

**A peer-reviewed version of this preprint was published in PeerJ on 20 January 2020.**

[View the peer-reviewed version](https://peerj.com/articles/cs-251) (peerj.com/articles/cs-251), which is the preferred citable publication unless you specifically need to cite this preprint.

Hao Z, Lv D, Ge Y, Shi J, Weijers D, Yu G, Chen J. 2020. *R*Idiogram: drawing SVG graphics to visualize and map genome-wide data on the idiograms. PeerJ Computer Science 6:e251  
<https://doi.org/10.7717/peerj-cs.251>

# ***R*ldeogram: drawing SVG graphics to visualize and map genome-wide data on the idiograms**

Zhaodong Hao<sup>1,2</sup>, Dekang Lv<sup>3</sup>, Ying Ge<sup>3</sup>, Jisen Shi<sup>1</sup>, Dolf Weijers<sup>2</sup>, Guangchuang Yu<sup>Corresp., 4</sup>, Jinhui Chen<sup>Corresp. 1</sup>

<sup>1</sup> Key Laboratory of Forest Genetics & Biotechnology of Ministry of Education, Co-Innovation Center for Sustainable Forestry in Southern China, Nanjing Forestry University, Nanjing, China

<sup>2</sup> Laboratory of Biochemistry, Wageningen University and Research, Wageningen, Netherlands

<sup>3</sup> Institute of Cancer Stem Cell, Dalian Medical University, Dalian, China

<sup>4</sup> Institute of Bioinformatics, School of Basic Medical Sciences, Southern Medical University, Guangzhou, China

Corresponding Authors: Guangchuang Yu, Jinhui Chen  
Email address: gcyu1@smu.edu.cn, chenjh@njfu.edu.cn

**Background:** Owing to the rapid advances in DNA sequencing technologies, whole genome from more and more species are becoming available at increasing pace. For whole-genome analysis, idiograms provide a very popular, intuitive and effective way to map and visualize the genome-wide information, such as GC content, gene and repeat density, DNA methylation distribution, etc. However, most available software programs and web servers are available only for a few model species, such as human, mouse and fly. As boundaries between model and non-model species are shifting, tools are urgently needs to generate idiograms for a broad range of species are needed to help better understanding fundamental genome characteristics.

**Results:** The R package *Rldeogram* allows users to build high-quality idiograms of any species of interest. It can map continuous and discrete genome-wide data on the idiograms and visualize them in a heat map and track labels, respectively.

**Conclusion:** The visualization of genome-wide data mapping and comparison allow users to quickly establish a clear impression of the chromosomal distribution pattern, thus making *Rldeogram* a useful tool for any researchers working with omics.

# 1 *RIdeogram*: drawing SVG graphics to visualize and 2 map genome-wide data on the idiograms

3  
4

5 Zhaodong Hao<sup>1,2</sup>, Dekang Lv<sup>3</sup>, Ying Ge<sup>3</sup>, Jisen Shi<sup>1</sup>, Dolf Weijers<sup>2</sup>, Guangchuang Yu<sup>4</sup> and Jinhui  
6 Chen<sup>1</sup>

7

8 <sup>1</sup>Key Laboratory of Forest Genetics & Biotechnology of Ministry of Education, Co-Innovation  
9 Center for Sustainable Forestry in Southern China, Nanjing Forestry University, Nanjing  
10 Jiangsu, China

11 <sup>2</sup>Laboratory of Biochemistry, Wageningen University, Wageningen, Haarlem, The Netherlands

12 <sup>3</sup>Institute of Cancer Stem Cell, Dalian Medical University, Dalian, Liaoning, China

13 <sup>4</sup>Institute of Bioinformatics, School of Basic Medical Sciences, Southern Medical University,  
14 Guangzhou, Guangdong, China

15

16 Corresponding Author:

17 Guangchuang Yu<sup>4</sup>

18 1023 ShaTaiNan Road, Guangzhou, Guangdong, 510515, China

19 Email address: gcyu1@smu.edu.cn

20 Jinhui Chen<sup>1</sup>

21 159 LongPan Road, Nanjing, Jiangsu, 210037, China

22 Email address: chenjh@njfu.edu.cn

23

## 24 **Abstract**

25 **Background:** Owing to the rapid advances in DNA sequencing technologies, whole genome  
26 from more and more species are becoming available at increasing pace. For whole-genome  
27 analysis, idiograms provide a very popular, intuitive and effective way to map and visualize the  
28 genome-wide information, such as GC content, gene and repeat density, DNA methylation  
29 distribution, etc. However, most available software programs and web servers are available only  
30 for a few model species, such as human, mouse and fly. As boundaries between model and non-  
31 model species are shifting, tools are urgently needs to generate idiograms for a broad range of  
32 species are needed to help better understanding fundamental genome characteristics.

33 **Results:** The R package *RIdeogram* allows users to build high-quality idiograms of any species  
34 of interest. It can map continuous and discrete genome-wide data on the idiograms and visualize  
35 them in a heat map and track labels, respectively.

36 **Conclusion:** The visualization of genome-wide data mapping and comparison allow users to  
37 quickly establish a clear impression of the chromosomal distribution pattern, thus making  
38 *RIdeogram* a useful tool for any researchers working with omics.

39

## 40 Introduction

41 Recently, with the development of sequencing technologies, especially rapid advances in third  
42 generation sequencing (Pacific Biosciences and Oxford Nanopore Technologies), BioNano  
43 genome mapping and High-throughput chromatin conformation capture sequencing, many  
44 sequenced species have their genomes updated to chromosome level and more and more non-  
45 model species have their genomes sequenced (Jiao & Schneeberger 2017; Phillippy 2017). After  
46 the chromosome-level genome completion, an overview of some genome characteristics can help  
47 to better understand a species genome, such as gene and transposon distribution across the  
48 sunflower genome (Badouin et al. 2017).

49 A idiogram, also known as a karyotype, is defined as the phenotypic appearance of  
50 chromosomes in the nucleus of a eukaryotic cell and has been widely used to visualize the  
51 genome-wide data since the first web server, *Idiographica*, came online in 2007 (Kin & Ono  
52 2007). However, this web server (updated November, 2017) still only caters to four species:  
53 human, mouse, rat and fly. Recently, an R package called *chromoMap* was published in CRAN  
54 and this allows users to interactively visualize and map chromosome elements on the  
55 chromosome plot. In addition, there are two JavaScript libraries for chromosome visualization,  
56 one is *Ideogram.js* and the other is *karyotypeSVG*. We still lack the choice of available drawing  
57 tools for plotting idiograms conveniently and effectively with a wide range species, despite the  
58 widespread need to visualize features along entire chromosomes.

59 Scalable Vector Graphics (SVG) is a language for describing two-dimensional graphics  
60 applications and images. SVG graphics is defined in an eXtensible Markup Language (XML)  
61 text file which means that one can easily use any text editor or drawing software to create and  
62 edit SVG graphics. Most R graphics packages are built on two graphics systems, the traditional  
63 graphics system and the grid graphics system. Here, we developed an R package (*RIdeogram*) to  
64 draw high-quality idiograms without species limitations, that allows to visualize and map whole-  
65 genome information on the idiograms based on the SVG language.

66

## 67 Description

68 The package *RIdeogram* is written in R (R Core Team, 2018), one of the most popular  
69 programming languages widely used in statistical computing, data analytics and graphics.  
70 However, this new R graphics package is not built based on any existing graphics systems. We  
71 use the R environment to read the custom input files and calculate the drawing element positions  
72 in a coordinate system. Then, we use R to write all element information into a text file following  
73 the XML format which are used to define graphics by the SVG language. A list of the currently  
74 implemented commands is given in Table 1. In general, there are three main functions, *data*,  
75 *ideogram* and *convertSVG* implemented in the package *RIdeogram*. Users can use the function  
76 *data* to load the example data or the basic R function *read.table* to load the custom data from  
77 local files. Then, the function *ideogram* can be used to compute the information for all drawing  
78 elements based on the input files and generate a A4-sized SVG file containing a vector graphic  
79 which can be conveniently viewed and modified using the software Adobe Illustrator or

80 Inkscape. Alternatively, users can also use the function *convertSVG* to convert this SVG file into  
81 an adjustable image format (pdf, png, tiff, or jpg) with a user-defined resolution according to the  
82 practical requirements.

83 In general, there are two types of data, i.e., continuous and discrete data. For mapping and  
84 visualizing, *RIdeogram* considers the continuous data, such as gene density across the whole  
85 genome in 1-Mb windows, as overlaid features and maps them on the idiograms with dark/light  
86 colors representing high/low values. For the other data type that are scattered throughout the  
87 whole genome, such as the chromosomal distribution of members in one gene family,  
88 *RIdeogram* can add track labels next to the idiograms with three shapes (box, circle and triangle)  
89 available to represent different characteristics of these members, such as the subclade that one  
90 gene member belongs to. Users can also combine the shapes and colors to represent more than  
91 three distinct characteristic types.

92 *RIdeogram* is available through CRAN (<https://cran.r-project.org/web/packages/RIdeogram/>)  
93 and is developed on GitHub (<https://github.com/TickingClock1992/RIdeogram>). Further  
94 extensions in development and fixes can be seen in the issue listing page on the package's  
95 GitHub page. The new function that we are planning to implement in next version include, but  
96 are not limited to, linking genome regions on two adjacent idiograms with Bezier curves or  
97 straight lines for synteny visualization and enlarging the user-specified genome regions to  
98 display detailed characteristics, as we gather more from users.

99

## 100 Examples

101 Our first example use the data contained in this package. After the completion of genome  
102 sequencing, assembly and annotation, *RIdeogram* can be used to give some idea of how genes  
103 are distributed across the whole genome. The example data contained numbers of protein-coding  
104 genes calculated in 1-Mb windows which can be considered as continues data and positions of  
105 500 random selected non-coding RNAs, including ribosomal RNAs (rRNAs), transfer RNAs  
106 (tRNAs) and microRNAs (miRNAs), which can be considered as discrete data. *RIdeogram* maps  
107 the gene density information on the idiograms as overlaid features in a heat map and adds track  
108 labels next to the idiograms with green boxes, purple circles and orange triangles representing  
109 rRNAs, tRNAs and miRNAs, respectively (Fig. 1). Obviously, inter- and intra-chromosomal  
110 gene distributions are non-uniform. For instance, the chromosomal regions adjacent to the  
111 centromeres are gene-poor in chromosome 1, 9 and 16 while those are gene-rich in chromosome  
112 11, 14 and 17. This function can be applied to many different situations, such as single  
113 nucleotide polymorphism (SNP) density and candidate markers (Fig. S1 & Data S1, original data  
114 see Li et al. 2019), DNA methylation dynamics and potential activated genes (Fig. S2 & Data S2,  
115 original data see Huang et al. 2019) and transcription factor (TF) binding sites and candidate  
116 target genes (Fig. S3 & Data S3, original data see Shamimuzzaman & Vodkin 2013).

117 Besides visualizing some specific genome characteristics across the whole genome at the  
118 chromosome level as showed in Fig. 1, *RIdeogram* can also be used to compare two relevant  
119 genome features, such as gene and repeat density, which will provide some important

120 implications for better understanding the relevance of chromosomal distribution patterns of these  
121 two features. The example data implemented in this package also contained the information of  
122 long terminal repeat (LTR) distribution across the human genome. Since the transposable  
123 elements have been suggested to have a potential detrimental effect on gene expression (Hollister  
124 & Gaut 2009), the distributions of gene and LTR are supposed to be opposite across the whole  
125 genome as a result of natural selection. As expect, the region that has a relatively high gene  
126 content usually has a relatively low LTR density and vice versa (Fig. 2), indicating that LTR  
127 seems to avoid inserting in the regions with a high gene content in the genome. This similar  
128 phenomenon was also observed in the sunflower genome explained using two idiogram graphics,  
129 one showing the gene distribution and the other showing the LTR distribution (Badouin et al.  
130 2017). Using *RIdeogram*, users can integrate these two graphics into one, much easier for  
131 researchers to interpret and readers to understand. Apart from the differences, this function can  
132 also be used to show the similarities, like the similar genetic diversity patterns across the whole  
133 genome between two geographical groups of the same species (Fig. S4 & Data S4, original data  
134 see Chen et al. 2019).

135

## 136 Conclusion

137 The *RIdeogram* package provides an efficient and effective way to build idiograms with no  
138 species limitations and map genome-wide information on the idiograms for better visualizing and  
139 understanding the chromosomal distribution patterns of some particular genomic features.  
140 Meanwhile, this package is user-friendly and accessible for biologists without extensive  
141 computer programming expertise. In addition, *RIdeogram* can generate two types of images, a  
142 vector graphic or a bitmap file, both in high-quality and meeting conventional requirements for  
143 direct use in presentations or journal publications.

144

## 145 Acknowledgements

146 We thank Dr. Zhongjuan Zhang for her comments on the manuscript.

147

## 148 References

- 149 Badouin H, Gouzy J, Grassa CJ, Murat F, Staton SE, Cottret L, Lelandais-Briere C, Owens GL,  
150 Carrere S, Mayjonade B, Legrand L, Gill N, Kane NC, Bowers JE, Hubner S, Bellec A,  
151 Berard A, Berges H, Blanchet N, Boniface MC, Brunel D, Catrice O, Chaidir N, Claudel  
152 C, Donnadiou C, Faraut T, Fievet G, Helmstetter N, King M, Knapp SJ, Lai Z, Le Paslier  
153 MC, Lippi Y, Lorenzon L, Mandel JR, Marage G, Marchand G, Marquand E, Bret-  
154 Mestries E, Morien E, Nambeesan S, Nguyen T, Pegot-Espagnet P, Pouilly N, Raftis F,  
155 Sallet E, Schiex T, Thomas J, Vandecasteele C, Vares D, Vear F, Vautrin S, Crespi M,  
156 Mangin B, Burke JM, Salse J, Munos S, Vincourt P, Riaseberg LH, and Langlade NB.  
157 2017. The sunflower genome provides insights into oil metabolism, flowering and Asterid  
158 evolution. *Nature* 546:148-152. DOI: 10.1038/nature22380
- 159 Chen J, Hao Z, Guang X, Zhao C, Wang P, Xue L, Zhu Q, Yang L, Sheng Y, Zhou Y, Xu H, Xie  
160 H, Long X, Zhang J, Wang Z, Shi M, Lu Y, Liu S, Guan L, Zhu Q, Yang L, Ge S, Cheng  
161 T, Laux T, Gao Q, Peng Y, Liu N, Yang S, and Shi J. 2019. *Liriodendron* genome sheds

- 162 light on angiosperm phylogeny and species-pair differentiation. *Nature Plants* 5:18-25.  
163 DOI: 10.1038/s41477-018-0323-6
- 164 Hollister JD, and Gaut BS. 2009. Epigenetic silencing of transposable elements: a trade-off  
165 between reduced transposition and deleterious effects on neighboring gene expression.  
166 *Genome Research* 19:1419-1428. DOI: 10.1101/gr.091678.109
- 167 Huang H, Liu R, Niu Q, Tang K, Zhang B, Zhang H, Chen K, Zhu JK, and Lang Z. 2019. Global  
168 increase in DNA methylation during orange fruit development and ripening. *Proceedings*  
169 *of the National Academy of Sciences of the United States of America* 116:1430-1436.  
170 DOI: 10.1073/pnas.1815441116
- 171 Jiao WB, and Schneeberger K. 2017. The impact of third generation genomic technologies on  
172 plant genome assembly. *Current Opinion in Plant Biology* 36:64-70. DOI:  
173 10.1016/j.pbi.2017.02.002
- 174 Kin T, and Ono Y. 2007. Idiographica: a general-purpose web application to build idiograms on-  
175 demand for human, mouse and rat. *Bioinformatics* 23:2945-2946. DOI:  
176 10.1093/bioinformatics/btm455
- 177 Li X, Singh J, Qin M, Li S, Zhang X, Zhang M, Khan A, Zhang S, and Wu J. 2019. Development  
178 of an integrated 200K SNP genotyping array and application for genetic mapping,  
179 genome assembly improvement and genome wide association studies in pear (*Pyrus*).  
180 *Plant Biotechnology Journal* 17:1582-1594. DOI: 10.1111/pbi.13085
- 181 Phillippy AM. 2017. New advances in sequence assembly. *Genome Research* 27:xi-xiii. DOI:  
182 10.1101/gr.223057.117
- 183 Shamimuzzaman M, and Vodkin L. 2013. Genome-wide identification of binding sites for NAC  
184 and YABBY transcription factors and co-regulated genes during soybean seedling  
185 development by ChIP-Seq and RNA-Seq. *BMC Genomics* 14:477. DOI: 10.1186/1471-  
186 2164-14-477
- 187

**Table 1** (on next page)

Table 1: Functions and data contained in the package *RIdeogram*.



1 **Table 1. Functions and data contained in the package *RIdeogram*.**

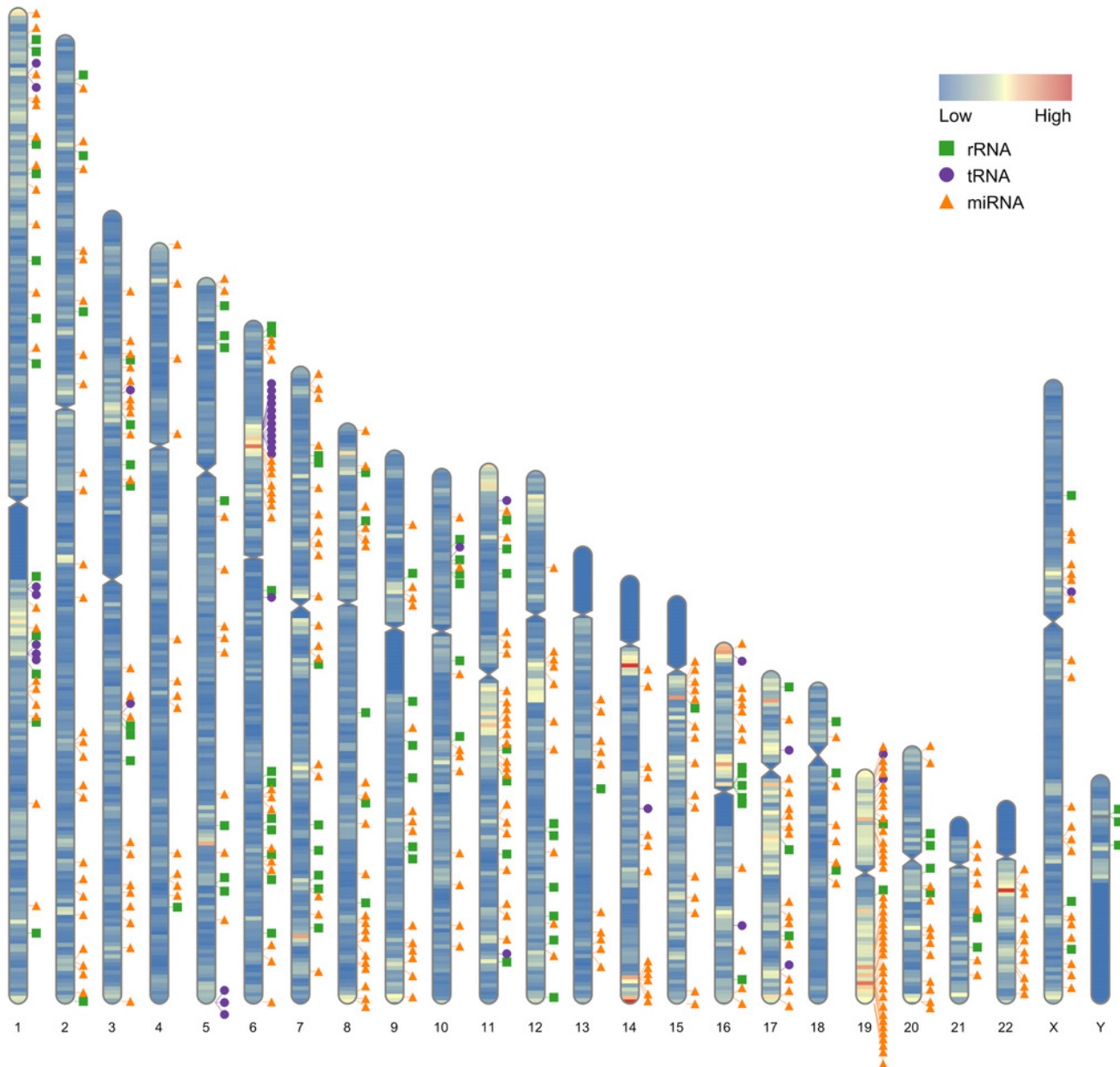
Function/Data name	Description
data(human_karyotype)	Load a data frame of the human karyotype information including the length information of each chromosome and the position information of each centromere.
data(gene_density)	Load a data frame of the gene number calculated in 1-Mb windows across the whole human genome.
data(LTR_density)	Load a data frame of the LTR number calculated in 1-Mb windows across the whole human genome.
data(Random_RNAs_500)	Load a data frame of the positions of 500 random selected RNAs throughout the whole human genome.
ideogram	Map and visualize the genome-wide data on the ideograms
convertSVG	Convert the output file from the SVG format to the format users chose.
svg2tiff	Convert the output file from the SVG format to the TIFF format.
svg2pdf	Convert the output file from the SVG format to the PDF format.
svg2jpg	Convert the output file from the SVG format to the JPG format.
svg2png	Convert the output file from the SVG format to the PNG format.

2

# Figure 1

Figure 1: Gene distribution across the whole human genome.

The overlaid heatmap shows the gene density and the tack labels refer to 500 random selected RNAs consisted of rRNAs (green boxes), tRNA (purple circles) and miRNA (orange triangles) locus across the human genome. Annotation information was downloaded from the GENCODE website (<https://www.gencodegenes.org>).



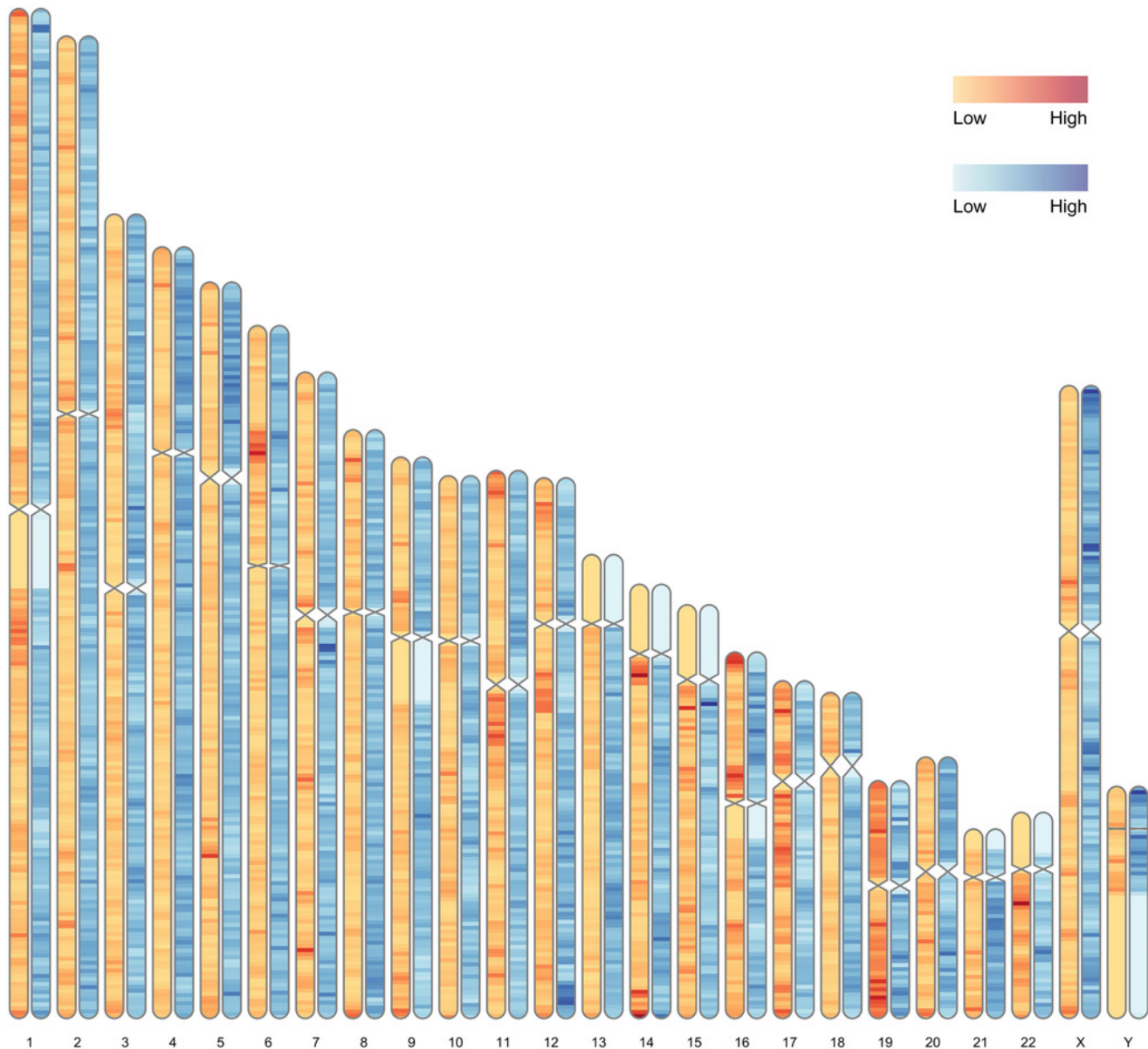
**Figure 1. Gene distribution across the whole human genome.**

The overlaid heatmap shows the gene density and the tack labels refer to 500 random selected RNAs consisted of rRNAs (green boxes), tRNA (purple circles) and miRNA (orange triangles) locus across the human genome. Annotation information was downloaded from the GENCODE website (<https://www.genencodegenes.org>).

## Figure 2

Figure 2: A comparison of chromosomal distribution of genes and LTRs in the human genome.

The gene number and LTR number are both counted in a 1-Mb window. Red color represents the gene number (range 0–135 per Mb) and blue color represents the LTR number (range 0–606 per Mb). The light and dark colors represent a low and high content, respectively. This plot shows that gene and LTR have an opposite distribution pattern along the human chromosomes.



**Figure 2. A comparison of chromosomal distribution of genes and LTRs in the human genome.**

The gene number and LTR number are both counted in a 1-Mb window. Red color represents the gene number (range 0–135 per Mb) and blue color represents the LTR number (range 0–606 per Mb). The light and dark colors represent a low and high content, respectively. This plot shows that gene and LTR have an opposite distribution pattern along the human chromosomes.