# Divergent and convergent evolution of housekeeping genes in human-pig lineage

Kai Wei [1] , Tingting Zhang [1] , Lei Ma [Corresp. 1]

[1] College of Life Science, Shihezi University, Shihezi, Xinjiang, China

Corresponding Author: Lei Ma
Email address: malei1979@hotmail.com

Housekeeping genes are ubiquitously expressed and maintain basic cellular function across tissue/cell types conditions. The present study aimed to develop a set of pig housekeeping genes and compare characteristics of structure, evolution and function of housekeeping genes in the human-pig lineage. Using RNA sequencing data, we identified a list of 3,136 pig housekeeping genes. Comparing to human homologous counterparts, we found pig housekeeping genes were longer and subjected to slight weaker purifying selection pressure and faster neutral evolution. Common housekeeping genes, shared by the two species, have stronger purifying selection than species-specific genes. But pig-specific and human-specific housekeeping genes have similar functions. Some species-specific housekeeping genes have evolved independently to form similar protein-active sites or structure, such as classical catalytic serine-histidine-aspartate triad and zinc finger features, implying that they have converged for maintaining the basic cellular function, which led to equivalent solutions for adapting to the environment. Human and pig housekeeping genes have varied in their structure and gene list, but they have converged on the maintenance of basic cellular functions essential for the existence of a cell, regardless of its specific role in the species. The results shed light on the evolutionary dynamics of housekeeping genes.

1 **Divergent and convergent evolution of housekeeping genes in**
2 **human-pig lineage**
3

4 **Kai Wei †, Tingting Zhang †, Lei Ma** [*]
5

6 College of Life Science, Shihezi University, Shihezi City, Xinjiang Province, China

7

8 †These authors contributed equally to this work

9 *Corresponding author

10   Lei Ma: malei1979@hotmail.com

## Abstract

12    Housekeeping genes are ubiquitously expressed and maintain basic cellular function across

13    tissue/cell types conditions. The present study aimed to develop a set of pig housekeeping genes

14    and compare characteristics of structure, evolution and function of housekeeping genes in the

15    human-pig lineage. Using RNA sequencing data, we identified a list of 3,136 pig housekeeping

16    genes. Comparing to human homologous counterparts, we found pig housekeeping genes were

17    longer and subjected to slight weaker purifying selection pressure and faster neutral evolution.

18    Common housekeeping genes, shared by the two species, have stronger purifying selection than

19    species-specific genes. But pig-specific and human-specific housekeeping genes have similar

20    functions. Some species-specific housekeeping genes have evolved independently to form

21    similar protein-active sites or structure, such as classical catalytic serine-histidine-aspartate triad

22    and zinc finger features, implying that they have converged for maintaining the basic cellular

23    function, which led to equivalent solutions for adapting to the environment. Human and pig

24    housekeeping genes have varied in their structure and gene list, but they have converged on the

25    maintenance of basic cellular functions essential for the existence of a cell, regardless of its

26    specific role in the species. The results shed light on the evolutionary dynamics of housekeeping

27    genes.

28    **Keywords:** Housekeeping genes; Gene structure; Basal cellular function; Convergent evolution;

29    Pig

## Background

31    Housekeeping genes are typically genes consistently expressed across tissues and developmental

32    stages for the maintenance of basic cellular functions (Butte et al.2001; Zhu et al.2003). They

33    have unique genomic features, including gene structure (Eisenberg and Levanon 2003;

34   Vinogradov 2004), nucleotide composition (Vinogradov 2003), and upstream sequence

35   conservation (Farré et al.2007; Belloraet al.2007). They are often considered as the minimally

36   essential gene set for normal cellular physiology (Butte et al.2001) and are widely used as

37   internal controls for gene expression experiments as well as computational biology studies

38   (Thellin et al.1999; Robinson and Oshlack 2010;Rubie et al.2005; Vandesompele et al.2002).

39

40   In previous studies, many human housekeeping gene sets have been identified. However, some

41   sets have little overlap. For example, only 155 genes were shared by three lists of microarray-

42   defined housekeeping genes, including 501, 425 and 567 genes, respectively (Warrington et

43   al.2000; Hsiao et al.2001; Eisenberg and Levanon 2003). The low overlap may be explained by

44   several reasons. First, their complex transcriptional organization may cause diverse definitions of

45   housekeeping genes (Gingeras 2007). Second, the expression of some housekeeping genes may

46   vary depending on experimental conditions (Greer et al.2010). The question of why these genes

47   vary across conditions awaits further investigations. Third, traditional techniques have their own

48   drawbacks. For instance, the microarray technology has limited dynamic range and sensitivity,

49   and also suffers from poor detectability and reproducibility for low-copy and transiently-

50   expressed genes (Marioni et al.2008; Fu et al.2009; Bradford et al.2010; Draghici et al.2006).

51

52   RNA sequencing (RNA-seq) data greatly improve the detectability of housekeeping genes. For

53   example, the amount of human housekeeping genes revisited by the RNA-seq data has increased

54   ten-fold the previous estimates based on microarray data (Eisenberg and Levanon 2013). With

55   advances in technology, large-scale RNA sequencing has provided new insights into the

56   definition of housekeeping genes. Some studies have suggested that transcripts should be used as

57   housekeeping units (Gingeras 2007; Gerstein et al.2007).

58

59   The comparative analysis of housekeeping genes between human and other animals is of great

60   interest. Human housekeeping genes are commonly used as control genes in the real-time

61   quantitative polymerase chain reaction (qRT-PCR) for other animals. However, whether human

62   genes can be used as references for other animals remains unclear. For instance, the most

63   commonly used human reference genes (e.g. *ACTB* and *GAPDH*) do not always apply to all

64   tissues of different organisms (Brattelid et al.2010; Kozera et al.2013). Therefore, to well define

65   a housekeeping genes set in another animal may be valuable.

66

67   As an important meat resources for humans, the pig (*Sus Scrofa*) is a well-studied organism. And

68   because of anatomical similarities with humans, the pig is often used as a biomedical model in

69   research as well (Lunney 2007; Rolandsson et al.2002; Lee et al.2009; Becker et al.2010).

70   Surveying pig housekeeping genes may help pave the way for a greater understanding basal

71   mechanisms that maintain cell function. In the present study, we identified housekeeping genes

72   in pig using the RNA-seq data, and then compared their structure and function with human

73   orthologs. In addition, we discussed the impact of selection pressure and convergent evolution on

74   functional conservation of housekeeping genes. The present study provided detailed information

75   of pig housekeeping genes and their functional features, and offered insights into evolutionary

76   dynamics on them.

77

## Materials and Methods

### Data preparation

In order to define housekeeping gene sets, the gene expression datasets were downloaded from Sequencing Read Achieve (SRA) database of National Center for Biotechnology Information (NCBI, Sep, 2016) (Kodama et al.2012). In addition, pig genomic annotation (*Sus Sscrofa*10.2) was downloaded from the Ensembl Genome Browser (Sep, 2016) (Kinsellaet al.2011). The RNA-seq dataset of 14 experiments were used to identify housekeeping genes, which were derived from 21 tissues (heart, spleen, liver, kidney, lung, musculus longissimus dorsi, occipital cortex, hypothalamus, frontal cortex, cerebellum, endometrium, mesenterium, greater omentum, backfat, gonad, ovary, placenta, testis, blood, uterine and lymph nodes), containing a total of 131 samples(Supplementary material1: Table S1 ). The SRA files were downloaded from the NCBI and then converted to fastq files using fastq-dump (Kodama et al.2012). RNA-seq reads were then filtered by IlluQC.pl (Patel and Jain 2012) while requiring an average read quality above 20, and then were aligned to pig genome sequence (Sus Sscrofa10.2) using Tophat (Trapnell et al.2009; Külahoglu et al.2014; Ghosh S, Chan et al.2016). The alignments were then fed to an assembler Cufflinks (Trapnel et al.2010) to assemble aligned RNA-seq reads into transcripts and estimate their abundances, which were measured in Fragments Per Kilobase of exon per Million fragments mapped (FPKM).

### To define housekeeping genes

Housekeeping genes were defined according to the following criteria: (i) the transcripts could be detected in all 21 tissues; (ii) the transcripts showed low expression variance across tissues: $P > 0.1$ (Kolmogorov-Smirnov test); (iii) no

100 exceptional expression in any single tissue; that is, the expression values were restricted within

101 the fourfold range of the average across tissues; and (iv) all transcripts of a housekeeping

102 candidate gene met the above criteria.

**Structure analysis of housekeeping genes**

103

104 The structure data of genes were taken from the Ensembl BioMart (Kinsella et al.2011). Human

105 housekeeping genes were derived from the reference (Eisenberg and Levanon 2013), considering

106 its similar type of data and stringency of the definition. We obtained 3,136 and 3,804

107 housekeeping genes of pig and human, respectively. Length of various parts of housekeeping

108 genes between them were compared by Mann-Whitney test (Table 1).

**Gene ontology analysis of housekeeping genes**

109

110 The analysis of functional annotations of housekeeping genes was performed using DAVID, ver.

111 6.7, available on their website (Huang da et al.2009; Huang da et al.2009). All expressed genes

112 in the data were used as background. Comparative analysis of housekeeping genes between

113 human and pig was performed. The false discovery rates (FDR) were calculated to estimate the

114 extent to which genes were enriched in GO categories (Ashburner et al.2000). Probabilities less

115 than 0.01 were used as the cut-off value and considered to show significant level of the

116 correlation. Heat map analysis was also conducted through DAVID outcomes to visualize a

117 matrix of enriched GO.

**Evolutionary feature analysis of housekeeping genes**

118

119 The number of non-synonymous substitutions per non-synonymous site (dN) and the number of

120 synonymous substitutions per synonymous site (dS) were estimated using the Nei-Gojobori

121 method embedded in MEGA 7.0 (Z-test, $P<0.05$)(Kumar et al.2016; Nei and Kumar 2000). From

122 the Scope row, select the Overall Average option. For the Gaps/Missing data treatment option,

123  select Pairwise Deletion. The genome sequence of orthologous genes were downloaded from

124  Ensembl BioMart. The dN/dS ratios were calculated to assess selection pressure (Hurst 2002;

125  Yang and Nielsen 2002; Dasmeh et al.2014). The information of active sites and zinc fingers of

126  proteins were obtained from UniProt Knowledgebase (UniProtKB) (Boutet et al.2016; Pundir et

127  al.2015). Species-specific housekeeping genes that have similar function were processed to

128  search their active sites or zinc fingers.

129

## Results

131  **Gene expression profile**

132  To identify the housekeeping genes in pig, we surveyed the expression distribution of 30,585

133  transcripts across 21 tissues of pig (see Methods, Figure 1, Supplementary material 1: Figure S1).

134  The detectability of RNA-seq data was high, and only 116 transcripts undetected in the present

135  study. The 226 transcripts showed tissue-specific expression(expressed in one tissue), whereas

136  6072 transcripts was found broadly expressed in all tissues (Figure 1). This finding was

137  consistent with the expression tissue-breadth of human genes (Zhu et al.2008; Eisenberg and

138  Levanon 2013).

139

140  **Identification of pig housekeeping genes**

141  To obtain the transcripts with the ubiquitous expression level across pig tissues, we selected the

142  transcripts detected in all tissues and then obtained 6072 candidates. The background differences

143  between different sequencing projects result in batch effect between samples, including

144  difference of sequencing depth and coverage. Therefore, we chose a single sequencing project to

145  assess the uniformity of gene expression, which contains a larger sample size. Furthermore, the

146 expression uniformity of those candidates in ERP002055 sequencing project was tested by the

147 Kolmogorov-Smirnov (K-S) test and then was accessed by the *P*-value of the test(Farajzadeh et

148 al.2013). Figure S2 of Supplementary material 1 represents the frequencies of the candidates

149 with the *P*-value being greater than the given cutoff. For about 67% of all candidates, the *P*-

150 values were above 0.1, implying their expression levels were not significantly varied across

151 tissues and had a high level of the expression uniformity. Therefore, we defined the cutoff of the

152 uniform level as $P > 0.1$ for the following analyses, which resulted in a list of 4068 unique

153 transcripts, belonging to 3754 genes. The housekeeping gene was further restricted into the gene

154 whose all transcripts passed the criteria. Altogether, the 3,136 genes passed the restriction

155 (Supplementary material 2), about a third of which were unannotated.

156

157 Figure 2 shows the overlap of pig housekeeping genes identified in the present study with

158 previously reported human housekeeping genes (Warrington et al.2000; Hsiao et al.2001;

159 Eisenberg and Levanon 2003; Eisenberg and Levanon 2013). In order to more accurately

160 describe the features, housekeeping genes were grouped into three sets of genes, namely,

161 common housekeeping genes observed both in pig and human, human-specific and pig-specific

162 housekeeping genes. We obtained 1,012 common, 2,792 human-specific and 2,124 pig-specific

163 housekeeping genes, respectively.

164

165

166 **Structure comparison of housekeeping genes between pig and human**

167 The comparison of length distribution of total intron, 5' untranslated region (UTR) and coding

168 sequence (CDS) in homologous housekeeping genes shows that pig genes dominates the fraction

169 of long length whereas human genes are prone to short length (Figure 3A - C). Furthermore,

170 Table 1 compares the average lengths of various structures of the housekeeping genes that

171 correspond to one another in pig and human. All structures of pig housekeeping genes were

172 significantly longer than human's (Table 1), which were consistent with the previous analyses of

173 pig genomes (Groenen et al.2012), implying that different purifying selection pressures were

174 applied between pig and human. Selective pressure may make gene as short as possible for

175 reducing the cost in the transcription process (Ucker and Yamamoto 1984; Castillo-Davis et

176 al.2002).

177

178 **Evolutionary dynamics of housekeeping genes**

179 Evolutionary features of housekeeping genes may provide a deeper understanding for the

180 evolutionary trend of housekeeping gene in different species. For the maintenance of essential

181 function, housekeeping genes are thought to evolve more slowly than other genes (Zhang and Li

182 2004). To survey that feature, the number of non-synonymous substitutions per non-synonymous

183 site (dN), the number of synonymous substitutions per synonymous site (dS) and dN/dS ratio

184 were calculated for pig and human housekeeping genes using mouse(*Mus musculus*) as outgroup

185 (Supplementary material 3 and 4), respectively. Generally, synonymous substitutions occurred

186 randomly and do not appear to change the gene function, but the non-synonymous substitutions

187 occurred nonrandomly, which may change the function of housekeeping genes and suffer strong

188 selection pressure (Nei and Kumar 2000, Kimura 1983).

189

190 The dN followed a power law distribution similar to that of the dN/dS (Figure 4A,

191 Supplementary material 1: Figure S3A), displaying a relatively large number of genes with a few

192 non-synonymous substitutions and a small fraction of genes with much more substitutions

193    (Figure 4A). In addition, most of the dN/dS ratios were lower than one, implying that purifying

194    selection have acted on housekeeping genes to ensure the stability of most of genes' function.

195    The less the dN/dS ratio is, the stronger purifying selection is. Furthermore, purifying selection

196    pressure on housekeeping genes were slightly stronger in human than in pig (Figure 4A, B).

197

198    The dN/dS ratios of common housekeeping genes showed no difference between pig and human,

199    but the ratios of species-specific housekeeping genes were significantly lower in human than in

200    pig (Mann-Whitney test, $P < 0.05$) (Figure 4B, Figure 5D). Furthermore, for both human and pig,

201    the dN/dS ratios of common genes were significantly lower than species-specific genes (Figure

202    5A for pig and Supplementary material 1: Figure S4 for human). This result suggested that

203    common housekeeping genes suffered more stringent purifying selection to remove alleles than

204    species-specific genes.

205

206    On the other side, these results of the dN/dS (or dN) also implied that human housekeeping

207    genes have evolved more stably than pig genes (Figure 5B-D). The dS of human species-specific

208    genes were prone towards lower values than pig genes (Figure 5C), showing that human

209    housekeeping genes have slower neutral evolution than pig housekeeping genes.

210

211    The dS followed an approximately normal distribution (Supplementary material 1：Figure S3B),

212    occurring to be around a central value (0.77 and 0.63 in pig and human housekeeping genes,

213    respectively). This finding implies the random tendency of synonymous substitutions. There was

214    no statistic difference in the synonymous substitutions between common and species-specific

215    genes within a species (Figure 5A for pig and Supplementary material 1: Figure S4 for human).

216

217 In addition, considering the mouse is close to human and pig in phylogeny, and may be more

218 close to human(Meredith et al. 2011). So, we also selected elephant (*Loxodonta africana*) as

219 outgroup to calculate dN,dS, and dN/dS for pig and human housekeeping genes,

220 respectively(Additional 5 and 6). Furthermore, all analyses of evolutionary dynamics were

221 performed to verify foregoing results using elephant as outgroup, and the results is similar to the

222 previous analysis of mouse as outgroup (Supplementary material 7).

223

224 **Associated function of housekeeping genes**

225 We then characterized the housekeeping genes that enriched molecular function, biological

226 process, cellular component, and disease, respectively, based on the Database for Annotation,

227 Visualization, and Integrated Discovery (DAVID) program. The heat map shown in Figure 6

228 illustrates the similar enrichment of housekeeping genes between pig and human. Briefly,

229 housekeeping genes were predominantly detected as the genes associated with Gene Ontology

230 (GO) terms related to basal metabolism that are indispensable for cellular physiology, indicating

231 housekeeping genes are essential for basic physiological processes (Figure 6).

232

233 It was worth noting that many pig housekeeping genes were enriched in human diseases,

234 especially in several cancers with high mortality rates: breast cancer, lung cancer and colorectal

235 cancer (Figure 6D). This finding may be beneficial for studies of human disease (Tu et al.2006),

236 given that pig may not have some human risk genes. For instance, alcohol-induced cirrhosis was

237 enriched in human housekeeping genes, but not in pig.

238

### Functional convergence

Interestingly, the functional enrichment analyses showed a coherent trend in pig and human

housekeeping genes although the low overlap of gene lists and the difference in gene structure

between the two species were found. For example, for biological process, pig and human showed

a slight difference in the GO term enrichment (Figure 6A). In addition, similar trends were also

observed in the active molecules that related to basic metabolism and gene expression (Figure 6B

and C).


The above analysis revealed that functions of housekeeping genes between pig and human were

consistent, implying that selection pressure may preclude the species-differentiation of

housekeeping genes for the maintenance of basal cellular functions, especially for species-

specific housekeeping genes. To confirm this conjecture, we performed functional enrichment

analysis for common and species-specific housekeeping genes, respectively. The heat map

shown in Figure 7 illustrates the more similarity between two species-specific terms than

between common and species-specific terms. These results indicated housekeeping genes

suffered strong selection pressure for maintaining normal life activities, and human and pig

species-specific housekeeping genes converged on the basal cellular function.


### Mechanistic convergence

To understand the mechanistic constraints on the function of housekeeping proteins, we analyzed

the evolutionary constraints on protein structure, active site feature and chemical reaction center.

We found some similar active site features in housekeeping peptidases (Figure 8, Table 2), which

reflected the intrinsic chemical constraints on enzymes, leading evolution to independently

converge on equivalent solutions repeatedly (Buller and Townsend 2013; Dodson and Wlodawer

263   1998). The chemical and physical constraints on enzyme catalysis have caused identical triad

264   arrangements in housekeeping peptidases in human-pig lineage, such as classical catalytic

265   Ser/His/Asp triad and non-classical variants (Table 2). However, the peptide sequences and

266   three-dimensional structure profiles of them were totally different (Figure 8A and B). Classical

267   Ser/His/Asp catalytic triad is a universal phenomenon in the serine protease class (E.C. 3.4.21),

268   where serine is the nucleophile, histidine is the general base or acid, and the aspartate helps

269   orient the histidine residue and neutralize the charge that develops on the histidine during the

270   transition states (Polgar 2005; Ekici et al.2008). Interestingly, almost all proteins in Table 2

271   contained histidine as an active site to provide a proton receptor (Wang et al.2006). In addition,

272   Cys/His and Glu/His/Asp in peptidases also evolved convergent; however, these active sites have

273   rarely been mentioned in previous reports to our knowledge.

274

275   **Structural convergence**

276   Moreover, many housekeeping proteins tended to form common zinc finger features involved in

277   the regulation of gene expression (Figure 9, Supplementary material 1: Table S2 and S3). For

278   example, $C_2H_2$ type is one of major zinc fingers in transcription factors (Wolfe et al.2000; Li et

279   al.2004). This analysis of housekeeping protein structure and function revealed several

280   interrelated and previously unrecognized relationships of structure–function constraints. These

281   fundamental constraints have promoted the convergent evolution of housekeeping genes,

282   especially for species-specific housekeeping genes and low homology genes.

283

## Discussion

In the present study, we defined a set of pig housekeeping genes with a wide range of expression

and low expression variation across tissues. The present set of housekeeping genes in pig showed

lower overlap with a human set. Some housekeeping genes of human were not in our list, such as

*GAPDH* and *ACTB* (Barber et al.2005;de Jonge et al.2007; Nygard et al.2007), thus whether

human housekeeping genes can be used as reference controls for other species remains to be

further verified.

After divergence from common ancestor, pig and human have accumulated difference in the

sequence and structure of housekeeping genes. On a molecular level, that can happen from

random mutation, for example, the synonymous substitution. The dS distribution followed an

approximately normal distribution, showing a random tend of synonymous substitutions. On the

other side, the divergence was also related to adaptive changes. Human housekeeping genes were

found to be shorter than pig genes (Figure 3A - C). The possible reason is food intake and stored

energy is less in.human than pig, so the shorter structure is good for human to consume less time

and cost in the process of gene expression (Ucker and Yamamoto 1984; Izban and Luse 1992).

In addition, the stronger purifying selection in human comparing to pig (Figure 4A) might result

in a lower degree of genetic redundancy as well (Zhang and Li 2004). In other words, human

housekeeping genes would have evolved more stably than pig, because advantageous and stable

living environment. Moreover, human and pig have evolved their own species-specific

housekeeping genes, which might lead to the formation of the two species, allowing

differentiated fixation of characteristics. In addition, purifying selection is stronger in common

than in species-specific housekeeping genes and show some differences in GO enrichment. This

may indicate common housekeeping genes were more indispensable than species-specific and

308     involve more functions for sustain life. Such as *GTF2H1* (general transcription factor IIH subunit

309     1) and *CXXC1* (CXXC finger protein 1) in common are crucial for regulation of many of gene

310     expression(Shiekhattar et al.1995; Andersen et al.2001), but in species-specific housekeeping

311     genes were not enrichment.

312

313     However, although human and pig have been divergent for millions of years, both species

314     independently converged towards similar features of housekeeping genes. One of the most

315     unexpected observations stemmed from species-specific housekeeping genes. The GO

316     enrichment analysis revealed that pig-specific and human-specific housekeeping genes have

317     similar functions. In addition, some housekeeping proteins evolved independently to have similar

318     active sites, sidechains, catalytic centers or binding sites to complete similar catalytic reaction or

319     molecular function (Buller and Townsend 2013; Polgar 2005; Ekici et al.2008; Brannigan et

320     al.1995; Chen et al. 2008; Klug 2010; Klug 1999; Hall 2005; Brown 2005), although these

321     proteins showed very low homology with each other. They have "converged" on the maintenance

322     of basic cellular functions, which led to equivalent solutions for adapting to the environment

323     (Nielsen 2005; Hurst 2009). Functional similarity across species may be caused by adaptive

324     evolution (Zhang and Li 2004; Kimura 1983), which drive different species-specific genes to

325     perform similar essential functions, regardless of its specific role in species.

326

327     As known, it is still under investigation to attain large-scale gene expression profile. The current

328     transcriptome sequencing data in pig may be inadequate to meet the requirement to define the

329     housekeeping genes. The accurate definition of housekeeping genes is still an unresolved issue.

330     Therefore, the present set of pig housekeeping genes had limitations, but it successfully offered

331    some instances, the characteristics of which were similar to those reported in previous studies.

332    As new technologies emerge, high-quality deep-sequencing transcriptome profiling data may

333    open up opportunities to improve the stringency in defining housekeeping genes and narrowing

334    the catalog of housekeeping genes that are expressed in a single cell (Tang et al.2009).

335    Furthermore, the advancement of statistical methods will greatly improve housekeeping gene

336    detection. More specifically, the concept of "housekeeping" or "maintenance" should be defined

337    in a hierarchical way related to cell types, growth stages, cell cycles as well as various

338    physiological conditions, and in terms of specific transcript variant (Zhu et al.2008). Thus, we

339    will be able to observe several sets of housekeeping genes in a single species. In addition, more

340    stringent sets of housekeeping genes will also provide powerful support for structural and

341    functional genomics, especially to analyze the cellular basal function of different species (Kumar

342    and Hedges 1998; Meredith et al.2011; Kumar et al.2002).

## Conclusions

344    The present study offered insight into the general aspects of housekeeping gene structure and

345    evolution. Diverging from the ancestor of human and pig, housekeeping genes have varied in

346    gene structure and gene list, but they have converged on the maintenance of basic cellular

347    function that are essential for the existence of a cell, regardless of their specific role in species.

348    The results in the present study will shed light on the evolutionary dynamics of the housekeeping

349    genes.

## Declarations

**Ethics approval and consent to participate**

352    We reused public data from the NCBI database and did not report on or involve the use of any

353    another animal data.

354

355

**Availability of data and material**

All data generated or analysed during this study are included in this published article and its

supplementary information files.

**Authors' contributions**

Kai Wei and Lei Ma designed the study. Kai Wei and Tingting Zhang performed the data

analyses and drafted the manuscript. Lei Mai revised the manuscript. All authors read and

approved the final manuscript.

**Competing interests**

The authors declare that they have no competing interests.

**Author detail**

[1] College of Life Science, Shihezi University, Shihezi City, Xinjiang Province, China

# References

Butte AJ, Dzau VJ, Glueck SB. 2001. Further defining housekeeping, or "maintenance," genes

Focus on "A compendium of gene expression in normal human tissues". Physiol. Genomics,

7(2):95-96.

381  Zhu J, He F, Song S, Wang J, Yu J. 2008. How many human genes can be defined as

382  housekeeping with current expression data? BMC Genomics, 9:172. doi: 10.1186/1471-2164-9-

383  172.

384  Eisenberg E, Levanon EY. 2003. Human housekeeping genes are compact. Trends Genet.

385  19(7):362-365. doi:10.1016/S0168-9525(03)00140-9.

386  Vinogradov AE. 2004. Compactness of human housekeeping genes: selection for economy or

387  genomic design? Trends Genet. 20(5):248-253. doi:10.1016/j.tig.2004.03.006.

388  Vinogradov AE. 2003. Isochores and tissue-specificity. Nucleic Acids Res. 31(17):5212-5220.

389  doi:10.1093/nar/gkg699.

390  Farré D, Bellora N, Mularoni L, Messeguer X, Albà MM. 2007.   Housekeeping genes tend to

391  show reduced upstream sequence conservation. Genome Biol. 8(7):R140. doi: 10.1186/gb-2007-

392  8-7-r140.

393  Bellora N, Farré D, Albà MM. 2007. Positional bias of general and tissue-specific regulatory

394  motifs in mouse gene promoters. BMC Genomics, 8:459. doi: 10.1186/1471-2164-8-459.

395  Thellin O, Zorzi W, Lakaye B, De Borman B, Coumans B, Hennen G, Grisar T, Igout A, Heinen

396  E. 1999. Housekeeping genes as internal standards: use and limits. J. Biotechnol. 75(2-3):291-

397  295. doi:10.1016/S0168-1656(99)00163-7.

398  Robinson MD, Oshlack A. 2010. A scaling normalization method for differential expression

399  analysis of RNA-seq data. Genome Biol. 11(3):R25. doi: 10.1186/gb-2010-11-3-r25.

400  Rubie C, Kempf K, Hans J, Su T, Tilton B, Georg T, Brittner B, Ludwig B, Schilling M. 2005.

401  Housekeeping gene variability in normal and cancerous colorectal, pancreatic, esophageal,

402  gastric and hepatic tissues. Mol. Cell Probes. 19(2):101-109. doi:10.1016/j.mcp.2004.10.001.

403  Vandesompele J, De Preter K, Pattyn F, Poppe B, Van Roy N, De Paepe A, Speleman F. 2002.

404  Accurate normalization of real-time quantitative RT-PCR data by geometric averaging of

405  multiple internal control genes. Genome Biol. 3(7):RESEARCH0034.1. doi:10.1186/gb-2002-3-

406  7-research0034.

407  Warrington JA, Nair A, Mahadevappa M, Tsyganskaya M. 2000. Comparison of human adult

408  and fetal expression and identification of 535 housekeeping/maintenance genes. Physiol.

409  Genomics, 2(3):143-147.

410  Hsiao LL, Dangond F, Yoshida T, Hong R, Jensen RV, Misra J, Dillon W, Lee KF, Clark KE,

411  Haverty P, Weng Z, Mutter GL, Frosch MP, MacDonald ME, Milford EL, Crum CP, Bueno R,

412  Pratt RE, Mahadevappa M, Warrington JA, Stephanopoulos G, Stephanopoulos G, Gullans SR.

413  2001. A compendium of gene expression in normal human tissues. Physiol. Genomics, 7(2):97-

414  104.

415  Gingeras TR. 2007. Origin of phenotypes: genes and transcripts. Genome Res. 17(6):682-690.

416  doi:10.1101/gr.6525007.

417  Greer S, Honeywell R, Geletu M, Arulanandam R, Raptis L. 2010. Housekeeping genes;

418  expression levels may change with density of cultured cells. J. Immunol. Methods, 355(1-2):76-

419  79. doi:10.1016/j.jim.2010.02.006.

420  Marioni JC, Mason CE, Mane SM, Stephens M, Gilad Y. 2008. RNA-seq: an assessment of

421  technical reproducibility and comparison with gene expression arrays. Genome Res. 18(9):1509-

422  1517. doi:10.1101/gr.079558.108.

423  Fu X, Fu N, Guo S, Yan Z, Xu Y, Hu H, Menzel C, Chen W, Li Y, Zeng R, Khaitovich P. 2009.

424  Estimating accuracy of RNA-Seq and microarrays with proteomics. BMC Genomics, 10:161. doi:

425  10.1186/1471-2164-10-161.

426  Bradford JR, Hey Y, Yates T, Li Y, Pepper SD, Miller CJ. 2010. A comparison of massively

427  parallel nucleotide sequencing with oligonucleotide microarrays for global transcription profiling.

428  BMC Genomics, 11:282. doi: 10.1186/1471-2164-11-282.

429  Draghici S, Khatri P, Eklund AC, Szallasi Z. 2006. Reliability and reproducibility issues in DNA

430  microarray measurements. Trends Genet.  22(2):101-109.  doi:10.1016/j.tig.2005.12.005

431  Eisenberg E, Levanon EY. 2013. Human housekeeping genes, revisited. Trends Genet.

432  29(10):569-574. doi:10.1016/j.tig.2013.05.010.

433  Gerstein MB, Bruce C, Rozowsky JS, Zheng D, Du J, Korbel JO, Emanuelsson O, Zhang ZD,

434  Weissman S, Snyder M. 2007. What is a gene, post-ENCODE? History and updated definition.

435  Genome Res. 17(6):669-681. doi:10.1101/gr.6339607.

436  Brattelid T, Winer LH, Levy FO, Liestol K, Sejersted OM, Andersson KB. 2010. Reference gene

437  alternatives to Gapdh in rodent and human heart failure gene expression studies. BMC Mol. Biol.

438  11:22. doi: 10.1186/1471-2199-11-22.

439  Kozera B, Rapacz M. 2013. Reference genes in real-time PCR. J. Appl. Genet. 54(4):391-406.

440  doi:10.1007/s13353-013-0173-x.

441  Lunney JK. 2007. Advances in swine biomedical model genomics. Int J Biol Sci. 3(3):179-184.

442  Rolandsson O, Haney MF, Hagg E, Biber B, Lernmark A. 2002. Streptozotocin induced diabetes

443  in minipig: a case report of a possible model for type 1 diabetes? Autoimmunity, 35(4):261-264.

444  Lee L, Alloosh M, Saxena R, Van Alstine W, Watkins BA, Klaunig JE, Sturek M, Chalasani N.

445  2009. Nutritional model of steatohepatitis and metabolic syndrome in the Ossabaw miniature

446  swine. Hepatology, 50(1):56-67. doi:10.1002/hep.22904.

447  Becker ST, Rennekampff HO, Alkatout I, Wiltfang J, Terheyden H. 2010. Comparison of

448  vacuum and conventional wound dressings for full thickness skin grafts in the minipig model.

International journal of oral and maxillofacial surgery. 39(7):699-704. doi:10.1016/j.ijom.2010.03.016.

Kodama Y, Shumway M, Leinonen R. 2012. The Sequence Read Archive: explosive growth of sequencing data. Nucleic Acids Res. 40(Database issue):D54-56. doi: 10.1093/nar/gkr854.

Farajzadeh, L, Hornshoj H, Momeni J, Thomsen B, Larsen K, Hedegaard J, Bendixen C, Madsen LB. 2013. Pairwise comparisons of ten porcine tissues identify differential transcriptional regulation at the gene, isoform, promoter and transcription start site level. Biochem. Biophys. Res. Commun. 438(2):346-352. doi:10.1016/j.bbrc.2013.07.074.

Patel RK, Jain M. 2012. NGS QC Toolkit: A Toolkit for Quality Control of Next Generation Sequencing Data. PloS one., 7(2):: e30619. doi:10.1371/journal.pone.0030619.

Trapnell C, Pachter L, Salzberg SL. 2009. TopHat:discovering splice junctions with RNA-Seq. Bioinformatics. 25(9): 1105-1111. doi:10.1093/bioinformatics/btp120.

Külahoglu C, Bräutigam A. 2014. Quantitative Transcriptome Analysis Using RNA-seq. Methods Mol. Biol. 1158:71-91. doi: 10.1007/978-1-4939-0700-7_5.

Ghosh S, Chan KK. 2016. Analysis of RNA-Seq Data Using TopHat and Cufflinks. Methods Mol. Biol.1374:339-361. doi: 10.1007/978-1-4939-3167-5_18.

Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren J, Salzberg SL, Wold BJ, Pachter L. 2010. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. Nat. Biotechnol. 28(5):511-515. doi: 10.1038/nbt.1621.

Kumar S, Stecher G, Tamura K. 2016. MEGA7: Molecular Evolutionary Genetics Analysis version 7.0 for bigger datasets. Mol. Biol. Evol. 33(7):1870-1874. doi:10.1093/molbev/msw054.

471 Nei M, Kumar S. 2000. Molecular Evolution and Phylogenetics, Oxford University Press,Oxford,

472 pp.52-72.

473 Hurst LD. 2002. The Ka/Ks ratio: diagnosing the form of sequence evolution. Trends Genet.

474 18(9):486. doi:10.1016/S0168-9525(02)02722-1.

475 Yang Z, Nielsen R. 2002. Codon-substitution models for detecting molecular adaptation at

476 individual sites along specific lineages. Mol. Biol. Evol. 19(6):908-917.

477 doi:10.1093/oxfordjournals.molbev.a004148.

478 Huang da W, Sherman BT, Lempicki RA. 2009. Systematic and integrative analysis of large

479 gene lists using DAVID bioinformatics resources. Nat. protoc. 4(1):44-57. doi:

480 10.1038/nprot.2008.211.

481 Huang da W, Sherman BT, Lempicki RA. 2009. Bioinformatics enrichment tools: paths toward

482 the comprehensive functional analysis of large gene lists. Nucleic Acids Res. 37(1):1-13. doi:

483 10.1093/nar/gkn923.

484 Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, Davis AP, Dolinski K,

485 Dwight SS, Eppig JT, Harris MA, Hill DP, Issel-Tarver L, Kasarskis A, Lewis S, Matese JC,

486 Richardson JE, Ringwald M, Rubin GM, Sherlock G. 2000. Gene ontology: tool for the

487 unification of biology. The Gene Ontology Consortium. Nat. Genet. 25(1):25-29. DOI:

488 10.1038/75556.

489 Zhang L, Li WH. 2004. Mammalian housekeeping genes evolve more slowly than tissue-specific

490 genes. Mol. Biol. Evol. 21(2):236-239. doi: 10.1093/molbev/msh010.

491 Kimura M. 1983. The Neutral Theory of Molecular Evolution. Cambridge Univ. Press,

492 Cambridge, U.K.

493    Kumar S, Hedges SB. 1998. A molecular timescale for vertebrate evolution. Nature,

494    392(6679):917-920. doi:10.1038/31927.

495    Meredith RW, Janecka JE, Gatesy J, Ryder OA, Fisher CA, Teeling EC, Goodbla A, Eizirik, E,

496    Simao TL, Stadler T, Rabosky DL, Honeycutt RL, Flynn JJ, Ingram CM, Steiner C, Williams TL,

497    Robinson TJ, Burk-Herrick A, Westerman M, Ayoub NA, Springer MS, Murphy WJ. 2011.

498    Impacts of the Cretaceous Terrestrial Revolution and KPg extinction on mammal diversification.

499    Science, 334(6055):521-524. doi: 10.1126/science.1211028.

500    Barber RD, Harmer DW, Coleman RA, Clark BJ. 2005. GAPDH as a housekeeping gene:

501    analysis of GAPDH mRNA expression in a panel of 72 human tissues. Physiol. Genomics,

502    21(3):389-395. doi:10.1152/physiolgenomics.00025.2005

503    de Jonge HJ, Fehrman RS, de Bont ES, Hofstra RM, Gerbens F, Kamps WA, de Vries EG, van

504    der Zee AG, te Meerman GJ, ter Elst A. 2007. Evidence based selection of housekeeping genes.

505    PloS one, 2(9):e898. doi:10.1371/journal.pone.0000898.

506    Freilich S, Massingham T, Bhattacharyya S, Ponsting H, Lyons PA, Freeman TC, Thornton JM.

507    2005. Relationship between the tissue-specificity of mouse gene expression and the evolutionary

508    origin and function of the proteins. Genome Biol. 6(7):R56. doi:10.1186/gb-2005-6-7-r56.

509    Zhu J, He F, Hu S, Yu J. 2008. On the nature of human housekeeping genes. Trends Genet.

510    24(10):481-484. doi:10.1016/j.tig.2008.08.004.

511    Groenen MA, Archibald AL, Uenishi H, Tuggle CK, Takeuchi Y, Rothschild MF, Rogel-

512    Gaillard C, Park C, Milan D, Megens HJ, et al. 2012. Analyses of pig genomes provide insight

513    into porcine demography and evolution. Nature, 491(7424):393-398. doi: 10.1038/nature11622.

514   Ucker DS, Yamamoto KR. 1984. Early events in the stimulation of mammary tumor virus RNA

515   synthesis by glucocorticoids. Novel assays of transcription rates. J. Biol. Chem. 259(12):7416-

516   7420.

517   Izban MG, Luse DS. 1992. Factor-stimulated RNA polymerase II transcribes at physiological

518   elongation rates on naked DNA but very poorly on chromatin templates. J. Biol. Chem.

519   267(19):13647-13655.

520   Nielsen R. 2005. Molecular Signatures of Natural Selection. Annu. Rev. Genet. 39:197-218. doi:

521   10.1146/annurev.genet.39.073003.112420.

522   Hurst LD. 2009. Genetics and the understanding of selection. Nat. Rev. Genet. Doi:10(2):83-93.

523   10.1038/nrg2506.

524   Tang F, Barbacioru C, Wang Y, Nordman E, Lee C, Xu N, Wang X, Bodeau J, Tuch BB,

525   Siddiqui A, Lao K, Surani MA. 2009. mRNA-Seq whole-transcriptome analysis of a single cell.

526   Nat. methods, 6(5):377-382. doi:10.1038/nmeth.1315.

527   Dasmeh P, Serohijos AW, Kepp KP, Shakhnovich EI. 2014. The influence of selection for

528   protein    stability    on    dN/dS    estimations.    Genome    Biol.    Evol.    6(10):2956-67.    doi:

529   10.1093/gbe/evu223.

530   Tu Z, Wang L, Xu M, Zhou X, Chen T, Sun F. 2006. Further understanding human disease genes

531   by    comparing    with    housekeeping    genes    and    other    genes.    BMC    Genomics,    7:31.    doi:

532   10.1186/1471-2164-7-31.

533   Kumar S, Subramanian S. 2002. Mutation rates in mammalian genomes. Proc. Natl. Acad. Sci.

534   USA. 99(2):803-808. doi:10.1073/pnas.022629899

535  Nygard AB, Jorgensen CB, Cirera S, Fredholm M. 2007. Selection of reference genes for gene

536  expression studies in pig tissues using SYBR green qPCR. BMC Mol. Biol. 8:67. doi:

537  10.1186/1471-2199-8-67.

538  Castillo-Davis CI, Mekhedov SL, Hartl DL, Koonin EV, Kondrashov FA. 2002. Selection for

539  short introns in highly expressed genes. Nat. Genet.   31(4):415-418. doi:10.1038/ng940.

540  Buller AR, Townsend CA. 2013. Intrinsic evolutionary constraints on protease structure, enzyme

541  acylation, and the identity of the catalytic triad. Proc Natl Acad Sci USA. 110(8):E653-661. doi:

542  10.1073/pnas.1221050110.

543  Polgar L. 2005. The catalytic triad of serine peptidases. Cell Mol. Life Sci. 62(19-20):2161-2172.

544  doi: 10.1007/s00018-005-5160-x

545  Ekici OD, Paetzel M, Dalbey RE. 2008. Unconventional serine proteases: variations on the

546  catalytic    Ser/His/Asp    triad    configuration.    Protein    Sci.    17(12):2023-2037.    doi:

547  10.1110/ps.035436.108.

548  Brannigan JA, Dodson G, Duggleby HJ, Moody PC, Smith JL, Tomchick DR, Murzin AG. 1995.

549  A protein catalytic framework with an N-terminal nucleophile is capable of self-activation.

550  Nature, 378(6555):416-419. doi:10.1038/378416a0.

551  Chen L, Wang H, Zhang J, Gu L, Huang N, Zhou JM, Chai J. 2008. Structural basis for the

552  catalytic  mechanism  of  phosphothreonine  lyase.  Nat.  Struct.  Mol.  Biol.  15(1):101-102.

553  doi:10.1038/nsmb1329.

554  Wang LJ, Sun N, Terzyan S, Zhang XJ, Benson DR. 2006. A Histidine/Tryptophan $\pi$-Stacking

555  Interaction Stabilizes the Heme-Independent Folding Core of Microsomal Apocytochrome

556  b5Relative to that of Mitochondrial Apocytochrome b5. Biochemistry 45 (46): 13750 -13759.

557  doi: 10.1021/bi0615689.

558    Wolfe SA, Nekludova L, Pabo CO. 2000. DNA recognition by Cys2His2 zinc finger proteins.

559    Annu. Rev. Biophys. Biomol. Struct. 29:183-212.

560    doi: 10.1146/annurev.biophys.29.1.183.

561    Li L, He S, Sun JM, Davie JR. 2004. Gene regulation by Sp1 and Sp3. Biochemistry and cell

562    biology , 82(4):460-471. doi: 10.1139/o04-045.

563    Klug A. 2010. The discovery of zinc fingers and their applications in gene regulation and

564    genome manipulation. Q. Rev. Biophys. 43(1):1-21. doi:10.1017/S0033583510000089.

565    Klug A. 1999. Zinc finger peptides for the regulation of gene expression. J. Mol. Biol.

566    293(2):215-218. doi: 10.1006/jmbi.1999.3007.

567    Hall TM. 2005. Multiple modes of RNA recognition by zinc finger proteins. Curr. Opin. Struct.

568    Biol. 15(3):367-373. doi:10.1016/j.sbi.2005.04.004.

569    Brown RS. 2005. Zinc finger proteins: getting a grip on RNA. Curr. Opin. Struct. Biol. 15(1):94-

570    98. doi:10.1016/j.sbi.2005.01.006.

571    Boutet E, Lieberherr D, Tognolli M, Schneider M, Bansal P, Bridge AJ, Poux S, Bougueleret L,

572    Xenarios, I. 2016. UniProtKB/Swiss-Prot, the Manually Annotated Section of the UniProt

573    KnowledgeBase: How to Use the Entry View. Methods Mol. Biol. 1374:23-54. doi:10.1007/978-

574    1-4939-3167-5_2.

575    Pundir S, Magrane M, Martin MJ, O'Donovan C. 2015. Searching and Navigating UniProt

576    Databases. Curr. Protoc. Bioinformatics, 50:1.27.1-10. doi: 10.1002/0471250953.bi0127s50.

577    Dodson G, Wlodawer A. 1998. Catalytic triads and their relatives. Trends Biochem. Sci.

578    23(9):347-352. doi:10.1016/S0968-0004(98)01254-7

579    Meredith RW, Janečka JE, Gatesy J, Ryder OA, Fisher CA, Teeling EC, Goodbla A, Eizirik E,

580    Simão TL, Stadler T, Rabosky DL, Honeycutt RL, Flynn JJ, Ingram CM, Steiner C, Williams TL,

581  Robinson TJ, Burk-Herrick A, Westerman M, Ayoub NA, Springer MS, Murphy WJ. 2011.

582  Impacts of the cretaceous terrestrial revolution and KPg extinction on mammal diversification.

583  Science, 334(6055):521-524. doi: 10.1126/science.1211028

584  Shiekhattar R, Mermelstein F, Fisher RP, Drapkin R, Dynlacht B, Wessling HC, Morgan DO,

585  Reinberg D. 1995. Cdk-activating kinase complex is a component of human transcription factor

586  TFIIH. Nature, 374(6519):283-287. doi:10.1038/374283a0.

587  Lee JH, Voo KS, Skalnik DG. 2001. Identification and characterization of the DNA binding

588  domain of CpG-binding protein. J. Biol. Chem.      276(48):44669-44676.  doi:

589  10.1074/jbc.M107179200.

590  Esteve-Codina A, Kofler R, Palmieri N, Bussotti G, Notredame C, Perez-Enciso M. 2011.

591  Exploring the gonad transcriptome of two extreme male pigs with RNA-seq. BMC Genomics,

592  12:552. doi:10.1186/1471-2164-12-552.

593  Martinez-Montes AM, Fernández A, Pérez-Montarelo D, Alves E, Benitez RM, Nuñez Y, Óvilo

594  C, Ibañez-Escriche N, Folch, JM, Fernández AI. 2016. Using RNA-Seq SNP data to reveal

595  potential causal mutations related to pig production traits and RNA editing. Anim. Genet.

596  48(2):151-165. doi:10.1111/age.12507.

597  Wang T, Jiang A, Guo Y, Tan Y, Tang G, Mai M, Liu H, Xiao J, Li M, Li X. 2013. Deep

598  sequencing of the transcriptome reveals inflammatory features of porcine visceral adipose tissue.

599  Int. J. Biol. Sci. 9(6):550-556. doi:10.7150/ijbs.6257.

600  Pérez-Montarelo D, Madsen O, Alves E, Rodriguez MC, Folch JM, Noguera JL, Groenen MA,

601  Fernández AI. 2014. Identification of genes regulating growth and fatness traits in pig through

602  hypothalamic   transcriptome   analysis.   Physiol.   Genomics,   2014,   46(6):195-206.

603  doi:10.1152/physiolgenomics.00151.2013.

604 Jiang S, Wei H, Song T, Yang Y, Peng J, Jiang S. 2013. Transcriptome comparison between

605 porcine subcutaneous and intramuscular stromal vascular cells during adipogenic differentiation.

606 PloS one. 8(10):e77094. doi:10.1371/journal.pone.0077094.

607 Li M, Tian S, Jin L, Zhou G, Li Y, Zhang Y, Wang T, Yeung CKL, Chen L, Ma J. et al. 2013.

608 Genomic analyses identify distinct patterns of selection in domesticated pigs and Tibetan wild

609 boars. Nat. Genet. 45(12):1431-1438. doi:10.1038/ng.2811.

610 Samborski A, Graf A, Krebs S, Kessler B, Bauersachs S. 2013. Deep sequencing of the porcine

611 endometrial transcriptome on day 14 of pregnancy. Biol. Reprod. 88(4):84.

612 doi:10.1095/biolreprod.113.107870.

613 Zhang X, Huang L, Wu T, Feng Y, Ding Y, Ye P, Yin Z. 2015. Transcriptomic Analysis of

614 Ovaries from Pigs with High And Low Litter Size. PloS one. 10(10):e0139514.

615 doi:10.1371/journal.pone.0139514.

616 Endale Ahanda ML, Fritz ER, Estelle J, Hu ZL, Madsen O, Groenen MA, Beraldi D,

617 Kapetanovic R, Hume DA, Rowland RR, Lunney JK, Rogel-Gaillard C, Reecy JM, Giuffra E.

618 2012. Prediction of altered 3'- UTR miRNA-binding sites from RNA-Seq data: the swine

619 leukocyte antigen complex (SLA) as a model region. PloS one. 7(11):e48607.

620 doi:10.1371/journal.pone.0048607.

621 Liu H, Nguyen YT, Nettleton D, Dekkers JC, Tuggle CK. 2016. Post-weaning blood

622 transcriptomic differences between Yorkshire pigs divergently selected for residual feed intake.

623 BMC Genomics, 17:73. doi: 10.1186/s12864-016-2395-x.

624 Rahman KM, Camp ME, Prasad N, McNeel AK, Levy SE, Bartol FF, Bagnell CA. 2016. Age

625 and Nursing Affect the Neonatal Porcine Uterine Transcriptome. Biol. Reprod. 2016, 94(2):46.

626 doi:10.1095/biolreprod.115.136150.

627　Miller LC, Bayles DO, Zanella EL, Lager KM. 2016. Effects of Pseudorabies Virus Infection on

628　the Tracheobronchial Lymph Node Transcriptome. Bioinform. Biol. Insights. 9(Suppl 2):25-36.

629　doi: 10.4137/BBI.S30522.

630　Samborski, A, Graf A, Krebs S, Kessler B, Reichenbach M, Reichenbach HD, Ulbrich SE,

631　Bauersachs S. 2013. Transcriptome changes in the porcine endometrium during the

632　preattachment phase. Biol. Reprod. 89(6):134. doi: 10.1095/biolreprod.113.112177.

**Figure 1**(on next page)

The number of tissues where a given transcript was detected.

The expression breadth (horizontal axis) denotes the number of tissues where a given transcript was detected. The zero value of the expression breadth indicates undetected transcripts.

**Figure 2**(on next page)

Overlap of housekeeping genes between pig and human.

Overlap of pig housekeeping gene set identified in the present study(A) with three human gene sets identified by microarray data (Warrington et al.2000; Hsiao et al.2001; Eisenberg and Levanon 2003) and (B)with a human set identified by RNA-seq data (Eisenberg and Levanon 2013).

**Figure 3**(on next page)

Comparison of length distribution of homologous housekeeping gene structures between pig and human.

nt, nucleotide(s); 5'UTR, 5'untranslated region (UTR); CDS, coding sequence.

# Figure 4(on next page)

Purifying selection on housekeeping genes.

(A) The distribution of the dN/dS ratio. (B) The dN/dS ratios of total (all HK), common (co-HK) and species-specific (sp-HK) housekeeping genes were compared between pig and human (Mann-Whitney test, * denoted $P < 0.05$), respectively.

# Figure 5(on next page)

Comparison of evolutionary features of housekeeping genes.

(A) The dN, dS and dN/dS of all, common and species-specific of pig housekeeping genes were compared based on the Mann-Whitney test, respectively. All such means which share a common English letter are similar; otherwise, they differ significantly at $p < 0.05$. (B) - (D) Distributions of dN, dS and dN/dS of species-specific housekeeping genes in pig and human.

**Figure 6**(on next page)

Functional enrichment analysis for housekeeping genes.

Housekeeping genes were enriched in GO categories of (A)biological process, (B) cellular component, (C) molecular function, (D) molecular functions . The basal cellular function between pig and human showed high consistency. (A) (1) Biological process categories included the basal metabolism, (2) regulation of metabolic processes, (3) cellular transport, (4) cell cycle, (5)gene expression and regulation. (B) (1) Cellular component categories included organelle, (2) nuclear, (3) micromolecular complex. (C) (1) Molecular function categories included catalytic activity, (2) transcription factor activity, (3)binding activity, (4) transporter activity. (D) (1) Disease categories included tumour, (2) cancer, (3) chromosomal damage and repair, (4) other disease.
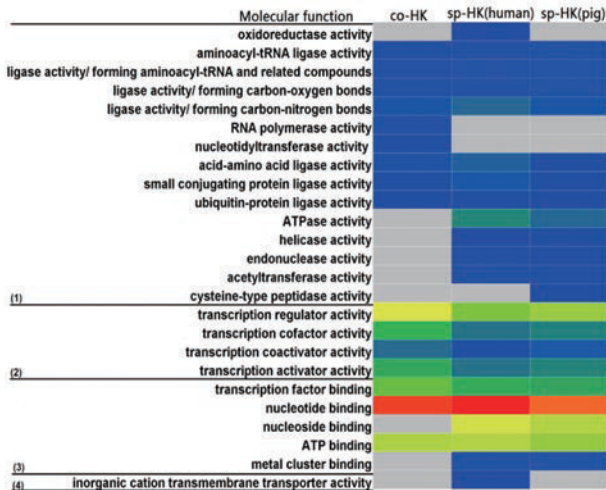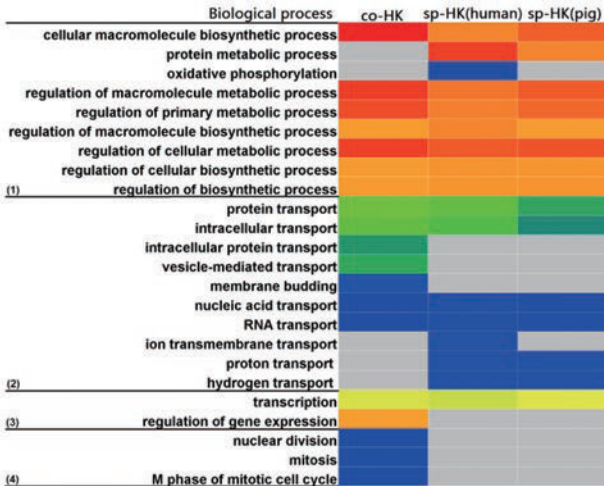
# Figure 7(on next page)

Comparison of functional enrichment analysis.

When we compared functional enrichment, common housekeeping genes (co-HK) showed significant difference with species-specific housekeeping genes (sp-HK), but the sp-HKgenes between pig and human showed very high consistency. (A) (1) Biological process categories included the basal metabolism and regulation, (2) cellular transport, (3 )gene expression and regulation, (4) nuclear division. (B) (1) Molecular function categories included catalytic activity, (2 )transcription factor activity, (3) binding activity, (4) transporter activity.
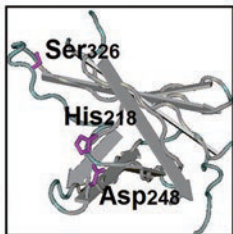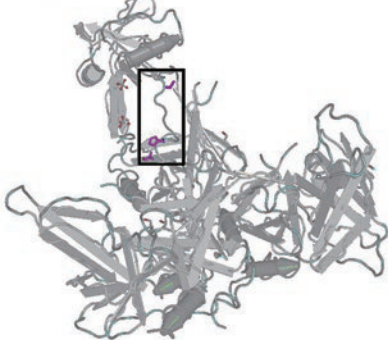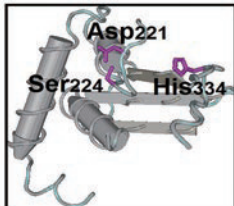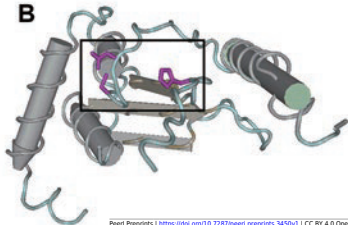
# Figure 8(on next page)

Structures of the "classical" Ser/His/Asp triad configuration.

(A) Serine protease HTRA4 from pig. (B) OTU domain-containing protein 5 from human. A zoomed-in view of the catalytic domain is shown to the right of each structure. The side chains of Ser/His/Asp triad are shown in principle.

A
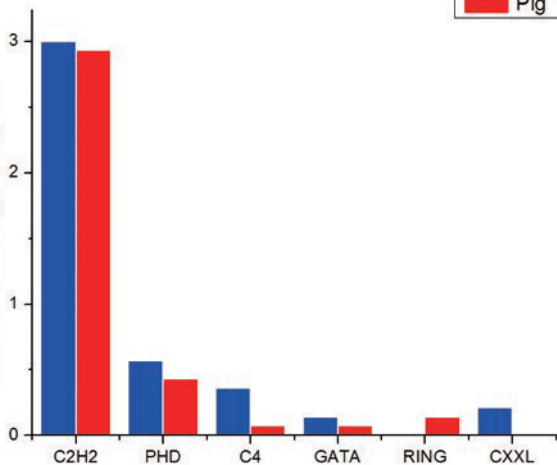
Ser₃₂₆

His₂₁₈

Asp₂₄₈

B

Asp₂₂₁

Ser₂₂₄        His₃₃₄

**Figure 9**(on next page)

Convergent evolution of regulatory proteins towards forming common zinc finger.

The number of zinc fingers per gene was standardized through dividing the number of each type of zinc finger by the number of proteins containing the zinc finger.

# Table 1(on next page)

Comparison of housekeeping genes between pig and human

[a] The length is measured in nucleotides. [b] The value gives the average and standard error of mean. [c] The $p$-value was calculated based on the Mann-Whitney test. UTR, untranslated region; CDS, coding sequence.

1 **Table 1 Comparison of housekeeping genes between pig and human**

| Structure | Pig | Human | $P$-value [c] |
|---|---|---|---|
| Total intron length [a] | 28,108±173 [b] | 21,062±297 | $1.5e^{-105}$ |
| 5' UTR length | 156±3 | 125±1.5 | $3.7e^{-34}$ |
| 3' UTR length | 658±13 | 549±5 | $1.4e^{-73}$ |
| Average exon length per gene | 261±3 | 227±1 | $1.8e^{-6}$ |
| CDS length | 2,181±10 | 1,460±5 | $8.7e^{-234}$ |
| Transcript length | 3,312±13 | 2,200±5 | $7.7e^{-7}$ |
| Number of exons | 9.2±0.1 | 8.8±0.2 | $1.7e^{-4}$ |

2 [a] The length is measured in nucleotides. [b] The value gives the average and standard error of mean.

3 [c] The $p$-value was calculated based on the Mann-Whitney test. UTR, untranslated region; CDS,

4 coding sequence.

5

**Table 2**(on next page)

Active site of convergently related peptidases.

*a* the number following amino acid represents the position of the amino acid in protein.

1   **Table 2 Active site of convergently related peptidases**

| Species | Gene | Protein | Nucleophile [a] | General base | Other active site residues |
|---|---|---|---|---|---|
| Pig | BLMH | Bleomycin hydrolase | Cys73 | His372 | Asn396 |
| | AFG3L2 | AFG3-like protein 2 | Glu575 | His574 | Asp649 |
| | HTRA4 | Serine protease HTRA4 | Ser326 | His218, | Asp248 |
| | CAPN7 | Calpain-7 | Cys290 | His458 | Asn478 |
| Human | OTUD5 | OTU domain-containing protein 5 | Ser224 | His334 | Asp221 |
| | SENP6 | Sentrin-specific protease 6 | Cys1030 | His765 | Asp917 |
| | USP14 | Ubiquitin carboxyl-terminal hudrolase 14 | Cys114 | His435 | |
| | LONP1 | Lon protease homolog, mitochondrial | Ser855 | Lys898 | |

2   [a] the number following amino acid represents the position of the amino acid in protein.

3