

Widespread RNA and DNA Sequence Differences in the Human Transcriptome

Mingyao Li,^{1*} Isabel X. Wang,^{8*} Yun Li,^{6,7} Alan Bruzel,⁸ Allison L. Richards,⁴ Jonathan M. Toung,⁵ Vivian G. Cheung^{2,3,8†}

¹Departments of Biostatistics and Epidemiology, University of Pennsylvania School of Medicine, Philadelphia, PA 19104, USA.

²Department of Genetics, University of Pennsylvania School of Medicine, Philadelphia, PA 19104, USA. ³Department of Pediatrics, University of Pennsylvania School of Medicine, Philadelphia, PA 19104, USA. ⁴Cell and Molecular Biology Graduate Program, University of Pennsylvania School of Medicine, Philadelphia, PA 19104, USA. ⁵Genomics and Computational Biology Graduate Program, University of Pennsylvania School of Medicine, Philadelphia, PA 19104, USA.

⁶Departments of Genetics, University of North Carolina School of Medicine, Chapel Hill, NC 27599, USA. ⁷Department of Biostatistics, University of North Carolina School of Medicine, Chapel Hill, NC 27599, USA. ⁸Howard Hughes Medical Institute, Chevy Chase, MD 20815, USA.

*These authors contributed equally to this work.

†To whom correspondence should be addressed. E-mail: vcheung@mail.med.upenn.edu

The transmission of information from DNA to RNA is a critical process. We compared RNA sequences from human B cells of 27 individuals to the corresponding DNA sequences from the same individuals and uncovered more than 10,000 exonic sites where the RNA sequences do not match that of the DNA. All 12 possible categories of discordances were observed. These differences were nonrandom as many sites were found in multiple individuals and in different cell types including primary skin cells and brain tissues. Using mass spectrometry, we detected peptides that are translated from the discordant RNA sequences and thus do not correspond exactly to the DNA sequences. These widespread RNA-DNA differences in the human transcriptome provide a yet unexplored aspect of genome variation.

DNA carries genetic information that is passed onto messenger RNA (mRNA) and proteins that perform cellular functions and it is assumed that the sequence of mRNA reflects that of the DNA. This assumed precision is important since mRNA serves as the template for protein synthesis. Hence, genetic studies have mostly focused on DNA sequence polymorphism as the basis of individual differences in disease susceptibility. Studies of mRNA and proteins analyze their expression and not sequence differences among individuals.

There are, however, known exceptions to the one-to-one relationship between DNA and mRNA sequences. These include errors in transcription (1, 2) and RNA-DNA differences that result from RNA editing (3–7). Errors are rare since proof-reading and repair mechanisms ensure the fidelity of transcription (8–10). RNA editing is carried out by enzymes that target mRNA post-transcriptionally: ADARs

that deaminate adenosine to inosine which is then recognized by the translation machineries as a guanosine (A-to-G), and APOBECs which edit cytidine to uridine (C-to-U). Previously, sequence comparisons and computational predictions have identified many A-to-G editing sites (6, 7, 11–13). In contrast, C-to-U changes are rare; apolipoprotein B is one of a small number of known target genes of human APOBEC1 (14, 15).

We obtained sequences of DNA and RNA samples from immortalized B-cells of 27 unrelated Centre d'Etude du Polymorphisme Humain (CEPH) (16) individuals, who are part of the International HapMap (17, 18) and the 1000 Genomes (19) Projects. When we compared the DNA and RNA sequences of the same individuals, we found 28,766 events at over 10,000 exonic sites that differ between the RNA and the corresponding DNA sequences. Each of these differences were observed in at least two individuals; many of these were seen in B-cells, as well as in primary skin cells and brain tissues from a separate set of individuals and in expressed sequence tags from cDNA libraries of various cell types. About 43% of the differences are transversions and therefore cannot be the result of typical deaminase-mediated RNA editing. By mass spectrometry, we also found peptide sequences that correspond to the RNA variant sequences, but not the DNA sequences, suggesting that the RNA forms are translated into proteins.

We compared the DNA and RNA sequences from B-cells of 27 unrelated CEPH individuals (table S1). We chose these samples because much information is available on them including dense DNA genotypes obtained using different technologies (20, 21). The genomes of B-cells from the CEPH collection are stable as evidenced by Mendelian inheritance of genetic loci that allowed the construction of

microsatellite- to SNP-based human genetic maps (20, 21). More recently the International HapMap Consortium (17, 18) obtained millions of SNP genotypes and the 1000 Genomes Project (19) sequenced the DNA of these individuals. Comparison of sequence data from these two projects showed high concordance (~99%). Here, we used the DNA genotypes and sequences from the two projects for our analyses. First, we considered sites that are monomorphic in the human genome. A monomorphic site is one where there is no evidence for sequence variation at that locus in dbSNP, the HapMap and the 1000 Genomes Projects. Different studies have analyzed these 27 and hundreds of additional individuals for DNA variants; thus if a site has not been identified as polymorphic, most likely all individuals have the same sequences at these sites. But to be certain, for these sites in the 27 individuals, we compared their DNA sequences from the 1000 Genomes Project with the sequences of the human reference genome and carried out traditional Sanger sequencing (22). To be included in our analysis, we required that each site to be covered by at least four reads in the 1000 Genomes Project and that the sequences from 1000 Genomes are the same as the reference genome. To ensure the integrity of the aliquots of B-cells we used for analyses, we carried out Sanger sequencing of their DNA and found perfect concordance of sequences with data from 1000 Genomes (thus also the reference genome sequences) (table S2). Second, we considered SNPs. For each individual, a SNP locus is included only if it is homozygous and the HapMap as well as the 1000 Genomes projects reported the same sequence. We have high confidence in those sequences since despite using different technologies (microarray-based genotyping in HapMap and high-throughput sequencing in 1000 Genomes), identical sequences were obtained in the two projects.

We sequenced the RNA of B-cells from the same 27 individuals using high-throughput sequencing technology from Illumina (23). The resulting RNA sequence reads were mapped to the Gencode genes (24) in the reference human genome. In total, we generated ~1.1 billion reads of 50 base pairs (bp) (~41 million reads and 2 Gb sequence per individual), of which ~69% of the reads mapped uniquely to the transcriptome (see Methods in (39)). To be confident of the base calls, for each individual, we focused our analysis on high quality reads (quality score ≥ 25) and sites that are covered by at least 10 uniquely mapped reads. Another study (25) had carried out RNA-Sequencing of the same individuals but at a lower coverage; at these sites we compared our sequences with those from their study, the concordance rate of the sequences is >99.5%. This is reassuring given that the samples were prepared and sequenced in different laboratories

For each of 27 individuals, we compared the mRNA sequences from B-cells with the corresponding DNA sequences (fig. S1). The comparison revealed many sites where the mRNA sequences differ from the corresponding DNA sequences of the same individual. To ensure that these are actual differences and to minimize the chance of sequencing errors, we required that at least 10% of the reads covering a site to be different from the DNA sequence and at least two individuals show the same RNA-DNA difference at the site. We call each occurrence of a difference between RNA and DNA sequences, an event, and the chromosomal location where such a difference occurs, a site. Each person can contribute an event to the site, thus there could be multiple events at a site.

Among our 27 subjects, we identified 28,766 events where the RNA sequences do not match those of the corresponding DNA sequences. These events are found in 10,210 exonic sites (table S10) in the human genome and reside in 4,741 known genes (36% of 13,214 genes that are covered by 10 or more RNA-Seq reads in at least one part of the gene, in two or more individuals). With gene orientation information in Gencode, we observed all 12 possible categories of base differences between RNA and its corresponding DNA (Fig. 1A). All 12 types of differences were found in each of the 27 samples; the relative proportion of each type is similar across individuals. There are 6,698 A-to-G events which can be the result of deamination by ADAR. There are 1,220 C-to-T differences which can also be mediated by a deaminase. However it is important to note that APOBEC1 and its complementation factor AICF that deaminate cytidine are not expressed in our B-cells (FPKM (26) ~ 0 for both genes) thus it is likely that an unknown deaminase or other mechanism is involved. Even for relatively well characterized proteins such as APOBEC1, a recent RNA-Seq study of Apobec1^{-/-} mice uncovered many previously unknown targets (27). In addition, we found 12,507 transversions (43%); which cannot result from classic deaminase-mediated editing. Since we do not know the mechanism by which these differences between RNA and DNA sequences arise, we refer to them as RNA-DNA Differences or RDD. An example of RDD is a C-to-A difference on chromosome 12 (at position 54,841,626 bp) in the myosin light chain gene *MYL6* where 16 of our subjects have C/C in their DNA but A/C in their RNA sequences. Another example is an A-to-C difference on chromosome 6 (at position 44,328,823 bp) in the gene, *HSP90AB1*, that encodes a heat shock protein, where 8 individuals who have homozygous A/A DNA genotype but have A/C in their RNA. Additional examples are shown in Table 1. These sites where RNA sequences differ from the corresponding DNA sequences appear to be non-random since the identical differences were found in multiple individuals: 8,163 (80.0%) of the sites were found in at least 50% of the informative

individuals (i.e. with RNA-Seq coverage ≥ 10 and DNA-Seq coverage ≥ 4 at the site). Some sites were found in all or nearly all informative individuals. For example, the DNA sequences of all 19 informative individuals at position 49,369,615 bp of chromosome 3 in the *GPXI* gene are G/G whereas their RNA sequences are G/A. (The remaining individuals were not included because available data did not meet our inclusion criteria; although the data suggest the same RDD in all remaining individuals: G/G in DNA, and G/A in RNA).

Computational and experimental validations also upheld these observed RNA-DNA differences. First, for 120 sites (10 sites per RDD type; randomly selected and all examples cited in this paper; see Table 1 and table S3), we looked for evidence of RDD in the human EST database by BLAST alignment (28) and manual inspection of each result. For 81 of the 120 sites, we found EST clones that contain the RDD alleles. The numbers of sites found in human ESTs are similar across different RDD types (average 67.5%, range: 60 to 90%). Second, we examined previously identified A-to-G editing sites (6). Fourteen of the A-to-G sites that we identified were found in their data despite the fact that different cell types were studied. Even the levels of editing at these sites are similar between the two studies (see fig. S2). Twelve additional sites were found in both studies but were filtered because they did not meet our selection criteria.

Next, we validated our findings experimentally by Sanger sequencing of both DNA and RNA at 12 randomly selected sites in B-cells (2 to 9 individuals/ site), primary skin (foreskin; 8 to 10 individuals/ site) and brain cortex (6 to 10 individuals/ site). We regrew the B-cells from our subjects and extracted DNA and mRNA from the same aliquots of cells. From sequencing the paired DNA and RNA samples and analysis of each chromatogram by two individuals independently, we confirmed 57 events in 11 sites (see Table 2, fig. S3). In *EIF2AK2*, in all the 8 individuals whose samples were sequenced, three sites were found within 10 nucleotides (see below). RDD was not found in one site in *NDUFC2*. Sanger sequencing is not very sensitive or quantitative thus we do not expect to validate all sites, especially those with low levels of RDD.

To assess if RDD shows cell type specificity, we looked for evidence of RNA-DNA sequence differences using primary human cells. We studied the same sites as above by Sanger sequencing of DNA and RNA samples from primary skin fibroblasts and brain (cortex) of a separate set of normal individuals (for each site, we examined the DNA and RNA of 6 to 10 samples per cell type). We identified 55 RDD events in primary skin cells and 62 events in brain cortex (Table 2). The results suggest that most sites are shared across cell types (Table 2); although there are exceptions, for example, an A-to-G difference in *EIF2AK2* (chr2: 37,181,512) which was

only found in B-cells and brain cortex but not in primary skin cells. We also queried the EST database for evidence of RDD (Table 1, table S3). The RNA alleles are seen in a wide range of tissues from embryonic stem cells to brain and testis; they are also found in tumors such as lung carcinoma and neuroblastoma.

Validation at the sequence level is important but does not address all concerns such as the difficulty in aligning sequences that are highly similar and errors introduced by enzymes in reverse transcription steps. We believe that such artifacts are unlikely considering the consistent patterns across sequencing methods and the fact that we observed all 12 types of nucleotide differences. An alternate and independent validation would be to ask whether the RNA variants in RDD sites are translated to proteins. To do so, first we searched mass spectrometry data from human ovarian cancer cells (29) and leukemic cells for putative RDD sites. Since the levels of most RDDs are less than 100%, both DNA and the RDD-forms of the mRNAs should be available to be translated (from here-on, we refer to mRNAs that correspond identically to the DNA sequences as DNA-forms and those that contain a RDD as RNA-forms). In the ovarian cancer and leukemic cells, we indeed found examples of proteins with peptides encoded by both DNA and RNA forms of mRNA (table S4). Encouraged by the search results and cognizant of possible genome instability and thus DNA mutations in cancer cells, we carried out mass spectrometry analysis of our B-cells.

We analyzed the proteome of our B-cells using liquid chromatography-tandem mass spectrometry and detected peptides for 3,217 proteins. Despite advances in mass spectrometry, far less than 50% of peptides can be detected in most studies (30, 31). We identified 327 peptides that cover RDD sites: 299 of them are encoded by the DNA-forms and 28 by RNA-forms of RDD containing mRNAs (FDR < 1%; tables S5 and S9). For 17 RDD sites, peptides that correspond to both DNA and RNA forms were identified (Table 3). By BLAST alignment, we ensured that these 28 peptides are unique to the genes that contain the RDD sites. In addition, we sequenced the DNA of the B-cells used for mass spectrometry and validated that the DNA sequences are the same as the reference genome but differ from the RNA sequences and thus do not encode the RNA-forms of the peptides (table S2). It is easier to detect more abundant proteins by mass spectrometry; for most RDD sites, the unaltered DNA forms are more abundant than variant RNA forms of mRNA (see below) thus it is not surprising to find more peptides that corresponds to the DNA rather than the RNA sequences. However, the counts of peptides corresponding to the DNA and RNA forms of RDD sites should not be taken as a measure of the proportions of DNA versus RNA forms of mRNA that are translated because

differences in the amino acid sequences of the DNA and RNA forms of the peptides affect the ability of mass spectrometry to detect them. In addition, when a peptide is not detected, it does not mean that it is absent from the proteome, it could be a result of sampling.

The proteomic data provide an independent validation that mRNA sequences are not always identical to DNA sequences, and demonstrate that RNA-forms of genes are translated to proteins. They also show that there are peptides in human cells that are not exactly encoded by the DNA sequences. An example of a protein variant that results from RDD is RPL28 (T-to-A, chr19: 60,590,467). The RDD led to a loss of a STOP codon. We identified peptides corresponding to the 55 amino acid extension of RPL28 protein in the ovarian cancer cells and in our B-cells (Fig. 2). Previously identified cases of RNA editing leading to proteins not encoded by genomic DNA, such as the apolipoprotein B (3, 4), serotonin and glutamate receptors (32–34) in humans and plant ribosomal protein S12 (35), also support our hypothesis that RDD leads to protein isoforms that do not correspond to the DNA sequences of the encoding genes.

Using our selection criteria, we found that in each person among the Gencode genes, there are on average 1,065 exonic events that differ in the RNA and DNA sequences. But the number of events varied among individuals (range: 282 to 1,863) by up to 6-fold across our 27 subjects (Fig. 1B). The degree of sequence coverage and sequencing errors in DNA or RNA samples do not explain these individual differences (39). Thus there is likely a biological basis for the individual variation in the number of editing and RDD events. We found no significant correlation between *ADAR* expression with the number of RDDs or the numbers of A-to-G events ($P > 0.5$). Thus, either *ADAR* expression does not affect the number of editing or RDD events, or our sample size is not sufficient to detect the correlation.

The 10,210 sites which showed RNA and DNA sequence differences are not evenly distributed across the genome: chromosome 19 has the most whereas chromosome 13 has the fewest number of sites. This pattern is observed after correction for differences in size and gene density among chromosomes. RDD sites are significantly ($P < 10^{-10}$) enriched in genes that play a role in helicase activity, protein and nucleotide binding (table S6).

We also noted that the 10,210 sites which showed RNA and DNA sequence differences are not evenly distributed within genes. About 44% (4,453 sites) of them are located in coding exons (10% were found in the last exons), 4% (386 sites) are in the 5' UTRs, and 39% (3,977 sites) are in the 3' UTRs (see table S7, those remaining cannot be classified because of differences in gene structures across isoforms). The results suggest that there are more sites in the 3' ends than the 5' ends of genes; a pattern that was also observed in

deamination-mediated RNA editing (27, 36). Seventy-one percent of the coding sites result in non-synonymous amino acid changes, including 2.1% that lead to the gain or loss of a stop codon if translated into proteins. Relative to other structural features in genes, we found 4% of RDD sites are within 2 nucleotides of exon borders and 5% are within 30 nucleotides of poly(A) signals (table S7). Among RDD types, the numbers of sites near splice junctions are quite similar but the numbers near poly(A) sites are more different. C-to-A and G-to-A differences are found more often near poly(A) sites.

Sites also tended to cluster; for example, 2,613 sites (26%) are within 25 bp, and 1,059 sites (10%) are adjacent to each other. Statistical analysis using a runs test supports that the locations of the sites are not random (median $P = 0.22$). We did not find obvious patterns or associations with motifs shared across the sites, except for the A-to-G and A-to-C differences that show a preference for a cytidine 5' to the adenosine; as previously observed in *ADAR* mediated A-to-G changes (7, 34).

We examined the percentage of mRNAs that differs in sequence from the corresponding DNA. For each site to determine the RDD level, we counted the number of reads with a different nucleotide from that in the corresponding DNA sequence. The distribution of the level is bimodal (Fig. 1C); the average level is 20% (median = 13%). However, for some sites, RDD was detected in nearly 100% of the RNA sequences such as the A-to-C difference in the gene that encodes an mRNA decapping enzyme, *DCP1A* (chr3:53297343). This level is correlated with the frequency and types of RNA-DNA differences. Sites found in more than 50% of the informative individuals tend to have higher levels of RNA editing or RDD than other sites ($P < 10^{-5}$; fig. S5). The levels also differ across individuals. For example at a G-to-A site in the gene *RHOT1*, which encodes a RAS protein that plays a role in mitochondrial trafficking (chr17:27526465), in one person, the level was 90% while in another person, it was only 18%. We identified 437 sites with 10 or more informative individuals where the individuals with the highest levels and the lowest levels differ by 2 fold or more (range: 2 to 8.6 fold).

In conclusion, we have uncovered thousands of exonic sites where the RNA sequences do not match those of the DNA sequences; including transitions and transversions. These findings challenge the long-standing belief that in the same individuals, DNA and RNA sequences are nearly identical. To increase the confidence in our results, we obtained the DNA, RNA and protein sequences from different individuals and cell types using a range of technologies (fig. S1b). The samples included cell lines and primary cells from healthy individuals and tumors. We used data from public resources such as EST databases, The HapMap and 1000 Genomes Projects as well as those that we

generated with traditional Sanger sequencing, high-throughput sequencing technologies and mass spectrometry. Table 4 showed the DNA, RNA and peptide sequences at 15 confirmed sites which illustrate that the RNA and peptide sequences are the same but differ from the corresponding DNA sequences. The results support our observation that in an individual, DNA and RNA sequences from the same cells are not always identical and some of the variant RNA sequences are translated into proteins. The consistent pattern of the observations suggests that the RDDs have biological significance and are not just “noise.” At nearly all RDD sites, we observed only one RDD type across cell types and in different individuals. If the DNA sequence is A/A, and the RNA is A/C in one sample, in other samples, we see the same A-to-C difference, but not other types of differences. These results suggest that there are unknown aspects of transcription and/or post-transcriptional processing of RNA. These differences may now be studied along with those in other genomes and organisms such as the mitochondrial genomes of trypanosomes and chloroplasts of plants, where RNA editing and modifications are relatively common (35, 36).

The underlying mechanisms for these events are largely unknown. For most of the cases, we do not know yet whether a different base was incorporated into the RNA during transcription or if these events occur post-transcriptionally. About 23% of the sites are A-to-G differences; some of these are likely mediated by ADAR, but other, currently unknown, mechanisms can be involved. If it is a co-transcriptional process, then the signal can be in the DNA or the RNA such as secondary structures or modified nucleotides. In addition, as some of the RDDs are found near splice and poly(A) sites; it is possible that this may be a facet of systematic RNA processing steps such as splicing and cleavage (37, 38).

Our findings supplement previous studies demonstrating RNA-DNA differences in the human genome, and show that these differences go beyond A-to-G transition. These findings impact our understanding of genetic variation; in addition to DNA sequence variation, we identify individual variation in RNA sequences. For monomorphic DNA sequences that show RDD there is an overall increase in genetic variation. Thus, this variation contributes not only to individual variation in gene expression but also diversifies the proteome since some identified sites lead to nonsynonymous amino acid changes. We speculate that this RNA sequence variation likely affects disease susceptibility and manifestations. To date, mapping studies have focused on identifying DNA variants as disease susceptibility alleles. Our results suggest that the search may need to include RNA sequence variants that are not in the DNA sequences.

References and Notes

1. R. T. Libby, J. A. Gallant, *Mol Microbiol* **5**, 999 (1991).

2. J. F. Sydow, P. Cramer, *Curr Opin Struct Biol* **19**, 732 (2009).
3. S. H. Chen *et al.*, *Science* **238**, 363 (1987).
4. L. M. Powell *et al.*, *Cell* **50**, 831 (1987).
5. B. L. Bass, H. Weintraub, *Cell* **55**, 1089 (1988).
6. J. B. Li *et al.*, *Science* **324**, 1210 (2009).
7. A. Athanasiadis, A. Rich, S. Maas, *PLoS Biol.* **2**, e391. (2004).
8. M. J. Thomas, A. A. Platas, D. K. Hawley, *Cell* **93**, 627 (1998).
9. D. Wang *et al.*, *Science* **324**, 1203 (2009).
10. N. Zenkin, Y. Yuzenkova, K. Severinov, *Science* **313**, 518 (2006).
11. M. Sakurai, T. Yano, H. Kawabata, H. Ueda, T. Suzuki, *Nat. Chem. Biol.* **6**, 733 (2010).
12. K. Nishikura. (Annual Reviews, 2010), vol. 79, pp. 321-349.
13. E. Levanon *et al.*, *Nat Biotechnol.* **22**, 1001 (2004).
14. S. G. Conticello, *Genome Biol* **9**, 229 (2008).
15. A. Chester, J. Scott, S. Anant, N. Navaratnam, *Biochim. Biophys. Acta-Gene Struct. Expression* **1494**, 1 (2000).
16. J. Dausset *et al.*, *Genomics* **6**, 575 (1990).
17. International HapMap Consortium, *Nature* **426**, 789 (2003).
18. International HapMap Consortium, *Nature* **437**, 1299 (2005).
19. R. M. Durbin *et al.*, *Nature* **467**, 1061 (2010).
20. H. Cann, *Curr Opin Genet Dev.* **2**, 393 (1992).
21. T. C. Matise *et al.*, *Am J Hum Genet* **73**, 271 (2003).
22. F. Sanger, S. Nicklen, A. R. Coulson, *Proc Natl Acad Sci U S A* **74**, 5463 (1977).
23. D. R. Bentley *et al.*, *Nature* **456**, 53 (2008).
24. J. Harrow *et al.*, *Genome Biol.* **7 Suppl 1**, S4.1 (2006).
25. S. B. Montgomery *et al.*, *Nature* **464**, 773 (2010).
26. C. Trapnell *et al.*, *Nat Biotechnol* **28**, 511 (2010).
27. B. Rosenberg, C. Hamilton, M. Mwangi, S. Dewell, F. Papavasiliou, *Nat Struct Mol Biol.* **18**, 230 (2011).
28. S. Altschul, W. Gish, W. Miller, E. Myers, D. Lipman, *J Mol Biol.* **215**, 403 (1990).
29. K. Sodek *et al.*, *Mol Biosyst.* **4**, 762 (2008).
30. A. Michalski, J. Cox, M. Mann, *Journal of proteome research* **10**, 1785 (2011).
31. L. M. de Godoy *et al.*, *Genome Biol* **7**, R50 (2006).
32. C. M. Burns *et al.*, *Nature* **387**, 303 (1997).
33. H. Lomeli *et al.*, *Science* **266**, 1709 (1994).
34. S. Maas, S. Patt, M. Schrey, A. Rich, *Proc. Natl. Acad. Sci. U. S. A.* **98**, 14687 (2001).
35. C. Phreaner, M. Williams, R. Mulligan, *Plant Cell.* **8**, 107 (1996).
36. H. Hundley, A. Krauchuk, B. Bass, *Rna.* **14**, 2050 (2008).
37. S. Rueter, C. Burns, S. Coode, P. Mookherjee, R. Emeson, *Science.* **267**, 1491 (1995).

38. S. Rueter, T. Dawson, R. Emeson, *Nature*. **399**, 75 (1999).

39. Supporting online material.

Acknowledgments: Dedicated to the memory of Dr. Tom Kadesch who gave us important suggestions, taught us salient and subtle points on gene expression, and inspired us with his enthusiasm. Dr. Kadesch died during the preparation of this manuscript. We thank Drs. D. Epstein, H. Kazazian, D. Puppione and L. Simpson for suggestions and discussions. We thank Drs. C. Gunter, R. Nussbaum and J. Puck for comments on the manuscript, M. Morley for help with data analysis, W. Ankener for sample processing and J. Devlin for results on Sanger sequencing. The mass spectrometry analysis was carried out at the Wistar Proteomic Facility; we thank K. Speicher for help and suggestions. Funded by grants from the National Institutes of Health (to VGC and ML) and support from the Howard Hughes Medical Institute. The RNA-Seq data have been deposited to NCBI GEO under the accession number GSE25840.

Supporting Online Material

www.sciencemag.org/cgi/content/full/science.1207018/DC1

Materials and Methods

Figs. S1 to S5

Tables S1 to S10

References (40–44)

3 March 2011; accepted 10 May 2011

Published online 19 May 2011; 10.1126/science.1207018

Fig. 1. Characteristics of RDD sites. **(A)** Frequency of the 12 categories of RNA-DNA differences identified in B-cells of 27 normal individuals. **(B)** Number of RNA editing/ RDD events in 27 normal individuals. **(C)** Bimodal distribution of the levels of 28,766 RNA editing/ RDD events.

Fig. 2. Identification of peptides coded by both RNA and DNA sequences. **(A)** The RNA-form of a RDD leads to loss of a stop codon in RPL28 and extension of 55 amino acids. Peptides detected by mass spectrometry are shown in red. Extended protein sequence due to RDD is underlined. **(B & C)** MS-MS data confirm the detection of peptides encoded by the RDD containing RPL28 mRNA. The representative spectra of one peptide (SLIGTASEPR) from ovarian cancer cells **(B)** and cultured B-cells **(C)** are shown.

Table 1. Selected examples of sites that show RNA-DNA Differences in B-cells and EST clones.

Gene	Chr	Position (bp)*	Type	No. of informative individuals ^{†^}	No. of individuals with RDD [^]	Average level ^{‡^} [range]	EST
<i>HSP90AB1</i>	6	44,328,823	A-to-C	11	8	0.39 [0.15, 0.79]	BQ355193 (head neck), BX413896 (B cell)
<i>AZIN1</i>	8	103,910,812	A-to-G	17	10	0.22 [0.12, 0.37]	CD359333 (testis), BF475970 (prostate)
<i>CNBP</i>	3	130,372,812	A-to-T	18	16	0.13 [0.10, 0.21]	EL955109 (eye), BJ995106 (hepatoblastoma)
<i>MYL6</i>	12	54,841,626	C-to-A	16	16	0.35 [0.12, 0.60]	EC496428 (prostate), BG030232 (breast adenocarcinoma)
<i>RBM23</i>	14	22,440,217	C-to-G	11	5	0.18 [0.11, 0.35]	BQ232763 (testis, embryonic)
<i>RPL23</i>	17	34,263,515	C-to-T	12	8	0.16 [0.10, 0.22]	BP206252 (smooth muscle), CK128791 (embryonic stem cell)
<i>BLNK</i>	10	97,957,645	G-to-A	14	7	0.14 [0.11, 0.17]	BF972964 (leiomyosarcoma), BE881159 (lung carcinoma)
<i>C17orf70</i>	17	77,117,583	G-to-C	2	2	0.26 [0.24, 0.28]	AA625546 (melanocyte), AA564879 (prostate)
<i>HMG2</i>	1	26,674,349	G-to-T	7	4	0.22 [0.14, 0.43]	BX388386 (neuroblastoma), BE091398 (breast)
<i>CANX</i>	5	179,090,533	T-to-A	9	8	0.20 [0.13, 0.30]	EL950052, DB558106
<i>EIF3K</i>	19	43,819,430	T-to-C	19	14	0.16 [0.10, 0.27]	AI250201 (ovarian carcinoma), AI345393 (lung carcinoma)
<i>RPL37</i>	5	40,871,072	T-to-G	6	6	0.27 [0.16, 0.45]	CF124792 (T cell), DW459229 (liver)

* hg18 build of the human genome

[^] B-cells

[†] RNA-Seq ≥ 10 reads, DNA-Seq ≥ 4 reads

[‡] Calculated by tallying RNA-Seq reads that contain RDD and those that do not.

Table 2. Sanger sequencing of RDD sites.

Gene	Chr	Position (bp) [#]	Type	Location	Amino acid change	B-Cells*		Primary Skin Fibroblast*		Brain (cortex)*	
						# informative individuals	# individuals showing RDD	# informative individuals	# individuals showing RDD	# informative individuals	# individuals showing RDD
<i>EIF2AK2</i>	2	37,181,512	A-to-G	3' UTR	Not applicable	8	8	8	0	10	10
	2	37,181,517	A-to-G	3' UTR	Not applicable	8	8	8	3	10	10
	2	37,181,520	A-to-G	3' UTR	Not applicable	8	8	8	3	10	10
	2	37,181,538	A-to-G	3' UTR	Not applicable	8	8	8	6	10	10
<i>AZIN1</i> [†]	8	103,910,812	A-to-G	coding, exonic	S to G	2	2	10	0	9	8
<i>DPP7</i>	9	139,128,755	C-to-T	coding, exonic	Synonymous (P)	9	2	8	1	10	0
<i>PPWD1</i>	5	64894960	G-to-A	coding, exonic	E to K	2	2	8	8	8	8
<i>HLA-DQB2</i>	6	32,833,537	G-to-A	coding exonic	G to S	2	2	10	10	ne	ne
	6	32,833,545	G-to-A	coding, exonic	R to H	2	2	10	10	ne	ne
	6	32,833,550	C-to-T	coding, exonic	Synonymous (I)	2	2	10	10	ne	ne
<i>BLCAP</i> [‡]	20	35,580,977	A-to-G	coding, exonic	Q to R	6	4	10	4	6	6
<i>NDUFC2</i>	11	77,468,303	C-to-G	coding, exonic	L to V	10	0	10	0	10	0

* In all cases, matched DNA and RNA samples from the same individuals were sequenced
[†] Also reported by Li, Levanon et al, 2009(6). [‡] Known site that we used as positive control.

ne = not expressed

Table 3 Peptides encoded by both DNA and RNA-forms of mRNA at RDD sites.

Protein	Position (bp) [#]	RDD	AA change	DNA-form [†]	RNA-form [†]
AP2A2	chr11:976858	T-to-G	Y-to-D	<u>Y</u> LALESMTLASSEFSHEAVK	<u>D</u> LALESMTLASSEFSHEAVK
DFNA5*	chr7:24705225	T-to-A	L-to-Q	VFP <u>L</u> LLCITLNGLCALGR	VFP <u>Q</u> LLCITLNGLCALGR
ENO1	chr1:8848125	T-to-C	L-to-P	EG <u>L</u> ELLK	EG <u>P</u> ELLK
ENO3	chr17:4800624	T-to-G	V-to-G	LAQSNWGW <u>G</u> MVSHR	LAQSNWGW <u>G</u> MVSHR
FABP3	chr1:31618424	T-to-A	W-to-R	MVDAFLGT <u>W</u> K	MVDAFLG <u>T</u> R
FH*	chr1:239747217	T-to-A	I-to-K	I <u>E</u> YDTFGELK	<u>K</u> EYDTFGELK
HMGB1	chr13:29935772	T-to-A	Y-to-N	MSS <u>Y</u> AFFVQTCR	MSS <u>N</u> AFFVQTCR
NACA	chr12:55392932	G-to-A	D-to-N	<u>D</u> IELVMSQANVSR	<u>N</u> IELVMSQANVSR
NSF	chr17:42161411	T-to-C	V-to-A	LLDY <u>V</u> PIGPR	LLDY <u>A</u> PIGPR
POLR2B	chr4:57567852	T-to-A	L-to-Q	IISDG <u>L</u> K	IISDG <u>Q</u> K
RAD50*	chr5:131979610	T-to-G	L-to-R	W <u>L</u> QDNLTLR	W <u>R</u> QDNLTLR
RPL12	chr9:129250509	A-to-G	N-to-D	HSG <u>N</u> ITFDEIVNIAR	HSG <u>D</u> ITFDEIVNIAR
RPL32*	chr3:12852658	G-to-T	A-to-S	<u>A</u> AQLAIR	<u>S</u> AQLAIR
RPS3AP47*	chr4:152243651	C-to-A	T-to-K	EVQ <u>I</u> NDLK	EVQ <u>K</u> NDLK
SLC25A17	chr22:39520485	A-to-G	E-to-G	TTHMVLL <u>E</u> IIK	TTHMVLL <u>G</u> IIK
TUBA1*	chr2:219823379	A-to-G	E-to-G	EDMAAL <u>E</u> K	EDMAAL <u>G</u> K
TUBB2C	chr9:139257297	G-to-A	G-to-D	LHFFMP <u>G</u> FAPLTSR	LHFFMP <u>D</u> FAPLTSR

* DNA sequences of these and others were verified by Sanger sequencing (see Table S2).

^RDD in *RPL28* leads to the loss of a stop codon, the resulting additional peptides are found only in the RNA-form of the mRNA.

hg 18 build of the human genome

† For each peptide, the amino acid that differs between the DNA and RNA forms are underlined.

Table 4. Corresponding DNA, RNA and peptide sequences at selected sites.

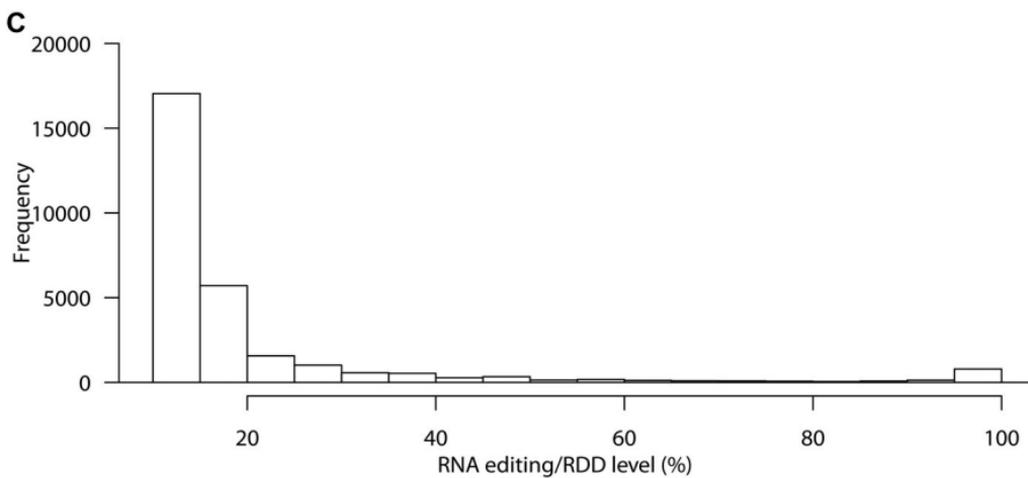
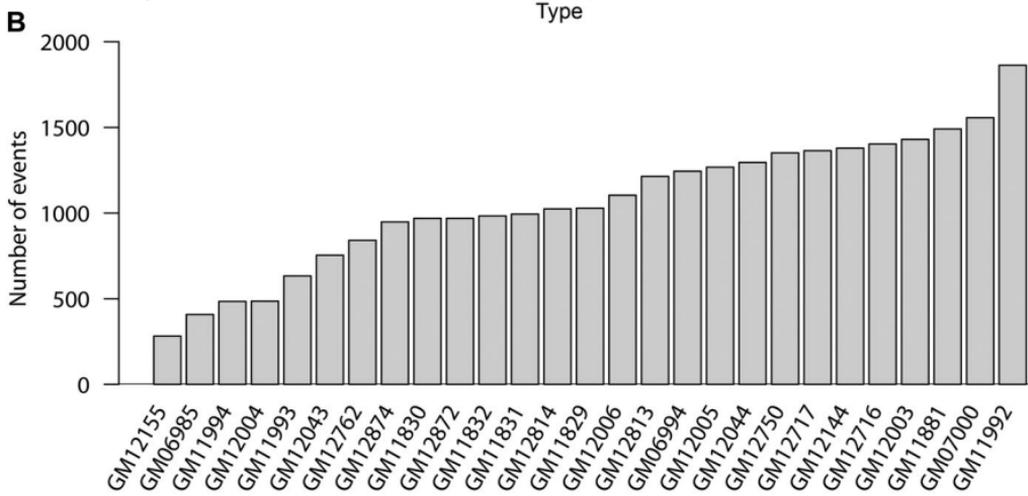
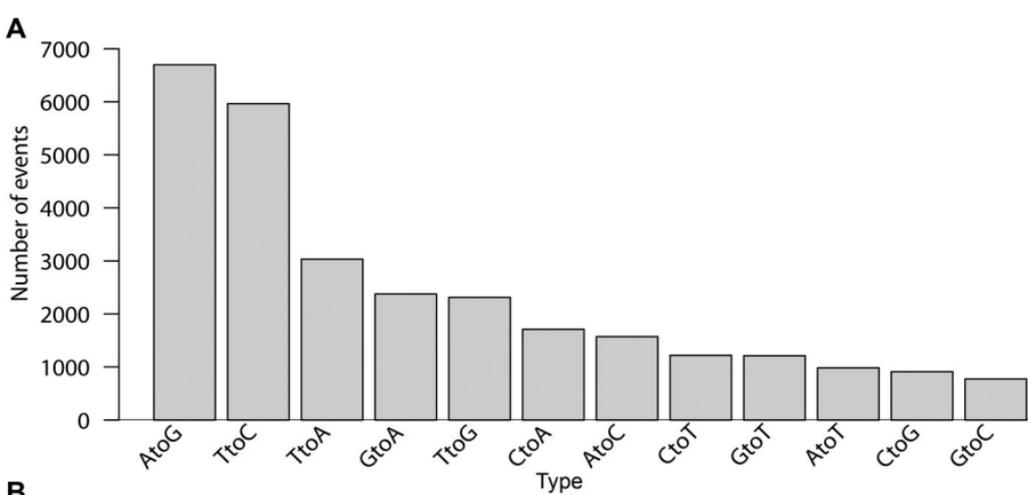
RDD	Gene	Location	DNA*†	RNA†	Peptide (DNA-form,LC-MS/MS)	Peptide (RNA-form, LC-MS/MS)
TtoG	<i>CD22</i>	chr19:40514815	<u>C</u> TG	C <u>G</u> G	ND	MHLLGPWLL <u>L</u> R
TtoA	<i>DFNA5</i>	chr7:24705225	<u>C</u> TG	C <u>A</u> G	VF <u>P</u> LLLCITLNGLCALGR	VF <u>P</u> QLLCITLNGLCALGR
TtoC	<i>ENO1</i>	chr1:8848125	<u>C</u> TG	C <u>C</u> G	EG <u>L</u> ELLK	EG <u>P</u> ELLK
TtoA	<i>FH</i>	chr1:239747217	<u>A</u> T <u>A</u>	<u>A</u> <u>A</u> <u>A</u>	<u>I</u> EYDTFGELK	<u>K</u> EYDTFGELK
TtoA	<i>HMGB1</i>	chr13:29935772	<u>I</u> <u>A</u> T	<u>A</u> <u>A</u> T	MSS <u>Y</u> AFFVQTCR	MSS <u>N</u> AFFVQTCR
AtoC	<i>HMGB1</i>	chr13:29935469	<u>A</u> <u>A</u> <u>A</u>	<u>A</u> <u>A</u> <u>C</u>	ND	TMSA <u>K</u> <u>N</u>
AtoC	<i>ITPR3</i>	chr6:33755773	<u>G</u> <u>A</u> C	<u>G</u> <u>C</u> C	ND	DGVEDHSPLMYHISLV <u>A</u> LLAACAEGK
TtoG	<i>RAD50</i>	chr5:131979610	<u>C</u> T <u>A</u>	C <u>G</u> <u>A</u>	W <u>L</u> QDNLT <u>L</u> R	W <u>R</u> QDNLT <u>L</u> R
GtoT	<i>ROD1</i>	chr9:114026264	<u>G</u> <u>G</u> <u>A</u>	<u>G</u> T <u>A</u>	ND	NLFIEA <u>V</u> CSVK
GtoT	<i>RPL32</i>	chr3:12852658	<u>G</u> <u>C</u> T	<u>I</u> <u>C</u> T	<u>A</u> AQLA <u>I</u> R	<u>S</u> AQLA <u>I</u> R
AtoG	<i>RPS25P8</i>	chr11:118393375	<u>A</u> <u>A</u> C	<u>G</u> <u>A</u> C	ND	EVP <u>D</u> YK
CtoA	<i>RPS3AP47</i>	chr4:152243651	<u>A</u> <u>C</u> <u>A</u>	<u>A</u> <u>A</u> <u>A</u>	EVQ <u>I</u> NDLK	EVQ <u>K</u> NDLK
GtoT	<i>SUPT5H</i>	chr19:44655806	C <u>A</u> <u>G</u>	C <u>A</u> <u>T</u>	ND	TPMYGSQTPL <u>H</u> DGSR
TtoC	<i>TOR1AIP1</i>	chr1:178144365	<u>I</u> <u>C</u> <u>A</u>	<u>C</u> <u>C</u> <u>A</u>	ND	QPSVLSP <u>G</u> YQK
AtoG	<i>TUBA1</i>	chr2:219823379	<u>G</u> <u>A</u> <u>G</u>	<u>G</u> <u>G</u> <u>G</u>	EDMAA <u>L</u> <u>E</u> K	EDMAA <u>L</u> <u>G</u> K

* DNA sequences are monomorphic according to dbSNP, 1000 Genomes and HapMap Projects; all individuals should have the reference allele. We verified this by Sanger sequencing of the B-cells used for mass spectrometry.

† RDD sites are underlined.

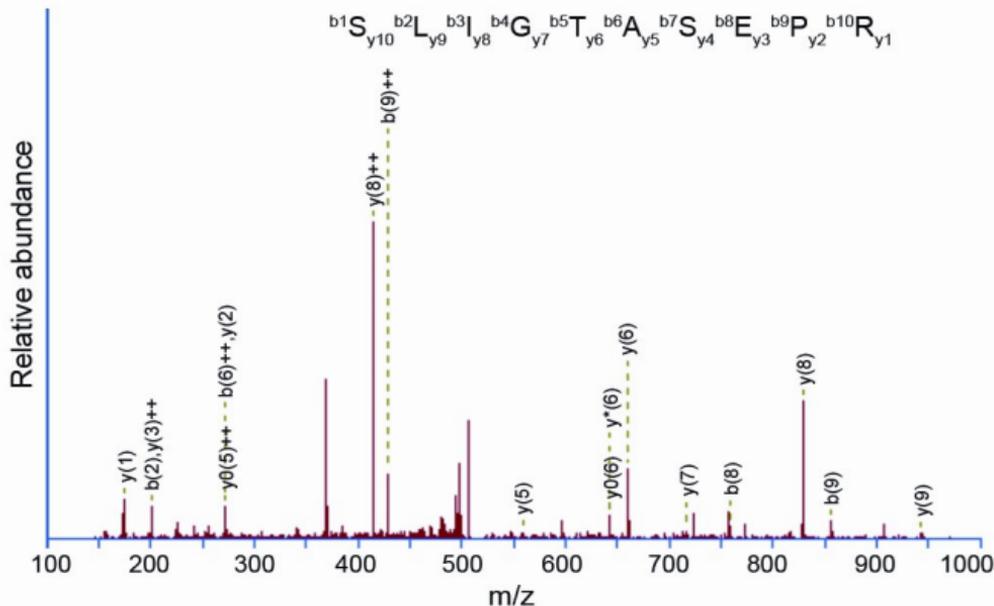
LC-MS/MS = liquid chromatography and tandem mass spectrometry

ND = not detected by mass spectrometry; however this does not mean that the peptides are absent in the B-cell proteome. It is likely a result of sampling.



A

1 **MSAHLQWMVV RNCSSFLIKR NKQTYSTEPN NLKARNSEFRY NGLIHRKTVG**
 51 **VEPAADGKGV VVVIKRRSER VFLRSLIGTA SEPRVLLLSG SNKRSLASD**
 101 **PPVSGTRSPG SSQLLGTWGP RSGES**

B**C**