

Targeted RNA-seq improves efficiency, resolution, and accuracy of allele specific expression for human term placentas

Weisheng Wu,¹ Jennie L. Lovett,² Kerby Shedden,³ Beverly I. Strassmann,^{2,4} and Claudius Vincenz^{4,*}

¹BRCF Bioinformatics Core, University of Michigan, Ann Arbor, MI 48109, USA,

²Department of Anthropology, University of Michigan, Ann Arbor, MI 48109, USA,

³Department of Statistics, University of Michigan, Ann Arbor, MI 48109, USA and

⁴Research Center for Group Dynamics, Institute for Social Research, University of Michigan, Ann Arbor, MI 48106, USA

*Corresponding author: University of Michigan, 5255 ISR, 426 Thompson Street, Ann Arbor, MI 48106, USA. Email: vincenz@umich.edu

Abstract

Genomic imprinting is an epigenetic mechanism that results in allele-specific expression (ASE) based on the parent of origin. It is known to play a role in the prenatal and postnatal allocation of maternal resources in mammals. ASE detected by whole transcriptome RNA-seq (wht-RNAseq) has been widely used to analyze imprinted genes using reciprocal crosses in mice to generate large numbers of informative SNPs. Studies in humans are more challenging due to the paucity of SNPs and the poor preservation of RNA in term placentas and other tissues. Targeted RNA-seq (tar-RNAseq) can potentially mitigate these challenges by focusing sequencing resources on the regions of interest in the transcriptome. Here, we compared tar-RNAseq and wht-RNAseq in a study of ASE in known imprinted genes in placental tissue collected from a healthy human cohort in Mali, West Africa. As expected, tar-RNAseq substantially improved the coverage of SNPs. Compared to wht-RNAseq, tar-RNAseq produced on average four times more SNPs in twice as many genes per sample and read depth at the SNPs increased fourfold. In previous research on humans, discordant ASE values for SNPs of the same gene have limited the ability to accurately quantify ASE. We show that tar-RNAseq reduces this limitation as it unexpectedly increased the concordance of ASE between SNPs of the same gene, even in cases of degraded RNA. Studies aimed at discovering associations between individual variation in ASE and phenotypes in mammals and flowering plants will benefit from the improved power and accuracy of tar-RNAseq.

Keywords: targeted RNA-seq; quantification of ASE; allele-specific expression; human placenta; genomic imprinting

Introduction

Genomic imprinting is an epigenetic phenomenon that results in allele-specific expression (ASE) based on parent of origin. Many imprinted genes are found in the nutritive tissue of placental and marsupial mammals as well as flowering plants (Tucci *et al.* 2019; Batista and Köhler 2020). Under the kinship hypothesis, genomic imprinting evolved due to a conflict of interest between the genes an offspring inherited from its mother versus its father over the number of resources to be allocated to the current offspring (Moore and Haig 1991). One consequence of imprinting is that, for specific genomic regions, the paternal and maternal genomes are not equivalent. Evidence that both are required for normal development derives from a series of mouse studies in the 1980s that generated gynogenotes or androgenotes. In uniparental disomies (UPD), nonequivalency was limited to certain genomic regions that later were identified as clusters of imprinted genes (Tucci *et al.* 2019). In humans, about 100 imprinted genes have been identified (Babak *et al.* 2015; Baran *et al.* 2015). The highest proportion of imprinted genes were

expressed in embryonic, extra-embryonic and brain tissues (Babak *et al.* 2015), and impacted neurological development, placentation, and fetal growth (Peters 2014). Regulation of imprinting is governed by imprinting control regions (ICRs) through epigenetic mechanisms involving DNA methylation, lncRNAs, histone modifications, and high-order chromatin organization (Farhadova *et al.* 2019; Thamban *et al.* 2020).

High throughput sequencing technologies including RNA-seq and DNA methylation sequencing have been widely used to study genomic imprinting (Li and Li 2019). Transcriptome wide ASE is determined by combining quantification of whole transcriptome RNA-seq (wht-RNAseq) reads with identification of heterozygous SNPs in DNA (Wang and Clark 2014). Animal studies gain additional power from reciprocal cross breeding of closely related strains, which produces higher SNP densities and phased reference genomes that pinpoint the parent of origin of each allele at every SNP. These crosses permit imprinting to be distinguished from sequence dependent allelic expression bias (Wang *et al.* 2013, 2019; Babak *et al.* 2015; Chen *et al.* 2016). In humans, fewer SNPs are present than in crossbred animal

Received: January 25, 2021. Accepted: May 12, 2021

© The Author(s) 2021. Published by Oxford University Press on behalf of Genetics Society of America.

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

models due to lower genetic diversity; nonetheless RNA-seq has been successfully employed in many human tissues (Metsalu et al. 2014; Hamada et al. 2016; Mozaffari et al. 2018; Zink et al. 2018; Jadhav et al. 2019; Pilvar et al. 2019). Frequently, the parent's genotype is not available in human studies and ASE is determined without the parent of origin of the bias (Babak et al. 2015; Baran et al. 2015; Gulyás-Kovács et al. 2018).

To quantify ASE from RNA-seq, best practice protocols have been proposed to accommodate several technical factors (Castel et al. 2015). For example, appropriate alignment methods should be used to reduce the tendency for mapping bias to favor the reference alleles (Stevenson et al. 2013; Van De Geijn et al. 2015). A few studies showed that the accuracy of ASE quantified from RNA-seq was especially limited when the read depth on the measured SNPs was insufficient (Fontanillas et al. 2010; Heap et al. 2010; Nothnagel et al. 2011), which could lead to low power to predict imprinting and poor agreement of ASE between the SNPs from the same genes (Zou et al. 2019) and even from the same exons (DeVeale et al. 2012).

The consensus is that most imprinted genes in humans and mice have been identified with these genome-wide approaches (Babak et al. 2015; Baran et al., 2015; Zink et al. 2018). However, the well-documented population variability in imprinting and potential phenotypic effects are poorly understood (Zink et al. 2018; Vincenz et al. 2020). Thus, there is a need to quantify ASE with high precision in a population setting in a cost-efficient manner. In a previous study, we used wnt-RNAseq to measure ASE in 91 known imprinted genes in human term placentas collected from a cohort in Mali. We showed that departures from mono-allelic RNA expression were prevalent in many imprinted genes in this cohort. The number of reads we obtained from imprinted genes was limited because many highly expressed genes in placenta are not imprinted and constituted a large fraction of the total reads (Vincenz et al. 2020).

To overcome this limitation, we employed a tar-RNAseq approach to focus sequencing resources on the genes of interest. We enriched RNA against a targeting panel designed to cover exonic regions of 520 genes and quantified ASE on the informative SNPs in this panel. We hypothesized that this tar-RNAseq dataset could substantially increase the coverage of SNPs and genes of interest and improve the accuracy of ASE determination, compared to our previously published wnt-RNAseq dataset. In order to test this hypothesis, we performed the same ASE analysis and compared the results for genes common to both datasets. For 75 genes, reported in the literature to be imprinted, we had at least one well-covered heterozygous SNP in both our tar-RNAseq and wnt-RNAseq datasets. In support of our hypothesis, we show that tar-RNAseq covered many more informative SNPs and greatly improved the SNP read depth, which allowed us to measure ASE at more sites in more genes. We also obtained two results that were not expected by the deeper sequencing of a targeted approach. First, tar-RNAseq produced much higher concordance of ASE between the SNPs from the same genes resulting in improved quantification of the gene-level expression bias. Second, tar-RNAseq in combination with rRNA depletion permitted efficient ASE determination from degraded RNA whereas higher RNA integrity was required for wnt-RNAseq. These improvements have enabled us to quantify the inter-individual variability of ASE in our cohort with high resolution and accuracy, which will be critical for querying associations between genomic imprinting and growth phenotypes. Our results provide performance metrics for this approach on samples collected in the

field, which can be applied to design ASE studies in other populations or species.

Materials and methods

IRB

Informed consent or assent was obtained from participants depending on whether they were adults or children. Institutional Review Boards (IRB) approval was obtained from the University of Michigan IRBMED (HUM00043670) and from La Faculté de Médecine de Pharmacie et d'Odontostomatologie (FMPOS) de Bamako in Mali (N°2016/68/CD/FMPOS).

Capture design

Sample collection, nucleic acid purification, and wnt-RNAseq were described previously (Vincenz et al. 2020). The capture region for tar-RNAseq included exonic regions for all genes with reports in the literature indicating imprinted expression or allelic methylation. The criteria for inclusion were nonstringent to avoid the exclusion of imprinted genes at the cost of including some nonimprinted genes. Furthermore, the targeted genes included genes relevant to diseases that are of interest in this cohort ($n = 67$), and genes with consistent high placental expression in the wnt-RNAseq dataset ($n = 71$, Supplementary Table S1). ERCC spike-in controls were also targeted except for nine transcripts that spanned the expression range. A genetic variation identified in the cohort was taken into account by including 100bp capturing oligonucleotides containing the alternate allele for all SNPs spaced more than 50bp apart ($n = 4634$). The targeting regions were evaluated by the NimbleDesign Software and oligonucleotides covering 2,797,406 bases were synthesized using the Roche SeqCap RNA Developer platform (Supplementary .bed file).

Genotyping

Genotyping of F2 umbilical cord tissues ($n = 227$) and F1 saliva samples ($n = 189$) was performed with targeted DNA sequencing (Roche: SeqCap EZ Choice). The region genotyped for the tar-RNAseq samples overlapped the region of the wnt-RNAseq samples (1.4 Mb) and included more genes for a total of 3.9 Mb. The analysis presented here is limited to the 75 genes that had RNAseq data in both datasets. This subset of genes mapped to 0.58 Mb and 0.48 Mb of the regions genotyped in the tar-RNAseq and the wnt-RNAseq samples, respectively. Library preparation and hybridization captures were performed at the University of Michigan Advanced Genomics Core following the manufacturer's protocols.

Tar-RNAseq Library preparation and sequencing

The University of Michigan Advanced Genomics Core prepared KAPA RNA HyperPrep Kit (Roche KK8540) libraries or KAPA RNA HyperPrep with RiboErase Kit (Roche KK8560) libraries from 1000 ng, DNaseI digested, placental total RNA using conditions adapted to each sample's RNA quality. Initially, only samples with poorly defined rRNA bands on Agilent traces ($RIN < 2.5$) were depleted of rRNA with RiboErase prior to fragmentation. Later, samples with intermediate RNA ($RIN < 6.0$) were also processed with RiboErase as the cost was not prohibitive. Fragmentation conditions were established based on each sample's Agilent Bioanalyzer DV₂₀₀ quality metric which reflects the percentage of RNA fragments above 200 nucleotides: DV₂₀₀ < 55 at 65°C for 1 minute; 55 > DV₂₀₀ < 70 at 65°C for 4 minutes; DV₂₀₀ > 70 and RIN < 3.8 at 85°C for 4 minutes; DV₂₀₀ > 70 and

RIN > 3.8 94°C for 4 minutes. ERCC exogenous RNA controls (ThermoFisher Scientific 4456739) were included in all library preparations according to the manufacturer's guidelines. Six indexed cDNA libraries were pooled for each capture reaction totaling 1 µg of cDNA. In cases where Kapa RNA HyperPrep plus RiboErase libraries were multiplexed with Kapa RNA HyperPrep (nonrRNA-depleted) libraries, the amount of the rRNA-depleted library was adjusted to 10-fold less than nondepleted RNA libraries in these mixtures. Libraries were sequenced on an Illumina NovaSeq (S4). RNAs with RIN ≤ 3.8 and DV₂₀₀ ≤ 70 were generally selected for KAPA RNA Prep Plus RiboErase library preparations. In total, 236 RNA samples from 227 F2 individuals were sequenced.

Pyrosequencing

Allelic expression of select heterozygous SNPs was validated by pyrosequencing. cDNA synthesis by RT was performed immediately after DNaseI digestion of placental RNA with the ProtoScript[®]II First Strand cDNA Synthesis Kit (E6560, New England Biolabs) and random hexamer primers. Qiagen PyroMark Assay Design 2.0 software was used to design pyrosequencing primers and amplicons were generated with PyroMark PCR Kit (978705, Qiagen) and sequenced using a PyroMark Q96 MD workstation. Nineteen SNPs were pyrosequenced in 3 placentas. The Pearson correlation coefficient for major allele frequency between RNAseq and pyrosequencing was 0.98 ($P = 9.6 \times 10^{-19}$).

DNA sequencing analysis

Illumina adapter contamination and read ends with base quality <20 were removed using Trimmomatic (Bolger et al. 2014). Reads shorter than 36 nt after trimming were discarded. Trimmed reads were aligned to hg38 reference genome using BWA (Li and Durbin 2009). Read deduping and base quality score recalibration were performed using MarkDuplicates and BaseRecalibrator, respectively, from GATK (DePristo et al. 2011; Van der Auwera et al. 2013). SNPs and short INDELs were called using HaplotypeCaller, GenomicsDBImport, and GenotypeGVCFs from GATK. Resulting variants underwent GATK-recommended hard-filtering for SNPs and INDELs separately. Furthermore, we applied a series of filters in order to remove less-confident genotypes that included the following: (1) variants with genotyping quality <20 or total read depth <20; (2) variants falling in the regions with 100mer-alignability score <1 using the Umap multi-read mappability track (Karimzadeh et al. et al. 2018); (3) variants falling in the ENCODE Blacklist regions (Amemiya et al. 2019) or the genomic SuperDups regions (Bailey et al. 2002); (4) variants with known alternate allele mapping bias identified in a previous study (Panousis et al. 2014; Castel et al. 2015); (5) variants that had more than one alternate allele; (6) heterozygous SNPs whose reference allele frequency was <0.2 or >0.8; (7) homozygous SNPs whose reference allele frequency was >0.05; (8) homozygous reference sites whose reference allele frequency was <0.95; (9) SNPs where >5% of reads supported an allele that was neither reference nor alternate; (10) SNPs exhibiting excess heterozygosity (GATK-calculated metrics ExcessHet >54.69); and (11) SNPs having a nearby INDEL within 150 bp. PhaseByTransmission in GATK was used to phase the variants in a subset of the samples (45%) where both parents were genotyped. The phased variants were filtered by requiring the transmission probability score to be no lower than 20, and then combined with the variants phased by HaplotypeCaller. Eight F2 samples were excluded from the final phasing results due to excessive Mendelian violations indicative of nonpaternity or tube error.

RNA sequencing analysis

Previously published wht-RNAseq data (Vincenz et al. 2020) was used with the allelic read counts recalculated using deduped alignments, and the same workflow was also used for the analysis of the tar-RNAseq data. Illumina adapter contamination and read ends with base quality <20 were removed using Trimmomatic. Reads shorter than 36 nt after trimming were discarded. HISAT2 (Kim et al. 2015) was used to first build a new reference for each individual to incorporate the genomic variants identified from the corresponding DNA sample, and second to align the paired trimmed reads onto this reference with splice sites from GENCODE GTF (Harrow et al. 2012). Alignments were filtered and deduped using WASP (Van De Geijn et al. 2015) to reduce biases. Properly paired alignments with the highest mapping quality were selected as confident alignments and used for downstream analyses.

StringTie (Pertea et al. 2016) was used to quantify the relative expression at the transcript level. Alignments were split into sense-strand and antisense-strand alignments. ASEReadCounter from GATK was used to calculate allele-specific RNA read depth in both strands at each heterozygous SNP of the paired DNA sample. SNPs were annotated with the coordinates of the exons to which they mapped and overlapping exons in the same gene were merged into one interval. SNPs covered by at least 10 reads and mapped to unique genes and transcripts expressed at >0.1 TPM in a placental reference RNA-seq dataset were retained (Majewska et al. 2017). The SNP level imprinting codes were generated after considering all genes affected by the SNP using VEP, Variant Effect Predictor (McLaren et al. 2016). Targeted enrichment was measured as one minus the off-target aligned base ratio computed by CollectHsMetrics in GATK. Maternal contamination was assessed and removed as previously described (Vincenz et al. 2020).

For the comparison between tar-RNAseq and wht-RNAseq, we only used the SNPs in the genes that carried at least one SNP in at least one sample in both datasets and only paternally expressed (PEGs), maternally expressed (MEGs), and complexly expressed (CEGs) genes were considered, which limited the comparison to 75 genes (Supplementary Table S1, column D). SNP-level Pat-Freq was calculated as the ratio of the paternal allele read count to the total read count. For gene-level Pat-Freq, we summed the paternal allele read counts and total read counts from all the SNPs of the gene and calculated their ratio. To determine the ASE correlations between SNPs at the gene or exon level, Pearson correlation coefficients were calculated across all pairwise combinations of SNPs mapping to the same gene or exon.

RNAseq library preparation cost evaluation

A comparison between the cost of wht-RNAseq and tar-RNAseq (Supplementary Table S3) on a per sample basis was made to weigh increases in sequencing depth and coverage against costs of adding a target capture step to the library preparation method. Library preparation service costs reflect pricing as of May 7, 2020 at the University of Michigan Advanced Genomics Core. All other reagents reflect pricing at the time of purchase. The custom Roche SeqCap Target Enrichment System employed for tar-RNAseq has been discontinued but similar products are currently offered by multiple suppliers (e.g., agilent.com, arborbiosci.com, illumina.com, idtdna.com, qiagen.com, thermofisher.com, twist-bioscience.com).

Data availability

The data for the wht-RNAseq study are registered in dbGap as “Placental Transcriptome and Stunting.” The wht-RNAseq data and the corresponding genotypes obtained through targeted sequencing, the FASTQ files with the sequences from RNAseq, and the SNP level file with the allele-specific counts were deposited in dbGaP as phs001782.v1.p1. The person level nonmolecular data are available at the same site. The analogous data for the tar-RNAseq will be made available through dbGaP once the manuscript has been accepted.

Supplementary material is available at *figshare*: <https://doi.org/10.25387/g3.12251810>.

Results

Congruent ASE values from targeted RNA-seq and whole transcriptome RNA-seq

We compared allelic count distributions in tar-RNAseq and wht-RNAseq datasets for a set of 75 genes that have been reported to be imprinted in the literature and that contained at least one SNP in both datasets. The SNPs in these datasets were derived from 227 and 40 term placentas in tar-RNAseq and wht-RNAseq, respectively. For both datasets, exonic SNPs were identified by DNA-seq of umbilical cord tissue. The total size of the regions genotyped for the tar-RNAseq samples was 3.9 Mb yielding 3647 exonic SNPs, and for the wht-RNAseq samples was 1.4 Mb yielding 2517 exonic SNPs that mapped to the 75 genes of interest (Materials and Methods). SNPs were annotated as in our earlier study (Vincenz *et al.* 2020) as PEGs or MEGs based on the parental bias reported in the literature, and as CEGs, for genes with complex imprinting patterns or conflicting literature data.

To verify the parent of origin for the expression bias, we phased the SNPs in the subset of samples for which parental genotypes were known and calculated paternal allele frequency (Pat-Freq) (Figure 1A). By combining transmission- and read-based phasing, we were able to phase on average 265 and 68 SNPs per sample in 110 and 28 samples from tar-RNAseq and wht-RNAseq, respectively. In both datasets, the distributions of Pat-Freq in all three categories agreed with the previously reported imprinting directions (Figure 1A). The interquartile range of Pat-Freq in CEGs was smaller for tar-RNAseq than for wht-RNAseq. To determine allelic bias for all data, phased and unphased, we

calculated ASE as $|0.5 - (\text{Reference reads}/\text{Total reads})|$ (Figure 1B) (Castel *et al.* 2015). Both tar-RNAseq and wht-RNAseq data showed the strongest allelic bias for PEGs, reduced allelic bias for MEGs, and close to biallelic expression for many CEGs. The ASE distribution in the tar-RNAseq data exhibited more biallelic expression in all groups (Figure 1B).

The agreement between datasets is further illustrated by the correlation of gene-level Pat-Freq between tar-RNAseq and wht-RNAseq data for the placenta that was assayed both ways (Figure 2). The mean Pearson correlation coefficient was 0.95 overall, 0.90 for PEGs, 1.00 for MEGs, and 0.58 for CEGs. Thus, the correlation between the two datasets for gene-level Pat-Freq was strong.

RNA preparations from the same samples were processed with or without RiboErase to assess the reproducibility of ASE estimates with our workflows. The correlation of Ref-Freq between the samples in a pair was high ($0.96, P < 2 \times 10^{-303}, n = 2$) (Supplementary Figure S4). Thus, similar to the well-documented high reproducibility of whole Exome targeting (Cherukuri *et al.* 2015), tar-RNAseq delivers repeatable ASE estimates and tolerates the inclusion of RiboErase. Furthermore, validation of a subset of 19 SNPs in 3 samples revealed high correlation of Ref-Freq

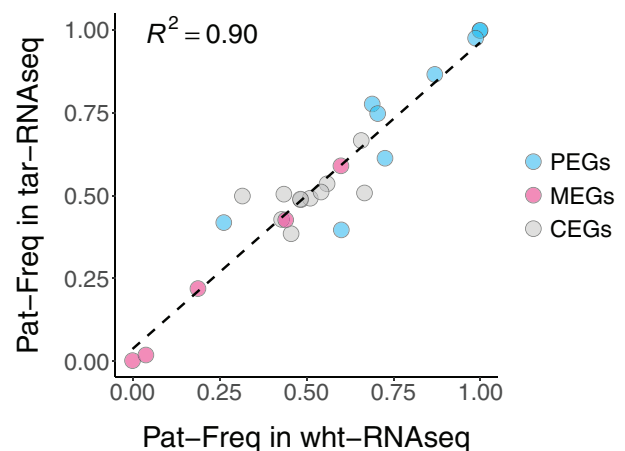


Figure 2 Correlation of gene-level average Pat-Freq between wht-RNAseq and tar-RNAseq. The scatter plot shows the Pat-Freq in each gene that had data in a sample sequenced with both wht-RNAseq and tar-RNAseq. PEGs, MEGs, and CEGs are denoted by blue, pink, and gray, respectively. The linear regression is shown by a dashed line.

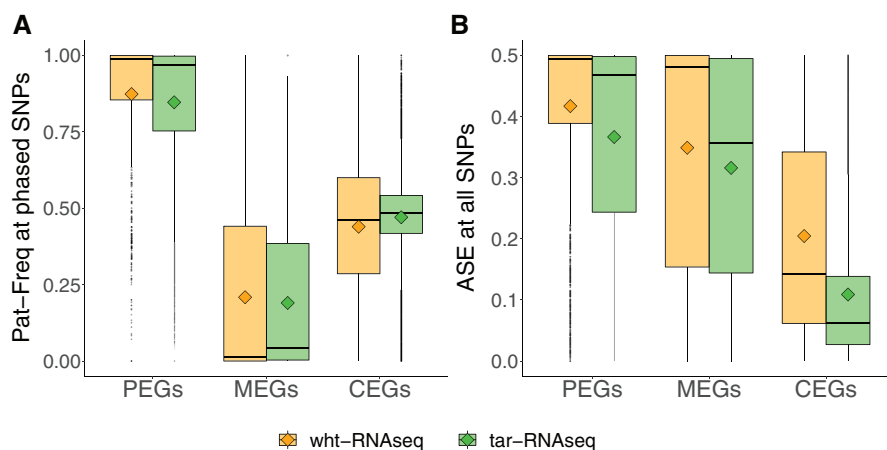


Figure 1 Distribution of parent of origin expression and ASE in PEGs, MEGs, and CEGs. The box plots show the distributions of (A) Pat-Freq and (B) ASE at the SNPs in PEGs, MEGs, or CEGs. The mean of each distribution is indicated by a diamond. Orange and green colors denote wht-RNAseq and tar-RNAseq, respectively.

between ASE determined by pyrosequencing and tar-RNAseq ($r = 0.98$, $P < 1 \times 10^{-18}$), similar to what is observed using wht-RNAseq in high RIN samples in mice (0.91 , $P < 1 \times 10^{-16}$) (Perez et al., 2015).

Targeted RNA-seq improved SNP coverage

Even though ASE measurements between the two datasets were congruent, we obtained substantial improvement in SNP coverage from tar-RNAseq compared to wht-RNAseq. After removing the SNPs that had fewer than 10 total read counts, wht-RNAseq was able to cover only 79 (or 20%) of the SNPs, on average, across the samples, while tar-RNAseq covered up to 337 (or 80%) of the SNPs (Figure 3A) in the 75 genes that were common between the two datasets. We observed the same pattern when analyzing all SNPs in each dataset (Supplementary Figure S1). The improvement in coverage was achieved for tar-RNAseq with a mean of only 80×10^6 reads per sample (SD 52×10^6)—far less than the mean number of reads per sample of 269×10^6 for wht-RNAseq (SD 110×10^6). Stated in terms of the number of total SNPs covered per billion bases sequenced, tar-RNAseq produced 54 SNPs/ 10^9 bases while wht-RNAseq delivered only 2 SNPs/ 10^9 bases. This improved use of sequencing resources was expected from the enrichment of the RNA fragments of interest. In our tar-RNAseq dataset, the percentage of bases that aligned to the targeted region was 85%, on average, indicative of successful enrichment (Materials and Methods). In principle, increasing sequencing depth could overcome the coverage deficits of wht-

RNAseq. We calculate that the per sample cost would increase 35-fold, which for most projects is prohibitive especially in the context of a population study.

While both datasets had at least one SNP in the 75 genes analyzed here, each gene had informative SNPs in more samples in tar-RNAseq than in wht-RNAseq. Specifically, after quality filtration, 24 genes per sample had at least one SNP in wht-RNAseq with each gene carrying, on average, 3 SNPs; these values increased to 55 genes with 6 SNPs per gene in tar-RNAseq (Figure 3, B and C). In addition, the average total read count for the final SNPs in tar-RNAseq was four times higher than in wht-RNAseq (Figure 3D). In sum, more informative SNPs were obtained by tar-RNAseq than wht-RNAseq, which allowed us to measure the inter-individual variability of ASE in more genes.

Arguably, the comparison could be confounded by the fact that different samples were analyzed between our tar-RNAseq and wht-RNAseq datasets as only one individual was sequenced by both technologies. However, it is unlikely that the difference in the samples between the two datasets was responsible for the large difference in the SNP coverage described above. Three sample-level QC metrics that could contribute to SNP coverage are number of genotyped SNPs, RIN, and maternal contamination. We show that RIN and maternal contamination did not differ significantly between the two datasets (Supplementary Table S2). We show that the choice of sequencing approach is more important than the number of genotyped SNPs through linear regressions in which the dependent variables were five different

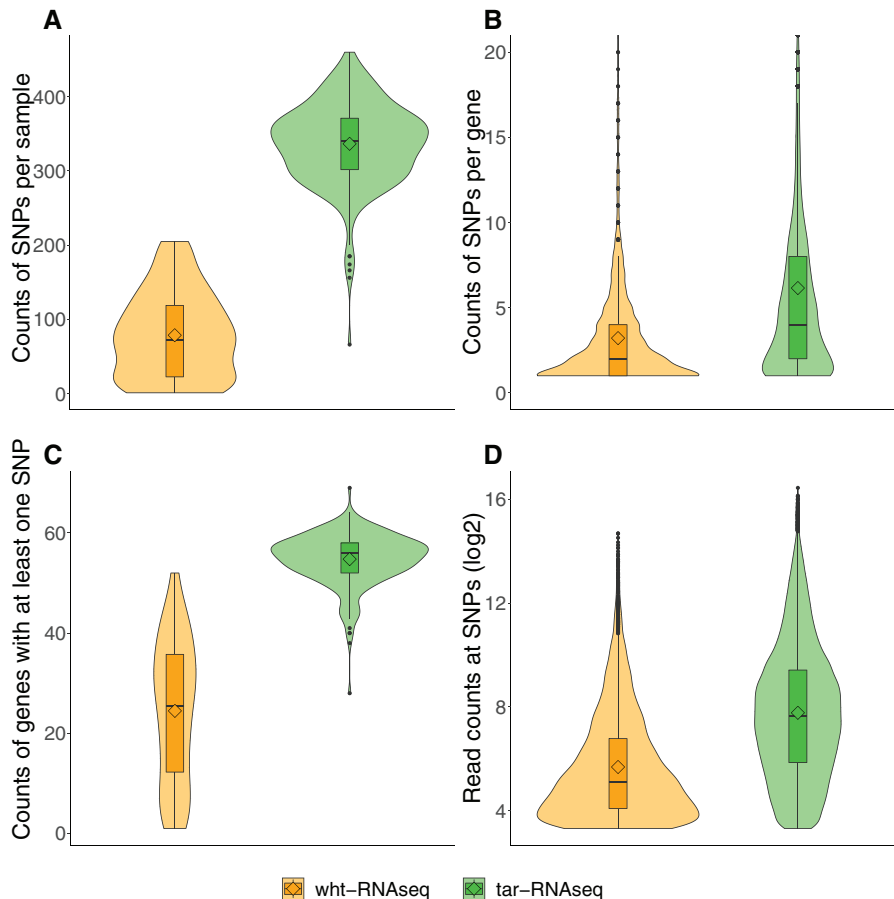


Figure 3 Comparison of SNP coverage between wht-RNAseq and tar-RNAseq. The violin plots show the distributions of (A) the counts of SNPs per sample, (B) the counts of SNPs per gene, (C) the counts of genes with at least one SNP, and (D) the read counts at SNPs. Orange and green colors denote wht-RNAseq and tar-RNAseq, respectively.

measures relevant to SNP coverage and the independent variables were tar-RNAseq (vs wht-RNAseq) and counts of genotyped exonic hetSNPs per sample (Supplementary Figure S3). Tar-RNAseq yielded a huge improvement in SNP coverage relative to wht-RNAseq at all observed numbers of genotyped SNPs. Moreover, even in the samples that had about 500 hetSNPs using wht-RNAseq, the SNP coverage was lower than in the samples that had about 300 hetSNPs using tar-RNAseq (Supplementary Figure S3). The mean number of genotyped SNPs was 418 for tar-RNAseq and 386 for wht-RNAseq, but evidently this difference could not underlie the improvement in SNP coverage using tar-RNAseq. In sum, our findings are not sensitive to the difference in the samples used in the two sequencing approaches.

Targeted RNA-seq improved concordance of ASE from the same genes

Low SNP read coverage can limit the concordance of ASE between the SNPs from the same gene (Zou et al. 2019). To determine the relationship between read coverage and concordance of ASE in our data, we calculated Pearson correlation coefficients for the pairwise combinations of SNPs mapping to the same gene in each sample. The mean correlation coefficient was 0.54 for wht-RNAseq and 0.90 for tar-RNAseq. At every read depth threshold, including the highest, the concordance was always much stronger in tar-RNAseq than in wht-RNAseq data (Figure 4). The concordance of the SNPs from the same exon showed the same pattern (Supplementary Figure S2). In wht-RNAseq, deduping improved SNP concordance but not to the level observed with deduped tar-RNAseq. Thus, hybridization capture improved this variable well beyond what would be expected from the increase in sequencing depth alone.

Targeted RNA-seq in combination with rRNA depletion permitted assessment of ASE even in degraded samples

RNA degradation contributed to the reduced SNP coverage in the wht-RNAseq samples and inefficient rRNA removal is a factor known to interfere with the complexity of sequencing libraries (Stark et al. 2019). In tar-RNAseq, ribosomal RNA should, in principle, have been removed by the hybridization reaction. However, we were able to rescue samples having low RIN by using rRNA depletion to improve SNP coverage in tar-RNAseq. Thus, we observed a strong positive correlation between RIN and SNP coverage fraction in the wht-RNAseq but not in the tar-RNAseq data (Figure 5). The combination of tar-RNAseq and rRNA depletion routinely produced high SNP coverage in samples with substantial RNA degradation ($DV_{200} \sim 50\%$).

Targeted RNA-seq improved the measurement of relative expression

Although the focus of our efforts was on ASE and not relative expression (Perez et al. 2015), we compared the performance of the two approaches with particular regard to degraded samples. As expected, compared with wht-RNAseq, tar-RNAseq yielded higher TPM values as the genes of interest constituted a larger fraction of the total sequenced transcripts (Supplementary Figure S5). Expression measurements were more consistent for low RIN samples from tar-RNAseq than for wht-RNAseq (Supplementary Figures S5A and S6B). Below a RIN of 3, many genes were expressed at TPM close to 0 in the wht-RNAseq dataset—visible as a transition in the heat maps of gene expression versus RIN (Supplementary Figure S5, A and B). No such transition was seen for the tar-RNAseq heat map and the relative expression of

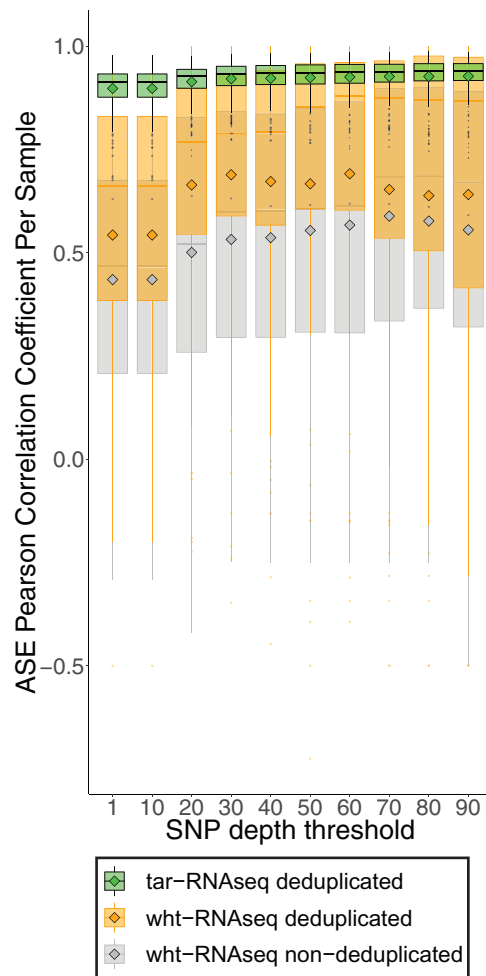


Figure 4 Concordance of ASE for SNPs in the same gene. The Pearson correlation coefficients were calculated from pairwise combinations of the SNPs from the same genes and their distributions are shown in box plots, stratified by escalated depth filtering thresholds. The mean of each distribution is indicated by a diamond. Data from tar-RNAseq (with deduping), wht-RNAseq (with deduping), and wht-RNAseq (without deduping) are denoted by green, orange, and gray colors, respectively.

highly or lowly expressed genes was not influenced by RNA integrity. These observations indicate that tar-RNAseq produced more reliable measurements of expression and was less affected by sample quality.

Discussion

ASE analysis has been performed on a variety of wht-RNAseq datasets including simulated sequences (Raghupathy et al. 2018), RNA from cells cultured *in vitro* (Lappalainen et al. 2013 2013; Gutierrez-Arcelus et al. 2013), and RNA from inbred mice (Perez et al. 2015). In humans, ASE analysis of RNA from many tissues was performed as part of the GTEx project (Babak et al. 2015; Baran et al. 2015) (<https://gtexportal.org/home/>). Placentas were not included in GTEx, but human placental tissue has been analyzed by other groups using wht-RNAseq and analyzed for ASE (Hamada et al. 2016; Hanna et al. 2016; Pilvar et al. 2019). The goal of these studies in regard to ASE was to identify imprinted genes through a transcriptome wide approach and to categorize them by imprinting status. The ENCODE study aimed to go beyond categorization and pursued a more quantitative approach that entailed calculation of the significance values for the parent of

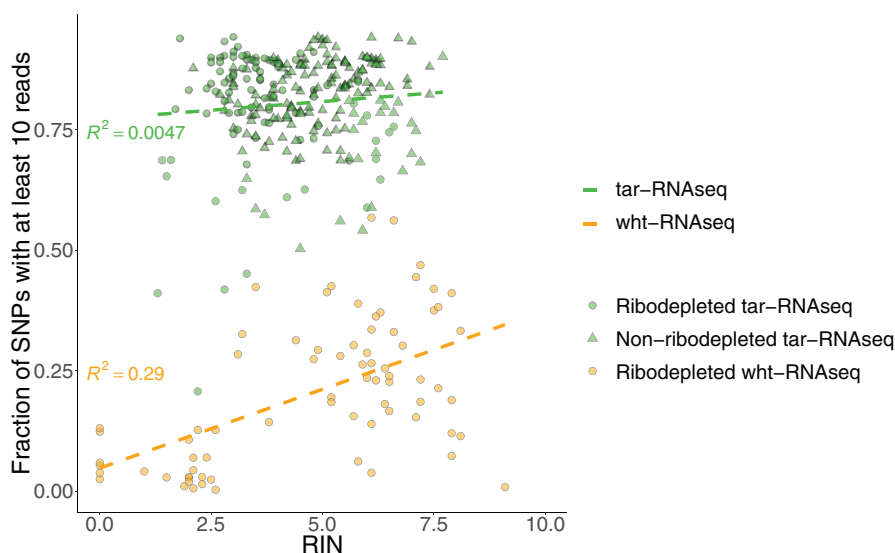


Figure 5 Combination of tar-RNAseq and RiboErase rescued most low RIN samples. The scatter plots show the fraction of SNPs covered by no fewer than 10 reads versus RIN. The ribodepleted tar-RNAseq, nonribodepleted tar-RNAseq, and ribodepleted wht-RNAseq samples are denoted by green circles, green triangles, and orange circles, respectively. The linear regression for tar-RNAseq and wht-RNAseq data points is shown by a green and an orange dashed line, respectively.

origin effect for individual SNPs (Zink et al. 2018). In contrast with the foregoing studies, our goal was to generate quantitative gene-level ASE estimates with high precision and accuracy. Such estimates are required for investigation of the functional significance of inter-individual variation in ASE (Vincenz et al., 2020).

Toward that end, we compared allelic count distributions in tar-RNAseq and wht-RNAseq datasets for a set of 75 genes that had been reported in the literature to be imprinted and that had at least one SNP in both datasets. We found that the two methods produced similar allelic expression biases. However, wht-RNAseq was able to cover only 20% of the SNPs, on average, across the samples, whereas tar-RNAseq covered 80% of the SNPs, with improvements in SNP coverage of 27-fold per billion bases sequenced. In humans, the paucity of SNPs makes it imperative to cover all the SNPs in the genes of interest. Using tar-RNAseq, we were able to obtain sufficient coverage at four times as many SNPs in twice as many genes in a sample, on average. Moreover, the mean number of SNPs per gene doubled, and the mean read depth per SNP increased fourfold, without increasing library preparation costs, making tar-RNAseq more cost effective (Supplementary Table S3). A complete dataset would have sufficient reads at every SNP in every person, which is a goal that was more closely achieved by tar-RNA seq than by wht-RNAseq. Having a richer dataset will enable us to determine the interindividual variation in ASE for more genes across more individuals, so that we can better query the association between genomic imprinting and growth phenotypes in our cohort study.

Maternal contamination is a potential confounder unique to placental tissue and is a limiting factor in molecular analyses (Konwar et al. 2019). The degree of contamination can be directly determined from the RNAseq data by quantitating nonfetal alleles (Hamada et al. 2016). The greater SNP sequencing coverage and depth of SNPs in tar-RNAseq enabled us to quantify maternal contamination for each gene in each placenta with greater sensitivity.

Importantly, targeted RNA-seq had some additional nonanticipated benefits. Gene-level ASE estimates are imprecise due to poor concordance of SNP-level ASE over a gene. Some

discordance between the SNPs can be due to differing imprinting status between the transcript variants from the same gene (Perez et al. 2015), but poorly identified technical variables also contribute (Babak et al. 2015; Baran et al. 2015). It is known that selecting SNPs with higher sequencing coverage leads to improved concordance (Zou et al. 2019), which we also saw in wht-RNAseq. Removing PCR duplicates further improved concordance but not nearly to the levels achieved by tar-RNAseq. In contrast with relative expression analyses, ASE is only based on the read count ratio between the alleles and removing read duplicates reduces the technical noise. The greatly improved concordance of SNPs strongly argues in favor of using tar-RNAseq for applications that require accurate gene-level ASE estimates. In future efforts, it might be possible to gain additional power for SNPs with low read depth by using unique molecular indices (UMI) (Islam et al. 2014) in conjunction with tar-RNAseq as there are reports that PCR duplicates can affect ASE quantification in such circumstances (Castel et al. 2015).

Degraded RNA is found in many human samples, including term placentas, post-mortem samples of stomach and kidney, and formalin-fixed paraffin-embedded (FFPE) samples (Walker et al. 2016; Zhang et al. 2017; Konwar et al. 2019) (<https://gtexportal.org/home/>). We were successful in combining random primed library preparation with rRNA depletion to generate libraries from degraded RNA for tar-RNAseq. Little to no loss of coverage was observed with degraded RNA from most samples, and we salvaged samples with DV_{200} as low as 50%. The critical role for efficient rRNA depletion in preparing libraries from degraded RNA is well known (Stark et al. 2019) and is in part due to the inability to target the poly(A) tail in degraded samples. The hybridization reaction with the capturing oligos should, in principle, be sufficient to remove rRNA. However, our results show that in degraded samples, removal of rRNA prior to library preparation improved SNP coverage. Improved data quality has previously been reported for gene expression analysis in FFPE samples with RNAseq and rRNA depletion (Zhao et al. 2014) or when capturing the whole exome (Pennock et al. 2019). However, to our knowledge, our study is the first to report the unexpected synergy between

rRNA depletion and tar-RNAseq. It was also more efficient and cost effective to focus on a targeted region of only 4 Mb instead of the human exome of 64 Mb (the total length of Roche SeqCap EZ Exome Probes). Tar-RNAseq also improved relative expression estimates for degraded samples. Importantly, we document that cost savings is only one of the advantages of tar-RNAseq and other synergies may become the predominant motive to use this technology as sequencing costs continue to fall.

A strength of our study is that the placenta samples were collected from healthy women who were of similar ages and belonged to the same cohort and ethnicity, using a standardized protocol. Moreover, we compared the same 75 genes using both wht-RNAseq and tar-RNAseq. A limitation of our study is that only one sample was sequenced using both methods. However, we examined three parameters that could potentially differ between samples and influence data yield (RIN, maternal contamination, number of genotyped SNPs) and showed that our conclusion about the superiority of tar-RNAseq was not sensitive to any of these parameters. We also note that library preparation reagents for wht-RNAseq libraries and tar-RNAseq libraries were from different manufacturers. Although we did not try to estimate ASE in relation to the cellular composition of the fetal compartment of the placenta, a recent single cell study showed that placental samples collected using a protocol similar to ours were mostly comprised of trophoblast and syncytiotrophoblast cells of the fetus (Yuan *et al.* 2020).

In conclusion, we compared tar-RNAseq and wht-RNAseq in a study of ASE in 75 known imprinted genes in placental tissue collected from a healthy human cohort. Tar-RNAseq covered more SNPs of interest and at greater depth. In previous research on humans, discordant ASE values for SNPs of the same gene have limited the ability to accurately quantify ASE. We show that Tar-RNAseq improved the reliability of ASE detection by greatly increasing the concordance of ASE measurement between the SNPs from the same gene. In combination with rRNA depletion, tar-RNAseq performed well even in cases of degraded RNA. The advantages of tar-RNAseq go beyond the savings on sequencing costs alone and include higher accuracy in ASE estimates in samples with varying RNA quality, as is typical for field collections. Targeted sequencing will benefit the study of associations between individual variation in ASE and phenotypes in humans or in other species where growth phenotypes are of interest, such as domesticated animals. The data we presented here originated from field samples and provide metrics to inform the design of such projects.

Acknowledgments

The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health, the John Templeton Foundation, or the National Science Foundation. The authors thank the study participants who made this research possible as well as three Malian gynecologists for their advice and helpful suggestions: Prof. Amadou Dolo, Prof. Niani Mounkoro, and Prof. Mamadou Traoré. They also thank our field manager, Zachary Dolo, and the Malian medical team: Dr. Gouro Dicko, Dr. Akoro Dolo, Madeleine Goita, Aissa Dolo, Younus Dolo, Jeremy Sagara, and Safoura Guindo. For permission to carry out this study in Mali, they are grateful to the

Centre National de la Recherche Scientifique et Technologique and the Comité d'Éthique de la Faculté de Médecine de Pharmacie et d'Onto-Stomatologie of the University of Sciences, Techniques, and Technologies of Bamako (authorization N°2016/68/CD/FMPOS). We also acknowledge the contribution of the Advanced Genomics Core at the University of Michigan.

Funding

This research was supported by the Eunice Kennedy Shriver National Institute of Child Health & Human Development of the National Institutes of Health (R01HD088521 and R21HD077465 to B.I.S.); the John Templeton Foundation (52269 to B.I.S.); and the National Science Foundation program in Biological Anthropology (NSF BCS-1354814 to B.I.S.).

Conflicts of interest

None declared.

Literature cited

- Amemiya HM, Kundaje A, Boyle AP. 2019. The ENCODE blacklist: identification of problematic regions of the Genome. *Sci Rep.* 9: 9354.
- Babak T, DeVeale B, Tsang EK, Zhou Y, Li X, *et al.* 2015. Genetic conflict reflected in tissue-specific maps of genomic imprinting in human and mouse. *Nat Genet.* 47:544–549.
- Bailey JA, Gu Z, Clark RA, Reinert K, Samonte RV, *et al.* 2002. Recent segmental duplications in the human genome. *Science.* 297: 1003–1007.
- Baran Y, Subramaniam M, Biton A, Tukiainen T, Tsang EK, *et al.*; GTEx Consortium. 2015. The landscape of genomic imprinting across diverse adult human tissues. *Genome Res.* 25:927–936.
- Batista RA, Köhler C. 2020. Genomic imprinting in plants-revisiting existing models. *Genes Dev.* 34:24–36. [10.1101/gad.332924]
- Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics.* 30: 2114–2120.
- Castel SE, Levy-Moonshine A, Mohammadi P, Banks E, Lappalainen T. 2015. Tools and best practices for data processing in allelic expression analysis. *Genome Biol.* 16:195.
- Chen Z, Hagen DE, Wang J, Elsik CG, Ji T, *et al.* 2016. Global assessment of imprinted gene expression in the bovine conceptus by next generation sequencing. *Epigenetics.* 11:501–516.
- Cherukuri PF, Maduro V, Fuentes-Fajardo KV, Lam K, Adams DR, *et al.*; NISC Comparative Sequencing Program. 2015. Replicate exome-sequencing in a multiple-generation family: improved interpretation of next-generation sequencing data. *BMC Genomics.* 16:998.
- DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, *et al.* 2011. A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat Genet.* 43:491–501.
- DeVeale B, van der Kooy D, Babak T. 2012. Critical evaluation of imprinted gene expression by RNA-seq: a new perspective. *PLoS Genet.* 8:e1002600.
- Farhadova S, Gomez-Velazquez M, Feil R. 2019. Stability and liability of parental methylation imprints in development and disease. *Genes (Basel).* 10:999.
- Fontanillas P, Landry CR, Wittkopp PJ, Russ C, Gruber JD, *et al.* 2010. Key considerations for measuring allelic expression on a genomic

- scale using high-throughput sequencing. *Mol Ecol.* 19:212–227. doi:10.1111/j.1365-294X.2010.04472.x
- Gulyás-Kovács A, Keydar I, Xia E, Fromer M, Hoffman G, et al. 2018. Unperturbed expression bias of imprinted genes in schizophrenia. *Nat Commun.* 9:7.
- Gutierrez-Arcelus M, Lappalainen T, Montgomery SB, Buil A, Ongen H, et al. 2013. Passive and active DNA methylation and the interplay with genetic variation in gene regulation. *eLife.* 2:e00523.
- Hamada H, Okae H, Toh H, Chiba H, Hiura H, et al. 2016. Allele-specific Methylome and transcriptome analysis reveals widespread imprinting in the human placenta. *Am J Hum Genet.* 99:1045–1058.
- Hanna CW, Peñaherrera MS, Saadeh H, Andrews S, McFadden DE, et al. 2016. Pervasive polymorphic imprinted methylation in the human placenta. *Genome Res.* 26:756–767.
- Harrow J, Frankish A, Gonzalez JM, Tapanari E, Diekhans M, et al. 2012. GENCODE: the reference human genome annotation for the ENCODE project. *Genome Res.* 22:1760–1774.
- Heap GA, Yang JHM, Downes K, Healy BC, Hunt KA, et al. 2010. Genome-wide analysis of allelic expression imbalance in human primary cells by high-throughput transcriptome resequencing. *Hum Mol Genet.* 19:122–134.
- Islam S, Zeisel A, Joost S, La Manno G, Zajac P, et al. 2014. Quantitative single-cell RNA-seq with unique molecular identifiers. *Nat Methods.* 11:163–166.
- Jadhav B, Monajemi R, Galalova KK, Ho D, Draisma HHM, et al.; GoNL Consortium. 2019. RNA-Seq in 296 phased trios provides a high-resolution map of genomic imprinting. *BMC Biol.* 17:50.
- Karimzadeh M, Ernst C, Kundaje A, Hoffman MM. 2018. Umap and Bismap: quantifying genome and methylome mappability. *Nucleic Acids Res.* 46:e120. [10.1093/nar/gky677]
- Kim D, Langmead B, Salzberg SL. 2015. HISAT: a fast spliced aligner with low memory requirements. *Nat Methods.* 12:357–360.
- Konwar C, Del Gobbo G, Yuan V, Robinson WP. 2019. Considerations when processing and interpreting genomics data of the placenta. *Placenta.* 84:57–62.
- Lappalainen T, Sammeth M, Friedländer MR, 't Hoen PAC, Monlong J, et al.; Geuvadis Consortium. 2013. Transcriptome and genome sequencing uncovers functional variation in humans. *Nature.* 501:506–511.
- Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics.* 25:1754–1760.
- Li Y, Li J. 2019. Technical advances contribute to the study of genomic imprinting. *PLoS Genet.* 15:e1008151.
- Majewska M, Lipka A, Pauksztó L, Jastrzebski JP, Myszczyński K, et al. 2017. Transcriptome profile of the human placenta. *Funct Integr Genomics.* 17:551–563.
- McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GRS, et al. 2016. The ensembl variant effect predictor. *Genome Biol.* 17:122.
- Metsalu T, Viltrop T, Tiirats A, Rajashekar B, Reimann E, et al. 2014. Using RNA sequencing for identifying gene imprinting and random monoallelic expression in human placenta. *Epigenetics.* 9:1397–1409.
- Moore T, Haig D. 1991. Genomic imprinting in mammalian development: a parental tug-of-war. *Trends Genet.* 7:45–49.
- Mozaffari SV, Stein MM, Magnaye KM, Nicolae DL, Ober C. 2018. Parent of origin gene expression in a founder population identifies two new candidate imprinted genes at known imprinted regions. *PLoS One.* 13:e0203906.
- Nothnagel M, Wolf A, Herrmann A, Szafranski K, Vater I, et al. 2011. Statistical inference of allelic imbalance from transcriptome data. *Hum Mutat.* 32:98–106.
- Panousis NI, Gutierrez-Arcelus M, Dermitzakis ET, Lappalainen T. 2014. Allelic mapping bias in RNA-sequencing is not a major confounder in eQTL studies. *Genome Biol.* 15:467.
- Pennock ND, Jindal S, Horton W, Sun D, Narasimhan J, et al. 2019. RNA-seq from archival FFPE breast cancer samples: molecular pathway fidelity and novel discovery. *BMC Med Genomics.* 12:195.
- Perez JD, Rubinstein ND, Fernandez DE, Santoro SW, Needleman LA, et al. 2015. Quantitative and functional interrogation of parent-of-origin allelic expression biases in the brain. *eLife.* 4:41.
- Pertea M, Kim D, Pertea GM, Leek JT, Salzberg SL. 2016. Transcript-level expression analysis of RNA-seq experiments with HISAT, StringTie and Ballgown. *Nat Protoc.* 11:1650–1667.
- Peters J. 2014. The role of genomic imprinting in biology and disease: an expanding view. *Nat Rev Genet.* 15:517–530.
- Pilvar D, Reiman M, Pilvar A, Laan M. 2019. Parent-of-origin-specific allelic expression in the human placenta is limited to established imprinted loci and it is stably maintained across pregnancy. *Clin Epigenet.* 11:94.
- Raghupathy N, Choi K, Vincent MJ, Beane GL, Sheppard KS, et al. 2018. Hierarchical analysis of RNA-seq reads improves the accuracy of allele-specific expression. *Bioinformatics.* 34:2177–2184.
- Stark R, Grzelak M, Hadfield J. 2019. RNA sequencing: the teenage years. *Nat Rev Genet.* 20:631–656.
- Stevenson KR, Coolon JD, Wittkopp PJ. 2013. Sources of bias in measures of allele-specific expression derived from RNA-seq data aligned to a single reference genome. *BMC Genomics.* 14:536.
- Thamban T, Agarwal V, Khosla S. 2020. Role of genomic imprinting in mammalian development. *J Biosci.* 45:20.
- Tucci V, Isles AR, Kelsey G, Ferguson-Smith AC, Erice Imprinting Group. 2019. Genomic imprinting and physiological processes in Mammals. *Cell.* 176:952–965.
- Van De Geijn B, Mcvicker G, Gilad Y, Pritchard JK. 2015. WASP: Allele-specific software for robust molecular quantitative trait locus discovery. *Nat Methods.* 12:1061–1063.
- Van der Auwera GA, Carneiro MO, Hartl C, et al. 2013. From fastQ data to high-confidence variant calls: the genome analysis toolkit best practices pipeline. *Curr Protoc Bioinformatics.* 43:11.10.1–11. [10.1002/0471250953.bi1110s43]
- Vincenz C, Lovett JL, Wu W, Shedden K, Strassmann BI. 2020. Loss of imprinting in human placentas is widespread, coordinated, and predicts birth phenotypes. *Mol Biol Evol.* 37:429–441. [10.1093/molbev/msz226]
- Walker DG, Whetzel AM, Serrano G, Sue LI, Lue LF, et al. 2016. Characterization of RNA isolated from eighteen different human tissues: results from a rapid human autopsy program. *Cell Tissue Bank.* 17:361–375.
- Wang X, Clark AG. 2014. Using next-generation RNA sequencing to identify imprinted genes. *Heredity (Edinb).* 113:156–166.
- Wang X, Miller DC, Harman R, Antczak DF, Clark AG. 2013. Paternally expressed genes predominate in the placenta. *Proc Natl Acad Sci USA.* 110:10705–10710. doi:10.1073/pnas.1308998110
- Wang Y, Gao S, Zhao Y, Chen W-H, Shao J-J, et al. 2019. Allele-specific expression and alternative splicing in horse×donkey and cattle×yak hybrids. *Zool Res.* 40:293–304.
- Yuan V, Hui D, Yin Y, Peñaherrera M, Beristain A, et al. 2020. Cell-specific characterization of the placental methylome. *BMC Genomics.* 22:6. [10.21203/rs.3.rs-38223/v2]
- Zhang P, Lehmann BD, Shyr Y, Guo Y. 2017. The utilization of formalin fixed-paraffin-embedded specimens in high throughput genomic studies. *Int J Genomics.* 2017:1926304.

- Zhao W, He X, Hoadley KA, Parker JS, Hayes DN, et al. 2014. Comparison of RNA-Seq by poly (A) capture, ribosomal RNA depletion, and DNA microarray for expression profiling. *BMC Genomics*. 15:419.
- Zink F, Magnusdottir DN, Magnusson OT, Walker NJ, Morris TJ, et al. 2018. Insights into imprinting from parent-of-origin phased methylomes and transcriptomes. *Nat Genet*. 50:1542–1552.
- Zou J, Hormozdiari F, Jew B, Castel SE, Lappalainen T, et al. 2019. Leveraging allelic imbalance to refine fine-mapping for eQTL studies. *PLoS Genet*. 15:e1008481.

Communicating editor: B. J. Andrews