

 Open access • Posted Content • DOI:10.1101/454363

A Whole-Genome Association Approach for Large-scaled Inter-species Trait

— [Source link](#) 

Qi Wu, Huizhong Fan, Lei Chen, Yibo Hu ...+1 more authors

Institutions: Chinese Academy of Sciences, Northwestern Polytechnical University

Published on: 26 Oct 2018 - bioRxiv (Cold Spring Harbor Laboratory)

Topics: Genome-wide association study, Candidate gene, Genetic association, Genomics and Genetic marker

Related papers:

- [Genome-Wide Associations of Gene Expression Variation in Humans](#)
- [Finding genetic variants in plants without complete genomes](#)
- [On the Prospects of Whole-genome Association Mapping in *Saccharomyces cerevisiae*](#)
- [Linkage analysis in the next-generation sequencing era.](#)
- [Genomic properties of structural variants and short tandem repeats that impact gene expression and complex traits in humans](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/a-whole-genome-association-approach-for-large-scaled-inter-vqte1gy0qe>

1 **A Whole-Genome Association Approach for Large-scaled** 2 **Inter-species Trait**

3 Qi Wu^{1,#}, Huizhong Fan^{1,2,#}, Lei Chen³, Yibo Hu^{1,4}, Fuwen Wei^{1,2,4*}

4

5 ¹*Key Laboratory of Animal Ecology and Conservation Biology, Institute of Zoology, Chinese*
6 *Academy of Sciences, Beijing, China.*

7 ²*University of Chinese Academy of Sciences, Beijing, China.*

8 ³*Northwestern Polytechnical University, Xi'an, China.*

9 ⁴*Center for Excellence in Animal Evolution and Genetics, Chinese Academy of Sciences, Kunming,*
10 *China*

11 [#]These authors contributed equally to this work.

12 ^{*}Correspondence should be addressed to Fuwen Wei (weifw@ioz.ac.cn).

13

1 **Abstract:**

2 Genome wide association studies (GWAS) have provided an avenue for the association between
3 common genetic variants and complex traits. However, using SNP as a genetic marker, GWAS has
4 been confined to detect genetic basis traits only for within species but not for the large-scale
5 inter-species traits. Here, we propose a practical statistical approach that is using kmer frequencies
6 as the genetic markers to associate genetic variants with large scale inter-species traits. We applied
7 this new approach to the trait of chromosome number in 96 mammalian proteomes, and we
8 prioritized 130 genes including *TP53* and *BAD*, of which 6 were candidate genes. These genes
9 were proved to be associated with cellular reaction of DNA double-strand breaks caused by
10 chromosome fission/fusion. Our study provides a new effective genomic strategy to perform
11 association studies for large-scaled inter-species traits, using the chromosome number as a case.
12 We hope this approach could provide exploration for broadly widely traits.

13

14 **Key words:** large-scaled inter-species trait, chromosome number, kmer frequency,
15 alignment-free, genome-wide association study

16

1 **Introduction**

2 With the development of sequencing technologies, GWAS have been proven powerful and
3 increasingly affordable tool for scanning and detecting candidate genes responsible for complex
4 disease in human or important economically traits in plants or domestic animals [1-4]. Generally
5 speaking, it is the central aim of genetics that to understand the association between heritable
6 variations and phenotypic changes [5]. At present, all of the GWAS are focused on the genetic
7 variants in hundreds or thousands of individuals under species. It is an interesting question to ask
8 that whether GWAS is available between genetic variants and large-scale inter-species traits. In
9 practice, the major problem is lack of an approach to bridge the genome data to the heritable traits.
10 This entails two close related issues: firstly, a marker is needed to quantify various difference of
11 genome/proteome sequences among species. Secondly, a statistical method is needed to assess the
12 association between the makers and the large-scale inter-species traits.

13

14 For the first issue, we are able to identify the probable maker to represent the characteristics of
15 different genomes with the help of alignment-free (AF) sequence comparison analysis [6-9].
16 Broadly speaking, AF methods focuses on the subsequences (words or kmer) of defined or varied
17 length from the sequence, which their types [6], frequencies [7], or position [10] of the kmers on
18 the sequence were considered as characteristics. Various compositional heterogeneity of sequence,
19 including CpG islands, transposon elements, centromere as well as telomere in the genome could
20 be delicately reflected in the characteristics of kmers [11]. Moreover, all types of variations such
21 as substitution, insertion/deletion, duplication as well as rearrangements could also be recorded in
22 the characteristic of subsequences [12]. For quantifying the overall feature of the sequence, AF
23 methods compares and analyzes the characteristics between the concerned sequences . In practice,
24 many excellent AF sequence comparison studies have been successfully performed to solve the
25 problems of whole-genome phylogeny construction. The results of AF comparison studies are in
26 correspond with those using morphological data or molecular data [7,9,12,13]. Therefore, it is
27 reasonable to deduce that the kmer characteristics of the genome can describe the sequence
28 variations which determine the large-scaled inter-species traits. This suggests the possibility of

1 using kmer characteristics as the genetic markers to preform GWAS analysis on inter-species
2 traits.

3

4 Numerous strategies and statistical approaches have been developed to meet the technical
5 challenges and the great opportunities provided by GWAS [4,14,15]. One of such approaches is
6 gene-based GWAS analysis [14]. In GWAS analysis, gene-based association analysis could
7 combine the significance of all SNPs in a gene and effectively assess the association between gene
8 and phenotype traits in hundreds or thousands of individuals. As such, it is possible to combine the
9 significance of all kmer frequencies in a gene and assess the significance of the association
10 between the gene and the inter-species traits in the evolutionary studies. Thus, the candidate genes
11 that have significant effects on the large-scale inter-species traits during the evolution could be
12 identified. Moreover, by mapping the significant kmers to the gene sequences, the potential
13 functional domain regions in the genes associated with the traits could also be identified.

14

15 In this study, we proposed a practical statistical approach that use kmer frequency as the genetic
16 markers to quantify the overall statistical feature of the genomes/proteomes among different
17 species, and to associate genetic variants with large scale inter-species traits. The proteome data of
18 96 mammals was used as samples. The chromosome number of these species and the kmer
19 frequencies of each proteome were used as the phenotypic variables and the genetic markers
20 respectively. A total of 130 genes including *TP53* and *BAD* were identified significantly associated
21 with the chromosome variations. The annotation results of these genes show that they are related
22 to the cellular reaction of DNA double-strand breaks caused by chromosome fission/fusion. Our
23 study provides a new approach to perform GWAS on the large-scale inter-species trait, and it was
24 applied to analysis the genetic variants that are associated with chromosome number during the
25 evolution of species. In fact, besides chromosome number, any phenotypic traits, including those
26 used for characteristic of advanced taxa, could theoretically be studied with our approach,
27 identifying genetic variants for trait determination.

28

1 **Results**

2 **Chromosome number statistics**

3 In mammals, the chromosome number changed variously among mammalian species [16], ranging
4 from $2n=102$ (*Octomys mimax*) [17,18] to $2n=6$ (female in *Muntiacus muntjak*) [19], with the
5 putative number of $2n=46$ in mammalian common ancestor [16]. In present work, we collected the
6 chromosome number from 96 mammals (Table S1) ranging from $2n=6$ (female in *Muntiacus*
7 *crinifrons*) to $2n=84$ (*Ceratotherium simum simum*).

8

9 **Kmer frequencies calculation and quality control**

10 The AF-based software of CVTree [11] was used to construct the phylogenic tree of the 96
11 mammals. The species was selected on the basis of the genome availability in public database and
12 the chromosome number availability in the previous studies (*Materials and Methods*, Table S1).
13 The set of these 96 species contain the representative of four orders (Cetartiodactyla, Carnivora,
14 Rodentia and Primates). The phylogenic analysis based on AF sequence alignment showed that
15 $k=7$ oligopeptide (heptapeptide) was able to obtain the best evolutionary tree under four branches
16 (Fig. 1). The human protein sequences were used for heptapeptide-grouping because of the quality
17 of the genome and the completeness of the annotation. According to the quantify control criteria, a
18 total of 82816 groups of heptapeptides which included 2529430 types of unique heptapeptides
19 were picked up for the further analyses.

20

21 **Assessment of kmer contribution to chromosome number variance**

22 Using random forest method, we estimated the contribution of nearly 2.5 million heptapeptides to
23 the variance of the chromosome number within the mammalian lineage. The results indicated that
24 the total heptapeptides explained 18.1% of the chromosome number variance , which verified the
25 substantial role of the heptapeptide copy number for the chromosome number variation.

26

1 **Association analysis**

2 For the association analysis, the Manhattan and QQ plots shown in Fig.2. Table 1 present the
3 significant genes detected by our new approach, including their starting and ending positions in
4 the human genome, the ensemble IDs, the normal-P values and the q values after FDR test. In total,
5 we discovered 6 proteins from 5 genes associated with chromosome number (normal P-values <
6 0.001 and FDR q-value < 0.05). Of these five genes, the *BAD* had a normal-P value of 1.32E-06,
7 showing the strongest evidence of association with the chromosome number. Moreover, a total of
8 130 proteins which belong to 84 genes were identified to be associated with the chromosome
9 number (P values < 0.001 without FDR test) (Supplemental Table S2).

10

11 **Relations between oligopeptide copy number to domain function**

12 Of the 130 significant proteins, the positions covered by the heptapeptides are mainly intrinsic
13 disorder (ID) comparing to other pfam domain regions (Figure 3A). The comprehensive roles of
14 ID in protein interaction network has been discovered in the previous studies [20,21]. Some
15 investigations indicated a crucial role of ID for P53-induced apoptosis [21]. P53 is one of the most
16 well known and intensively-studied tumor suppressors, which plays the central role in apoptosis
17 responding to a broad variety of cell stresses including the DNA double-strand breaks (DSB).
18 Using P53 as an example, we found ID plays a role of bridge between oligopeptides and proteins
19 (Figure 3B, 3C). In Figure 3B both dog and Chinese hamster are species with the chromosome
20 number which is distinct from mammal ancestor, while human chromosome number are quite
21 close to mammal ancestor. It could see that the majority of kmers (except “SGTAKSV”) have a
22 consistent trend of frequencies in dog and Chinese hamster whose chromosome number has
23 deviated from mammal ancestor. There were 4 kmers mapped in the ID region of the gene *P53*, all
24 of which appeared also in ID region of other proteins in human proteome (Figure 3C). These
25 proteins with the same ID region may contribute to the function of P53.

26

27 **Function enrichment results**

1 To further analyze the function of the significant genes, we conducted the KEGG pathway
2 analysis by using the Genetrial2 platform. The results showed that these 130 proteins were
3 enriched in cancer (hsa05200) and Glycerophospholipid metabolism (hsa00564) (Table S3). The
4 Related genes of *BAD*, *TP53* and *BRCA2* have been proved to play an essential role in cell
5 apoptosis induced by DSB (Figure 4A); and the genes of *GDKQ*, *PLD1* and *CHPT1* are related
6 with the biosynthesis and the metabolism of Phosphatidic acid, which also takes part in the
7 pathways of apoptosis (Figure 4B).

8

9 **Discussion**

10 We performed a gene-based GWAS analysis on 96 mammal species in this study. However, a
11 critical constraint in GWAS is the sample size, which is also a potential constraint for the
12 association study in large-scaled inter-species trait. Typically, a GWAS study needs a very large
13 sample size to detect loci varies from phenotype to phenotype, because each SNP makes such a
14 small contribution to the whole variability in phenotype of the trait [15,22]. In the case of
15 large-scaled inter-species trait, the species number should also be large. However, only 96
16 sequences with chromosome numbers were available in the public databases. With the further
17 development of sequencing technologies, more and more high-quality genome could be acquired.
18 this will further provide more data for our study.

19

20 In this study, the kmer frequencies were used to describe the overall feature of genome, which is
21 similar to the AF sequence comparison studies. Generally, the kmer used for AF methods can be
22 either oligonucleotide or oligopeptide so that both genome/transcriptome and proteome data could
23 be used for AF phylogeny study. However, computational complexity was extremely high with
24 oligonucleotide of kmer due to the large size of the mammal genomes. So, many AF phylogeny
25 studies always choose to select oligopeptide in proteome as kmer for complicated multicellular
26 organisms [23]. Here, the oligopeptide frequencies were extracted through CVTree program [7]
27 for association study from the mammalian proteome data. In addition, another essential issue is to
28 determine what length of kmer could reflect the overall feature of proteome. In CVTree, each

1 length of oligopeptide could obtain a distance tree and be filtered out by the criterion based on the
2 accepted phylogenic topologies. As mentioned above, the genetic variation determining the
3 variation of large-scaled inter-species traits was also recorded in kmer characteristics. So, it is
4 reasonable that the length of kmer whose frequencies could obtain the best tree should also obtain
5 the most amount of information of variation to determine the trait. Therefore, we used the length of
6 kmer which determined the best tree for further trait association study.

7

8 The connections between the oligopeptide copy number and the targeted proteins, and between the
9 target proteins and the biology context of the trait determination needs further explanation. The
10 heptapeptides are enriched mainly in low complexity or intrinsic disorder (ID) [24]. These
11 unstructured regions in protein have been discovered to be far more important than previous
12 considered [25]. Critical functional involvement includes phosphorylation [26], alternative
13 splicing [27], protein-protein interaction in interactomes [28], and so force. Therefore, the process
14 can be modeled in three steps. Firstly, the changes of the heptapeptide copy number caused the
15 changes of the ID region appearing in different proteins. Secondly, different ID region distribution
16 in different proteins produces different protein-protein interaction or other function alteration of
17 the target proteins. Finally, the altered function of the related proteins in the proteome caused the
18 variation of the traits.

19

20 Further explanation on the chromosome number determination is required, particularly on the
21 function enrichment result. The changes of chromosome number need to be fixed in population,
22 during which the DNA double-strand breaks (DSB) may increase in the cells of individuals in the
23 population. As we know, the increasing DSB may trigger apoptosis. Therefore, for the survive of
24 individuals with different chromosome number, a mechanism which can evade apoptosis under
25 high level of DSB is required. The critical role of DSBs during cancer development [29,30],
26 peculiar the most recent advances in the relation of DSBs and chromothripsis [31], has offered
27 some indirect evidences for identifying the role of DSBs on molecular mechanism about the
28 changes of chromosomes number in evolution.

1

2 To better explain the whole-genome role in trait determination, we introduce an analogy involving
3 the comparison of statistical physics. For system consist of large number of microscopic particles,
4 although the motion of single particle could be accurately described with mechanics principle, the
5 macroscopic phenomenon must be quantified with the average statistical feature of particles [32].
6 As far as the genome case is concerned, it is the similar one. Although the function of single gene
7 obeying to central dogma, the macroscopic heritable trait was an average result of a large number
8 of genes in the complex interaction network. In physics, it has been an implicit view of point that
9 dynamics of single particle could not explain the macroscopic feature of a system with numerous
10 particles [32]. In biology, however, it is a rather radical or revolutionary opinion that traits in
11 individual level must be determined by global feature of the whole genome instead of the local
12 effects of individual genes.

13

14 **Conclusion**

15 By using kmer frequencies as the genomic marker to replace SNPs, we extended genome-wide
16 association studies for genetic analysis of the large-scale inter-species trait. Such an approach can
17 be used not only for identifying traditional complex traits within species, but also the large-scaled
18 inter-species complex traits. In this work we took chromosome number variation within mammals
19 as a case to show the availability and the advantages of this approach. We offered a biological
20 explanation for how the overall feature of sequences could participate in the determination of traits,
21 which is the kmers associated with certain phenotype may cover a stretch of continuous
22 polypeptide sequences in protein sequences. These stretch of sequence may work as some
23 functional domain such as ID region which changes its location and the copy number to influence
24 the protein-protein interactions. So that the pattern of the protein interaction network may be
25 changed. And finally the phenotype may be affected. This study has provided a practical method
26 for genetic analysis of the large-scale inter-species trait.

27

1 **Materials and methods**

2 **The acquisition and filtration of protein sequence**

3 We selected 96 mammals with annotated proteome data. The newly artiodactyla genome project
4 offered a dataset of 40 newly sequenced genomes (unpublished). Using its annotation system, 8
5 genome assemblies from NCBI were re-annotated to obtain a refined proteome sequence. 48
6 proteome sequences were directly downloaded from NCBI.

7

8 For kmer frequency analysis, protein sequence from the 96 proteomes were filtered according to
9 two criteria: 1) the length of protein should be longer than 150 amino acid residues, and 2) For
10 multiple proteins from one gene caused by alternative splicing, we only selected the longest one.

11

12 **Determination for the length of oligopeptide**

13 The alignment-free methods for evolutionary tree construction were used to determine the best k
14 value of the orthologous oligopeptide for copy number counting. Although not applied as widely
15 as alignment method such as BLAST, alignment-free methods have long been developed to
16 analyze sequences [16,33]. Different researchers use different terms to description a class of
17 similar methods. Here we used the well-known term of kmer frequency in computational biology.
18 One disparity is that kmer in computational biology usually denotes oligonucleotides, while in
19 alignment-free method it could denotes either oligopeptide or oligonucleotide. These methods
20 counted the frequencies of kmers for DNA sequence (or protein sequence), defining composition
21 vector with the frequencies to describe the statistical feature of the sequence, constructing a
22 random background to exclude from the vector, and then defining a measurement as the
23 evolutionary distance to construct a distance matrix for concerned sequences (species). Finally NJ
24 tree could be used to describe the evolutionary relation of these sequences.

25

26 We used the CVTree to optimize k value for the mammalian proteome data above. Considering the

1 restriction of the computer power, CVTree sets the k value of protein sequence from 3 to 14.
2 Different k values could generate different NJ trees. The random background was defined as a
3 Markovian chain which could be excluded by subtraction parameter. The optimized criterion is
4 designed with three aspects: 1) the four superorders within Eutheria should be correctly clustered,
5 which are Euarchontoglires, Laurasiatheria, Xenarthra and Afrotheria. 2) the monophyletic orders
6 within Euarchontoglires and Laurasiatheria should be correctly clustered, including Carnivora,
7 Cetartiodactyla, Chiroptera and Perissodactyla in Laurasiatheria; and Glires and Primates in
8 Euarchontoglires, and 3) the phylogeny of the monophyletic orders within their respective
9 superorders should be consistent with the general accepted phylogenomic tree. All of the kmer
10 frequency for each species and the distance matrix for NJ tree were the output results by the
11 CVTree program.

12

13 **Definition of protein-kmer matrix and quality control**

14 In GWAS, it is a general strategy to group SNP within one gene (or transcript) together to perform
15 linear regression. Here we used the similar methods with the aim of grouping oligopeptides into
16 protein sequences to identify the association between protein and the phenotype. We used the
17 protein sequences of the unfiltered human proteome downloaded from the Ensembl Genes version
18 90 as the reference to group oligopeptides. Each group corresponds to the oligopeptides appearing
19 in one human protein sequence. Theoretically, the grouping should be random. But this would
20 produce a vast number of groups which is impossible for association computation.

21

22 In CVTree results, the frequency value of -1.0 or 0.0 was used to denote missing data. We
23 removed these data for oligopeptide groups with the following criterion: 1) All kmers with the
24 missing data < 90% and all transcripts with kmers <5 were filtered out, 2) All the species with the
25 missing data >10% were ruled out for each transcript.

26

27 **Phenotypic preparation**

1 The mammalian karyotype has been extensively investigated since 1960s [16]. Since ever, some
2 species were renamed, or their taxonomic places were changed. This makes some difficulties to
3 get the chromosome number information of the species with the available proteome data. We
4 manually checked related literature for these data. There are also some species with intra-species
5 variation of the chromosome number, the numbers are rounded to the nearest integer average for
6 these cases. We can see the details of chromosome number in Supplementary Table S1. To
7 eliminate the influence of species difference, the phenotypic values of the chromosome number
8 was adjusted by the evolutionary relation of the 96 mammals. In other words, the residuals after
9 fitting the above distance matrix from CVTree were treated as the phenotypic values of traits for
10 the further association studies.

11

12 **Variance component analysis**

13 The analysis is conducted through the R script programmed with R language version 3.4.3. As
14 the core part of the program, RandomForest (version 4.6-14) was used for the percentage of
15 variance explained calculated which can be directly provided from the R script outcome.

16 **Association analysis**

17 For each group (protein or transcript) of oligopeptides, we constructed principal components (PCs)
18 according to the following rules [14]. We treated each oligopeptide (kmer) within a group as a
19 single variable, and calculated the variance-covariance matrix. Then we calculated the eigenvalues
20 and the eigenvectors of the covariance matrix. PCs were obtained by multiplying the eigenvectors
21 by the kmer frequency matrix. Thus, PCs were mutually independent within a protein sequence,
22 and could be treated as the independent variables for regression analysis and significance tests.
23 After these steps, multiple regression analysis for the adjusted chromosome number on the
24 oligopeptides (kmers) was turned to a simple regression by using the PC. A simple correlation
25 coefficient corrected by chromosome numbers and PCs could be used to calculate the P-value of
26 the PC.

27

1 In the Fisher's combination test, χ^2 was constructed by combining K independent p values for PCs,
2 as follows:

$$3 \quad Z_F = -2 \sum_{i=1}^K \log P_i$$

4 This test statistic followed a χ^2 distribution with the degree of freedom $2 \times K$ under the null model,
5 and a new p-value was calculated from this Chi-square distribution for each transcript.

6

7 The false discovery rate test was adopted for multiple tests of all groups (proteins) in the proteome.
8 Proteins whose nominal q value less than 0.05 was considered as core associated gene set at the
9 proteome-wide significance level. Besides, proteins with nominal p value less than 0.001 was
10 categorized under the peripheral associated gene set.

11

12 **Case study of ID region in P53**

13 In the association study, one of P53 transcripts (Ensembl number ENSP00000426252) was used
14 and showing significant association with chromosome number, in which 6 heptapeptide were
15 selected after filtering absences in the 96 mammals. So, the selected three species with the
16 representative chromosome number, human (2n=46), dog (2n=78), Chinese hamster (2n=22)
17 showed the subtracted heptapeptide frequencies. The heptapeptide locations were identified in
18 SMART (URL see below). The frequency was subtracted with the random background in CVTree,
19 and the result values were compared among the three species. Not that the heptapeptide of
20 "AAPAPAP" was absent in dog and Chinese hamster, so only five were shown in Figure 3B. For
21 all other proteins, the domain identification of proteins was all performed in SMART.

22

23 **Function enrichment and analysis**

24 We used the peripheral associated gene set for a function enrichment with Genetrial2 and for a
25 protein interaction network analysis with STRING. For Genetrial2 analysis, the enrichment

1 algorithm used was Over-representation analysis, the adjustment method was FDR adjustment
2 offered by Genetrial2 with significance level of 0.1. We excluded results of pathways which
3 include no gene from the core associated gene set. For STRING analysis, the low confidence for
4 minimum required interaction score (>0.150) was used. The subcellular location information was
5 obtained in GeneCards.

6

7 **Urls used in Methods**

8 NCBI, <https://www.ncbi.nlm.nih.gov/assembly>. Ensembl <http://www.ensembl.org/index.html>.
9 SMART, <https://smart.embl-heidelberg.de/>. Genetrial2, <https://genetrial2.bioinf.uni-sb.de/>.
10 STRING, <https://string-db.org/>. GeneCards, <http://www.genecards.org/>.

11

12 **Authors' contributions**

13 Qi Wu and Fuwen Wei designed the research. Qi Wu and Lei Chen prepared the data. Qi Wu and
14 Huizhong Fan performed the research and interpreted the data. Qi Wu, Huizhong Fan, Yibo Hu
15 and Fuwen Wei wrote the manuscript.

16

17 **Competing interests**

18 The authors declared no competing interests of this work.

19

20 **Acknowledgments**

21 We thank Prof. Qiang Qiu for the helpful discussion of the project and Dr. Xinhai Li for the help
22 and assistance in random forest method application. This work is supported by grants from
23 National Key Program of Research and Development, Ministry of Science and Technology Grant
24 (No. 2016YFC0503200), the Strategic Priority Research Program of the Chinese Academy of

1 Sciences (No. XDPB0202), the Key Research Program of the Chinese Academy of Sciences (No.
2 KJZD-EW-L07) and the NSFC Major Scientific Research Program (No. 91746119).

3

4 **References**

- 5 [1] Klein RJ, Zeiss C, Chew EY, Tsai JY, Sackler RS, Haynes C, et al. Complement factor H
6 polymorphism in age-related macular degeneration. *Science*. 2005 Apr 15;308(5720):385–9.
- 7 [2] Visscher PM, Wray NR, Zhang Q, Sklar P, McCarthy MI, Brown MA, et al. 10 Years of
8 GWAS Discovery: Biology, Function, and Translation. *Am J Hum Genet*. 2017 Jul 6;101(1):5–
9 22.
- 10 [3] Hayes B, Goddard M. Genome-wide association and genomic selection in animal breeding.
11 *Genome*. 2010 Nov;53(11):876–83.
- 12 [4] Bartoli C, Roux F. Genome-Wide Association Studies In Plant Pathosystems: Toward an
13 Ecological Genomics Approach. *Front Plant Sci*. 2017;8:763.
- 14 [5] Fisher RA. The correlation between relatives on the supposition of Mendelian inheritance.
15 *Trans R Soc Edinb*. 1918 Oct 1; 52:399–433.
- 16 [6] Zieleszinski A, Vinga S, Almeida J, Karlowski WM. Alignment-free sequence comparison:
17 benefits, applications, and tools. *Genome Biol*. 2017 Oct 3;18(1):186.
- 18 [7] Xu Z, Hao B. CVTree update: a newly designed phylogenetic study platform using
19 composition vectors and whole genomes. *Nucleic Acids Res*. 2009 Jul;37(Web Server
20 issue):W174-178.
- 21 [8] Chan CX, Bernard G, Poirion O, Hogan JM, Ragan MA. Inferring phylogenies of evolving
22 sequences without multiple sequence alignment. *Sci Rep*. 2014 Sep 30;4:6504.
- 23 [9] Bernard G, Chan CX, Chan Y-B, Chua X-Y, Cong Y, Hogan JM, et al. Alignment-free
24 inference of hierarchical and reticulate phylogenomic relationships. *Brief Bioinform*. 2017 Jun
25 30;
- 26 [10] Deng M, Yu C, Liang Q, He RL, Yau SS-T. A novel method of characterizing genetic
27 sequences: genome space with biological distance and applications. *PLoS One*. 2011 Mar
28 2;6(3):e17293.
- 29 [11] Karlin S, Campbell AM, Mrazek J. Comparative DNA analysis across diverse genomes. *Annu*

- 1 Rev Genet. 1998;32:185–225.
- 2 [12] Chan CX, Bernard G, Poirion O, Hogan JM, Ragan MA. Inferring phylogenies of evolving
3 sequences without multiple sequence alignment. *Sci Rep.* 2014 Sep 30;4:6504.
- 4 [13] Wu W, Zhang M, Zhu L, Wu Q. Construction and discussion of whole-genome phylogeny of
5 mammals by alignment-free method. *Acta Theriologica Sinica.* 2018;38(2):174-182.
- 6 [14] Xia J, Fan H, Chang T, Xu L, Zhang W, Song Y, et al. Searching for new loci and candidate
7 genes for economically important traits through gene-based association analysis of Simmental
8 cattle. *Sci Rep.* 2017 Feb 7;7:42048.
- 9 [15] Yang J, Benyamin B, McEvoy BP, Gordon S, Henders AK, Nyholt DR, et al. Common SNPs
10 explain a large proportion of the heritability for human height. *Nat Genet.* 2010 Jul;42(7):565-9.
- 11 [16] O’Brian S. *Atlas of Mammalian Chromosomes.* Hoboken: John Wiley & Sons Inc.; 2006.
- 12 [17] Gallardo MH, Bickham JW, Honeycutt RL, Ojeda RA, Kohler N. Discovery of tetraploidy in
13 a mammal. *Nature.* 1999 Sep 23;401(6751):341.
- 14 [18] Svartman M, Stone G, Stanyon R. Molecular cytogenetics discards polyploidy in mammals.
15 *Genomics.* 2005 Apr;85(4):425–30.
- 16 [19] Wurster DH, Benirschke K. Indian muntjac, *Muntiacus muntjak*: a deer with a low diploid
17 chromosome number. *Science.* 1970 Jun 12;168(3937):1364–6.
- 18 [20] Peng, Z., et al. (2013) Resilience of death: intrinsic disorder in proteins involved in the
19 programmed cell death. *Cell Death Differ.* 20(9):1257-1267.
- 20 [21] Xue, B., et al. (2013) Intrinsically disordered regions of p53 family are highly diversified in
21 evolution. *Biochim Biophys Acta.* 834(4): 725–738.
- 22 [22] Boyle EA, Li YI, Pritchard JK. An Expanded View of Complex Traits: From Polygenic to
23 Omnigenic. *Cell.* 2017 Jun 15;169(7):1177–86.
- 24 [23] Sims GE, Jun S-R, Wu GA, Kim S-H. Whole-genome phylogeny of mammals: evolutionary
25 information in genic and nongenic regions. *Proc Natl Acad Sci U S A.* 2009 Oct
26 6;106(40):17077–82.
- 27 [24] Dunker AK, Oldfield CJ, Meng J, Romero P, Yang JY, Chen JW, et al. The unfoldomics
28 decade: an update on intrinsically disordered proteins. *BMC Genomics.* 2008 Sep 16;9 Suppl
29 2:S1.

- 1 [25] Haerty W, Golding GB. Low-complexity sequences and single amino acid repeats: not just
2 “junk peptide sequences.” *Genome*. 2010 Oct;53(10):753–62.
- 3 [26] Iakoucheva LM, Radivojac P, Brown CJ, O’Connor TR, Sikes JG, Obradovic Z, et al. The
4 importance of intrinsic disorder for protein phosphorylation. *Nucleic Acids Res*.
5 2004;32(3):1037–49.
- 6 [27] Romero PR, Zaidi S, Fang YY, Uversky VN, Radivojac P, Oldfield CJ, et al. Alternative
7 splicing in concert with protein intrinsic disorder enables increased functional diversity in
8 multicellular organisms. *Proc Natl Acad Sci U S A*. 2006 May 30;103(22):8390–5.
- 9 [28] Cumberworth A, Lamour G, Babu MM, Gsponer J. Promiscuity as a functional trait:
10 intrinsically disordered regions as central players of interactomes. *Biochem J*. 2013 Sep
11 15;454(3):361–9.
- 12 [29] Norbury CJ, Hickson ID. Cellular responses to DNA damage. *Annu Rev Pharmacol Toxicol*.
13 2001;41:367–401.
- 14 [30] Khanna KK, Jackson SP. DNA double-strand breaks: signaling, repair and the cancer
15 connection. *Nat Genet*. 2001 Mar;27(3):247–54.
- 16 [31] Schipler A, Mladenova V, Soni A, Nikolov V, Saha J, Mladenov E, et al. Chromosome
17 thripsis by DNA double strand break clusters causes enhanced cell lethality, chromosomal
18 translocations and 53BP1-recruitment. *Nucleic Acids Res*. 2016 Sep 19;44(16):7673–90.
- 19 [32] Su R. *Statistics physics*. 2nd ed. Beijing: Higher Education Press; 2002.
- 20 [33] Alberts B, Johnson A, Lewis J, Raff M, Roberts K, Walter P. *Molecular biology of the cell*. in
21 *Molecular biology of the cell* 204 (Garland Science, 2015).

22

23 **Figure legends**

24 **Figure 1** Optimization for the length (k value) of orthologous oligopeptide. The NJ tree is given
25 by the value of k=7. Although the inner topology may be ambiguous or conflict, the 4 orders and
26 the 4 superorders of placental mammals shown with different colors were clustered right. And the
27 relation of the eight lineages were acceptable. The exception for position of superorder Xenarthra

1 could be seen.

2

3 **Figure 2** Overall results of the EGWAS. **A.** the plots of Manhattan, the 25 chromosomes
4 (including 22 autosomes, X, Y and the non-chromosome) are color coded. **B.** the
5 Quantile-Quantile plots of $-\log_{10}$ (p values) for chromosome, the observed negative logarithms of
6 the p values in association study are plotted against their expected values under the null
7 hypothesis of no association with the phenotype.

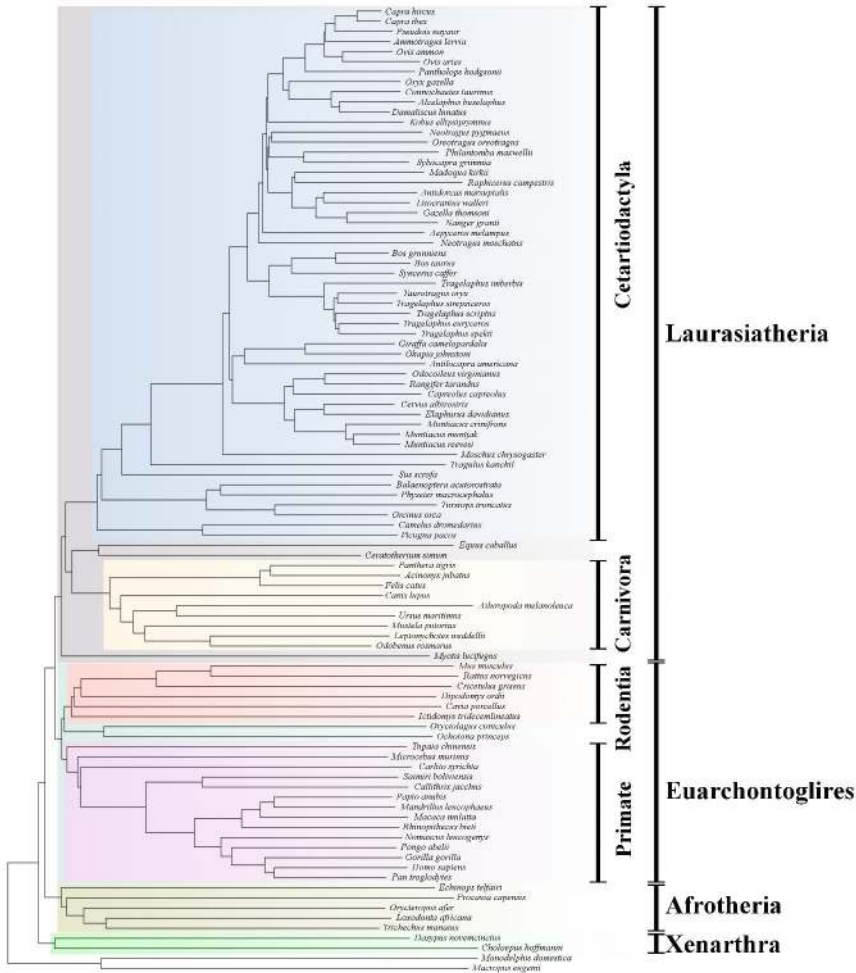
8

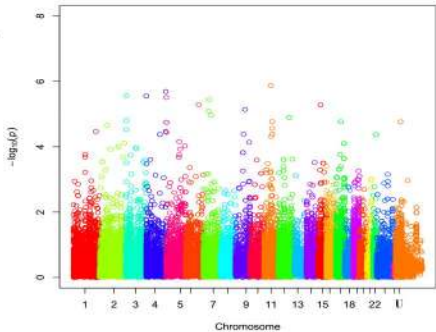
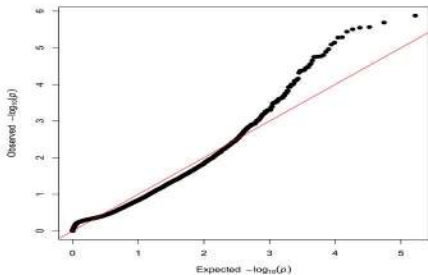
9 **Figure 3** The relation from oligopeptide to function of proteome with the ID as a bridge. **A.** The
10 ratio of ID in all domains mapped with significantly associated heptapeptide. **B, C, And D**
11 Showed the detailed result of the P53 protein (ENSP00000426252). The coverage of six
12 heptapeptides on the functional domain in human protein are indicated in the center, involving a
13 ID region (purple) and the C-terminal of the P53 PFAM domain. The position of the domains and
14 the oligopeptides in proteins were indicated. The four heptapeptides in the ID region were also
15 appeared in other 72 proteins (supplementary Table S4), six of which covered by at least two of
16 the heptapeptides were shown around. None of the 72 was in the associated gene set, but some of
17 them have had bound evidence to be functional related with P53 such as CDKN1C40.

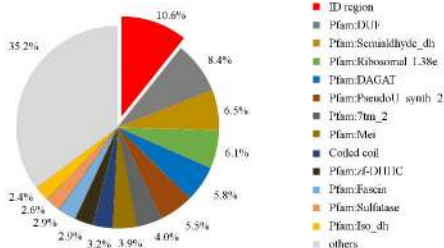
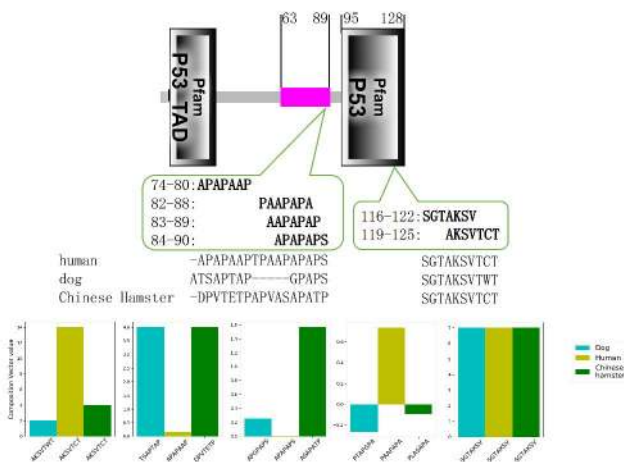
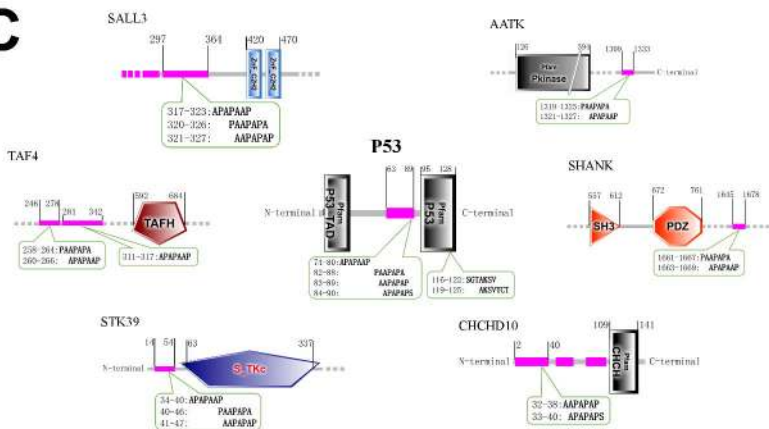
18

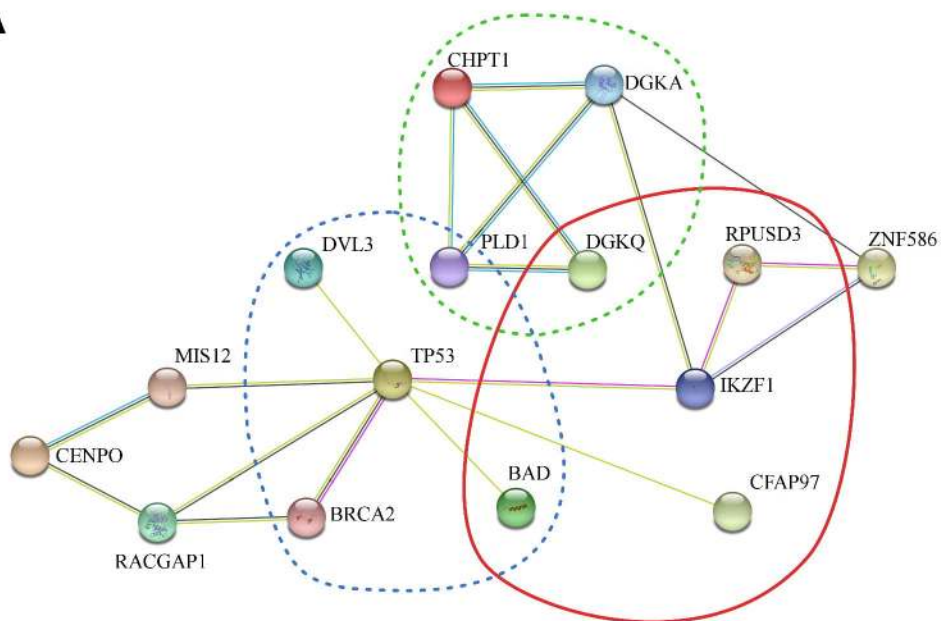
19 **Figure 4** Intensive function analysis for the candidate associated genes. **A.** The refined protein
20 interaction network of related genes modified from STRING website. The balls indicated proteins
21 while the lines connecting balls indicated interaction between proteins. The colors of balls and
22 lines were default legends from STRING result. The genes in the red solid rings in the right
23 indicated the five genes in the core associated gene set. The region marked with blue dotted lines
24 in the left indicated the five genes in the pathways related to cancer. The region marked with green
25 dotted lines above indicated the four genes in the pathways of Glycerophospholipid metabolism.
26 **B.** The related genes and apoptosis modified from KEGG. The black circle denoted to
27 Glycerophospholipid molecules with names and molecule structures marked nearby respectively.
28 The gene encoding enzymes for related chemical reactions were illustrated in blue boxes. The

- 1 majority molecule related to apoptosis in cancer was phosphatidate (phosphatidic acid). TP53 and
- 2 BAD may play essential role during the process to evade apoptosis.
- 3



A**B**

A**B****C**

A**B**