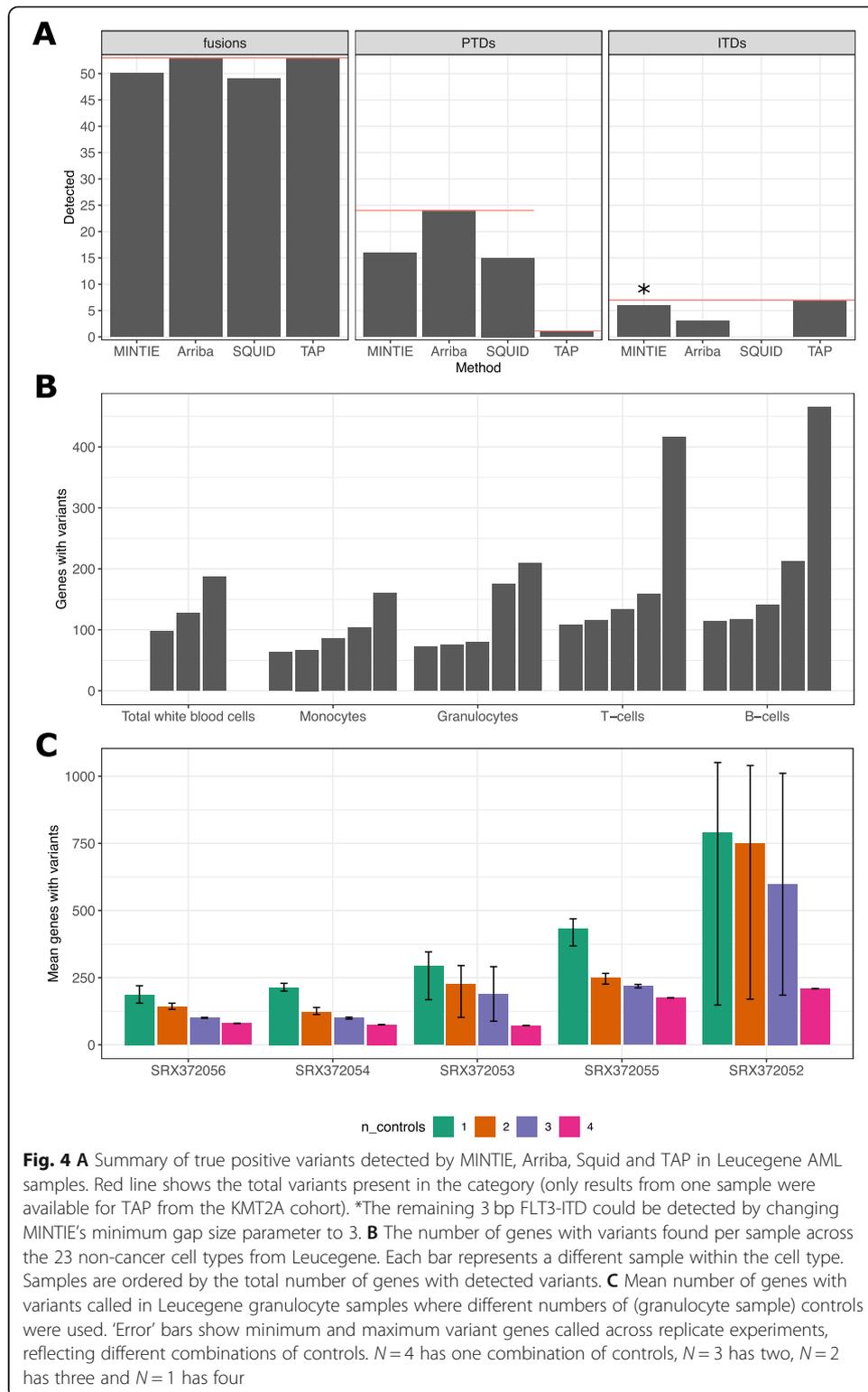
To aid with variant interpretation and prioritisation, MINTIE outputs classifications for the type of transcribed variant found: fusion, intra-genic rearrangement, deletion, insertion, novel/extended exon, novel exon junction and retained intron. We explored the classification accuracy of MINTIE's results using the simulation. Figure 3B shows that the variant types MINTIE assigned to variants called within the simulated genes are mostly correct (86.5%) across the 15 categories.

### **MINTIE identifies known fusions and ITDs in Leucegene samples**

In order to demonstrate the sensitivity of MINTIE in real samples, we applied MINTIE to a set of 77 AML samples with known fusions, internal tandem duplications and partial tandem duplications (Additional file 1: Note S1). This included the NUP98-NSD1 fusions that were validated from Lavallée et al. [35] ( $N = 7$ ). We also tested MINTIE on samples containing CBF-MYH11 fusions, RUNX1-RUNX1T1 fusions and FLT3-ITDs known to be in the core binding factor (CBF) AML data [27], as well as a cohort containing KMT2A-PTDs [32], identified by Audemard et al. [33]. Although we benchmarked just five variants across three variant types, which represents a small subset of the types of events MINTIE can detect, it enabled MINTIE's sensitivity to be confirmed using well known and validated events. The Leucegene data was also used to explore the impact of control choice, read coverage in real data and rate of background events.

In all these analyses, we used 13 normal Leucegene samples as controls (5 granulocytes, 5 monocytes and 3 total white blood cells). In the AML data sets, a median of 592 variant genes were reported per sample (range 261–2265). MINTIE detected 50/53 total fusions (Fig. 4A); manual inspection revealed that two CBF-MYH11 fusions were missed due to low expression and one NUP98-NSD1 fusion was not assembled. Additionally, all FLT3-ITDs were detected in the CBF and NUP98-NSD1 cohorts (default parameters needed to be adjusted for one patient sample in order to detect a small 3 bp FLT3-ITD). MINTIE was able to find 16/24 KMT2A-PTDs in the samples identified by Audemard et al. [33] when using normal samples as controls and 19/24 when using no controls (Additional file 1: Table S1). Some of the variants that were successfully detected (using normals as controls) had estimated expressed variant allele frequencies (VAFs) down to 0.0580 (median 0.2220, max 0.7127), suggesting that only a small amount of the total gene expression needs to come from the variant isoform for detection. The variants that were missed had low coverage ( $< 17$  reads). Of the 8 missed variants, 6 were due to insufficient assembly and 2 were filtered out by the counts per million (CPM) filter due to low counts in their respective ECs.

We also analysed these Leucegene samples with three selected state-of-the-art methods, which are designed for these variant types: Arriba (a representative fusion finder), SQUID (a representative large TSV caller) and TAP (a general-purpose caller). Arriba found all fusions and KMT2A-PTDs but missed more than half the FLT3-ITDs (Fig. 4A). Squid performed similarly to MINTIE, finding 49/53 fusions and 15/24 KMT2A-PTDs, but failed to identify any of the FLT3-ITDs. This highlights how current tools designed to detect larger rearrangements such as fusions and PTD may be insensitive to smaller events such as ITDs. TAP was able to detect all events



(Fig. 4A). However, due to difficulties in installing and running TAP, detection rates were taken from their publication [15], and results from only one KMT2A-PTD sample were available. These results show that MINTIE has sensitivity within the range of existing state-of-the-art tools designed specifically for detecting these known, well

defined variants. The validation cohort incorporates only three of the 15 variant types from our simulation study and the PTDs and ITDs were focused on a single gene in each case (KMT2A and FLT3 respectively). Thus, they represent only a fraction of MINTIE's utility. Nevertheless, these results are consistent with our simulation and suggest that MINTIE is likely to have good sensitivity over a range of events.

#### **MINTIE detects a low number of background variants**

We applied MINTIE to a set of 23 non-cancer healthy adult sorted blood cell samples obtained from Leucegene [34]. We expected these samples to have low numbers of transcriptional variants compared to cancer samples and thus provide an estimate to the background rate of detected variants found by MINTIE. This set was composed of several cell types: T cells ( $N = 5$ ), B cells ( $N = 5$ ), granulocytes ( $N = 5$ ), monocytes ( $N = 5$ ) and total white blood cells ( $N = 3$ ). Each sample was run with MINTIE using all other samples of the same cell type as controls.

Figure 4B shows the total number of genes containing transcriptional variants identified across the non-cancer samples. In general, MINTIE had a low background variant rate, with a median of 122 affected genes across samples (range 61–1397). We identified 19,893 expressed genes across the non-cancer samples ( $> 1$  CPM in at least one sample); thus, the range translates to a median of 0.6% of expressed genes identified as containing variants. 83.3% of the 2363 identified variant genes were protein coding (compared to 74.39% of expressed genes that were protein coding). 56.2% of variants were classified as novel splicing variants and may represent true variation in the samples. The number of variants called was variable between samples, even within samples of the same cell type. We found only a weak correlation between library size and number of transcriptional variants detected (spearman correlation coefficient = 0.2105). A moderate fraction of variants (23.3%) were of unknown type and may be due to poor alignment of the de novo assembled transcripts to the genome. We note that one B cell sample contained the largest number of variant calls (1710 variants). Upon manual inspection, we noted that 1259 of these variants were clustered in a 1 MB region at the end of the long arm of chromosome 14 containing immunoglobulin genes, likely indicating that these genes, which are commonly rearranged, were highly expressed in this sample.

#### **Well-matched controls reduce the number of background variants**

We found that using correctly matched controls was important for filtering out the normal range of transcriptional diversity in each sample. We compared the number of variant genes observed in normal total white blood cells (TWBCs) when using different controls: the same cell type (TWBCs), as well as granulocytes, monocytes, T cells and B cells. Using the other TWBCs as controls resulted in the fewest number of variant genes found across all three samples ( $N = 412$ ), followed by monocytes ( $N = 715$ ) and granulocytes ( $N = 763$ ) (Additional file 1: Fig. S4). This corresponds with the expression similarity (Additional file 1: Fig. S5). B cells and T cells cluster together in gene expression, distinct from TWBCs, and using them as controls identified the most variant genes at 1564 and 1630 respectively. These results confirm that comparing against

controls of similar expression profiles combined with enough control samples resulted in fewer total number of variants found.

To further explore the effect of control number and similarity on MINTIE's variant calls, we performed an experiment using 1–4 controls in different combinations using the Leucegene granulocyte samples. We observed fewer variant genes called as the number of controls were increased (Fig. 4C). One sample (SRX372052) displayed notably higher variability in variant genes called, compared to other samples, when using 1–3 controls. On inspection of expression similarity through a PCA plot (not shown), the sample was revealed to cluster more closely with two samples than the others. When the less similar samples were randomly chosen as controls, variant numbers were significantly higher. These results suggest that using more controls reduces the number of background variants called and confirms that sample similarity will impact the number of variants called.

We also ran the cohort of 24 Leucegene AML samples containing KMT2A-PTDs against three different sets of controls: (i) 13 normal (non-cancer) controls (5 granulocytes, 5 monocytes and 3 total white blood cells), (ii) a reduced set of 3 normal controls containing one sample from each cell type, (iii) a set of 13 AML samples from a different cohort and (iv) no controls. MINTIE was able to detect 19/24 variants using no controls, 16/24 variants using the normal controls, 15/24 variants using the reduced normal set and 11/24 variants using AML controls (Additional file 1: Table S1). The lower detection rate versus other cancers was due to high read counts in the controls for the ECs corresponding to the KMT2A variant. While there was no evidence of KMT2A rearrangements in the control samples used, some of the assembled contigs contained other, more common variants in the same EC (such as an extra A insertion in a small A repeat region proximal to the breakpoint), and thus were not found to be significant at the DE step. Differences were observed in total variants found using different control sets (Additional file 1: Fig. S6), with a median of 645.5 variant genes (range 264–2265) found per sample for the 13 normal controls compared to 913 median variant genes (range 490–2552) for the reduced set of 3 non-cancer control samples and a median of 211.5 variant genes (range 130–2852) found when using the 13 AML controls. Using no controls yielded significantly higher variant numbers (range 4750–11796, median 9255.5), indicating that using even small numbers of controls acts as an effective filtering strategy. If we assume most variants are background events, unrelated to cancer, this suggests that using samples with well-matched transcriptomes as controls reduces the number of detected background variants.

### **MINTIE identifies novel transcriptomic variants in B-ALL**

In order to test MINTIE's ability to discover previously undetected events, we applied MINTIE to a set of 91 samples across a cohort of 87 paediatric B-ALL patients from the Royal Children's Hospital, Melbourne Australia [36]. As the control group, we used a subset of 34 samples from the cohort with already known driver fusions, identified by molecular testing, and validated in the RNA-seq with fusion calling (Additional file 1: Note S2). A median of 48 variant genes was found per sample (range 15–633, Additional file 1: Fig. S7). We filtered on variants found in 379 recurrently altered genes (not considering immunoglobulin genes) in over 2500 paediatric cancer samples

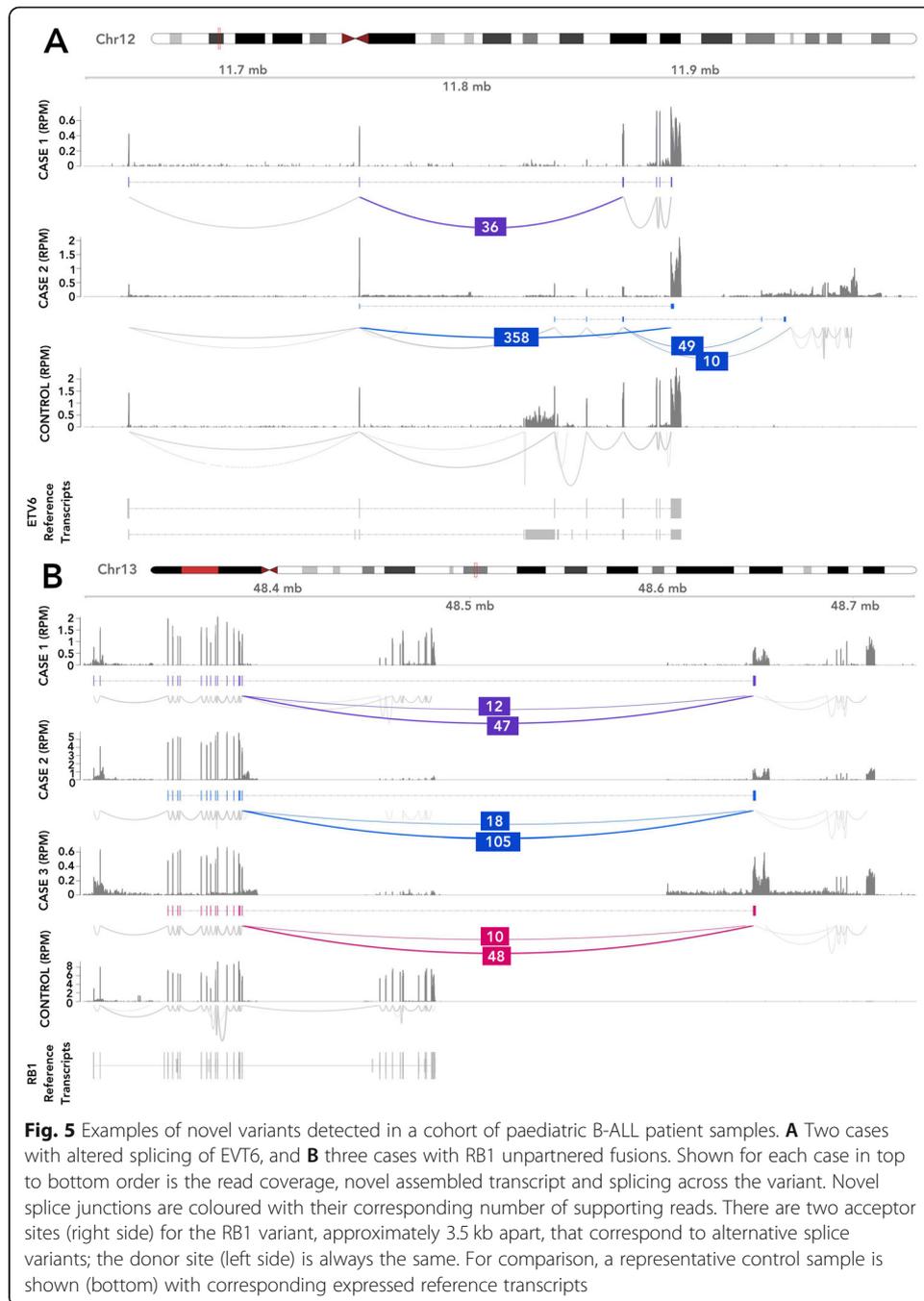
reported in two landscape papers [37, 38] and found 339 variants across 131 of these genes. The top recurrent genes were ETV6 (26 variants across 8 samples), IKZF1 (14 variants across 7 samples) and IKZF2 (13 variants across 3 samples) (Additional file 1: Fig. S8).

One sample contained an in-frame, single-exon tandem duplication in IKZF1. In-frame deletions of IKZF1 are recurrent in paediatric ALL and associated with poor prognosis [19]. They are believed to act in a dominant negative manner and we hypothesise that this might also be the case for the IKZF1-PTD we detected. We also found one sample with an in-frame tandem duplication of 4 exons in PAX5, which correlates with previously reported amplification of PAX5 exons in paediatric ALL [39]. A subtype of B-ALL is defined by PAX5 alterations, with a range of diverse variant types across patients (fusions, amplifications and mutations) [40].

In addition, we found three samples with exon skipping and cryptic exons of ETV6 (Fig. 5A). In one sample, case 1, we observed a unique alternative splicing event between exons 2–5, resulting in lowered expression of intervening exons, indicating a potential deletion. Exon skipping due to an intragenic deletion in ETV6 has previously been reported in paediatric ALL [41]. Karyotyping of the sample indicated a complex rearrangement of ETV6, but no fusion involving ETV6 was found. The same exon skipping event was observed in a single TCGA Lung Squamous Cell Carcinoma sample, found with SeqOthello [42]. In the other two B-ALL samples, case 2, which were from diagnosis and relapse of the same patient, we observed two novel ETV6 isoforms: splicing between exons 2–8 and splicing between exon 5 and two downstream novel intergenic exons. Expression of exons 6 and 7 was absent.

A B-ALL subtype classifier based on gene expression (<https://github.com/Oshlack/ALLSorts> v0.1) was run on the four diagnosis samples harbouring ETV6 splicing, IKZF1-PTD and PAX5-PTD variants. All samples were found to have a strong probability of belonging to their respective class, ETV6-RUNX1-like, IKZF1 N159Y and PAX5alt (data not shown), suggesting these events may be drivers.

MINTIE also detected an unpartnered recurrent fusion involving the tumour suppressor gene RB1 from the end of exon 17 to an intergenic region approximately 165 kb downstream (Fig. 5B) in three other samples. In all three samples, two splice variants of the fusion were seen, which involved different novel downstream exons (Additional file 1: Table S2). Read coverage dropout suggested a 188 kb genomic deletion involving the 3' end of RB1 in one sample; however, this was not clear in the two other samples. RB1 focal deletions have been reported in both B-ALL and T-ALL [38] from DNA-based assays, supporting this hypothesis. We further expanded our search for this variant by querying TCGA using SeqOthello [42]. One of the fusion splice variants involving RB1's exon 17 and the same downstream intergenic region was detected in two other samples, a breast invasive carcinoma and an oesophageal carcinoma. MINTIE successfully detected the variant in both samples when we obtained the RNA-seq data and processed it through the pipeline, using 3 normal controls from TGA. These results suggest that focal RB1 deletions occur in multiple cancer types, and form a stable transcript that can be detected from RNA-seq alone using MINTIE.



### MINTIE identifies novel splicing and transcribed structural variants in rare disease

To demonstrate the utility of MINTIE in rare disease, we analysed RNA sequencing data from patients with rare muscle disorders from Cummings et al. [14]. In the study, diagnosis was made by identifying altered transcripts using RNA-seq in pathogenic genes with an identified variant of unknown significance. From this data set, there were 10 patients that had RNA-seq available with 13 splicing variants that could potentially be detected by our method. We used 10 muscle RNA-seq samples from GTEx [43] as controls and ran MINTIE transcriptome wide using RNA-seq alone without knowledge of which genes had DNA mutations. MINTIE identified the correct novel splice variants for 9/13 events (Table 1).

**Table 1** List of variants detectable by MINTIE from Cummings et al. [14] rare muscle disease data

Patient	Gene	Junction 1	Junction 2	DNA variant	RNA variant	Detected variant
E2	NEB	chr2:151687733-151690727		SNP	Skipped exon	Y
	NEB	chr2:151687733-151688175		SNP	Extended exon	Y
	NEB	chr2:151687733-151688007		SNP	Extended exon	Y
C9	NEB	chr2:151531896-151533330		SNP	Extended exon	Y
N25	NEB	chr2:151498352-151498491		SNP	Novel exon	Y
N31	COL6A1	chr21:45989778-45989894	chr21:45989965-45990258	SNP	Novel exon	Y
N32	COL6A1	chr21:45989778-45989894	chr21:45989965-45990258	SNP	Novel exon	Y
C11	RYR1	chr19:38467720-38468966		SNP	Truncated exon	Y
C1	POMGNT1	chr1:46194651-46195811		SNP	Skipped exon	N (low CPM)
	POMGNT1	chr1:46189357-46189458		SNP	Retained intron	Y
E4	TTN	chr2:178580609-178581906	chr2:151498552-151499298	INDEL	Skipped exon	N (not assembled)
N22	TTN	chr2:178777317-178777460		SNP	Truncated exon	N (low CPM)
C3	DMD	chrX:31729748-31819975		Inversion-deletion	Skipped exon	N (EC not significant)

Two variants were filtered out due to low expression (exon skipping in POMGNT1 and exon truncation in TTN), while an exon skipping event in TTN which was only supported by 2 junction reads, failed to be assembled. A transcript harbouring the DMD skipped exon in patient C3 was not found to be significant in the DE step, as a sizable number of reads supporting the variant transcript were found in the controls. Manual inspection revealed that the assembled transcript for this DMD case contained two variants, one being the target variant (the skipped exon), and a secondary variant (a homozygous SNP). The secondary variant was not in the reference transcriptome but was found in 9/10 controls. Cases such as this are rare but represent one potential weakness of an EC-based method, whereby multiple variants assembled into the same transcript cannot be resolved at the EC-level.

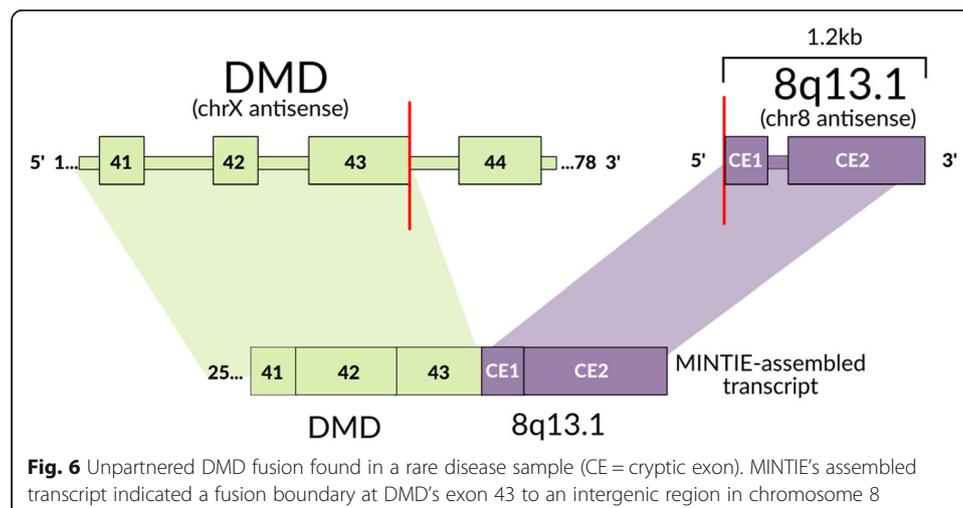
Patient samples C2 and C4 were identified with large-scale inversions in the DMD gene using DNA sequencing that manifested as lower exonic coverage across regions of the gene. Cummings et al. did not report an exon skipping RNA-seq feature, and thus, we could not check specific boundaries in the RNA-seq. However, MINTIE did find novel transcriptomic variants in DMD in both samples. Running BLAT [44] on these soft-clip sequences revealed that they aligned on the opposite strand to the gene sequence, supporting an inversion variant. This is an example where MINTIE can provide insight into DNA-level structural variants solely from the RNA-seq.

We also investigated whether MINTIE could detect any variants of interest that were missed by the Cummings et al. study. Prioritising the same set of NMD

candidate genes as the study's authors, and stringent filtering on MINTIE's outputs, together with manual validation (see the "Methods" section), we identified five variants: four deletions and a fusion across eight patients (Additional file 1: Table S3). All deletions were found to be rare UTR variants found in the general population (< 5% MAF), of which, two were reported in ClinVar (DMD's 3' UTR deletion and VAPB's 3' UTR deletion) as variants of unknown clinical significance. The VAPB mutation was reported to be associated with Spinal Muscular Atrophy (ClinVar Submission Accession SCV000434613). Most interestingly, we identified a fusion transcript in DMD in a patient where exome sequencing suggested that DMD was a strong candidate gene [14]; however, no corresponding RNA variant was found in the study. The unpartnered fusion joined the end of exon 43 of DMD on the X chromosome to an intergenic region on chromosome 8 with two cryptic exons (Fig. 6). Upon manual inspection, we noted 10 split-reads with three different soft-clip sequences at the exon 43 boundary; of these, 5 sequences corresponded to the breakpoint on chromosome 8, matching the assembly as expected. However, two other sequences were found to correspond to regions 650 bp upstream (2 split-reads) and 70 kb downstream (3 split-reads), suggesting multiple potential fusion transcripts. We confirmed the presence of this rearrangement with the original authors through private correspondence, who independently discovered it using whole genome sequencing and recently reported it in Waddell et al. [45].

### Computational Performance

Computer resources are often seen as a limitation to the application of de novo assembly in variant detection. In MINTIE, we used SOAPdenovo-Trans as the default assembler due to its computational efficiency compared to other assemblers. MINTIE also leverages pseudo-alignment to minimise compute times. To benchmark performance, we selected three samples (containing the minimum, median and maximum variants called) from the Leucegene KMT2A-PTD cohort and ran MINTIE against the set of 13 controls previously described (Table 2). Tests were run on a server with 32 cores and 190 GB of available memory with MINTIE's bpipe



**Table 2** Compute times for three selected samples from the Leucegene KMT2A-PTD cohort, including the sample with the least variants called, the median, and the most. Reads were 100 bp and samples were run against 13 AML controls on a server with 32 cores and 190 GB of available memory

Sample	Variants	Read pairs	Time	Max mem (GB)
08H012	128 (min)	67,789,336	3 h 24 min	101.47
07H155	214 (median)	98,946,255	6 h 42 min	89.85
10H007	2093 (max)	113,017,564	7 h 04 min	102.64

pipeline concurrency limited to 32. The variant numbers also corresponded with the number of read pairs, ranging from 67 million to 113 million. Compute times were between 3 h and 24 min and 7 h and 4 min, where the run time scaled with the number of reads. Maximum memory usage did not appear to be significantly correlated with the read number, as all samples were between roughly 90–100 GB max memory usage (we note that some samples may utilise > 190 GB memory, primarily due to the assembly step). The first three steps, dedupe, trim and assemble steps account for approximately over half the run time and most of the memory usage. For instance, in one sample, these stages took 4 h and 4 min and used, at maximum, 102.64GB of memory. The rest of the pipeline finished in 3 h and used at maximum 16.95 GB of memory. These benchmarking results demonstrate that MINTIE's utility is not hindered by computational considerations, and many samples can be processed on a medium sized server in a reasonable time frame.

## Discussion

Here, we present the MINTIE pipeline, a caller for detecting a broad range of transcribed variants  $\geq 7$  bp in size. MINTIE excels at detecting unusual and rare upregulated variants. MINTIE frames variant calling as a problem of detecting sample specific expression of transcripts with novel sequence, without relying on read alignment to the reference genome. MINTIE achieves this by performing de novo assembly to reconstruct all transcripts expressed in a sample. It then leverages the unique approach of using equivalence class counts to reliably detect differential expression of unannotated sequence compared to a set of controls. Once novel transcripts are identified, MINTIE performs several steps to filter out transcripts found in the reference annotation and finally annotates the remaining variants into subtypes.

We demonstrated MINTIE's ability to find a diverse range of fusions, structural and splice transcriptomic variants by simulating 1500 variants, where MINTIE achieved > 86% recall. We further benchmarked fusion callers, transcribed structural variant callers and splice assemblers and found that MINTIE detected the full range of simulated variant types, while the other methods could not. We validated the ability of MINTIE to detect fusions and duplications in a set of validated AML cancer samples. Even for fusion finding, where many reliable methods exist, MINTIE's sensitivity was comparable. We propose that MINTIE could serve as either a first-pass scanner that is able to detect a large range of variants including fusions or as a transcriptome-wide discovery tool for rare variants if nothing is found using conventional fusion finders.

Using the Leucegene data, we showed that using several well-matched controls are important for reducing the number of background variants. For many cases, other

disease samples with similar expression profiles can be used as controls. In this instance, we advise using samples with known driver variants, so that the unique variants to be detected are unlikely to be present in the controls. Identical variants that are present in the set of controls will limit detection sensitivity in the case sample.

MINTIE has the potential to discover a wide range of important events that have thus far evaded detection due to the limitation of conventional analysis approaches. We demonstrated this by running MINTIE on a cohort of 87 B-ALL samples, where we found interesting events in disease relevant genes, including 3 samples with the same RB1 unpartnered fusion, IKZF1 and PAX5 PTDs, and 2 samples with altered splicing of ETV6. On rare disease samples, we showed that MINTIE could detect novel splice variants and, importantly, was able to discover diagnostically relevant genomic rearrangements from RNA-seq alone.

It is important to consider some of the limitations of MINTIE. The methodology has been optimised to detect many variant types and thus some specialised approaches geared at detecting fewer variant types outperform MINTIE in accuracy and runtime. MINTIE also requires sufficient expression of the variant relative to the controls. A further consideration is that MINTIE is, to some extent, reliant on appropriate control samples. MINTIE is a method based on RNA sequencing and can only identify structural rearrangements that affect transcribed regions. Similarly, lowly expressed genes or variants that knock out the expression of a gene will be difficult to detect. Finally, MINTIE utilises de novo transcriptome assembly to reconstruct transcript sequences. While this has the advantage of allowing complex variants to be detected, the run-time and RAM requirements are generally higher than many of the alignment-based approaches and the shortcomings of this approach also affect the method, such as difficulty with repetitive regions. We chose an assembler with good performance while minimising compute time. Slower, more computationally intensive assemblers may reconstruct more accurate transcripts. Although we have not benchmarked the effects of different transcript assemblers, MINTIE allows an assembled transcriptome from any source to be passed to the program bypassing default assembly by SOAPdenovo-Trans. Looking forward, long read transcriptome sequencing technology is readily being adopted and MINTIE could be used with polished long read data in tandem with short reads, avoiding the disadvantages of assembly.

Here, we have presented a novel, fast, transcriptome-wide approach for detecting rare, transcribed variants, and have demonstrated the use of the tool through simulations, real cancer and non-cancer samples and have demonstrated the detection of several alterations, some of which are clinically relevant. We propose that MINTIE will be used to identify novel drivers and loss-of-function variants in cancer and in rare diseases. Unlike existing methods, which are targeted at specific types of rearrangements or splicing, or specific genes, MINTIE is more agnostic to variant type. We have tested MINTIE on several types of variants, but there are potentially many other types which MINTIE could detect, for example, retrotransposon insertions into genes or complex events involving multiple rearrangements, which could be generated by chromoplexy or chromothripsis. Little is known about the frequency or characteristics of transcribed variants, but we now have a method to uncover and study them.

## Methods

### Pipeline detail

#### *De novo assembly*

RNA-seq reads of the case sample are de-duplicated using fastuniq [46] v1.1 and trimmed using Trimmomatic [47] v0.39 and then de novo assembled using SOAPdenovo-Trans [30] v1.03 with k-mers of 29, 49 and 69 (by default). (Optionally, Trinity [48] or rna-SPAdes [49] may be used for assembly, or the user may provide their own assembly.) The different k-mer assemblies are merged, sequences are de-duplicated and assembled transcripts longer than 150 bp (by default) are retained.

#### *Quantification*

All assembled transcripts are merged with the CHESSE [50] v2.2 transcriptome reference sequence, which is then indexed by Salmon [51] v0.14. Salmon is then run on this index for the case and control samples in single-end mode (otherwise short assembled transcripts may not be correctly counted). Sequence bias is corrected (--seqBias), equivalence class (EC) counts are extracted (--dumpEq), mapping validation (--validate-Mappings) and hard filtering (--hardFilter) are enabled.

#### *Differential expression*

ECs are matched between all samples by membership (i.e. an EC composed of transcripts A, B and C in sample 1 is given an identifier, and will be considered the same EC as one composed of the same transcripts in sample 2). Library sizes are calculated prior to any filtering. ECs containing *only* assembled transcripts (i.e. the EC does not contain any transcripts from the reference) are kept and, after light filtering (CPM > 0.1 in the case sample), the case sample is compared to the controls using edgeR [52] v3.26.5. Dispersions are estimated using the classic (non-GLM) model [53], and transcripts are fitted using edgeR's GLM quasi-likelihood fit function [54] with the robust flag set to true [55]. edgeR's GLM likelihood-ratio test [56] is used to perform differential expression. All transcripts associated with significant ECs are retained (FDR < 0.05 and logFC > 2 by default; we used a logFC > 5 for all analyses unless otherwise indicated).

#### *Annotation*

All significant transcripts are aligned using GMAP [57] (2020-06-04) to the hg38 genome with all alternative reference contigs removed. GMAP's --max-intronlength-ends parameter is set to 500 kb to prevent soft-clipping at long terminal introns, and --chimera-margin is set to correspond to MINTIE's clip length (default = 20). All aligned transcripts are extracted and compared to the CHESSE [50] v2.2 reference GTF to identify fusions, TSVs and novel splicing. We require that at least 30 base-pairs and 30% of the transcript sequence is aligned somewhere in the genome to retain a candidate transcript. Transcripts not overlapping any known reference exons are discarded. Broadly, we set the criteria for novel variants as follows: (i) the transcript contains an insertion or deletion of at least 7 base-pairs, (ii) the transcript has a soft-clip or hard-clip (based on the GMAP's alignment) that is least 20 base-pairs in length, (iii) the transcript contains a splice junction not seen in the reference transcriptome and (iv) the

transcript contains a block of at least 20 base-pairs that is not normally transcribed in the reference genome. Passing category (i), (ii), (iii) and/or (iv) will retain the transcript as potentially harbouring a candidate variant. We discard any variants matching category (iv) where no splicing occurs (i.e. the transcript is aligned as an unspliced block) unless the transcript's start and end are contained within flanking exons, indicating a potential retained intron.

Annotated variants are further refined by matching novel blocks with novel splice junctions and identifying valid donor/acceptor sites for novel junctions (defined as observing the sequences AG/GT or AC/CT before the splice junction). By default, the motif checking allows tolerance for one mutation on either end of a junction (this can be made more or less lenient by the user). As extended/novel exons should create novel exon boundaries, we consolidate these variants and discard any extended/novel exons not supported by a corresponding novel junction. Novel and extended exons must be accompanied by at least one novel junction at either (or both) ends of the block and have a valid donor or acceptor motif. Alignment gaps (of at least 7 bp) within single exons are checked to capture potential novel intron, and exon ends are checked for truncation by at least 20 bp. A single transcript may have multiple annotated variants, but only those matching specified filters are retained. The variant sizes listed above for filtering are defaults only and may be adjusted in MINTIE.

#### Selecting appropriate control samples

Control samples ideally have the same transcriptional profile as the case samples but without the variants of interest. Controls can be from normal samples from the same tissue type as the case sample if they still have a similar transcriptional program as the case. As RNA sequencing data from the given tissue may be difficult to obtain, other samples from the same cancer type (and same tissue type) can be used as controls. We recommend using at least one control, but ideally 3–10 (using large numbers of controls, > 20 for example, will result in increased processing time). In some cases, no appropriate controls may be available, in which case MINTIE can be run without controls. In this mode, differential expression is not performed; however, novel ECs are still selected based on read quantification.

#### Simulations

We simulated 1500 variants using code available in the MINTIE source code repository ([https://github.com/Oshlack/MINTIE/tree/master/simu/run\\_simu.py](https://github.com/Oshlack/MINTIE/tree/master/simu/run_simu.py)). Simulations were generated by extracting sequences from the transcripts listed in the hg38 UCSC RefSeq reference and simulating reads from the resulting sequence. One hundred variants from 15 variant types were generated (five fusion types: canonical, extended exon, novel exon, with insertion and unpartnered, five TSV types: insertions, deletions, ITDs, PTDs and inversions, and five novel splice variants: extended exons, novel exons, truncated exons, skipped exons and retained introns). Only transcripts from genes that did not overlap any other genes were used in the simulation. Additionally, each transcript had to have at least 3 exons to be considered as a simulation transcript.

All fusions were simulated by selecting the first two and the last two exons from two random transcripts from different genes and inserting the intervening sequence.

Canonical fusions contained no intervening sequence, while fusions with extended exons inserted 30–199 bp of intronic sequence from the end of the second exon of the first transcript. Similarly, fusions with novel exons contained intronic sequence 30–199 bp downstream with a size of 30–199 bp. Non-canonical fusions with insertions were generated by inserting 7–49 bp of randomly-generated sequence between the two fusion transcripts.

Small TSVs were generated by inserting, duplicating or deleting sequence within randomly selected exons from randomly selected transcripts. These small variant types were between 7 and 49 base-pairs and had to reside at least 10 bp within the exon. Inversions and partial-tandem duplications were generated by selecting 1–3 random exons within a transcript and either inverting or duplicating their sequence in tandem. Extended and novel exons were generated by adding intronic sequence downstream (directly, or with a 30–199 bp gap) of a randomly selected exon. To ensure that novel or extended exons did not overlap exons from other transcripts (or downstream exons of the same transcript), each candidate exon was checked for these potential overlaps (which would otherwise result in obfuscation of the variant, or the wrong variant type being created). Novel junction (skipped exon) variants were created by selecting a random pair of exons and checking whether an existing junction existed between them, creating a transcript with this junction if not. Two randomly selected neighbouring exons were both truncated at their facing ends (end and start respectively) by 30–199 bp. Retained introns included sequence from a randomly selected intron from a given transcript that was > 30 bp. The presence of correct splicing motifs was not considered for the simulation.

In addition to each variant gene, the sequence of the unaltered wild-type gene was added to the simulated case sample's reference. An additional 100 unaltered background genes were also added to the case sample. A control sample reference was also generated, which included the unaltered wildtype sequence only for all simulated transcripts. ART-illumina [31] v2.5.8 was run on the corresponding references with 100 bp paired-end reads with a fragment size of 300 and coverage of 50x (transcripts thus have an effective coverage of 100x, given the bi-allelic reference containing variant and wild-type transcripts).

The simulated variant fastq files were run through the MINTIE pipeline with default k-mers (29, 49 and 69). Minimum read length was set to 80 (for trimming), minimum assembled transcript length was set to 100, and motif checking was disabled. To ensure that setting a fixed low dispersion value did not adversely over-state the results, we also reran the simulations at four extra dispersion levels: 0.3, 0.5, 0.6 and 0.7 (the original being 0.1). Increasing dispersion had minimal effect on the results until a value between 0.6 and 0.7, after which very few ECs are found to be differentially expressed (Additional file 1: Fig. S1B). This indicated that dispersion is unlikely to play a significant role in the observed results, unless the value is abnormally high (> 0.6). All other parameters were left at default. Seven other tools were run on the resulting simulation output: TAP [15] (commit 8940e45), Barnacle [17] v1.0.4, SQUID v1.5 [16], JAFFA [24] v1.0.7, Arriba [26] v1.1.0, CICERO [18] v1.3.0, KisSplice [11] v2.4.0 and StringTie [12] v1.3.6. TAP was run using the tap2.py pipeline with k-mer lengths of 49 and 79, and --max\_diff\_splice set to 4 to bypass motif checking. Other parameters were kept at default. The transcript reference was generated with PAVfinder's extract\_transcript\_

sequence.py script, using the UCSC RefGene sequence reference as input. Barnacle was run with its default hg19 reference. SQUID was run using STAR's alignments. JAFFA was run using its direct mode and default reference. References used in generating simulations were used whenever possible. Default parameters were used for each tool unless otherwise specified. Samtools [58] v1.8, BLAT [44] v36.1, STAR [59] v2.5.3a and bwa [60] v0.7.17 were used by all methods that required these tools. CICERO was run through the St. Jude Cloud Genomics Platform, using fastq files as input (STAR alignments were performed through St. Jude's Rapid RNA-seq pipeline, which incorporates CICERO).

Variants were counted as true positives for all tools run for the benchmarking analysis if a variant was reported in the same gene name as the simulated variant (for MINTIE, Arriba and TAP). Gene coordinates were used for all other tools where gene names were not reported (SQUID, StringTie and KissSplice) or where using specific references proved prohibitive (Barnacle). Only transcripts marked as novel were considered for StringTie. Sequence results from KisSplice were aligned with GMAP [57] (2019-05-12), and aligned coordinates were extracted in order to obtain variant positions. As long as a variant coordinate was found within the expected gene region, this was counted as a hit. Only one gene/coordinate match was needed (of a partnered fusion pair) to be called a true positive. As Barnacle and CICERO were run using a hg19 reference, we used UCSC liftover (<https://genome.ucsc.edu/cgi-bin/hgLiftOver>) to convert the hg38 simulation coordinates to hg19 (6 TSVs/NSVs and 9 fusion gene coordinates were partially deleted in the liftover; the TSVs/NSVs could not be counted, while the fusion were only affected at one coordinate and could still be counted). All code used to process each tool's results can be found in the paper analysis code (see Code Availability).

To perform experiments with varying coverage, the simulated fastq files were down-sampled using seqtk v1.0 (<https://github.com/lh3/seqtk>) at 40x, 20x and 10x (resulting in variant coverage of half of each of these values, due to the heterogeneous nature of the simulation). The varying dispersion experiments were generated by manually changing the fixed dispersion parameter (0.1) in the MINTIE source code, to values 0.7, 0.6, 0.5, 0.3. No other parameters were modified. Detection limit experiments were generated by generating nine simulations of 300 variants each (deletions, insertions and ITDs), each with the same random seed (123), only varying by the variant size, ranging from 1 to 9.

### Running MINTIE on Leucegene samples

MINTIE was run on four Leucegene cohorts: the NUP98-NSD1 fusion cohort [35], the CBF AML cohort [61], the KMT2A-PTD cohort [32] and the non-cancer peripheral blood (normals) data set [34], using k-mer sizes of 49 and 79 for assembly. We used the patient IDs identified by Audemard et al. [62] to identify the samples containing KMT2A-PTDs, as these were not disclosed in the original paper. For all AML data sets, we used 13 Leucegene normals as the controls (5 monocytes, 5 granulocytes and 3 TWBCs). We also ran the KMT2A-PTD data set against a set of randomly selected AML samples from the CBF cohort (see Additional file 1: Note S1 for the control samples used). The normal samples were run individually against the remaining samples

from that cell type. (The pooled peripheral CD34+ blood cells were not considered for this analysis.) In the granulocyte 1-4 control experiments, all combinations of cases and controls were run for the five total samples, resulting in four total runs for 1 vs. 1, three total runs for 1 vs. 2, two total runs for 1 vs. 3 and a single run for 1 vs. 4.

Salmon [51] (v0.14) was run separately on the Leucegene normals, the KMT2A-PTD cohort and the selected CBF AML control samples in paired-end mode with default parameters and the `--seqBias` flag using the same CHES transcriptome reference used in all other analyses. Counts were summarised at the gene level using `tximport` [63] (v1.12.3) with the `countsFromAbundance='lengthScaledTPM'` parameter. The `voom` function from `Limma` [64] v3.42.0 was used to normalise libraries, and the `base prcomp` function (R v3.6.2) was used to perform principal component analysis.

Squid and Arriba were run on the Leucegene samples using the same parameters and workflow as on the simulations. As the experiment emphasised sensitivity and compared fewer tools, variants were checked with greater stringency, requiring, in the case of fusions, both variant genes or loci to be called in one variant call (or contig), or as an unknown soft-clip variant containing either fusion gene (in MINTIE's case). This only happened in two cases, and these variants were manually validated by BLATing the contig sequence to ensure the correct fusion was being reported. In cases where a variant was contained in a single gene and caller output listed gene 1 and gene 2, these both had to match the single gene. A variant filter was also used where classifications were reported (MINTIE, Arriba and TAP) to ensure variants were plausibly classified. Plausible classifications for fusions included 'fusion' or 'unknown' for MINTIE, 'translocations' or 'inversions' for Arriba and 'fusions' for TAP. Plausible classifications for ITDs included 'insertion', 'unknown' or 'intergenic rearrangement' for MINTIE, 'ITD' or 'duplication' for Arriba and 'ITD' for TAP. Plausible classifications for PTDs included 'intergenic rearrangement' or 'unknown' for MINTIE, and 'duplication' for Arriba and TAP.

#### Running MINTIE on B-ALL and TCGA samples

B-ALL samples used in this study were obtained from the Royal Children's Hospital; data is described in Brown et al. [36]. Seven samples were excluded due to significantly different sequencing characteristics (shorter read lengths and using an unstranded protocol). Patient samples were selected as controls if they had a positive molecular test and the driver fusions were confirmed in the RNA-seq data. If there were multiple samples obtained from the same patient, they were not considered as controls. This resulted in 34 controls (see Additional file 1: Note S2 for a list). MINTIE was run with k-mer sizes 49 and 79 and a  $\log_{FC} > 5$  cutoff. The RB1 variants were detected in two samples initially, and a third was found after rerunning the pipeline at a lower  $\log_{FC}$  cutoff ( $> 2$ ). For the TCGA samples that were identified to contain the variant (using `SeqOthello`), we obtained the RNA-seq sequence data and ran these samples through the MINTIE pipeline. We used single random TCGA control of the same disease type, sex and similar read length.  $\log_{FC}$  cutoff was 2 and k-mer sizes of 29 and 39 were used due to short read lengths (samples had 51 bp and 76 bp reads for the BRCA and ESCA samples respectively).

### Running MINTIE on rare disease samples

MINTIE was run on each of the 52 available samples described in Cummings et al. [14] against a set of 10 muscle sample controls from GTEx [43]. Motif checking was turned off due to the presence of SNP-induced splice sites described in the study; all other parameters were set as default. We manually inspected variants called and transcripts assembled in order to determine whether a given variant was called. In order to look for novel candidate variants, we considered all available samples that did not have RNA variants described in the Cummings et al. paper. We applied the following additional filters to the results: (i) VAF > 0.1, (ii) total reads in controls < 10, (iii) < 10 variants found on the de novo assembled transcript, (iv) the variant was classified as a TSV or (any type of) fusion and (v) the affected variant affected a gene within the NMD gene list used by the Cummings et al. authors. Each variant was then manually curated considering contig alignment, read coverage and variant location.

### Abbreviations

SVs: Structural variants; TSVs: Transcribed structural variants; NSVs: Novel splice variants; ECs: Equivalence classes; ITD: Internal tandem duplication; PTD: Partial tandem duplication; WGS: Whole genome sequencing; TCGA: The cancer genome atlas; AML: Acute myeloid leukaemia; B-ALL: B cell acute lymphoblastic leukaemia; VAF: Variant allele frequency; NTS: Non-templated sequence

### Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s13059-021-02507-8>.

**Additional file 1:** Supplementary Figures, Tables and Notes

**Additional file 2:** Review history

### Acknowledgements

We would like to thank Jinze Liu and Xiaofei Zhang for querying TCGA with SeqOthello for us and help from Christoffer Flensburg on interpreting variants and software suggestions and Beryl Cummings, Leigh Waddell, Sandra Cooper and Daniel MacArthur for correspondence about the rare disease data. Some results presented in this manuscript used data generated by the TCGA Research Network: <https://www.cancer.gov/tcga>. Computing resources were provided by MCRI, The Peter MacCallum Cancer Centre and The University of Melbourne Science IT.

### Peer review information

Tim Sands was the primary editor of this article and managed its editorial process and peer review in collaboration with the rest of the editorial team.

### Authors' contributions

MC: formal analysis, methodology, software, visualisation, writing—original draft preparation, writing—review and editing; BS: visualisation; IJM: writing—review and editing, PGE: writing—review and editing; AO: conceptualization, supervision, methodology, writing—original draft preparation, writing—review and editing; NMD: conceptualization, supervision, methodology, software, writing—original draft preparation, writing—review and editing. The authors read and approved the final manuscript.

### Funding

This work was funded by NHMRC project grant GNT1140626 to AO, IJM, PGE and NMD.

### Availability of data and materials

- MINTIE software: <https://github.com/Oshlack/MINTIE> [65]. The archived code used in the paper analyses can be obtained from Zenodo [66].
- Code for generating paper analyses and figures: <https://github.com/Oshlack/MINTIE-paper-analysis>. An archived version can be found in Zenodo [67]. Figures can be viewed at <https://oshlacklab.com/MINTIE-paper-analysis>.
- The 1500 variant simulation (15 variant types), including the downsampled versions, are available in Zenodo [68], or can be reproduced using the `run_simu.py` script under <https://github.com/Oshlack/MINTIE/tree/master/simu>, using the same random seeds as the `fullsimu_params.ini` file. The downsampled simulations can be reproduced by using `seqtk v1.0` with random seeds of 1587, 6471 and 9505 for the 40x, 20x and 10x simulations respectively.
- The 2700 simulation for testing INDEL and ITD detection limits can be found in Zenodo [69] and can be reproduced using the code mentioned above (with variant sizes and numbers adjusted as mentioned in the "Methods" section).
- The Leucegene cohorts are available on the Sequence Read Archive: the NUP98-NSD1 fusion cohort [35] (GSE49642, GSE67039 and GSE52656), the CBF AML cohort [61] (GSE49642, GSE52656, GSE62190, GSE66917 and GSE67039), the

KMT2A-PTD cohort [32] (accessions GSE49642, GSE52656, GSE66917 and GSE67039) and the non-cancer peripheral blood (normals) data set [34] (accession GSE51984).

- The paediatric B-ALL samples are available from the European Genome-Phenome Archive (accession number EGAS00001004212).
- The rare disease and control data are available from dbGaP under accession IDs phs000655.v3.p1 and phs000424.v6.p1 (SRR809444, SRR810225, SRR811771, SRR813656, SRR813983, SRR809595, SRR810249, SRR812773, SRR813802 and SRR815020) respectively.
- TCGA samples were obtained from dbGaP under accession ID phs000178
  - Case IDs: 71c5ab4f-ce13-432d-9a90-807ec33cf891 (BRCA) and eae803bf-172f-492f-a381-8f4c040232a2 (ESCA)
  - Control IDs: b5e182ff-159a-44af-881e-8f21bbe96193 (BRCA) and 241a9e1b-3f1a-4eca-90be-d5d48fedce6d (ESCA)

## Declarations

### Ethics approval and consent to participate

Samples analysed in the B-ALL cohort were approved for collection and analysis by the Royal Children's Hospital Human Research Ethics Committee (HREC 34127).

### Competing interests

The authors declare no competing interests.

### Author details

<sup>1</sup>Peter MacCallum Cancer Centre, Melbourne, VIC, Australia. <sup>2</sup>Murdoch Children's Research Institute, Parkville, Australia. <sup>3</sup>Sir Peter MacCallum Department of Oncology, The University of Melbourne, Parkville, Australia. <sup>4</sup>School of BioSciences, University of Melbourne, Parkville, Australia. <sup>5</sup>Walter and Eliza Hall Institute, Parkville, Australia. <sup>6</sup>Faculty of Medicine, Dentistry and Health Sciences, University of Melbourne, Parkville, Australia. <sup>7</sup>Children's Cancer Institute, UNSW, Sydney, Australia. <sup>8</sup>Department of Paediatrics, University of Melbourne, Parkville, Australia.

Received: 28 July 2020 Accepted: 27 September 2021

Published online: 22 October 2021

## References

1. Saito M, et al. Development of Lung Adenocarcinomas with Exclusive Dependence on Oncogene Fusions. *Cancer Res.* 2015;75:2264–72.
2. Patch A, et al. Whole-genome characterization of chemoresistant ovarian cancer. *Nature.* 2015;489–94. <https://doi.org/10.1038/nature14410>.
3. Grimwade D, et al. Refinement of cytogenetic classification in AML Younger adult patients treated in UKMRC. *Blood.* 2010;116:354–66.
4. Li Y, et al. Patterns of somatic structural variation in human cancer genomes. *Nature.* 2020;578:112–21.
5. Sanchis-Juan A, et al. Complex structural variants in Mendelian disorders: identification and 27 breakpoint resolution using short- and long-read genome sequencing. *Genome Med.* 2018;10:95.
6. Holt JM, et al. Identification of pathogenic structural variants in rare disease patients through genome Sequencing. *bioRxiv.* 2019;627661. <https://doi.org/10.1101/627661>.
7. Calabrese C, et al. Genomic basis for RNA alterations in cancer. *Nature.* 2020;578:129–36.
8. Haas BJ, et al. Accuracy assessment of fusion transcript detection via read-mapping and de novo fusion transcript assembly-based methods. *Genome Biol.* 2019;20:1–16.
9. Kumar A, et al. Substantial interindividual and limited intraindividual genomic diversity among tumors from men with metastatic prostate cancer. *Nat Med.* 2016;22:1–13.
10. Trapnell C, et al. Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks. *Nat Protoc.* 2012;7:562–78.
11. Sacomoto GAT, et al. KISSPLICE: de-novo calling alternative splicing events from RNA-seq data. *BMC Bioinformatics.* 2012;13:1–12.
12. Pertea M, et al. StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat Biotechnol.* 2015; 33:290–5.
13. Gonorazky HD, et al. Expanding the Boundaries of RNA Sequencing as a Diagnostic Tool for Rare Mendelian Disease. *Am J Hum Genet.* 2019;104:1007.
14. Cummings BB, et al. Improving genetic diagnosis in Mendelian disease with transcriptome sequencing. *Sci Transl Med.* 2017;9:eaal5209.
15. Chiu R, Nip KM, Chu J, Birol I. TAP: a targeted clinical genomics pipeline for detecting transcript variants using RNA-seq data. *BMC Med Genet.* 2018;11:79.
16. Ma C, Shao M, Kingsford C. SQUID: Transcriptomic structural variation detection from RNA-seq. *Genome Biol.* 2018;19:1–16.
17. Swanson L, et al. Barnacle: detecting and characterizing tandem duplications and fusions in transcriptome assemblies. *BMC Genomics.* 2013;14:550.
18. Tian L, et al. CICERO: a versatile method for detecting complex and diverse driver fusions using cancer RNA sequencing data. *Genome Biol.* 2020;21:126.
19. Mullighan CG, et al. Deletion of IKZF1 and Prognosis in Acute Lymphoblastic Leukemia. *N Engl J Med.* 2009;360:470–80.
20. Bolouri H, et al. The molecular landscape of pediatric acute myeloid leukemia reveals recurrent structural alterations and age-specific mutational interactions. *Nat Med.* 2017. <https://doi.org/10.1101/125609>.
21. Sudmant PH, et al. An integrated map of structural variation in 2,504 human genomes. *Nature.* 2015; 526:75–81.
22. STAR-Fusion: Fast and Accurate Fusion Transcript Detection from RNA-Seq | bioRxiv. <https://www.biorxiv.org/content/10.1101/120295v1.abstract>.

23. Kim D, Salzberg SL. TopHat-Fusion: An algorithm for discovery of novel fusion transcripts. *Genome Biol.* 2011;12:1–15.
24. Davidson NM, Majewski IJ, Oshlack A. JAFFA: High sensitivity transcriptome-focused fusion gene detection. *Genome Med.* 2015;7:43.
25. Melsted P, et al. Fusion detection and quantification by pseudoalignment. *bioRxiv.* 2017;166322:10.1101/166322.
26. Uhrig S, et al. Accurate and efficient detection of gene fusions from RNA sequencing data. *Genome Res.* 2021;gr.257246:119. <https://doi.org/10.1101/gr.257246.119>.
27. Qiu Y, Ma C, Xie H, Kingsford C. Detecting transcriptomic structural variants in heterogeneous contexts via the Multiple Compatible Arrangements Problem. *Algorithms Mol Biol.* 2020;15:9.
28. Audoux J, et al. DE-kupl: exhaustive capture of biological variation in RNA-seq data through k-mer decomposition. *Genome Biol.* 2017;18:243.
29. Xie Y, et al. SOAPdenovo-Trans: De novo transcriptome assembly with short RNA-Seq reads. *Bioinformatics.* 2014;30:1660–6.
30. O'Leary NA, et al. Reference sequence (RefSeq) database at NCBI: current status, 29 taxonomic expansion, and functional annotation. *Nucleic Acids Res.* 2016;44:D733–45.
31. Huang W, Li L, Myers JR, Marth GT. ART: a next-generation sequencing read simulator. *Bioinformatics.* 2012;28:593–4.
32. Lavallée V-P, et al. The transcriptomic landscape and directed chemical interrogation of MLL-rearranged acute myeloid leukemias. *Nat Genet.* 2015;47:1030–7.
33. Audemard É, et al. Target variant detection in leukemia using unaligned RNA-Seq reads. *bioRxiv.* 2018;295808. <https://doi.org/10.1101/295808>.
34. Pabst C, et al. GPR56 identifies primary human acute myeloid leukemia cells with high repopulating potential in vivo. *Blood.* 2016;127:2018–27.
35. Lavallée VP, et al. Identification of MYC mutations in acute myeloid leukemias with NUP98-NSD1 translocations. *Leukemia.* 2016;30:1621–4.
36. Brown LM, et al. The application of RNA sequencing for the diagnosis and genomic classification of pediatric acute lymphoblastic leukemia. *Blood Adv.* 2020;4:1–3.
37. Gröbner SN, et al. The landscape of genomic alterations across childhood cancers. *Nature.* 2018;555:321–7.
38. Ma X, et al. Pan-cancer genome and transcriptome analyses of 1,699 paediatric leukaemias and solid tumours. *Nat Publ Group.* 2018. <https://doi.org/10.1038/nature25795>.
39. Mullighan CG, et al. Genome-wide analysis of genetic alterations in acute lymphoblastic leukaemia. *Nature.* 2007;446:758–64.
40. Gu Z, et al. PAX5-driven subtypes of B-progenitor acute lymphoblastic leukemia. *Nat Genet.* <https://doi.org/10.1038/s41588-018-0315-5>.
41. Zhang J, et al. Key pathways are frequently mutated in high-risk childhood acute lymphoblastic leukemia: a report from the Children's Oncology Group. *Blood.* 2011;118:3080–7.
42. Yu Y, et al. SeqOthello: Query over RNA-seq experiments at scale. *bioRxiv.* 2018;258772. <https://doi.org/10.1101/258772.30>.
43. The GTEx Consortium. The Genotype-Tissue Expression (GTEx) pilot analysis: Multitissue gene regulation in humans. *Sci.* 2015;348:648–60.
44. Kent WJ. BLAT—The BLAST-Like Alignment Tool. *Genome Res.* 2002;12:656–64.
45. Waddell LB, et al. WGS and RNA Studies Diagnose Noncoding DMD Variants in Males With High Creatine Kinase. *Neurol Genet.* 2021;7:e554.
46. Xu H, et al. FastUniq: A Fast De Novo Duplicates Removal Tool for Paired Short Reads. *PLoS One.* 2012;7:e52249.
47. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics.* 2014;30:2114–20.
48. Haas BJ, et al. De novo transcript sequence reconstruction from RNA-Seq: reference generation and analysis with Trinity. *Nat Protoc.* 2014;8:1494–512.
49. Bushmanova E, Antipov D, Lapidus A, Pribelski AD. rnaSPAdes: a de novo transcriptome assembler and its application to RNA-Seq data. *GigaScience.* 2019;8:giz100.
50. Perteu M, et al. CHESS: a new human gene catalog curated from thousands of large-scale RNA sequencing experiments reveals extensive transcriptional noise. *Genome Biol.* 2018;19:332825.
51. Patro R, Duggal G, Love MI, Irizarry RA, Kingsford C. Salmon provides fast and bias-aware quantification of transcript expression. *Nat Methods.* 2017;14:021592.
52. Robinson MD, McCarthy DJ, Smyth GK. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinforma Oxf Engl.* 2010;26:139–40.
53. Chen Y, Lun ATL, Smyth GK. Differential Expression Analysis of Complex RNA-seq 31 Experiments Using edgeR. In: Datta S, Nettleton D, editors. *Statistical Analysis of Next Generation Sequencing Data*: Springer, Cham; 2014. p. 51–74. [https://doi.org/10.1007/978-3-319-07212-8\\_3](https://doi.org/10.1007/978-3-319-07212-8_3).
54. Lund SP, Nettleton D, McCarthy DJ, Smyth GK. Detecting differential expression in RNA-sequence data using quasi-likelihood with shrunken dispersion estimates. *Stat Appl Genet Mol Biol.* 2012;11:5.
55. Phipson B, Lee S, Majewski IJ, Alexander WS, Smyth GK. Robust hyperparameter estimation protects against hypervariable genes and improves power to detect differential expression. *Ann Appl Stat.* 2016;10:946–63.
56. McCarthy DJ, Chen Y, Smyth GK. Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Res.* 2012;40:4288–97.
57. Wu TD, Reeder J, Lawrence M, Becker G, Brauer MJ. GMAP and GSNAP for Genomic Sequence Alignment: Enhancements to Speed, Accuracy, and Functionality. In: Mathé E, Davis S, editors. *Statistical Genomics: Methods and Protocols*: Humana Press, New York, NY; 2016. p. 283–334. [https://doi.org/10.1007/978-1-4939-3578-9\\_15](https://doi.org/10.1007/978-1-4939-3578-9_15).
58. Li H, et al. The Sequence Alignment/Map format and SAMtools. *Bioinforma Oxf Engl.* 2009;25:2078–9.
59. Dobin A, et al. STAR: Ultrafast universal RNA-seq aligner. *Bioinformatics.* 2013;29:15–21.
60. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinforma Oxf Engl.* 2009;25:1754–60.
61. Lavallée VP, et al. RNA-sequencing analysis of core binding factor AML identifies recurrent ZBTB7A mutations and 32 defines RUNX1-CBFA2T3 fusion signature. *Blood, Am J Hematol.* 2016;128:872–5.

62. Audemard EO, et al. Targeted variant detection using unaligned RNA-Seq reads. *Life Sci Alliance*. 2019;2:e201900336.
63. Sonesson C, Love MI, Robinson MD. Differential analyses for RNA-seq: transcript-level estimates improve gene-level inferences. *F1000Research*. 2016;4:1521.
64. Ritchie ME, et al. Limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res*. 2015;43:e47.
65. Cmero M, et al. MINTIE v0.2.0 code for Genome Biology paper. (GitHub, 2020). <https://github.com/Oshlack/MINTIE>.
66. Cmero M, et al. MINTIE v0.2.0 code for Genome Biology paper. (Zenodo, 2020). doi:<https://doi.org/10.5281/zenodo.5516712>.
67. Cmero M, et al. Oshlack/MINTIE-paper-analysis. *Genome Biol*. 2021. <https://doi.org/10.5281/zenodo.5516708>.
68. Cmero M, et al. 1,500 simulated transcriptomic variants for MINTIE paper. (2020) doi:<https://doi.org/10.5281/zenodo.4876713>.
69. Cmero M, et al. 2,700 simulated small INDELS and ITDs for MINTIE paper. (2021) doi:<https://doi.org/10.5281/zenodo.4876678>.

### **Publisher's Note**

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.