

The colorectal microRNAome

Jordan M. Cummins[†], Yiping He[†], Rebecca J. Leary[†], Ray Pagliarini[†], Luis A. Diaz, Jr.[†], Tobias Sjoblom[†], Omer Barad[‡], Zvi Bentwich[‡], Anna E. Szafranska[§], Emmanuel Labourier[§], Christopher K. Raymond[¶], Brian S. Roberts[¶], Hartmut Juhl^{||}, Kenneth W. Kinzler[†], Bert Vogelstein^{†,††}, and Victor E. Velculescu^{†,††}

[†]The Sidney Kimmel Comprehensive Cancer Center and Howard Hughes Medical Institute, Johns Hopkins University Medical Institutions, Baltimore, MD 21231; [‡]Rosetta Genomics, 10 Plaut Street, Science Park, Rehovot 76706, Israel; [§]Ambion Diagnostics, 2130 Woodward Street, Austin, TX 78744; [¶]Rosetta Inpharmatics, 401 Terry Avenue North, Seattle, WA 98109; and ^{||}Indivumed GmbH, Center for Cancer Research at the Israelitic Hospital, Orchideenstieg 14, 22297 Hamburg, Germany

Contributed by Bert Vogelstein, December 27, 2005

MicroRNAs (miRNAs) are a class of small noncoding RNAs that have important regulatory roles in multicellular organisms. The public miRNA database contains 321 human miRNA sequences, 234 of which have been experimentally verified. To explore the possibility that additional miRNAs are present in the human genome, we have developed an experimental approach called miRNA serial analysis of gene expression (miRAGE) and used it to perform the largest experimental analysis of human miRNAs to date. Sequence analysis of 273,966 small RNA tags from human colorectal cells allowed us to identify 200 known mature miRNAs, 133 novel miRNA candidates, and 112 previously uncharacterized miRNA* forms. To aid in the evaluation of candidate miRNAs, we disrupted the *Dicer* locus in three human colorectal cancer cell lines and examined known and novel miRNAs in these cells. These studies suggest that the human genome contains many more miRNAs than currently identified and provide an approach for the large-scale experimental cloning of novel human miRNAs in human tissues.

colorectal cancer | dicer | microRNA

MicroRNAs (miRNAs) are ≈22-nt noncoding RNAs that are processed from larger (≈80-nt) precursor hairpins by the RNase III enzyme Dicer into miRNA:miRNA* duplexes (1–3). One strand of these duplexes associates with the RNA-induced silencing complex (RISC), whereas the other is generally degraded (1). The miRNA–RISC complex targets messenger RNAs for translational repression or mRNA cleavage. There has been considerable debate about the total number of miRNAs that are encoded in the human genome. Initial estimates, relying mostly on evolutionary conservation, suggested there were up to 255 human miRNAs (4). More recent analyses have demonstrated there are numerous nonconserved human miRNAs and suggest this number may be significantly larger (5).

Both cloning and bioinformatic approaches have been used to identify miRNAs. Direct miRNA cloning strategies identified many of the initial miRNAs and demonstrated that miRNAs are found in many species (6–16). However, the throughput of this approach is low, and cloning approaches have appeared to approach saturation (8). Bioinformatic strategies have recently been used to identify potential miRNAs predicted on the basis of various sequence and structural characteristics (4, 7). However, such gene predictions may not point to all legitimate miRNAs, especially those that are not phylogenetically conserved, and all *in silico* predictions require independent experimental validation.

To increase the efficiency of discovery of small RNA species, we have developed an approach called miRNA serial analysis of gene expression (miRAGE). This approach combines aspects of direct miRNA cloning and SAGE (17). Similar to traditional cloning approaches, miRAGE starts with the isolation of 18- to 26-base RNA molecules to which specialized linkers are ligated, and which are reverse-transcribed into cDNA (Fig. 1A). However, subsequent steps, including amplification of the complex mixture of cDNAs using PCR, tag purification, concatenation, cloning, and sequencing, have been performed by using SAGE methodology optimized for small RNA species. This approach

has the advantage of generating large concatemers that can be used to identify as many as 35 tags in a single sequencing reaction, whereas existing cloning protocols analyze on average approximately five miRNAs per reaction (8).

Results and Discussion

Genome-Wide miRNA Analysis with miRAGE. Using miRAGE, we analyzed 273,966 cDNA tags obtained from four human colorectal cancers and two matching samples of normal colonic mucosae. Comparing these tags to the existing miRNA database identified 68,376 tags matching known miRNA sequences. These represent the largest collection of human miRNA sequences identified to date, because all previous human miRNA cloning analyses in aggregate have analyzed <2,000 miRNA molecules. The expression level of the miRNAs detected by miRAGE ranged over 4 orders of magnitude (from 23,431 observations for miR-21 to 20 miRNAs that were observed only once), suggesting this approach can detect miRNAs present at varied expression levels. The identified miRNA tags matched 200 of the mature miRNAs present in the public miRBase database (2) (Table 2, which is published as supporting information on the PNAS web site), and 52 of these were expressed at significantly different levels between tumor cells and normal colonic epithelium ($P < 0.05$, Fisher exact test; Table 3, which is published as supporting information on the PNAS web site). Importantly, of the already catalogued miRNAs, these results provide novel experimental evidence for 62 miRNAs whose presence in this database was based solely on phylogenetic predictions.

In addition to detecting known or predicted miRNAs, 1,411 of the miRAGE tags represented 100 previously unrecognized miRNA* forms of known miRNAs (Table 4, which is published as supporting information on the PNAS web site). miRNA* molecules correspond to the short-lived complementary strand present in initial miRNA duplexes, and their biologic role, if any, has yet to be elucidated. Although miRNA* have been inferred to exist for all miRNAs, only 24 human miRNAs* have previously been reported in the public database. These analyses therefore provide substantially greater evidence for the presence of these molecules in human cells.

Evaluation of Novel miRNAs. We next focused on evaluating whether the miRAGE tags not matching known miRNAs might represent novel miRNA species. As a first step, miRAGE tags were compared with existing gene databases to exclude sequences matching known

Conflict of interest statement: K.W.K. and V.E.V. receive research funding from Genzyme, and K.W.K., V.E.V., and Johns Hopkins University own Genzyme stock, which is subject to certain restrictions under Johns Hopkins University policy. K.W.K. and V.E.V. are also paid consultants to Genzyme. The terms of this arrangement are being managed by Johns Hopkins University in accordance with its conflict of interest policies.

Abbreviations: miRNA, microRNA; SAGE, serial analysis of gene expression; miRAGE, microRNA SAGE; *Dicer*^{ex5}, *Dicer* exon 5-disrupted lines; qRT-PCR, quantitative RT-PCR.

^{††}To whom correspondence may be addressed. E-mail: vogelbe@welch.jhu.edu or velculescu@jhmi.edu.

© 2006 by The National Academy of Sciences of the USA

RT (Invitrogen) for 50 min at 45°C. Subsequently, the procedures for amplifying, isolating, purifying, concatenating, cloning, and sequencing tags are nearly identical to those performed in Long-SAGE and Digital Karyotyping, except that miRAGE PCR products range in size from 110 to 118 bp, and miRAGE tags (not ditags) were released from linkers with XhoI endonuclease (NEB). The sequencing of concatemer clones was performed by contract sequencing at Agencourt (Beverly, MA). Resulting sequence files were trimmed by using PHRED sequence analysis software (Codon-Code, Dedham, MA), and 18- to 26-bp tags were extracted by using the SAGE2000 software package, which identifies the fragmenting enzyme site between tags, extracts intervening tags, and records them in a database.

Bioinformatic Analyses of miRAGE Tags. Step 1: Grouping and comparing miRAGE tags to known RNAs. All tags sharing a common set of 11 of 12 core internal sequence elements were assembled into groups containing all related members. The tag with the most counts in each group was further analyzed. Grouping facilitated analysis by (i) eliminating rare sequencing errors and (ii) removing trivial miRNA variants, because miRNAs are known to display both 5' and 3' variation. The tags were subsequently compared to databases of known RNA sequences (miRNAs, mRNAs, rRNAs, etc.), using BLAST, and those tags matching known sequences were removed from further analysis. The tags obtained by miRAGE were compared with public databases on September 1, 2005. Subsequent additions and changes to these databases are not reflected in the data analysis.

Step 2: Secondary structure analysis and hairpin stability scoring of candidate miRNAs. To determine potential miRNA precursor structures, each tag was compared to the human genome sequence. For tags with perfect matches, a total of 75 bp (60 + 15 bp) of flanking genomic sequence around each tag was extracted. Because there are two possible precursors for each tag (i.e., the tag can be located on the 5' or 3' arm of a putative hairpin), pairs of theoretical precursors were extracted from the human genome at the position of each tag and were carried through the following analysis. Secondary structure and free energy of folding were determined for each pair of precursor structures by using MFOLD 3.2 (26, 27) and compared to values obtained for known miRNAs. The values used for thermodynamic evaluation were the free energy of folding of each precursor sequence ($\Delta G_{\text{folding}}$) and the difference of $\Delta G_{\text{folding}}$ between the two possible precursors ($\Delta\Delta G_{\text{folding}}$). Analysis of an arbitrary set of 126 known miRNAs using these thermodynamic analyses revealed that the highest $\Delta G_{\text{folding}}$ was -22.6 , and there were no miRNAs with a $\Delta G_{\text{folding}} > -29.0$, which had a $\Delta\Delta G_{\text{folding}} < 5$. Therefore, for a candidate miRNA precursor structure to be considered legitimate, it would have to have either (i) $\Delta G_{\text{folding}} \leq -29$ or (ii) $-29 < \Delta G_{\text{folding}} \leq -22$ and $\Delta\Delta G_{\text{folding}} > 5$. In cases where both precursors fulfilled these criteria, the member of each pair with the lowest $\Delta G_{\text{folding}}$ was further considered. Precursors that had not been excluded up to this point were subsequently analyzed to determine whether they conformed to generally acceptable miRNA base-pairing standards (base-pairing involving at least 16 of the first 22 nucleotides of the miRNA and the other arm of the hairpin) (18).

Step 3: Determination of hairpin conservation. We classified all candidate miRNAs as either "conserved" or "nonconserved" by using the University of California at Santa Cruz phastCons database (28). This database has scores at each nucleotide in the human genome that correspond to the degree of conservation of that particular nucleotide in chimpanzee, mouse, rat, dog, chicken, pufferfish, and zebrafish. The algorithm is based on a phylogenetic hidden Markov model using best-in-genome pairwise alignment for each species (based on BLASTZ), followed by multialignment of the eight genomes. A hairpin was defined as conserved if the average phastCons conservation score over the seven species in any 15-nt sequence in the hairpin stem is at least 0.9 (5, 29).

Determination of Homology of Candidate miRNAs to Existing miRNAs.

One hundred random 22 mers were generated and compared to the miRBase database using the SSEARCH search algorithm, and expect values were obtained for each. *E* values for randomly generated sequences ranged from 0.07 to 23. All 133 miRNA candidates were subsequently analyzed, and tags with *E* values < 0.05 were deemed to have homology to existing miRNAs.

miRNA Microarray Expression Analysis. Five micrograms of total RNA from human placenta, prostate, testes, and brain (Ambion, Austin, TX) were size-fractionated (< 200 nt) by using the *mir*Vana kit (Ambion) and labeled with Cy3 (placenta and testes) and Cy5 (prostate and brain) fluorescent dyes. Pairs of labeled samples were hybridized to dual-channel microarrays. Microarray assays were performed on a μ ParaFlo microfluidics chip with each of the detection probes containing a nucleotide sequence of coding segment complementary to a specific microRNA sequence and a long nonnucleotide molecule spacer that extended the detection probe away from the substrate. The melting temperature of the detection probes was balanced by incorporation of varying number of modified nucleotides with increased binding affinities. The maximal signal level of background probes was 180. A miRNA detection signal threshold was defined as twice the maximal background signal.

Quantitative RT-PCR (qRT-PCR) Expression Analysis. qRT-PCRs were performed by using SuperTaq Polymerase (Ambion) and the *mir*Vana qRT-PCR miRNA Detection Kit (Ambion) following the manufacturer's instructions. Reactions contained custom-designed oligonucleotide DNA primers (Integrated DNA Technologies) specific for 36 novel putative miRNAs or *mir*Vana qRT-PCR Primer Sets specific for hsa-miR-16, hsa-miR-24, hsa-miR-143, or human 5S rRNA as positive controls. For each set of primers, 100 ng of FirstChoice human colon Tumor/Normal Adjacent Tissue RNA (Ambion); a pool containing 50 ng of HCT116, RKO, and DLD-1 cell lines total RNA; a pool containing 50 ng of FirstChoice Total RNA from human brain, cervix, thymus, and skeletal muscle (Ambion); and a no-template negative control were tested. All RNAs were treated with TURBO DNase. qRT-PCR was performed on an ABI7000 thermocycler (Applied Biosciences), and end-point reaction products were also analyzed on a 3.5% high-resolution agarose gel (Ambion) stained with ethidium bromide to discriminate between the correct amplification products (≈ 90 bp) and the potential primer dimers.

Targeted Disruption of the Human *Dicer* locus. The strategy for creating knockouts with AAV vectors was performed as described (30, 31). The targeting construct pAAV-Neo-*Dicer* was made by PCR, by using bacterial artificial chromosome clone CITB 2240H23 (Invitrogen) as the template for the homology arms. A targeted insertion was made in exon 5, which is part of the helicase domain. Details of the vector design and sequences of all PCR primers are available from the authors upon request. Stable G418-resistant clones were initially selected in the presence of Geneticin (Invitrogen), then routinely propagated in the absence of selective agents.

Determination of Differential Expression. Tag numbers from the different libraries were normalized and compared by using a Fisher exact test (significance threshold $P = 0.05$) with Bonferroni correction (32).

We thank Levy Kopelovich for helpful discussions and David Wong for sequence analysis and primer design. This work was supported by the Strang Cancer Prevention Center, the Division of Cancer Prevention of the National Cancer Institute, the Ludwig Trust, the Pew Charitable Trusts, and National Institutes of Health Grant CA057345.

