

Ab initio reconstruction of cell type-specific transcriptomes in mouse reveals the conserved multi-exonic structure of lincRNAs

Mitchell Guttman^{1,2,6}, Manuel Garber^{1,6}, Joshua Z Levin¹, Julie Donaghey¹, James Robinson¹, Xian Adiconis¹, Lin Fan¹, Magdalena J Koziol^{1,3}, Andreas Gnirke¹, Chad Nusbaum¹, John L Rinn^{1,3}, Eric S Lander^{1,2,4} & Aviv Regev^{1,2,5}

Massively parallel cDNA sequencing (RNA-Seq) provides an unbiased way to study a transcriptome, including both coding and noncoding genes. Until now, most RNA-Seq studies have depended crucially on existing annotations and thus focused on expression levels and variation in known transcripts. Here, we present Scripture, a method to reconstruct the transcriptome of a mammalian cell using only RNA-Seq reads and the genome sequence. We applied it to mouse embryonic stem cells, neuronal precursor cells and lung fibroblasts to accurately reconstruct the full-length gene structures for most known expressed genes. We identified substantial variation in protein coding genes, including thousands of novel 5' start sites, 3' ends and internal coding exons. We then determined the gene structures of more than a thousand large intergenic noncoding RNA (lincRNA) and antisense loci. Our results open the way to direct experimental manipulation of thousands of noncoding RNAs and demonstrate the power of *ab initio* reconstruction to render a comprehensive picture of mammalian transcriptomes.

A critical task in understanding mammalian biology is defining a precise map of all the transcripts encoded in a genome. Although much is known about protein coding genes in mammals, recent studies have suggested that the mammalian genome also encodes many thousands of large noncoding RNA (ncRNA) genes^{1–4}. Recently, we used a chromatin signature, combining histone-3 Lys4 trimethylation modifications (H3K4me3), known to mark promoter regions, and histone-3 Lys36 trimethylation modifications (H3K36me3), known to mark the entire transcribed regions (K4-K36 region; see **Supplementary Fig. 1**), to discover the genomic regions encoding ~1,600 lincRNAs in four mouse cell types⁴ and ~3,300 lincRNAs across six human cell types⁵.

Defining the complete gene structure of these lincRNAs is a prerequisite for experimental and computational studies of their function. We previously gained initial insights by hybridizing total RNA to tiling microarrays defined across the K4-K36 region⁴. This provided a coarse list of putative exonic locations but could not define the precise gene structures and exon connectivity.

Advances in RNA-Seq have opened the way to unbiased and efficient assays of the transcriptome of any mammalian cell^{6–10}. Recent studies in mouse and human cells have mostly focused on using RNA-Seq to study known genes^{6–8,10,11} and have depended on existing annotations. They were thus of limited utility for discovering the complete gene structure of lincRNAs or other noncoding transcripts.

An alternative strategy is to use an *ab initio* reconstruction approach^{9,12–14} to learn the complete transcriptome of an individual sample from solely the unannotated genome sequence and millions

of relatively short sequence reads. A complete *ab initio* transcriptome reconstruction of a sample will (i) identify all expressed exons; (ii) enumerate all the splicing events that connect them; (iii) combine them into transcriptional units; (iv) determine all isoforms, including alternative ends and (v) discover novel transcripts. A successful *ab initio* method should be applicable to large and complex mammalian genomes and should be able to reconstruct transcripts of variable sizes, expression levels and protein coding capacity.

Despite early successes in yeast⁹, *ab initio* reconstruction of a mammalian transcriptome has remained an elusive and substantial computational challenge. There has been important recent progress, including (i) efficient gapped aligners (for example, TopHat¹³) that can map short reads that span splice junctions ('spliced reads'); (ii) use of such gapped alignments to discover splicing events^{9,13}; (iii) exon identification methods¹⁴; and (iv) genome-independent assembly of unmapped reads to sequence contigs (for example, Abyss¹²). Each of these methods provides an important component toward reconstruction, but none can reconstruct the complete transcriptome of a mammalian cell, due to scaling issues⁹, limitations in handling splicing¹⁴ or inability to identify transcripts with moderate coverage¹².

Here we present Scripture, a comprehensive method for *ab initio* reconstruction of the transcriptome of a mammalian cell that uses gapped alignments of reads across splice junctions (exploiting recent increases in read length) and reconstructs reads into statistically significant transcript structures. We applied Scripture to RNA-Seq data from mouse embryonic stem cells (ESC), mouse neural progenitor

¹Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA. ²Department of Biology, Massachusetts Institute of Technology, Cambridge, Massachusetts, USA. ³Department of Pathology, Beth Israel Deaconess Medical Center, Boston, Massachusetts, USA. ⁴Department of Systems Biology, Harvard Medical School, Boston, Massachusetts, USA. ⁵Howard Hughes Medical Institute, Massachusetts Institute of Technology, Cambridge, Massachusetts, USA. ⁶These authors contributed equally to this work. Correspondence should be addressed to M. Guttman (mguttman@mit.edu), M. Garber (mgarber@broadinstitute.org) or A.R. (aregev@broad.mit.edu).

Received 10 March; accepted 6 April; published online 2 May 2010; corrected after print 9 July 2010; doi:10.1038/nbt.1633

cells (NPC) and mouse lung fibroblasts (MLF) and correctly identified the complete annotated full-length gene structures for most expressed, known, protein coding genes. The reconstruction of the three transcriptomes revealed substantial variation in protein coding genes between cell types, including thousands of novel 5' start sites, 3' ends or additional coding exons. Many of these variant structures are supported by independent data. We also discovered the gene structure and expression level of over 2,000 noncoding transcripts, including hundreds of transcripts from previously identified lincRNA loci, over a thousand more lincRNAs with similar properties and hundreds of multi-exonic antisense ncRNAs. We show that lincRNAs have no significant coding potential and that they are evolutionary conserved. Our results open the way to direct experimental manipulation of this new class of genes and highlight the power of RNA-Seq along with an *ab initio* reconstruction to provide a comprehensive picture of cell-specific transcriptomes.

RESULTS

RNA-Seq libraries

We used massively parallel (Illumina) sequencing to sequence cDNA libraries from poly(A)⁺ mRNA from ESC, NPC and MLF cells, with

76-base paired-end reads. For the ESC library, we generated a total of 152 million paired-end reads. Using a gapped aligner¹³, 93 million of these were alignable (497 Mb aligned bases, 262-fold average coverage of known protein coding genes expressed in ESC). We obtained similar numbers for the NPC and MLF libraries (Online Methods). In ESC, 76% of these reads mapped within the exonic regions of known protein coding genes, 9% were in introns of known protein coding genes, and 15% mapped in intergenic regions. We found a strong correlation between expression levels of protein coding genes as measured by RNA-Seq and Affymetrix expression arrays ($r = 0.88$ for all genes; **Supplementary Fig. 2**).

Scripture: a method for transcriptome reconstruction

We next developed Scripture, a genome-guided method to reconstruct the transcriptome using only an RNA-Seq data set and an (unannotated) reference genome sequence. Scripture consists of five steps (**Fig. 1**, **Supplementary Note 1** and Online Methods). (i) We use reads aligned to the genome, including those with gapped alignments¹³ spanning exon-exon junctions ('aligned spliced reads', **Fig. 1a,c**). 'Spliced' reads provide direct information on the location of splice

Figure 1 Scripture: a method for *ab initio* transcriptome reconstruction from RNA-Seq data. **(a)** Spliced and unspliced reads. A typical expressed four-exon gene (*1500032D16Rik*, top; exons, gray boxes) with coverage from different type of reads. Unspliced reads (black bars) fall within a single exon, whereas spliced reads (bars broken into 'dumbbells') span exon-exon junctions (thin horizontal lines connect the alignment of a read to the exons it spans). The coverage track (bottom) shows the aggregate coverage of both spliced and unspliced reads. **(b–g)** A schematic description of Scripture. **(b)** A cartoon example. Reads (black bars) originate from sequencing a contiguous RNA molecule. Shown are transcripts from two different genes (blue and red boxes), one with seven exons (blue boxes) and one with three exons (red boxes), which are adjacent in the genome (black line). The grayscale vertical shading in subsequent panels is shown for visual tracking. **(c)** Spliced reads. Scripture is initiated with a reference genome sequence and spliced aligned reads (dumbbells) with gaps in their alignment (thin horizontal lines). Scripture uses splice site information to orient spliced reads (arrowheads). **(d)** Connectivity graph construction. Scripture builds a connectivity graph by drawing an edge (curved arrow) between any two bases that are connected by a spliced read gap. Edges are color coded to relate to the original RNA and eventual transcript. **(e)** Path scoring. Scripture scans the graph with fixed-sized windows and uses coverage from all reads (spliced and unspliced; bottom track) to score each path for significance (*P*-values shown as edge labels). **(f)** Transcript graph construction. Scripture merges all significant windows and uses the connectivity graph to give significant segments a graph structure (three graphs, in this example). **(g)** Refinement with paired-end data. Scripture uses paired-end (dashed curved lines) to join previously disconnected graphs (gene 1, bold dashed line), find breakpoint regions within contiguous segments (detectable in this example by the lack of dashed lines between genes 1 and 2) and eliminate isoforms that result in paired-end reads mapping at a distance with low likelihood.

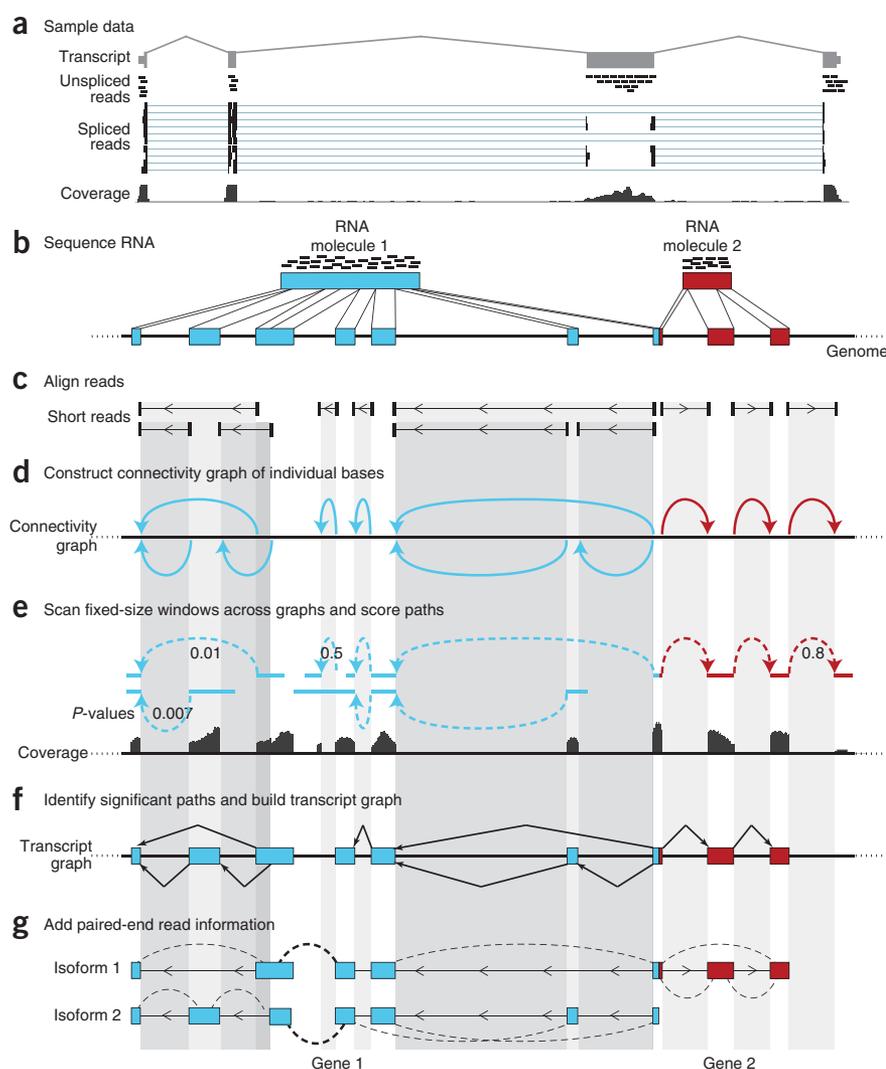
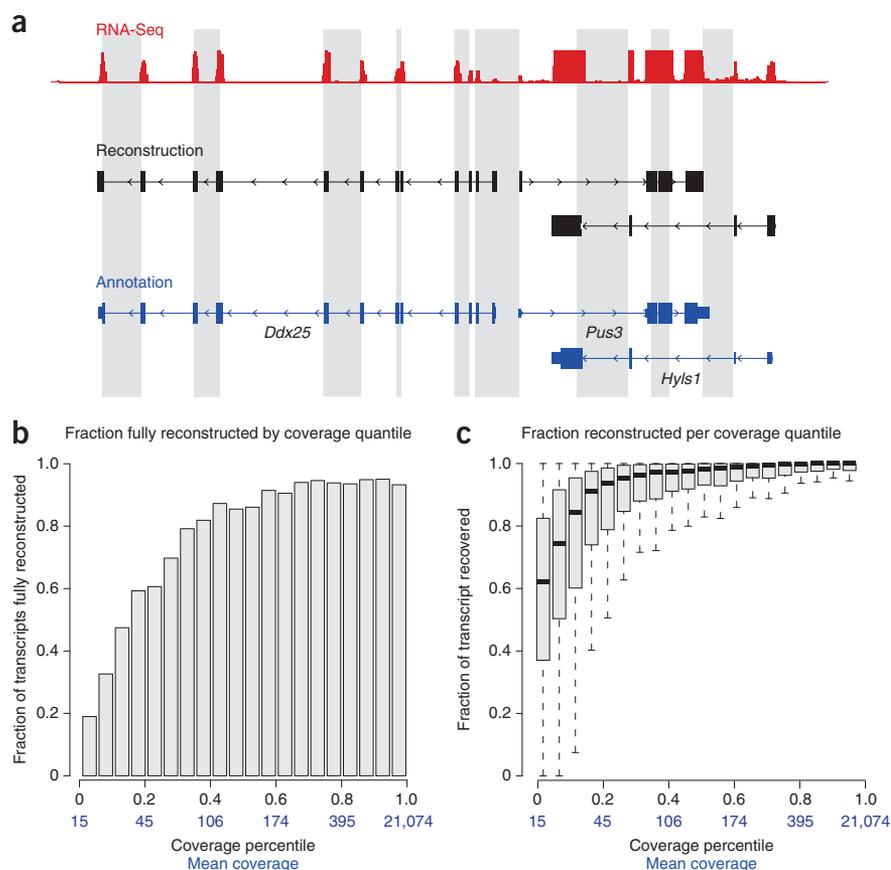


Figure 2 Scripture correctly reconstructs full-length transcripts for most annotated protein coding genes. (a) A typical Scripture reconstruction on mouse chromosome 9. Top, RNA-Seq read coverage (from both unspliced and spliced reads); middle, three transcripts reconstructed by Scripture, including exons (black boxes) and orientation (arrow heads); bottom, RefSeq annotations for this region. All three transcripts are fully reconstructed from 5' to 3' ends, capturing all internal exons; Scripture correctly reconstructed the overlapping transcripts *Pus3* and *Hlys1*. (b) Fraction of genes fully reconstructed in different expression quantiles (in 5% increments) in ESC. Each bar represents a 5% quantile of read coverage for genes expressed; mean read coverage is noted in blue. The height of each bar is the fraction of genes in that quantile that were fully reconstructed. For example, ~20% of the transcripts at the bottom 5% of expression levels were fully reconstructed; ~94% of the genes at the top 95% of expression are fully reconstructed. (c) Portion of gene length reconstructed in different expression quantiles in ESC. Shown is a box plot of the portion of each transcript's length that was covered by a Scripture reconstruction in each 5% coverage quantile. Black line in each box, median; rectangle, 25%–75% coverage quantiles; whiskers, extreme coverage values within expression quantile. For example, at the bottom 5% of expression, Scripture reconstructed a median length of 60% of the full length transcript.



junctions within the transcript, and ~30% of 76-base reads are expected on average to span an exon–exon junction. From the aligned spliced reads, we construct a ‘connectivity graph’ (Fig. 1d), where two bases in the genome are connected if they are immediate neighbors either in the genomic sequence itself or within a spliced read. We use agreement with splicing motifs at each putative junction to orient the connection (edge) in the connectivity graph^{9,13} (Fig. 1d). (ii) To infer transcripts, we use a statistical segmentation approach⁴ and both spliced and unspliced reads to identify paths in the connectivity graph with mapped read enrichment compared to the genomic background (Fig. 1e). This is done by scoring a sliding window using a test statistic for each region, computing a threshold for genome-wide significance, and using the significant windows to define intervals. (iii) From the paths, we construct a ‘transcript graph’ connecting each exon in the transcript (Fig. 1f). Each path through the graph is directed and represents one oriented (strand-specific) isoform of the gene. Alternative spliced isoforms are identified by considering all possible paths in the transcript graph. (iv) We augment the transcript graph with connections based on paired-end reads and their distance constraints, allowing us to join transcripts or remove unlikely isoforms (Fig. 1g, below). (v) We generate a catalog of transcripts defined by the paths through the transcript graph.

Paired-end reads aid in transcriptome reconstruction

Paired-end information, consisting of reads that came from the two ends of the sequenced RNA fragment, provides valuable additional information in the reconstruction.

First, the presence of paired ends linking two regions shows that they appear in the same transcript; such a connection might not otherwise be apparent because low expression levels or unalignable sequence might prevent a continuous chain of overlapping sequence

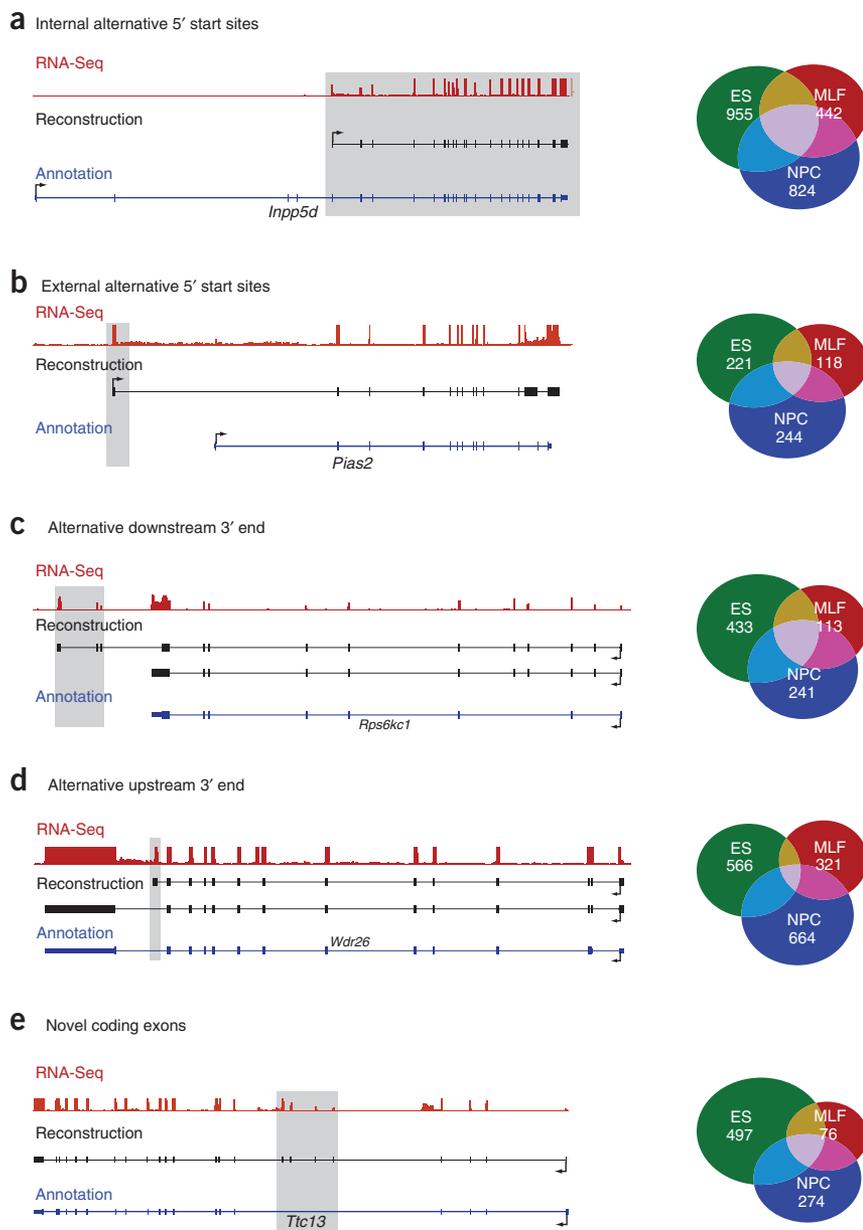
reads (spliced or unspliced) across the transcript. We thus augment the transcript graphs with paired-end information, where available, to (indirectly) link nodes in the graph. We use these indirect links (Fig. 1g) to add edges between disconnected graphs, add internal nodes (exons) that might have been missed within a path (transcript) and add extra support for existing edges. This refines the structure of our transcripts and increases our confidence in them, especially in weakly expressed transcripts, which are more likely to have coverage gaps.

Second, the distribution of library insert sizes constrains the distance between the paired-end reads; these distance constraints can be used to infer the relative likelihood of some potential transcripts (for example, those in which the paired ends would be much closer or much further than expected). We infer the distribution of insert sizes for a given library from the position of read pairs on transcripts from those genes for which there is only a single transcript model (that is, no detectable alternative splicing; Online Methods). For example, in the ESC library, this distribution matches well with the experimentally determined sizes. Using this distribution, we assign likelihoods to each connection, filtering unlikely ones (Online Methods).

Reconstruction of full-length gene structures

We applied Scripture to our mouse ESC RNA-Seq data set and compared our reconstructions to protein coding gene annotations¹⁵. Scripture identified 16,389 nonoverlapping, multi-exonic transcript graphs that correspond to 15,352 known multi-exonic genes (Online Methods). Of reconstructed genes, 88.4% are covered by a single graph (no fragmentation of the reconstructed transcript) and 8.0% are covered by two transcript graphs (fragmentation of the transcript to two separate pieces in the reconstruction). Focusing on the 13,362 genes with significant expression ($P < 0.05$ compared with

Figure 3 Alternative 5' ends, 3' ends and novel coding exons in transcripts reconstructed by Scripture. Representative examples (tracks, left) and summary counts (Venn diagrams, right) numbers represent those unique to each cell type compared to other two) of five categories of variation discovered in Scripture transcripts compared to the known annotations. In each representative example, shown is the coverage by RNA-Seq reads (top track), the reconstructed annotation (middle track) and the known annotation (bottom track). The novel regions in the reconstruction are marked by gray shading. In each proportional Venn diagram we show the number of transcripts in this class in each cell type (ESC, green; NPC, blue; MLF, red) and their overlap. (a) Internal alternative 5' start sites. (b) External alternative 5' start sites. (c) Alternative downstream 3' end (extended termination). (d) Alternative upstream 3' end (early termination). (e) Novel coding exons.



background coverage; see Online Methods), Scripture reconstructed the full-length structure of the longest known splice isoform (from 5' to 3' end, including all exons and splice junctions; **Fig. 2a**) for 10,355 of them (~78%). All of our reconstructed transcripts for known multi-exonic transcripts also had the correct orientation (strand), allowing us to reconstruct genes that overlap one another on opposite strands (**Fig. 2a**).

Complete transcript structures were recovered across a very broad range of expression levels (**Fig. 2b,c**) for both single and multi-exonic genes. For example, Scripture accurately reconstructed the full-length transcript of ~73% of the known protein coding genes at the second quintile of expression, and ~94% of the genes from the top quintile. Furthermore, the average proportion of bases reconstructed for each transcript was high (**Fig. 2c**). Even for the bottom 5% of expressed genes, we recovered on average 62% of each of these transcripts' bases (**Fig. 2c**). For single-exon genes, we recovered on average 80% of the transcribed bases. We obtained similar results in the other two cell types (19,835 and 20,407 transcript graphs for 14,212 and 13,351 known genes in NPC and MLF, respectively). Most of the genes that were not fully reconstructed are those with low expression; it should be possible to reconstruct most of these by generating more RNA-Seq data. The few highly expressed genes that were not fully reconstructed are either the result of alignment artifacts caused by recent processed pseudogenes or stem from novel transcriptome variations, missing from the current annotation (explored in detail below).

Novel transcriptome variations in annotated protein coding genes

Given that most of the Scripture reconstructions of protein coding genes were accurate, we next investigated the differences between the reconstructed transcriptome and the known gene annotations (**Supplementary Table 1**). We focused on transcripts with (i) novel 5' start sites; (ii) novel 3' ends; and (iii) previously unidentified exons within the transcriptional units of known protein coding genes.

In each category, we first discuss below the reconstructed transcripts in ESC and then consider the results for the NPC and MLF.

1. Alternative 5' start sites are supported by H3K4me3 marks. We found 1,804 transcripts in ESC that match the annotated 3' end but have an alternative 5' start site, derived from an extra exon not overlapping the annotated first exon. We distinguish between internal alternative 5' start sites (1,397 cases; **Fig. 3a**), which occur downstream of the annotated start, and external alternative 5' start sites (407 cases; **Fig. 3b**), which occur upstream of the annotated start. Ninety percent of the internal 5' start sites and 75% of the external 5' start sites contained an H3K4me3 modification, a mark of the promoter region of genes¹⁶ (**Supplementary Fig. 3**). These alternative start sites are on average 21 kb upstream of the annotated site, substantially revising the annotated promoters. Notably, ~60% of the transcripts with an alternative start site (internal or external) had no reconstructed isoform starting at the annotated 5' start site.

We obtained similar results from NPC and MLF (**Fig. 3a,b**, right; **Supplementary Table 1**). Altogether, we identified 2,813 internal

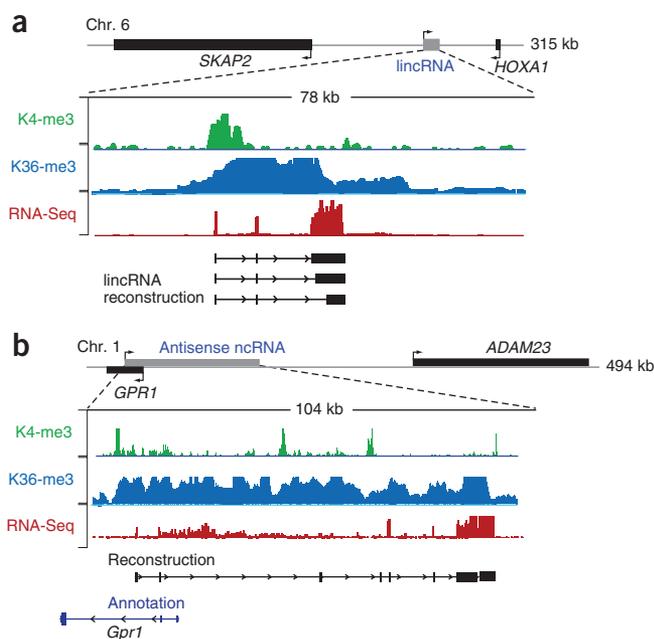


Figure 4 Noncoding transcripts reconstructed by Scripture. (a) A representative example of a lincRNA expressed in ESC. Top: mouse genomic locus containing the lincRNA and its neighboring protein coding genes. Bottom: magnified view of the lincRNA locus showing the coverage of H3K4me3 (green track), H3K36me3 (blue track) and RNA-Seq reads (red track) overlapping the transcribed lincRNA locus, as well as its Scripture reconstructed transcript isoforms (black). (b) A representative example of a multi-exonic antisense ncRNA expressed in ESC. Top: mouse genomic locus containing the antisense transcript. Bottom: magnified view of the antisense locus showing the coverage of H3K4me3 (green track), H3K36me3 (blue track) and RNA-Seq reads (red track) overlapping the transcribed antisense locus, as well as its Scripture reconstructed gene structure (black) and the annotated overlapping transcript (blue).

5' start sites (2,302 supported by H3K4me3 in their respective tissues), and 807 external 5' start sites in at least one cell type. In particular, 33% of these novel 5' ends are likely active in ESCs but not in NPCs or MLFs.

2. Alternative 3' untranslated regions are supported by polyadenylation motifs. There are 551 (~4%) ESC-reconstructed transcripts with an alternative 3' end downstream of the annotated 3' end (mean distance 30 kb downstream, Fig. 3c). Of these, 275 (~50%) showed evidence of a polyadenylation motif within the novel 3' exon, which is only slightly lower than for annotated 3' ends (60%) and much higher than for randomly chosen size-matched exons (6%). The frequency of the polyadenylation motif supports the accuracy of the reconstruction.

To conservatively distinguish between upstream (early) termination and incomplete reconstruction, we designated novel 3' ends only in those cases that did not overlap any of the known exons in the annotated transcript and that contained complete 5' start sites. We identified 759 transcripts with upstream 3' ends in ESC (Fig. 3d); 44% of them contained a polyadenylation motif, supporting their biological relevance. For most (90%) of these transcripts, Scripture also reconstructed an isoform that contained the annotated 3' end.

We obtained similar results for NPC and MLF (Fig. 3c,d, right; Supplementary Table 1). Altogether, we identified 940 downstream 3' ends and 1,850 upstream 3' ends in at least one cell type.

3. Additional coding exons are highly conserved and preserve ORFs. We found 534 transcripts in ESC with at least one extra,

previously unannotated internal coding exon spliced into annotated protein coding transcripts (Fig. 3e). These transcripts contained 588 novel internal exons, ranging in length from 6 bp to 3.5 kb (median, 111 bp; 20–80% quantiles, 60–224 bp). Of these extra exons, 322 (54.5%) were present in all versions of the reconstructed transcript in ESC. Most (83%) of these novel exons maintain the reading frame of the transcript and are as highly conserved as known coding exons (Supplementary Fig. 4), consistent with their coding capacity. We validated the presence of the novel exons within five of five tested transcripts, using reverse transcription followed by PCR (RT-PCR) followed by Sanger sequencing (Online Methods).

We obtained similar results in MLF (124 transcripts, 144 exons) and NPC (325 transcripts, 363 exons) (Fig. 3e, right). A majority of exons (~70%) were present in all versions of the reconstructed transcript within a cell type. Altogether, we identified 960 novel internal exons in at least one cell type (Fig. 3e, right).

Gene structures of previously identified lincRNA loci

We next turned to identifying the gene structures of transcripts expressed from known lincRNAs loci. We had previously identified 317 lincRNA loci on the basis of K4-K36 domains in ESC cells⁴. When applied to ESC RNA-Seq data, Scripture reconstructed multi-exonic gene structures for 250 (78.8%) of them (Fig. 4). This is comparable to the proportion (78.5%) reconstructed for protein coding genes with K4-K36 domains in ESC. Scripture reconstructed 87% (160 of 183) of ESC lincRNAs for which we previously identified an RNA hybridization signal from tiling microarrays. We discuss possible reasons for the few remaining discrepancies in Supplementary Note 2.

The reconstructed lincRNA transcripts in ESC have on average 3.7 exons, an average exon size of 350 bp and an average mature spliced size of 3.2 kb (in comparison, protein coding genes have on average 9.7 exons, exon length of 291 bp and length of 2.9 kb). The Scripture-identified strand information for each lincRNA is consistent with that inferred from the location of H3K4me3 modification and with the orientation determined from a strand-specific RNA-Seq library which we generated independently (Online Methods). Most lincRNAs likely represent 5' complete transcripts based on overlap with H3K4me3 (82%) and 3' complete transcripts based on presence of a polyadenylation motif (~50%, comparable to 60% for protein coding genes and far above background of 6%).

Similarly, Scripture successfully reconstructed lincRNA gene structures for K4-K36 lincRNA loci in MLF and NPC (232 of 289 in MLF and 224 of 270 in NPC). Most are likely 5' complete (69% in MLF and 81% in NPC based on overlap with H3K4me3) and many may be 3' complete based on detectable 3' polyadenylation sites (18% in MLF and 37% in NPC). In addition, we successfully reconstructed another 116 lincRNAs previously identified only in mouse embryonic fibroblasts but which were now reconstructed in at least one of the other three cell types. Altogether, we identified gene structures for 609 previously defined lincRNA loci in at least one of the three cell types.

Discovery of novel lincRNAs

In addition to the previously identified lincRNAs, we found another 1,140 multi-exonic transcripts that map to intergenic regions (591 in ESC, 318 in MLF, and 528 in NPC; Fig. 5). Most of these transcripts do not seem to encode proteins, and are designated as noncoding, on the basis of their codon substitution frequency (CSF) scores^{17,18} (Online Methods) across the mature (spliced) RNA transcript (88%; Fig. 5a) and on the lack of an open reading frame (ORF) larger than 100 amino acids (80%; Fig. 5b). Careful review of the

remaining ~12% revealed 66 loci that are likely to be novel protein coding genes (high CSF score, ORF >200 amino acids and very high evolutionary conservation; **Supplementary Fig. 5**).

Most of the novel lincRNA loci were not identified in our previous study owing to the stringent criteria we imposed when using chromatin maps to identify lincRNAs. Specifically, we required that a K4-K36 domain extend over at least 5 kb and be well separated from the nearest known gene locus⁴. Indeed, most novel intergenic transcripts (76%) were enriched for a K4-K36 domain (a comparable proportion as that for expressed protein coding genes) but failed to meet the size and distance criteria or could not be identified at a genome-wide significance level (without knowing their locus a priori). On average, the genomic loci of the novel lincRNAs are closer to neighboring genes and have smaller sizes (~3.5 kb average), and the transcripts are shorter (859 bp). Of the lincRNAs that did not have a chromatin signature that reached genome-wide significance, ~40% showed chromatin modifications enriched at a nominal significance level (compared to 57% for protein coding genes).

On average, the lincRNAs are expressed at levels that are readily detectable, albeit somewhat lower than those of protein coding genes. The median expression level of the reconstructed lincRNAs, as estimated by reads per kilobase of exonic sequence per million aligned reads (RPKM; see Online Methods) was approximately one-third of the expression of protein coding genes (**Fig. 5d**), with ~25% of lincRNAs having expression levels higher than the median level for protein coding genes (**Fig. 5d**). The novel lincRNAs identified in this study are expressed at somewhat lower levels than those from chromatin identified loci, consistent with the fact that chromatin enrichment is positively correlated with expression levels (**Fig. 5d**).

We compared the novel lincRNA genes to a collection of ~35,000 mouse cDNA and found evidence that ~43% of our lincRNAs were present in this collection¹. This is comparable to the reported fraction (40%) of known transcripts covered by the same cDNA catalog¹. The remaining lincRNAs are found in our study but not in the comparison catalog. These were likely previously missed owing to the different cell types and limited coverage of the previous study¹.

Most lincRNAs are evolutionarily conserved

The reconstructed full-length gene structures of lincRNAs allow us to accurately assess their evolutionary sequence conservation in each exon and in small windows. To this end, we identified the orthologous sequences for each lincRNA across 29 mammals and estimated conservation by a metric (ω ; Online Methods) reflecting the total contraction of the branch length of the evolutionary tree connecting them¹⁹. We calculated ω over the entire lincRNA transcript, as well as over individual exons.

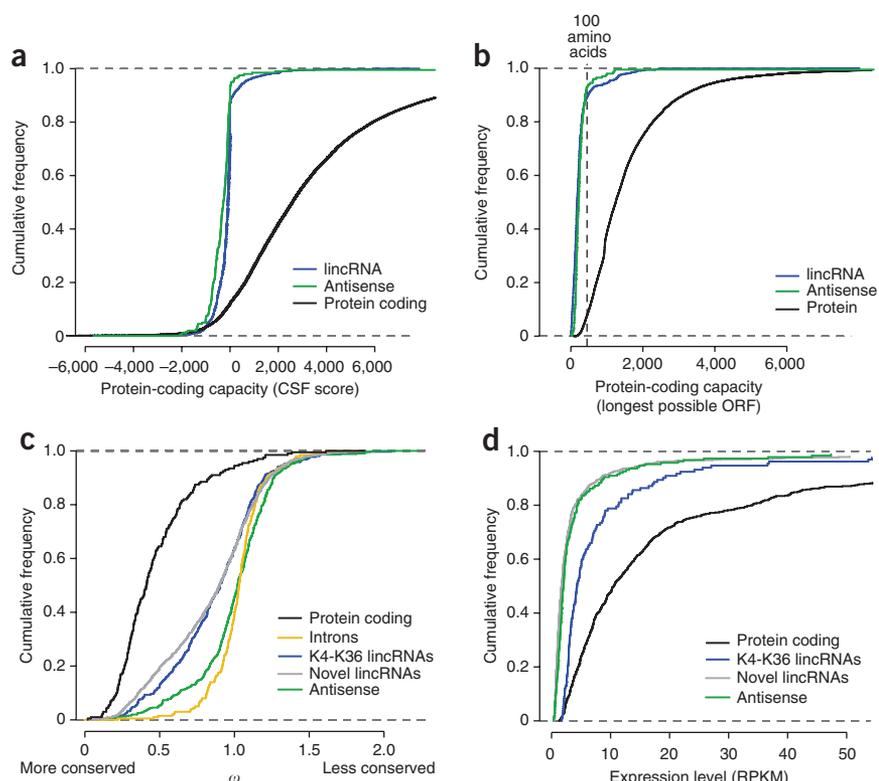


Figure 5 Protein coding capacity, conservation levels and expression of lincRNAs and multi-exonic antisense transcripts. (**a,b**) Coding capacity of protein coding, lincRNAs and multi-exonic antisense transcripts. Shown is the cumulative distribution of CSF scores (**a**) and maximal ORF length (**b**) for protein coding transcripts, lincRNAs and multi-exonic antisense transcripts. (**c**) Conservation levels for exons from protein coding transcripts, lincRNAs, multi-exonic antisense transcripts and introns. Shown is the cumulative distribution of sequence conservation across 29 mammals for exons from protein coding exons, introns, exons from previously annotated lincRNA loci, exons from newly annotated lincRNA transcripts and exons from multi-exonic antisense transcripts. (**d**) Expression levels of protein coding, lincRNAs and multi-exonic antisense transcripts. Shown is the cumulative distribution of expression levels, in reads per kilobase of exonic sequence per million aligned reads (RPKM) in ESC for protein coding transcripts, transcripts from previously annotated lincRNA loci, transcripts from newly annotated lincRNA loci and multi-exonic antisense transcripts.

On the basis of our high-resolution gene structures, the lincRNA sequences show greater conservation than random genomic regions or introns (**Fig. 5c**), comparable to eight known functional lincRNAs^{20–22}, and lower than protein coding exons. The results are consistent with our previous estimates of conservation⁴. Interestingly, conservation levels are indistinguishable between the chromatin-defined lincRNAs⁴ and the novel ones identified only in this study (**Fig. 5c**), consistent with membership in the same class of functional large ncRNA genes. These conservation levels are considerably higher than those reported for a previous catalog of large noncoding RNAs¹.

We also determined the specific regions within each lincRNA that are under purifying selection and thus likely to be functional, by computing ω within short windows (Online Methods). On average, 22% of the bases within the lincRNAs lie within conserved patches (comparable to the value of 25% for the eight known functional lincRNAs, much higher than the 7% for intronic bases and lower than the 77% for protein coding bases, **Supplementary Fig. 6**). These conserved patches provide a critical starting point for functional studies²³.

Variations in lincRNA expression and isoforms

A substantial fraction (~41%) of the novel lincRNAs reconstructed in at least one cell type show evidence for expression in at least two of the

three cell types. This is comparable to the 45% of the previously identified lincRNAs present in at least two out of the three cell types. In contrast, 80% of expressed protein coding genes are expressed across two of the three cell types. This is not merely a result of the lower overall expression of lincRNAs, as the fraction of cell type-specific lincRNAs is higher than that of tissue specific protein coding genes in every expression quantile (Supplementary Fig. 7). Thus, lincRNAs are likely to be more tissue specific than protein coding genes.

A substantial portion of lincRNA loci also produce alternative spliced isoforms. For example, within ESC we identified two or more alternative spliced isoforms for 25% of lincRNA genes, comparable to the 30% for protein coding genes (15% of lincRNAs in MLF have alternative spliced isoforms, and 14.7% in NPC). Altogether, 28.8% of the 1,749 lincRNA loci had evidence for alternative isoforms in any of the three cell types.

Identification of hundreds of large antisense transcripts

Scripture reconstructed hundreds of transcripts that overlap known protein coding gene loci but are transcribed in the opposite orientation and likely represent antisense transcripts. To determine orientation, we required that any identified antisense transcript be multi-exonic (Online Methods).

Using these criteria, we identified 201 antisense multi-exonic transcripts in ESC (Fig. 4b); these transcripts had an average five exons per transcript and an average transcript size of 1.7 kb. On average, the antisense transcripts overlapped the genomic locus of the sense protein coding gene by 1,023 bp (83% of the locus), and most (64%) overlapped at least one sense exon, but this overlap was substantially lower (766 bp, 48% of the transcript exons). Some of these antisense transcripts (79, ~40%) were identified by a previous cDNA sequencing study^{1,24}, but most (122, ~60%) were previously unidentified. Most (~85%) antisense transcripts were non-protein coding by both ORF analysis (Fig. 5b) and CSF scores (Fig. 5a). Four of the newly identified antisense transcripts had a large, conserved open reading frame and are likely novel, previously unannotated protein coding genes.

We validated the reconstructed ESC antisense transcripts by three independent sets of experimental data. (i) Most of the antisense loci carried an H3K4me3 mark at their 5' end (Fig. 4b), consistent with their independent and antisense transcription (for example, 64% of the 164 transcripts where it was possible to detect an independent H3K4me3 mark because the 5' end of the antisense transcript did not overlap the 5' ends of the sense gene). (ii) We generated and sequenced a strand-specific library in ESC (17.5 million Illumina reads; Online Methods), and found a significant ($P < 0.05$) number of reads on the antisense strand in >90% of cases (the remaining are likely missed in this limited sequencing owing to lower expression). (iii) We confirmed five of five tested antisense transcripts using RT-PCR to unique exons of the antisense transcript (Online Methods) followed by Sanger sequencing.

We obtained similar results for antisense transcripts in MLF and NPC (112 and 202 multi-exonic antisense transcripts, respectively). Altogether, we identified 469 antisense transcripts expressed in at least one cell type, only 125 of which (27%) were previously identified in large-scale sequencing of mouse cDNAs²⁴. The remaining 344 (73%) were unidentified by the previous study, likely reflecting the distinct cell types used in that study and the limited coverage of previous catalogs.

The 469 antisense transcripts are expressed at levels comparable to those of the novel lincRNAs (Fig. 5d) but show substantially lower sequence conservation. Indeed, the antisense ncRNAs showed very little evolutionary conservation as estimated by the ω metric for the

portions that do not overlap protein coding exons on the sense strand, suggesting that the antisense ncRNAs are a distinct class from the lincRNAs (Fig. 5c).

DISCUSSION

Despite the availability of the genome sequence of many mammals, a comprehensive understanding of the mammalian transcriptome has been an elusive goal. In particular, the computational tools needed to reconstruct all full-length transcripts from the wealth of short read data were largely missing. A recent study proposed to overcome this limitation experimentally by using very long reads (for example, 454 sequencing) as a scaffold for short read reconstruction²⁵. This is applicable, albeit at a substantial cost, for highly expressed genes but would require extraordinary depth to cover more weakly expressed ones.

Here we present Scripture, a new computational method to reconstruct a mammalian transcriptome with no prior knowledge of gene annotations. Scripture relies on longer reads that span splice junctions to connect discontinuous (spliced) segments and resolve multiple splice isoforms, and uses paired-end information to refine these transcripts. Scripture can identify short but strongly expressed transcripts as well as transcripts with much lower expression for which there is aggregate evidence along the entire transcript length. Although Scripture does rely on a reference genome sequence, many of its components can also be used in the development of methods for assembly of transcripts from read data only.

We applied Scripture to RNA-Seq data from pluripotent ESCs and differentiated lineages and showed that we can accurately reconstruct most expressed, annotated protein coding genes, at a broad range of expression levels, as well as uncover many new isoforms in the protein coding transcriptome. This variation may have key regulatory roles, defining new cell type-specific promoters, untranslated regions and protein coding exons. We used Scripture's sensitivity and resolution to reconstruct the gene structures and strand information of hundreds of lincRNAs and multi-exonic antisense transcripts, many of which are only moderately expressed.

Scripture identified over a thousand lincRNAs across the three cell types studied. Most of the lincRNAs identified were not previously found by classical large-scale cDNA sequencing¹. Many of these lincRNAs could not be reliably identified solely on the basis of chromatin structure owing to their proximity to protein coding genes or their short genomic lengths. Overall, we found that the ratio of expressed protein coding to noncoding genes in these cell types was ~10:1 but that the total number of RNA molecules was more heavily biased toward the protein coding fraction (~30:1), results similar to previous observations²⁶.

Scripture identifies precise gene structures for most previously found lincRNA loci (as well as for the newly discovered ones), a prerequisite for further studies. For example, we used these to identify the specific regions within each lincRNA that are under purifying selection (conservation), a starting point for experimental and computational investigation.

Taken together, our results highlight the power of *ab initio* reconstructions to annotate a genome, to discover transcriptional variation within known protein coding genes and to provide a rich catalog of precise gene structures for noncoding RNAs. The next step is clearly to apply this approach to a wide range of mammalian cell types, to obtain a comprehensive picture of the mammalian transcriptome.

METHODS

Methods and any associated references are available in the online version of the paper at <http://www.nature.com/naturebiotechnology/>.

Accession codes. NCBI Gene Expression Omnibus (GEO), GSE20851.

Note: Supplementary information is available on the Nature Biotechnology website.

ACKNOWLEDGMENTS

We thank M. Wernig (MIT) for providing NPC; M. Lin and M. Kellis (MIT) for CSF code; the Broad Sequencing Platform for sample sequencing; L. Gaffney for assistance with graphics; and C. Burge, J. Merkin, R. Bradley and members of Lander and Regev laboratories—in part, M. Yassour, T. Mikkelsen and I. Amit—for discussions. A.R. and J.L.R. were supported by the Merkin Family Foundation for Stem Cell Research at the Broad Institute. M. Guttman was supported by a Vertex scholarship. Work was supported by a Burroughs Wellcome Fund Career Award at the Scientific Interface, a US National Institutes of Health PIONEER award, a US National Human Genome Research Institute (NHGRI) R01 grant and the Howard Hughes Medical Institute (A.R.), and NHGRI and the Broad Institute of MIT and Harvard (E.S.L.).

AUTHOR CONTRIBUTIONS

M. Guttman and M. Garber conceived the project, designed research, implemented Scripture, performed computational analysis and wrote the paper. A.G., C.N. and J.Z.L. oversaw cDNA sequencing, provided molecular biology advice and helped to edit the manuscript. J.D. constructed cDNA libraries, performed validation experiments and helped to edit the manuscript. J.R. implemented components of Scripture and provided computational support and technical advice. X.A., L.F. and M.J.K. constructed cDNA libraries. J.L.R. provided reagents and helped edit the manuscript. E.S.L. designed research direction and wrote the paper. A.R. provided cDNA sequencing guidance, conceived the project, designed research direction and wrote the paper.

COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

Published online at <http://www.nature.com/naturebiotechnology/>.

Reprints and permissions information is available online at <http://npg.nature.com/reprintsandpermissions/>.

- Carninci, P. *et al.* The transcriptional landscape of the mammalian genome. *Science* **309**, 1559–1563 (2005).
- Kapranov, P. *et al.* RNA maps reveal new RNA classes and a possible function for pervasive transcription. *Science* **316**, 1484–1488 (2007).
- Bertone, P. *et al.* Global identification of human transcribed sequences with genome tiling arrays. *Science* **306**, 2242–2246 (2004).
- Guttman, M. *et al.* Chromatin signature reveals over a thousand highly conserved large non-coding RNAs in mammals. *Nature* **458**, 223–227 (2009).
- Khalil, A.M. *et al.* Many human large intergenic noncoding RNAs associate with chromatin-modifying complexes and affect gene expression. *Proc. Natl. Acad. Sci. USA* **106**, 11667–11672 (2009).
- Cloonan, N. *et al.* Stem cell transcriptome profiling via massive-scale mRNA sequencing. *Nat. Methods* **5**, 613–619 (2008).
- Wang, E.T. *et al.* Alternative isoform regulation in human tissue transcriptomes. *Nature* **456**, 470–476 (2008).
- Mortazavi, A., Williams, B.A., McCue, K., Schaeffer, L. & Wold, B. Mapping and quantifying mammalian transcriptomes by RNA-Seq. *Nat. Methods* **5**, 621–628 (2008).
- Yassour, M. *et al.* Ab initio construction of a eukaryotic transcriptome by massively parallel mRNA sequencing. *Proc. Natl. Acad. Sci. USA* **106**, 3264–3269 (2009).
- Pan, Q., Shai, O., Lee, L.J., Frey, B.J. & Blencowe, B.J. Deep surveying of alternative splicing complexity in the human transcriptome by high-throughput sequencing. *Nat. Genet.* **40**, 1413–1415 (2008).
- Maher, C.A. *et al.* Transcriptome sequencing to detect gene fusions in cancer. *Nature* **458**, 97–101 (2009).
- Birrol, I. *et al.* De novo transcriptome assembly with ABySS. *Bioinformatics* **25**, 2872–2877 (2009).
- Trapnell, C., Pachter, L. & Salzberg, S.L. TopHat: discovering splice junctions with RNA-Seq. *Bioinformatics* **25**, 1105–1111 (2009).
- Denoeud, F. *et al.* Annotating genomes with massive-scale RNA sequencing. *Genome Biol.* **9**, R175 (2008).
- Pruitt, K.D., Tatusova, T. & Maglott, D.R. NCBI reference sequences (RefSeq): a curated non-redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids Res.* **35**, D61–D65 (2007).
- Mikkelsen, T.S. *et al.* Genome-wide maps of chromatin state in pluripotent and lineage-committed cells. *Nature* **448**, 553–560 (2007).
- Lin, M.F., Deoras, A.N., Rasmussen, M.D. & Kellis, M. Performance and scalability of discriminative metrics for comparative gene identification in 12 *Drosophila* genomes. *PLOS Comput. Biol.* **4**, e1000067 (2008).
- Lin, M.F. *et al.* Revisiting the protein-coding gene catalog of *Drosophila melanogaster* using 12 fly genomes. *Genome Res.* **17**, 1823–1836 (2007).
- Garber, M. *et al.* Identifying novel constrained elements by exploiting biased substitution patterns. *Bioinformatics* **25**, i54–i62 (2009).
- Brown, C.J. *et al.* A gene from the region of the human X inactivation centre is expressed exclusively from the inactive X chromosome. *Nature* **349**, 38–44 (1991).
- Rinn, J.L. *et al.* Functional demarcation of active and silent chromatin domains in human HOX loci by noncoding RNAs. *Cell* **129**, 1311–1323 (2007).
- Willingham, A.T. *et al.* A strategy for probing the function of noncoding RNAs finds a repressor of NFAT. *Science* **309**, 1570–1573 (2005).
- Zhao, J., Sun, B.K., Erwin, J.A., Song, J.J. & Lee, J.T. Polycomb proteins targeted by a short repeat RNA to the mouse X chromosome. *Science* **322**, 750–756 (2008).
- Katayama, S. *et al.* Antisense transcription in the mammalian transcriptome. *Science* **309**, 1564–1566 (2005).
- Wu, J. Q. *et al.* Dynamic transcriptomes during neural differentiation of human embryonic stem cells revealed by short, long, and paired-end sequencing. *Proc. Natl. Acad. Sci. USA* **107**, 5254–5259 (2010).
- Ramsköld, D., Wang, E.T., Burge, C.B. & Sandberg, R. An abundance of ubiquitously expressed genes revealed by tissue transcriptome sequence data. *PLOS Comput. Biol.* **5**, e1000598 (2009).

ONLINE METHODS

Cell culture. Mouse ESCs (V6.5) were cultured with irradiated mouse embryonic fibroblasts (GlobalStem; GSC-6002C) on 0.2% gelatin-coated plates in a culture medium consisting of Knockout DMEM (Invitrogen; 10829018) containing 10% FBS (GlobalStem; GSM-6002), 1% penicillin-streptomycin (Invitrogen 15140-163), 1% non-essential amino acids (Invitrogen 11140-076), 1% L-glutamine, 4 μ l β -mercaptoethanol and 0.01% leukemia inhibitory factor (LIF; Millipore; ESG1106). ESCs were passaged once on gelatin without mouse embryonic fibroblasts before RNA extraction. V6.5 ESCs were differentiated into NPCs through embryoid body formation for 4 d and selection in ITSFn medium²⁷ for 5–7 and maintained in fibroblast growth factor-2 (FGF-2) and epidermal growth factor-2 (EGF-2) (R&D Systems) as described²⁷. The cells uniformly express Nestin and Sox2 and can differentiate into neurons, astrocytes and oligodendrocytes. Mouse lung fibroblasts (ATCC) were grown in DMEM with 10% FBS and penicillin/streptomycin at 37 °C, 5% CO₂.

RNA extraction and library preparation. RNA was extracted using the protocol outlined in the RNeasy kit (Qiagen). Extracts were treated with DNase (Ambion 2238). Polyadenylated RNAs were selected using Ambion's MicroPoly(A)Purist kit (AM1919M) and RNA integrity confirmed using Bioanalyzer (Agilent). We used a cDNA preparation procedure that combines fragmentation of mRNA to a peak size of ~750 nucleotides by heating⁶ followed by random-primed reverse transcription⁸. We previously found^{9,28} that this protocol provides relatively uniform coverage of the whole transcript, thus assisting in *ab initio* reconstruction. Specifically, a 'regular' RNA sequencing library (non-strand specific) was created as previously described²⁸, with the following modifications. Poly(A)⁺ RNA (250 ng) was fragmented by heating at 98 °C for 33 min in 0.2 mM sodium citrate, pH 6.4 (Ambion). Fragmented RNA was mixed with 3 μ g random hexamers (Invitrogen), incubated at 70 °C for 10 min, and placed on ice briefly before starting cDNA synthesis. First-strand cDNA synthesis was performed using Superscript III (Invitrogen) for 1 h at 55 °C, and second-strand using *E. coli* DNA polymerase and *E. coli* DNA ligase at 16 °C for 2 h. cDNA was eluted using the Qiagen MiniElute kit with 30 μ l of the manufacturer's EB buffer. DNA ends were repaired using dNTPs and T4 polymerase (NEB), followed by purification using the MiniElute kit. Adenine was added to the 3' end of the DNA fragments using dATP and Klenow exonuclease (NEB; M0212S) to allow adaptor ligation, and fragments were purified using MiniElute. Adaptors were ligated and incubated for 15 min at room temperature (25 °C). Phenol/chloroform/isoamyl alcohol (Invitrogen 15593-031) extraction followed to remove the DNA ligase. The pellet was then resuspended in 10 μ l EB buffer. The sample was run on a 3% agarose gel (Nusieve 3:1 agarose, Lonza) and a 160–380 base pair fragment was cut out and extracted. PCR was performed with Phusion High-Fidelity DNA Polymerase with the manufacturer's GC buffer (New England Biolabs) and 2 M betaine (Sigma). PCR conditions were 30 s at 98 °C; 16 cycles of 10 s at 98 °C, 30 s at 65 °C, 30 s at 72 °C; 5 min at 72 °C; forever at 4 °C. Products were run on a polyacrylamide gel for 60 min at 120 V. The PCR products were cleaned up with Agencourt AMPure XP magnetic beads (A63880) to completely remove primers and product was submitted for Illumina sequencing.

The strand-specific library was created from 100 ng of poly(A)⁺ RNA using the previously published RNA ligation method²⁹ with modifications from the manufacturer (Illumina; data not shown). The insert size was 110 to 170 bp.

RNA-Seq library sequencing. All libraries were sequenced using the Illumina Genome Analyzer (GAII). We sequenced three lanes for ESC, corresponding to 152 million reads; two lanes for MLE, corresponding to 161 million reads; and two lanes for NPC, corresponding to 180 million reads.

Alignments of reads to the genome. All reads were aligned to the mouse reference genome (NCBI 37, MM9) using the TopHat aligner¹³. Briefly, TopHat uses a two-step mapping process, first using Bowtie³⁰ to align all reads that map directly to the genome (with no gaps), and then mapping all reads that were not aligned in the first step using gapped alignment. TopHat uses canonical and non-canonical splice sites to determine possible locations for gaps in the alignment.

Generation of connectivity graph. Given a set of reads aligned to the genome, we first identified all spliced reads as those whose alignment to the reference genome contained a gap. These reads and the reference genome were used to construct connectivity graphs. Each connectivity graph contains all bases from a single chromosome. The nodes in the graph are bases and the edges connect each base to the next base in the genome as well as to all bases to which it is connected through a spliced read (Fig. 1). In the analysis presented, we identified as an edge any two bases in the chromosome that were connected by two or more spliced reads. The connectivity graph thus represents the contiguity that exists in the RNA but that is interrupted by intron sequences in the reference genome.

Identification of splice site motifs and directionality. We restricted our analysis to spliced reads that mapped connecting donor/acceptor splice sites, either canonical (GT/AG) or noncanonical (GC/AG and AT/AC). We oriented each mapped spliced read using the orientation of the donor/acceptor sites it connected.

Construction of transcript graphs. The spliced edges in the connectivity graph reflect bases that were connected in the original RNA but are not contiguous in the genome. To construct a transcript graph, we use a statistical segmentation strategy to traverse the graph topology directly and determine 'paths' through the connectivity graph that represent a contiguous path of significant enrichment over the background distribution (see below). In this segmentation process, we scan variably sized windows across the graph and assign significance to each window. We then merge significant paths into a 'transcript graph'. Specifically, for a window of fixed size, we slide the window across each base in the connectivity graph (after augmenting it with the unspliced reads). If a window contains only contiguous unspliced reads, then it represents an unspliced part of the transcript. However, if the window hits an edge in the connectivity graph connecting two separate parts of the genome (based on two or more spliced reads), then the path follows this edge to a noncontiguous part of the genome, denoting a splicing event. Similarly, when alternative splice isoforms are present, if a base connects to multiple possible places, then we compute all windows across these alternative paths. Using a simple recursive procedure, we can compute all paths of a fixed size across the graph.

Identification of significant segments. To assess the significance of each path, we first define a background distribution. We estimate a genomic defined background distribution by permuting the read alignments in the genome and counting the number of reads that overlap each region and the frequency by which they each occur. Specifically, if we are interested in computing the probability of observing alignment *a* (of length *r*) at position *i* (out of a total genome size of *L*) we can permute the alignments and ask how often read *a* overlaps position *i*. Under this uniform permutation model, the probability that read *a* overlaps position *i* is simply *r/L*. Extending this reasoning, we can compute the probability of observing *k* reads (of average length *r*) at position *i* as the binomial probability. Given the many reads and the large genome size, the binomial formula can be well approximated by a Poisson distribution where $\lambda = np$ (that is, the number of reads times the number of possible positions).

Given a distribution for the real number of counts over each position, we scan the genome for regions that deviate from the expected background distribution. First, consider a fixed window size *w*. We slide this window across each position (allowing for overlapping windows), and compute the probability of each observed window based on a Poisson distribution with $\lambda = wnp$. Since we are sliding this window across a genome of size *L*, we correct our nominal significance for multiple testing by computing the maximum value observed for a window size (*w*) across a number of permutations of the data. This distribution controls the family wise error rate, defined as the probability of observing at least one such value in the null distribution³¹. Notably, we can estimate this maximum permutation distribution well by a distribution known as the scan statistic distribution³², which depends on the size of the genome that we scan, the window size used and our estimate of the Poisson λ parameter. This method provides us with a general strategy to determine a multiple testing-corrected *P*-value for a

specified region of the genome in any given sample. We use this method to compute a corrected significance cutoff for any given region.

Finally, to identify significant intervals, we scan the genome using variably sized windows, computing significance values for each and filtering by a 0.05 significance threshold. For each window size, we merge the significant regions that pass this cutoff into consecutive intervals. We trim the ends of the intervals as needed, because we are computing significant windows (rather than regions) and it is possible that an interval need not be fully contained within a significant region. Trimming is performed by computing a normalized read count for each base in the interval compared to the average number of reads in the genome. We then trim the interval to the maximum contiguous subsequence of this value. We test this trimmed interval using the scan procedure and retain it only if it passes our defined significance level.

We work with a range of different window sizes in order to detect paths (intervals) with variable support. Small windows have the power to identify short regions of strong enrichment (for example, a short exon that is highly expressed), whereas long windows capture long contiguous regions with often lower and more 'diffuse' enrichment (for example, a longer, lower-expression transcript, whose 'moderate evidence' aggregates along its entire length).

Estimation of library insert size. We estimated the insert size distribution by taking all reconstructed transcripts for which we only reconstructed a single isoform and computing the distribution of distances between the paired-end reads that aligned to them.

Weighting of isoforms using paired end edges. Using the size constraints imposed by the length of the paired ends, we assigned weights to each path in the transcript graph. We classified all paired ends overlapping a given path and assigned them to all possible paths that they overlapped. We then assigned a probability to each paired end of the likelihood that it was observed from this transcript given the inferred insert size for the pair in that path. We used an empirically determined distribution of insert sizes, estimated from single isoform graphs. We then scaled each value by the average insert size. We refer to this scaled value as our insert distribution. For each paired end in a path, we computed I , the inferred insert size (the distance between nodes following along the full path) minus the average insert size. We then determined the probability of I as the area in our insert distribution between $-I$ and I . This value is the probability of obtaining the observed paired-end insert distance given this distribution of paired-end reads. We use this probability to compute a weighted score for each path by summing all paired ends that fall within the path weighted by the probability of observing insert size they span in the path. Paired ends that support multiple isoforms equally will count equally for all, but paired ends with biases toward some isoforms and against others will provide weighted evidence for each isoform. We assign this weight to each isoform path. This score is normalized by the number of paired ends overlapping the path. We filter out paths with little support (normalized score < 0.1) from paired reads.

Determination of expression levels from RNA-Seq data. Expression levels are computed as previously described⁸. Briefly, the expression of a transcript is computed in reads per kilobase of exonic sequence per million aligned reads (RPKM) defined as $RPKM_{\text{transcript}} = 10^9 r/Rt$, where r is the number of reads mapped to the exonic region of the transcript, t is the total exonic length of the transcript and R is the total number of reads mapped in the experiment.

Array expression profiling in ESC cells. Microarray hybridization data was obtained from our previous studies including ESC¹⁶, NPC¹⁶ and MLF⁴.

Comparisons to known annotation. The reconstructed transcripts were compared to the RefSeq genome annotation¹⁵ (NCBI release 39). To determine whether a known annotation of a protein coding gene from RefSeq was fully reconstructed, we first compared the 5' and 3' ends of the reconstructed versus the annotated transcript. If these overlapped, we further verified that all exons in the annotated transcript matched those in the reconstructed version. To score the portion of an annotated transcript

covered by our reconstructions, we found the reconstructed transcript whose exons covered the largest fraction of the annotated transcript and reported the portion of the annotation that it covered.

ChIP-seq profiles in ESC cells and determination of K4 and K36 regions. To determine regions enriched in chromatin marks from chromatin immunoprecipitation followed by sequencing (ChIP-seq) data, we applied a previously described method⁴ to ESC, MLF and NPC data^{4,16}.

Determination of external and internal 5' start sites. We identified alternative 5' start sites by comparing the 5' exon of our reconstructed transcripts to the location of the 5' exon of the annotated gene overlapping it. If the reconstructed 5' start site resided upstream of the annotated 5', we termed it an external start site. For novel 5' ends that were downstream of the annotated 5' end (internal), we required a few more criteria to avoid reconstruction biases due to low coverage. First, we required that the novel internal 5' end did not overlap any of the known exons within the known gene. Second, we required that the reconstructed gene contain a completed 3' end. To determine the presence of H3K4me3 modifications overlapping the promoter regions defined by these novel start sites, we computed regions of enriched H3K4me3 genome-wide (as previously described) and intersected the location of the novel 5' exon (both internal and external) with the location of an H3K4me3 peak.

Determination of premature/extended 3' end. To determine novel 3' ends, we compared the locations of the 3' exon of our reconstructed 3' ends and those of annotated genes. If the reconstruction extended past the annotated 3' end, we classified it as an extended 3' end. If the reconstruction ended before the annotated 3' end, we required that it not overlap any known exon and have a fully reconstructed 5' start site.

Determination of sequence conservation levels. We used the SiPhy¹⁹ algorithm and software package (http://www.broadinstitute.org/genome_bio/siphy/) to estimate ω , the deviation ('contraction' or 'extension') of the branch length compared to the neutral tree based on the total number of substitutions estimated from the alignment of the region of interest across 20 placental mammals (build MM9, <http://hgdownload.cse.ucsc.edu/goldenPath/mm9/multiz30way/>). For global (whole transcript) conservation, we estimated ω for each protein coding, lincRNA and antisense transcript exon and compared it to similarly sized regions within introns. To identify local regions of conservation within a transcript, we computed ω for all 12-mers within the transcript sequence and assigned a P -value for each 12-mer based on the chi-squared distribution, as previously described¹⁹. We then took all 12-mers showing significance at $P < 0.05$, collapsed overlapping 12-mers and identified constrained regions within the transcript (**Supplementary Fig. 6**, for example).

ORF determination. We estimated maximal supported open reading frames (ORFs) for each transcript built by scanning for start codons and computing the length (in nucleotides) until the first stop codon was reached.

Codon substitution frequency scores. To further estimate the coding potential of novel transcripts, we evaluated whether evolutionary sequence substitutions were consistent with the preservation of the reading frame of any detected peptide. Briefly, if a transcript encodes a protein, we expect a reduction in frame-shifting indels, nonsynonymous changes and, in general, any substitution that affects the encoded protein. To assess this, we used the CSF method as previously described^{17,18}.

RT-PCR validations. Primers were obtained for a randomly selected set of predicted lincRNA, protein coding genes, antisense transcripts and intron primers (**Supplementary Table 2**); all began with an M13 primer sequence. RNA from ESC cells was extracted using Qiagen's RNeasy kit (74106). A one-step RT-PCR reaction was run using Invitrogen's One-Step RT-PCR kit (12574-018) following the manufacturer's instructions, with the following PCR protocol: 55 °C for 30 min, 94 °C for 2 min; 40 cycles of 94 °C for 15 s, 64 °C for 30 s, 68 °C for 1 min; 68 °C for 5 min; 4 °C forever. Samples were separated on a

3% agarose gel, and all bands were cut out and gel extracted using the QIAquick Gel Extraction kit (Qiagen 28706). DNA (30 ng) was mixed with 3.2 pmol M13 forward or M13 reverse primer and sequenced in both directions.

Data availability. The sequencing data from this study are available at the NCBI Gene Expression Omnibus (GEO) under accession code GSE20851 and as **Supplementary Data**. The Scripture method is implemented as a standalone Java application and is available as **Supplementary Software** and at <http://www.broadinstitute.org/software/Scripture/>, along with all assembled transcripts in both GFF and BED file formats. All transcript graphs are also available in the dot graph language.

27. Conti, L. *et al.* Niche-independent symmetrical self-renewal of a mammalian tissue stem cell. *PLoS Biol.* **3**, e283 (2005).
28. Berger, M. F. *et al.* Integrative analysis of the melanoma transcriptome. *Genome Res.* **20**, 413–427 (2010).
29. Lister, R. *et al.* Highly integrated single-base resolution maps of the epigenome in *Arabidopsis*. *Cell* **133**, 523–536 (2008).
30. Langmead, B., Trapnell, C., Pop, M. & Salzberg, S.L. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol.* **10**, R25 (2009).
31. Ewens, W.J. & Grant, G.R. *Statistical Methods in Bioinformatics: An Introduction* 2nd edn. (Springer, 2005).
32. Glaz, J., Naus, J.I. & Wallenstein, S. *Scan Statistics* (Springer, 2001).

Corrigendum: Safety signal dampens reception for mipomersen antisense

Jim Kling

Nat. Biotechnol. 28, 295–297 (2010); published online 8 April 2010; corrected after print 9 July 2010

In the version of this article initially published, some of the oligos in Table 1 are described as phosphorothioate modified. In fact, all antisense oligonucleotides are phosphorothioate-modified oligos. In addition, Lucanix, which is not an antisense oligo, has been removed from the table. The error has been corrected in the HTML and PDF versions of the article.

Corrigendum: *Ab initio* reconstruction of cell type-specific transcriptomes in mouse reveals the conserved multi-exonic structure of lincRNAs

Mitchell Guttman, Manuel Garber, Joshua Z Levin, Julie Donaghey, James Robinson, Xian Adiconis, Lin Fan, Magdalena J Koziol, Andreas Gnirke, Chad Nusbaum, John L Rinn, Eric S Lander & Aviv Regev

Nat. Biotechnol. 28, 503–510 (2010); published online 02 May 2010; corrected after print 9 July 2010

In the version of this article initially published, the fourth sentence in the Online Methods section “RNA extraction and library preparation,” that read in part “procedure that combines a random priming step with a shearing step^{8,9,28} and results in fragments of ~700 bp in size,” should have read, “procedure that combines fragmentation of mRNA to a peak size of ~750 nucleotides by heating⁶ followed by random-primed reverse transcription⁸.” The error has been corrected in the HTML and PDF versions of the article.

Erratum: US biodefense contracts continue to lure biotechs

Catherine Shaffer

Nat. Biotechnol. 28, 187–188 (2010); published online 8 March 2010; corrected after print DD Month 9 July 2010

In the version of this article initially published, in Table 1, the Emergent BioSolutions’ anthrax countermeasures in development listed AV-7909 as being in phase 2 under a \$447.6 million BARDA contract; AV-7909 is in phase 1 and the BARDA contract is for \$29.7 million. AIGIV is in phase 1/3, not phase 1/2. Finally, a third product was omitted; anthrax monoclonal is in preclinical testing under a \$24 million BARDA contract. The \$447.6 million BARDA contract was for procurement and product enhancements on BioThrax. Also, on p.188, column 2, line 7, the vaccine requires five injections, not six as originally stated. The error has been corrected in the HTML and PDF versions of the article.

Erratum: Single base-resolution methylome of the silkworm reveals a sparse epigenomic map

Hui Xiang, Jingde Zhu, Quan Chen, Fangyin Dai, Xin Li, Muwang Li, Hongyu Zhang, Guojie Zhang, Dong Li, Yang Dong, Li Zhao, Ying Lin, Daojun Cheng, Jian Yu, Jinfeng Sun, Xiaoyu Zhou, Kelong Ma, Yinghua He, Yangxing Zhao, Shicheng Guo, Mingzhi Ye, Guangwu Guo, Yingrui Li, Ruiqiang Li, Xiuqing Zhang, Lijia Ma, Karsten Kristiansen, Qihong Guo, Jianhao Jiang, Stephan Beck, Qingyou Xia, Wen Wang & Jun Wang

Nat. Biotechnol. 28, 516–520 (2010); published online 02 May 2010; corrected after print 9 July 2010

In the version of this article initially published, references 4 and 7 were inadvertently interchanged. The error has been corrected in the HTML and PDF versions of the article.

Erratum: Up for grabs

Michael Eisenstein

Nat. Biotechnol. 28, 544–546 (2010); published online 7 June 2010; corrected after print 9 July 2010

In the version of the article originally published, it was stated that the Cohen-Boyer patents generated hundreds of billions of dollars in licensing revenue. It should have read hundreds of millions of dollars. The error has been corrected in the HTML and PDF versions of the article.