

*Zajac et al.*

Base preferences in non-templated nucleotide incorporation by MMLV-derived reverse transcriptases

## **Supplementary Information**

Table S1:	2
Primers and oligonucleotides used in this work.	
Table S2:	3
p-Values for the performed optimizations.	
Table S3:	3
Transcript and spike sequences used for querying the Illumina sequencing reads.	
Table S4:	4
Correlation between the performed experiments.	
Table S5:	4
Composition of the ribo base portion of the TSO.	
Table S6:	5
Composition of DNA base in position 4, i.e. the DNA base adjacent to the ribobase stretch.	
Table S7:	6
Number of guanosines at the template-switching interface.	
Figure S1:	7
Regression analysis of TSO length data.	
Figure S2:	8
Hit distributions along the analyzed transcripts.	
Figure S3:	11
Composition of the ribo base portion of the TSO.	
Figure S4:	12
Composition of DNA base in position 4, i.e. the DNA base adjacent to the ribobase stretch.	
Figure S5:	13
Number of guanosines at the template-switching junction.	
Figure S6:	14
Barcode complexity for different transcripts.	
Figure S7:	15
Barcode complexity for different RNA spikes.	

Table S1

Name	Sequence	Length [bp]	Ordered from
<b>STRT Bio-T30VN</b>	5'-Bio-aagcagtggtatcaaccgagagtCGACTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTTN	58	Eurofins MWG Operon
<b>STRT v2-7</b>	5'-aagcagtggtatcaaccgagagtGCAGUGCUGGACATrGrGrG	40	Eurofins MWG Operon
<b>STRT v2-7-N2</b>	5'-aagcagtggtatcaaccgagagtUNNGGACATrGrGrG	35	Eurofins MWG Operon
<b>STRT v2-7-N4</b>	5'-aagcagtggtatcaaccgagagtUNNNNGGACATrGrGrG	37	Eurofins MWG Operon
<b>STRT v2-7-N6</b>	5'-aagcagtggtatcaaccgagagtUNNNNNNGGACATrGrGrG	39	Eurofins MWG Operon
<b>STRT v2-7-N8</b>	5'-aagcagtggtatcaaccgagagtUNNNNNNNNGGACATrGrGrG	41	Eurofins MWG Operon
<b>STRT v2-7-N10</b>	5'-aagcagtggtatcaaccgagagtUNNNNNNNNNNGGACATrGrGrG	43	Eurofins MWG Operon
<b>STRT v2-7-N12</b>	5'-aagcagtggtatcaaccgagagtUNNNNNNNNNNNGGACATrGrGrG	45	Eurofins MWG Operon
<b>STRT N10-rN3</b>	5'-aagcagtggtatcaaccgagaguNNNNNNNNNNrMrN	36	Integrated DNA Technologies
<b>STRT N10-rG3</b>	5'-aagcagtggtatcaaccgagaguNNNNNNNNNNrGrGrG	36	Integrated DNA Technologies
<b>STRT N12-rG3</b>	5'-gcagtggtatcaaccgagaguNNNNNNNNNNrGrGrG	36	Integrated DNA Technologies
<b>STRT-PCR</b>	5'-Bio-aagcagtggtatcaaccgagagt	23	Eurofins MWG Operon

**Primers and oligonucleotides used in this work.**

The sequences, lengths and vendors are indicated. All sequences are written from 5'. The lowercase letters indicate the STRT amplification handle. The underlined bases for the T30 oligonucleotide denote a SalI recognition sequence. Similarly, the underlined bases for the STRT v2-7 oligonucleotide indicate a BtsI recognition site. The barcode (used in the STRT protocol to uniquely label each cell's RNA) is shown in bold and underlined typeface. Ribo bases are preceded by an 'r' and are shown in italics. An 'N' is a DNA degenerate position and an 'rN' is a RNA degenerate position. The oligonucleotides were unmodified except for the T30 oligonucleotide and STRT-PCR primer that carried a 5'-biotin.

Table S2

	Comparison	p-value
TSO amount	10 nM vs. 200 nM	1.06E-04
	40 nM vs. 200 nM	3.44E-04
	1 $\mu$ M vs. 200 nM	8.77E-03
	2.5 $\mu$ M vs. 200 nM	2.06E-03
	5 $\mu$ M vs. 200 nM	7.09E-02
	1 $\mu$ M vs. 2.5 $\mu$ M	1.54E-02
	1 $\mu$ M vs. 5 $\mu$ M	0.973
	2.5 $\mu$ M vs. 5 $\mu$ M	0.133
SS enzyme	SSII vs. SSIII	1.80E-02
	SSII vs. Cycled SSIII	8.28E-03
	SSIII vs. Cycled SSIII	1.63E-02
	Cycled SSIII vs. NTC	0.861
SSII amount	200 U vs. 10 U	9.69E-03
	50 U vs. 10 U	0.470
	5 U vs. 10 U	6.95E-03
	1 U vs. 10 U	3.07E-03
	50 U vs. 200 U	1.07E-02
	5 U vs. 50 U	7.22E-03
	1 U vs. 50 U	2.95E-03

**p-Values for the performed optimizations.**

Student's unpaired t-test with a two-tailed distribution was used.

Table S3

	Transcript / spike	Sequence
Transcripts	MALAT1	AGGCATTGAGGCAGCCAGCGCAGGGGCTTC
	RPLP1	CCTTTCTCAGCTGCCGCCAAGGTGCTCGG
	MT2A	ACCACGCCTCCTCCAAGTCCCAGCGAACCC
	AHSG	CCTTTCCCAGCAGAGCACCTGGGTTGGTCC
	CNIH4	AGGAGCGGCGGCGACGGAGGAGGAGGATGG
Spikes	MC28	GGAATTCTCCAGATTACTTCCATTTCCGCC
	MJ-500-37	GGAATTCTGGACATTAATTAGGGCTGAAAG

**Transcript and spike sequences used for querying the Illumina sequencing reads.**

The human transcript sequences were obtained from definitions in the refFlat.txt file for hg19 from the UCSC Genome Browser. The spike sequences were provided by Life / Ambion. Apart from AHSG that starts with base number three in the transcript defined by refFlat.txt for hg19 and MALAT1 where the main template-switching peak is 1300 bp from the 5'-end, the other query sequences coincide with the defined 5'-end of the transcripts.

Table S4

<b>Reaction</b>	<b>RNA10G3</b>	<b>RNA12G3</b>	<b>RNA10N3</b>
RNA10G3	<b>1.000</b>	0.999	0.999 (0.484)
RNA12G3		<b>1.000</b>	1.000 (0.506)
RNA10N3			<b>1.000</b>

**Correlation between the performed experiments.**

R2 values are shown. These correlations are based on the five investigated transcripts: MALAT1, RPLP1, MT2A, AHSG and CNIH4. For the correlations involving RNA10N3, the MALAT1 gene was omitted as it generated an unusually low number of reads in this reaction. The R<sup>2</sup> values featuring MALAT1 are shown in parentheses.

Table S5

<b>Spike / Transcript</b>	<b>Guanosine percentage [%]</b>			<b>Percentage [%]</b>				
	<b>Position 1</b>	<b>Position 2</b>	<b>Position 3</b>	<b>AGG</b>	<b>CGG</b>	<b>GGG</b>	<b>TGG</b>	<b>NGG</b>
MC28	98	80	46	18	11	38	11	79
MJ-500-37	97	79	54	16	10	43	9	77
MALAT1	91	85	58	15	11	44	8	77
RPLP1	93	84	62	11	9	52	7	79
MT2A	91	87	63	10	12	53	4	79
<i>Average</i>	<i>94</i>	<i>83</i>	<i>57</i>	<i>14</i>	<i>11</i>	<i>46</i>	<i>8</i>	<i>78</i>

**Composition of the ribo base portion of the TSO.**

Table S6

		Position 4 percentage [%]				
	Spike / Transcript	A	C	G	T	Order
<i>ERCC10N3</i>	MC28	20	27	31	22	G > C > T > A
	MJ-500-37	18	25	35	22	
<i>ERCC10G3</i>	MC28	25	19	33	23	G > A > T > C
	MJ-500-37	24	20	35	22	
<i>RNA10N3</i>	MALAT1	16	33	32	19	G > C > T > A
	RPLP1	15	31	35	20	
	MT2A	14	35	38	14	
<i>RNA10G3</i>	MALAT1	21	22	36	22	G > C/T/A
	RPLP1	20	20	38	23	
	MT2A	19	20	40	21	
<i>RNA12G3</i>	MALAT1	15	22	38	25	G > T > C > A
	RPLP1	16	20	40	25	
	MT2A	12	22	40	26	

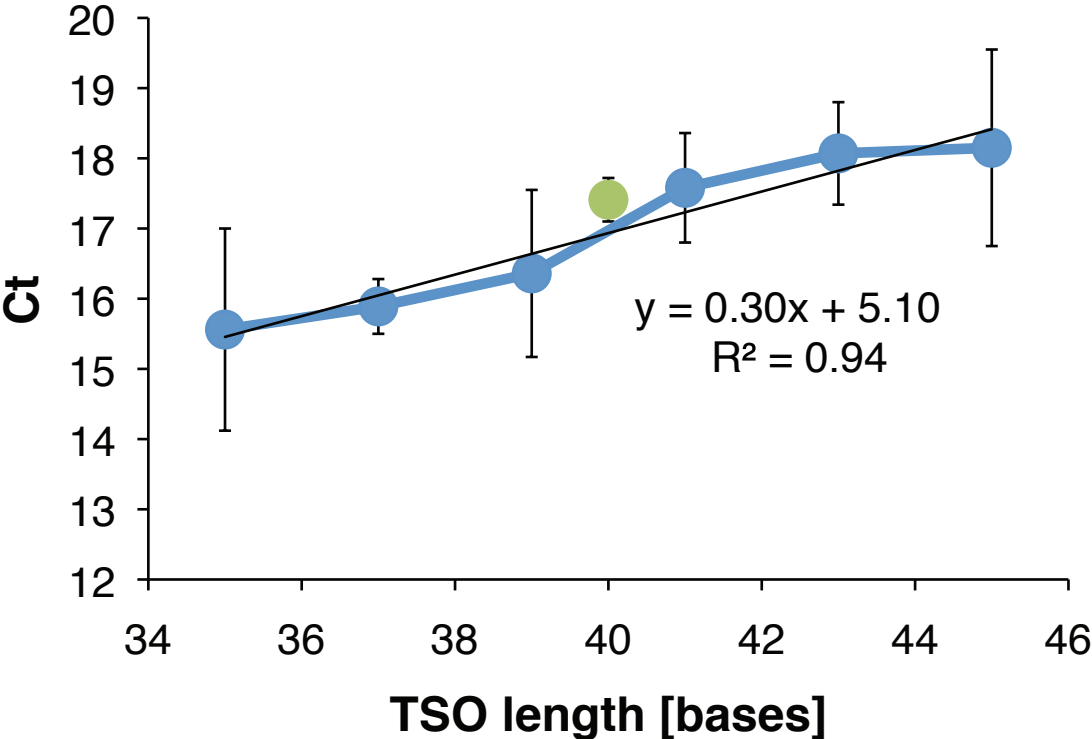
Composition of DNA base in position 4, i.e. the DNA base adjacent to the ribobase stretch.

Table S7

Number of guanosines	Percentage [%]											
	ERCC10N3		ERCC10G3		RNA10N3		RNA10G3		RNA12G3			
	MC28	MJ-500-37	MC28	MJ-500-37	RPLP1	MT2A	RPLP1	MT2A	RPLP1	MT2A		
2	18.86	20.31	0.65	0.76	20.14	22.56	0.48	0.31	0.51	0.39		
3	35.98	32.33	4.29	6.62	27.96	27.57	16.71	23.01	16.99	26.12		
4	28.18	28.54	62.88	59.93	32.95	32.22	56.85	58.40	57.73	57.67		
5	14.37	16.69	28.08	28.71	17.72	17.65	24.49	18.16	23.61	15.35		
6	2.23	1.71	3.53	3.48	1.25	0.00	1.47	0.02	1.09	0.47		
7	0.35	0.30	0.52	0.43	0.00	0.00	0.00	0.11	0.07	0.01		
8	0.04	0.12	0.05	0.07	0.00	0.00	0.00	0.00	0.00	0.00		

Number of guanosines at the template-switching interface.

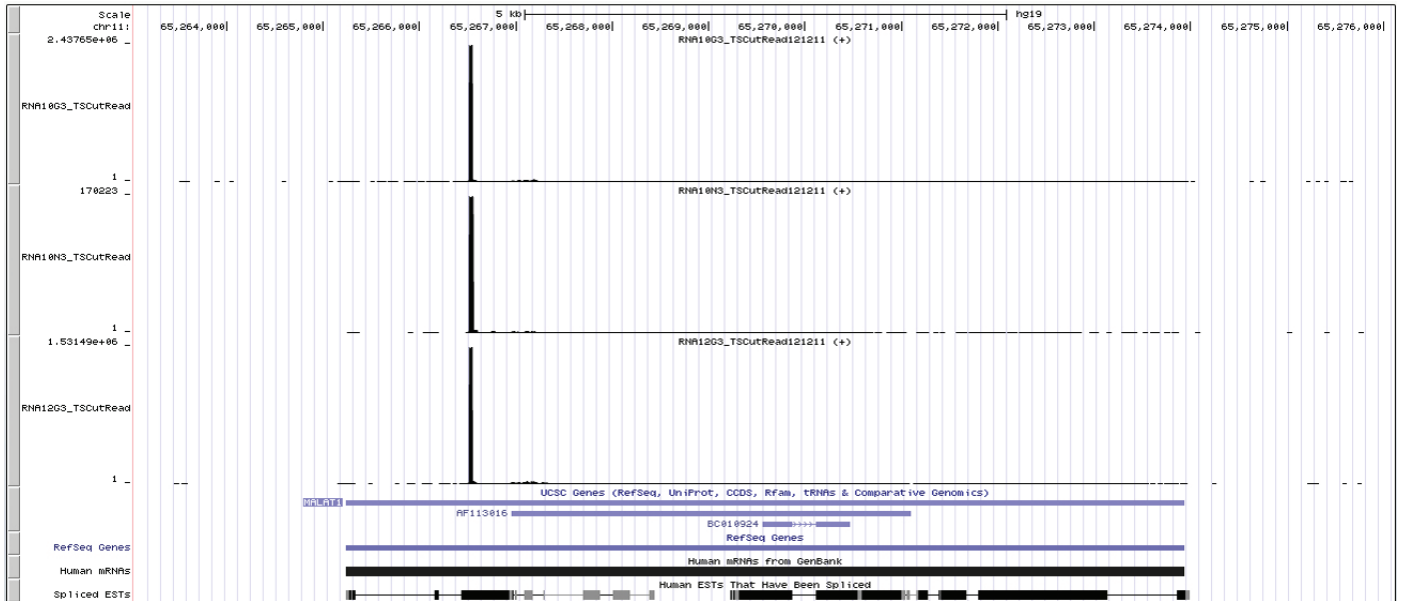
Figure S1



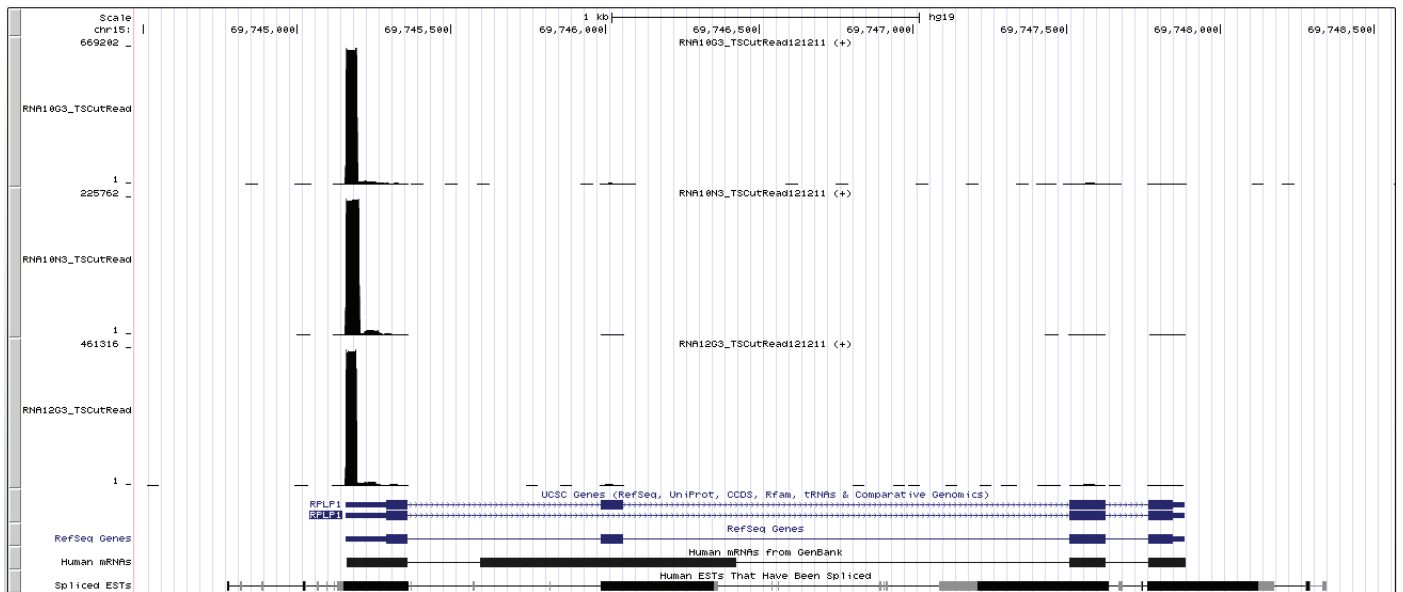
Regression analysis of TSO length data.

Figure S2

### MALAT1



### RPLP1



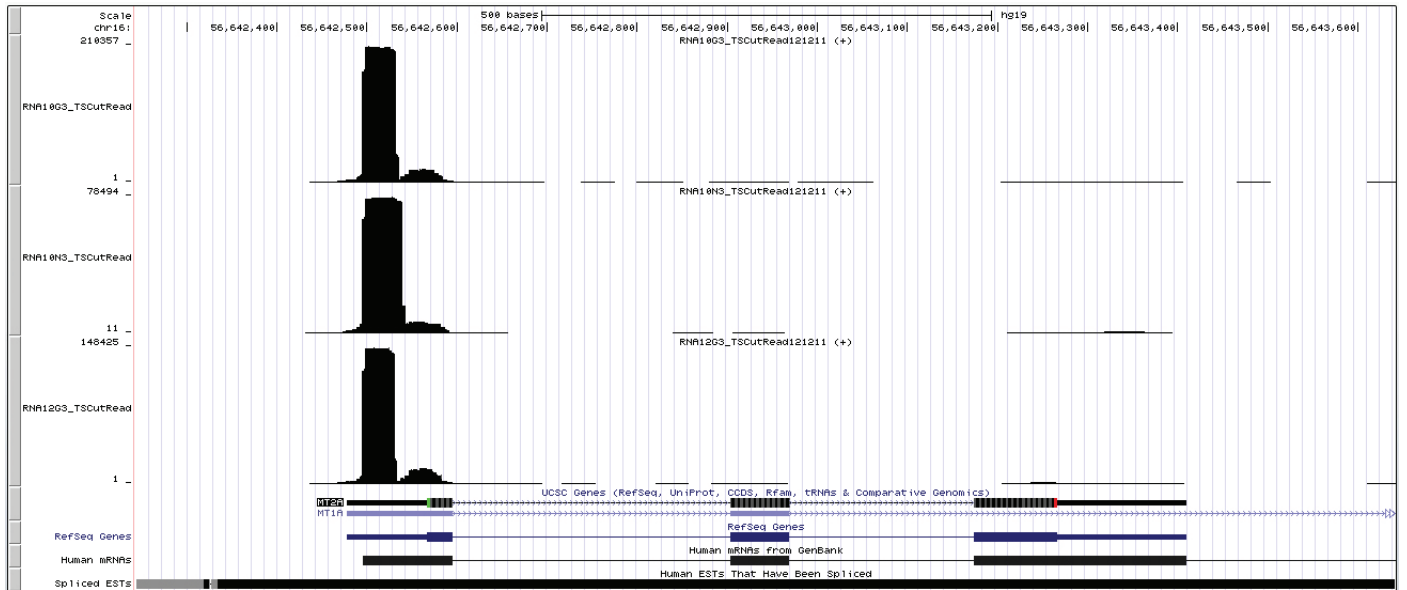
#### Hit distributions along the analyzed transcripts.

Images from Genome Browser are shown. The upper tracks correspond to the RNA10G3 reaction, the middle tracks to the RNA10N3 reaction and the bottom tracks to the RNA12G3 reaction. The hit distributions are highly similar.



Figure S2 (continued)

**MT2A**



**AHSG**

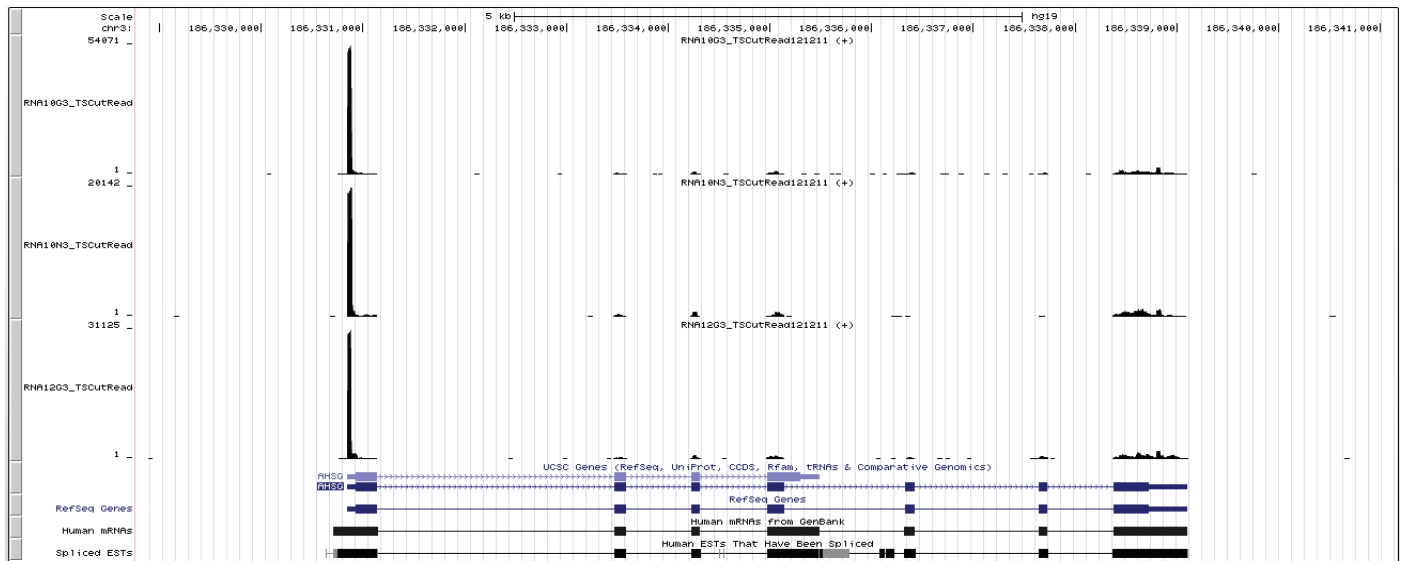


Figure S2 (continued)

### CNIH4

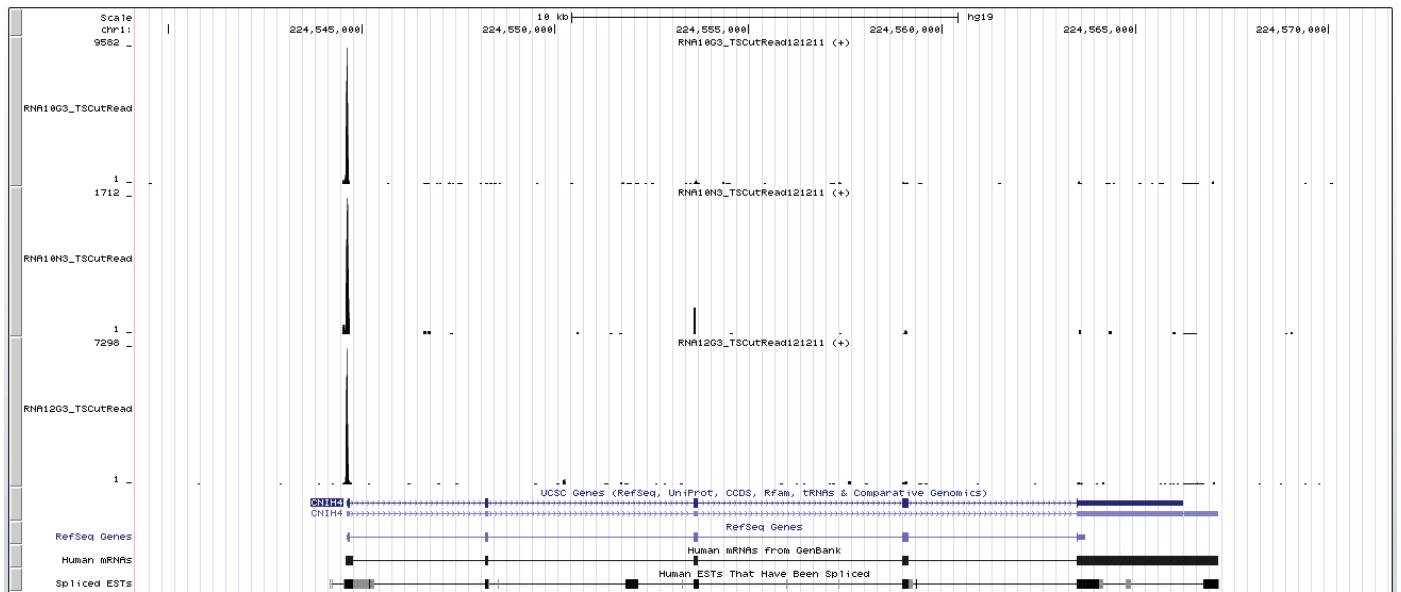
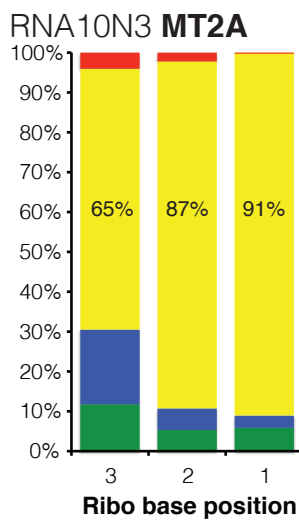
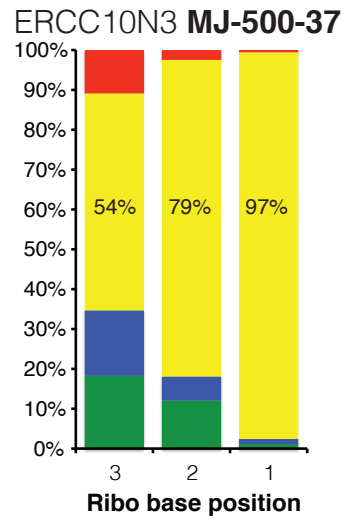
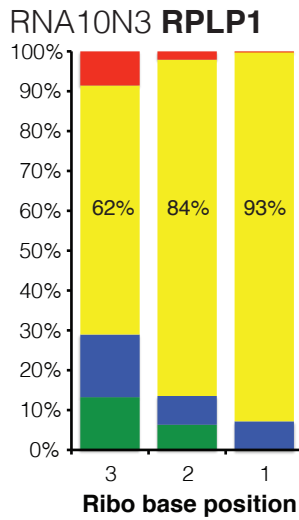
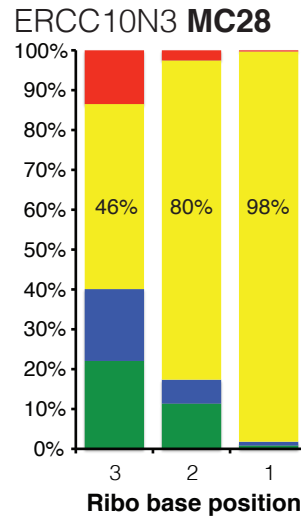
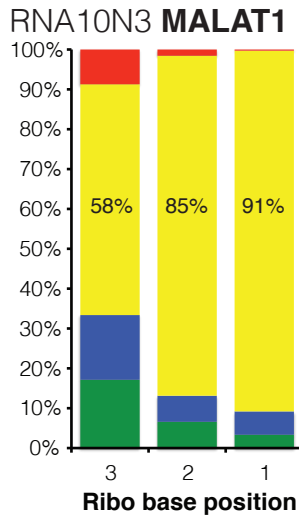


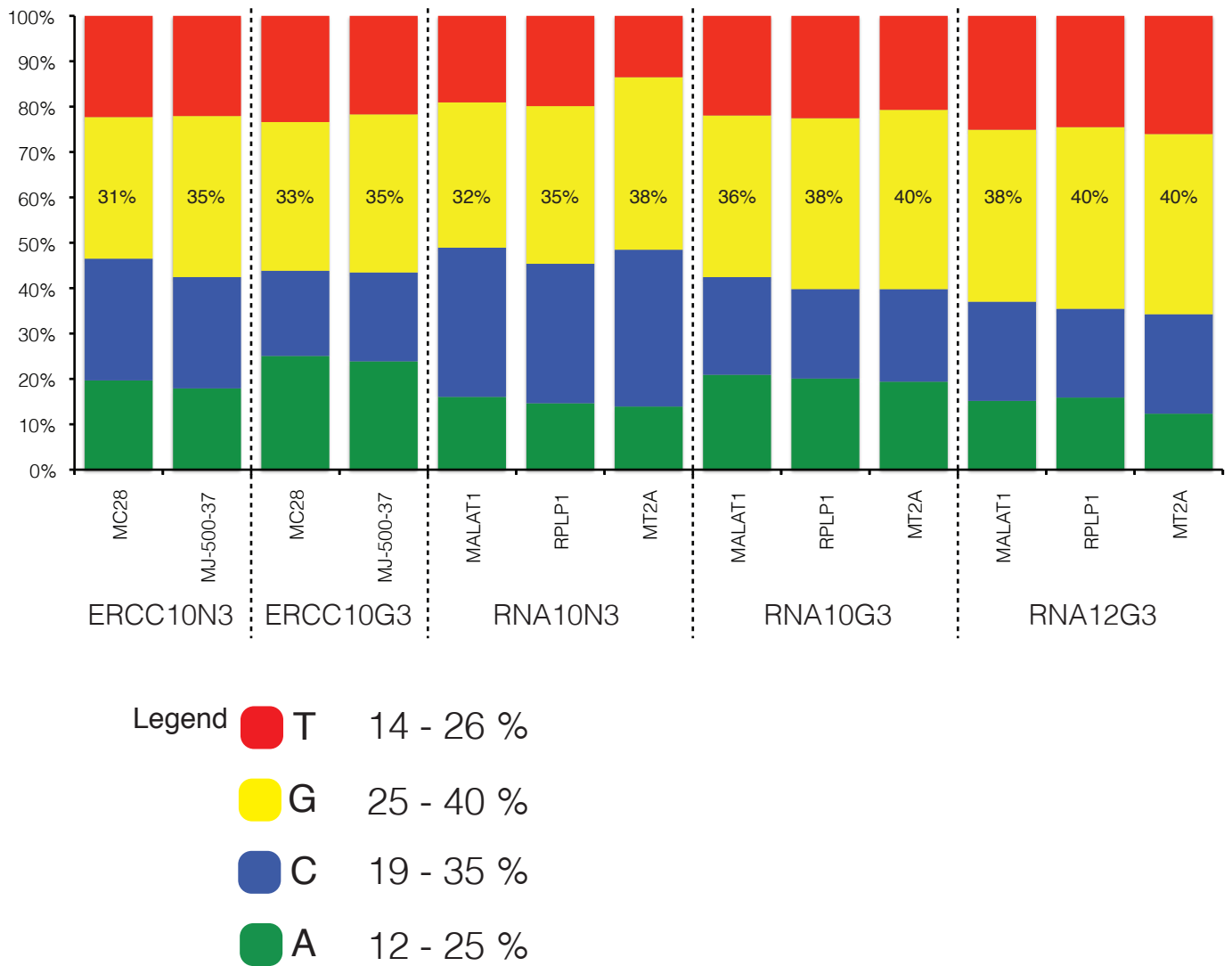
Figure S3



**Composition of the ribo base portion of the TSO.**

Position 1 represents the ribo base at the template-switching junction. In all cases, G is the preferred nucleotide in the three ribo positions of the TSO. The preference for G decreases as the distance from the template-switching site increases.

Figure S4

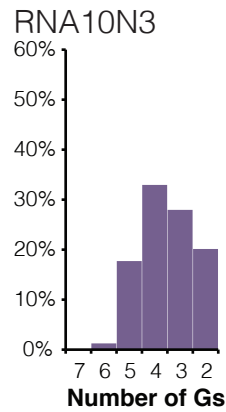
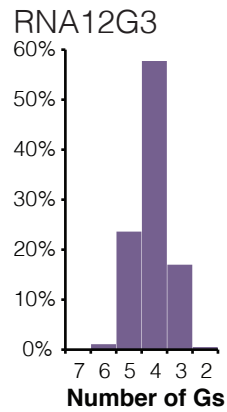
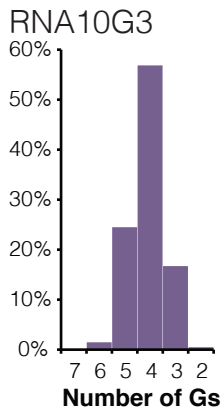


**Composition of DNA base in position 4, i.e. the DNA base adjacent to the ribo base stretch.**

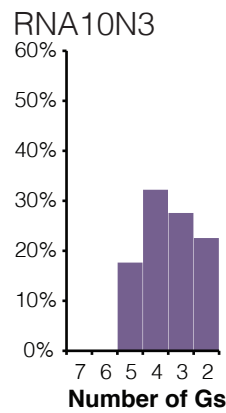
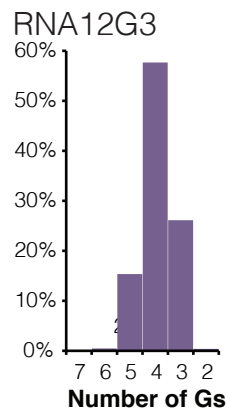
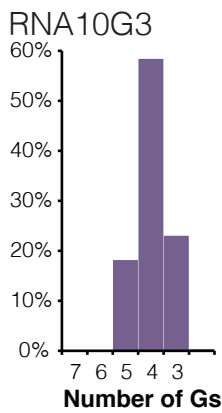
Position 4 is the first DNA position in the TSO as seen from the template-switching site. As such, it is located adjacent to the ribo base portion. G is the preferred nucleotide in this position.

Figure S5

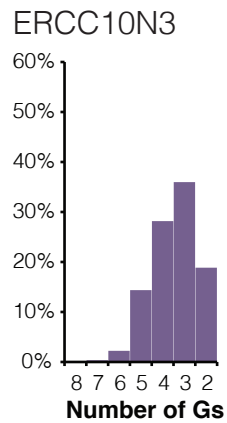
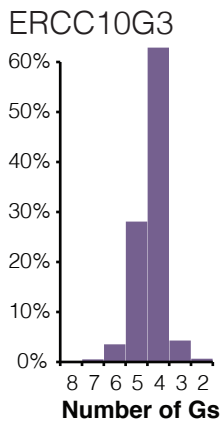
**RPLP1**



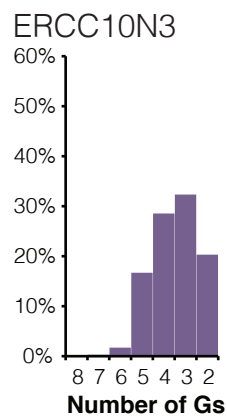
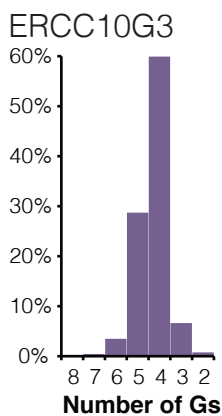
**MT2A**



**MC28**



**MJ-500-37**

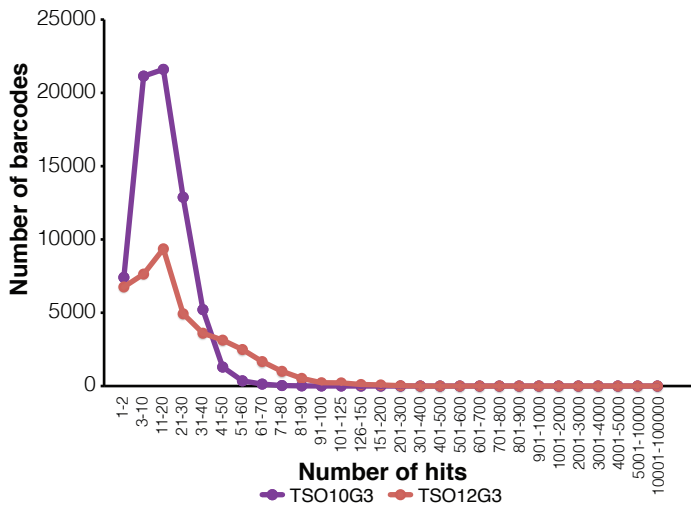


**Number of guanosines at the template-switching junction.**

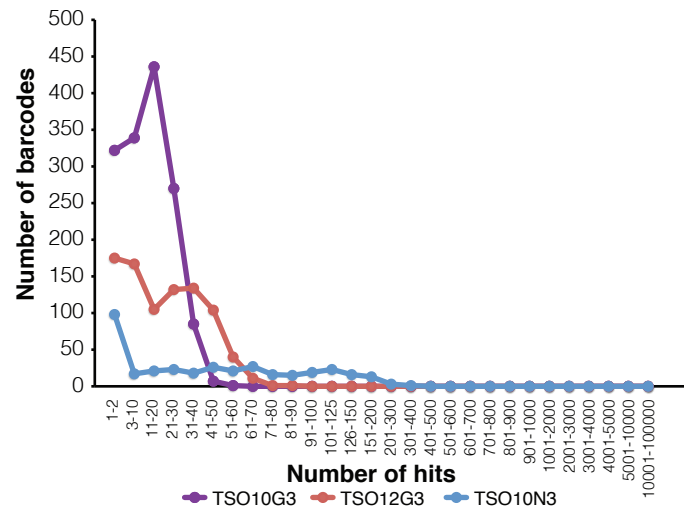
The graphs show the distribution of the number of guanosines observed in the sequencing data for the analyzed transcripts and RNA spike molecules. In general, three to four Gs are observed most frequently.

Figure S6

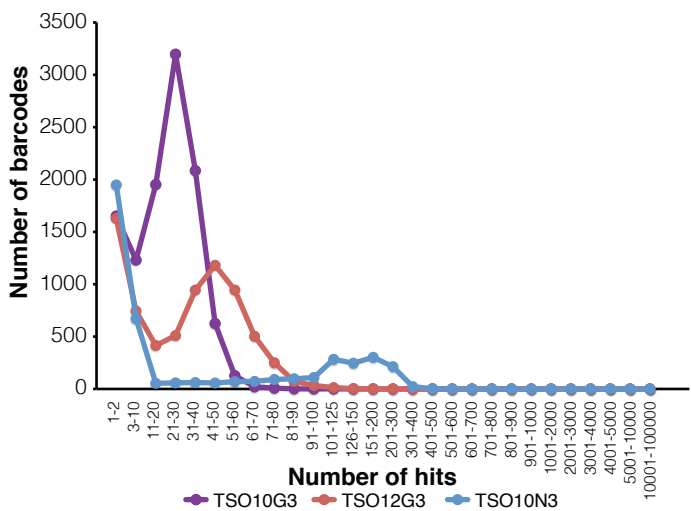
**MALAT1**



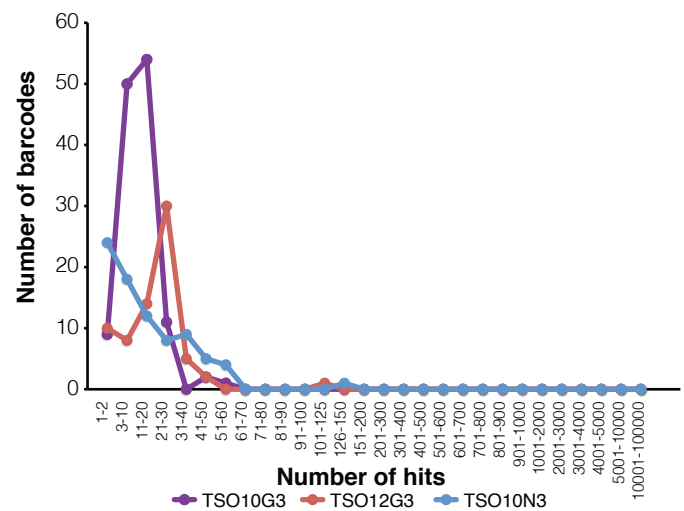
**AHSG**



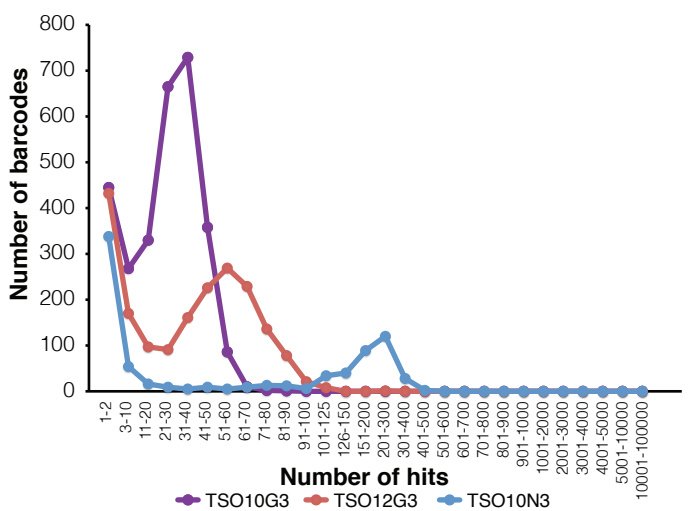
**RPLP1**



**CNIH4**



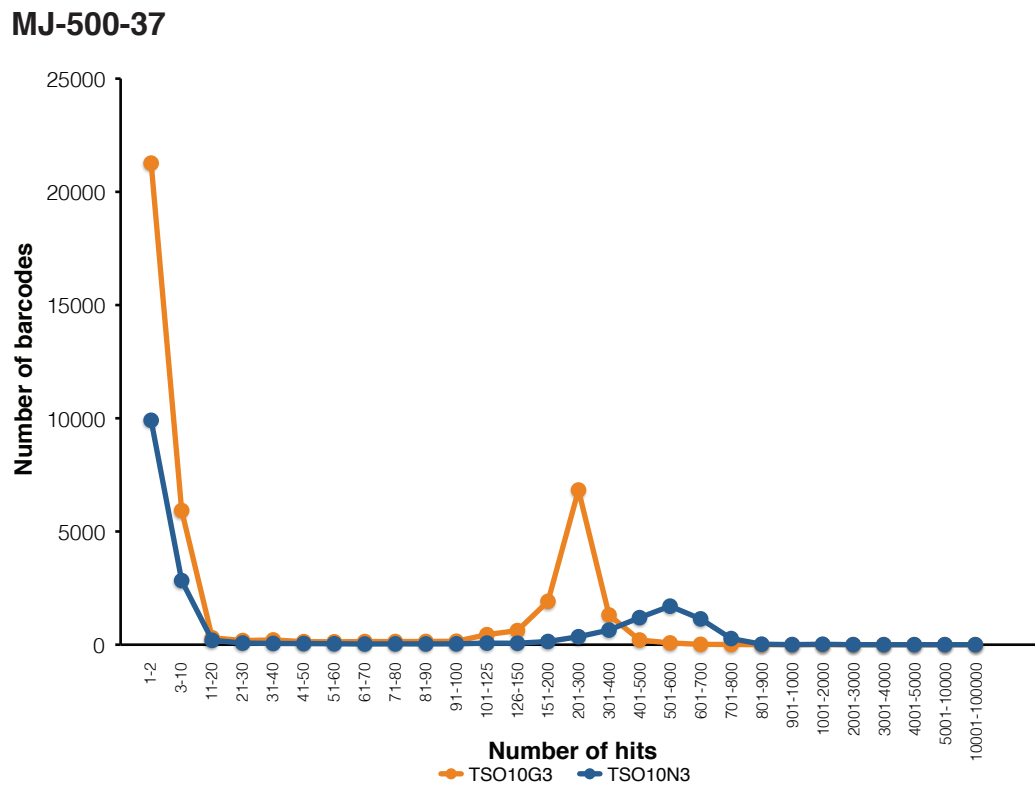
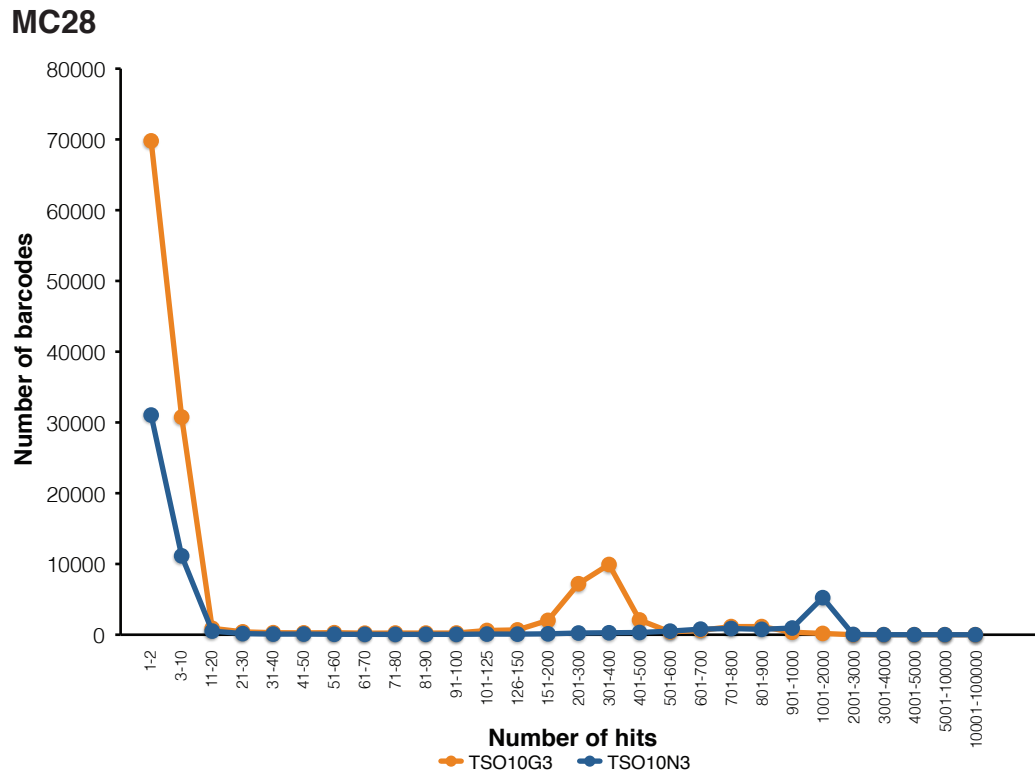
**MT2A**



**Barcode complexity for different transcripts.**

The graphs show the barcode distribution in the three reactions employing total RNA for the analyzed transcripts. The reaction with TSO10G3 exhibits the highest apparent complexity, followed by the reaction with TSO12G3 and with the TSO10N3 reaction showing the lowest apparent complexity.

Figure S7



#### Barcode complexity for different RNA spikes.

The graphs show the barcode distribution in the two reactions employing ERCC spikes for the analyzed spike molecules. The reaction with TSO10G3 exhibits higher apparent complexity than the reaction with TSO12G3.