1

## Prometheus: omics portals for interkingdom comparative genomic analyses

3

4 Gunhwan Ko[†], Insu Jang[†], Namjin Koo[†], Seong-Jin Park, Sangho Oh, Min-Seo Kim, Jin-Hyuk Choi,

5 Hyeongmin Kim, Young Mi Sim, Iksu Byeon, Pan-Gyu Kim, Kye Young Kim, Gukhee Han, Jong-Cheol

6 Yoon, Yunji Hong, Kyung-Lok Mun, Banghyuk Lee, Deayeon Ko, Wangho Song, and Yong-Min Kim*

7

8 Korean Bioinformation Center, Korea Research Institute of Bioscience and Biotechnology, Daejeon

9 34141, Republic of Korea

10

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

12

13 *To whom correspondence should be addressed. Tel: +82-42-879-8534, Fax: +82-42-870-8519, E-mail:

14 ymkim@kribb.re.kr

15

16

17

18

19

20

21

22

23 **Running Title:** Portal for interkingdom comparative genomics

24

25 **Keywords**

26 Comparative genomics, Molecular evolution, Interkingdom analysis, Domain architectures-based gene

27 search, Web-based platform

28

**Abstract**

Functional analyses of genes are crucial for unveiling biological responses, for genetic engineering, and for developing new medicines. However, functional analyses have largely been restricted to model organisms, representing a major hurdle for functional studies and industrial applications. To resolve this, comparative genome analyses can be used to provide clues to gene functions as well as their evolutionary history. To this end, we present Prometheus (http://prometheus.kobic.re.kr), web-based omics portal that contains more than 17,215 sequences from prokaryotic and eukaryotic genomes. This portal supports interkingdom comparative analyses via a domain architecture-based gene identification system, Gene Search, and users can easily and rapidly identify single or entire gene sets in specific pathways. Bioinformatics tools for further analyses are provided in Prometheus or through BioExpress, a cloud-based bioinformatics analysis platform. Prometheus suggests a new paradigm for comparative analyses with large amounts of genomic information.

**Introduction**

The completion of the Human Genome Project (2003) was not an end but rather a new beginning for further functional genomic analyses. The ENCyclopedia Of DNA Elements (ENCODE) was launched to begin investigating the functions of the identified human genes[1]. In addition, large-scale functional studies, such as interactome or network analyses, were performed in model organisms, including Arabidopsis thaliana, Saccharomyces cerevisiae, and Drosophila melanogaster. These efforts accumulated network information on various interactomes and gene functions. These vast amounts of biological information enabled functional studies that contributed to the unveiling of biological responses, the cloning of genes of interest, and the development of molecular markers for model organisms or medicines in humans[2,3]. Thus, the trend of functional analyses has been transferred from candidate gene research to genome-wide research. However, this flood of information has largely been restricted to model organisms, and it has been challenging for researchers to apply these data to newly sequenced genomes.

Since next-generation sequencing (NGS) technology was developed in the mid-2000s, an enormous amount of genomic information has been analyzed and amassed in public databases. As the numbers of sequenced genomes increased, many tools and pipelines were developed to investigate gene functions, identify gene families, and perform comparative genomic analyses. However, the application of comparative analyses is restricted to functional gene annotations and newly sequenced genome analyses. Newly sequenced genomes are initially compared to those that have previously been analyzed, including genomes of closely related species, to provide information on genome structure changes and gene repertoires. Such comparisons can also predict gene paralogues, which are genes related by duplication events, or orthologues, which are those related by speciation events[4-6]. As orthologues tend to be more similar in function that paralogues[7], they are widely used for functional gene annotations[8]. Moreover, recent gene-of-interest studies that include multigenome orthologues offer insight into their mechanisms for adapting to the environment[9,10]. However, these comparative genomic analyses were performed at genome-, genus-, or kingdom-wide  levels, thereby restricting comparisons to the species, family, or order level[11-13]. To understand the evolution of genes of interest more precisely, interkingdom analyses are needed, particularly because many genes in eukaryotic genomes have universal common ancestries in Bacteria and Archaea[14].

Here, we report an omics portal for interkingdom comparative genomic analyses named Prometheus (http://prometheus.kobic.re.kr). We collected 17,215 sequences from 16,730 species and constructed four primary databases to provide basic genome information, with more detailed information on individual genes provided in secondary databases. Researchers can then access detailed information on genes of interest, such as gene structure, domain architecture, subcellular localization, orthologues, and paralogues, as well as their sequences. In particular, Prometheus provides Gene Search to identify genes

3

1   of interest based on their domain architectures from prokaryotes to eukaryotes and performs various

2   comparative analyses, such as comparison of chromosome sequences, sequence alignment, and

3   phylogenetic analyses. Furthermore, researchers can perform various bioinformatics analyses with these

4   and their own sequencing data in a cloud-based platform, BioExpress. Prometheus suggests a new

5   paradigm for genome research, from single genes of interest to entire gene pathways.

6

7   **Materials and Methods**

8   **Web interface**

9   Prometheus furnishes data search, configuration of data analyses, data visualization, and storage

10  of users' own data. The interface is implemented using a Hypertext Markup Language (HTML),

11  cascading style sheets (CSS) and uses a jQuery JavaScript library (jQuery) to modify web page contents.

12  To visualize data, dynamic web interface is constructed by Asynchronous JavaScript and XML (Ajax)

13  using JavaScript Object Notation (JSON) data format. Furthermore, genome browser was constructed

14  using Scalable Vector Graphics (SVG) and phylogenetic viewer is constructed using JavaScript. Web

15  interface of Prometheus supports a cross-browsing.

16

17  **Construction of taxonomy combined heatmap of photolyase/cryptochrome family**

18  Sequences for the photolyase/cryptochrome family of genes from different species in previous

19  study[15] were collected and domain architectures were investigated using InterProScan v5.0[16]. Each of the

20  subtypes reported in previous studies were investigated using Gene Search in Prometheus. The numbers

21  of each of the subfamily genes were calculated for individual species and visualized as a heatmap using R

22  scripts. The taxonomic tree was constructed using phyloT in iTOL[17], an online tool that generates

23  phylogenetic trees based on the NCBI taxonomy. Finally, the taxonomic tree and heatmap were combined

24  using Adobe Illustrator.

25

26  **Bioinformatics analysis using a cloud-based analysis system, BioExpress**

27  LAST[18], BLAST[19], Clustal Omega[20], MUSCLE[21] and InterPro[16] programs are run in hybrid-

28  cluster system, BioExpress. To support further genomic analyses using personal data such as RNA-Seq,

29  Chip-Seq, or genome resequencing data, Prometheus links to BioExpress and users can perform further

30  various genomic analyses using personal data in My Gene and various analysis pipelines in BioExpress.

31  BioExpress is constructed by Hadoop to support high-speed analysis of a large amount of data. To

32  maintain a large data of user, Prometheus stores the data divided by optimized block size using Hadoop

33  Distributed File System (HDFS) into various computer servers. These storage system can maintain three

34  copies of user data and provides stable data storage by reducing risk of data loss. Web server of

1    Prometheus transmits task, progress and result of data analysis to BioExpress server using apache thrift

2    library-based Remote Procedure Call (RPC) and received result as JSON format data. The result of

3    genomic analyses is stored in HDFS and downloaded in web browser using HTTP protocol. In case of

4    large amount of data, users can download their data using KoDS (KOBIC Data Transfer Solution). KoDS

5    is a high-speed file transfer software using TCP/IP protocol and transferred user data is stored in HDFS.

6

7    **Construction of Database**

8        The database consisted of primary and secondary data tables in the Prometheus was constructed

9    using MySQL database management system. In database, primary data tables were created through data is

10   opened in five public databases and secondary data tables were constructed by parsing results of

11   bioinformatics tools such as InterProScan[16], OrthoMCL[22], MultiLoc2[23] and TargetP[24]. Detailed methods

12   for construction of database are described in Supplemental Note Section 1.

13

14   **Results**

15   **Concept and construction of Prometheus**

16       Prometheus provides an integrated pipeline for interkingdom comparative genomic analyses and

17   comprises four major sections, Genome Archive, Gene Search, BioExpress, and Genome Analysis. Users

18   can identify genes of interest using Gene Search and investigate their domain architectures using

19   InterPro[16] in Genome Analysis. Furthermore, users can obtain additional species information via

20   accessing the Korean Bioresource Information System (KOBIS) or perform further analyses by accessing

21   the cloud-based BioExpress (Figure 1).

22       To establish Prometheus, 17,215 sequences from 16,730 species were collected and stored in four

23   primary databases. The genomic information in Genome Archive (Figure 2A and Supplementary Table 1)

24   is arranged by taxonomic rank (obtained from NCBI), which users can access by clicking the species

25   name in the taxonomic tree or using a key word search. This General Information provides details on

26   genome assembly, annotation, and taxonomy. In eukaryotic genomes, distinct versions of genome

27   assembly and annotation were provided, and so each version is stored separately (Figure 2B). In

28   prokaryotic genomes, genomic information is separated by strain to support metagenomics analyses.

29   Genomes were classified according to criteria from RefSeq, which provided most of the genomic data

30   (Supplementary Table 1), to construct the database and to visualize the genomic information. In total, 435

31   eukaryotic genomes, 15,984 prokaryotic genomes, and 311 archaea genomes were collected and

32   assembled into the four primary databases containing information on assembled genomes, general feature

33   formats (GFFs), coding sequences (CDSs), and protein sequences, for a grand total 213,478,449 records

34   (Supplementary Table 2). Taxonomic information in Genome Archive is stored in a taxonomy database

5

1   and general information of genome assembly and annotation is stored in a genome report databank.

2   Totally, 51 database were constructed with 1,163,053,603 records (Supplementary Tables 3, 4, and 5).

3   General information on individual genomes is obtained using Genome Browser (Figure 2C), with

4   zoom in/out functions ranging from 100% to 1,000% and a gene search function by position or gene

5   name. Users can access and download the individual gene's information (CDS and peptide sequence) by

6   key word search or by clicking within the Genome Browser. Detailed information on individual genes is

7   provided in Gene Viewer (Figure 2D), and users can access the Genome Browser or result pages in Gene

8   Search. Bioinformatics analyses, including InterPro[16], OrthoMCL[22], MultiLoc2[23], and TargetP[24] were

9   performed using protein sequences of each species to construct six secondary databases, which are

10  presented in separate sections within the Gene Viewer (Figure 2D). In summary, genomic information

11  collected from major public databases is contained in the Genome Archive, and integrative information of

12  genomes or individual genes from individual databases is accessed via Gene Viewer.

13

14  **Analyses of transcriptional factors and tricarboxylic acid (TCA) cycle in Gene Search**

15  The major function of Prometheus is to perform interkingdom comparative analyses. To support

16  this objective, secondary databases containing information on domain architectures and

17  orthologues/paralogues of individual genes were constructed. We validated the utility of Prometheus by

18  performing an interkingdom investigation of transcription factors (TFs) and genes involved in the TCA

19  cycle using Gene Search (Figure 3, Supplementary Tables 6, and 7). A pipeline (iTAK v1.7)[25] was used to

20  identify plant TFs and classify protein kinases. TFs, transcriptional regulators (TRs), and kinases were

21  identified by consensus rules mainly summarized from PlnTFDB[26], PlantTFDB[18] with families from

22  PlantTFact[27], and AtFDB[28]. Domain architectures of each TF were investigated using InterProScan, and

23  their domain architectures depicted by InterPro entry (IPR) terms were used for further analyses using

24  Gene Search. To provide additional information about identified genes, the number of domain subtypes

25  are depicted in a summary table in Gene Search and as a header of sequence data in a FASTA file

26  (Supplementary Figure 1). Users can categorize identified genes into each subtype. We identified and

27  validated 79,960 genes from 15 gene families using the iTAK v1.7 (Figure 3A and Supplementary Table

28  6)[25]. The accuracy of our Gene Search ranged from 86.03% to 99.98%, with an average accuracy of

29  96.41%. High rates of accuracy were observed for genes encoding TFs containing significant IPR terms,

30  such as FAR1, MADS-box, NAC or Dof domains, whereas those for TFs without significant IPR terms,

31  such as B3-type TFs or CAMTA, showed lower rates. Thus, these data suggest that specific IPR terms or

32  exact domain architectures are required to enhance the accuracy of Gene Search.

33  Genes involved in the TCA cycle were further investigated with Gene Search to demonstrate the

34  potential for applying comparative genomics at the pathway level. As the TCA cycle is a fundamental

1    metabolic pathway for survival in prokaryotes and eukaryotes, we selected this for an interkingdom

2    comparative genomic analysis. A total of 435,044 genes were identified from 20 individual genes in the

3    TCA cycle using Gene Search, and the ratios of species harboring each gene in the TCA cycle were

4    shown as heatmaps (Figure 3B and Supplementary Table 7). These results showed that some genes, such

5    as those encoding isocitrate dehydrogenase (IDH1 and IDH2) and malate dehydrogenase (MDH1 and

6    MDH2) evolved in a lineage-specific manner. Furthermore, the results show the lineage-specific rates of

7    functionally redundant genes, such as those encoding succinate dehydrogenase and succinyl-CoA

8    synthase. This investigation of the TCA cycle also provided information on the gene repertoires and the

9    evolution of the TCA cycle in each kingdom. Thus, Prometheus provides information for evolutionary

10   studies of single genes or those in specific pathways, including the distributions and rates of genes, as

11   well as repertoires of gene orthologues in pathways. In addition, Prometheus provides the domain

12   architectures of genes as well as their CDSs and/or peptide sequences.

13

14   **Tools for comparative analyses and personalized management system via My Genes in Prometheus**

15        To support comparative analyses in Prometheus, essential tools such as LAST (a program for

16   comparing sequences at the chromosome level)[29], Basic Local Alignment Search Tool (BLAST)[19], and

17   InterPro are provided in Genome Analysis (Supplementary Figure 2). Users can monitor the progress of

18   analysis in a personalized page, My Genes (Supplementary Figure 3), and download the result files from

19   each program via a file menu. In the case of data from InterProScan, the result file is shown in a graphic

20   format and results are downloaded in a tsv file format (Supplementary Figure 4). Thus, users can

21   investigate domain architectures of genes of interest and perform interkingdom identification using Gene

22   Search.

23        We performed a comparative analysis of genes in the photolyase/cryptochrome family using a

24   gene set from a previous study[15] as a control (Figure 4 and Supplementary Table 8). The domain

25   architectures of photolyase/cryptochrome subfamilies are the same and family IPR terms are different

26   (Figure 4A), enabling a more accurate identification of each subfamily. The results also indicated lineage-

27   specific distributions of photolyase/cryptochrome gene families in each kingdom. Furthermore, the gene

28   repertoires of each subgroup of these families are shown in a combined taxonomy heatmap (Figure 4B),

29   demonstrating lineage-specific evolution and the expansion of subgroups at the species level. These data

30   demonstrate that Gene Search and bioinformatics tools in Genome Analysis in Prometheus support

31   interkingdom comparative analyses. In summary, Prometheus provides the bioinformatics tools essential

32   for comparative analyses, and users can combine these tools with interkingdom comparative analyses in

33   Gene Search to unveil gene function or the evolution of genes/gene families.

34

7

1

**Further genomic analyses using BioExpress with personalized data via My Genes**

Personal data, such as RNA-Seq or Chip-Seq data, and downloaded data from Prometheus, such as genome sequences (genome, CDS, and peptide) in Genome Archive or FASTA files from Gene Search can be uploaded and stored in My Genes (Supplementary Figure 3) and further analyzed using the BioExpress platform (http://bioexpress.kobic.re.kr/bioexpress.en/). BioExpress is a cloud-based analysis platform, and programs for bioinformatics analysis are modularized and shown as icons (Supplementary Figure 5). Users can construct their own analysis pipelines by selecting and linking each modularized program using arrows.

We performed a transcriptomic analysis in BioExpress using the genome of *Hibiscus syriacus*[6] and RNA-Seq data. For this, TopHat[30] and Cufflink[31] programs were used, and genes differentially expressed in tissues from a previous study were identified and visualized as a heatmap (Supplementary Figure 6). Thus, users can perform bioinformatics analyses with personal data in My Genes by linking to BioExpress. This combination of Prometheus and BioExpress can provide convenient and user-friendly analysis conditions for non-bioinformatician scientists.

**Discussion**

Since NGS technology was developed and applied to biology, vast amounts of genomic data have accumulated. With these data, comparative analyses of species or genes can be performed to unveil gene function or evolution. For instance, the evolution of pungency in peppers was discovered by a comparative analysis with tomato and potato genomes[5]. However, only a small number of biologists can perform these comparative analyses using bioinformatics tools. Indeed, the accessibility of bioinformatics analysis is currently a major hurdle for ongoing biologic research. Thus, we constructed Prometheus, a web-based omics portal for interkingdom comparative genomic analyses. Biologists can identify genes or gene families of interest using the domain architectures in Gene Search. Genes from multigene families containing various domain architectures can be detected, such as for the photolyase/cryptochrome family[15] and the nucleotide-binding leucine-rich repeat gene family[32]. Additional subtype information of identified genes is provided in the headers for their sequences in FASTA files.

The goal of combining kingdom-wide gene identification with subtype information is to provide evolutionary insight by detecting lineage-specific subtypes or subtype distribution patterns, as exemplified by the analysis of gene subtypes involved in the TCA cycle. Moreover, users can perform comparative analyses of single genes as well as sets of genes involved in specific signaling pathways. We found that genes containing specific domains showed high rates of accuracy in domain architecture-based Gene Search in Prometheus. However, the accuracy was reduced for genes without specific IPR terms,

8

1   which is a limitation of domain architecture-based gene search systems using InterPro[16] or the pfam[33]

2   database. Nevertheless, this limitation will be minimized as Prometheus is updated with new releases of

3   these databases.

4       To support comparative analyses, Prometheus incorporates various tools, such as LAST[29], Clustal

5   Omega[20], and Phylogeny viewer, in Genome Analysis. This is a valuable addition, as there are currently

6   few web sites for comparative analyses with large restrictive or functionally important gene families, such

7   as TFs. For TFs in plants, there are two representative web sites, PlnTFDB[26] and PlantTFDB[18], but their

8   gene repertoires differ due to their rules for indemnification of TFs[25]. Prometheus clears this particular

9   hurdle via its domain architecture-based Gene Search system, thereby providing biologists with a

10  powerful comparative analysis platform with various tools for further studies.

11      Prometheus provides information to users on individual genomes assigned by taxonomy in

12  Genome Archive via Genome Browser. Here, users can download the genomic and peptide sequences and

13  CDSs as well as upload their own data for further analyses in Prometheus or the cloud-based BioExpress

14  platform. Furthermore, user can access detailed information of genes of interest in the Gene Viewer page.

15  The connection with BioExpress enables Prometheus to provide various bioinformatics tools and allows

16  biologists to analyze their own data in same platform. Thus, unlike other comparative genomics portals or

17  platforms, Prometheus provides tools not only for comparative analyses but also for genomic analyses,

18  such as transcriptome or resequencing analyses. In conclusion, Prometheus is an integrated platform for

19  interkingdom comparative genomic analyses with additional support for other genomic analyses with the

20  user's own data. Thus, Prometheus offers biologists a new paradigm for comparative genome analyses

21  and evolution studies. The platform and InterPro version will be updated annually with newly sequenced

22  genomes to ensure that broad and precise data are available to researchers. Furthermore, newly developed

23  tools for comparative genomic analyses will continue to be added to support various analyses. Finally,

24  visualization of domain subtype architectures identified by Gene Search is now being developed and will

25  be available for updates in the near future.

26

27  **Availability of data and materials**

28  The raw sequence reads of RNA-Seq from *Hibiscus syriacus* has been deposited at DDBJ/ENA/GenBank

29  under accession SRP087036 (PRJNA341314). Detailed methods to perform comparative analysis of

30  Prometheus are provided in Tutorial section (http://prometheus.kobic.re.kr)

31

32  **Acknowledgements**

4

5    **Competing interests**

6    The authors declare that they have no competing financial interests.

7

**References**

1. 1. Maher, B. 2012, ENCODE: The human encyclopaedia. *Nature*, **489**, 46-48.

2. 2. Arabidopsis Interactome Mapping, C. 2011, Evidence for network evolution in an Arabidopsis interactome map. *Science*, **333**, 601-607.

3. 3. Rolland, T., Tasan, M., Charloteaux, B., et al. 2014, A proteome-scale map of the human interactome network. *Cell*, **159**, 1212-1226.

4. 4. Kim, S., Park, M., Yeom, S. I., et al. 2014, Genome sequence of the hot pepper provides insights into the evolution of pungency in Capsicum species. *Nat Genet*, **46**, 270-278.

5. 5. Tomato Genome, C. 2012, The tomato genome sequence provides insights into fleshy fruit evolution. *Nature*, **485**, 635-641.

6. 6. Kim, Y. M., Kim, S., Koo, N., et al. 2017, Genome analysis of Hibiscus syriacus provides insights of polyploidization and indeterminate flowering in woody plants. *DNA Res*, **24**, 71-80.

7. 7. Altenhoff, A. M., Studer, R. A., Robinson-Rechavi, M. and Dessimoz, C. 2012, Resolving the ortholog conjecture: orthologs tend to be weakly, but significantly, more similar in function than paralogs. *PLoS Comput Biol*, **8**, e1002514.

8. 8. Mi, H., Muruganujan, A., Casagrande, J. T. and Thomas, P. D. 2013, Large-scale gene function analysis with the PANTHER classification system. *Nat Protoc*, **8**, 1551-1566.

9. 9. Wang, X., Guo, H., Wang, J., et al. 2016, Comparative genomic de-convolution of the cotton genome revealed a decaploid ancestor and widespread chromosomal fractionation. *The New phytologist*, **209**, 1252-1263.

10. 10. Peng, T., Lin, J., Xu, Y. Z. and Zhang, Y. 2016, Comparative genomics reveals new evolutionary and ecological patterns of selenium utilization in bacteria. *The ISME journal*.

11. 11. Goodstein, D. M., Shu, S., Howson, R., et al. 2012, Phytozome: a comparative platform for green plant genomics. *Nucleic Acids Res*, **40**, D1178-1186.

12. 12. Herrero, J., Muffato, M., Beal, K., et al. 2016, Ensembl comparative genomics resources. *Database : the journal of biological databases and curation*, **2016**.

13. 13. Tyner, C., Barber, G. P., Casper, J., et al. 2017, The UCSC Genome Browser database: 2017 update. *Nucleic Acids Res*, **45**, D626-D634.

14. 14. Theobald, D. L. 2010, A formal test of the theory of universal common ancestry. *Nature*, **465**, 219-222.

15. 15. Kim, Y. M., Choi, J., Lee, H. Y., Lee, G. W., Lee, Y. H. and Choi, D. 2014, dbCRY: a Web-based comparative and evolutionary genomics platform for blue-light receptors. *Database : the journal of biological databases and curation*, **2014**, bau037.

16. 16. Hunter, S., Apweiler, R., Attwood, T. K., et al. 2009, InterPro: the integrative protein signature

11

database. *Nucleic Acids Res*, **37**, D211-215.

17.  Letunic, I. and Bork, P. 2016, Interactive tree of life (iTOL) v3: an online tool for the display and annotation of phylogenetic and other trees. *Nucleic Acids Res*, **44**, W242-245.

18.  Jin, J., Zhang, H., Kong, L., Gao, G. and Luo, J. 2014, PlantTFDB 3.0: a portal for the functional and evolutionary study of plant transcription factors. *Nucleic Acids Res*, **42**, D1182-1187.

19.  Altschul, S. F., Gish, W., Miller, W., Myers, E. W. and Lipman, D. J. 1990, Basic local alignment search tool. *J Mol Biol*, **215**, 403-410.

20.  Sievers, F. and Higgins, D. G. 2014, Clustal Omega, accurate alignment of very large numbers of sequences. *Methods in molecular biology*, **1079**, 105-116.

21.  Edgar, R. C. 2004, MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res*, **32**, 1792-1797.

22.  Li, L., Stoeckert, C. J., Jr. and Roos, D. S. 2003, OrthoMCL: identification of ortholog groups for eukaryotic genomes. *Genome research*, **13**, 2178-2189.

23.  Blum, T., Briesemeister, S. and Kohlbacher, O. 2009, MultiLoc2: integrating phylogeny and Gene Ontology terms improves subcellular protein localization prediction. *BMC Bioinformatics*, **10**, 274.

24.  Emanuelsson, O., Brunak, S., von Heijne, G. and Nielsen, H. 2007, Locating proteins in the cell using TargetP, SignalP and related tools. *Nat Protoc*, **2**, 953-971.

25.  Zheng, Y., Jiao, C., Sun, H., et al. 2016, iTAK: A Program for Genome-wide Prediction and Classification of Plant Transcription Factors, Transcriptional Regulators, and Protein Kinases. *Mol Plant*, **9**, 1667-1670.

26.  Perez-Rodriguez, P., Riano-Pachon, D. M., Correa, L. G., Rensing, S. A., Kersten, B. and Mueller-Roeber, B. 2010, PlnTFDB: updated content and new features of the plant transcription factor database. *Nucleic Acids Res*, **38**, D822-827.

27.  Dai, X., Sinharoy, S., Udvardi, M. and Zhao, P. X. 2013, PlantTFcat: an online plant transcription factor and transcriptional regulator categorization and analysis tool. *BMC Bioinformatics*, **14**, 321.

28.  Yilmaz, A., Mejia-Guerra, M. K., Kurz, K., Liang, X., Welch, L. and Grotewold, E. 2011, AGRIS: the Arabidopsis Gene Regulatory Information Server, an update. *Nucleic Acids Res*, **39**, D1118-1122.

29.  Kielbasa, S. M., Wan, R., Sato, K., Horton, P. and Frith, M. C. 2011, Adaptive seeds tame genomic sequence comparison. *Genome research*, **21**, 487-493.

30.  Trapnell, C., Pachter, L. and Salzberg, S. L. 2009, TopHat: discovering splice junctions with RNA-Seq. *Bioinformatics*, **25**, 1105-1111.

31.  Ghosh, S. and Chan, C. K. 2016, Analysis of RNA-Seq Data Using TopHat and Cufflinks.

1    *Methods in molecular biology*, **1374**, 339-361.

2    32.    Seo, E., Kim, S., Yeom, S.-I. and Choi, D. 2016, Genome-Wide Comparative Analyses Reveal

3           the Dynamic Evolution of Nucleotide-Binding Leucine-Rich Repeat Gene Family among

4           Solanaceae Plants. *Frontiers in plant science*, **7**.

5    33.    Mistry, J. and Finn, R. 2007, Pfam: a domain-centric method for analyzing proteins and

6           proteomes. *Methods in molecular biology*, **396**, 43-58.

7

8

13

1    **Figure 1**. Concept and construction progresses of Prometheus.

2    Schematic showing the workflow for the construction of Prometheus (left), detailed information for each

3    stage (middle), and the functions available with Prometheus (right).

4

5    **Figure 2**. Construction of primary and secondary databases.

6    (A) Screenshot of the Genome Archive page. (B) The numbers of species and genome versions used for

7    the construction of Prometheus. (C) Screenshot of the Genome Browser of Prometheus. A region of the

8    human genome (HGP 38) is shown enlarged in the small box. (D) Screenshot of the Gene Viewer for

9    providing detailed information of individual genes. Gene structure, domain architecture, subcellular

10   localization, and orthologous and paralogous genes are shown in each panel.

11

12   **Figure 3**. Identification of TFs and genes in the TCA using Gene Search.

13   (A) Validation of identified TFs using iTAK pipeline. (B) Human TCA cycle genes were investigated and

14   used for further analysis. The ratios of each gene are shown as heatmaps.

15

16   **Figure 4**. Interkingdom comparative analysis of the photolyase/cryptochrome gene family.

17   (A) Domain architectures of the photolyase/cryptochrome gene family. (B) Distribution of

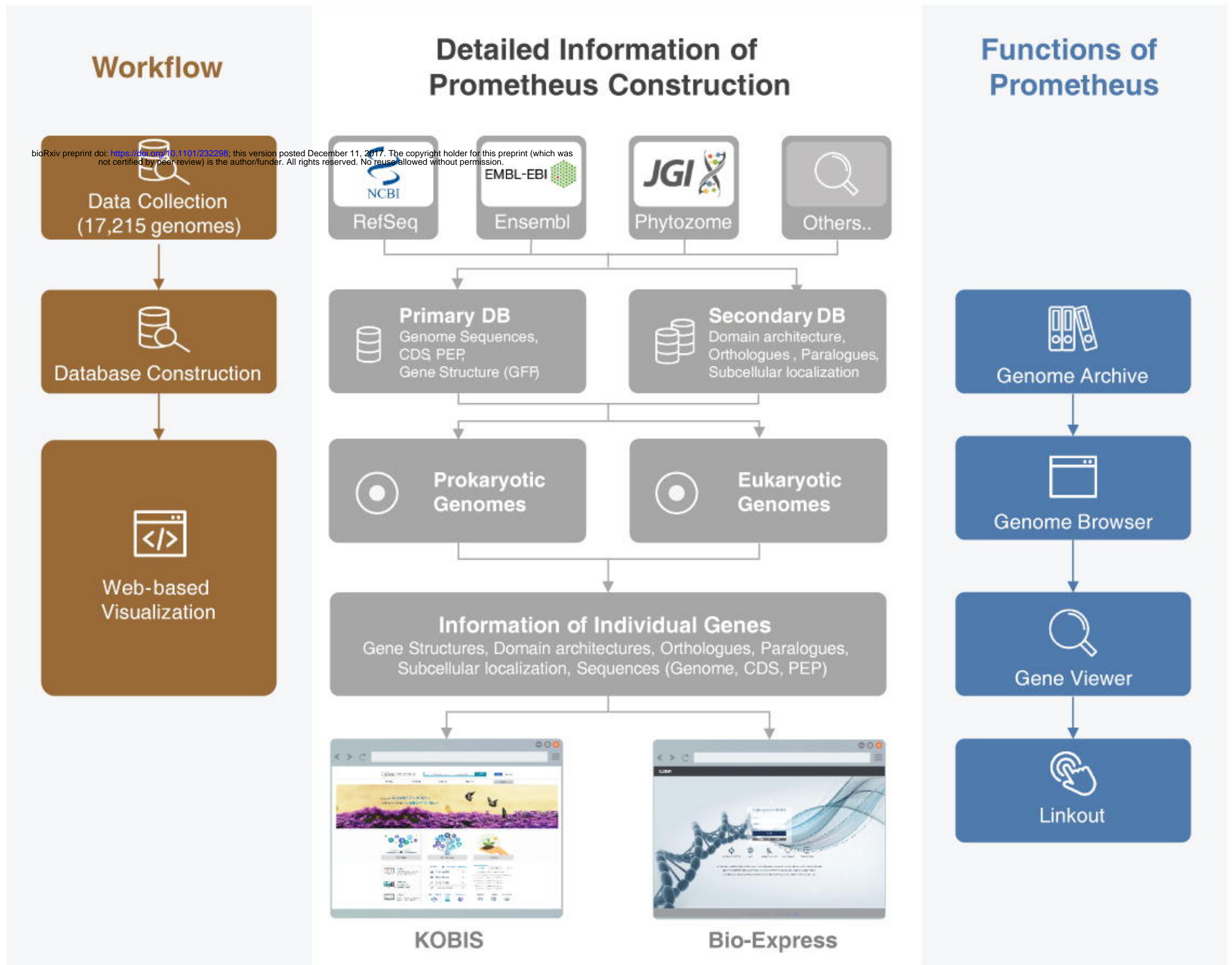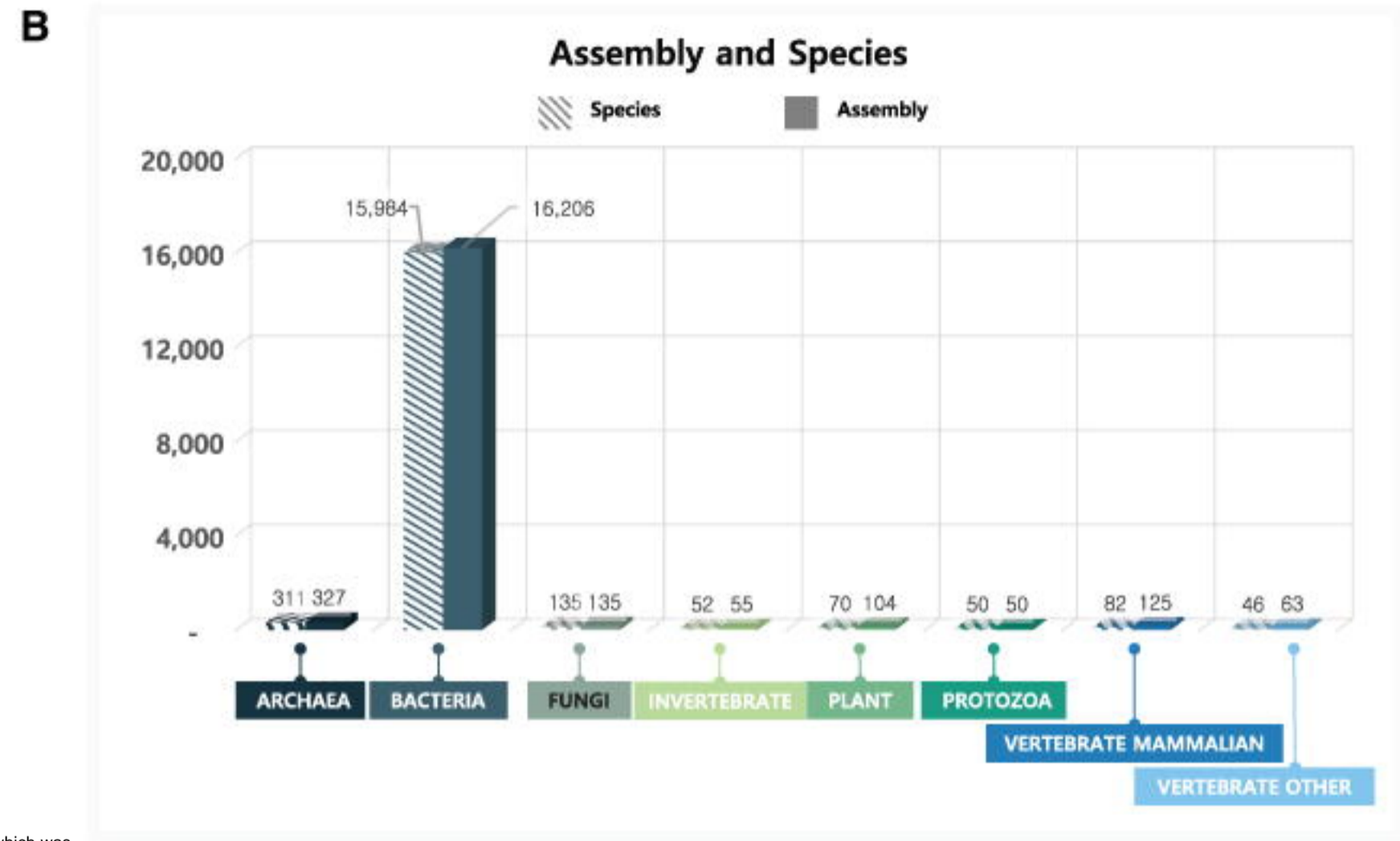18   photolyase/cryptochrome genes and taxonomic tree.

# Figure1

**Workflow**

Data Collection
(17,215 genomes)

Database Construction

Web-based
Visualization

**Detailed Information of
Prometheus Construction**

NCBI
RefSeq

EMBL-EBI
Ensembl

JGI
Phytozome

Others..

**Primary DB**
Genome Sequences,
CDS, PEP,
Gene Structure (GFF)

**Secondary DB**
Domain architecture,
Orthologues, Paralogues,
Subcellular localization

**Prokaryotic
Genomes**

**Eukaryotic
Genomes**

**Information of Individual Genes**
Gene Structures, Domain architectures, Orthologues, Paralogues,
Subcellular localization, Sequences (Genome, CDS, PEP)

KOBIS

Bio-Express

**Functions of
Prometheus**

Genome Archive

Genome Browser

Gene Viewer

Linkout

# Figure2

# Figure3



**A**

**Family site**

FAR1 (98.8%) | MADS (99.5%) | Homeobox (99.9%) | GARP-ARR-B (94.2%) | AUX-IAA (98.7%) | CAMTA (91.04%) | LIM (86.0%)

4700 / 58 | 4425 / 24 | 6405 / 1 | 692 / 43 | 2190 / 28 | 376 / 37 | 314 / 51

**Domain only**

NAC (96.7%) | B3-Type TF (94.5%) | C2C2-Dof (99.9%) | C2C2-CO (99.9%) | C2C2-GATA (93.8%) | CCAAT (93.4%) | BES1 (99.9%)

12007 / 401 | 8590 / 488 | 3516 / 2 | 1448 / 2 | 608 / 40 | 5733 / 404 | 1027 / 1

iTAK   Prometheus

**B**

Citric acid cycle

# Figure4