# Cross-platform transcriptomic profiling of the response to recombinant human erythropoietin

Guan Wang  ( ✉ g.wang2@brighton.ac.uk )
University of Brighton

**Traci Kitaoka**
Illumina (United States)

**Ali Crawford**
Illumina (United States)

**Qian Mao**
BGI Group (China)

**Andrew Hesketh**
University of Brighton

**Fergus Guppy**
University of Brighton

**Garrett Ash**
Yale University

**Jason Liu**
Yale University

**Mark Gerstein**
Yale University

**Yannis Pitsiladis**
University of Brighton

1 **Cross-platform transcriptomic profiling of the response to recombinant**

2 **human erythropoietin**

3

4 Guan Wang[1,2]*, Traci Kitaoka[3], Ali Crawford[3], Qian Mao[4], Andrew Hesketh[5], Fergus M.

5 Guppy[5,6], Garrett I. Ash[7,8], Jason Liu[9], Mark B. Gerstein[9-12], Yannis P. Pitsiladis[6]*

6

7 [1]Sport and Exercise Science and Sports Medicine Research and Enterprise Group, University of

8 Brighton, Brighton, UK.

9 [2]Centre for Regenerative Medicine and Devices, University of Brighton, Brighton, UK.

10 [3]Illumina, San Diego, CA, USA.

11 [4]BGI, Shenzhen, China.

12 [5]School of Pharmacy and Biomolecular Sciences, University of Brighton, Brighton, UK.

13 [6]Centre for Stress and Age-related Disease, University of Brighton, Brighton, UK.

14 [7]Veterans Affairs Connecticut Healthcare System, West Haven, CT, USA.

15 [8]Center for Medical Informatics, Yale University, New Haven, CT, USA.

16 [9]Program in Computational Biology & Bioinformatics, Yale University, New Haven, CT, USA.

17 [10]Department of Molecular Biophysics & Biochemistry, Yale University, New Haven, CT, USA.

18 [11]Department of Computer Science, Yale University, New Haven, CT, USA.

19 [12]Department of Statistics and Data Science, Yale University, New Haven, CT, USA.

20 *Corresponding authors. Email: g.wang2@brighton.ac.uk; y.pitsiladis@brighton.ac.uk.

21     **Abstract**

22     RNA-seq has matured and become an important tool for studying RNA biology. Here we

23     compared two RNA-seq (Illumina sequencing by synthesis and MGI DNBSEQ$^{TM}$) and two

24     microarray platforms (Illumina Expression BeadChip and GeneChip$^{TM}$ Human Transcriptome

25     Array 2.0) in healthy individuals administered recombinant human erythropoietin for

26     transcriptome-wide quantification of differential gene expression. The results show that total

27     RNA sequencing combined with DNB-seq produced a multitude of genes of biological relevance

28     and significance in response to recombinant human erythropoietin, in contrast to other platforms.

29     Through data triangulation linking genes to functions, genes representing the processes of

30     erythropoiesis as well as non-erythropoietic functions of erythropoietin were unveiled. This

31     study provides a knowledge base of genes characterising the responses to recombinant human

32     erythropoietin through cross-platform comparison and validation.

33

## Introduction

High-throughput technologies in gene discovery, quantification and functional investigation have advanced our understanding of complex traits and facilitated disease diagnosis, prevention and treatment over the past decade[1, 2]. Although technologies continue to evolve for discerning and characterising genes and gene-protein interactions both *ex-* and *in-vivo*, uncovering coding transcriptomes of bulk cells can capture global gene expression patterns that may directly pinpoint important biological processes at the molecular level. Which tool to use will ultimately depend on the fundamental research question. Here, we performed RNA-seq and microarray analyses in healthy individuals administered recombinant human erythropoietin (rHuEPO) to assess their discriminatory capacity, and importantly, to explore the implications of the findings to better understanding the systemic responses to rHuEPO; a first of its kind in the investigation of transcriptome-wide responses to rHuEPO in humans. This study primarily differs from previous cross-platform gene-expression studies in 1) systematic comparisons between two RNA-seq platforms (MGI DNBSEQ-G400RS and Illumina NextSeq 500), 2) comparisons with the benchmarking microarrays (GeneChip[TM] Human Transcriptome Array 2.0 and Illumina HumanHT-12 v4 Expression BeadChip), 3) the use of a relatively large number of the same experimental samples across all four platforms, and 4) the adoption of a data triangulation approach across platforms to prioritise the functional genes of diagnostic potential.

Eighteen endurance-trained Caucasian males at sea level (Glasgow, Scotland; age: 26.0±4.5 yrs, weight: 74.8±7.9kg, height: 179.8±5.4cm) underwent 4 weeks of rHuEPO injections (50 IU/kg every 2 days)[3]. Whole blood samples collected from the 18 subjects across 8 time points — before (−14- and −1-day prior to the first injection; Base1 and Base2), during (2-, 14- and 28-day

57    into the administration; EPO3, EPO4 and EPO5) and post rHuEPO administration (2-, 14- and

58    28-day after the last injection; Post6, Post7 and Post8) — were analysed on the Illumina

59    HumanHT-12 v4 Expression BeadChip previously[4]. In the current study, 50 samples from 10 of

60    the 18 subjects collected at Base1, Base2, EPO3, EPO4 and Post7 were analysed on the two

61    RNA-seq platforms as well as on the GeneChip[TM] array for quantifying differential gene

62    expression (DGE). This experimental design aimed to identify the gene expression response to

63    rHuEPO through robust quantification processes, and to generate results with wide applications

64    ranging from developing effective therapeutics targeting clinical disorders associated with EPO

65    dysfunctions to facilitating sensitive testing strategies against blood doping in sport.

66

67    **Results**

68    **Total RNA DNB-seq (MGI) identifies a wealth of mRNA genes in response to rHuEpo**

69    We identified 16,738 genes (MGI RNA-seq), 16,581 genes (Illumina RNA-seq), 29,517

70    transcript clusters (GeneChip), and 10,622 transcripts (BeadChip) for the DGE analyses (Table

71    1). Both MGI and Illumina RNA-seq generated good base call quality, with an average quality

72    score of >34 across the read lengths and across the samples (Supplementary Fig. 1). No sample

73    contamination/swaps (Supplementary Fig. 2) and no other significant surrogate variables of

74    batch effects were detected in these sequencing datasets. Genome mapping using HISAT2[5]

75    (against the reference genome assembly GRCh38.p12[6]) showed the overall alignment rates of

76    94.2% (MGI; 197.9M total reads) and 95.0% (Illumina; 110.4M total reads) (Supplementary

77    Table 1). RseQC[7] revealed a large proportion of the sequences aligned to introns in the MGI

78    RNA-seq data (37.4% versus 8.7% Illumina on average; Supplementary Fig. 3 and

79    Supplementary Table 2), a result coinciding with the differing sequencing library preparation

80    methods used (total RNA-seq with rRNA depletion and globin mRNA reduction, MGI versus

81    mRNA enrichment, Illumina). RseQC also showed that among a total of 186.6M (MGI) and

82    104.6M (Illumina) averaged reads observed, 52.4% (~ 97.8M reads) and 74.8% (~ 78.2M reads)

83    of the reads were effectively mapped to the coding sequences (exons), respectively

84    (Supplementary Table 2). The average Salmon[8] transcriptome mapping rates, following

85    selective-alignment-based lightweight mapping, were 38.8% (38.3M aligned reads; MGI) and

86    81.9% (45.1M aligned reads; Illumina) (Supplementary Table 3). The seeming discrepancies

87    observed in total reads and alignment rates across the software tools (HISAT2, RseQC and

88    Salmon) were expected given their specific usage. Overall, these data suggested high quality

89    sequences obtained from both sequencing platforms. For the purposes of cross-platform

90    comparison, the relative abundance estimates of transcripts after Salmon transcriptome mapping

91    were summed to gene level, and genes were considered expressed when the gene-level

92    abundance estimates were equal to or more than 5 in at least 4 samples; resulting in the exclusion

93    of 17,198 and 18,347 genes from the MGI and Illumina RNA-seq datasets, respectively (Table

94    2). Gene annotation resulted in 3,852 (MGI) and 2,860 (Illumina) un-defined gene mappings

95    removed from the sequencing datasets (Table 2). As a result, 16,738 and 16,581 protein-coding

96    genes identified from MGI and Illumina sequencing, respectively, were used for the downstream

97    DGE analyses (Table 1 and 2).

98

99    Initial quality control metrics revealed variability in eight out of the fifty GeneChip[TM] arrays

100   (Supplementary Fig. 4A). Two of the eight samples were then repeated for chip scanning, and

101   the other six samples were repeated from the target preparation step (Supplementary Fig. 4B).

102   Raw intensity values obtained from the GeneChip and BeadChip analyses correspond to 67,480

103    and 47,286 coding and non-coding transcriptomic features, respectively (Table 1). The process

104    of normalisation and filtering unveiled 29,517 transcript clusters (GeneChip) and 10,622

105    transcripts (BeadChip) as identified features (Table 1), with the detailed filtering steps and the

106    resulting number of features summarised in Table 2. Briefly, 18,494 and 6,900 probes were

107    removed as undetected and low-quality probes, respectively, from the BeadChip dataset (Table

108    2). While 8,166 probes were removed due to low average expression (cutoff value: 5.1) in the

109    BeadChip dataset, no such probes were necessarily excluded from the GeneChip dataset

110    (Supplementary Fig. 5 and Table 2). No significant surrogate variables representing the

111    underlying biases, potentially arising from library preparation and/or scanning, thereby

112    confounding the biological effects being studied, were observed in the two microarray datasets.

113

114    Unsupervised principal component analysis (PCA) revealed substantial variance, estimated using

115    the top 500 genes ranked by expression variance across all samples. Variances explained by the

116    principal component 1 and the principal component 2 were: 69% vs. 5% (MGI RNA-seq), 44%

117    vs. 9% (Illumina RNA-seq), 58% vs. 14% (GeneChip), and 78% vs. 7% (Beadchip)

118    (Supplementary Fig. 6). Gene clustering of the top 30 genes of high variance showed a good

119    distinction across biological conditions in all datasets (Supplementary Fig. 7). Nevertheless, a

120    more distinctive pattern across the conditions was observed following MGI RNA-seq compared

121    to Illumina RNA-seq and GeneChip (Supplementary Fig. 7, A versus B, C). In contrast with the

122    discrimination pattern presented in the 143 BeadChip samples, a higher expression level of the

123    examined top 30 genes was detected by MGI RNA-seq (Supplementary Fig. 7, A versus D). The

124    DESeq2[9] and limma[10] DGE analyses yielded 1,552, 582, 252 and 2,372 transcriptomic features

125    exceeding the pre-defined thresholds following MGI RNA-seq, Illumina RNA-seq, GeneChip[TM]

126     and BeadChip, respectively (thresholds for RNA-seq: a fold change of 1.2 and $s$-value of 0.005;

127     for microarray: a fold change of 1.2 and BH adjusted $p$-value of 0.05; note that the probability

128     thresholds bound to the fold change of 1.2) (Table 3). A significant proportion of these findings

129     were unique to MGI RNA-seq at EPO4 (66.8%) and Post7 (54.5%) (Supplementary Table 4).

130     Notably, substantial sub-proportions of the gene features identified from MGI RNA-seq

131     exceeded an absolute fold change of 2 (12.4% at EPO4 and 18.0% at Post7) and captured even

132     smaller changes between 1.2 and 2 (54.4% at EPO4 and 36.5% at Post7), when compared to the

133     Illumina RNA-seq and GeneChip[TM] gene features (ranging from 0% to 19.4%; Supplementary

134     Table 4 and Fig. 1). Furthermore, strong correlations between the two RNA-seq platforms on the

135     commonly identified genes were observed ($r = 0.74$ at EPO4 and $r = 0.85$ at Post7, $P < 2E-16$;

136     Fig. 2, a and d), whereas the correlations ranged from very weak ($r = 0.2$) to moderate ($r = 0.7$)

137     when compared RNA-seq to GeneChip[TM] ($P < 0.0003$; Fig. 2, b, c, e and f). Overall, MGI RNA-

138     seq, the total RNA DNB-seq, resulted in an increased sensitivity in identifying coding genes in

139     response to EPO compared to the Illumina mRNA-seq and GeneChip[TM] (Fig. 1, Fig. 2 and

140     Supplementary Data 1).

141

142     **Pathway analysis links the differentially expressed genes to erythropoiesis and non-**

143     **erythropoietic functions of EPO**

144     To explore the biological functions of the gene features identified from sequencing and

145     microarray, we performed a standard GSEA run (v4.0.3) subject to 1,000 phenotype

146     permutations[11, 12] on all datasets, using the MSigDB (v7.2)[11, 13] hallmark (H)[14] and Gene

147     Ontology (C5; BP: GO biological process)[15, 16] collections of functional gene sets. As expected,

148     heme metabolism emerged as the most significantly enriched pathway in all datasets following

149  the analysis on the 50 hallmark gene sets (FDR: MGI = 0.011 EPO4, Illumina = 0.033 EPO4,

150  GeneChip ≤ 0.017 EPO4/Post7 and BeadChip ≤ 0.004 EPO3/4/5/Post7/8; Supplementary Table

151  5). Leading edge genes, those contributing the most to the enrichment score of the heme

152  metabolism pathway constituting 200 genes, included 144 (MGI; EPO4), 105 (Illumina; EPO4),

153  125/96 (GeneChip; EPO4/Post7) and 101/103/103/97/84 (BeadChip; EPO3/4/5/Post7/8) genes

154  found in these datasets (Supplementary Data 2). Fifty-six leading edge genes overlapped across

155  all platforms and across conditions (pathway FDR < 0.1) (Supplementary Data 2). Of the 56

156  genes, 51 and 34 genes were also identified by the standard DGE analyses for the EPO4 and

157  Post7 conditions, respectively, across two or three of the MGI RNA-seq, Illumina RNA-seq and

158  GeneChip[TM] platforms (Supplementary Data 3). GSEA was able to detect the associated genes

159  that have fallen off the detection thresholds in the standard DGE analyses of Illumina RNA-seq

160  and GeneChip datasets (Supplementary Data 3). In addition, subsets of 10, 51, 51, 36, and 19 of

161  the 56 leading edge genes were found in the BeadChip DGE results across EPO3, EPO4, EPO5,

162  Post7 and Post8 conditions, respectively (Supplementary Data 4). The data suggest the

163  effectiveness of all four detection platforms and the effectiveness of GSEA in capturing the most

164  context-relevant biological pathway in response to rHuEPO. Next, GSEA was conducted on

165  7,530 GO biological processes included in the MSigDB C5 collection, and identified a total of

166  212, 134, and 33 biological pathways from MGI RNA-seq (EPO4), GeneChip (EPO4) and

167  BeadChip (EPO4, EPO5, Post7 and Post8) datasets, respectively, exceeding the pathway FDR <

168  0.1 and nominal $P < 0.05$. No significantly enriched GO biological processes were identified

169  from GSEA in the Illumina RNA-seq datasets. From the MGI RNA-seq dataset, these included

170  biological processes, resembling EPO cytoprotective functions and the downstream signal

171  transduction pathways[17-19], typically involved in response to oxidative stress (e.g. positive

172     regulation of mitophagy, hydrogen peroxide metabolic process, and nucleotide-excision repair,

173     DNA damage recognition), heme formation (e.g. porphyrin-containing compound metabolic

174     process), erythrocyte development, mTOR (target of rapamycin) signaling, regulation of energy

175     metabolism (e.g. regulation of generation of precursor metabolites and energy), low density

176     lipoprotein clearance, and nervous system development (Fig. 3). Key pathways characterising the

177     responses to EPO, such as autophagy of mitochondrion, positive regulation of cell cycle arrest,

178     iron ion homeostasis, tetrapyrrole metabolic process, erythrocyte development, and ventricular

179     system development also were identified from the GeneChip dataset (Supplementary Fig. 8). In

180     addition, other biological processes, including cyclic GMP mediated signaling, positive

181     regulation of cardiac muscle cell proliferation, and gamma-aminobutyric acid transport, were

182     observed, to name a few (Supplementary Fig. 8). In the BeadChip datasets, particular pathways

183     identified that were common to those observed on both the MGI RNA-seq and the GeneChip

184     platforms included hemoglobin metabolic process, erythrocyte development, and hydrogen

185     peroxide metabolic process (Supplementary Fig. 9). Further pathways of negative regulation of

186     necrotic cell death, negative regulation of TORC1 signaling, cellular response to monoamine

187     stimulus, monoamine transport, gas transport, lipid transport, drug transmembrane transport,

188     synaptic signaling, synapse organisation, and multicellular organism development, were found in

189     the BeadChip datasets (Supplementary Fig. 9). Leading edge genes from top 34, 14, and 16

190     pathways defined by the normalised enrichment score (NES) > 1.90 and from 38, 66, and 12 the

191     most enriched pathways representing a biological theme where NES < 1.90 were further

192     investigated in the MGI RNA-seq, GeneChip and BeadChip datasets, respectively. Out of a total

193     of 308 leading edge genes identified in the MGI RNA-seq EPO4 dataset overlapping with the

194     DESeq2 EPO4 DGE results, 135 genes also were found to be significantly expressed in the Post7

195      condition following the DESeq2 analysis (Supplementary Data 5). Of the 135 genes, top 10

196      genes filtered based on the GSEA ranks and pathway NESs — *BPGM*, *ALAS2*, *PKD1L3*,

197      *SLC4A1*, *AP2A1*, *IGF2*, *FAM210B*, *DYRK3*, *FECH* and *SLC25A37* — characterise erythrocyte

198      development, heme formation, metal ion homeostasis, cellular response to PH, LDL particle

199      clearance, glucose and energy metabolism, and TOR signaling (Supplementary Data 5). Fifty-

200      seven leading edge genes of the GeneChip EPO4 dataset were common to the limma EPO4 DGE

201      genes, while 15 (of the 57) were also present in the Post7 DGE results (Supplementary Data 6).

202      These 15 genes — *ALAS2*, *SLC4A1*, *FOXO3*, *TMOD1*, *FECH*, *SLC6A8*, *SLC25A39*, *SNCA*,

203      *FAM210B*, *EPB42*, *SLC25A37*, *YBX3*, *BPGM*, *STRADB*, and *BCL2L1* — correlate with heme

204      formation, bicarbonate transport, muscle atrophy, lens fiber cell development, gamma-

205      aminobutyric acid transport, erythrocyte development, cellular hyperosmotic response, negative

206      regulation of signal transduction in the absence of ligand and cellular response to amino acid

207      stimulus (Supplementary Data 6). Among 376 leading edge genes identified from the BeadChip

208      datasets, 76 also were observed in the limma DGE analysis. Top 10 genes (of the 76 genes) —

209      *KCNJ10*, *YBX3*, *SNCA*, *OR2W3*, *IRX1*, *OR2W5*, *CAMK2A*, *ACP4*, *NCDN* and *HOXC10* — are

210      involved in regulation of neuronal synaptic plasticity and necrotic cell death, sensory perception

211      of smell, proximal/distant pattern formation, and cell fate specification, and were enriched in the

212      GSEA Post7 dataset as well as were significantly expressed across EPO4, EPO5, Post7 and

213      Post8 conditions following the limma DGE analysis (Supplementary Data 7). Among the above

214      135, 15 and 76 leading edge genes identified on the three platforms, *BPGM*, *ALAS2*, *SLC4A1*,

215      *FAM210B*, *EPB42*, *SNCA*, *YBX3* and *TMOD1* were detected by all three platforms

216      (Supplementary Data 8). *FECH*, *SLC25A37*, *FOXO3*, *BCL2L1*, and *SLC25A39* were common

217      between MGI RNA-seq and GeneChip, *SLC6A8* between GeneChip and BeadChip, and *SLC7A5*,

218  *PINK1*, *DMTN*, *TRIM58*, *SESN3*, *GATA1*, *FURIN*, *HBQ1*, *EIF2AK1*, and *HBM* between MGI

219  RNA-seq and BeadChip (Supplementary Data 8). One hundred and twelve leading edge genes

220  (top 5: *PKD1L3*, *AP2A1*, *DYRK3*, *IGF2* and *TAL1*) were uniquely identified by MGI RNA-seq, 1

221  (*STRADB*) by GeneChip and 57 by BeadChip (top 5: *KCNJ10*, *OR2W3*, *IRX1*, *OR2W5* and

222  *CAMK2A*) (Supplementary Data 8). Further, 43 leading edge genes identified from one or more

223  of the three platforms were detected by Illumina RNA-seq across EPO4 and Post7 conditions

224  following the DESeq2 DGE analysis (Supplementary Data 9).

225

226  To follow up on the DGE and GSEA results, we performed additional analysis using the

227  Reactome database to examining the pathway components inferred from the 43 genes in pathway

228  diagrams and to confirming the gene functions attributed to rHuEPO across the experimental

229  conditions. By overlaying the gene expression values on Reactome pathway diagrams (release

230  73)[20], 13 and 8 significantly expressed interaction networks represented by 29 and 13 of the 43

231  genes, or their interactors (IntAct score $\geq$ 0.556), were identified in the MGI RNA-seq and

232  BeadChip datasets, respectively (pathway FDR < 0.05; see Supplementary Data 10 for pathway

233  entities and statistics and Supplementary Data 11 for the corresponding pathway overviews).

234  Notably, pathway components in the entire cascade of $O_2/CO_2$ exchange in erythrocytes were the

235  most significantly altered in the MGI RNA-seq datasets as opposed to findings obtained from the

236  Illumina RNA-seq, GeneChip and BeadChip platforms, including the pathway genes *SLC4A1*,

237  *HBB*, *CA1*, *AQP1*, *RHAG*, *HBA1* and *CYBSR1* (pathway FDR $\leq$0.003, Supplementary Data 10)

238  across EPO4, up-regulation and Post7, down-regulation (see Fig. 4 for the enhanced high-level

239  pathway diagram of the Post7 dataset). Finally, by overlapping a total of 172 and 91 significantly

240  expressed Reactome pathway genes and their interactors (IntAct score > 0.9 of high confidence

241      interactions, Supplementary Data 10) emerged from the 13 and 8 networks with the genes

242      identified from the standard DGE analyses in the MGI RNA-seq and BeadChip datasets, 80 and

243      41 genes were further confirmed, respectively (Supplementary Data 12). These 80 and 41 genes

244      represent the candidate genes that warrant further studies to rule out potential confounding

245      factors that mimic the EPO effect in terms of developing robust anti-doping gene signatures, or

246      to verify the role of the genes in EPO production and function for therapeutic purposes. The

247      subsets of the top 10 genes (sorted by the standard DGE $s$-value $< 0.005$ or FDR $< 0.05$)

248      accompanied by their corresponding GSEA and Reactome pathways are presented in

249      Supplementary Table 6 and 7.

250

251      **Discussion**

252      Taken together, cross-platform comparison in 10 subjects administered rHuEPO (50 IU every 2

253      days for 4 weeks) was conducted following gene expression quantification on MGI DNBSEQ-

254      G400RS, Illumina NextSeq 500 and GeneChip$^{TM}$ HTA2.0 platforms. To initiate a direct

255      comparison, only the coding gene features were extracted and compared across platforms. There

256      was a 2.28-fold increase in genes significantly expressed following MGI RNA-seq, as compared

257      to the combined number of genes identified on the other two platforms (Fig. 1 and

258      Supplementary Table 4). Furthermore, among 1,126 genes identified at EPO4, 25.5% of the

259      genes overlapped between MGI RNA-seq and the other two platforms, and 66.8% of the genes

260      were unique to MGI RNA-seq; among 674 genes identified at Post7, the corresponding figures

261      were 21.5% and 54.5%, respectively (Supplementary Table 4). Among genes with an absolute

262      fold change less than 2, Illumina RNA-seq captured a much higher proportion of the identified

263      gene features compared to GeneChip (10.3% vs 1.7% EPO4; 26.4% vs 0.6% Post7;

264     Supplementary Table 4). The experimental effect of EPO was largely captured by MGI RNA-seq

265     (PC1: 69% vs PC2: 5%), followed by GeneChip (PC1: 58% vs PC2: 14%) and Illumina RNA-

266     seq (PC1: 44% vs PC2: 9%) by examining the top 500 genes showing the highest variability

267     across samples (Supplementary Fig. 6, A to C). These observations support the supreme

268     performance of total RNA DNB-seq on MGI DNBSEQ-G400RS, followed by mRNA-seq on

269     Illumina NextSeq 500 and GeneChip[TM] HTA2.0 in this study. Nevertheless, genes characterised

270     by Illumina HumanHT-12 v4 Expression BeadChip in the 18 subjects represented a total of 85%

271     of variance captured by PC1 (78%) and PC2 (7%) (Supplementary Fig. 6D), suggesting

272     increased statistical power owing to the larger sample size (i.e. 143 samples from 18 subjects in

273     contrast to 50 samples from 10 subjects comprising the sample sets analysed on the other

274     platforms). Following on, quantitative pathway analysis by GSEA identified the heme

275     metabolism pathway enriched in all datasets across the four high-throughput gene quantification

276     platforms when analysing the MSigDB Hallmark collection of functional gene sets, while 212,

277     134, and 33 enriched pathways were unveiled from MGI RNA-seq, GeneChip and BeadChip

278     datasets, respectively, by examining the MSigDB C5 collection of 7,530 biological processes.

279     The pathway results underpinned the biological relevance of the gene expression findings,

280     particularly with a wealth of functional information emerged from the MGI RNA-seq dataset

281     (Fig. 3). Pathways of interest were prioritised by focusing on pathways of NES > 1.9 and

282     representative pathways of the biological themes where the NES < 1.9. Three hundred and eight,

283     57 and 376 leading edge genes were extracted from the pathways of interest, eventually led to

284     135, 15 and 76 genes also confirmed by the standard DGE analyses of MGI RNA-seq, GeneChip

285     and BeadChip datasets, respectively (Supplementary Data 5-7). Among these genes, 43 were

286     further validated in the list of genes resulted from the Illumina RNA-seq DGE analysis

287 (Supplementary Data 9). Despite strong positive correlations observed at the gene level between

288 MGI RNA-seq and Illumina RNA-seq, the lack of significantly expressed pathways following

289 GSEA in the Illumina RNA-seq datasets is in line with the generally weaker signals being picked

290 up by Illumina RNA-seq in this study (Fig. 2). To better understand the interacting networks or

291 signaling cascades represented by the 43 genes, we explored the Reactome database and

292 generated a total of 21 pathway overviews detailing the pathway entities/genes, their expression

293 levels, and their interactions with other entities within the pathway or across different pathways

294 (Supplementary Data 11). Finally, by extracting the significantly altered genes involved in these

295 Reactome networks and by matching these genes to the results of the standard DGE analyses, we

296 concluded with the lists of 80 and 41 genes that are of biological relevance to rHuEPO, identified

297 on the MGI RNA-seq and BeadChip platforms, respectively (Supplementary Data 12). They

298 represent the top biological pathways enriched in metabolism of porphyrins, $O_2/CO_2$ exchange in

299 erythrocytes, response to oxidative stress induced cellular senescence, and tissue damage caused

300 by amyloid deposition.

301

302 This comprehensive profiling of rHuEPO gene expression based on both RNA-seq and

303 microarrays has generated a robust set of genes of biological significance in relation to

304 erythropoiesis as well as non-erythropoietic effects of rHuEPO. It also establishes a knowledge

305 base of genes capturing a wide range of magnitude of changes attributable to rHuEPO by RNA-

306 seq, highlighting advantages of total RNA-seq combined with DNB-seq in quantifying gene

307 transcription. The adoption of a data triangulation approach by cross-platform comparisons and

308 by linking genes to their functions reinforces the biological findings and mitigates gene

309 expression perturbations caused by normal physiological changes such as seasonal changes and

310     lifestyle related changes. The longitudinal nature of the current investigation in healthy

311     individuals would help facilitate detailed studies of erythroid disorders and help formulate target

312     therapeutics, through disrupting and examining the mechanisms of the putative genes involved in

313     erythropoiesis and non-erythropoietic functions of rHuEPO. Finally, this study underpins the

314     follow-up studies needed to develop sensitive and robust gene signatures of blood doping in

315     sport.

316

317     **Methods**

318     **Subjects**

319     In a previously funded research project by the World Anti-Doping Agency (grant no.:

320     08C19YP), we collected whole blood samples from 18 endurance-trained Caucasian males at sea

321     level from Glasgow, Scotland (26.0±4.5 yrs, 74.8±7.9 kg, 179.8±5.4 cm), who underwent 4-

322     week 50 IU·kg$^{-1}$ body mass of rHuEPO every second day[3]. Daily oral iron supplementation (100

323     mg of elemental iron, ferrous sulphate tablets, Almus, Barnstable, UK) was given during the 4

324     weeks of rHuEPO administration[3]. Whole blood samples were collected at baseline (2 weeks and

325     1 day before rHuEPO; denoted by Base1 and Base2, respectively), during the rHuEPO

326     administration (2 days, 2 and 4 weeks following the 1st injection; denoted by EPO3, EPO4 and

327     EPO5, respectively) and for 4 weeks after the rHuEPO administration (1, 2 and 4 weeks after the

328     last injection; denoted by Post6, Post7 and Post8, respectively) for gene expression profiling on

329     the HumanHT-12 v4.0 Expression BeadChip (Illumina, San Diego, CA, USA)[4]. In the current

330     study (grant no.: ISF15E10YP), samples from 10 out of the 18 subjects collected at Base1,

331     Base2, EPO3, EPO4 and Post7 were analysed on a new microarray platform (GeneChip$^{TM}$

332     Human Transcriptome Array 2.0 or HTA2.0, Thermo Fisher Scientific, Waltham, MA, USA)

333　　and on two RNA-seq platforms (NextSeq500, Illumina, San Diego, CA, USA, and DNBSEQ-

334　　G400RS, MGI Tech, Shenzhen, China) for cross-platform gene expression comparisons for

335　　robust detection of EPO gene signatures. The studies were approved by the University of

336　　Glasgow Ethics Committee (Scotland, UK) and the University of Brighton Ethics Committee

337　　(England, UK) and were performed in accordance with the "Declaration of Helsinki". Written

338　　informed consent was obtained from all subjects.

339

340　　**RNA collection and preparation**

341　　Three milliliters of whole blood was collected from an antecubital vein using Tempus™

342　　Blood RNA tubes (Thermo Fisher Scientific, Waltham, MA, USA). Each Tempus™ tube

343　　contains 6 mL of RNA stabilising reagent and was vigorously mixed immediately after

344　　collection for 10 s. The blood samples were incubated at room temperature for approximately 3

345　　hours and then stored at −20°C or −80°C before subsequent analysis or transportation to the

346　　analytical lab. Total RNA was isolated from the whole blood according to the manufacturer's

347　　instructions (Tempus™ Spin RNA Isolation Kit, Thermo Fisher Scientific, Waltham, MA,

348　　USA). The purified total RNA was eluted in 90 μL elution buffer and stored in three aliquots at

349　　−80°C until further analysis. Initial RNA quantity and purity was assessed by the Nanodrop™

350　　ND-2000 Spectrophotometer (Thermo Fisher Scientific, Waltham, MA, US). RNA integrity was

351　　assessed using the Agilent 2100 Bioanalyser (Agilent Technologies, Santa Clara, CA, USA)

352　　prior to the RNA-seq and GeneChip analyses.

353

354　　**Microarray analysis with HumanHT-12 v4.0 Expression BeadChip**

355    Detailed sample preparation for the Illumina microarray experiment are available elsewhere[4].

356    Briefly, 500 ng of total RNA was used for complimentary RNA (cRNA) synthesis using the

357    Illumina[TM] TotalPrep RNA Amplification Kit (Thermo Fisher Scientific, Waltham, MA, USA).

358    Seven hundred and fifty nanograms of the purified labelled cRNA samples were hybridised to

359    the HumanHT-12 v4.0 Expression BeadChip arrays containing > 47,000 probes, following the

360    manufacturer's recommended procedures (Illumina, San Diego, CA, USA). The Bead arrays

361    were scanned on the Illumina BeadArray Reader. In this current study, the raw intensity values

362    were exported using the Illumina GenomeStudio software (v2.0; Gene Expression Module). The

363    bioconductor "limma" package[10] was used for background correction, data normalisation (using

364    the "neqc" function)[21] and differential gene expression analysis (DGE)[22] for paired samples

365    (using the "treat" function) in the 18 subjects across all 8 time points (i.e. Base1, Base2, EPO3,

366    EPO4, EPO5, Post6, Post7 and Post8). Notably, only probes expressed in at least 7 samples at a

367    detection $p<0.05$ were kept. Probes were annotated to illuminaHumanv4.db[23] and only probes

368    with "good" and "perfect" matching quality were retained followed by removing probes with

369    "NA" or multiple mappings.  Probes with low expression values below 5.1 were excluded prior

370    to the DGE analysis (assessed using the limma "plotSA" function). Transcripts were considered

371    significantly expressed for a fold change of 1.2 bounded to a 5% false discovery rate (FDR)

372    (thereby, giving more weight to fold change for gene ranking). These are common cut-off values

373    being used for declaring biologically and statistically significant findings in a DGE analysis[24].

374

375    **Microarray analysis with GeneChip[TM] HTA2.0**

376    One hundred nanograms of total RNA was processed using the GeneChip[TM] WT Plus Reagent

377    Kit according to the manufacturer's instructions (Thermo Fisher Scientific, Waltham, MA, US)

378      for 10 out of the 18 subjects at the selected time points (i.e. Base1, Base2, EPO3, EPO4 and

379      Post7). Single-stranded cDNA (ss-cDNA) was synthesised by the reverse transcription of cRNA.

380      Two hundred microlitres of hybridisation cocktail (containing approximately 5.2 μg fragmented

381      and labelled ss-cDNA) was loaded onto the GeneChip[TM] HTA2.0 (Thermo Fisher Scientific,

382      Waltham, MA, US). The GeneChip[TM] arrays were incubated in the GeneChip[TM] Hybridization

383      Oven 645 for 16 hours, washed and stained on the GeneChip[TM] Fluidics Station 450. The arrays

384      were then scanned using the GeneChip[TM] Scanner 3000 7G. The Applied Biosystems[TM]

385      Transcriptome Analysis Console (version:4.0.1.36; Thermo Fisher Scientific, Waltham, MA,

386      US) was used to perform initial data QC and data visualisation. The relative log expression box

387      plots were plotted following the quality assessment steps illustrated in *ref*[25]. The Bioconductor

388      "oligo" package[26] was used to read in the raw intensity CEL files, and the "rma" function was

389      used for background correction, normalisation, and data summarisation to the gene level (defined

390      by the argument "core"). Probes were annotated to hta20transcriptcluster.db[27] and probes with

391      "NA" or multiple mappings were removed. The "limma" package was then used to perform the

392      usual DGE analysis for paired samples (the analysis setting is identical to that used in the

393      Illumina microarray analysis illustrated above). Transcript clusters (loosely equal to genes) were

394      considered significantly expressed at a fold change of 1.2 bounded to a 5% FDR.

395

**396      RNA-seq on Illumina NextSeq500**

397      Five hundred nanograms of total RNA was used for sequencing according to the Illumina TruSeq

398      Stranded mRNA sample prep guide - high sample protocol (Illumina, San Diego, CA, USA).

399      Briefly, mRNA molecules were purified using the poly-T oligo attached magnetic beads

400      following which the mRNA was fragmented and primed for cDNA synthesis. A single "A" base

| 401 | was subsequently added to the 3-prime end of the synthesised blunt-ended cDNA and ligated |
| 402 | with index adapters for hybridisation onto a flow cell. The DNA fragments with adapters on both |
| 403 | ends were amplified via polymerase chain reaction to generate the final double-stranded cDNA |
| 404 | (ds-cDNA) library followed by library validation and normalisation and pooling of the samples. |
| 405 | Samples were pooled and then sequenced at 2x75 bp read length to a depth of approximately 64 |
| 406 | M reads per sample on the Illumina NextSeq 500 (Illumina, San Diego, CA, USA). 10 out of the |
| 407 | 18 subjects at the selected time points (i.e. Base1, Base2, EPO3, EPO4, and Post7) were |
| 408 | analysed. Raw sequences were examined by FastQC[28] for basic quality checks (e.g. per base |
| 409 | sequence quality, adaptor content, and per base N content), FastQ Screen[29] for mapping against |
| 410 | multiple reference genomes for detecting sample swaps or sample contamination that may have |
| 411 | resulted from sources other than humans (i.e. in this case, mapping against human, |
| 412 | mouse and rat genomes were conducted), HISAT2[5] for alignment to the reference genome |
| 413 | assembly (GRCh38.p12[6]) using the Ensembl 94 annotation[30] prior to RseQC[7] for read |
| 414 | distribution analysis, Salmon[8] for aligning to the transcriptome and transcripts quantification |
| 415 | (using selective alignment with the *decoy aware* target transcriptome to eliminate potential |
| 416 | spurious mapping to unannotated genomic locus over a *k*-mer length of 31, along with --SeqBias |
| 417 | and --gcBias flags switched on to correct for any unwanted effects), bioconductor package |
| 418 | "tximport"[31] for summarising transcript-level estimates to genes based on Ensembl release 94[30], |
| 419 | and DESeq2[9] for paired sample DGE analysis. Pre-filtering was performed to keep genes that |
| 420 | have at least 5 reads across 4 samples prior to the DGE analysis. Ensembl IDs were mapped to |
| 421 | gene symbols using the bioconductor package "org.Hs.eg.db"[32] and un-defined mappings were |
| 422 | removed (i.e. gene with "NA" or multiple mappings). MultiQC[33] was used to aggregate the |
| 423 | analysis results from the FastQC, FastQ Screen and RseQC runs from multiple samples. |

424 Unsupervised principal component analysis (PCA) for top 500 genes of high variance and gene

425 clustering analysis for the top 30 genes were performed following the DESeq2 vignette on data

426 quality assessment procedures[34]. The bioconductor package "SVA"[35] was used to assess

427 surrogate variables that may represent other variations in the data for further correction.

428 Shrinkage estimator "apeglm" was used for the shrinkage of log fold change estimates and for

429 ranking genes by effect size[36]. Genes exceeding a fold change of 1.2 bounded to the default $s$-

430 value < 0.005 were reported.

431

432 **RNA-seq on MGI DNBSEQ-G400RS**

433 Four hundred nanograms of total RNA was used for sequencing on the MGI DNBSEQ-G400RS

434 instrument (MGI, Shenzhen, China). Total RNA was first treated with Globin-Zero Gold Kit

435 (Illumina, San Diego, CA, USA) for rRNA depletion and globin mRNA reduction. The ds-

436 cDNA library preparation is in line with the Illumina RNA-seq protocol described in the above

437 section. The ds-cDNAs were then heat denatured and circularised by the splint oligo sequence to

438 generate the single strand circle DNA followed by rolling circle replication to create DNA

439 nanoballs (DNB) for processing on the MGI DNBSEQ-G400RS. The same 50 samples used for

440 the GeneChip[TM] and Illumina RNA-seq profiling were again analysed on this platform. These

441 samples were sequenced on 6 flowcells at 2x100 bp read length aimed at a sequencing depth of

442 64 M reads. Raw sequences were processed for quality assessment, alignment, transcripts

443 quantification and DGE analysis as described above in the "**RNA-seq on Illumina NextSeq500**"

444 section. The same cut-offs as in the Illumina RNA-seq section for defining a significant result

445 were applied (i.e. a fold change of 1.2 bounded to the default $s$-value < 0.005).

446

**Gene set enrichment analyses in GSEA and Reactome**

The pathway enrichment analysis was performed in accordance with recommendations from *ref*

[37], where appropriate. Specifically, normalised RNA-seq counts (outputted from DESeq2

"counts" function with the argument "normalized=TRUE") and normalised microarray gene

expression values were subjected to gene set enrichment analysis using GSEA (v4.0.3)[11, 12] by

examining the Molecular Signatures Database (MSigDB)[11, 13] Hallmark (H; containing 50 gene

sets)[14] and Gene Ontology (C5; BP: subset of GO biological processes containing 7,573 gene

sets)[15, 16] collections of functional gene sets. Low count genes (by removing genes with counts

below 5 in at least 4 samples) and genes with unidentified mappings from RNA-seq, and control

probes, low-quality probes and probes with unidentified mappings from microarray analyses

were excluded from the expression datasets prior to the GSEA. A standard GSEA run was

applied for each dataset by performing 1,000 phenotype permutations and by collapsing the

Ensembl IDs and probe IDs to gene symbols by mapping to their corresponding chip platforms

available from the MSigDB database (i.e. Human_ENSEMBL_Gene_ID_MSigDB.v7.2.chip for

RNA-seq, Human_AFFY_hta_2_0_MSigDB.v7.2.chip for GeneChip[TM] HTA2.0 and

Human_Illumina_HumanHT_12_v4_Array_MSigDB.v7.2.chip for Illumina BeadChip). Other

main parameters used in a GSEA run included the default ranking metric "Signal2Noise", gene

set size filters (15-200 for H, and 10-500 for C5) and collapsing mode ("Sum_of_probes" for

RNA-seq, and "Max_probe" for microarray). Default values were used for other fields of the

GSEA run. EnrichmentMap App[38] was used for creating biological networks of the GSEA

pathways (pathway FDR<0.1, nominal P<0.05 and Jaccard Overlap coefficient >0.375 with

combined constant k=0.5) and AutoAnnotate App[39] for gene sets annotation and clustering

(MCL Cluster annotation) in Cytoscape (v3.8.0)[40]. The most significantly enriched gene set was

470      used to label a gene set cluster, characterised by the normalised enrichment score (NES). Raw

471      counts from the RNA-seq (outputted from DESeq2 "counts" function by setting

472      "normalized=FALSE"), and normalised and log2 transformed gene expression values from

473      microarray analyses were uploaded onto Reactome (v73)[20] for quantitative pathway analysis

474      (ReactomeGSA) using the PADOG algorithm[41, 42] for gene expression visualisation in pathway

475      diagrams. Protein-protein interactors derived from the IntAct database[43] with the IntAct score $\geq$

476      0.556 (of medium to high confidence interactions) were included in the analysis to improve the

477      Reactome pathway coverage.  For consistency, these expression datasets were collapsed to gene

478      symbols using the "Collapse Dataset" tool in the GSEA software prior to the ReactomeGSA.

479

480      **Cross-platform DGE comparison**

481      Direct comparisons for the coding gene features identified across MGI DNBSEQ-G400RS,

482      Illumina NextSeq 500 and GeneChip[TM] HTA2.0 platforms in the 10 subjects (comprised of 50

483      samples) were carried out on the differentially expressed genes following the formal

484      DESeq2/limma DGE analyses. A sankey diagram was plotted for visualisation of the DGE

485      results using the "ggalluvial" package[44]. The cross-platform correlations were computed using

486      the "ggscatter" function in the package "ggpubr"[45]. "ggplot2"[46] and "cowplot"[47] packages were

487      used for creating publication-quality figures, where appropriate. Leading edge genes from the

488      significantly expressed GSEA pathways (derived from MGI DNBSEQ-G400RS, GeneChip[TM]

489      HTA2.0 and HumanHT-12 v4.0 Expression BeadChip; including all pathways with the NES>1.9

490      or the representative pathway of a gene set cluster when the NES<1.9) were extracted and

491      compared to the DGE genes to generate common sets of genes identified by both the GSEA and

492      DGE analyses. These genes were then overlapped with the DGE results obtained from the

493     Illumina NextSeq 500 platform for confirmation. The interaction networks among pathway genes

494     were defined by expression overlay with the Reactome pathway diagrams, focusing on the

495     networks represented by the confirmed genes above. The final lists of genes were obtained by

496     extracting all significantly altered genes and their interactors involved in these Reactome

497     networks and by matching them back to the differentially expressed genes resulted from the

498     formal DESeq2/limma analyses.

499

500     **Data availability:** Note all data will be made available and deposited into appropriate

501     repositories at publication, including raw RNA-seq data, raw microarray data, and code required

502     to reproduce all analyses. Full data access may be provided to reviewers on request during

503     manuscript reviewing.

504

505

506

**Fig. 1. Sankey diagram showing the flow of the differentially expressed gene features stratified by platform, biological condition and absolute $\log_2$-transformed fold changes.** M/I/H: MGI RNA-seq/Illumina RNA-seq/HTA2.0; M/I: MGI RNA-seq/Illumina RNA-seq; M/H: MGI RNA-seq/HTA2.0; M: MGI RNA-seq; I: Illumina RNA-seq; and H: HTA2.0. abs(lfc): absolute $\log_2$-transformed fold change. The colour coded band represents a detection platform or a combination of the detection platforms. The

512    wider the band, the higher number of the identified features on a platform or across platforms. The x-axis represents the number of

513    identified features captured on each platform. Note, for M/I/H, M/I, and M/H, that biological magnitude of the features used for

514    stratification is based on the MGI RNA-seq DGE results. Thirty-four identified non-protein coding transcript clusters on the GeneChip

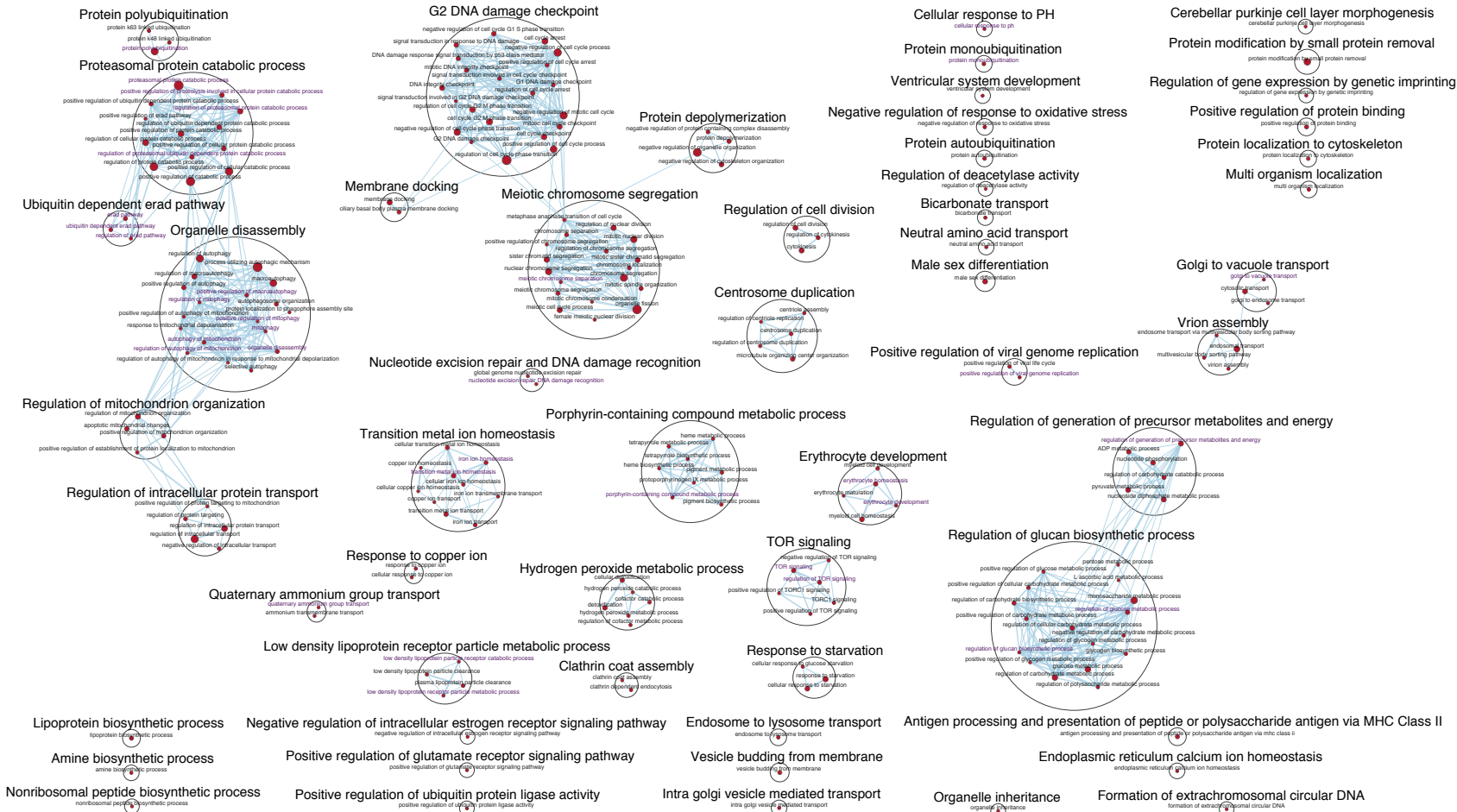515    are removed for the purposes of cross-platform comparison.

516



517

**Fig. 2. Cross-platform gene expression correlation analyses of log$_2$-transformed fold changes of all identified gene features. a-c** Genes identified when compared the level of expression between EPO4 and Base1 among the platform pairs in Illumina-MGI RNA-seq (**a**), GeneChip$^{TM}$ HTA2.0-MGI RNA-seq (**b**), GeneChip$^{TM}$ HTA2.0-Illlumina RNA-seq (**c**). **d-f** Genes identified when compared the level of expression between Post7 and Base1 among the

523    platform pairs in Illumina-MGI RNA-seq (**d**), GeneChip^(TM) HTA2.0-MGI RNA-seq (**e**),

524    GeneChip^(TM) HTA2.0-Illlumina RNA-seq (**f**). Genes identified as differentially expressed by

525    each pair are plotted in blue; genes that are only differentially expressed in Illumina RNA-seq,

526    MGI RNA-seq or GeneChip^(TM) HTA2.0 are plotted in yellow, grey and dijon, respectively; genes

527    not identified as differentially expressed by a pair are plotted in red. For simplicity, the

528    maximum expression value of a gene was used when multiple mapping of transcripts to the same

529    gene occurred.  *FOXO3B* is only differentially expressed in GeneChip^(TM) HTA2.0 when

530    compared to the MGI RNA-seq findings in (**b**), thus it has been removed from the correlation

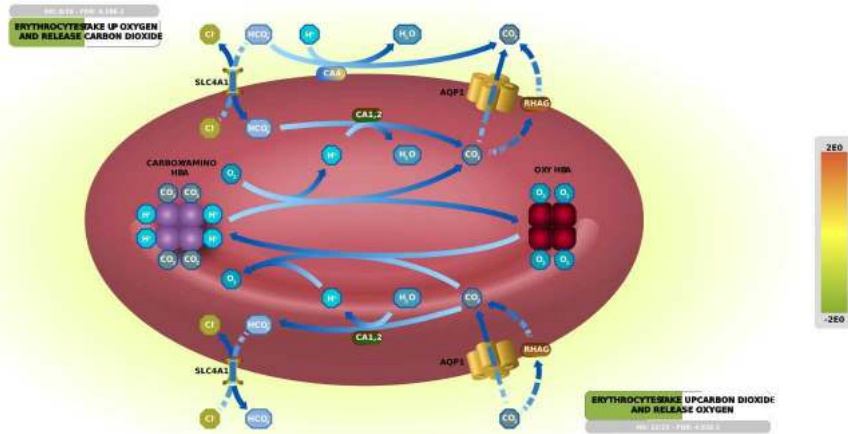531    analysis. R: Pearson's *r*. LogFC: $\log_2$-transformed fold change.

532

533



534

**Fig. 3. Biological network of the MGI RNA-seq dataset following Gene Ontology (biological process) gene set enrichment analysis in GSEA (v4.0.3) and visualisation in Cytoscape (3.8.0).** Each circle (node) represents a gene set and two nodes are

537    connected by lines (edges) indicating shared genes. The size of a node and width of an edge are proportional to the number of genes

538    enriched in a gene set and the number of genes shared between gene sets, respectively. Gene sets that are similar were annotated and

539    clustered to form a biological theme using the AutoAnnotate App in Cytoscape. The most significantly enriched gene set is used to

540    label a gene set cluster, defined by NES. Red node: gene set enriched in EPO4. Purple node label: top gene sets with NES > 1.90. The

541    enrichment map was created with pathway FDR < 0.1, nominal $P < 0.05$ and Jaccard Overlap coefficient > 0.375 with combined

542    constant k = 0.5.

543

544



545

546 **Fig. 4. Enhanced high-level Reactome pathway diagram for $O_2/CO_2$ exchange in**

547 **erythrocytes[48] by expression overlay with the MGI RNA-seq Post7 dataset.** This high-level

548 diagram represents two subpathways, namely erythrocyte take up oxygen and release carbon

549 dioxide and erythrocyte take up carbon dioxide and release oxygen. The green band indicates the

550 proportion of the pathway that is represented in the MGI RNA-seq Post7 dataset, and the colour

551 (green) represents the down-regulation of the pathway genes. The grey bar contains the

552 information for the number of pathway entities in the query dataset, the total number of the

553 pathway entities, and the FDR corrected over-representation probability.

554

**References:**

1. Shendure, J. et al., DNA sequencing at 40: past, present and future. *Nature*. **550**, 345-353 (2017).

2. Manolio, T. A. et al., Genomic Medicine Year in Review: 2019. *Am. J. Hum. Genet.* **105**, 1072-1075 (2019).

3. Durussel, J. et al., Haemoglobin mass and running time trial performance after recombinant human erythropoietin administration in trained men. *PLoS One*. **8**, e56151 (2013).

4. Durussel, J. et al., Blood transcriptional signature of recombinant human erythropoietin administration and implications for antidoping strategies. *Physiol. Genomics*. **48**, 202-209 (2016).

5. Kim, D., Paggi, J. M., Park, C., Bennett, C., Salzberg, S. L., Graph-based genome alignment and genotyping with HISAT2 and HISAT-genotype. *Nat. Biotechnol.* **37**, 907-915 (2019).

6. Schneider, V. A. et al., Evaluation of GRCh38 and de novo haploid genome assemblies demonstrates the enduring quality of the reference assembly. *Genome Res.* **27**, 849-864 (2017).

7. Wang, L., Wang, S., Li, W., RSeQC: quality control of RNA-seq experiments. *Bioinformatics*. **28**, 2184-2185 (2012).

8. Patro, R., Duggal, G., Love, M. I., Irizarry, R. A., Kingsford, C., Salmon provides fast and bias-aware quantification of transcript expression. *Nat. Methods*. **14**, 417-419 (2017).

9. Love, M. I., Huber, W., Anders, S., Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol.* **15**, 550 (2014).

10. Ritchie, M. E. et al., limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res.* **43**, e47 (2015).

578    11. Subramanian, A. et al., Gene set enrichment analysis: a knowledge-based approach for

579    interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci. U. S. A.* **102**, 15545-15550

580    (2005).

581    12. Mootha, V. K. et al., PGC-1alpha-responsive genes involved in oxidative phosphorylation

582    are coordinately downregulated in human diabetes. *Nat. Genet.* **34**, 267-273 (2003).

583    13. Liberzon, A. et al., Molecular signatures database (MSigDB) 3.0. *Bioinformatics*. **27**, 1739-

584    1740 (2011).

585    14. Liberzon, A. et al., The Molecular Signatures Database (MSigDB) hallmark gene set

586    collection. *Cell Syst.* **1**, 417-425 (2015).

587    15. Ashburner, M. et al., Gene ontology: tool for the unification of biology. The Gene Ontology

588    Consortium. *Nat. Genet.* **25**, 25-29 (2000).

589    16. The Gene Ontology Consortium. The Gene Ontology Resource: 20 years and still GOing

590    strong. *Nucleic Acids Res.* **47**, D330-d338 (2019).

591    17. Ghezzi, P., Brines, M., Erythropoietin as an antiapoptotic, tissue-protective cytokine. *Cell

592    Death Differ.* **11 Suppl 1**, S37-44 (2004).

593    18. Jelkmann, W., Regulation of erythropoietin production. *J. Physiol.* **589**, 1251-1258 (2011).

594    19. Maiese, K., Erythropoietin and diabetes mellitus. *World J. Diabetes*. **6**, 1259-1273 (2015).

595    20. Jassal, B. et al., The reactome pathway knowledgebase. *Nucleic Acids Res.* **48**, D498-d503

596    (2020).

597    21. Shi, W., Oshlack, A., Smyth, G. K., Optimizing the noise versus bias trade-off for Illumina

598    whole genome expression BeadChips. *Nucleic Acids Res.* **38**, e204 (2010).

599   22. Phipson, B., Lee, S., Majewski, I. J., Alexander, W. S., Smyth, G. K., Robust hyperparameter

600   estimation protects against hypervariable genes and improves power to detect differential

601   expression. *Ann. Appl. Stat*. **10**, 946-963 (2016).

602   23. Dunning, M., Lynch, A., Eldridge, M., illuminaHumanv4.db: Illumina HumanHT12v4

603   annotation data (chip illuminaHumanv4). R package version 1.26.0.  (2015).

604   24. Vaes, E., Khan, M., Mombaerts, P., Statistical analysis of differential gene expression

605   relative to a fold change threshold on NanoString data of mouse odorant receptor genes. *BMC*

606   *Bioinformatics*. **15**, 39 (2014).

607   25. Klaus, B., Reisenauer, S., An end to end workflow for differential gene expression using

608   Affymetrix microarrays. *F1000Res*. **5**, 1384 (2016).

609   26. Carvalho, B. S., Irizarry, R. A., A framework for oligonucleotide microarray preprocessing.

610   *Bioinformatics*. **26**, 2363-2367 (2010).

611   27. MacDonald, J. W., hta20transcriptcluster.db: Affymetrix hta20 annotation data (chip

612   hta20transcriptcluster). R package version 8.7.0.  (2017).

613   28. Andrews, S., FastQC:  A Quality Control Tool for High Throughput Sequence Data [Online].

614   Available online at: http://www.bioinformatics.babraham.ac.uk/projects/fastqc/.  (2010).

615   29. Wingett, S. W., Andrews, S., FastQ Screen: A tool for multi-genome mapping and quality

616   control. *F1000Res*. **7**, 1338 (2018).

617   30. Ensembl Archive Release 94 (October 2018). Available at:

618   http://oct2018.archive.ensembl.org/index.html.  (2018).

619   31. Soneson, C., Love, M. I., Robinson, M. D., Differential analyses for RNA-seq: transcript-

620   level estimates improve gene-level inferences. *F1000Res*. **4**, 1521 (2015).

621    32. Carlson, M., org.Hs.eg.db: Genome wide annotation for Human. R package version 3.8.2.,

622    (2019).

623    33. Ewels, P., Magnusson, M., Lundin, S., Käller, M., MultiQC: summarize analysis results for

624    multiple tools and samples in a single report. *Bioinformatics*. **32**, 3047-3048 (2016).

625    34. Love, M. I., Anders, S., Huber, W., Analyzing RNA-seq data with DESeq2. Available at:

626    [https://www.bioconductor.org/packages/devel/bioc/vignettes/DESeq2/inst/doc/DESeq2.html -](https://www.bioconductor.org/packages/devel/bioc/vignettes/DESeq2/inst/doc/DESeq2.html)

627    [references](references).  (2020).

628    35. Leek, J. T. et al., sva: Surrogate Variable Analysis. R package version 3.38.0.  (2020).

629    36. Zhu, A., Ibrahim, J. G., Love, M. I., Heavy-tailed prior distributions for sequence count data:

630    removing the noise and preserving large differences. *Bioinformatics*. **35**, 2084-2092 (2019).

631    37. Reimand, J. et al., Pathway enrichment analysis and visualization of omics data using

632    g:Profiler, GSEA, Cytoscape and EnrichmentMap. *Nat. Protoc.* **14**, 482-517 (2019).

633    38. Merico, D., Isserlin, R., Stueker, O., Emili, A., Bader, G. D., Enrichment map: a network-

634    based method for gene-set enrichment visualization and interpretation. *PLoS One*. **5**, e13984

635    (2010).

636    39. Kucera, M., Isserlin, R., Arkhangorodsky, A., Bader, G. D., AutoAnnotate: A Cytoscape app

637    for summarizing networks with semantic annotations. *F1000Res*. **5**, 1717 (2016).

638    40. Shannon, P. et al., Cytoscape: a software environment for integrated models of biomolecular

639    interaction networks. *Genome Res*. **13**, 2498-2504 (2003).

640    41. Tarca, A. L., Draghici, S., Bhatti, G., Romero, R., Down-weighting overlapping genes

641    improves gene set analysis. *BMC Bioinformatics*. **13**, 136 (2012).

642    42. Griss, J. et al., ReactomeGSA - Efficient Multi-Omics Comparative Pathway Analysis. *Mol.*

643    *Cell. Proteomics*.  (2020).

644     43. Orchard, S. et al., The MIntAct project--IntAct as a common curation platform for 11

645     molecular interaction databases. *Nucleic Acids Res.* **42**, D358-363 (2014).

646     44. Brunson, J. C., ggalluvial: Layered Grammar for Alluvial Plots. *J. Open Source Softw*. **5**, 49

647     (2017).

648     45. Kassambara, A., ggpubr: 'ggplot2' based publication ready plots. R package version 0.4.0.

649     (2020).

650     46. Wickham, H., Elegant graphics for data analysis. Springer-Verlag New York.  (2016).

651     47. Wilke, C. O., cowplot: streamlined plot theme and plot annotations for 'ggplot2'. R package

652     version 1.1.1.  (2020).

653     48. O2/CO2 exchange in erythrocytes. Reactome, released 2012-06-12,

654     doi:10.3180/REACT_120969.1 (11/11/20).

655

661 **Author contributions:**

662 Conceptualization: GW, YPP; Formal analysis: GW; Funding acquisition: YPP; Investigation:

663 GW, TK; Methodology: GW, YPP; Project administration: AC, QM; Resources: GIA, JL, MBG;

664 Supervision: YPP; Validation: GW; Visualization: GW; Writing – original draft: GW; Writing –

665 review & editing: AH, FMG, GIA, JL, MBG, YPP, GW.

666 **Competing interests:**

667 The authors declare no competing interests.

668 **Materials & Correspondence:**

669 Correspondence and material requests should be addressed to GW and YPP.

670

671

672 **Figure legends:**

673

674

675 **Fig. 1. Sankey diagram [44] showing the flow of the differentially expressed gene features**

676 **stratified by platform, biological condition and absolute $\log_2$-transformed fold changes.**

677 M/I/H: MGI RNA-seq/Illumina RNA-seq/HTA2.0; M/I: MGI RNA-seq/Illumina RNA-seq;

678 M/H: MGI RNA-seq/HTA2.0; M: MGI RNA-seq; I: Illumina RNA-seq; and H: HTA2.0.

679 abs(lfc): absolute $\log_2$-transformed fold change. The colour coded band represents a detection

680 platform or a combination of the detection platforms. The wider the band, the higher number of

681 the identified features on a platform or across platforms. The x-axis represents the number of

682 identified features captured on each platform. Note, for M/I/H, M/I, and M/H, that biological

683 magnitude of the features used for stratification is based on the MGI RNA-seq DGE results.

684 Thirty-four identified non-protein coding transcript clusters on the GeneChip are removed for the

685 purposes of cross-platform comparison.

686 **Fig. 2. Cross-platform gene expression correlation analyses of $\log_2$-transformed fold**

687 **changes of all identified gene features. a-c** Genes identified when compared the level of

688 expression between EPO4 and Base1 among the platform pairs in Illumina-MGI RNA-seq (**a**),

689 GeneChip^TM HTA2.0-MGI RNA-seq (**b**), GeneChip^TM HTA2.0-Illlumina RNA-seq (**c**). **d-f**

690 Genes identified when compared the level of expression between Post7 and Base1 among the

691 platform pairs in Illumina-MGI RNA-seq (**d**), GeneChip^TM HTA2.0-MGI RNA-seq (**e**),

692 GeneChip^TM HTA2.0-Illlumina RNA-seq (**f**). Genes identified as differentially expressed by

693 each pair are plotted in blue; genes that are only differentially expressed in Illumina RNA-seq,

694 MGI RNA-seq or GeneChip^TM HTA2.0 are plotted in yellow, grey and dijon, respectively; genes

695 not identified as differentially expressed by a pair are plotted in red. For simplicity, the

696 maximum expression value of a gene was used when multiple mapping of transcripts to the same

697    gene occurred.  FOXO3B is only differentially expressed in GeneChip[TM] HTA2.0 when

698    compared to the MGI RNA-seq findings in (**b**), thus it has been removed from the correlation

699    analysis. R: Pearson's r. LogFC: $\log_2$-transformed fold change.

700    **Fig. 3. Biological network of the MGI RNA-seq dataset following Gene Ontology (biological**

701    **process) gene set enrichment analysis in GSEA (v4.0.3) and visualisation in Cytoscape**

702    **(3.8.0) [40].** Each circle (node) represents a gene set and two nodes are connected by lines (edges)

703    indicating shared genes. The size of a node and width of an edge are proportional to the number

704    of genes enriched in a gene set and the number of genes shared between gene sets, respectively.

705    Gene sets that are similar were annotated and clustered to form a biological theme using the

706    AutoAnnotate App [39] in Cytoscape. The most significantly enriched gene set is used to label a

707    gene set cluster, defined by NES. Red node: gene set enriched in EPO4. Purple node label: top

708    gene sets with NES > 1.90. The enrichment map was created with pathway FDR < 0.1, nominal

709    $P < 0.05$ and Jaccard Overlap coefficient > 0.375 with combined constant k = 0.5.

710    **Fig. 4. Enhanced high-level Reactome pathway diagram for $O_2/CO_2$ exchange in**

711    **erythrocytes [48] by expression overlay with the MGI RNA-seq Post7 dataset.** This high-level

712    diagram represents two subpathways, namely erythrocyte take up oxygen and release carbon

713    dioxide and erythrocyte take up carbon dioxide and release oxygen. The green band indicates the

714    proportion of the pathway that is represented in the MGI RNA-seq Post7 dataset, and the colour

715    (green) represents the down-regulation of the pathway genes. The grey bar contains the

716    information for the number of pathway entities in the query dataset, the total number of the

717    pathway entities, and the FDR corrected over-representation probability.

718

719

720      **Table 1.** Summary of the number of transcriptomic features available for the DGE analysis

721      across the four gene expression detection platforms.

722

| | MGI DNBSEQ-G400RS | Illumina NextSeq 500 | GeneChip™ HTA2.0 | Illumina HumanHT-12 v4 Expression BeadChip |
|---|---|---|---|---|
| Number of samples | 50 | 48 | 49 | 143 |
| Number of transcriptomic features following RNA-seq quantification (Salmon) or on the array | 175,775 | 175,775 | 285,543 | 47,286 |
| Number of identified features available for the DGE analysis | 16,738[g] | 16,581[g] | 29,517[tc] | 10,622[t] |

723
724      DGE: differential gene expression. g: protein-coding gene features. tc: protein coding and non-
725      protein coding transcript clusters, loosely corresponding to genes. t: coding and non-coding
726      transcripts. Two, one and one samples were removed from the DGE analyses due to human
727      processing errors, sample quality and sampling issue in the Illumina RNA-seq, GeneChip and
728      BeadChip datasets, respectively.

**Table 2. Transcript annotation and filtering of the RNA-seq and microarray data prior to the DGE analysis.**

| | MGI DNBSEQ-G400RS | Illumina NextSeq 500 | GeneChip™ HTA2.0 | Illumina HumanHT-12 v4 Expression BeadChip |
|---|---|---|---|---|
| Annotation database (N=the number of transcriptomic features) | Org.Hs.eg.db (N=175,775 transcripts following Salmon transcription quantification, aggregated into 37,788 genes using Ensembl 94 annotation) | Org.Hs.eg.db (N=175,775 transcripts following Salmon transcription quantification, aggregated into 37,788 genes using Ensembl 94 annotation) | hta20transcriptcluster.db (N=285,543 transcripts, corresponding to 67,480 protein-coding and non-protein coding transcript clusters) | illuminaHumanv4.db (N=47,286 coding and non-coding transcripts) |
| Undetected probes | - | - | - | 18,494 |
| Low quality probes | - | - | - | 6,900 |
| Low-expressed genes (RNA-seq) | 17,198 | 18,347 | - | - |
| "NA" mapping to stable gene symbols | 3,675 | 2,668 | 36,709 | 2,406 |
| Multiple mapping to stable gene symbols | 177 | 192 | 1,254 | 698 |
| Low-expressed probes (microarray)[1] | - | - | 0 | 8,166 |
| Identified features available for DGE analysis | 16,738[g] | 16,581[g] | 29,517[tc] | 10,622[t] |

730

731 NA: features with no gene symbols returned after annotation. DGE: differential gene expression. -: not applicable. [1]: low expressed
732 probes were further removed following assessing the average log expression and the mean-variance relationship after fitting the linear

733      model in limma microarray analysis. g: protein-coding gene features. tc: protein-coding and non-protein coding transcript clusters,
734      loosely corresponding to genes. t: coding and non-coding transcript features.

735 **Table 3.** Summary of the number of significantly expressed transcriptomic features across all platforms.

736

| | DGE thresholds | Base2 vs Base1 | EPO3 vs Base1 | EPO4 vs Base1 | Post7 vs Base1 | Up/down regulation |
|---|---|---|---|---|---|---|
| MGI DNBSEQ-G400RS (N = 50) | abs FC > 1.2 & $s$-value < 0.005 | 0 | 1 | 959 | 60 | Up |
| | | 0 | 0 | 81 | 451 | Down |
| Illumina NextSeq 500 (N = 48) | abs FC > 1.2 & $s$-value < 0.005 | 0 | 0 | 277 | 27 | Up |
| | | 0 | 0 | 20 | 258 | Down |
| GeneChip$^{TM}$ HTA2.0 (N = 49) | abs FC > 1.2 & FDR < 0.05 | 0 | 0 | 200 | 0 | Up |
| | | 0 | 0 | 1 | 51 | Down |
| Illumina HumanHT-12v4.0 Expression Beadchip (N = 143) | abs FC > 1.2 & FDR < 0.05 | 0 | 13 | 796 | 7 | Up |
| | | 0 | 0 | 1,315 | 254 | Down |

737
738 DGE: differential gene expression. abs FC: absolute fold change. FDR: false discovery rate. The number of protein-coding gene
739 features, and coding and non-coding transcript clusters and transcripts are reported following the RNA-seq, GeneChip and Beadchip
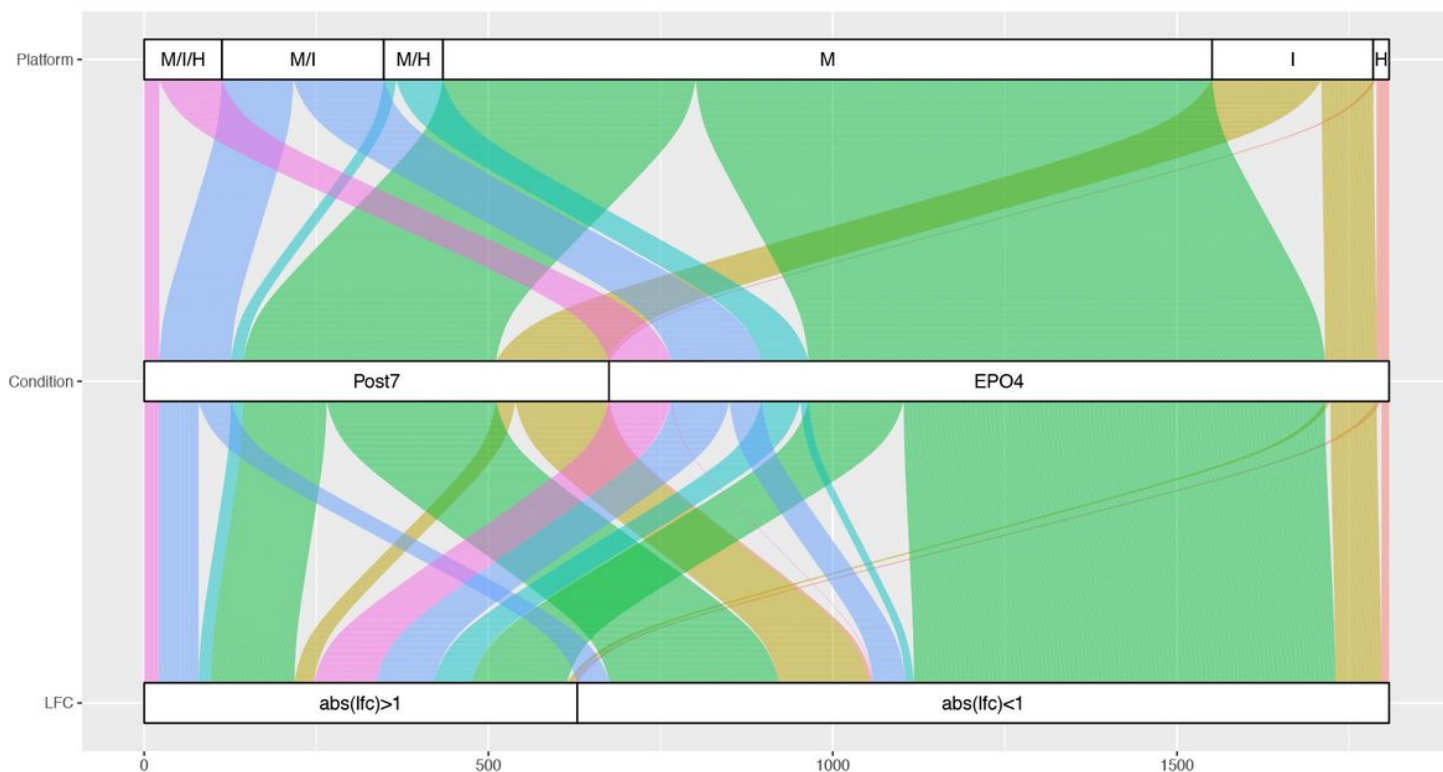740 DGE analyses, respectively.
741

# Figures



## Figure 1

Sankey diagram 44 showing the flow of the differentially expressed gene features stratified by platform, biological condition and absolute log2-transformed fold changes. M/I/H: MGI RNA-seq/Illumina RNA-seq/HTA2.0; M/I: MGI RNA-seq/Illumina RNA-seq; M/H: MGI RNA-seq/HTA2.0; M: MGI RNA-seq; I: Illumina RNA-seq; and H: HTA2.0. abs(lfc): absolute log2-transformed fold change. The colour coded band represents a detection platform or a combination of the detection platforms. The wider the band, the higher number of the identified features on a platform or across platforms. The x-axis represents the number of identified features captured on each platform. Note, for M/I/H, M/I, and M/H, that biological magnitude of the features used for stratification is based on the MGI RNA-seq DGE results. Thirty-four identified non-protein coding transcript clusters on the GeneChip are removed for the purposes of cross-platform comparison.
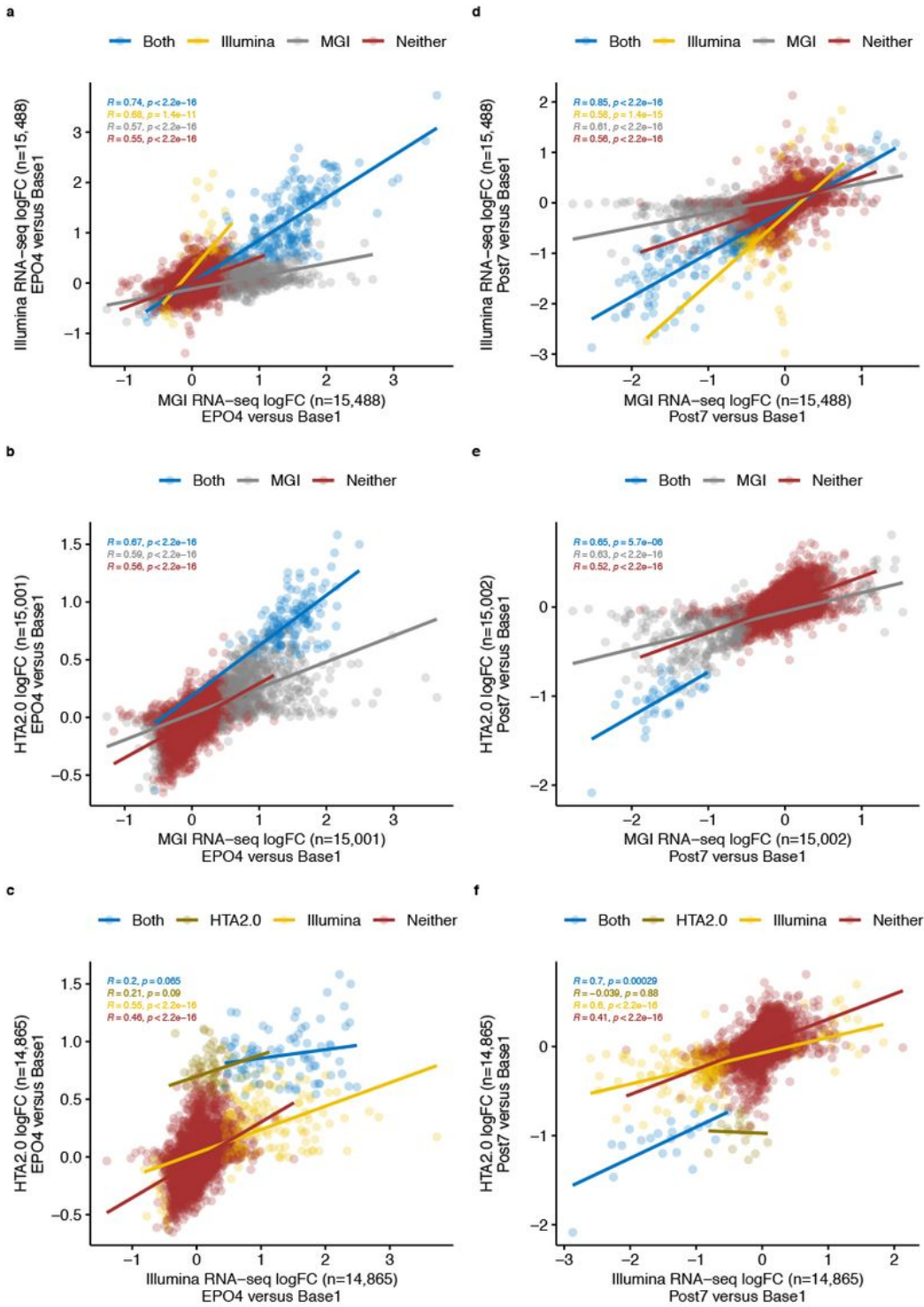
**Figure 2**

Cross-platform gene expression correlation analyses of log2-transformed fold changes of all identified gene features. a-c Genes identified when compared the level of expression between EPO4 and Base1 among the platform pairs in Illumina-MGI RNA-seq (a), GeneChipTM HTA2.0-MGI RNA-seq (b), GeneChipTM HTA2.0-Illlumina RNA-seq (c). d-f Genes identified when compared the level of expression between Post7 and Base1 among the platform pairs in Illumina-MGI RNA-seq (d), GeneChipTM HTA2.0-

MGI RNA-seq (e), GeneChipTM HTA2.0-IlIllumina RNA-seq (f). Genes identified as differentially expressed by each pair are plotted in blue; genes that are only differentially expressed in Illumina RNA-seq, MGI RNA-seq or GeneChipTM HTA2.0 are plotted in yellow, grey and dijon, respectively; genes not identified as differentially expressed by a pair are plotted in red. For simplicity, the maximum expression value of a gene was used when multiple mapping of transcripts to the same gene occurred. FOXO3B is only differentially e xpressed in GeneChipTM HTA2.0 when compared to the MGI RNA-seq findings in (b), thus it has been removed from the correlation analysis. R: Pearson's r. LogFC: log2-transformed fold change.
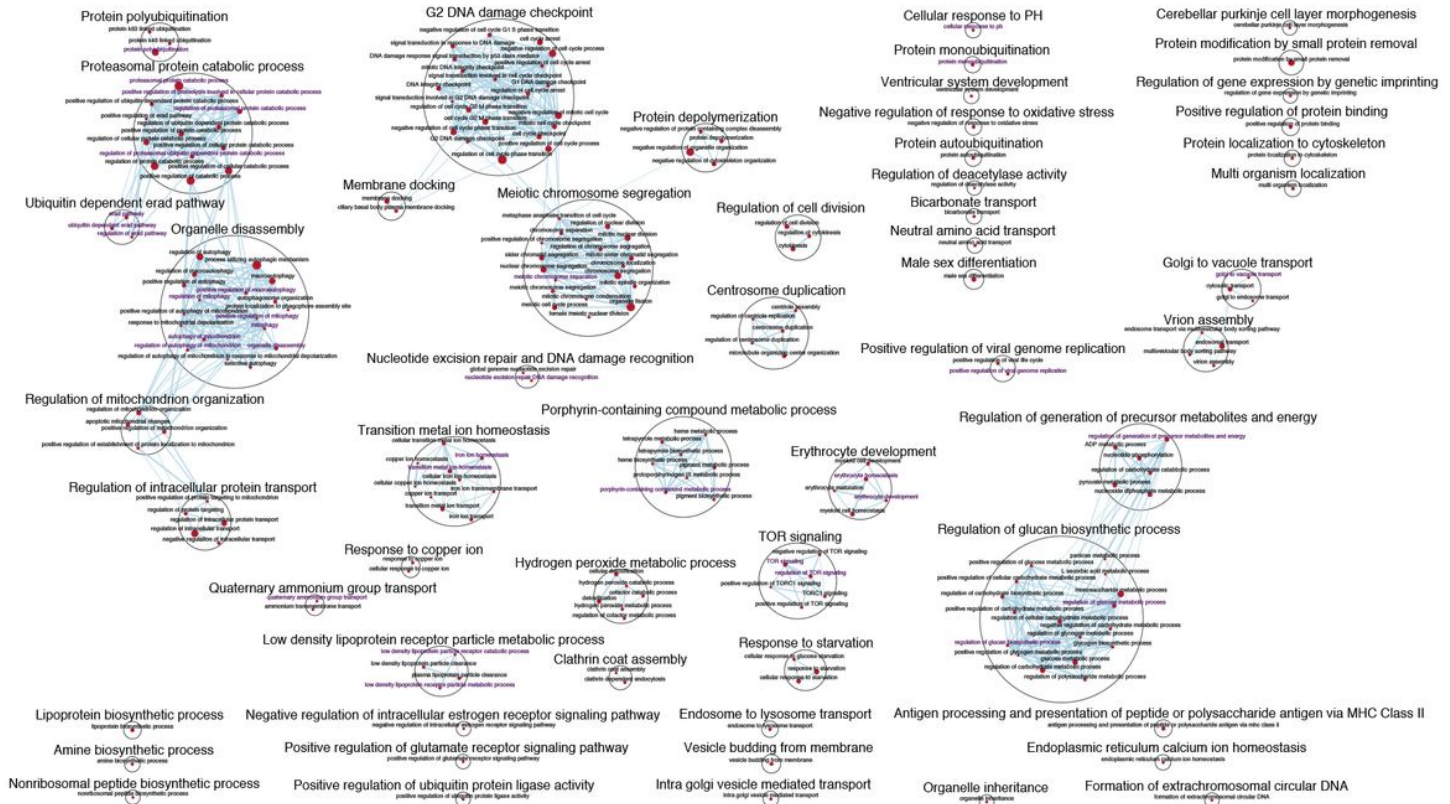


## Figure 3

Biological network of the MGI RNA-seq dataset following Gene Ontology (biological process) gene set enrichment analysis in GSEA (v4.0.3) and visualisation in Cytoscape (3.8.0) 40. Each circle (node) represents a gene set and two nodes are connected by lines (edges) indicating shared genes. The size of a node and width of an edge are proportional to the number of genes enriched in a gene set and the number of genes shared between gene sets, respectively. Gene sets that are similar were annotated and clustered to form a biological theme using the AutoAnnotate App 39 in Cytoscape. The most significantly enriched gene set is used to label a gene set cluster, defined by NES. Red node: gene set enriched in EPO4. Purple node label: top gene sets with NES > 1.90. The enrichment map was created with pathway FDR < 0.1, nominal P < 0.05 and Jaccard Overlap coefficient > 0.375 with combined constant k = 0.5.
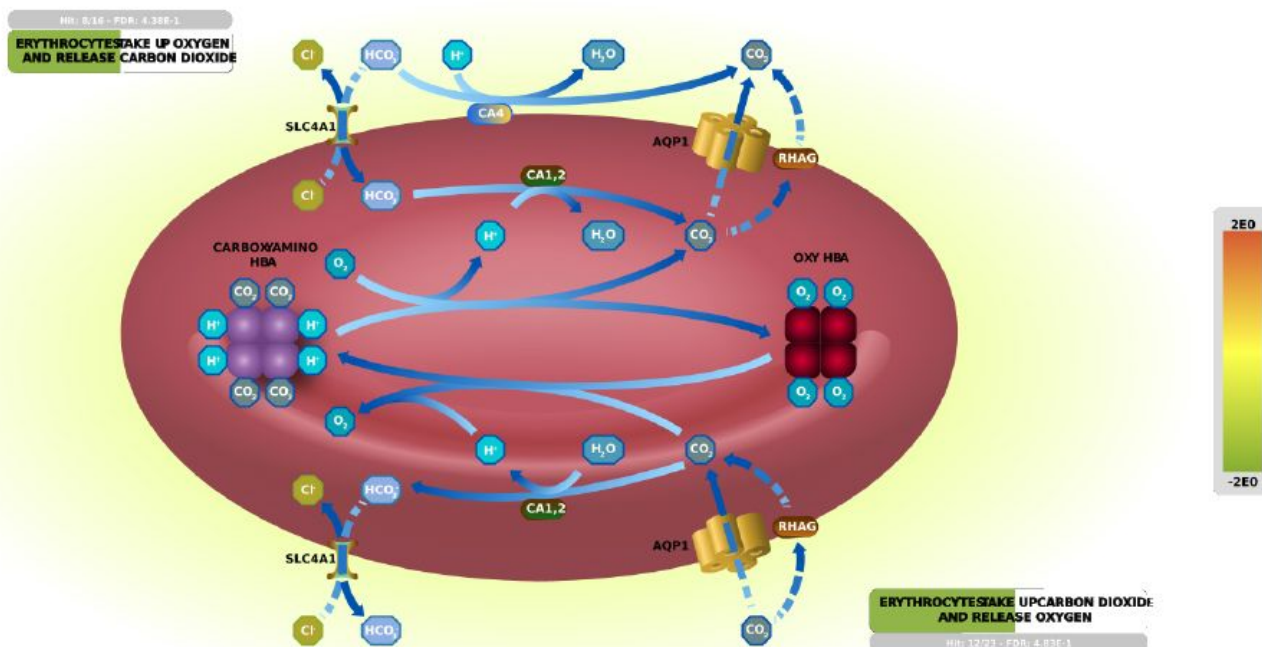
Figure 4

Enhanced high-level Reactome pathway diagram for O2/CO2 exchange in erythrocytes 48 by expression overlay with the MGI RNA-seq Post7 dataset. This high-level diagram represents two subpathways, namely erythrocyte take up oxygen and release carbon dioxide and erythrocyte take up carbon dioxide and release oxygen. The green band indicates the proportion of the pathway that is represented in the MGI RNA-seq Post7 dataset, and the colour (green) represents the down-regulation of the pathway genes. The grey bar contains the information for the number of pathway entities in the query dataset, the total number of the pathway entities, and the FDR corrected over-representation probability.

# Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- SupplementaryData112.zip
- WangetalSupplementaryInformation.pdf