

Induro™, a Novel Reverse Transcriptase for Nanopore Direct RNA Sequencing with Significantly Improved RNA 5' Coverage

Luo Sun, Yan Xu, Bradley W. Langhorst, Jian Sun, Nicole M. Nichols, Eileen T. Dimalanta, and Theodore B. Davis

New England Biolabs, Inc.



INTRODUCTION

Transcriptome analysis based on cDNA sequencing of direct or PCR amplified cDNA has limitations, including generation of cDNA chimeras during library preparation, artifacts associated with template switching, DNA contamination and size bias during PCR amplification. Recent advances from Oxford Nanopore Technologies® (ONT) have enabled direct RNA sequencing without PCR. This method not only interrogates full-length RNA molecules and generates higher resolution data on complex isoforms, but also allows for the detection of RNA modifications. This opens up new opportunities for studying RNA stability and regulation. To perform direct RNA sequencing, one approach is to first convert RNA into an RNA/cDNA hybrid molecule using a reverse transcriptase (RT). The RNA/cDNA hybrid molecule reduces RNA secondary structure formation, leading to higher sequencing throughput and longer sequencing reads without sacrificing RNA modification information. To accomplish this, an ideal reverse transcriptase would be highly processive and generate high cDNA yields.

Here we introduce Induro, a novel reverse transcriptase that delivers superior performance in the cDNA synthesis reaction. Using this novel RT, we generated a RNA/cDNA hybrid library starting with mouse brain polyA-enriched RNA. After ligation with ONT sequencing adapter following the SQK-RNA002 protocol, this library was sequenced on the GridION. When compared with other commonly used RTs, Induro produced longer average read lengths and better 5' end coverage of RNA transcripts. Similar results were also obtained by performing direct sequencing on a cDNA library generated using this RT enzyme. In summary, Induro is a significant improvement compared to other available RT enzymes. It can produce longer cDNA molecules, and exhibits higher cDNA yields, better full-length coverage, leading to more accurate transcriptome analysis.

METHODS

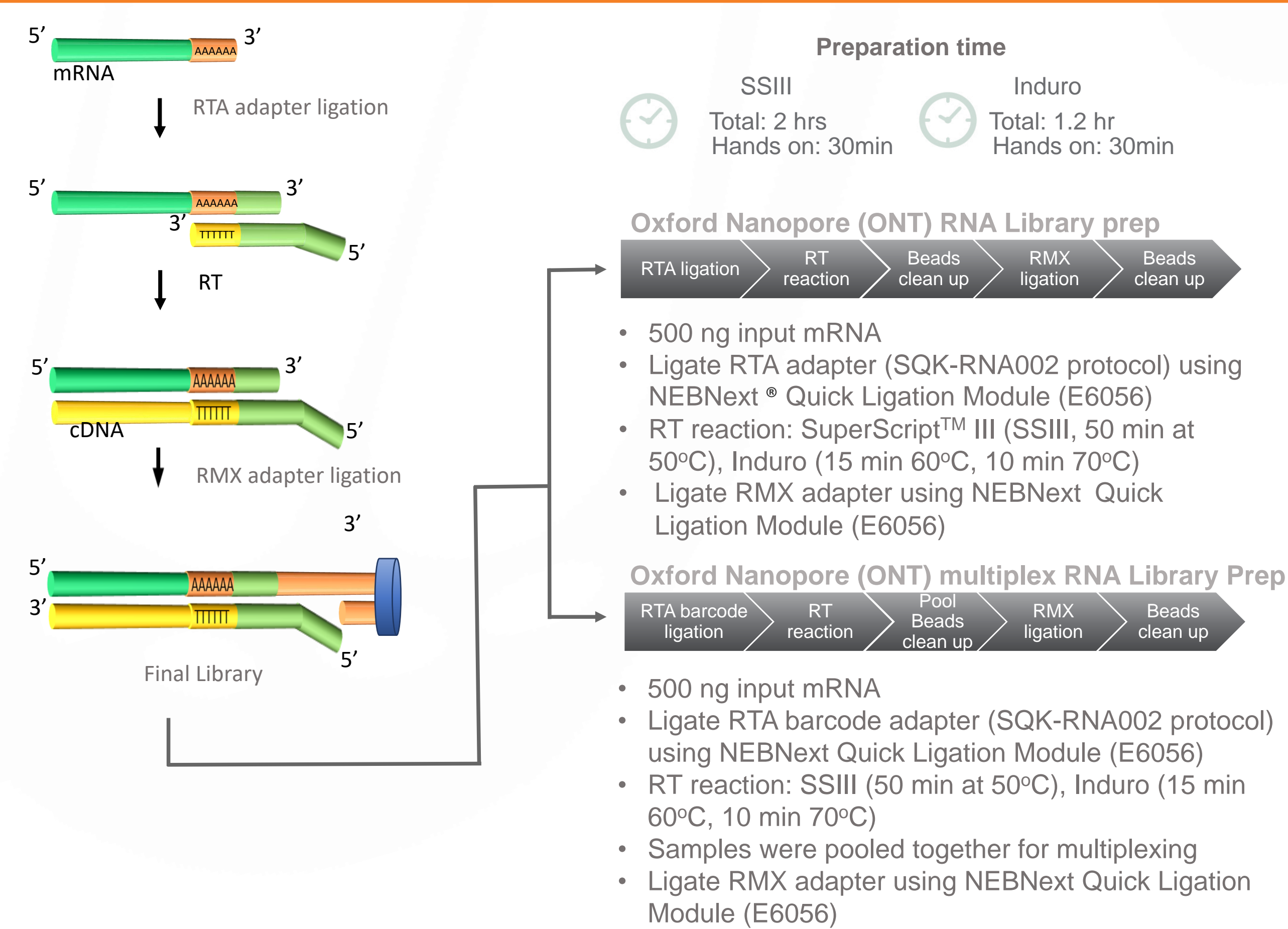


Figure 1. Direct RNA sequencing libraries were prepared following the Oxford Nanopore SQK-RNA002 protocol. 500ng mouse brain polyA (TaKaRa, cat # 636207) was used as input to make RNA/cDNA libraries with either Induro or SSIII. RNA libraries without cDNA synthesis (no_RT) were also generated. All three methods (no cDNA synthesis, or cDNA synthesis with Induro or SSIII) were used to make single-plex and multiplex libraries. Final Library yields were measured using Invitrogen™ Qubit™ dsDNA HS Assay Kit.

RESULTS

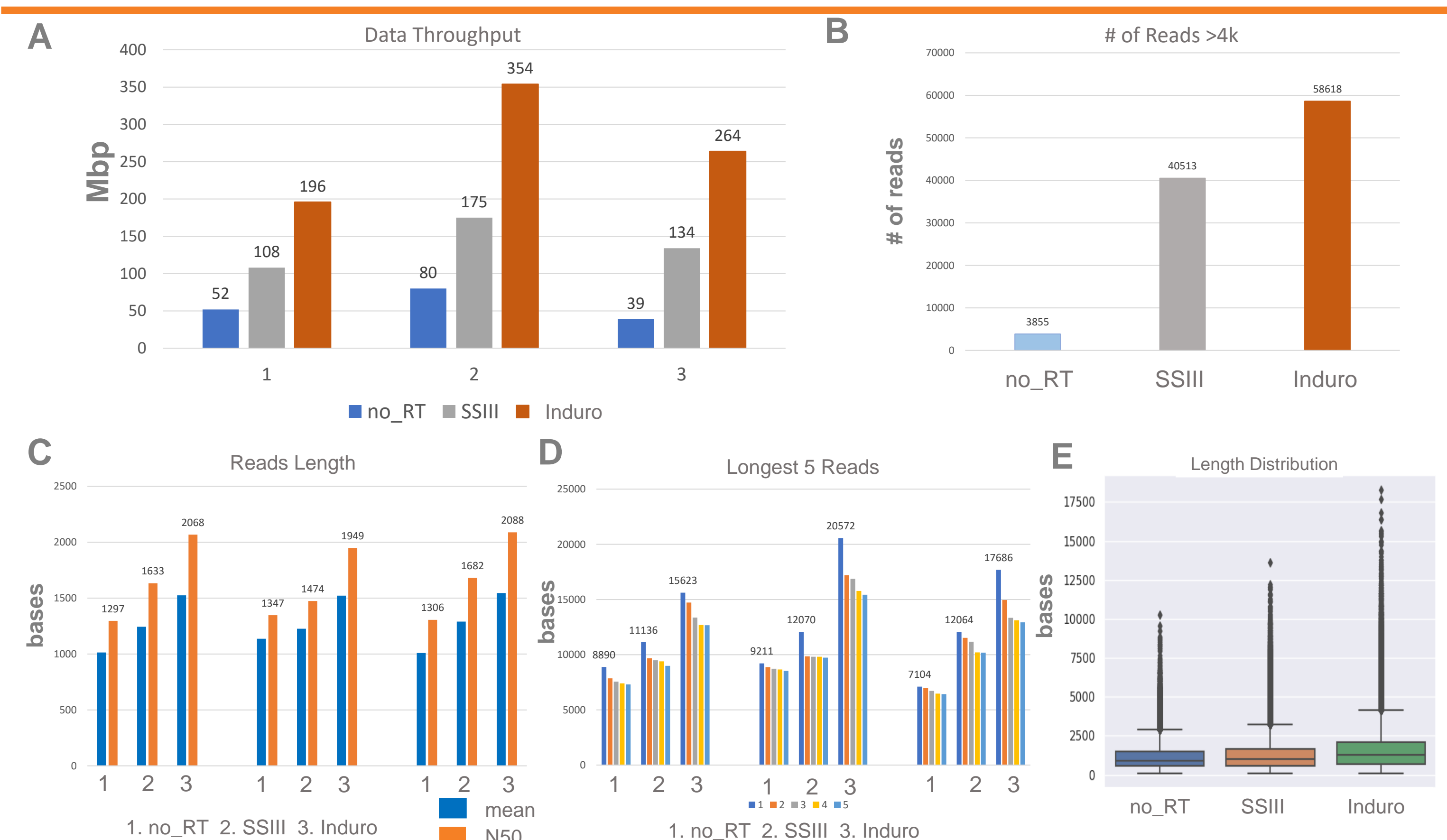


Figure 2. 3 sequencing runs were performed on pooled barcoded samples containing the direct RNA library (no_RT), the RNA/cDNA libraries made with either SSIII or Induro in the RT step. (A) Data throughput for each of the sequencing runs. The Induro libraries have significantly and consistently higher base-called data (orange bar). Qubit measurements of double stranded RNA/cDNA yield with Induro is higher than with SSIII, the ONT recommended RT (data not shown). (B) All reads from 3 multiplexed runs and two single-plex runs/sample were merged together for data analysis. The merged data was mapped to mouse mm10 genome using minimap2. Reads (358Mb for no_RT, 3.6Gb for SSIII, 3.2Gb for Induro) were filtered for reads with length >4kb. The number of reads >4kb long is 45% higher with the Induro library (orange bar, 2.6% of total reads) than from SSIII (gray bar, 1.3% of total reads). (C) Data analysis for mean and median read lengths shows that the Induro libraries have significantly longer mean length and N50 values. The no_RT libraries have the lowest mean length and N50 values. (D) Analysis of top 5 longest reads demonstrates that Induro libraries have the longest reads up to 20kb, significantly longer than reads from SSIII and no_RT libraries. (E) Boxplot of read length distribution shows Induro libraries have longer median read lengths. Also observed is the significantly longer tail for Induro libraries, indicating more long reads.

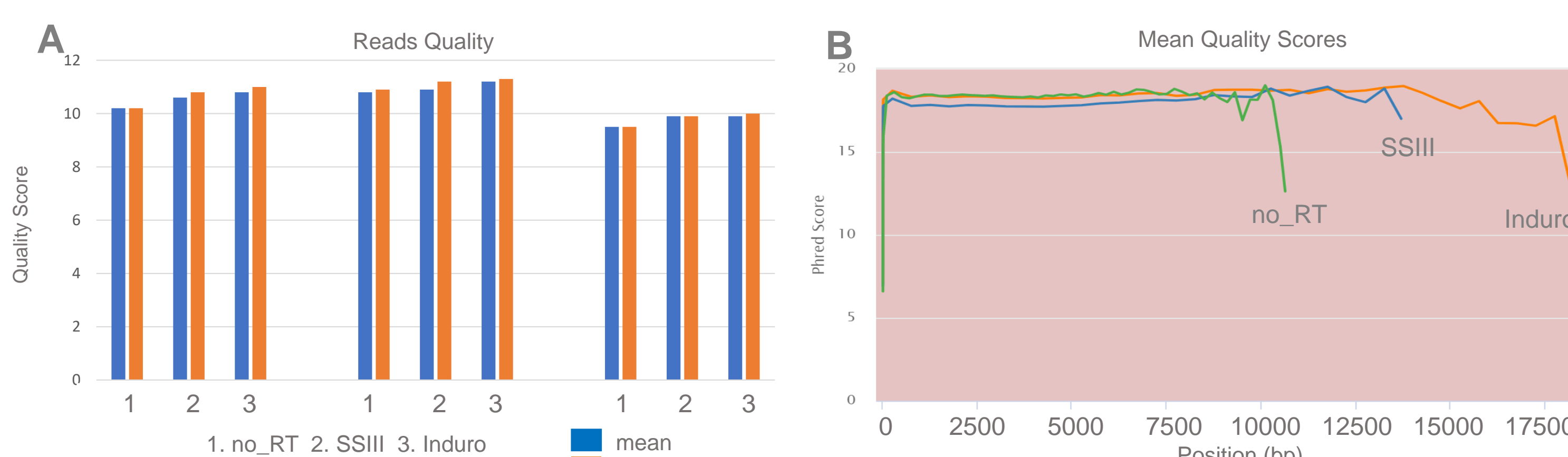


Figure 3. Read Quality (A) Read quality is similar between the 3 different libraries, and varied more by sequencing run than by library type. (B) Quality score as a function of base position is shown. Induro libraries maintained a higher mean quality score at lengths > 1kb compared to other samples.

RESULTS

no_RT	SSIII	Induro
98.9% reads aligned	99.3% reads aligned	99.5% reads aligned
Read Stats	Read Stats	Read Stats
Total reads 146,771	Total reads 314,671	Total reads 506,826
- Unaligned reads 1,542 (1.1%)	- Unaligned reads 2,093 (0.7%)	- Unaligned reads 2,710 (0.5%)
- Aligned reads 145,229 (98.9%)	- Aligned reads 312,578 (99.3%)	- Aligned reads 504,116 (99.5%)
--- Single-align reads 144,817 (98.7%)	--- Single-align reads 311,605 (99.0%)	--- Single-align reads 502,473 (99.1%)
--- Gapped-align reads 35 (0.02%)	--- Gapped-align reads 107 (0.03%)	--- Gapped-align reads 298 (0.06%)
--- Chimeric reads 377 (0.26%)	--- Chimeric reads 866 (0.28%)	--- Chimeric reads 1,345 (0.27%)
----- Trans-chimeric reads 377 (0.26%)	----- Trans-chimeric reads 863 (0.27%)	----- Trans-chimeric reads 1,340 (0.26%)
----- Self-chimeric reads 0 (0.00%)	----- Self-chimeric reads 3 (0.00%)	----- Self-chimeric reads 5 (0.00%)
Base Stats (of aligned reads)	Base Stats (of aligned reads)	Base Stats (of aligned reads)
Total bases 155,340,385	Total bases 394,326,079	Total bases 782,448,115
- Unaligned bases 10,349,985 (6.7%)	- Unaligned bases 19,964,900 (5.1%)	- Unaligned bases 35,776,076 (4.6%)
- Aligned bases 144,990,400 (93.3%)	- Aligned bases 374,361,179 (94.9%)	- Aligned bases 746,672,039 (95.4%)
--- Single-aligned bases 144,775,569 (93.2%)	--- Single-aligned bases 373,806,638 (94.8%)	--- Single-aligned bases 745,665,454 (95.3%)
--- Other-aligned bases 14,162 (0.01%)	--- Other-aligned bases 48,854 (0.01%)	--- Other-aligned bases 159,130 (0.02%)

Figure 4. Alignment Analysis
All reads from 3 multiplex runs were merged for no_RT, SSIII and Induro libraries. Merged data was mapped to the mouse genome mm10 using minimap2. Mapping quality was assessed using Alignqc. Result shows very similar alignment statistics for all 3 samples. Induro and SSIII libraries have a slightly higher percentage of reads and bases aligning to the genome. Chimeric read rates were similar for all 3 libraries.

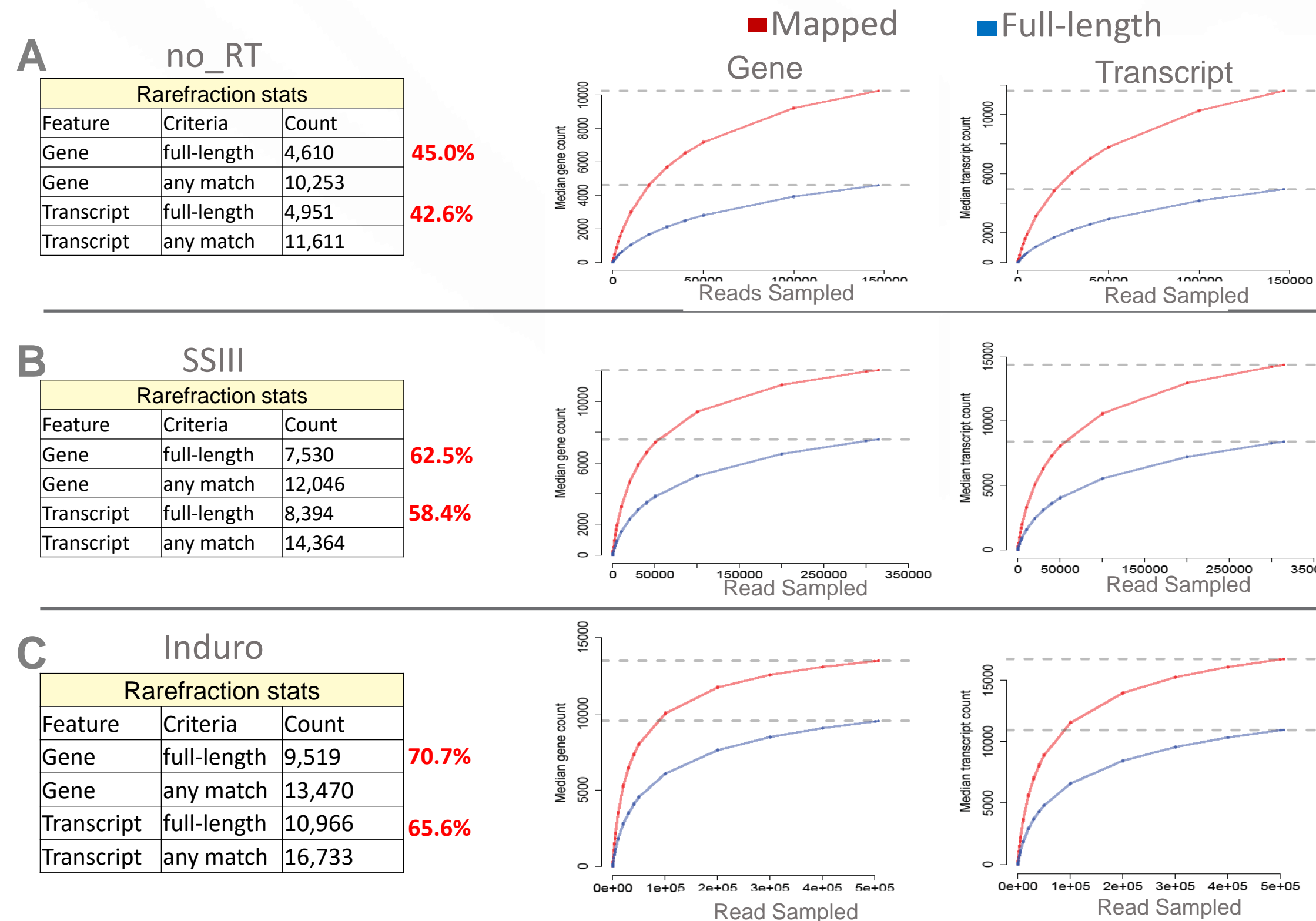


Figure 5. Annotation and rarefaction analysis. Mapped data was analyzed using Alignqc for the number of detected genes and transcripts for no_RT (A), SSIII (B), and Induro (C) libraries. The percentage of full-length genes detected is 70.7% from Induro vs 62.5% from SSIII and 45.0% from no_RT libraries. Similarly, the percentage of full-length transcripts detected was higher for Induro libraries (65.6%) compared to SSIII (58.4%) or no_RT (42.6%) libraries. This data demonstrates that higher full-length genes and transcripts were obtained with Induro.

no_RT	SSIII	Induro
Alignment stats	Alignment stats	Alignment stats
Best alignments sampled 901	Best alignments sampled 801	Best alignments sampled 701
Base stats	Base stats	Base stats
Bases analyzed 1,002,927	Bases analyzed 1,050,043	Bases analyzed 1,113,365
- Correctly aligned bases 897,968 (89.5%)	- Correctly aligned bases 945,093 (90.0%)	- Correctly aligned bases 1,003,544 (90.1%)
- Total error bases 104,959 (10.465%)	- Total error bases 104,950 (9.995%)	- Total error bases 109,821 (9.864%)
--- Mismatched bases 26,282 (2.621%)	--- Mismatched bases 27,913 (2.658%)	--- Mismatched bases 29,214 (2.624%)
--- Deletion bases 49,267 (4.912%)	--- Deletion bases 50,477 (4.807%)	--- Deletion bases 52,627 (4.727%)
--- Complete deletion bases 28,937 (2.885%)	--- Complete deletion bases 29,803 (2.838%)	--- Complete deletion bases 30,640 (2.752%)
--- Homopolymer deletion bases 20,330 (2.027%)	--- Homopolymer deletion bases 20,674 (1.969%)	--- Homopolymer deletion bases 21,987 (1.975%)
--- Insertion bases 29,410 (2.932%)	--- Insertion bases 26,560 (2.529%)	--- Insertion bases 27,980 (2.513%)
--- Complete insertion bases 18,786 (1.873%)	--- Complete insertion bases 16,025 (1.526%)	--- Complete insertion bases 16,745 (1.504%)
--- Homopolymer insertion bases 10,624 (1.059%)	--- Homopolymer insertion bases 10,535 (1.003%)	--- Homopolymer insertion bases 11,235 (1.009%)

Figure 6. Error statistics. A random down-sampled subset of sequencing data from no_RT, SSIII and Induro were mapped to the mouse genome mm10. Mapping data was analyzed for errors including mismatch, deletion, insertion and homopolymers using Alignqc. Similar error rates were seen with different libraries.

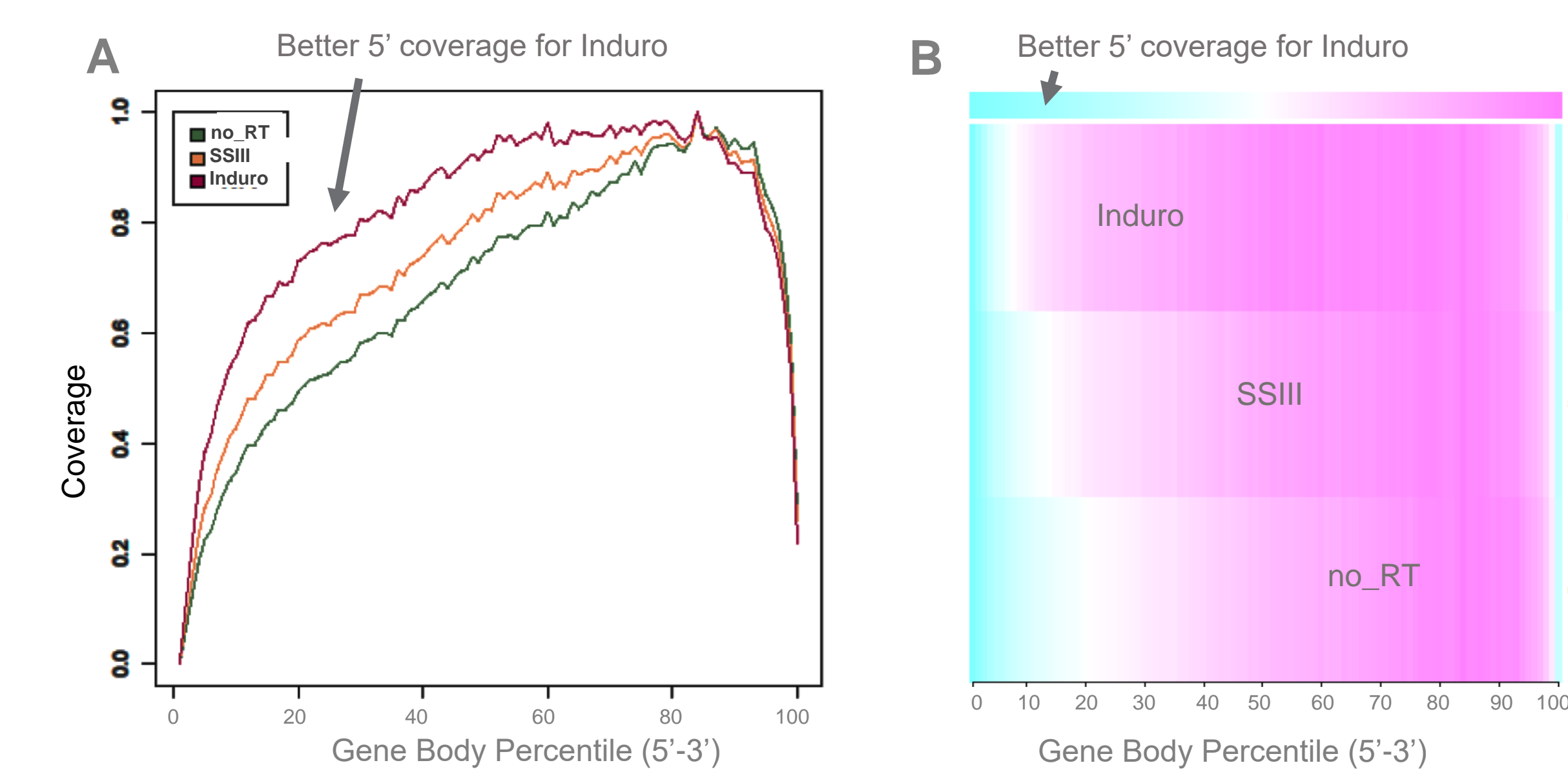


Figure 7. Gene body coverage percentile
Gene Body Coverage (A) and a heat map (B) were generated using RSeQC. In the heat map, the aqua color indicates low coverage and pink indicates high coverage. In both plots, it clearly shows that Induro libraries have significantly better 5' coverage than libraries made using SSIII or made without cDNA synthesis (no_RT).

Table 1. Direct cDNA Sequencing Using Different RT Enzymes

	Induro	SSIII	SSIV	Maxima	TGIRT	Marathon	ProtoScript II
Mean read length (bp)	2562	2118.7	2142	1893	2077	1971	1995
Median read length (bp)	1896	1665	1682	1441	1611	1573	1561
Number of reads	175400	97638	148859	343925	144746	170860	149240
Read length N50 (bp)	3658	2855	2926	2681	2859	2726	2703

Table 1. Impact of RT enzymes on the length of cDNA reads following recommended conditions was studied using direct cDNA sequencing on the GridION. In this experiment, after first strand cDNA was synthesized in the RT reaction, 2nd strand cDNA was synthesized using NEBNext Ultra II Non-Directional RNA Second Strand Synthesis Module. Then the cDNA products were purified and ligated to a barcoded adapter following the protocol of Nanopore EXP-NBD104 kit. cDNA library concentrations were measured on a Qubit after the barcoding step and then equal amounts of DNA were pooled and ligated with a sequencing adapter. Finally, the resulting DNA library was cleaned up with beads and run on a GridION with a R9.4.1 flow cell. The data was demultiplexed and analyzed. Results are consistent with direct RNA sequencing data in that reads from the Induro library has the longest mean, median and N50 values.

CONCLUSIONS

- Induro, a novel reverse transcriptase shows superior performance for generating libraries for direct RNA and direct cDNA Oxford Nanopore sequencing.
- Induro is very processive, requiring a short incubation time to generate full length cDNA.
- Longer read lengths, higher percentage of full-length genes and transcripts, improved 5' coverage can be achieved with Induro compared to SSIII without increasing error rates or decreasing mapping rates.

ACKNOWLEDGEMENTS

- We thank Dr. Maggie Heider for proof reading and Dr. Guoping Ren for the discussion and assistance in developing Induro.