

 Open access • Posted Content • DOI:10.1101/2020.09.28.316265

When DNA gets in the way in RNA-seq experiments, a sequel — [Source link](#)

Jasper Verwilt, María Dolores Giráldez, Wim Trypsteen, Ruben Van Paemel ...+3 more authors

Institutions: Ghent University

Published on: 29 Sep 2020 - bioRxiv (Cold Spring Harbor Laboratory)

Topics: DNA Contamination

Related papers:

- [When DNA gets in the way: A cautionary note for DNA contamination in extracellular RNA-seq studies.](#)
- [Accurate estimation of expression levels of homologous genes in RNA-seq experiments.](#)
- [Multi-perspective quality control of Illumina RNA sequencing data analysis.](#)
- [Elimination of PCR duplicates in RNA-seq and small RNA-seq using unique molecular identifiers](#)
- [Variant analysis pipeline for accurate detection of genomic variants from transcriptome sequencing data.](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/when-dna-gets-in-the-way-in-rna-seq-experiments-a-sequel-4iwb5re01b>

When DNA gets in the way in RNA-seq experiments, a sequel

Jasper Verwilt^{a,b,*,\$}, Maria D. Giraldez^{c,d,*,\$}, Wim Trypsteen^{a,b}, Ruben Van Paemel^{a,b},
Katleen De Preter^{a,b}, Pieter Mestdagh^{a,b}, Jo Vandesompele^{a,b}

Using a newly developed method dubbed SILVER-Seq—enabling extracellular RNA sequencing (exRNA-seq) directly from a small volume of human serum or plasma—Yan et al. recently reported in *Current Biology* a potential exRNA biomarker for the early diagnosis of Alzheimer’s disease [1]. After the publication of the initial paper describing the SILVER-Seq method [2], we reported our concern regarding potential DNA contamination in their datasets [3]. Although the authors replied they were able to successfully treat RNA samples with DNase to avoid such contamination, they did not address our observations of the majority of reads without evidence of being derived from RNA, nor documented verified absence of DNA after DNase treatment [4]. To assess whether the newly data generated may suffer from DNA contamination, we downloaded the publicly available sequencing data and evaluated two quality control metrics (i.e., fraction of exonic and splice reads), which were not reported in the paper. We found that both quality metrics were much lower than expected for RNA-seq data (6.28% exonic and 0.478% splice reads), in line with our previous findings on the first SILVER-Seq paper. These observations suggest the data and results presented by Yan et al. are affected by DNA contamination, an issue that may be inherent to the SILVER-Seq technology.

^aDepartment of Biomolecular Medicine, Ghent University, 9000 Ghent, Belgium;

^bOncoRNALab, Cancer Research Institute Ghent, 9000 Ghent, Belgium;

^cDigestive Diseases Unit, Virgen del Rocio University Hospital, 41013 Seville, Spain;

^dOncoDigest Group, Institute of Biomedicine of Seville (IBiS), 41013 Seville, Spain;

*Correspondence: jasper.verwilt@ugent.be; mdgiraldez-ibis@us.es

\$ shared 1st authors

24

25 RNA sequencing (RNA-seq) has transformed transcriptome characterization in a
26 wide range of biological contexts and is increasingly used to study samples with a
27 low RNA concentration, such as human biofluids. Biofluids contain microRNA and
28 other types of sncRNA, fragments of multiple RNA classes (e.g., mRNA, lncRNA,
29 tRNA, mtRNA) and circular RNA [5]. The presence of a variety of exRNA molecules
30 in the human bloodstream and other biofluids has opened up new avenues for the
31 development of minimally invasive biomarkers for a wide range of diseases.
32 However, the explosion in exRNA research has resulted in a growing field lacking
33 standardized protocols, consensus on data analysis, consistent findings and
34 sufficient experimental detail in many publications, which prevents researchers from
35 critically evaluating the quality of the presented results or reproducing the
36 experiments. Besides, performing exRNA-seq experiments without adequate quality
37 controls may result in several issues, one being sample contamination [6].

38

39 RNA-seq contaminants can be either external (originating from a different sample or
40 another species) or, although often overlooked, internal (originating from other
41 molecules from the same sample). Endogenous DNA contamination can be
42 particularly troubling as it can be hard to detect unless specific quality control
43 measurements are performed. RNA-seq experiments suffering from DNA
44 contamination can lead to biased results as it affects proper data quantification and
45 normalization. Due to the low concentration of RNA in human biofluids, DNA
46 contamination can be particularly vexing in exRNA-seq, preventing the reliable
47 detection of potential biomarkers.

48

49 DNase treatment is included in most standardized RNA-seq protocols but, in some
50 instances, it is not completely effective (not all DNA is removed) and can result in
51 impaired final libraries. This problem can be aggravated in protocols using crude
52 biofluids without RNA purification, which may contain DNase-inhibiting molecules
53 (one of them being actin, which has long been known for inhibiting DNase activity
54 [7]). Serum in particular contains not only cell-free DNA but also genomic DNA that
55 originates from lysis of white blood cells during *ex vivo* clotting [8], thus increasing
56 the risk of DNA contamination in RNA-seq experiments.

57

58 To evaluate whether the RNA-seq signal in the paper by Yan et al. [1] might be
59 affected by contaminating DNA, we replicated the pipeline used in the paper as
60 accurately as possible (no details were reported regarding parameters of sequencing
61 and data analysis) and calculated several quality control metrics. A step-by-step
62 overview of the used tools can be found in the Supplemental Methods section and
63 the full code is uploaded to GitHub
64 (https://github.com/jasperverwilt/exRNA_contamination).

65

66 In order to confirm, or refute, our suspicions, we were mainly interested in two data
67 quality metrics: the fraction of exonic reads (5% in case of sequencing pure DNA)
68 and the fraction of splice reads (0% in case of DNA). Considering all samples, we
69 observe exonic fractions ranging from 4.7% to 25.4%, with a median value of 6.28%
70 (Figure 1A); and splice fractions ranging from 0.206% to 1.27%, with a median value
71 of 0.478% (Figure 1B). In addition to the splice and exonic fractions, we checked the
72 strandedness of the data. If SILVER-Seq would employ a stranded library
73 preparation approach (which we do not know for sure, as it is unreported), and the

74 data turns out to be unstranded, contaminating DNA might be at play (since DNA is
75 double stranded, the reads can originate from both strands). With a strandedness of
76 100% for perfectly stranded data, and 50% for pure DNA, the observed median
77 strandedness was 49.2%, with individual values ranging from 47.9 to 70.4%
78 (Supplemental Figure 1). These results support the hypothesis of DNA
79 contamination.

80

81 The low median value of exonic and spliced reads prompts us to conclude that most
82 of the SILVER-Seq data generated by Yan et al. is affected by DNA contamination.
83 We deduct that SILVER-Seq is a stranded library prep method, given that some
84 samples showed a strandedness higher than 50%. The wide ranges of the exonic
85 and splice read fractions and variable strandedness level indicate that DNA is
86 differentially present in the samples, with some samples performing consistently
87 worse or better for all quality metrics: SRR10015490, for example, showed relatively
88 high values for all the metrics (Figure 1, Supplemental Figure 1).

89

90 Finally, the biogenesis of exRNA is not well established yet and some authors argue
91 that biofluids might be enriched in intron and antisense sequences compared with
92 cellular RNAs [9]. However, we are concerned that DNA contamination is the most
93 likely explanation here as: (a) some exRNA-seq studies have consistently reported a
94 high proportion of exonic reads and adequate strandedness [10]; (b) the inherent
95 challenge of avoiding DNA contamination, especially when working with crude
96 biosamples as input; and, (c) the high variability of the evaluated quality control
97 metrics across the reported samples. We would like to emphasize that our
98 observations do not undermine the potential utility of SILVER-Seq. Our letter is a call

99 for thorough reporting of methodology and analysis details including quality control
100 metrics in exRNA-seq studies. We hope that our plea helps to move the exRNA field
101 forward by promoting consistency among laboratories and increasing experimental
102 transparency and reproducibility.
103

104 **References**

105

- 106 1. Yan, Z., Zhou, Z., Wu, Q., Chen, Z.B., Koo, E.H., and Zhong, S. (2020).
107 Presymptomatic Increase of an Extracellular RNA in Blood Plasma Associates
108 with the Development of Alzheimer's Disease. *Curr. Biol.* *30*, 1771-1782.e3.
109 Available at: <https://doi.org/10.1016/j.cub.2020.02.084> [Accessed September
110 16, 2020].
- 111 2. Zhou, Z., Wu, Q., Yan, Z., Zheng, H., Chen, C.-J., Liu, Y., Qi, Z., Calandrelli,
112 R., Chen, Z., Chien, S., *et al.* (2019). Extracellular RNA in a single droplet of
113 human serum reflects physiologic and disease states. *Proc. Natl. Acad. Sci.*
114 *116*, 19200–19208.
- 115 3. Verwilt, J., Trypsteen, W., Van Paemel, R., De Preter, K., Giraldez, M.D.,
116 Mestdagh, P., and Vandesompele, J. (2020). When DNA gets in the way: A
117 cautionary note for DNA contamination in extracellular RNA-seq studies. *Proc.*
118 *Natl. Acad. Sci.* *117*, 18934–18936. Available at:
119 <http://www.pnas.org/lookup/doi/10.1073/pnas.2001675117> [Accessed August
120 26, 2020].
- 121 4. Zhou, Z., Wu, Q., Yan, Z., Zheng, H., Chen, C.J., Liu, Y., Qi, Z., Calandrelli, R.,
122 Chen, Z., Chien, S., *et al.* (2020). Reply to Verwilt et al.: Experimental evidence
123 against DNA contamination in SILVER-seq. *Proc. Natl. Acad. Sci. U. S. A.* *117*,
124 18937–18938. Available at: www.pnas.org/cgi/doi/10.1073/pnas.2008585117
125 [Accessed August 26, 2020].
- 126 5. Hulstaert, E., Morlion, A., Cobos, F.A., Verniers, K., Nuytens, J., Eynde, E.
127 Vanden, Yigit, N., Anckaert, J., Geerts, A., Hindryckx, P., *et al.* (2019). Charting
128 extracellular transcriptomes in The Human Biofluid RNA Atlas. *bioRxiv*,

- 129 823369.
- 130 6. Nieuwenhuis, T.O., Yang, S.Y., Verma, R.X., Pillalamarri, V., Arking, D.E.,
131 Rosenberg, A.Z., McCall, M.N., and Halushka, M.K. (2020). Consistent RNA
132 sequencing contamination in GTEx and other data sets. *Nat. Commun.* *11*, 1–
133 10. Available at: <https://doi.org/10.1038/s41467-020-15821-9> [Accessed
134 September 17, 2020].
- 135 7. Blikstad, I., Markey, F., Carlsson, L., Persson, T., and Lindberg, U. (1978).
136 Selective assay of monomeric and filamentous actin in cell extracts, using
137 inhibition of deoxyribonuclease I. *Cell* *15*, 935–943. Available at:
138 <http://www.cell.com/article/0092867478902775/fulltext> [Accessed September
139 17, 2020].
- 140 8. Lee, T.H., Montalvo, L., Chrebtow, V., and Busch, M.P. (2001). Quantitation of
141 genomic DNA in plasma and serum samples: Higher concentrations of
142 genomic DNA found in serum than in plasma. *Transfusion* *41*, 276–282.
143 Available at: [https://onlinelibrary.wiley.com/doi/full/10.1046/j.1537-
144 2995.2001.41020276.x](https://onlinelibrary.wiley.com/doi/full/10.1046/j.1537-2995.2001.41020276.x) [Accessed September 17, 2020].
- 145 9. Qin, Y., Yao, J., Wu, D.C., Nottingham, R.M., Mohr, S., Hunicke-Smith, S., and
146 Lambowitz, A.M. (2016). High-throughput sequencing of human plasma RNA
147 by using thermostable group II intron reverse transcriptases. *RNA* *22*, 111–
148 128. Available at: <http://www.rnajournal.org/cgi/doi/10.1261/rna.054809.115>.
149 [Accessed September 16, 2020].
- 150 10. Everaert, C., Helsmoortel, H., Decock, A., Hulstaert, E., Van Paemel, R.,
151 Verniers, K., Nuytens, J., Anckaert, J., Nijs, N., Tulkens, J., *et al.* (2019).
152 Performance assessment of total RNA sequencing of human biofluids and
153 extracellular vesicles. *Sci. Rep.* *9*, 17574. Available at:

154 <http://www.nature.com/articles/s41598-019-53892-x> [Accessed January 20,
155 2020].

156

157

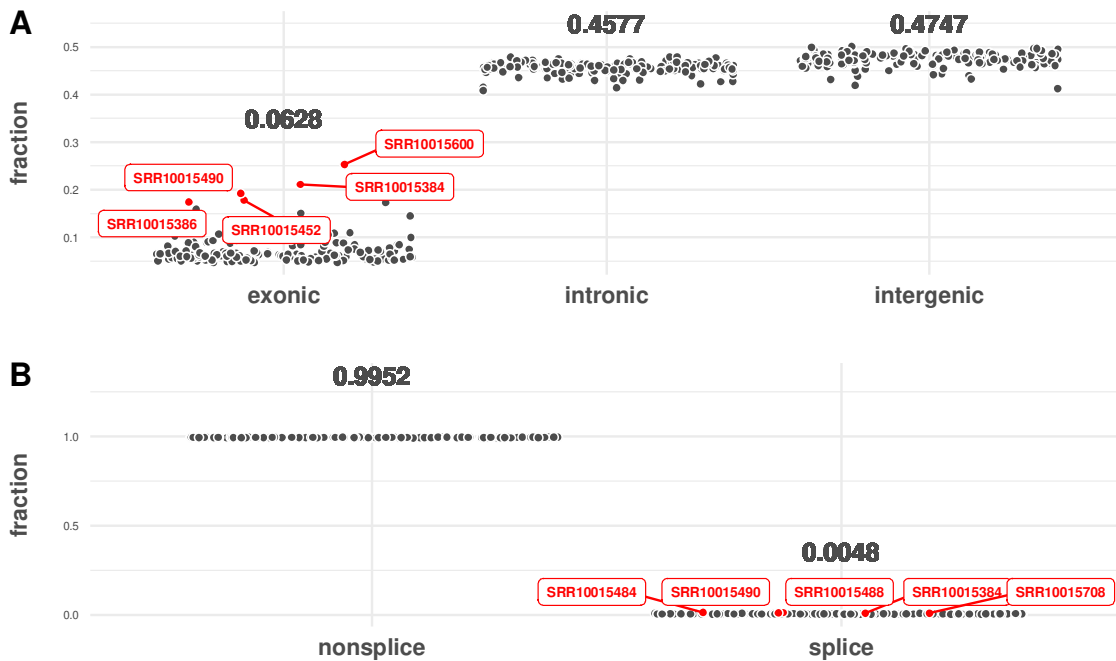
158

159 **Figures**

160

161 **Figure 1**

162



163

164

165 **Figure legends**

166

167 **Figure 1:** Regional coverage and splice read fractions of the data. (A) Fractions of
168 reads mapping to exonic, intronic and intergenic regions. The data points are
169 calculated values for individual samples. The median fractions over all samples are
170 printed. The five samples with the highest exonic coverage are annotated and
171 colored red. (B) Fractions of reads mapping to splice and nonsplice regions. The data
172 points are calculated values for individual samples. The median fractions over all

173 samples are printed. The five samples with the highest fraction of reads mapping to

174 splice junctions are annotated and colored red.

175

176