

Supporting Information

Shiroguchi et al. 10.1073/pnas.1118018109

SI Materials and Methods

Evaluation of Digital Counting for Synthesized DNAs Using Random Barcodes by Deep Sequencing. To evaluate the capability of random barcodes to digitally count DNA molecules, we attempted digital counting of synthesized DNA on the Illumina platform using single read sequencing.

Design of single-stranded DNA with random barcodes. We designed and purchased (Integrated DNA Technologies) nine single-stranded DNA molecules, each with adapter sequences at both ends that are compatible with the Illumina platform (Dataset S3). Ten, 16, or 20 random bases were positioned near the beginning of sequencing reads. We measured the concentration of each DNA by absorbance at 260 nm using the extinction coefficient provided by Integrated DNA Technologies.

Sequencing sample preparation and sequencing. We mixed all nine single-stranded DNA templates for a final concentration of 0.17 aM ssDNA-1, 0.67 aM ssDNA-2, 3.3 aM ssDNA-3, 33 aM ssDNA-4, 330 aM ssDNA-5, 660 aM ssDNA-6, 660 aM ssDNA-7, 660 aM ssDNA-8, and 660 aM ssDNA-9 in HF buffer [New England Biolabs(NEB)] with 0.4 mM dNTP, 1 μ M of each amplification primer (Integrated DNA Technologies) (Dataset S3), and 1 U Phusion DNA polymerase (NEB) in a final volume of 50 μ L. These DNA molecules were amplified by PCR (1 cycle of 98 °C for 30 s, 20 cycles of both 98 °C for 10 s and 60 °C for 30 s, and 1 cycle of 72 °C for 10 min). After column purification (Zymo Research), the concentration was measured by quantitative PCR (qPCR) (Fast SYBR Master Mix; Applied Biosystems) using conventional Illumina qPCR primers against a reference PCR product amplified from ssDNA-2 using the amplification primers mentioned above, the concentration of which was measured by NanoDrop (LMS). The samples were sequenced on an Illumina Genome Analyzer II in one lane with 36-base single reads. We obtained 26,566,253 reads resulting from 26,566,253 sequencing clusters, of which 25,738,380 mapped to the designed template sequences (Dataset S3).

Data analysis. We sorted the raw reads by the known sequence for each template (Dataset S3) and allowed up to two single-base mismatches and one insertion or deletion of a single base before the barcode region. We counted the number of molecules for each template by counting the number of uniquely observed barcodes.

Results. For all templates, the results were completely inconsistent with the estimated number of molecules added in the PCR mixture. For example, the measured number of molecules of ssDNA-8 [shown in Fig. 2B (red)] was 15.5-fold higher than the estimated number of the DNA molecules added to the original PCR mixture; we identified 155,161 unique barcodes compared with an estimated input of 10,000 molecules. The histogram of the distribution of the number of reads for the random barcodes [for example, Fig. 2B (red)] showed that all ssDNA templates have a remarkably high peak at a very low number of reads. When we assumed that two molecules with up to two mismatches in the barcode region were originally from same DNA molecule, the peak significantly decreased [for example, the peak in Fig. 2B (red) decreased about 30-fold]. This results in 23,890 remaining barcodes and a roughly 2.4-fold higher output than the estimated input (two other templates, ssDNA-1 and ssDNA-3 also result in a 5.8-fold and a 27-fold higher estimation, respectively, after the same process). Such error correction may reduce the number of potential false-positives, but it is still inaccurate and results in uncertainty; the maximum number of allowable nucleotide differences used for correction greatly affects the final results. This

finding suggests that newly generated sequences at the barcode position, from sequencing or PCR amplification errors that did not exist in the original PCR mixture, may have become new artifactual unique barcodes, which resulted in a significant number of false-positives. The data for the red histogram shown in Fig. 2B was generated by down-sampling reads such that the maximum number of reads for any given molecule, and the number of total molecules depicted in the figure corresponds to the parallel optimized barcode experiment described in the main text.

Efficacy of Digital Counting Strategy. The correlation plot in Fig. 2A demonstrates the efficacy of our digital counting strategy for the spike-in sequences. In addition, the histogram in Fig. 2B (green) for the optimized barcode also highlights an important difference between the conventional approach of counting the number of reads corresponding to a target sequence and digital counting based on unique barcode labeling. We know that the abundance of the spike-in molecule was accurately quantified by our digital counting scheme, and so the vast majority of unique barcodes observed for the spike-in molecule correspond to individual molecules of the spike-in in our original sample. Hence, the correct number of digital counts for most of those molecules is one. Because Fig. 2B is a histogram of the number of reads observed for each of the spike-in barcodes, it is equivalent to the conventional counting distribution for individual copies of the spike-in in the original sample. Therefore, according to the conventional counting strategy, the abundance of each individual copy of the spike-in in our original sample varies over three orders-of-magnitude (Fig. 2B). This discrepancy is most likely because of intrinsic noise in sample preparation (arising from amplification and sequencing), which in this case is catastrophic for the conventional method but does not affect the accuracy of our digital counting strategy.

Tenfold Down-Sampling of Spike-in Reads. We randomly down-sampled the spike-in reads in the replicate experiment by a factor of 10. For each read of each of the five spike-ins, there was a 10% chance that it was kept and counted, whereas the other 90% was discarded. Fig. S3 shows that there is little dropout between these two conditions (the data show the average dropout rate for the five spike-ins to be 1.6%) and that for the spike-in with the highest number of molecules, the change in the single-molecule coverage histogram is minimal.

Simulation Comparison of Performance of Optimized Barcoding Scheme vs. Random Barcoding Scheme. We assumed a theoretical library of 6,000 identical template molecules and barcoded each molecule *in silico* under two separate conditions: (i) with 6,000 optimized barcodes and (ii) with 6,000 random barcodes. The optimized barcodes mirror the technique described in the main text (Paired-end reads, 20 nucleotides), and the random barcodes mirror the technique described above (single-end reads, 16 nucleotides). We note that previous uses of the random barcode procedure have only involved single-end reads (1, 2). Each of the barcoded templates in both conditions were “amplified” without sequence-dependent bias and noise for simplicity to a copy number of 400. Then, using the averaged sequencing error profile for the Illumina platform characterized by Nakamura et al. (3), we simulated sequencing errors of the first 53 bases. Mapping dropout was also included in the simulation (we allowed at most two sequencing errors per transcript for mapping). If we allow one mismatch in both barcode regions simulating

actual experimental conditions for the optimized barcode simulation, we observe that there are no artificially created barcodes (i.e., false-positives) (Fig. S1). However, for the random barcode simulation, there is a substantial peak for low copy barcodes, which is because of sequencing errors generating new barcodes not originally present in the original sample, as we see that the total number of original molecules estimated is almost two orders-of-magnitude greater than the actual value (Fig. S1). From the *in silico* sequencing results in the random case, we cannot determine which barcodes were actually attached to the template molecules and which barcodes were generated by sequencing error because all possible random sequences can be potentially attached to the template molecule. As a result, sequencing error generates many false-positives. In the case of optimized barcodes, sequencing errors occurred as well. However, we can identify the sequences that were not in the set of used barcodes, and so sequencing errors in the barcode region can be identified and removed, which does not result in false-positives. If we assume resulting random barcode sequences that differ by one nucleotide are identical, the estimated output is 58,556 and still contains many apparent false-positives (the two nucleotide case results in 9,597 estimated output molecules). Although error correction reduces the number of false-positives, it is clearly not as accurate as our optimized barcode pair method and still results in uncertainty; the maximum number of allowable nucleotide differences used for correction greatly affects the final results.

Simulations Demonstrating the Utility of Digital RNA Sequencing for Differential Expression Analysis. We performed computer simulations to demonstrate that digital counting using optimized barcodes can substantially outperform conventional counting in differential expression analysis experiments, where expression profiles are compared under various conditions. We simulated differential expression analysis for three cases: a system with the same copy number distribution as the *Escherichia coli* fragment library described in Fig. 3, a system with approximately the same copy number distribution as the *E. coli* transcription unit library described in Fig. S2, and a system with the copy number distribution that was experimentally measured for microRNA in human stomach tissue (4). In each case, two gene expression profiles, each having the same copy number distribution, were generated randomly and compared with each other. For the *E. coli* fragment, transcription unit, and human microRNA libraries, each molecule in the simulation (83,863 and 352,088 total molecules for the *E. coli* fragments and human microRNA libraries, respectively, and 83,500 total molecules for ~4,000 transcription units for the transcription unit library) was assigned a barcode from a large pool (21,025 barcode pairs for the *E. coli* fragment and transcription unit libraries corresponding to 145 barcode sequences and 189,225 barcode pairs for the microRNA library, corresponding to 435 barcode sequences). However, the barcodes were not assigned randomly. Instead, they were assigned based on the barcode sampling bias distribution, which we experimentally measured for the *E. coli* transcriptome (Fig. 3B). After assigning barcodes, we simulated PCR amplification of each molecule. We know from our conventional counting results for the *E. coli* transcriptome that a significant amount of noise is incurred during amplification because not every PCR cycle for every molecule is 100% efficient (Fig. 3C). Imperfect efficiency and exponential amplification conspire to produce the majority of the variability in conventional counting of low copy transcripts observed in Fig. 3C. Hence, in our simulation, we assumed a PCR cycle efficiency (70%) that reproduces the experimentally measured amplification variance for single molecules observed in Fig. 3C. We note that the PCR simulation does not account for sequence-dependent or cycle-dependent amplification bias. Just as in our *E. coli* transcriptome RNA sequencing (RNA-Seq) experiment, we applied 18 cycles of PCR to every molecule in

each library. At the end of the simulation, we counted the number of unique barcodes assigned to each sequence (digital counts) and the number of PCR amplicons for each sequence (conventional counts). After running this simulation for both libraries in the two cases, we obtained a simulated fold-change for each gene and compared it to the original input fold-change for each gene. Fig. S4 shows the ratio of the simulated to input fold-changes for each fragment in the *E. coli* fragment library, each transcription unit in the *E. coli* transcription unit library, and each microRNA in the microRNA library for both digital and conventional counting as a function of the lower of the two input copy numbers for each gene. Although the digital counting result is imperfect because the barcode sampling is slightly biased and the set of barcode sequences is finite, digital counting vastly outperforms conventional counting in terms of accurately quantifying fold-changes between two gene expression profiles in both cases. We note that the copy-number distributions in the two cases differ significantly. Copy numbers fall between 1 and 20 for the *E. coli* fragment library and between 1 and ~20,000 for the microRNA library, but digital counting is advantageous in both cases. In addition, the performance of conventional counting vs. digital counting is significantly improved for the *E. coli* transcription unit simulation compared with the fragment library because the PCR noise for each transcription unit is averaged over several fragments. This finding is consistent with the improved correlation of conventional vs. digital counting depicted in Fig. S2 compared with the fragment correlation shown in Fig. 3C. The overall performance discrepancy between conventional and digital counting and the increase in conventional counting accuracy as a function of copy number are expected from the amplification noise that is evident in Fig. 3C. We conclude that our digital counting scheme will be advantageous for differential expression analysis, particularly for applications involving small fragments as in the microRNA case simulated here.

In-Depth Design of Barcode. Each of the 2,358 barcodes was analyzed for sequence characteristics that would contribute to either amplification or sequencing errors.

Initial filtering. All barcodes with less than 40% or greater than 60% GC-content or containing homopolymers greater than length four were deleted. All barcode sequences were compared with the PE 1.0 and 2.0 Illumina PCR primer sequences and discarded if there were more than 10 total base matches or more than five consecutive base matches in any possible alignment with either primer sequence (sense and antisense). Each barcode was compared with the final four, five, and six bases closest to the 3' end of the PCR primer sequence (sense and antisense), respectively. The final six bases for both the PE 1.0 and PE 2.0 are identical. If any of these regions contained more than three consecutive base matches with a given barcode in any possible alignment, that barcode was discarded. All barcode sequences were compared with all other barcode sequences (including cases of offset sequences), and a barcode was discarded if there were more than 15 total base matches or 10 consecutive base matches to any other barcode sequence (sense and antisense) in any possible alignment. The total number of hydrogen bonds in the longest consecutive matching region of each barcode to any other barcode sequence was calculated, and barcodes with greater than 26 total hydrogen bonds in that region were discarded. Each barcode was aligned with the entire *E. coli* genome (sense and antisense). If at any position in the genome a barcode contained more than 16 total matching bases, 12 consecutive matching bases, or 32 hydrogen bonds present in any given consecutively matching region, that barcode was discarded. Finally, all possible indels were generated for each barcode and the resulting sequence was compared with each original barcode sequence. If the resulting indel sequence could incur fewer than five point mutations and result in the exact sequence of any

barcode, the barcode sequence that generated the given indel was deleted.

Score filtering. In addition to these universal thresholds, a more in-depth analysis of barcode-barcode and barcode-*E. coli* genome hybridization was performed, particularly regarding barcode hybridization melting temperatures with respect to PCR amplification. Although for a given barcode sequence, when comparing complementarity to a large set of reference sequences through sequence alignment, there will be an alignment condition that results in a region in the barcode sequence where the absolute maximum number of consecutive base matches is achieved (as described above), there are other possible conditions when a barcode contains a region where the number of consecutive bases in the region of maximum consecutive base matches does not reach the absolute maximum value, as described above. We shall define this value for any given region as the score. When comparing a barcode to the sense and antisense sequences of all other barcodes, we determined all alignment conditions in the cases where the score of the region that contained the highest number of consecutive base matches (first score) were 10, 9, 8, and 7 bases, respectively. For each of these four first scores, the condition where the maximum score of the region where the second-highest number of consecutive base matches (second score) occurred was determined. For example, we shall denote a condition where the first score is 10 and the second score is 3 as “10-3” and define this condition as the duplex. If the sum of the first score and the second score compared with all other barcode sequences for any barcode was greater than 12, that barcode was discarded. If the distance between these two regions (maximum and second-most consecutive matches) was one base, and the sum of the first score and the second score was greater than 11, the barcode was discarded. The maximum value of the number of consecutive base matches for a region under all alignment conditions that contained the third-highest number consecutive base matches (third score) was also determined for all barcodes. Given the maximum third score, the respective maximum first score was determined; the maximum second score given both of these conditions was also determined. We shall define this condition as the triplex and denote it, for example, as 7-4-3. We manually deleted the barcodes with the following triplexes: 7-3-3, 6-5-4, 6-5-3, 6-4-4, 5-5-4, and 5-4-4. We also manually deleted barcodes with a triplex of 6-4-3, where both the distances between adjacent regions corresponding to the scores was one base.

The same analysis was done for all barcodes aligned against the entire *E. coli* genome (sense and antisense). Barcodes with a first score and second score sum of greater than 15 were discarded. Barcodes where the first score region and the second score region were separated by one nucleotide and had a first score and second score sum of more than 14 were discarded as well. Barcodes with the following triplexes were deleted: 8-4-4, 7-5-4, 6-6-4, 6-5-5, and 5-5-5. After filtering, a total of 150 barcodes remained.

In-Depth Design and Preparation of Adapter. *Adapter design.* To avoid sequencing errors resulting from cluster overlap (i.e., low sequence complexity) and to reduce potential ligation bias, an additional two- to five-bp extension—CT, ACT, GACT, or TGACT—was added to the 3'-end of each barcode. These sequences mimic the T-overhang in the conventional Illumina paired-end adapter and conserve the sequence of the last two bases. For each of the 150 final barcodes, we attached these four different adapter extensions to the 3' end of the barcode. We then obtained the same values as used in the initial filtering step (see above) for each of the four adapter candidates for all 150 barcodes. We determined the following four parameters of analysis: PCR primer matching (PC), 3' end of PCR primer matching (TP), barcode-barcode matching (BB), and barcode

E. coli genome matching (EC), and calculated the complementarity score of each category for all barcode-adapter candidates as follows:

$$PC = \{\text{Sum of [maximum total base matches to the PE 1.0 and PE 2.0 PCR primers (sense and antisense for a total of four terms)]} + 2 \cdot \{\text{Sum of [maximum consecutive base matches to the PE 1.0 and PE 2.0 PCR primers (sense and antisense for a total of four terms)]}^2\}$$

$$TP = \{\text{Sum of [maximum total base matches to the final four bases of the PCR primer sequence (sense and antisense for a total of two terms)]}^2 + 1.5 \cdot \{\text{Sum of [maximum total base matches to the final five bases of the PCR primer sequence (sense and antisense for a total of two terms)]}^2 + 2 \cdot \{\text{Sum of [maximum total base matches to the final six bases of the PCR primer sequence (sense and antisense for a total of two terms)]}^2\}$$

$$BB = \{\text{Sum of [maximum total base matches to all other barcode candidates (sense and antisense for a total of two terms)]}^2 + 2 \cdot \{\text{Sum of [maximum consecutive base matches to all other barcode candidates (sense and antisense for a total of two terms)]}^2\}$$

$$EC = \text{Maximum total base matches to entire } E. coli \text{ genome (sense only)} + 2 \cdot \{\text{maximum consecutive base matches to the entire } E. coli \text{ genome (sense only)}\}^2$$

The total complementarity score (TC) for each barcode candidate was calculated as follows:

$$TC = 3 \cdot PC + 15 \cdot TP + BB + EC$$

The TC value gives us a metric to determine the expected efficacy of each barcode candidate during PCR amplification. A low TC represents a lower chance of amplification errors caused by unwanted hybridization between barcodes and adapters, primers, or the sample. For each barcode, we selected the barcode-adapter candidate that had either the lowest or second-lowest TC among the four: this resulted in 150 final barcode-adapter sequences, of which 145 were randomly chosen and used.

Thirty-seven CT extensions and 36 of each of the other three extensions were used. Adapters were then designed in the same Y-shaped construct as the conventional Illumina paired-end adapter with a 22- to 25-bp extension that contained the barcode and a T-overhang (Fig. 1B). Both strands (A and B) of the adapter were ordered from Integrated DNA Technologies ([Dataset S1](#)).

Adapter generation. The 5'-end of strand B was phosphorylated in T4 DNA Ligase Reaction Buffer (NEB) containing 40- μ M strand B and 20 U T4 polynucleotide kinase (NEB) at 37 °C for 60 min in a 20 μ L volume, followed by a 25-min incubation at 70 °C, and a 5-min incubation at 90 °C for enzyme inactivation. The phosphorylated strand B was annealed to each respective strand A in NEB Buffer 2 (NEB). Each solution contained 20- μ M strand A and 20- μ M strand B in a total volume of 20 μ L. The solutions were first raised to 90 °C and cooled to 25 °C at a rate of 5 °C per minute (annealing temperature condition). Finally, equal volumes of all 145 annealed adapters were mixed.

Design and preparation of spike-in and “normalization” DNA. Fifteen-thousand random 30-bp sequences were generated such that even if a sequence accumulated 15 mutations, it would still be identifiable and distinguishable from all other generated sequences. Spike-in and normalization candidates with a maximum homopolymer length of greater than 3, or a GC-content less than 11 or greater than 19 were discarded. Spike-in and normalization candidates were also discarded if they exceeded a certain degree of complementarity or sequence identity (total matches and maximum consecutive matches) with (i) the Illumina paired-end sequencing primers, (ii) the 3'-end of the sequencing primers, (iii) the whole *E. coli* genome [K-12 MG1655 strain (U00096.2)], and (iv) all other generated spike-in candidates in the same fashion as barcode design. The final population consisted of 40 spike-in and normalization DNA candidates, of which three were chosen at random seven times (without replacement) and con-

catenated, with one deletion at the 60th base of strand A (corresponding to the 31st base of strand B) and an addition of a single A to the end of the sequence to form seven 90-base spike-in DNA sequences and one normalization DNA sequence. Both strands of the 5'-end phosphorylated DNA oligos were ordered from Integrated DNA Technologies ([Dataset S4](#)) and were annealed in 0.3× NEB Buffer 2 (NEB) with 50 μM of each strand using the annealing temperature gradient. All seven spike-in sequences were ligated to the barcoded adapter mixture in NEB Next Quick Ligation Reaction Buffer (NEB) with 6.7 μM annealed spike-in, 6.7-μM barcoded adapter, and 6 μM Quick Ligase (NEB) by incubating at 25 °C for 30 min. The product was run on a 5% polyacrylamide gel (Bio-Rad), and the targeted band (at ~270 bp) was removed from the gel. The gel slice was cut into small pieces and the embedded DNA was extracted into diffusion buffer (10 mM Tris-Cl pH 8.0, 50 mM NaCl, 0.1 mM EDTA) by overnight incubation at room temperature. Then, the extracted spike-ins were purified on a column (Zymo Research). Sequence analysis (GeneWiz) confirmed that the band contained the expected ligation product. The concentration of each spike-in was estimated by qPCR (Fast SYBR Master Mix; Applied Biosystems) using sequence-specific qPCR primers ([Dataset S4](#)) against a known-concentration Y-shaped Reference DNA (below). The concentrations of spike-ins for the second deepsequencing run were measured in parallel by digital PCR (Fluidigm) at the Molecular Genetics Core Facility of Children's Hospital Boston Intellectual and Developmental Disabilities Research Center. Each spike-in was measured a total of ten times on two separate chips (48.770).

Design and Preparation of Y-Shaped Reference DNA. From the original list of 150 barcode candidates, we chose two barcodes that were not present in the final list of 145 used. Then, we concatenated the two barcodes with the Y-shaped adapter sequences and a 90-bp targeted sequence mimic such that the targeted sequence mimic was between the barcodes, which were between the adapters ([Dataset S4](#)). The 90-bp targeted-sequence mimic was designed the same way as the spike-in and normalization DNAs. Both strands of the DNA oligos were ordered from Integrated DNA Technologies and their concentrations were measured by absorbance at 260 nm using the extinction coefficient provided by Integrated DNA Technologies. The DNA oligos were annealed in water with 5 μM of each strand using the annealing temperature gradient.

***E. coli* RNA Preparation and cDNA Generation.** *E. coli* [K-12 MG1655 strain (U00096.2)] was grown overnight at 30 °C in LB medium. The resulting culture was diluted 500-fold in fresh LB medium and grown at 30 °C for 3.5 h, such that the O.D. at 600 nm became 0.30–0.35. One milliliter of cells were quickly killed by addition of 0.1 mL stop solution [90% (vol/vol) ethanol and 10% (vol/vol) phenol]. The cells were collected by centrifugation (9,100 × g, 1.5 min, room temperature), suspended in 1 mL cooled PBS (Lonza), and centrifuged again (16,000 × g, 1.5 min, room temperature). The supernatant was removed and the cells were suspended in 0.1 mL of 1 mg/mL lysozyme in TE Buffer (pH 8.0) (Ambion). Next, 0.1 mL of lysis buffer (Genosys) was added and the mixture was vortexed for 5 s. After adding 0.2 mL of phenol chloroform pH 4.5 (Sigma) and vortexing three times for 5 s, the mixture was centrifuged (16,000 × g, 3 min, room temperature). The top layer of solution was taken and 0.15 mL of 100% 2-Propanol (Sigma) was added; the mixture was left on ice for 30 min. The solution was centrifuged (16,000 × g, 30 min, 4 °C) to precipitate the RNA. The RNA pellet was washed twice by centrifugation (16,000 × g, 5 min, 4 °C) with 0.75 mL of cold 70% (vol/vol) ethanol. After the second centrifugation, the supernatant was removed and the pellet was dried for 15 min at room temperature. Then, 88 μL of water was added and the

mixture was incubated for 15 min at room temperature, followed by resuspension. The resulting solution was mixed with 0.04U/μL DNase I (NEB) in DNase I Reaction Buffer (NEB) for a total volume of 100 μL, and the mixture was incubated at 37 °C for 30 min followed by addition of EDTA (Sigma) to a final concentration of 5 mM. The mixture was incubated at 75 °C for 10 min to inactivate DNase I, followed by column purification. Ribosomal RNA was removed using Ribo-Zero rRNA Removal Kit (Gram-negative bacteria) (Epicentre, Illumina). From this point, we followed the conventional Illumina protocol for mRNA sequencing sample preparation with a few modifications. The purified RNA was fragmented in 0.5× fragmentation buffer (Ambion) with ~500 ng RNA in a 100-μL reaction solution. The solution was incubated on ice for 1 min after the fragmentation buffer was added, followed by a 6-min incubation at 70 °C. The tube was placed on ice and incubated for 1 min followed by addition of 4-μL stop solution (Ambion). The fragmented RNA was purified with a column and eluted in 11.1 μL in water. One microliter of 50 μM Random Hexamer Primer (Applied Biosystems) was added to this solution and incubated at 65 °C for 5 min and then placed on ice. Four microliters 5× First Strand Buffer (Invitrogen), 2 μL 100 mM DTT (Invitrogen), 0.4 μL 25 mM dNTP Mix (Applied Biosystems), and 0.5 μL RNase inhibitor (Applied Biosystems) was added to the mixture. This mixture was then incubated at 25 °C for 2 min, followed by the addition of 1 μL SuperScript II (Invitrogen). The mixture was then incubated at 25 °C for 10 min, 42 °C for 50 min, and 70 °C for 15 min to synthesize the first strand of the cDNA and inactivate the enzyme, which was placed on ice and then purified on a column. The eluate from this column was used to generate the second strand of cDNA in NEB Next Second Strand Synthesis Reaction Buffer (NEB) with 0.3 U/μL DNA polymerase I (*E. coli*) (NEB), 1.25 U/μL *E. coli* DNA Ligase, and 0.25 U/μL RNase H in an 800 μL total volume solution at 16 °C for 2.5 h, followed by the column purification. The eluted double stranded cDNA was end-repaired in T4 DNA Ligase Buffer (NEB) with 0.4 mM Deoxynucleotide Solution Mix (NEB), 0.5 U/μL T4 DNA polymerase (NEB), 0.5 U/μL T4 Polynucleotide Kinase (NEB) in a 200 μL reaction solution by incubating at 20 °C for 30 min followed by column purification. The eluted end-repaired cDNA was dA-tailed in NEB 2 buffer (NEB) with NEB Next dA-tailing Reaction Buffer with 1mM dATP (NEB) and 0.3 U/μL Klenow Fragment (3' → 5' exo-) in a 50-μL solution by incubating at 37 °C for 30 min followed by column purification.

Sample-Adapter Ligation, Sequencing Sample Preparation, and Sequencing. The cDNA library was ligated to the barcoded adapter mixture and the conventional Illumina paired-end adapter (without phosphorothioate bond) (Integrated DNA Technologies) respectively in the NEB Next Quick Ligation Reaction Buffer (NEB) with 5.4 μL of cDNA produced above, 1.9 μM barcoded adapter (or conventional Illumina paired-end adapter), and 3.6 μM Quick Ligase (NEB), in a total volume of 10 μL by incubating at 25 °C for 15 min. The two solutions were separately run on a 5% polyacrylamide gel and the portion between 250 and 300 bp was cut. The gel slice was cut into small pieces and the embedded DNA was extracted into diffusion buffer (10 mM Tris-Cl pH 8.0, 50 mM NaCl, 0.1 mM EDTA) by overnight incubation at room temperature. Then, the extracted DNAs were column purified. The concentrations of the purified products were measured by qPCR (Fast SYBR Master Mix; Applied Biosystems) against a known-concentration Y-shaped reference sequence ([Dataset S4](#)) using our designed qPCR primers ([Dataset S4](#)) purchased from Integrated DNA Technologies. The sample ligated to the barcoded adapter and the conventional Illumina Paired-end adapter were amplified by PCR (1 cycle of 98 °C for 1 min, 18 cycles of 98 °C for 1 s, 65 °C

for 45 s, 72 °C for 40 s, and 1 cycle of 72 °C for 5 min) in HF buffer (NEB) with 0.63 mM dNTP, 0.5 mM of each amplification primer modified from the Illumina PCR primers PE 1.0 and 2.0 (Integrated DNA Technologies) (Dataset S4), 25 fM DNA sample, and Phusion DNA polymerase (NEB) in 20 µL, with spike-in DNAs (0.71 aM Spike-in 1, 1.0 aM Spike-in 2, 4.4 aM Spike-in 3, 18 aM Spike-in 4, 150 aM Spike-in 5, 480 aM Spike-in 6 for the first sequencing run, and 1.3 aM Spike-in 1, 6.9 aM Spike-in 3, 36 aM Spike-in 4, 150 aM Spike-in 5, 770 aM Spike-in 7 for the second sequencing run). Then, 10 pM Normalization DNA was added to both PCR products, and the DNA was purified twice on a column. The concentration of the purified product was measured by qPCR (Fast Fast SYBR Master Mix; Applied Biosystems) using the conventional Illumina qPCR primers (Integrated DNA Technologies) against a PCR product amplified from Y-shaped reference DNA using modified Illumina PCR primers whose concentration was measured by NanoDrop (LMS). The final concentration of each spike-in and normalization DNA in the purified products was measured by qPCR using sequence-specific qPCR primers described above and compared by normalization. The length distribution of the purified PCR product was measured by Bioanalyzer (Agilent). Samples with barcoded adapters were sequenced on an Illumina HiSeq 2000 with 2 × 100 (for the first sequencing run) and 2 × 50 (for the second) base paired-end reads in one lane. Sequence data has been deposited in the Gene Expression Omnibus (<http://www.ncbi.nlm.nih.gov/>, accession no. GSE34449).

Genome Viewer Files. We have generated genome viewer files (.sam format) and index files (.sai) for the IGV software for fragments from both the digital counting method and the con-

ventional counting method. This software allows visualization of the difference between both methods in a genome-wide manner. The supplementary files have been deposited in the Gene Expression Omnibus (<http://www.ncbi.nlm.nih.gov/>, accession no. GSE34449).

Spike-in Analysis. From the raw sequencing data, we isolated reads which contained barcode sequences that corresponded to our original list of 145 barcodes in both forward and reverse reads for each sequencing cluster that had at most one mismatch. We then aligned the first 28 bases (26 bases for the second sequencing run) of the targeted sequence of both the forward and reverse reads of each cluster to each spike-in sequence, which is known. Sequences with more than two mismatches were discarded. We then counted the number of unique tags present in each spike-in to determine the number of copies of each spike-in.

Comparison of Noise in Conventional vs. Digital Counting for the *E. coli* Transcriptome. We summed the total number of reads and the total number of digital counts for each base in each of the mapped sequences. For each transcription unit, we created bins that were 99-bp long and summed the total number of reads and the total number of digital counts present in each bin. Bins that yielded an average number of digital counts per base of greater than or equal to 1 were selected, and from these bins we calculated the average and sample SD of the summed reads and summed digital counts, respectively. We define the noise to be the sample SD divided by the mean and calculated this value for both reads and digital counts; the ratio of the noise for reads to digital counts was computed for each transcription unit (Fig. 3D).

- Casbon JA, Osborne RJ, Brenner S, Lichtenstein CP (2011) A method for counting PCR template molecules with application to next-generation sequencing. *Nucleic Acids Res* 39:e81.
- Kinde I, Wu J, Papadopoulos N, Kinzler KW, Vogelstein B (2011) Detection and quantification of rare mutations with massively parallel sequencing. *Proc Natl Acad Sci USA* 108:9530–9535.
- Nakamura K, et al. (2011) Sequence-specific error profile of Illumina sequencers. *Nucleic Acids Res* 39:e90.
- Ribeiro-dos-Santos A, et al. (2010) Ultra-deep sequencing reveals the microRNA expression pattern of the human stomach. *PLoS ONE* 5:e13205.

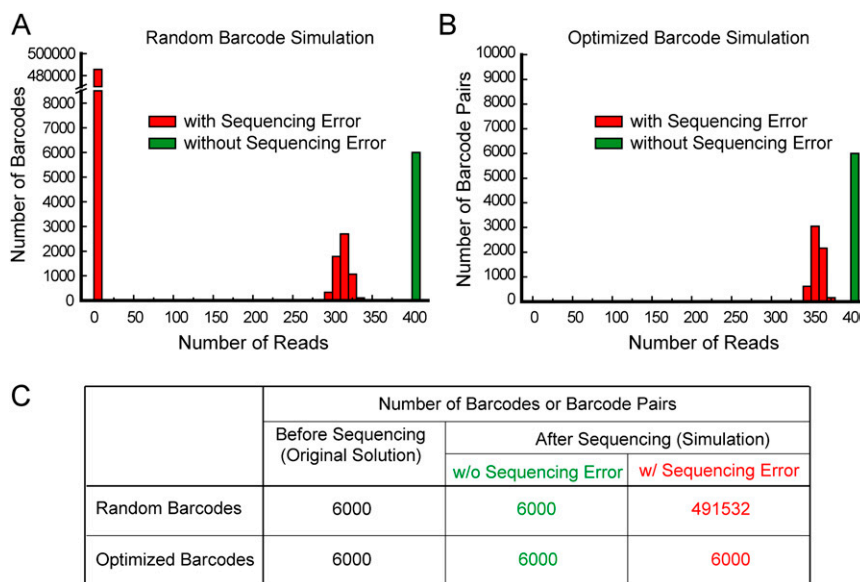


Fig. S1. Histogram of unique barcodes from parallel simulation of the theoretical library using optimized barcodes and random barcodes (*SI Materials and Methods*). (A) A library of 6,000 identical molecules was barcoded *in silico* by random barcodes (single-end 16 nucleotides). (B) Same as A, but with optimized barcodes (paired-end 20 nucleotides). Both conditions were then “amplified” 400-fold. After accounting for sequencing error in both the barcode region and the mapping region, the resulting number of unique barcodes and the respective number of reads were histogrammed for each case. (C) In the random barcode case, many artifactual barcode sequences were generated by sequencing error and thus have a low number of reads. Sequencing error occurred for the optimized barcode sequence as well. However, because the optimized barcode sequences were known (i.e., unused sequences are known) and designed to have minimally overlapping sequences, the sequencing errors were identified. Then, sequenced reads which contained errors were removed (*Materials and Methods*), and there is no population of artifactual barcode sequences with a small number of sequencing reads.

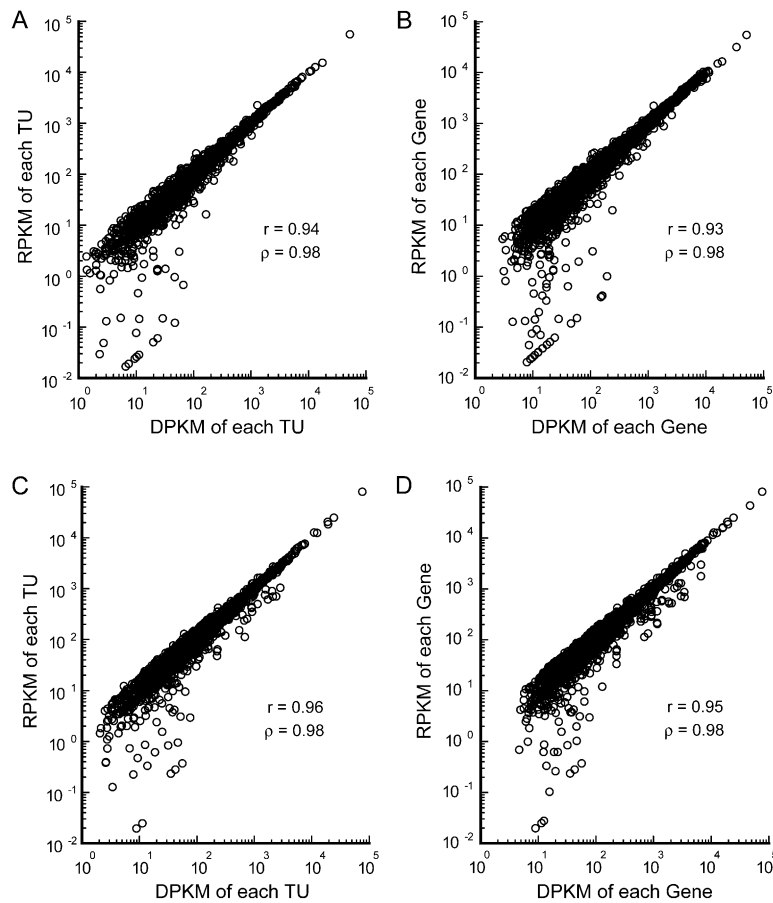


Fig. S2. (A) Comparison of uniquely mapped reads per kilobase of each transcription unit per million total uniquely mapped reads (RPKM) and uniquely mapped digital counts per kilobase of each transcription unit per million total uniquely mapped molecules (DPKM) for all detected transcription units. (B) The same plot as A, but for genes. (C and D) The same plots as A and B, respectively, from the second sequencing run. Genes and transcription units exhibit stronger correlations between RPKM and DPKM compared with individual fragment sequences (Fig. 3C) because of the averaging of reads and counts over long genes or transcription units. At higher copy numbers, analog reads and digital counts are well-correlated, but at lower copy numbers the data are not nearly as correlated (see [Datasets S5](#), [S6](#), and [S7](#) for specific values for each datapoint).

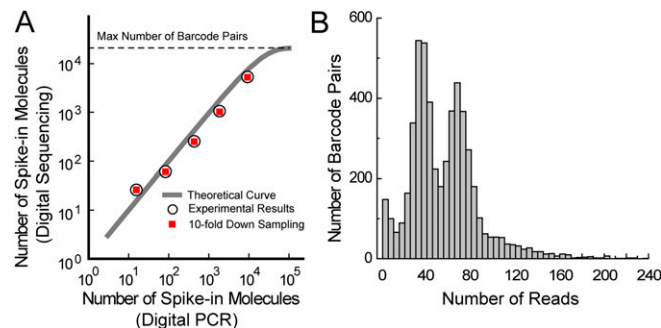


Fig. S3. Down-sampling of all spike-in data by a factor of 10 and the resulting digital counts obtained (*SI Materials and Methods*). (A) The correlation between the number of digital counts of all down-sampled spike-in molecules and the number of spike-in molecules measured by digital PCR. For each of the spike-in sequences, we still detected an average of 98% of the original number of unique barcodes. (B) The histogram of conventional reads for all unique barcodes in the down-sampled population of the highest copy spike-in (spike-in 7).

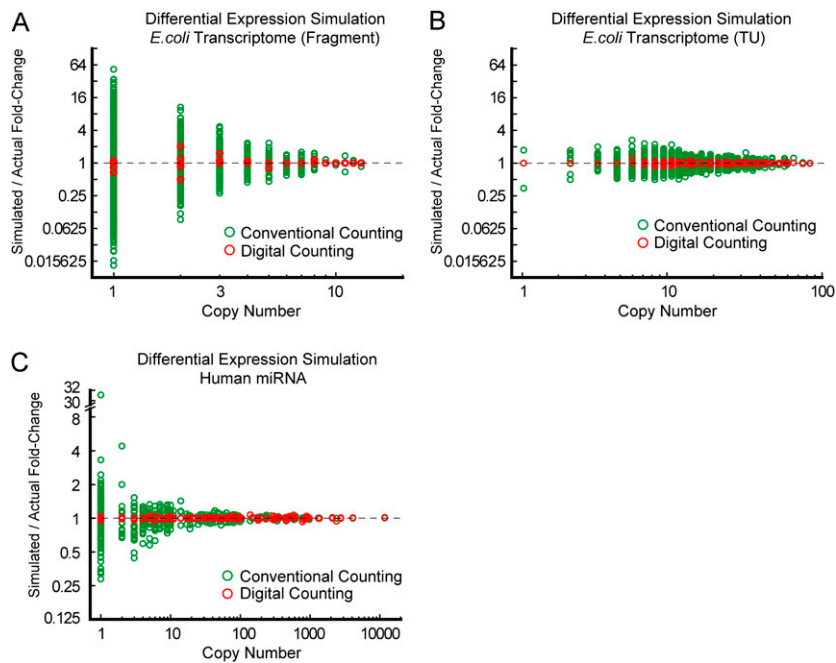


Fig. S4. Simulation demonstrating the superior performance of digital counting over conventional counting in differential expression analysis (*SI Materials and Methods*). RNA expression quantification was simulated using experimentally measured copy numbers, barcode sampling, and amplification noise distributions for two different libraries for each of three different systems [*E. coli* transcriptome fragments (*A*), *E. coli* transcription units (*B*), and human stomach microRNA (*C*)]. We plot the ratio of simulated to actual fold-change for each gene as a function of the lower of two copy numbers for the two compared libraries. Ideally, the value of this ratio is one for all genes. Because digital counting is almost completely immune to amplification noise, it gives consistently superior performance to conventional counting for differential expression, even at low copy numbers. We note that the discrepancy between conventional and digital counting is smaller for the *E. coli* transcription unit library in *B* than for the fragment library in *A* because amplification noise can be averaged over many fragments in the case of long transcription units.

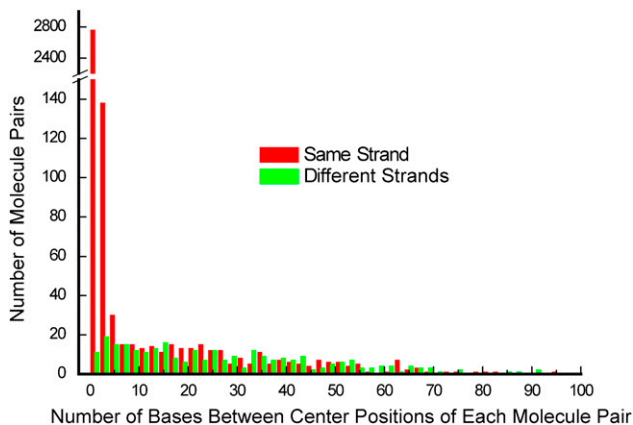


Fig. S5. Histogram of the number of bases between the center positions of all pairs of molecules mapped to the same transcription unit that contain the same barcode for pairs of molecules both mapped to the sense or antisense strand of the *E. coli* genome (red), and also for pairs of molecules mapping to different strands of the *E. coli* genome (green). We are able to distinguish the different strands of an original cDNA molecule because of the design of the paired-end sequencing adapters (Fig. 1*B*). The distribution for molecule pairs mapping to different strands is more uniform than the distribution for molecule pairs mapping to the same strand, which confirms that molecule pairs mapping to different strands on the *E. coli* genome are on average far apart from each other. From the distribution of molecule pairs mapping to the same strand, we determined that molecule pairs with fewer than four bases between center positions are identical molecules. This analysis allows us to keep 98% of all molecules but reduces the counting error rate.

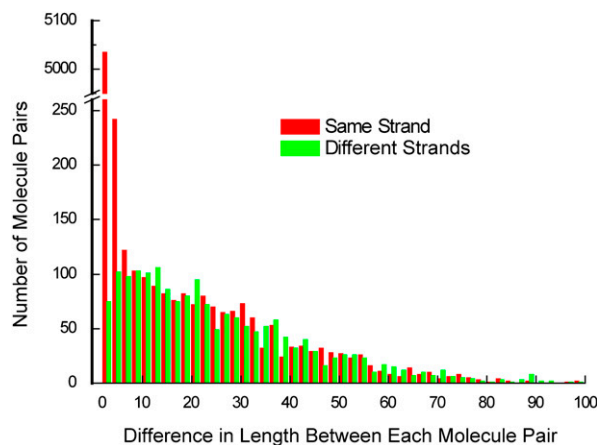


Fig. S6. Histogram of the difference in fragment length for all pairs of molecules mapped to the same transcription unit that contain the same barcode for pairs of molecules both mapped to the sense or antisense strand of the *E. coli* genome (red), and also for pairs of molecules mapping to different strands of the *E. coli* genome (green). The distribution for molecule pairs mapping to different strands is more uniform than the distribution for molecule pairs mapping to the same strand, which confirms that molecule pairs mapping to different strands on the *E. coli* genome are on average similar in size. From the distribution of molecule pairs mapping to the same strand, we determined that molecule pairs with a difference in fragment length of less than nine are identical molecules. This analysis allows us to keep 96% of all molecules but reduces the counting error rate.

Other Supporting Information Files

[Dataset S1 \(XLS\)](#)

[Dataset S2 \(XLS\)](#)

[Dataset S3 \(XLS\)](#)

[Dataset S4 \(XLS\)](#)

[Dataset S5 \(XLS\)](#)

[Dataset S6 \(XLS\)](#)

[Dataset S7 \(XLS\)](#)