

DARNED: a DAtabase of RNa EDiting in humans

Anmol Kiran and Pavel V. Baranov*

Biochemistry Department, University College Cork, Ireland

Associate Editor: Ivo Hofacker

ABSTRACT

Motivation: RNA editing is a phenomenon, which is responsible for the alteration of particular nucleotides in RNA sequences relative to their genomic templates. Recently, a large number of RNA editing instances in humans have been identified using bioinformatic screens and high-throughput experimental investigations utilizing next-generation sequencing technologies. However, the available data on RNA editing are not uniform and difficult to access.

Results: Here, we describe a new database DARNED (DAtabase of RNa EDiting) that provides centralized access to available published data related to RNA editing. RNA editing locations are mapped on the reference human genome. The current release of the database contains information on approximately 42 000 human genome coordinates corresponding to RNA locations that undergo RNA editing, mostly involving adenosine-to-inosine (A-to-I) substitutions. The data can be queried using a range of genomic coordinates, their corresponding functional localization in RNA molecules [Exons, Introns, CoDing Sequence (CDS) and UnTranslated Regions (UTRs)] and information regarding tissue/organ/cell sources where RNA editing has been observed. It is also possible to obtain RNA editing information for a specific gene or an RNA molecule using corresponding accession numbers. Search results provide information on the number of expressed sequence tags (ESTs) supporting edited and genomic bases, functional localization of RNA editing and existence of known single nucleotide polymorphisms (SNPs). Editing data can be explored in UCSC and Ensembl genome browsers, in conjunction with additional data provided by these popular genome browsers. DARNED has been designed for researchers seeking information on RNA editing and for the developers of novel algorithms for its prediction.

Availability: DARNED is accessible at <http://darned.ucc.ie>

Contact: p.baranov@ucc.ie; brave.oval.pan@gmail.com

Received on November 23, 2009; revised on May 3, 2010; accepted on May 25, 2010

1 INTRODUCTION

RNA editing is a biological phenomenon of changing RNA sequence relative to its corresponding DNA template by inserting, deleting or substituting one or a few nucleotides (Bass, 2002; Gott and Emeson, 2000; Maas and Rich, 2000). It is distinct from other nucleotide modifications in which nucleotides retain their identities during decoding or sequencing. For instance, a change of adenosine-to-inosine (A-to-I) is considered as RNA editing, since inosines are recognized as guanosines by polymerases and the translationary

machinery, while methylation of adenosines is not considered an editing event, since the identity of the nucleotide remains the same. RNA editing has been found across all kingdoms of life, including viruses (Casey, 2006; Kolakofsky *et al.*, 2005). There are several molecular mechanisms responsible for RNA editing. It can occur co-transcriptionally, for example, due to slippage of a growing nascent RNA chain relative to its template resulting in insertion or deletion of nucleotides (Baranov *et al.*, 2005; Iseni *et al.*, 2002; Penno *et al.*, 2006; Tamas *et al.*, 2008). Alternatively, it can be a result of enzymatic modification, e.g. deamination of cytidines and adenosines that leads to the formation of uridines and inosines, respectively (Bass, 2002).

Earlier, RNA editing was known to be particularly abundant in plant organelles (Covello and Gray, 1989) and was also identified in a few human genes, where it contributes to diversification of corresponding protein products, such as those encoding the ApoB protein and the glutamate and serotonin receptors (Bass, 2002; Chen *et al.*, 1987; Seeburg *et al.*, 1998). The abundance of RNA editing in humans, however, has only recently come to light. Most RNA editing instances found so far in humans involve A-to-I conversions that are carried out by ADARs (Adenosine Deaminase that Acts on RNAs). ADARs bind to double-stranded regions of RNA molecules and catalyze deamination of adenosines, converting them to inosine residues (Bass, 2002). Therefore, the majority of RNA editing instances occur in regions where long RNA stem-loop structures can be formed and are particularly abundant among Alu repeats (Athanasiadis *et al.*, 2004; Kim *et al.*, 2004; Levanon, K. *et al.*, 2005).

RNA editing plays a variety of functional roles in gene expression. RNA editing of a nucleotide within the CoDing Sequence (CDS) may change the identity of a particular encoded amino acid or could generate significantly altered protein sequences if editing leads to the formation of a stop codon (e.g. the ApoB case mentioned earlier). Editing occurring within introns can affect splicing, effectively resulting in the generation of alternatively spliced products (Maas and Gommans, 2009). Hyper-editing of UnTranslated Regions (UTRs) leads to retention of mRNA inside the nucleus, consequently down-regulating synthesis of the encoded protein by preventing transport of its mRNA to the cytoplasm (Sansam *et al.*, 2003). An interesting role of RNA editing is in diversification of miRNA sequences that consequently change their target repertoire (Blow *et al.*, 2006; Kawahara *et al.*, 2008). Abnormal RNA editing has also been associated with a number of diseases, including several neurological disorders and cancers (Klimek-Tomczak *et al.*, 2006; Vollmar *et al.*, 2004).

A large number of recently discovered RNA editing instances have been found by means of bioinformatic screens of discrepancies between genomic sequences and corresponding RNA products, and by applying statistical filters to distinguish RNA editing from

*To whom correspondence should be addressed.

sequencing errors or human polymorphisms (Athanasiadis *et al.*, 2004; Clutterbuck *et al.*, 2005; Kim *et al.*, 2004; Levanon *et al.*, 2004). Recent developments in high-throughput parallel sequencing technologies have facilitated the experimental analysis of a large amount of RNA editing candidates, making the past year particularly remarkable for the area (Li *et al.*, 2009). A surprisingly large number of confirmed RNA editing instances have been found in protein-coding regions. Interestingly, many proteins whose sequences are altered by means of RNA editing are not directly involved in synaptic functions, contrary to previous observations that RNA editing affects mostly proteins which function in the nervous system (Hoopengardner *et al.*, 2003; Jepson and Reenan, 2008; Levanon, E.Y. *et al.*, 2005). The data on RNA editing accumulated so far are spread across a number of research papers in the form of supplementary files. Even though there are a few databases dedicated to RNA editing, such as REDIdb (Picardi *et al.*, 2007) and DbRES (He *et al.*, 2007), they do not provide access to these new datasets on human RNA editing. According to a recent examination (during the preparation of this manuscript), REDIdb did not have any data related to human RNA editing, while DbRES contains information on only 34 instances. Therefore, there is a clear and timely need for a resource that would provide unified access to the accumulating information on RNA editing in humans.

To generate such a resource, we have designed a database, called DARNED (Database of RNa EDiting), which is freely accessible at <http://darned.ucc.ie>. DARNED is a convenient resource allowing centralized uniform access to growing information on RNA editing in humans. At present the database stores information on approximately 42 000 instances of RNA editing. While most of the current data are on A-to-I editing, a few citidine-to-uridine (C-to-U) instances are also included in the DARNED dataset.

2 DATA COLLECTION AND PROCESSING

For data collection, we searched PubMed for articles describing newly identified instances of RNA editing. We have identified four types of datasets: (i) those based on bioinformatics analysis of discrepancies between cDNA sequences and genomic sequences (Blow *et al.*, 2006; Clutterbuck *et al.*, 2005; Kim *et al.*, 2004; Levanon *et al.*, 2004; Levanon, K. *et al.*, 2005); (ii) analysis of single nucleotide polymorphism (SNP) data (Eisenberg *et al.*, 2005; Levanon, E.Y. *et al.*, 2005); (iii) analysis of miRNA (Blow *et al.*, 2006; Kawahara *et al.*, 2008); and (iv) results of high-throughput sequencing of RNA and DNA samples from the same tissues (Li *et al.*, 2009).

To describe RNA editing instances in a uniform format, we decided to map RNA editing locations to the reference human genome. Since original data have been formatted in different ways, we processed each dataset in a unique manner. For example, Levanon *et al.* (2004) generated a dataset with sequences flanking editing locations by 100 nt on both sides. To process these data, we used BLAT (Kent, 2002) to identify genomic locations matching these flanking sequences. Coordinates corresponding to the position 101 (within 201 nt sequences) have been chosen as genomic locations that are edited on the RNA level. (Kim *et al.*, 2004) have provided accession numbers of edited cDNAs also indicating matching coordinates from hg16 human genome assembly. To process these data, the coordinates of hg16 assembly were converted into coordinates of the hg18 assembly using UCSC liftOver tool

(Kuhn *et al.*, 2009). Further BLAT was used to generate alignments between cDNA sequences and genomic sequences. Adenosine-to-guanosine (A-to-G) mismatches were further identified using a custom python script and the coordinates of mismatches were extracted. Research articles dedicated to the analysis of SNP data have reported accession numbers of SNPs in dbSNP (Sherry *et al.*, 2001). We used these accession numbers to find genomic coordinates of SNPs that were further taken directly into the DARNED dataset as locations of RNA editing. Kawahara *et al.* (2008) have reported RNA editing locations within miRNA sequences. Since miRNA sequences are short, they are not suitable for the BLAT search. Therefore, corresponding pre-miRNA sequences have been downloaded from the miRBase (Griffiths-Jones *et al.*, 2006). Pre-miRNA sequences were then aligned to mature miRNA sequences using ClustalW2 (Larkin *et al.*, 2007) to identify RNA editing locations within miRNA sequences. Next, miRNA sequences were aligned to the reference genome using BLAT and coordinates corresponding to RNA locations were extracted in a manner similar to that described above. Genomic coordinates have been reported in a number of datasets. In cases where the coordinates corresponded to the hg18 assembly, the data were used directly, while for instances where coordinates corresponded to previous assemblies, they were converted to hg18 using the UCSC LiftOver tool. The total number of RNA editing instances currently stored in the database is about 42 000.

Information on C-to-U editing was obtained from the following research publications: Blanc and Davidson (2003), Keegan *et al.* (2001) and Kondo *et al.* (2004).

For all RNA instances, DARNED provides information on a number of expressed sequence tags (ESTs) supporting both edited and non-edited base identities, information on tissue or cell type sources and also on functional regions in RNAs. Tissue source information has been extracted directly from the datasets of the analyzed research articles. To obtain information on EST sequences supporting edited or non-edited bases, coordinates of matching ESTs and genomic blocks were obtained from the all_est MySQL table of UCSC genome browser. Base identities of RNA editing genomic locations were then compared with base identities in the corresponding locations of individual ESTs. Functional regions of RNA editing locations were identified using the refGene MySQL table of UCSC genome browser and first classified as intronic or exonic. Further, exonic regions were sub-classified as CDS, 5' or 3' UTRs. Locations that did not correspond to mRNA sequences were classified as 'others'.

3 DATABASE INTERFACE

The database can be freely accessed at <http://darned.ucc.ie>. The web interface is based on HTML and Python/CGI scripts that provide access to a local MySQL database. The search page of DARNED provides the interface for querying the database. Figure 1 illustrates a combination of the search form and an example of a search output. To search the database, a user needs to specify a region within the reference human genome (hg18 assembly) by specifying a chromosome ID and a coordinate range. The default is set at zero value, which will result in querying the entire chromosome length. The search can be limited by several filters. First, it can be limited to a tissue where supporting sequences (cDNA, EST or RNAseq short reads) have been identified. The user can then choose between

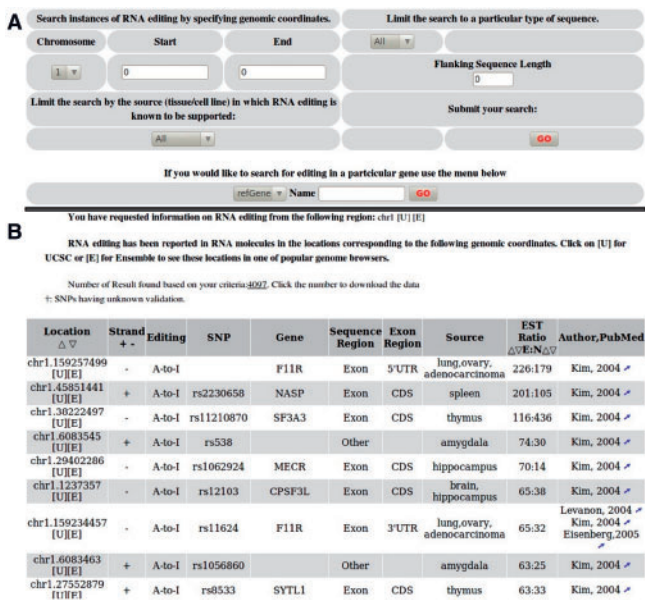


Fig. 1. (A) Database search form. (B) Search results table.

normal or cancerous tissues or human cell lines. After that choice has been made, the user is provided with additional menus that list available sources (a particular tissue type or a cell line). The data output can also be limited to a particular functional region within an RNA molecule (intronic, exonic or others) Exonic regions can be further specified as CDS, 5' or 3' UTRs. The last option in the search form is specification of flanking sequence length. This value is used to specify the length of genomic region in which RNA location will be displayed in UCSC genome browser custom tracks (Kuhn *et al.*, 2009) or Ensembl distributed annotation system (DAS; Hubbard *et al.*, 2009) for global exploration of RNA editing in the context of other available information for a chosen genomic region.

Alternatively, the database can be searched for RNA editing instances associated with a particular gene or a molecule. For this purpose, the user can choose the menu located at the bottom of the search page, where sequence type needs to be specified and then an accession number for corresponding sequence needs to be inserted into the adjacent box.

Search results are displayed in the form of a table (Fig. 1B). The text above the table contains links to the UCSC custom tracks [U] and Ensembl [E] corresponding to the entire region of the human genome specified in the search form. Each row of the table itself corresponds to a particular editing location. The table is organized into 10 columns: location, strand, editing, SNP, gene, sequence region, exon region, source, EST ratio and reference. The *location* column provides genomic reference coordinates with a chromosome index separated from a coordinate by a dot, e.g. chr1.239848 and links to custom tracks on UCSC [U] and Ensembl [E] genome browsers. Clicking on these links opens a separate window in a corresponding genome browser showing a coordinate of RNA editing within a flanking region specified in the search form. *Strand* indicates from which strand an edited RNA is transcribed. The menu option allows to group entries by strand. The *editing* column provides information on the type of editing associated with the particular chromosomal location, e.g. A-to-I or C-to-U. *SNP*

column indicates whether there is a known SNP in dbSNP, in which case RefSNP accession number is given that is hyperlinked to its corresponding entry in NCBI dbSNP. A number of SNPs in the dbSNP are not supported by genomic sequences and their validation is classified as 'unknown', such cases are likely to be RNA editing events misinterpreted as SNPs (Eisenberg *et al.*, 2005). To help users to compare DARNED and dbSNP data, SNPs with unknown validation are marked with † symbol. *Gene* column provides information on a gene if it is associated with a corresponding locus, the gene name is linked to a separate table (an example can be found in the Help page of DARNED). This table lists all the known transcript variants for this gene and functional classification of the edited locations in relation to the individual transcripts. Transcripts are listed by their accession numbers and hyperlinked to NCBI RefSeq database. *Sequence region* indicates whether the RNA editing location corresponds to intronic or exonic sequences and in the case of an exonic sequence, the next column *exonic region* also indicates a functional location within a corresponding mRNA, e.g. CDS, 5' or 3' UTR. Information on a tissue or a cell type where a corresponding RNA molecule has been identified is provided in the *source* column. *EST ratio* column provides information on the number of ESTs that are consistent with the edited base (E) and the number of ESTs that are consistent with the genomic base (N) separated by a colon. The EST ratio is not intended to indicate the efficiency of RNA editing, since (i) RNA editing could be tissue or developmentally regulated and (ii) available ESTs are likely to be skewed in representation of different sources. Clicking on a particular number of ESTs within a particular cell generates a new table where all ESTs are listed by their accession numbers and also by the source from which they were obtained. An example of such a table can be found in DARNED help page. All EST accession numbers are hyperlinked with corresponding entries in dbEST. Finally, the *PubMed* column contains citations of primary sources where RNA editing has been identified, with a link to its abstract on PubMed. The data in the table can be sorted using devoted links in the header of the table or can be grouped together according to a particular value in the table by clicking on that value within the table.

4 CUSTOM TRACKS AND FUNCTIONAL DISTRIBUTION OF RNA EDITING SITES

Custom tracks for the UCSC genome browser and Ensembl DAS have been designed to provide users with means to explore RNA editing datasets within their genomic locations in one of the two most popular human genome browsers. For each search, cgi script designs a new custom track in the bed format as a separate file for each individual chromosome. Figure 2 illustrates screen-shots of RNA editing custom tracks at different levels of magnification. Each RNA editing location is indicated as a differentially colored vertical dash. Differential coloring is used to indicate a type of functional RNA region where an edited base is located, i.e. red for intron, blue for CDS, dark green for 5' UTR and deep pink for 3' UTR.

The current distribution of known RNA editing instances in terms of their functional location within RNA molecules is shown in Figure 3. An interesting observation that can be made from this distribution is that RNA editing is significantly more frequent in 3' UTRs compared with 5' UTRs.

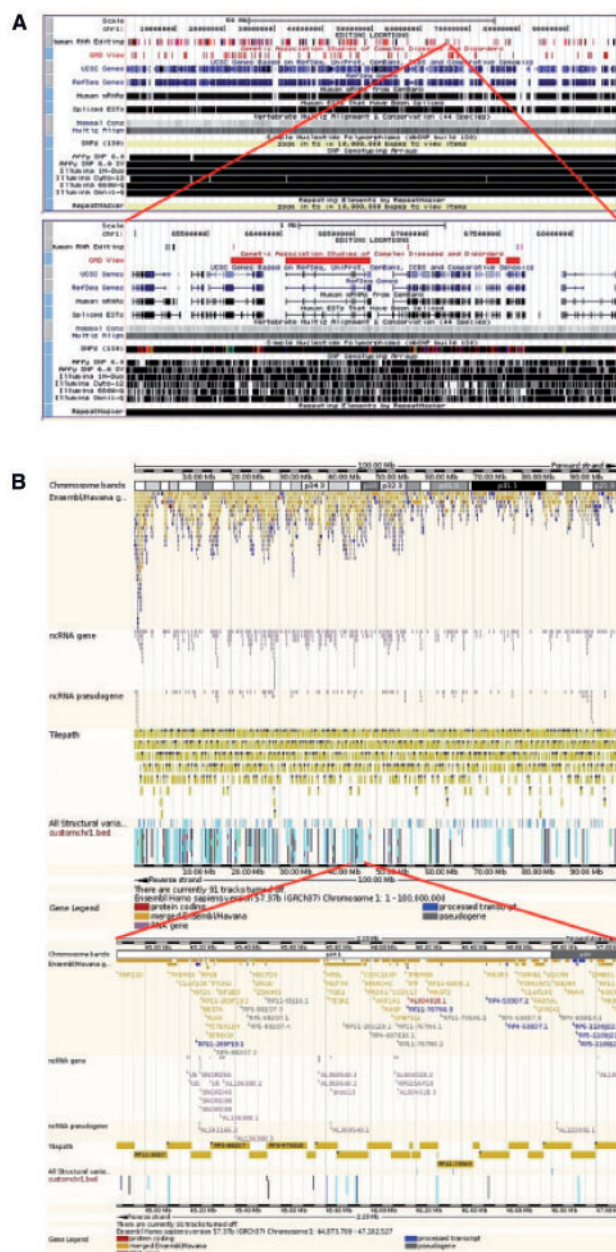


Fig. 2. (A) UCSC and (B) Ensembl custom tracks at different levels of magnification.

5 FUTURE PROSPECTS

DARNED will accommodate the novel RNA editing data that is anticipated to become available in the near future due to the progress in high-throughput sequencing technologies such as whole transcriptome shotgun sequencing (Mardis, 2008; Wang *et al.*, 2009). This anticipation is based on a recent frontier article on the analysis of ~36 K suspected RNA editing candidates in seven human tissues (Li *et al.*, 2009). These data are extremely useful in that they provide direct experimental analysis of RNA editing instances with reliable information on the efficiency of RNA editing and currently they form the largest bulk of data in DARNED.

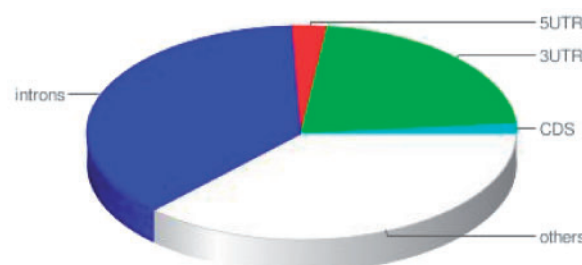


Fig. 3. Distribution of RNA editing instances in different functional locations of RNA molecules.

Using such data, we plan to design more detailed descriptions of RNA editing locations, which will include data on the efficiencies of RNA editing levels in individual human tissues based on quantities of corresponding short reads.

To keep the database synchronized with the developments in the RNA editing field, we constantly monitor newly published literature. We also plan to expand the Database to accommodate RNA editing data that is available for other model organisms.

ACKNOWLEDGEMENTS

We are grateful to Dr Michaël Bekaert for useful hints and tips during development of the database and to Prof. John Atkins and Kathy Barriscale for the help they provided during the preparation of the manuscript.

Funding: Science Foundation Ireland Principal Investigator Award 06/IN.1/B81 (to P.V.B.).

Conflict of Interest: none declared.

REFERENCES

- Athanasiadis, A. *et al.* (2004) Widespread A-to-I RNA editing of Alu-containing mRNAs in the human transcriptome. *PLoS Biol.*, **2**, e391.
- Baranov, P.V. *et al.* (2005) Transcriptional slippage in bacteria: distribution in sequenced genomes and utilization in IS element gene expression. *Genome Biol.*, **6**, R25.
- Bass, B.L. (2002) RNA editing by adenosine deaminases that act on RNA. *Annu. Rev. Biochem.*, **71**, 817–846.
- Blanc, V. and Davidson, N.O. (2003) C-to-U RNA editing: mechanisms leading to genetic diversity. *J. Biol. Chem.*, **278**, 1395–1398.
- Blow, M.J. *et al.* (2006) RNA editing of human microRNAs. *Genome Biol.*, **7**, R27.
- Casey, J.L. (2006) RNA editing in hepatitis delta virus. *Curr. Top. Microbiol. Immunol.*, **307**, 67–89.
- Chen, S.H. *et al.* (1987) Apolipoprotein B-48 is the product of a messenger RNA with an organ-specific in-frame stop codon. *Science*, **238**, 363–366.
- Clutterbuck, D.R. *et al.* (2005) A bioinformatic screen for novel A-I RNA editing sites reveals recoding editing in BC10. *Bioinformatics*, **21**, 2590–2595.
- Covello, P.S. and Gray, M.W. (1989) RNA editing in plant mitochondria. *Nature*, **341**, 662–666.
- Eisenberg, E. *et al.* (2005) Identification of RNA editing sites in the SNP database. *Nucleic Acids Res.*, **33**, 4612–4617.
- Gott, J.M. and Emeson, R.B. (2000) Functions and mechanisms of RNA editing. *Annu. Rev. Genet.*, **34**, 499–531.
- Griffiths-Jones, S. *et al.* (2006) miRBase: microRNA sequences, targets and gene nomenclature. *Nucleic Acids Res.*, **34**, D140–D144.
- He, T. *et al.* (2007) dbRES: a web-oriented database for annotated RNA editing sites. *Nucleic Acids Res.*, **35**, D141–D144.
- Hoopengardner, B. *et al.* (2003) Nervous system targets of RNA editing identified by comparative genomics. *Science*, **301**, 832–836.
- Hubbard, T.J. *et al.* (2009) Ensembl 2009. *Nucleic Acids Res.*, **37**, D690–D697.

- Iseni,F. *et al.* (2002) Chemical modification of nucleotide bases and mRNA editing depend on hexamer or nucleoprotein phase in Sendai virus nucleocapsids. *RNA*, **8**, 1056–1067.
- Jepson,J.E. and Reenan,R.A. (2008) RNA editing in regulating gene expression in the brain. *Biochim. Biophys. Acta*, **1779**, 459–470.
- Kawahara,Y. *et al.* (2008) Frequency and fate of microRNA editing in human brain. *Nucleic Acids Res.*, **36**, 5270–5280.
- Keegan,L.P. *et al.* (2001) The many roles of an RNA editor. *Nat. Rev.*, **2**, 869–878.
- Kent,W.J. (2002) BLAT—the BLAST-like alignment tool. *Genome Res.*, **12**, 656–664.
- Kim,D.D. *et al.* (2004) Widespread RNA editing of embedded alu elements in the human transcriptome. *Genome Res.*, **14**, 1719–1725.
- Klimek-Tomczak,K. *et al.* (2006) Editing of hnRNP K protein mRNA in colorectal adenocarcinoma and surrounding mucosa. *Br. J. Cancer*, **94**, 586–592.
- Kolakofsky,D. *et al.* (2005) Paramyxovirus mRNA editing, the “rule of six” and error catastrophe: a hypothesis. *J. Gen. Virol.*, **86**, 1869–1877.
- Kondo,N. *et al.* (2004) RNA editing of interleukin-12 receptor beta2, 2451 C-to-U (Ala 604 Val) conversion, associated with atopy. *Clin. Exp. Allergy*, **34**, 363–368.
- Kuhn,R.M. *et al.* (2009) The UCSC Genome Browser Database: update 2009. *Nucleic Acids Res.*, **37**, D755–D761.
- Larkin,M.A. *et al.* (2007) Clustal W and Clustal X version 2.0. *Bioinformatics*, **23**, 2947–2948.
- Levanon,E.Y. *et al.* (2004) Systematic identification of abundant A-to-I editing sites in the human transcriptome. *Nat. Biotechnol.*, **22**, 1001–1005.
- Levanon,E.Y. *et al.* (2005) Evolutionarily conserved human targets of adenosine to inosine RNA editing. *Nucleic Acids Res.*, **33**, 1162–1168.
- Levanon,K. *et al.* (2005) Letter from the editor: adenosine-to-inosine RNA editing in Alu repeats in the human genome. *EMBO Rep.*, **6**, 831–835.
- Li,J.B. *et al.* (2009) Genome-wide identification of human RNA editing sites by parallel DNA capturing and sequencing. *Science*, **324**, 1210–1213.
- Maas,S. and Gommans,W.M. (2009) Novel exon of mammalian ADAR2 extends open reading frame. *PLoS One*, **4**, e4225.
- Maas,S. and Rich,A. (2000) Changing genetic information through RNA editing. *Bioessays*, **22**, 790–802.
- Mardis,E.R. (2008) The impact of next-generation sequencing technology on genetics. *Trends Genet.*, **24**, 133–141.
- Penno,C. *et al.* (2006) Transcriptional slippage controls production of type III secretion apparatus components in *Shigella flexneri*. *Mol. Microbiol.*, **62**, 1460–1468.
- Picardi,E. (2007) REDIdb: the RNA editing database. *Nucleic Acids Res.*, **35**, D173–D177.
- Sansam,C.L. *et al.* (2003) Modulation of RNA editing by functional nucleolar sequestration of ADAR2. *Proc. Natl Acad. Sci. USA*, **100**, 14018–14023.
- Seeburg,P.H. *et al.* (1998) RNA editing of brain glutamate receptor channels: mechanism and physiology. *Brain Res.*, **26**, 217–229.
- Sherry,S.T. *et al.* (2001) dbSNP: the NCBI database of genetic variation. *Nucleic Acids Res.*, **29**, 308–311.
- Tamas,I. *et al.* (2008) Endosymbiont gene functions impaired and rescued by polymerase infidelity at poly(A) tracts. *Proc. Natl Acad. Sci. USA*, **105**, 14934–14939.
- Vollmar,W. *et al.* (2004) RNA editing (R/G site) and flip-flop splicing of the AMPA receptor subunit GluR2 in nervous tissue of epilepsy patients. *Neurobiol. Dis.*, **15**, 371–379.
- Wang,Z. *et al.* (2009) RNA-Seq: a revolutionary tool for transcriptomics. *Nat. Rev. Genet.*, **10**, 57–63.