1    **Tittle:**

2    **A technical assessment of the porcine ejaculated spermatozoa for a sperm specific RNA-**

3    **seq analysis**

4

5    **Authors and affiliations:**

6    Marta Gòdia[a,*]; email: marta.godia@cragenomica.es ORCID: 0000-0002-0439-4014

7    Fabiana Quoos Mayer[a,b,*]; email: bimmayer@gmail.com ORCID: 0000-0002-9324-8536

8    Julieta Nafissi[a,c]; email: julietanafissi@gmail.com

9    Anna Castelló[a,d]; email: anna.castello@uab.cat

10   Joan Enric Rodríguez-Gil[e]; email: juanenrique.rodriguez@uab.cat

11   Armand Sánchez[a,d]; email: armand.sanchez@cragenomica.es ORCID: 0000-0001-9160-1124

12   Alex Clop[a,#]; email: alex.clop@cragenomica.es ORCID: 0000-0001-9238-2728

13

14

15   [a]Animal Genomics Group, Centre for Research in Agricultural Genomics-CSIC-IRTA-UAB-

16   UB, Campus UAB, 08193 Cerdanyola del Valles, Catalonia, Spain. Tel: +34 93 5636600

17   [b]Permanent adress - Instituto de Pesquisas Veterinárias Desidério Finamor, Agricultural

18   Diagnostic and Research Departament, Secretariat of Agriculture, Livestock and Irrigation,

19   Eldorado do Sul, Rio Grande do Sul, Brazil

20   [c]Current address: Technology Institute (INTEC), Argentine University of Enterprise

21   (UADE), Buenos Aires, Argentina

22   [d]Unit of Animal Science, Department of Animal Science and Nutrition, Autonomous

23   University of Barcelona, Cerdanyola del Valles, Catalonia, Spain

24 ᵉUnit of Animal Reproduction, Department of Animal Medicine and Surgery, Autonomous

25 University of Barcelona, Cerdanyola del Valles, Catalonia, Spain

26 * These authors contributed equally to this work.

27 #Corresponding author.

30

31 **Abbreviations**

32 FPKM - Fragments Per Kilobase of transcript per Million mapped reads

33 KRT1 - Keratin 1

34 miRNA - micro RNA

35 miscRNA - miscellaneous RNA

36 Mt rRNA - mitochondrial ribosomal RNA

37 Mt tRNA - mitochondrial transference RNA

38 OAZ3 - Ornithine Decarboxylase Antizyme 3

39 ORT - Osmotic Resistance Test

40 piRNA - Piwi-interacting RNA

41 PRM1 - Protamine 1

42 PTPRC - Protein tyrosine phosphatase receptor type C

43 rRNA - ribosomal RNA.

44 snoRNA - small nucleolar RNA

45    snRNA - small nuclear RNA

46    SRR – Sperm Recovery Rate

47    tRNA - transfer RNA

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

**Abstract**

70

71 The study of the boar sperm transcriptome by RNA-seq can provide relevant information on

72 sperm quality and fertility and might contribute to animal breeding strategies. However, the

73 analysis of the spermatozoa RNA is challenging as these cells harbor very low amounts of

74 highly fragmented RNA, and ejaculates also contain other cell types with larger amounts of

75 non-fragmented RNA. Here we describe a strategy for a successful boar sperm purification,

76 RNA extraction and RNA-seq library preparation. Using these approaches our objectives

77 were: (i) to evaluate the sperm recovery rate (SRR) after boar spermatozoa purification by

78 density centrifugation using the non-porcine-specific commercial reagent BoviPure$^{TM}$; (ii) to

79 assess the correlation between SRR and sperm quality characters; (iii) to evaluate the

80 relationship between sperm cell RNA load and sperm quality traits and (iv) to compare

81 different library preparation kits for both total RNA-seq (SMARTer Universal Low Input

82 RNA and TruSeq RNA Library Prep kit) and small RNA-seq (NEBNext Small RNA and

83 TailorMix miRNA Sample Prep v2) for high throughput sequencing. Our results show that

84 pig SRR (~ 22 %) is lower than in other mammalian species and that it is not significantly

85 dependent of the sperm quality parameters analyzed in our study. Moreover, no relationship

86 between the RNA yield per sperm cell and sperm phenotypes was found. We compared a

87 RNA-seq library preparation kit optimized for low amounts of fragmented RNA with a

88 standard kit designed for high amount and quality of input RNA and found that for sperm, a

89 protocol designed to work on low quality RNA is essential. We also compared two small

90 RNA-seq kits and found not substantial differences in their performance. We propose the

91 methodological workflow described here for the RNA-seq screening of the boar spermatozoa

92 transcriptome.

93

## Introduction

RNA-seq is the current gold-standard technology for the high-throughput analysis of transcriptome profiles, which is essential to understand the molecular basis of phenotypes (Wang et al. 2009). Thus, if studied in livestock species, this information could contribute to designing animal breeding strategies. This method has been applied to map the transcriptome of multiple species and tissues including spermatozoa from human (Sendler et al. 2013), mouse (Fang et al. 2014), bovine (Card et al. 2013; Selvaraju et al. 2017) and horse (Das et al. 2013). Although these cells are considered transcriptionally and translationally inactive, they contain a wide population of coding and non-coding RNA molecules (Jodar et al. 2013), which functions have been related to spermatogenesis (Ostermeier et al. 2002), sperm chromatin reorganization (Martins and Krawetz 2005; Hamatani 2012), fertility potential (Jodar et al. 2015), early embryo development (Sendler et al. 2013) and trans-generational epigenetic inheritance (Rando 2016). Hence, the study of the sperm transcriptome is crucial for understanding its biology and its role in fertility, and can be thus of interest when applied to livestock research.

One of the main challenges for the study of the spermatozoa transcriptome is the extremely low RNA yield and high fragmentation of the transcripts typically present in these cells, as the standard RNA-seq chemistry normally requires a large amount (1 µg) of good quality RNA. To overcome this problem, new protocols to prepare high quality RNA-seq libraries from samples containing only tiny amounts (200 pg) of highly degraded RNA (e.g., paraffin embedded tissues) have been developed and already tested and compared in human sperm (Mao et al. 2014). A human mature sperm cell is estimated to contain a 600-fold lower amount of RNA than a somatic cell (Zhao et al. 2006). As a typical ejaculate contains somatic cells – mainly leukocytes, keratinocytes and other type of epithelial cells – as well as germ-line cells from different stages of spermatogenesis (Patil et al. 2013), the study of the spermatozoa transcriptome requires removing these RNA-rich cells for an unbiased analysis.

Somatic cells can be removed from sperm by the swim-up method (Jameel 2008), somatic cell lysis or by gradient centrifugation (Mao et al. 2013). Cell lysis approaches are efficient in eliminating somatic cells, but they also cause cell membrane damage and loss of mitochondrial sequences, thus risking to lose the sperm transcripts present in the cell's midpiece (Mao et al. 2013). Gradient centrifugation has been employed in the purification of sperm cells from several mammalian species using different commercial solutions, such as Percoll® (Ostermeier et al. 2002), PureSperm® (Sendler et al. 2013), EquiPure™ (Das et al. 2013) and BoviPure™ (Samardzija et al. 2006; Selvaraju et al. 2017). These gradients allow the motile mature spermatozoa to separate from somatic cells along with immature sperm cells (Mao et al. 2013). Typically, these commercial reagents are primarily used to improve sperm quality for artificial insemination, since they select progressive motile and morphologically normal spermatozoa (Samardzija et al. 2006). Although gradient centrifugation is convenient for these purposes, it significantly decreases the final number of recovered spermatozoa (Samardzija et al. 2006), adding yet another layer of complexity for the experimental analysis of the spermatozoa transcriptome. The sperm recovery rate (SRR) in gradient-based methods is mainly related to sperm motility (Samardzija et al. 2006), even though additional factors are likely to be involved since the number of recovered cells is lower than the expected based solely on initial motility values. The boar sperm is particularly sensitive to a wide spectrum of manipulations (Feugang 2017) and the use of a reagent not optimized for the porcine sperm may have detrimental effects on the SRR. Taking all this information into account, one of the main aims of this work was to evaluate the influence of different boar sperm quality traits on SRR after gradient density purification.

The levels of several RNA transcripts in sperm have been associated to semen quality traits and male fertility in many mammalian species including human (Jodar et al. 2012), cattle (Bissonnette et al. 2009) and pigs (Curry et al. 2011), among others. Likewise, abnormal levels of histone or protamine chromatin proteins in sperm have also been linked to spermatozoa defects (Carrell et al. 2007;

Hammoud et al. 2011) and it has been suggested that it could be related to spermatogenesis defects and alterations in RNA amounts in sperm cells (Carrell et al. 2007) and even sperm quality (Aoki et al. 2005). Thus, we searched for statistical relationship between RNA yield extracted per sperm cell and semen quality traits in swine. To determine the purity (lack of DNA or somatic RNA) of this sperm RNA, we developed three real-time quantitative PCRs (qPCRs) and tested its efficiency. Finally, we performed high throughput sequencing of a selection of these RNAs using two library preparation kits for total RNA-seq analysis ($N = 6$) and two kits for small RNA-seq analysis ($N = 3$) to compare their performances.

## Results

### *Spermatozoa Recovery Rate*

SRR was calculated in 285 samples and was in general low and with high variability between samples. The average SRR was 21.76 % with a standard deviation of 15.07 %. To shed light into the biological causes of this variance, we tested the dependence between SRR and sperm phenotypes. Significant covariates were adjusted for the given parameters: head abnormalities, tail abnormalities and distal droplets were adjusted for farm; motility 90 min for age; viability 0 min, viability 90 min, acrosomes 0 min and ORT for batch; acrosomes 90 min and neck abnormalities for farm and batch, and distal droplets for farm, age and batch. SRR and the sperm quality characters did not present normal distribution nor a linear relationship. Thus, a multivariate nonparametric test of independence was applied (Székely and Rizzo 2009). When considering the Bonferroni corrected *P*-value, SRR was found independent of all the sperm quality parameters (Table 1).

[Table 1 near here].

### *RNA yield*

Total RNA was extracted from 190 samples. The RNA yields averaged 1.6 fg per sperm cell, with ranges from 0.4 to 4.8 fg. The RNA Integrity Number (RIN) values, measured on 70 samples, was low (RIN < 2.6) and with undetectable ribosomal RNA profiles, which indicates the absence of RNA of somatic cell origin. The amount of RNA extracted per sperm cell was not significantly associated to the covariates farm, age or batch. The test of independence indicated null relationship between the total RNA extracted per sperm cell and the sperm quality phenotypes studied (Table 2).

[Table 2 near here].

### *qPCR controls*

The standard curves of the qPCR assays showed a good efficiency (97-97.9%). The three qPCR control assays displayed single peaks after the dissociation curve analysis, thus confirming that a single amplicon was generated in each reaction. The minus reverse transcription controls showed no

amplification of *PRM1* and *PTPRC*. 70 RNA samples were subjected to qPCR, and all presented quantification cycles (Cq) ranging between 14.6 and 21.3 for the sperm-specific gene *PRM1*. In contrast, the average Cq for *PTPRC* was 35.4 in 49 sperm samples and undetectable in the other 21 samples. The $\Delta$Cq $_{PTPRC-PRM1,}$ calculated as the Cq for *PTPRC* minus the Cq for *PRM1* in the sperm samples, ranged from 14.3 to 21.3. The intergenic region was undetectable in 66 samples and had Cqs > 36 in the other 4 and the $\Delta$Cq $_{Genomic-PRM1}$ ranged from 18.4 to 21. As a comparison, the liver RNA showed a *PRM1* and *PTPRC* Cqs of 38 and 24, respectively.

### *RNA-seq library preparation, sequencing and mapping statistics*

Four of the six samples that were chosen for total RNA-seq analysis (Sample_1 to Sample_6) presented $\Delta$Cq $_{PTPRC-PRM1}$ ranging from 17.4 to 19.1 and undetectable levels of *PTPRC* in the other two samples. Likewise, the $\Delta$Cq $_{Genomic-PRM1}$ ranged from 19.4 to 21 in three of the six samples and was undetectable in the other three.

The SMARTer and the TruSeq kits produced libraries with significantly different concentrations, which ranged between 53 and 120.7 nM (total RNA yield between 0.8 and 1.8 pmol) and between 0.5 and 2.9 nM (0.01 – 0.09 pmol), respectively (*P*-value = 0.03) (Table 3). All the libraries generated a similar percentage of high quality RNA-seq reads which mapped unambiguously to the swine reference genome, (SMARTer : 74.9 – 85.8 % and TruSeq: 70.8 – 82.3 %) (*P*-value = 0.13) (Table 3). Likewise, SMARTer yielded a higher percentage of reads uniquely mapping to annotated genes (37.8 - 48.4 %) when compared to the TruSeq libraries (28.8 - 38.5 %) (*P*-value = 0.02) (Table 3). The proportion of PCR duplicates was significantly higher for the TruSeq (89.3 - 97.9%) than for the SMARTer samples (75.9 - 80.3 %) (*P*-value = 0.03) (Table 3). We identified on average, 8,562 and 2,522 transcripts for the SMARTer and the TruSeq, respectively (*P*-value = $1.89 \times 10^{-4}$). The SMARTer datasets presented a mean FPKM of 363 and median FPKM of 4.8, whereas the TruSeq libraries showed a mean FPKM of 3,410 and median FPKM of 12.6. 32.5% and 46.1% of the genes

215  were identified at intermediate or high abundance levels (FPKM ≥ 10) for both the SMARTer and

216  the TruSeq, respectively.

217  [Table 3 near here].

218  Short RNA samples (Sample_7 to Sample_9) presented undetectable RNA levels for *PTPRC* and for

219  the intergenic region with the qPCR assay. Sequencing and mapping of short RNAs with the

220  NEBNext and the TailorMix displayed similar results. The proportion of reads mapping to annotated

221  features was similar for both protocols (77.4 - 82.9 %) (Table 4). Most of these reads mapped to

222  miRNAs (27.0 – 32.4 %) (Table 4), followed by mitochondrial tRNAs (22.1 – 27.3 %) and protein

223  coding genes (12.6 – 15.5 %) (Table 4). The remaining mapped reads corresponded to snRNAs,

224  piRNAs and tRNAs, among others (Table 4). Some of the most abundant miRNAs have been already

225  identified in swine sperm or in other mammalian species and include miR10b and miR34c, among

226  others (Capra et al. 2017; Chen et al. 2017; Jodar et al. 2013).

227  [Table 4 near here].

228  Further analysis of the transcriptome profile was carried with the SMARTer datasets using the

229  totality of the reads generated in each library (between 18.5 and 26.9 million reads per sample).

230  Genes related to somatic cell contamination, *PTPRC* and *KRT1* were absent (mean FPKM = 0.3 and

231  0, respectively) in these samples (Supplementary Table 1). On the contrary, the sperm specific *PRM1*

232  and *OAZ3* were among the most abundant transcripts with mean FPKMs of 15,368 and 22,670,

233  respectively (Supplementary Table 1). The pattern of relative expression of these four genes in

234  porcine white blood cells and in ear tissue was inverted when compared to sperm. Whilst *PRM1* and

235  *OAZ3* were absent, *PTPRC* and *KRT1* were abundant in the white blood cells and in the ear RNA-seq

236  datasets, respectively (Figure 1). We also quantified the amount of other previously reported somatic

237  and sperm specific gene biomarkers (Jodar et al. 2016). The abundance of the epithelial *CDH1*,

238  keratinocyte *KRT10*, leukocyte *IL8*, whole blood *HBB* and prostate *KLK3* genes ranged between 0

239　and 9 FPKM. On the contrary, the sperm-specific genes *PRM2, TNP1, ODF1* and *SMCP* showed

240　average FPKMs ranging between 779 and 7,742 (Supplementary Table 1).

241　[Figure 1 near here].

242

243　**Discussion**

244　Although spermatozoa are considered transcriptionally inactive, there is growing evidence that the

245　sperm RNA populations are related to spermatogenesis, fertility potential, chromatin reorganization,

246　embryo development and transgenerational epigenetic inheritance (Bohacek and Mansuy 2015; Jodar

247　et al. 2013). Since RNA load in spermatozoa is considerably lower than in somatic cells, an adequate

248　separation of these populations is imperative to study the spermatozoa transcriptome. The application

249　of purification methods decreases the final number of recovered sperm cells, and consequently the

250　cell availability for RNA extraction.

251　The present study is the first to analyze the performance of porcine SRR. The purification of the boar

252　sperm using gradient centrifugation with the non-porcine-specific reagent BoviPure$^{TM}$ yielded highly

253　purified spermatozoa as demonstrated by qPCR ($\Delta$Cq $_{PTPRC\text{-}PRM1}$ > 16) for the vast majority (97 %) of

254　the 70 samples. Nonetheless, the SRR was not only much lower but also more variable (21.76 ±

255　15.07 %) than that described in other species such as cattle (mean SRR = 31 %), human (69 %) and

256　horse (63 %) (Allamaneni et al. 2005; Samardzija et al. 2006; Das et al. 2010). These differences

257　may be due to the unique characteristics of the boar sperm. For example, the motility of the pig

258　sperm after ejaculation is slower than in other species, (e.g., horse and cattle), while it is also very

259　prone to be altered by a myriad of environmental incidences (Rodríguez-Gil and Bonet 2015). In

260　light of these singularities, we addressed the question of which sperm quality factors are influencing

261　SRR. The multivariate non-parametric test of independence revealed that none of the studied sperm

262　traits were related to SRR. This is somehow unexpected, particularly for motility and cell viability,

263　since a positive effect between these two traits and SRR have been previously described in cattle

11

(Samardzija et al. 2006). The differences in the physico-chemical properties between the ejaculate and the extender media in which the semen quality phenotypes are measured, and the BoviPure[TM] reagent during centrifugation, may divergently affect semen quality. This would explain the lack of dependency between the semen quality measures and SRR. A complementary hypothesis is that the time and speed of the density gradient centrifugation step enables all the boar's motile sperm, either fast or slow, to end up reaching the bottom of the tube, and be thus recovered. This would imply that the sperm recovery with BoviPure[TM] is not preferentially biased toward specific sperm sub-populations and therefore the molecular analysis of the recovered sperm robustly reflects that of the whole ejaculated mature sperm. Finally, SRR may be also affected by the composition and the physico-chemical characteristics of the ejaculates, which have been shown to  be affected by diet (Byrne et al. 2017), or abstinency in humans (Agarwal et al. 2016), which is related to the frequency of ejaculates or the time from prior ejaculate in pigs.

Two determinant parameters for a successful transcriptome analysis are both the RNA quality and yield. The RNA extraction method becomes a critical step when working with spermatozoa, since these cells have low amount of highly fragmented RNA. In the present work, we chose the Trizol method for RNA extraction after having tested other protocols involving commercial kits, which yielded even lower RNA yields (data not shown). The average amount of RNA extracted per sperm cell was 1.6 fg, a similar value to previously reported data in domestic swine (Yang et al. 2009), but lower than human (Goodrich et al. 2013; Pessot et al. 1989) and mice (Pessot et al. 1989). The low amount of RNA recovered and low RIN value is in fact an indication that the removal of somatic cells, with their large amount of non-fragmented RNA, during the cell purification steps, was highly efficient. The observed variability in RNA yields between samples could be due to inter-sample differences in the epididymosomes secreted by epididymal epithelial cells, which have been involved in post-testicular spermatogenesis and are known to contain a repertoire of RNAs (Belleannée et al. 2013), yet this mechanisms remains to be elucidated.

Spermatogenesis is a highly regulated process with many genes tightly controlling the different maturation steps (Legrand and Hobbs 2017) and playing a role in the sperm's fertility potential (Jodar et al. 2015). Our study in 190 samples suggests that the sperm quality parameters that we assessed are independent of the amount of RNA recovered - as a proxy of RNA load - per sperm cell. qPCR assays were also developed with the aim to determine the presence of RNA from somatic origin and gDNA contamination in our samples. Most of our samples showed at most, only traces of *PTPRC* (68 samples displayed $\Delta$Cq $_{\text{PTPRC-PRM1}}$ > 16) and gDNA was only detected in 4 samples with $\Delta$Cq $_{\text{Genomic-PRM1}}$ > 18.4. In qPCR, the amplification curve is exponential and the template doubles at every cycle. This amplification follows this formula: $X_N = X_1 * 2^N$, where $N$ is the number of amplification cycles, $X_1$ is the number of molecules prior amplification and $X_N$ is the number of molecules after $N$ PCR cycles. If we assume similar assay sensitivities, we can conclude that for a $\Delta$Cq $_{\text{PTPRC-PRM1}}$ = 16, the number of molecules of *PRM1* is $2^{16}$ = 65,536 times more abundant than the number of molecules of *PTPRC*. Likewise, when $\Delta$Cq $_{\text{Genomic-PRM1}}$ = 19, the number of *PRM1* RNA molecules is $2^{19}$ = 524,288 more abundant than the number of gDNA template. Hence, the majority of the RNA samples we processed were considered of sufficient spermatozoa purity. These qPCR assays can be used to determine sperm purity in porcine RNA samples and help selecting the purest RNAs for further analysis to obtain a reliable an accurate spermatozoa transcriptome. We must bear in mind that *PRM1* is also expressed in round spermatids (Siffroi et al. 1998; Steger et al. 2000) but we did not find any round-shaped cells in our samples following visual inspection of smear tissue under the microscope (Supplementary Figure 1). Thus, the presence of these cells is unlikely.

The purification and RNA extraction from boar sperm have proven to be suitable for the sequencing of total and small RNA by RNA-seq. To test the suitability of our samples for total RNA-seq analysis, we prepared sequencing libraries from six purified boar RNAs from different pigs with the SMARTer Universal Low Input RNA kit (Clontech) and with the TruSeq RNA Library Prep (Illumina) in parallel. Despite the fact that both protocols use the preferable amplification with

random primers instead of poly-dT (Mao et al. 2014), the SMARTer libraries still outperformed the TruSeq in several standard RNA-seq quality control parameters. Nevertheless, this is expected as the SMARTer protocol and chemistry is optimized for samples with low amount (10 ng) of highly fragmented RNA as for example, formalin-fixed paraffin embedded tissues. First, although the SMARTer protocol required less input RNA and it included a lower number of cycles in the amplification steps, it consistently yielded a much higher amount and concentration of cDNA library (Table 3), which is crucial to obtain optimal sequencing results. Second, even though the RNA-seq metrics of the two kits were similar (Table 3), the significantly higher proportion of PCR duplicates in the TruSeq datasets indicates a lower library complexity and number of unique transcripts. This was also indicated by the significant difference between the number of uniquely identified transcripts in both kits. The SMARTer datasets yielded twice the number of transcripts than TruSeq (13,233 versus 6,642). The vast majority of the TruSeq transcripts, 6,452, were also detected in the SMARTer dataset (Figure 2.A). This suggests that a proportion of RNAs were not captured with the TruSeq library preparation protocol. With the SMARTer kit we detected transcripts with lower abundance than with TruSeq and ultimately, a more comprehensive profile of the sperm transcriptome. We need to point out that there are other protocols from different providers, including Illumina, that have been designed for the RNA-seq analysis of samples with low amount and quality RNA, which have not been tested in this study. Furthermore, cluster analysis of the transcript levels shows a major kit effect, clustering together the libraries generated with the same kit rather than the libraries generated from matched RNAs (Figure 2.B).

To evaluate the adequacy of our samples for the sequencing of small RNAs, we used three samples and two short RNA library prep kits, the NEBNext Small RNA Library Prep Set (New England Biolabs) and the TailorMix miRNA Sample Prep v2 (SeqMatic). Both protocols showed very similar RNA-seq metrics (Table 4) but the analysis of the NEBNext libraries evidenced first, a larger number of detected miRNAs and second, slightly more similar miRNA abundance between samples

with a Pearson correlation of expression of 0.95 and 0.90 for the NEBNext and TailorMix, respectively.

RNA-seq is the gold standard approach for the genomic analysis of gene expression. This technology has been already applied to ejaculated sperm of different animal species such as human (Sendler et al. 2013), mouse (Fang et al. 2014), cattle (Card et al. 2013; Selvaraju et al. 2017) or horse (Das et al. 2013). In pig so far, the spermatozoa transcriptome was explored in 2009, using medium throughput sequencing approaches (Yang et al. 2009). The authors generated an Expressed Sequence Tag (EST) library and sequenced circa 5,000 clones using Sanger sequencing chemistry. This resulted in the identification of 271 genes with known function or cellular localization. That study, of high quality at that time, yielded a low number of transcripts and did not offer a comprehensive view of the boar sperm transcriptome (Yang et al. 2009). Current RNA-seq technologies, offer higher throughput and thousands of transcripts are typically detected. This is clearly illustrated in our study with the SMARTer datasets, whereby 13,233 genes were identified. The small non-coding transcriptome of the porcine sperm has been recently described (Chen et al. 2017). Chen and co-authors found a rich population of miRNAs and lower abundances of other families of small non coding RNAs (e.g., rRNAs, tRNAs, snRNAs) but dit not detect Piwi-interacting RNAs (piRNAs), a type of small non coding RNAs with important - yet weakly explored - functions in sperm biology. In contrast, our pipeline allowed us to identify a similar catalogue of small non coding RNAs but also piRNAs (Table 4 and Figure 2.C).

[Figure 2 near here].

Sperm RNA is likely to be transcriptionally inactive and contain fragmented transcripts that are remnants from spermatogenesis (Ostermeier et al. 2002) and RNA molecules that may function in signaling for embryogenesis after fertilization (Sendler et al. 2013). Hence, the study of the sperm transcriptome and the identification of its alterations could help the scientific community to identify robust, easy and non-invasive markers for sperm defects and male fertility. The purification method

364 and RNA extraction protocol described here, together with control qPCRs to evaluate the sperm

365 RNA purity warrants high quality RNA-seq experiments in boar sperm.

366 In conclusion, we have concatenated a series of well-established protocols to first, purify

367 spermatozoa from porcine ejaculates by gradient centrifugation, then extract RNA from the purified

368 sperm cells and finally prepare sequencing libraries from these samples to successfully sequence the

369 boar sperm-specific transcriptome by RNA-seq. We have also developed three qPCR assays to assess

370 the purity of the sperm RNA and compared the quality control metrics of different total and small

371 RNA-seq library preparation protocols. In addition, we have evaluated the boar's SRR with the

372 BoviPure$^{TM}$ and found that the boar's SRR is lower than in other mammalian species and not

373 dependent on any of the sperm quality parameters measured in our study. This recovered sperm was

374 thereafter used for RNA extraction. RNA yield per sperm cell was also lower than other species.

375 Moreover, we found no relationship between the quantity of RNA per sperm cell and the sperm

376 quality traits included in the analysis. Despite these caveats of the pig sperm, we obtained sufficient

377 sperm-specific RNA for RNA-seq studies. Thus, we recommend the methodological workflow

378 described here for the high throughput analysis of the boar spermatozoa transcriptome.

379

380

381

382

383

384

385

386

387

388

## Material and Methods

### *Sperm phenotyping*

From March 2015 to January 2017, specialized professionals at the farms obtained fresh ejaculates from 285 Pietrain boars kept in commercial farms. The ages of the animals ranged from 9 months to 5 years old. Sperm was collected with the hand glove method and immediately diluted (1:2) in freshly prepared commercial extender for storage at 16 ºC (MR-A extender; Kubus, S.L.; Majadahonda, Spain). No animal experiment has been performed in the scope of this research.

The samples were maintained at 16 ºC for a maximum time of 2 h for the phenotypic evaluation and a maximum of 4 h for the sperm cell purification. The analysis of sperm motility was performed with the commercial Computer-Aided Sperm Analysis (CASA) system (Integrated Sperm Analysis System V1.0; Proiser, Valencia, Spain) at 5 min and 90 min after incubation of the samples at 37 ºC. The percentages of sperm cell viability, structurally altered acrosomes and morphological abnormalities were measured after staining the samples with the eosin-nigrosin technique after 5 and 90 min incubation at 37 ºC as previously described (Bamba 1988). The osmotic resistance test (ORT) was performed by incubation at 37 °C for 10 min of the sperm samples on iso- and hypo-osmotic solutions, as previously described (Rodríguez-Gil and Rigau 1995). Sperm cell count was performed using a Neubauer cell chamber with not less than 200 cells examined.

### *Spermatozoa purification*

The purification of the spermatozoa cells was performed using 3 mL of BoviPure[TM] (Nidacon; Mölndal, Sweden), a commercial suspension of colloidal silica particles coated with silane in an isotonic salt solution, diluted to a final ratio of 60% (v/v) with BoviDilute[TM] (Nidacon; Mölndal, Sweden) in 15 mL RNase-free tubes. The volume of sperm that was layered on top of the cushion varied according to its concentration, with a maximum of 1 billion cells and not exceeding 11 mL. In all cases, the minimum volume ratio of 25 % diluted BoviPure[TM]/semen recommended by the manufacturer was maintained. The purification was made by centrifugation at 300 x g for 20 min at

414    20 °C with slow acceleration and deceleration rates (Allegra X-15R, Beckman Coulter; Brea, USA).

415    After centrifugation, all the upper phases were removed and the cell pellet was transferred to a new

416    RNase-free 15 mL tube, washed with 10 mL of RNase-free PBS and centrifuged at 1,500 x g for 10

417    min at 20 °C. The supernatant was then removed and the pellet was gently resuspended in 1 mL of

418    RNase-free PBS. Optical microscopy was used to confirm somatic cell removal of the purified

419    spermatozoa and sperm cell number was assessed in a Neubauer cell chamber. The resuspended

420    pellets were transferred to 1.5 mL RNase-free tubes and centrifuged at 1,500 x g for 10 min at 20 °C.

421    The resulting pellet was stored at -80 °C in 1 mL of Trizol® until further use for RNA extraction.

422    SRR was calculated as the number of cells obtained after purification divided by the initial number

423    of cells subjected to purification.

424    *RNA extraction*

425    Total RNA was extracted from 190 purified sperm samples, each from a different boar. The starting

426    number of cells ranged between 48 and 200 million (mean = 143 million) according to availability.

427    First, the cells were pre-lysated using a 5 mL sterile syringe with a 25 G needle for 5 min on ice,

428    followed by 2 min of vigorous vortex. Then, 200 µL of chloroform were added to the samples and

429    these were incubated for 3 min at room temperature first, and centrifuged at 12,000 x g for 15 min

430    afterwards. Supernatants were transferred to new RNase-free tubes and 500 µL of isopropanol were

431    added for RNA precipitation. Samples were then centrifuged at 12,000 x g for 10 min and the

432    supernatants were carefully removed. To wash the pellet, 500 µL of 75% (v/v) ethanol solution were

433    added and the samples were centrifuged at 13,000 x g for 5 min. The pellets were dried out at room

434    temperature for 10 min and resuspended in 30 µL of ultrapure water. All the centrifugations were

435    performed at 4 °C.

436    All RNA samples were subjected to DNase treatment with the Turbo DNA-free™ kit (Life

437    Technologies, USA) following the manufacturer's instructions. The RNA samples were then

438    quantified with Qubit™ RNA HS Assay kit (Invitrogen; Carlsbad, USA). To analyze overall RNA

439  fragmentation, RNA integrity number (RIN) was assessed on a 2100 Bioanalyzer using the Agilent

440  RNA 6000 Pico kit (Agilent Technologies; Santa Clara, USA). cDNA was synthesized using 2 µL of

441  RNA (1.7 - 38 ng) and the High Capacity cDNA Reverse Transcription kit in a final volume of 20

442  µL (Applied Biosystems; Waltham, USA) following the manufacturer's protocol.

443  *qPCR controls*

444  To verify that the purified samples were free of somatic cells and genomic DNA (gDNA), three

445  qPCR assays were developed. One assay targets *Protamine 1 (PRM1)* gene, which transcript is

446  specific to later stages of spermatogenesis and ejaculated mature spermatozoa (Wykes et al. 1997)

447  (PRM1_forward primer: 5'-AGTAGCAAGACCACCGCACT-3'; PRM1_reverse: 5'-

448  AGAGGGTCTTGAAGGCTGGT-3'). The second assay targets the *Protein tyrosine phosphatase*

449  *receptor type C* (*PTPRC*) gene, which is used as a marker of somatic cell contamination, since it is

450  expressed on most somatic cells and absent in spermatozoa (Das et al. 2013; Shafeeque et al. 2014)

451  (PTPRC_forward: 5'-AGAACAAGGTGGATGTCTATGGCTAT-3'; PTPRC_reverse: 5'-

452  TGTACTGTGCCTCCACCTGAAC-3'). The third assay amplifies an intergenic region

453  (Sscrofa10.2; chr18:25,459,856 - 25,459,926) and was designed to monitor the presence of gDNA

454  contamination (Intergenic_forward: 5'-ACGCAGTCAGAAGCCTGTGA-3'; Intergenic_reverse: 5'-

455  TGGTGTACATGCTCCGAAGGT-3').

456  To evaluate the performance of our qPCR assays, standard curves with serial dilutions of control

457  cDNA were made. For *PRM1*, *PTPRC* and gDNA qPCRs, pig cDNA from spermatozoa, liver and

458  gDNA were used as input, respectively. Liver cDNA was generated as indicated for sperm but using

459  1 µg of RNA starting material. The standard curve was generated with five ten-fold serial dilutions

460  of the cDNA templates. The reactions were performed with 10 µL of SYBR® Select Master Mix

461  (Life Technologies, USA), 0.3 µM of each primer, 5 µL of cDNA (for the serial dilutions 1/5 to

462  1/50,000 and for the query a dilution 1/5) or DNA (from 2 pg/uL to 20 ng/uL ) and ultrapure water to

463  a final volume of 20 µL. The thermal profile was: 50 ºC for 2 min, 95 ºC for 10 min and 40 cycles of

464  95 °C for 15 sec and 60 °C for 1 min. Moreover, to assess the specificity of the qPCR reactions, a

465  melting profile (95 °C for 15 sec, 60 °C for 15 sec and a gradual increase in temperature with a ramp

466  rate of 1% up to 95 °C) was programmed following the thermal cycling protocol. A minus reverse

467  transcription control was also included for the two cDNA assays (*PRM1* and *PTPRC*). The reactions

468  for the standard curves were performed in triplicate. For the queried samples (N = 70), the reactions

469  were performed in triplicate. Moreover, a liver cDNA was also included to monitor the expression of

470  *PRM1* and *PTPRC* in a somatic cell type.

471  ***Statistical analysis***

472  R v.3.3.0 was utilized for statistical analysis. The Shapiro-Wilk test was used to assess normality of

473  the data. One-way analysis of variance (ANOVA) was used to assess the effects of farm (N = 3), age

474  (N = 3) and batch collection day (N = 59) on SRR, fg per sperm cell and sperm quality parameters.

475  Significantly correlated covariates were adjusted with the R package 'limma' (Ritchie et al. 2015),

476  considering age and farm as fixed effects and batch collection day as batch effect. The R package

477  'energy' (Rizzo and Szekely 2016) was applied to assess a multivariate nonparametric test of

478  independence covariates between sperm quality phenotypes and SRR and fg of RNA per sperm cell.

479  The nominal significance threshold was set to a *P*-value $\leq$ 0.0041 after Bonferroni correction for

480  multiple testing ($\alpha$ = 0.05/12 = 0.0041). To determine whether the RNA-seq quality metrics of the

481  SMARTer Universal Low Input RNA and the TruSeq RNA library prep kits were significantly

482  different, we used the t-test for normally distributed data for reads mapped to the genome, number of

483  uniquely mapped reads and the number of detected genes, and the Wilcoxon test for the non-

484  normally distributed data, i.e. library concentration and proportion of PCR duplicates. The tests were

485  carried with R.

486  ***RNA-seq library preparation***

487  Purified RNA from six ejaculates from different boars (Sample_1 to Sample_6) was subjected to

488  ribosomal RNA depletion with the Ribo-Zero Gold rRNA Removal Kit (Illumina). Depleted RNA

489  was then used to prepare long RNA-seq libraries with two different protocols in parallel. On the one

490  hand, SMARTer Universal Low Input RNA library Prep kit (Clontech) was used, following the

491  manufacturer's instructions. On the other hand, TruSeq RNA Library Prep kit (Illumina) was

492  employed, adhering to the manufacturer's protocol with slight modifications adapted to a low

493  amount of starting RNA yield (100 ng). The concentration of the 12 RNA libraries was quantified

494  with the High Sensitivity DNA kit on a 2100 Bioanalizer (Agilent Technologies). The libraries were

495  sequenced in a HiSeq2000 system (Illumina) to generate 75 bp long paired-end reads.

496  RNA from 3 additional sperm samples (Sample_7 to Sample_9) was used to compare two short

497  RNA library Prep kits: NEBNext Small RNA Library Prep Set (New England Biolabs) and

498  TailorMix miRNA Sample Prep v2 (SeqMatic). One sample was prepared with the NEBNext kit,

499  another with the TailorMix, and a third sample with both kits to allow a more direct comparison of

500  results. Libraries were prepared following the company's instructions. The three samples were

501  quantified with the High Sensitivity DNA kit on a 2100 Bioanalizer and sequenced on a HiSeq2000

502  to generate 50 bp single-end reads.

503  *Bioinformatics analysis*

504  Read quality of the 12 long RNA-seq datasets was checked with FastQC v.0.11.2 (Andrews 2010).

505  Reads were then filtered using Trimmomatic v.0.33 (Bolger et al. 2014) for read quality and adaptor

506  contamination, with a minimum Phred quality score of 20 and length of over 30 bp. Trimmed reads

507  were mapped to the pig reference genome (Sscrofa 10.2) with STAR v.2.5.3a (Dobin et al. 2013)

508  using the default parameters and including the Ensembl v.83 pig reference annotation

509  (ftp://ftp.ensembl.org/pub/release-83/gtf/sus_scrofa). Transcript abundance was quantified as

510  Fragments Per Kilobase of transcript per Million mapped reads (FPKM) with RSem v.1.3.0 (Li and

511  Dewey 2011) with default parameters. FPKM is a normalized measure of gene expression based on

512  the number of reads mapping to a given gene corrected by the length of that gene and the sample

513  sequencing depth. To compare the performance of the SMARTer and the TruSeq protocols we

evaluated the proportion of PCR duplicates as it is a measure of the complexity of each library. To allow a fair comparison of the two protocols we analyzed the same number of reads in all the samples. More in detail, we randomly sub-selected 2,336,549 reads per sample since this number corresponds to the lowest sequencing depth obtained. The read selection was carried with seqtk v.1.2 (Shen et al. 2016). The proportion of PCR duplicates was calculated with Picard Tools v.1.110 (http://picard.sourceforge.net) MarkDuplicates. Graphs were performed with R: venn diagram with the R package 'VennDiagram' (Chen and Boutros 2011) and cluster dendrogram with the R package 'cluster' (Maechler et al. 2017).

We also evaluated the absence of RNA from somatic cell origin in our samples. For this we used the SMARTer libraries, which showed better outcomes, and the totality of the reads obtained for each of these libraries (between 18.5 and 26.9 million reads per sample). We also included two publicly available (http://www.ncbi.nlm.nih.gov/sra) boar RNA-seq datasets, one from whole blood cells (ERR1898477), which contains a large abundance of leukocytes and a second one from ear biopsy (SRR3437133), which contains a high proportion of keratinocytes, a specialized type of epithelial cells. We screened the presence of the two sperm specific genes, *PRM1* and *OAZ3* (*Ornithine Decarboxylase Antizyme 3*) (Jodar *et al.*, 2016), and two genes of somatic cell origin *PTPRC* (expressed in most somatic cells) and *KRT1* (*Keratin 1*), which is specific from keratinocytes. For data visualization, SMARTer mapped bam files were indexed with SAMtools v.1.3.0 (Li et al. 2009) and uploaded into the IGV viewer (Thorvaldsdóttir et al. 2013). We used a manual script to extract RNA levels of tissue-specific genes as described in (Jodar et al. 2016) as an ultimately control for RNA purity.

The 3 small RNA-seq datasets were analyzed for read quality with FastQC v.0.11.2 (Andrews 2010) and reads were sub-sampled to 887,406 reads per library with seqtk v.1.2 (Shen et al. 2016). Library adaptors and indexes were trimmed using Cutadapt v.1.0 (Martin 2011) and filtered for read quality, with a minimum quality score of 20, and minimum length of 10 bp with Trimmomatic v.0.33 (Bolger

539    et al. 2014). Trimmed reads were mapped to the pig reference genome (Sscrofa 10.2) with Bowtie 1

540    v.1.2.0 (Langmead 2010) with default parameters but allowing 0 mismatches (-n) in 'seed' region of

541    10 bp (-l). The proportion of PCR duplicates was calculated with SAMtools v.1.3.0 (Li et al. 2009)

542    rmdup for single-end reads. RNA levels of small non-coding RNAs were calculated with Bedtools

543    v.2.17.0 (Quinlan and Hall 2010) intersect against the boar Ensembl v.83 'gtf' annotation, miRBase

544    database (Griffiths-Jones et al. 2006), piRNA database (Rosenkranz 2016), and tRNA v.2.0 database

545    (Chan and Lowe 2016).

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

589 **Disclosure of interest:**

590 The authors report no conflict of interest

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

**Notes on contributors**

614 **Notes on contributors**

615 ACl, AS, FQM and MG conceived and designed the experiment. FQM and ACa designed the

616 primers and carried out the qPCR analyses. FQM, MG and JN performed sperm purifications and

617 RNA extractions. MG performed statistic and bioinformatics analysis. FQM and MG analyzed the

618 data. JERG carried the phenotypic analysis. FQM, MG and AC drafted the manuscript. All authors

619 read and approved the final manuscript.

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

**References**

Agarwal, A., Gupta, S., Du Plessis, S., Sharma, R., Esteves, S. C., Cirenza, C., Eliwa, J., Al-Najjar, W., Kumaresan, D., Haroun, N., Philby, S., and Sabanegh, E. (2016) Abstinence Time and Its Impact on Basic and Advanced Semen Parameters. Urology. **94**:102–110.

Allamaneni, S. S. R., Agarwal, A., Rama, S., Ranganathan, P., and Sharma, R. K. (2005) Comparative study on density gradients and swim-up preparation techniques utilizing neat and cryopreserved spermatozoa. Asian J Androl. **7**(86):86–92.

Andrews, S. (2010) FastQC: A quality control tool for high throughput sequence data. http://www.bioinformatics.babraham.ac.uk/projects/fastqc/.

Aoki, V. W., Moskovtsev, S. I., Willis, J., Liu, L., Mullen, J. B. M., and Carrell, D. T. (2005) DNA integrity is compromised in protamine-deficient human sperm. J Androl. **26**(6):741–748.

Bamba, K. (1988) Evaluation of acrosomal integrity of boar spermatozoa by bright field microscopy using an eosin-nigrosin stain. Theriogenology. **29**(6):1245–1251.

Belleannée, C., Calvo, É., Caballero, J., and Sullivan, R. (2013) Epididymosomes Convey Different Repertoires of MicroRNAs Throughout the Bovine Epididymis. Biol Reprod. **89**(2):30.

Bissonnette, N., Lévesque-Sergerie, J. P., Thibault, C., and Boissonneault, G. (2009) Spermatozoal transcriptome profiling for bull sperm motility: A potential tool to evaluate semen quality. Reproduction. **138**(1):65–80.

Bohacek, J., and Mansuy, I. M. (2015) Molecular insights into transgenerational non-genetic inheritance of acquired behaviours. Nat Rev Genet **16**(11):641–652.

Bolger, A. M., Lohse, M., and Usadel, B. (2014) Trimmomatic : a flexible trimmer for Illumina sequence data. Bioinformatics. **30**(15):2114–2120.

Byrne, C. J., Fair, S., English, A. M., Holden, S. A., Dick, J. R., Lonergan, P., and Kenny, D. A. (2017) Dietary polyunsaturated fatty acid supplementation of young post-pubertal dairy bulls alters the fatty acid composition of seminal plasma and spermatozoa but has no effect on semen volume or

664    sperm quality. Theriogenology. **90**:289–300.

665    Capra, E., Turri, F., Lazzari, B., Cremonesi, P., Gliozzi, T. M., Fojadelli, I., Stella, A., and Pizzi, F.

666    (2017) Small RNA sequencing of cryopreserved semen from single bull revealed altered miRNAs

667    and piRNAs expression between High- and Low-motile sperm populations. BMC Genomics.

668    **18**(1):14.

669    Card, C. J., Anderson, E. J., Zamberlan, S., Krieger, K. E., Kaproth, M., and Sartini, B. L. (2013)

670    Cryopreserved bovine spermatozoal transcript profile as revealed by high-throughput ribonucleic

671    acid sequencing. Biol Reprod. **88**(2):49.

672    Carrell, D. T., Emery, B. R., and Hammoud, S. (2007) Altered protamine expression and diminished

673    spermatogenesis: What is the link? Hum Reprod Update. **13**(3):313–327.

674    Chan, P. P., and Lowe, T. M. (2016) GtRNAdb 2.0: An expanded database of transfer RNA genes

675    identified in complete and draft genomes. Nucleic Acids Res. **44**(D1):D184–D189.

676    Chen, C., Wu, H., Shen, D., Wang, S., Zhang, L., Wang, X., Gao, B., Wu, T., Li, B., Li, K., and

677    Song, C. (2017) Comparative profiling of small RNAs of pig seminal plasma and ejaculated and

678    epididymal sperm. Reproduction. **153**(6):785–796.

679    Chen, H., and Boutros, P. C. (2011) VennDiagram: A package for the generation of highly-

680    customizable Venn and Euler diagrams in R. BMC Bioinformatics. **12**:35.

681    Curry, E., Safranski, T. J., and Pratt, S. L. (2011) Differential expression of porcine sperm

682    microRNAs and their association with sperm morphology and motility. Theriogenology. **76**(8):1532–

683    1539.

684    Das, P. J., McCarthy, F., Vishnoi, M., Paria, N., Gresham, C., Li, G., Kachroo, P., Sudderth, A. K.,

685    Teague, S., Love, C. C., Varner, D. D., Chowdhary, B. P., and Raudsepp, T. (2013) Stallion Sperm

686    Transcriptome Comprises Functionally Coherent Coding and Regulatory RNAs as Revealed by

687    Microarray Analysis and RNA-seq. PLoS ONE. **8**(2):e56535.

688    Das, P. J., Paria, N., Gustafson-Seabury, A., Vishnoi, M., Chaki, S. P., Love, C. C., Varner, D. D.,

689   Chowdhary, B. P., and Raudsepp, T. (2010) Total RNA isolation from stallion sperm and testis

690   biopsies. Theriogenology. **74**(6):1099–1106.e2.

691   Dobin, A., Davis, C. A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P., Chaisson, M.,

692   and Gingeras, T. R. (2013) STAR: Ultrafast universal RNA-seq aligner. Bioinformatics **29**(1):15–21.

693   Fang, P., Zeng, P., Wang, Z., Liu, M., Xu, W., Dai, J., Zhao, X., Zhang, D., Liang, D., Chen, X., Shi,

694   S., Zhang, M., Wang, L., Qiao, Z., and Shi, H. (2014) Estimated Diversity of Messenger RNAs in

695   Each Murine Spermatozoa and Their Potential Function During Early Zygotic Development. Biol

696   Reprod. **90**(5):1–11.

697   Feugang, J. M. (2017) Novel agents for sperm purification, sorting, and imaging. Mol Reprod Dev.

698   **84**(9):832–841.

699   Goodrich, R. J., Anton, E., and Krawetz, S. A. (2013) Isolating mRNA and small noncoding RNAs

700   from human sperm. Methods Mol Biol. **927**:385–396.

701   Griffiths-Jones, S., Grocock, R. J., van Dongen, S., Bateman, A., and Enright, A. J. (2006) miRBase:

702   microRNA sequences, targets and gene nomenclature. Nucleic Acids Res. **34**(Database issue):D140-

703   4.

704   Hamatani, T. (2012) Human spermatozoal RNAs. Fertil Steril. **97**(2):275–281.

705   Hammoud, S. S., Nix, D. A., Hammoud, A. O., Gibson, M., Cairns, B. R., and Carrell, D. T. (2011)

706   Genome-wide analysis identifies changes in histone retention and epigenetic modifications at

707   developmental and imprinted gene loci in the sperm of infertile men. Hum Reprod. **26**(9):2558–

708   2569.

709   Jameel, T. (2008) Sperm swim-up: A simple and effective technique of semen processing for

710   intrauterine insemination. J Pak Med Assoc. **58**(2):71–74.

711   Jodar, M., Kalko, S., Castillo, J., Ballescà, J. L., and Oliva, R. (2012) Differential RNAs in the sperm

712   cells of asthenozoospermic patients. Hum Reprod. **27**(5):1431–1438.

713   Jodar, M., Selvaraju, S., Sendler, E., Diamond, M. P., and Krawetz, S. A. (2013) The presence, role

714 and clinical use of spermatozoal RNAs. Hum Reprod Update. **19**(6):604–624.

715 Jodar, M., Sendler, E., Moskovtsev, S. I., Librach, C. L., Goodrich, R., Swanson, S., Hauser, R.,

716 Diamond, M. P., and Krawetz, S. A. (2015) Absence of sperm RNA elements correlates with

717 idiopathic male infertility. Sci Transl Med. **7**(295):295re6.

718 Jodar, M., Sendler, E., Moskovtsev, S. I., Librach, C. L., Goodrich, R., Swanson, S., Hauser, R.,

719 Diamond, M. P., and Krawetz, S. A. (2016) Response to Comment on 'Absence of sperm RNA

720 elements correlates with idiopathic male infertility'. Sci Transl Med. **8**(353):353tr1.

721 Langmead, B. (2010) Aligning short sequencing reads with Bowtie. Current Protocols in

722 Bioinformatics (SUPP.32):Chapter Unit-11.7.

723 Legrand, J. M. D., and Hobbs, R. M. (2017) RNA processing in the male germline: Mechanisms and

724 implications for fertility. Semin Cell Dev Biol. https://doi.org/10.1016/j.semcdb.2017.10.006.

725 Li, B., and Dewey, C. N. (2011) RSEM: accurate transcript quantification from RNA-Seq data with

726 or without a reference genome. BMC Bioinformatics. **12**(1):323.

727 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., and

728 Durbin, R. (2009) The Sequence Alignment/Map format and SAMtools. Bioinformatics.

729 **25**(16):2078–2079.

730 Maechler, M., Rousseeuw, P., Struyf, A., Hubert, M., and Hornik, K. (2017) Cluster Analysis Basics

731 and Extensions. R package version 2.0.6. http://cran.r-project.org/web/packages/cluster/index.html.

732 Mao, S., Goodrich, R. J., Hauser, R., Schrader, S. M., Chen, Z., and Krawetz, S. A. (2013)

733 Evaluation of the effectiveness of semen storage and sperm purification methods for spermatozoa

734 transcript profiling. Syst Biol Reprod Med. **59**(5):287–295.

735 Mao, S., Sendler, E., Goodrich, R. J., Hauser, R., and Krawetz, S. A. (2014) A comparison of sperm

736 RNA-seq methods. Syst Biol Reprod Med. **60**(5):308–315.

737 Martin, M. (2011) Cutadapt removes adapter sequences from high-throughput sequencing reads.

738 EMBnet.journal. **17**(1):10–12.

739  Martins, R. P., and Krawetz, S. A. (2005) Towards understanding the epigenetics of transcription by

740  chromatin structure and the nuclear matrix. Gene Ther Mol Biol. **9**(B):229–246.

741  Ostermeier, G. C., Dix, D. J., Miller, D., Khatri, P., and Krawetz, S. A. (2002) Spermatozoal RNA

742  profiles of normal fertile men. Lancet. **360**(9335):772–777.

743  Patil, P. S., Humbarwadi, R. S., Patil, A. D., and Gune, A. R. (2013) Immature germ cells in semen -

744  correlation with total sperm count and sperm motility. J Cytol. **30**(3):185–189.

745  Pessot, C. A., Brito, M., Figueroa, J., Concha, I. I., Yañez, A., and Burzio, L. O. (1989) Presence of

746  RNA in the sperm nucleus. Biochem Biophys Res Commun. **158**(1):272–278.

747  Quinlan, A. R., and Hall, I. M. (2010) BEDTools: A flexible suite of utilities for comparing genomic

748  features. Bioinformatics. **26**(6):841–842.

749  Rando, O. J. (2016) Intergenerational transfer of epigenetic information in sperm. Cold Spring Harb

750  Perspect Med. **6**(5):a022988.

751  Ritchie, M. E., Phipson, B., Wu, D., Hu, Y., Law, C. W., Shi, W., and Smyth, G. K. (2015) limma

752  powers differential expression analyses for RNA-sequencing and microarray studies. Nucleic Acids

753  Res. **43**(7):e47.

754  Rizzo ML and Szekely GJ. (2016) energy: E-Statistics: Multivariate Inference via the Energy of

755  Data. version 1.7-2 http://cran.r-project.org/package=energy.

756  Rodríguez-Gil, J. E., and Bonet, S. (2015) Current knowledge on boar sperm metabolism:

757  Comparison with other mammalian species. Theriogenology. **85**(1):4–11.

758  Rodríguez-Gil, J. E., and Rigau, T. (1995) Effects of slight agitation on the quality of refrigerated

759  boar sperm. Anim Reprod Sci. **39**(2):141–146.

760  Rosenkranz, D. (2016) piRNA cluster database: A web resource for piRNA producing loci. Nucleic

761  Acids Res. **44**(D1):D223–D230.

762  Samardzija, M., Karadjole, M., Matkovic, M., Cergolj, M., Getz, I., Dobranic, T., Tomaskovic, A.,

763  Petric, J., Surina, J., Grizelj, J., and Karadjole, T. (2006) A comparison of BoviPure and Percoll on

764     bull sperm separation protocols for IVF. Anim Reprod Sci. **91**(3–4):237–247.

765     Selvaraju, S., Parthipan, S., Somashekar, L., Kolte, A. P., Krishnan Binsila, B., Arangasamy, A., and

766     Ravindra, J. P. (2017) Occurrence and functional significance of the transcriptome in bovine (Bos

767     taurus) spermatozoa. Sci Rep. **7**:42392.

768     Sendler, E., Johnson, G. D., Mao, S., Goodrich, R. J., Diamond, M. P., Hauser, R., and Krawetz, S.

769     A. (2013) Stability, delivery and functions of human sperm RNAs at fertilization. Nucleic Acids Res.

770     **41**(7):4104–4117.

771     Shafeeque, C. M., Singh, R. P., Sharma, S. K., Mohan, J., Sastry, K. V. H., Kolluri, G., Saxena, V.

772     K., Tyagi, J. S., Kataria, J. M., and Azeez, P. A. (2014) Development of a new method for sperm

773     RNA purification in the chicken. Anim Reprod Sci. **149**(3–4):259–265.

774     Shen, W., Le, S., Li, Y., and Hu, F. (2016) SeqKit: A Cross-Platform and Ultrafast Toolkit for

775     FASTA/Q File Manipulation. PLOS ONE. **11**(10):e0163962.

776     Siffroi, J. P., Alfonsi, M. F., and Dadoune, J. P. (1998) Electron microscopic in situ hybridization

777     study of simultaneous expression of TNP1 and PRM1 genes in human spermatids. Ital J Anat

778     Embryol. **103**(4 Suppl 1):65–74.

779     Steger, K., Pauls, K., Klonisch, T., Franke, F. E., and Bergmann, M. (2000) Expression of

780     protamine-1 and -2 mRNA during human spermiogenesis. Mol Hum Reprod **6**(3):219–225.

781     Székely, G. J., and Rizzo, M. L. (2009) Brownian distance covariance. Ann Appl Stat. **3**(4):1236–

782     1265.

783     Thorvaldsdóttir, H., Robinson, J. T., and Mesirov, J. P. (2013) Integrative Genomics Viewer (IGV):

784     High-performance genomics data visualization and exploration. Brief Bioinform. **14**(2):178–192.

785     Wang Z, Gerstein M, Snyder M. (2013) RNA-Seq: a revolutionary tool for transcriptomics. Nat Rev

786     Genet. **10**(1):57-63.

787     Wykes, S. M., Visscher, D. W., and Krawetz, S. A. (1997) Haploid transcripts persist in mature

788     human spermatozoa. Mol Hum Reprod. **3**(1):15-19.

789   Yang, C. C., Lin, Y. S., Hsu, C. C., Wu, S. C., Lin, E. C., and Cheng, W. T. K. (2009) Identification

790   and sequencing of remnant messenger RNAs found in domestic swine (Sus scrofa) fresh ejaculated

791   spermatozoa. Anim Reprod Sci. **113**(1–4):143–155.

792   Zhao, Y., Li, Q., Yao, C., Wang, Z., Zhou, Y., Wang, Y., Liu, L., Wang, Y., Wang, L., and Qiao, Z.

793   (2006) Characterization and quantification of mRNA transcripts in ejaculated spermatozoa of fertile

794   men by serial analysis of gene expression. Hum Reprod. **21**(6):1583–1590.

795

**Figure legends**

Figure 1. Read mapping depth of sperm and somatic-specific genes in the porcine sperm, whole blood and ear RNA-seq datasets. A) Corresponds to the sperm-specific gene *PRM1* (Ensembl gene ID: ENSSSCG00000021337). B) Plot for the sperm-specific gene *OAZ3* (ENSSSCG00000027091). C) Read depth along the somatic cell specific gene *PTPRC* (ENSSSCG00000010908). D) Depth for the Keratinocyte specific gene *KRT1* (ENSSSCG00000000251). The number of reads produced for the sperm datasets are between 18.5 and 26.5 million. The white blood cells and the ear RNA-seq libraries include 18.3 and 21.7 million reads, respectively. The scale provided on the upper left axis of each graph indicates the raw number of reads mapped to the gene.

Figure 2. Comparison of the RNA-seq results from both the total and the small library preparation kits. A) Venn diagram representing the 13,233 different transcripts detected in the SMARTer datasets and 6,642 for the TruSeq. The majority of the TruSeq transcripts detected (6,452) were detected in both kits. B) Cluster dendrogram of the RNA transcript levels from both kits. S: Sample. C) Size distribution of mapped sequencing reads from small RNAs in both kits.