

生物信息学

天津医科大学
生物医学工程与技术学院

2015-2016 学年上学期 (秋)
2014 级生信班

第五章 基因组功能注释分析

伊现富 (Yi Xianfu)

天津医科大学 (TJMU)
生物医学工程与技术学院

2015 年 10 月



5.1,5.2 基因组功能注释分析基础

- ① 基础知识：组装版本，坐标系统，常用格式，逻辑运算模式
- ② 准备工作：坐标转换，格式转换，逻辑运算
- ③ 扩展知识：文本文件与文本编辑器

5.3 基因组功能的高级注释

- ① 高级注释：变异位点注释，富集分析，序列标识
- ② 扩展知识：box plot，解析图表

5.3 Galaxy 分析平台

- ① Galaxy 分析平台：简介，使用
- ② 扩展知识：数据处理三段论

教学提纲

- 1 引言
- 2 基因组组装版本
- 3 基因组坐标系统
- 4 基因组注释常用格式
- 5 文本文件与文本编辑器
- 6 基因组坐标的逻辑运算
- 7 总结与答疑
- 8 引言
- 9 变异位点的注释
- 10 基因集富集分析

- 11 序列标识
- 12 box plot
- 13 解析图表
- 14 总结与答疑
- 15 引言
- 16 Galaxy 分析平台
- 17 Galaxy 使用演示
- 18 数据处理三段论
- 19 总结与答疑
- 20 复习思考题

教学提纲

1

引言

2

基因组组装版本

3

基因组坐标系统

4

基因组注释常用格式

5

文本文件与文本编辑器

6

基因组坐标的逻辑运算

7

总结与答疑

8

引言

9

变异位点的注释

10

基因集富集分析

11

序列标识

12

box plot

13

解析图表

14

总结与答疑

15

引言

16

Galaxy 分析平台

17

Galaxy 使用演示

18

数据处理三段论

19

总结与答疑

20

复习思考题



基因组注释 (genome annotation)

从原始的基因组核酸序列中挖掘有用的生物学信息并阐释其生物学含义，包括基因组结构注释和基因组功能注释两大部分。

基因组结构注释 (structural annotation)

在基因组序列中寻找基因等功能元件并明确其基本结构。

基因组功能注释 (functional annotation)

在结构注释的基础上，将进化保守性 (evolutionary conservation) 和基因本体论 (gene ontology) 等元数据 (meta-data) 与功能元件对应起来，找到其生物学功能。



基因组注释 (genome annotation)

从原始的基因组核酸序列中挖掘有用的生物学信息并阐释其生物学含义，包括基因组结构注释和基因组功能注释两大部分。

基因组结构注释 (structural annotation)

在基因组序列中寻找基因等功能元件并明确其基本结构。

基因组功能注释 (functional annotation)

在结构注释的基础上，将进化保守性 (evolutionary conservation) 和基因本体论 (gene ontology) 等元数据 (meta-data) 与功能元件对应起来，找到其生物学功能。



基因组注释 (genome annotation)

从原始的基因组核酸序列中挖掘有用的生物学信息并阐释其生物学含义，包括基因组结构注释和基因组功能注释两大部分。

基因组结构注释 (structural annotation)

在基因组序列中寻找基因等功能元件并明确其基本结构。

基因组功能注释 (functional annotation)

在结构注释的基础上，将进化保守性 (evolutionary conservation) 和基因本体论 (gene ontology) 等元数据 (meta-data) 与功能元件对应起来，找到其生物学功能。



基因组注释

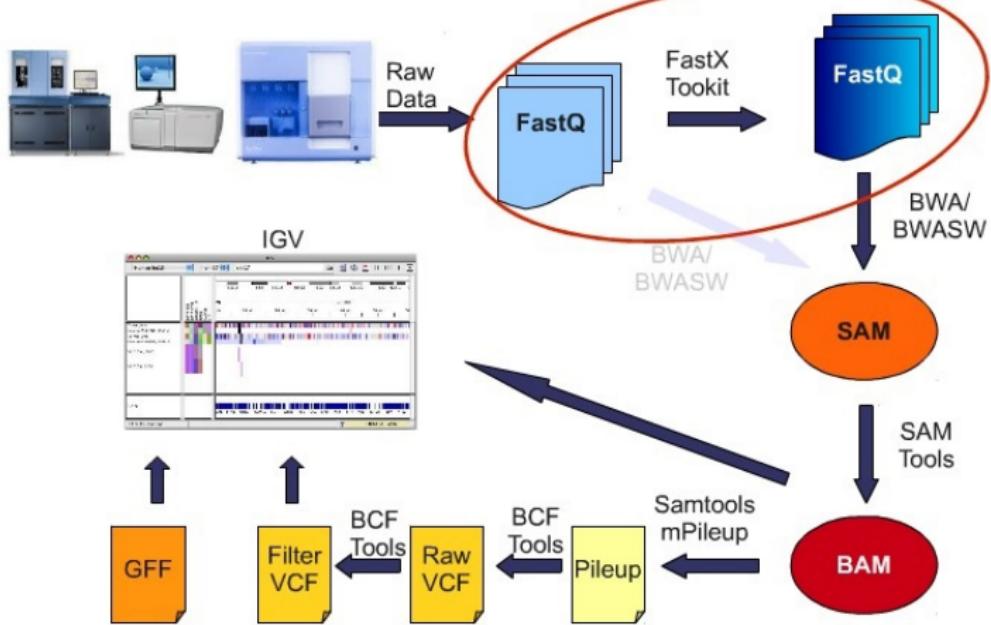
- 结构注释 ← 实验手段，单个基因
 - 限制性酶切位点分析、开放阅读框分析、启动子分析、CpG 岛识别
 - 重复序列分析、基因识别
 - mRNA 选择性剪接分析
- 功能注释 ← 组学时代，复杂疾病
 - 变异位点的注释
 - 基因集富集分析
 - 生物学通路分析
 - 相互作用网络分析
 - 分子进化分析



- 基因组组装版本
- 基因组坐标系统
- 注释常用格式
- 文本编辑器
- 坐标的逻辑运算



Sequence to Variation Workflow



教学提纲

1

引言

2

基因组组装版本

3

基因组坐标系统

4

基因组注释常用格式

5

文本文件与文本编辑器

6

基因组坐标的逻辑运算

7

总结与答疑

8

引言

9

变异位点的注释

10

基因集富集分析

11

序列标识

12

box plot

13

解析图表

14

总结与答疑

15

引言

16

Galaxy 分析平台

17

Galaxy 使用演示

18

数据处理三段论

19

总结与答疑

20

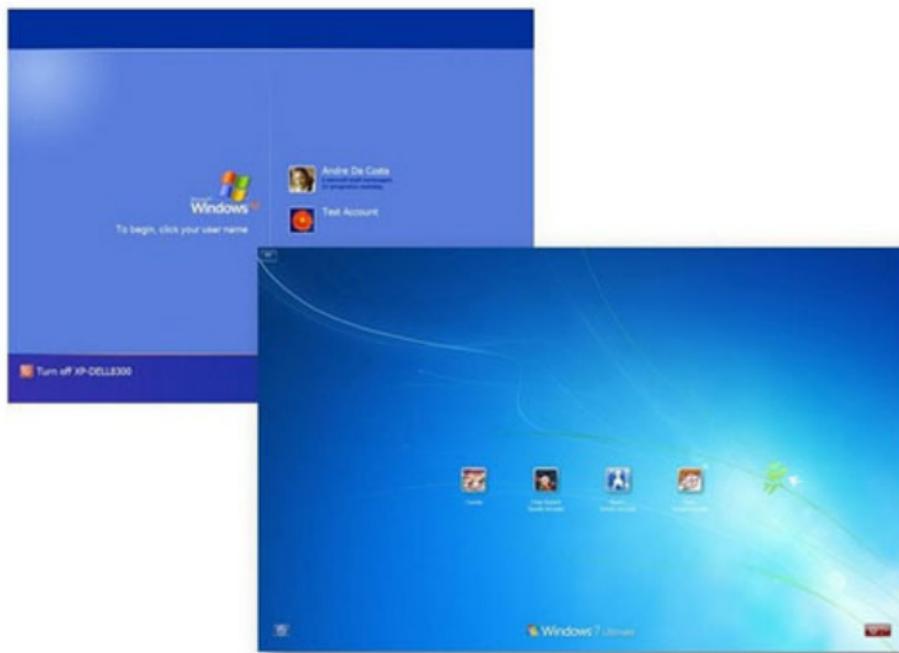
复习思考题



- These sequences were mapped to human and mouse genomes sequences ([hg18 and mm9](#), respectively) using BLASTN.
- We used DNA sequences from the human and mouse genome assemblies [hg18 and mm9](#).
- Currently there are 25,000 genes annotated in the human ([hg18](#)) and mouse ([mm9](#)) genome, which comprise less than 3% of the genome (UCSC genome browser; <http://genome.ucsc.edu/>).
- The [GRCh37/hg19](#) and [GRCm38/mm10](#) assemblies at the UCSC genome browser (<http://genome.ucsc.edu/>) were used for mapping the chromosomal defect and gene annotations.
- The genome assemblies from which the sequences obtained were Dec 2011 ([GRCm38/mm10](#)), Feb 2009 ([GRCh37/hg19](#)) and Nov 2004 ([Baylor3.4/rn4](#)) for mouse, human and rat respectively.



组装版本 | XP vs. Win7



基因组序列不是确定的吗？也需要版本升级？

SPECIES	UCSC	DATE	NCBI
Human	hg19	Feb. 2009	Genome Reference Consortium GRCh37
	hg18	Mar. 2006	NCBI Build 36.1
	hg17	May 2004	NCBI Build 35
	hg16	Jul. 2003	NCBI Build 34
Mouse	mm10	Dec. 2011	Genome Reference Consortium GRCm38
	mm9	Jul. 2007	NCBI Build 37
	mm8	Feb. 2006	NCBI Build 36
	mm7	Aug. 2005	NCBI Build 35

human: *Homo sapiens*; mouse: *Mus musculus*

hg: human genome; GRC: Genome Reference Consortium



基因组序列不是确定的吗？也需要版本升级？

SPECIES	UCSC	DATE	NCBI
Human	hg19	Feb. 2009	Genome Reference Consortium GRCh37
	hg18	Mar. 2006	NCBI Build 36.1
	hg17	May 2004	NCBI Build 35
	hg16	Jul. 2003	NCBI Build 34
Mouse	mm10	Dec. 2011	Genome Reference Consortium GRCm38
	mm9	Jul. 2007	NCBI Build 37
	mm8	Feb. 2006	NCBI Build 36
	mm7	Aug. 2005	NCBI Build 35

human: *Homo sapiens*; mouse: *Mus musculus*

hg: human genome; GRC: Genome Reference Consortium

基因组序列不是确定的吗？也需要版本升级？

SPECIES	UCSC	DATE	NCBI
Human	hg19	Feb. 2009	Genome Reference Consortium GRCh37
	hg18	Mar. 2006	NCBI Build 36.1
	hg17	May 2004	NCBI Build 35
	hg16	Jul. 2003	NCBI Build 34
Mouse	mm10	Dec. 2011	Genome Reference Consortium GRCm38
	mm9	Jul. 2007	NCBI Build 37
	mm8	Feb. 2006	NCBI Build 36
	mm7	Aug. 2005	NCBI Build 35

human: *Homo sapiens*; mouse: *Mus musculus*

hg: human genome; GRC: Genome Reference Consortium



- UCSC liftOver tool：支持 BED 和 “chrN:start-end” 格式的输入
- Galaxy (基于 UCSC liftOver tool)：支持 BED、GFF 和 GTF 格式的输入
- CrossMap：支持 SAM/BAM、Wiggle/BigWig、BED、GFF/GTF 和 VCF 格式的输入，输出对应格式
- NCBI Remap：支持 BED、GFF、GTF 和 VCF 等格式的输入
- Ensembl assembly converter (2015 年退休, CrossMap 继位)：支持 BED、GFF、GTF 和 PSL 格式的输入，但输出都是 GFF 格式的
- pyliftover：仅支持点坐标 (point coordinates) 的转换，无法对区段 (ranges) 坐标进行转换



教学提纲

1

引言

2

基因组组装版本

3

基因组坐标系统

4

基因组注释常用格式

5

文本文件与文本编辑器

6

基因组坐标的逻辑运算

7

总结与答疑

8

引言

9

变异位点的注释

10

基因集富集分析

11

序列标识

12

box plot

13

解析图表

14

总结与答疑

15

引言

16

Galaxy 分析平台

17

Galaxy 使用演示

18

数据处理三段论

19

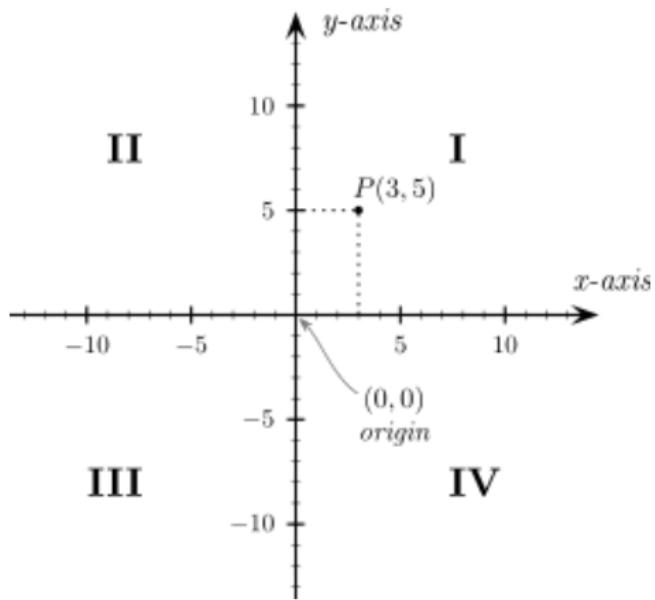
总结与答疑

20

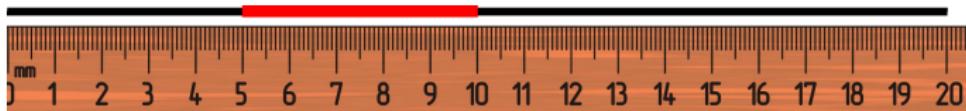
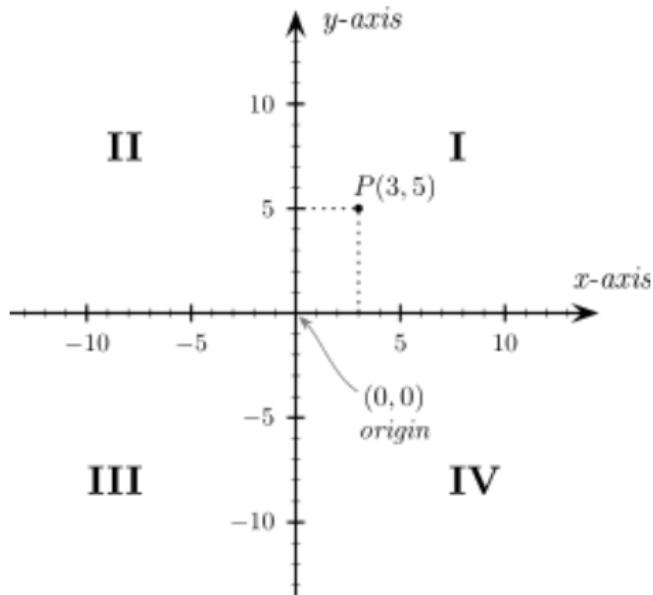
复习思考题



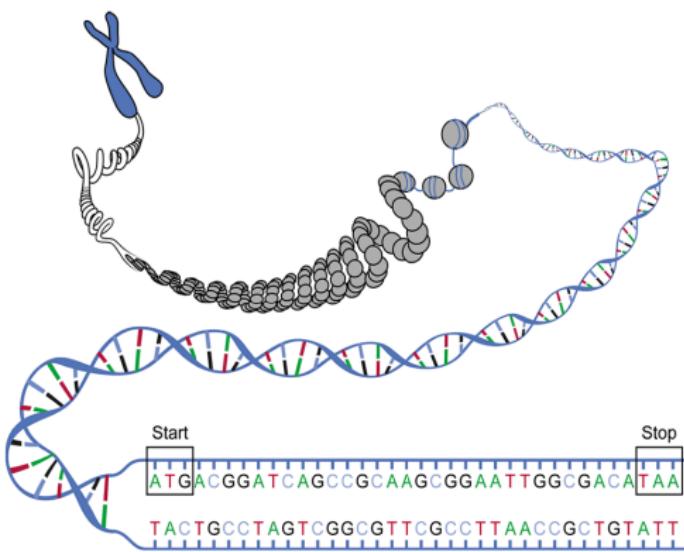
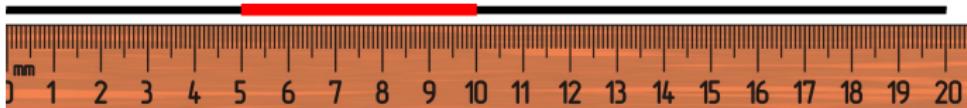
坐标系统 | 坐标轴



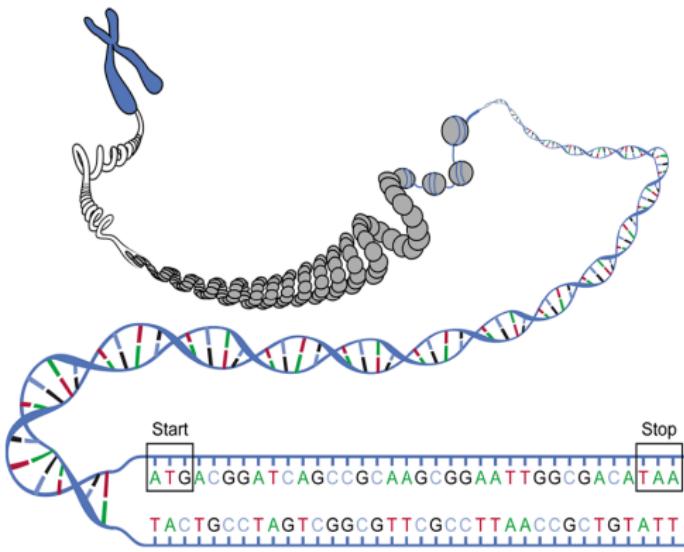
坐标系统 | 坐标轴



坐标系统 | 坐标轴



坐标系统 | 坐标轴



hg19

- SNP, rs1800468: "chr19 41860587"; "chr19:41860587"
- gene, *SAMD11*: "chr1 861121 879961"; "chr1:861121-879961"

坐标系统 | 两大系统

序列

0-based index	0	1	2	3	4	5	6	7
Sequence	A	A	T	T	G	G	C	C
1-based index	1	2	3	4	5	6	7	8

TG 的坐标

- 0-based, half-open: [3,5)
- 1-based, fully-closed: [4,5]

实例

- 0-based: BED、BAM、PSL、dbSNP、Table Browser
- 1-based: GFF、VCF、SAM、Wiggle、DAS、Genome Browser

坐标系统 | 两大系统

序列

0-based index	0	1	2	3	4	5	6	7
Sequence	A	A	T	T	G	G	C	C
1-based index	1	2	3	4	5	6	7	8

TG 的坐标

- 0-based, half-open: [3,5)
- 1-based, fully-closed: [4,5]

实例

- 0-based: BED、BAM、PSL、dbSNP、Table Browser
- 1-based: GFF、VCF、SAM、Wiggle、DAS、Genome Browser

坐标系统 | 两大系统

序列

0-based index	0	1	2	3	4	5	6	7
Sequence	A	A	T	T	G	G	C	C
1-based index	1	2	3	4	5	6	7	8

TG 的坐标

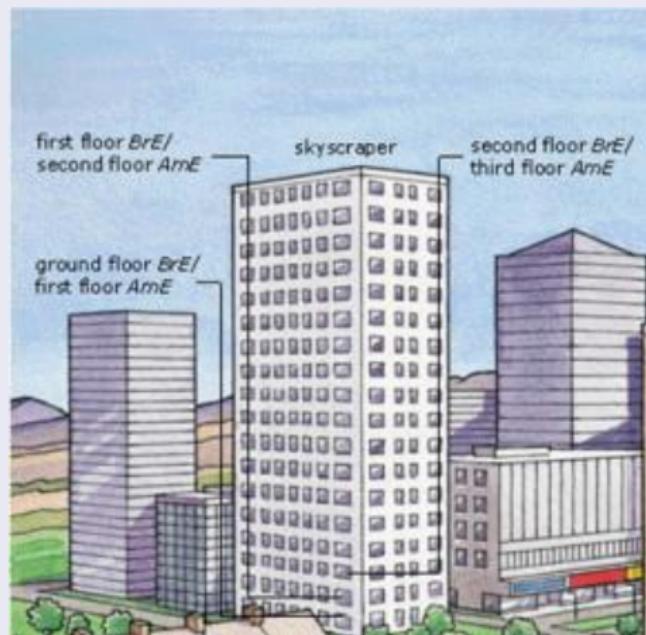
- 0-based, half-open: [3,5)
- 1-based, fully-closed: [4,5]

实例

- 0-based: BED、BAM、PSL、dbSNP、Table Browser
- 1-based: GFF、VCF、SAM、Wiggle、DAS、Genome Browser

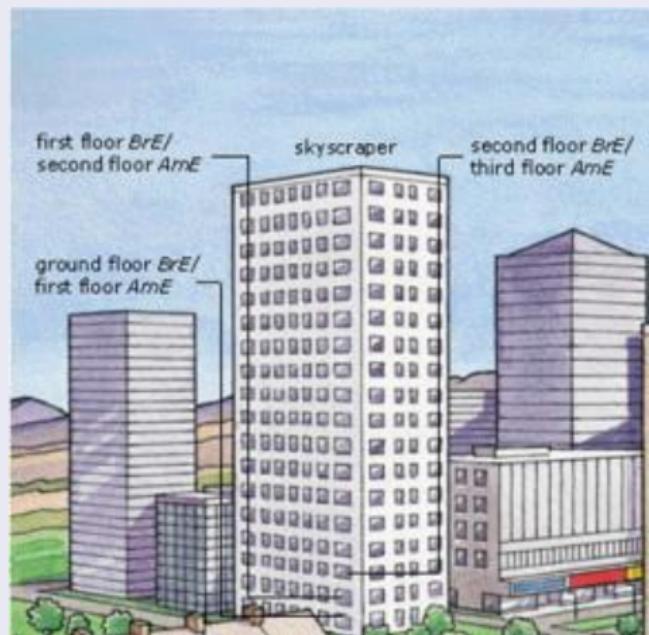
坐标系统 | 类比

first floor

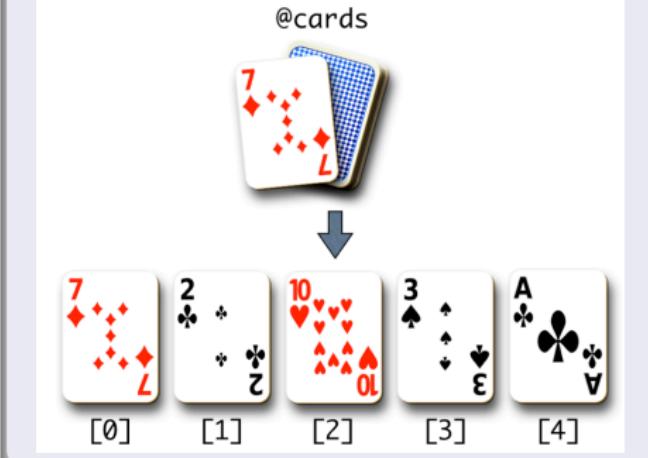


坐标系统 | 类比

first floor



数组



教学提纲

1

引言

2

基因组组装版本

3

基因组坐标系统

4

基因组注释常用格式

5

文本文件与文本编辑器

6

基因组坐标的逻辑运算

7

总结与答疑

8

引言

9

变异位点的注释

10

基因集富集分析

11

序列标识

12

box plot

13

解析图表

14

总结与答疑

15

引言

16

Galaxy 分析平台

17

Galaxy 使用演示

18

数据处理三段论

19

总结与答疑

20

复习思考题



格式 | 文件格式



EPS



FLA



HTML



GIF



IND



PHP



PPT



PDF



XLS



DOC



MOV



CSS



WAV



JPG



ZIP



MP3

格式 | FASTA

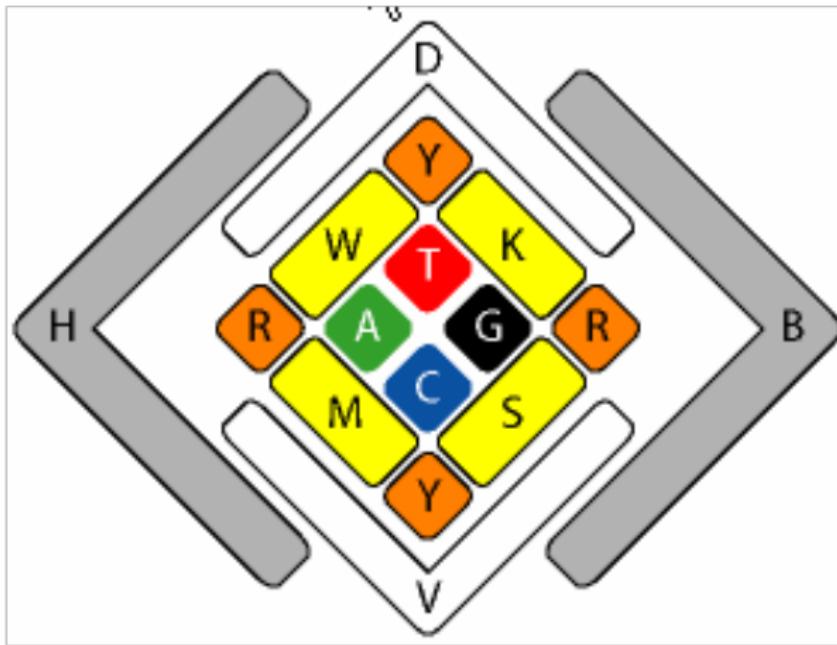
>gi|183121|gb|M29645.1|HUMGFI Human insulin-like growth factor II mRNA, complete cds
CAGGGGCCAAGAGTCACCAACCGAGCTGTGTGGAGGAGGTGGATTCCAGCCCCAGCCCCAGGGCTCT
GAATCGCTGCCAGCTCAGCCCCCTGCCAGCCTGCCACAGCCTGAGCCCCAGCAGGCCAGAGAGCCA
GTCCTGAGGTGAGCTGCTGTGGCTGTGGCCAGGGCACCCAGCGCTCCAGAACGTGAGGCTGGCAGCCA
GCCCCAGCCTCAGCCCCAACTGCGAGGCAGAGAGACACCAATGGGAATGCAATGGGAAGTCGATGCTG
GTGCTTCTCACCTTCTGGCTTCGCCTCGTGCATTGCTGCTTACCGCCCCAGTGAGACCCGTGCG
GCGGGGAGCTGGGGACACCCCTCCAGTTCTGTGTGGGGACCGCGCTCTACTTCAGCAGGCCGCAAG
CCGTGTGAGCCGTCGCAGCGTGGCATCGTGAGGAGTGCTGTTCCGCACTGTGACCTGGCCCTCTG
GAGACGTACTGTGCTACCCCGCCAAGTCCGAGAGGGACGTGTCGACCCCTCCGACCGTGCTCCGGACA
ACTTCCCCAGATAACCCGTGGCAAGTTCTCAATATGACACCTGGAAGCAGTCCACCCAGCGCCTGCG
CAGGGGCCTGCCTGCCCTCTGCGTGCCGCCGGGTACGTGCTGCCAAGGAGCTGAGGCCTTCAGG
GAGGCCAAACGTACCGTCCCTGATTGCTTACCCACCCAAAGACCCGCCACGGGGCGCCCCCCCAG
AGATGGCCAGCAATCGGAAGT GAGCAAACACTGCCAAGTCTGAGCCGGCGCCACCATCCTGAGCCT
CCTCTGACCACGGACGTTCCATCAGGTTCCATCCGAAATCTCGGTTCCACGTCCCCCTGGGCTT
CTCCTGACCCAGTCCCCGTGCCCGCCTCCCCGAAACAGGCTACTCTCCTCGGCCCCCTCCATGGGCTG
AGGAAGCACAGCAGCATCTCAAACATGTACAAATGATTGGCTTAAACACCTTCACATACCT



- 每一行最好不要超过 80 个字符
- 序列中的换行符不会影响序列的连续性
- 使用标准的 IUB/IUPAC 核酸代码和氨基酸代码
- 允许小写字母的存在，但会转换成大写
- 单个 “-” 代表不明长度的空位
- 在氨基酸序列中允许出现 “U” 和 “*”
- 任何数字都应该被去掉或转换成字母
- 不明核酸和氨基酸分别用 “N” 和 “X” 表示



格式 | FASTA | IUB/IUPAC 核酸



格式 | FASTA | IUB/IUPAC 核酸

Code	Meaning	Code	Meaning
A	Adenine	Y	Pyrimidine (C, T, or U)
C	Cytosine	K	T, U, or G (keto)
G	Guanine	W	T, U, or A (weak)
T	Thymine	B	C, T, U, or G (not A)
U	Uracil	D	A, T, U, or G (not C)
R	Purine (A or G)	H	A, T, U, or C (not G)
S	C or G (strong)	V	A, C, or G (not T, not U)
M	C or A (amino)	N	Any base (A, C, G, T, or U)
X	masked	-	gap of indeterminate length



格式 | FASTA | IUB/IUPAC 氨基酸

1	3	Meaning	1	3	Meaning
A	Ala	Alanine	B	Asx	Aspartic acid or Asparagine
C	Cys	Cysteine	D	Asp	Aspartic acid
E	Glu	Glutamic acid	F	Phe	Phenylalanine
G	Gly	Glycine	H	His	Histidin
I	Ile	Isoleucine	K	Lys	Lysine
L	Leu	Leucine	M	Met	Methionine
N	Asn	Asparagine	P	Pro	Proline
Q	Gln	Glutamine	R	Arg	Arginine
S	Ser	Serine	T	Thr	Threonine
U	Sec	Selenocysteine	V	Val	Valine
W	Trp	Tryptophan	X	Xaa	Any amino acid
Y	Tyr	Tyrosine	Z	Glx	Glutamine or Glutamic acid
*		translation stop	-		gap of indeterminate length
O	Pyl	Pyrrolysine			



格式 | FASTA | FASTA vs. Sequence

FASTA

```
>gi|183121|gb|M29645.1|HUMGFI Human insulin-like growth factor II mRNA, complete cds  
CAGGGGCCAAGAGTCACCACCGAGCTGTGAGGGAGGTGATTCCAGCCCCAGCCCCAGGGCTCT  
GAATCGCTGCCAGCTCAGCCCCCTGCCAGCCTGCCACAGCCTGAGCCCAGCAGGCCAGAGAGCCCA  
GTCCTGAGGTGAGCTGCTGTGGCTGTGGCCAGGGCACCCAGCCTCCAGAACTGAGGCTGGCAGCCA  
GCCCCAGCCTCAGCCCCAACCTGCAGGGCAGAGAGACACCAATGGGAATGCCAATGGGAAGTCGATGCTG  
GTGCTTCTCACCTTCTTGGCTTCGCCCTCGTGTGCATTGCTGCTTACGCCAGTGAGACCCCTGTGCG  
GCGGGGAGCTGGGGACACCCCTCAGTTGTCTGTGGGGACCCGGCTTACTTCAGCAGGCCAGAAG  
CCGTGTAGCCGTCGAGCCGTGGCATCGTGTGAGGAGTGTGTTCCGAGCTGTGACCTGGCCCTCG  
GAGACGTACTGTGCTACCCCCGCCAAGTCCAGAGAGGGACGTGTCACCCCTCGACCGTGCTTCCGGACA  
ACTTCCCCAGATAACCCGTGGCAAGTTCTCAATATGACACCTGGAAGCAGTCCACCCAGCGCCTGCG  
CAGGGGCTGCCCTCTCGTGTGCCGCCGGGTCACGTGCTGCCAAGGAGCTGAGGGCTTCAGG  
GAGGCCAAACGTACCGTCCCCCTGATTGCTCTACCCACCCAAGACCCGCCACGGGGCGCCCCCAG  
AGATGGCCAGCAATCGGAAGTGTGAGCAAAACTGCCAAGTCTGAGCCGGCGCCACCATCCTGAGCCT  
CCTCTGACCACGGACGTTCCATCAGGTTCCATCCGAAATCTCTGGTCCACGTCCCCCTGGGCTT  
CTCCTGACCCAGTCCCCGTCCCCGCCCTCCCGAAACAGGCTACTCTCTCGGCCCCCTCATGGGCTG  
AGGAAGCACAGCAGCATTTCAAACATGTACAAAATGATTGGCTTAAACACCTTACATACCT
```

Sequence

- GTACGACGGAGTGTATAAGATGGAAATCGGATACCAGATGAAATTGTGGATCAG
- MWTALPLLCAGAWLLSAGATAELTVNAIEKFHFTSWMKQHQKTYSSREYSHRLQVFAN

格式 | BED (Browser Extensible Data)

chr7	127471196	127472363	Pos1	0	+	127471196	127472363	255,0,0
chr7	127472363	127473530	Pos2	0	+	127472363	127473530	255,0,0
chr7	127473530	127474697	Pos3	0	+	127473530	127474697	255,0,0
chr7	127474697	127475864	Pos4	0	+	127474697	127475864	255,0,0
chr7	127475864	127477031	Neg1	0	-	127475864	127477031	0,0,255
chr7	127477031	127478198	Neg2	0	-	127477031	127478198	0,0,255
chr7	127478198	127479365	Neg3	0	-	127478198	127479365	0,0,255
chr7	127479365	127480532	Pos5	0	+	127479365	127480532	255,0,0
chr7	127480532	127481699	Neg4	0	-	127480532	127481699	0,0,255



BED#

BED12 包含全部 12 列

BED6 chrom, start, end, name, score, and strand

BED5 chrom, start, end, name, and score

BED4 chrom, start, end, and name

BED3 chrom, start, and end

例子

chr1	11873	14409	uc001aaa.3	0	+	11873	11873	0
3	354,109,1189,	0,739,1347,						



BED#

BED12 包含全部 12 列

BED6 chrom, start, end, name, score, and strand

BED5 chrom, start, end, name, and score

BED4 chrom, start, end, and name

BED3 chrom, start, and end

例子

chr1 11873 14409 uc001aaa.3 0 +



BED#

BED12 包含全部 12 列

BED6 chrom, start, end, name, score, and strand

BED5 chrom, start, end, name, and score

BED4 chrom, start, end, and name

BED3 chrom, start, and end

例子

chr1 11873 14409 uc001aaa.3 0



BED#

BED12 包含全部 12 列

BED6 chrom, start, end, name, score, and strand

BED5 chrom, start, end, name, and score

BED4 chrom, start, end, and name

BED3 chrom, start, and end

例子

chr1 11873 14409 uc001aaa.3



BED#

BED12 包含全部 12 列

BED6 chrom, start, end, name, score, and strand

BED5 chrom, start, end, name, and score

BED4 chrom, start, end, and name

BED3 chrom, start, and end

例子

chr1 11873 14409



格式 | GFF (General Feature Format)

```
##gff-version 3
ctg123 . operon      1300 15000  .  +  .  ID=operon001;Name=superOperon
ctg123 . mRNA        1300 9000   .  +  .  ID=mrna0001;Parent=operon001;Name=sonichedgehog
ctg123 . exon         1300 1500   .  +  .  Parent=mrna0001
ctg123 . exon         1050 1500   .  +  .  Parent=mrna0001
ctg123 . exon         3000 3902   .  +  .  Parent=mrna0001
ctg123 . exon         5000 5500   .  +  .  Parent=mrna0001
ctg123 . exon         7000 9000   .  +  .  Parent=mrna0001
ctg123 . mRNA        10000 15000  .  +  .  ID=mrna0002;Parent=operon001;Name=subsonicsquirrel
ctg123 . exon         10000 12000  .  +  .  Parent=mrna0002
ctg123 . exon         14000 15000  .  +  .  Parent=mrna0002
```



格式 | VCF (Variant Call Format)

```
##fileformat=VCFv4.0
##fileDate=20110705
##reference=1000GenomesPilot-NCBI37
##phasing=partial
##INFO=<ID=NS,Number=1,Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=.,Type=Float,Description="Allele Frequency">
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP membership, build 129">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FILTER=<ID=q10,Description="Quality below 10">
##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##FORMAT=<ID=HQ,Number=2,Type=Integer,Description="Haplotype Quality">
#CHROM POS ID REF ALT QUAL FILTER INFO
2 4370 rs6057 G A 29 . NS=2;DP=13;AF=0.5;DB;H2
2 7330 . T A 3 q10 NS=5;DP=12;AF=0.017
2 110696 rs6055 A G,T 67 PASS NS=2;DP=10;AF=0.333,0.667;AA=T;DB
2 130237 . T . 47 . NS=2;DP=16;AA=T
2 134567 microsat1 GTCT G,GTACT 50 PASS NS=2;DP=9;AA=G
```

	FORMAT	Sample1	Sample2	Sample3
2 4370 rs6057 G A 29 . NS=2;DP=13;AF=0.5;DB;H2	GT:GQ:DP:HQ	0 0:48:1:52,51 1 0:48:8:51,51	1/1:43:5:,,	
2 7330 . T A 3 q10 NS=5;DP=12;AF=0.017	GT:GQ:DP:HQ	0 0:46:3:58,50 0 1:3:5:65,3	0/0:41:3	
2 110696 rs6055 A G,T 67 PASS NS=2;DP=10;AF=0.333,0.667;AA=T;DB	GT:GQ:DP:HQ	1 2:21:6:23,27 2 1:2:0:18,2	2/2:35:4	
2 130237 . T . 47 . NS=2;DP=16;AA=T	GT:GQ:DP:HQ	0 0:54:7:56,60 0 0:48:4:56,51	0/0:61:2	
2 134567 microsat1 GTCT G,GTACT 50 PASS NS=2;DP=9;AA=G	GT:GQ:DP	0 1:35:4	0/2:17:2	1/1:40:3



教学提纲

1

引言

2

基因组组装版本

3

基因组坐标系统

4

基因组注释常用格式

5

文本文件与文本编辑器

6

基因组坐标的逻辑运算

7

总结与答疑

8

引言

9

变异位点的注释

10

基因集富集分析

11

序列标识

12

box plot

13

解析图表

14

总结与答疑

15

引言

16

Galaxy 分析平台

17

Galaxy 使用演示

18

数据处理三段论

19

总结与答疑

20

复习思考题



文本 | 纯文本 vs. 格式化文本

P8_Ain_Pro - 记事本
文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)
CLUSTAL X (1.83) multiple sequence alignment

RGDV_ABC75537 HSRQWIEITSALIECISEVGTCKCSFDTFQGLTINDISTLSNLHNQISUASUGFLNDPRTP
RGDV_AAY14576 HSRQWIEITSALIECISEVGTCKCSFDTFQGLTINDISTLSNLHNQISUASUGFLNDPRTP
RGDV_BAA02676 HSRQWIEITSALIECISEVGTCKCSFDTFQGLTINDISTLSNLHNQISUASUGFLNDPRTP
RGDV_AAY14579 HSRQWIEITSALIECISEVGTCKCSFDTFQGLTINDISTLSNLHNQISUASUGFLNDPRTP
RGDV_AAY14580 HSRQWIEITSALIECISEVGTCKCSFDTFQGLTINDISTLSNLHNQISUASUGFLNDPRTP
RGDV_AAO04253 HSRQWIEITSALIECISEVGTCKCSFDTFQGLTINDISTLSNLHNQISUASUGFLNDPRTP
RGDV_AAY14577 HSRQWIEITSALIECISEVGTCKCSFDTFQGLTINDISTLSNLHNQISUASUGFLNDPRTP
RGDV_AAY14578 HSRQWIEITSALIECISEVGTCKCSFDTFQGLTINDISTLSNLHNQISUASUGFLNDPRTP
WTU_P17380 HSRQHWEUETSALLEAISEYUVRUNGDTFSGLTTGDFNALSNNHTQLSUSSAGYUSDPRUP
RDVS_Q85451 HSRQHMDLTSALLEAISEYUVRUNGDTFSGLTTGDFNALSNNHTQLSUSSAGYUSDPRUP
RDVS_P17379 HSRQHMDLTSALLEAISEYUVRUNGDTFSGLTTGDFNALSNNHTQLSUSSAGYUSDPRUP
RDVA_Q85449 HSRQHMDLTSALLEAISEYUVRUNGDTFSGLTTGDFNALSNNHTQLSUSSAGYUSDPRUP
RDUF_Q85439 **** *:***** * ***** .. . *** *** *.*:*****.*****..*:***.*.*

RGDV_ABC75537 LQAHSCFUFNSTADRHAYHLQKNMFDSDUAPNUTDNFIATYIKPRFRSRTUSDULRQU
RGDV_AAY14576 LQAHSCFUFNSTADRHAYHLQKNMFDSDUAPNUTDNFIATYIKPRFRSRTUSDULRQU
RGDV_BAA02676 LQAHSCFUFNSTADRHAYHLQKNMFDSDUAPNUTDNFIATYIKPRFRSRTUSDULRQU
RGDV_AAY14579 LQAHSCFUFNSTADRHAYHLQKNMFDSDUAPNUTDNFIATYIKPRFRSRTUSDULRQU
RGDV_AAY14580 LQAHSCFUFNSTADRHAYHLQKNMFDSDUAPNUTDNFIATYIKPRFRSRTUSDULRQU
RGDV_AAO04253 LQAHSCFUFNSTADRHAYHLQKNMFDSDUAPNUTDNFIATYIKPRFRSRTUSDULRQU
RGDV_AAY14577 PQAHSCFUFNSTADRHAYHLQKNMFDSDUAPNUTDNFIATYIKPRFRSRTUSDULRQU

Ln 1, Col 1

Work on Word docs at the same time with others using SkyDrive.docx - M... Picture Tools Format
File Home Insert Page Layout References Mailings Review View
Clipboard Font Paragraph Styles
collaborate. Today, I'm thrilled to announce that we are taking a step towards improving collaboration – by bringing simultaneous editing to Word Web App in addition to Microsoft Word 2010, Microsoft Word for Mac 2011. (Word now joins Excel and OneNote on the web with simultaneous editing.)
It's no secret that we love Word around here. It's used to author all the specifications we write for our products, as well as all the blog posts we publish on this blog. These are professional documents that we produce all the time and there are countless millions of people that have been using Word to express and communicate thoughts and ideas.
Word has a long history of innovation. I remember when the red squiggle line arrived and I no longer had to manually spell check all the school papers I wrote. I also remember when Word introduced auto-correct and many common mistakes were corrected for me (still get corrected to this day). I remember when Outlook started to use Word as its default mail editor and all the power of Word arrived in the place that I write the most. In short, I have grown up with Word, first on the Mac, and now on the PC and am happy we can offer yet another feature to enable you to be more productive.
In fact, this post was written using Word 2010 on my ThinkPad, while my colleague Harrison has been making changes to this post such as inserting screen shots to demonstrate certain word features. And we even made a video of us doing this. It doesn't get more meta than this!
<insert video>
Here are some of the features that showcase how Word communicates to you about changes collaborators are making
Notifications of other collaborators
Word Web App Real time collaboration in Word on SkyDrive
Harrison Hoffman sign out
File Home Insert Page Layout References Mailings Review View
Clipboard Font Paragraph Styles
Other Sharron is editing this document
Work on Word docs at the same time with others using SkyDrive
Words 488 5 100% LiveSync 44

TIANJIN MEDICAL UNIVERSITY

三大类

- Windows: \r\n (CR+LF, 回车 + 换行) , 文件尾部直接 EOF (文件结束标志)
- Unix: \n (LF, 仅有换行) , 文件最后一行也会增加该字符, 然后才是 EOF
- Mac: \r (CR, 仅有回车)

识别与转换

- Windows: 文本编辑器, 如 Notepad++
- Unix: file 识别, fromdos & todos 转换



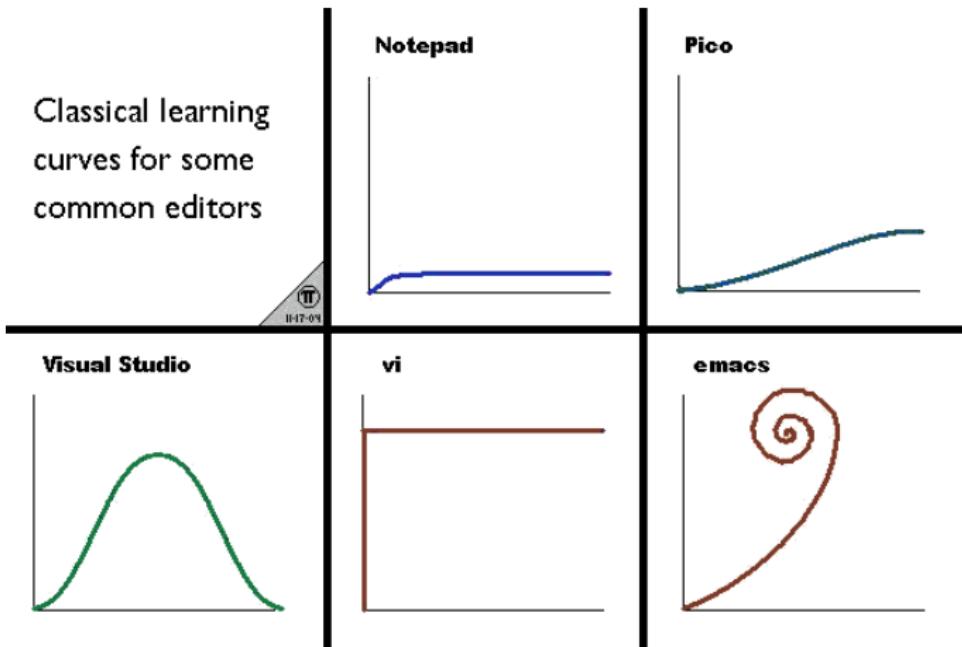
三大类

- Windows: \r\n (CR+LF, 回车 + 换行) , 文件尾部直接 EOF (文件结束标志)
- Unix: \n (LF, 仅有换行) , 文件最后一行也会增加该字符, 然后才是 EOF
- Mac: \r (CR, 仅有回车)

识别与转换

- Windows: 文本编辑器, 如 Notepad++
- Unix: file 识别, fromdos & todos 转换





文本 | 编辑器 | Notepad++

*C:\mediawiki\includes\Article.php - Notepad++

File Edit Search View Encoding Language Settings Macro Run Plugins Window ?

Article.php

```
2580 $modified = $current != '' && $protection != 'unprotected';  
2581  
2582 if ( $protect ) {  
2583     $comment_type = $modified ? 'modified' : 'unprotected';  
2584 } else {  
2585     $comment_type = 'unprotected';  
2586 }  
2587  
2588 $comment = $wgContLang->ucfirst( $comment_type );  
2589  
2590 # Only restrictions with the 'protect' key.  
# Otherwise, people who cannot normally edit  
$editrestriction = isset( $limit['protect'] ) ?  
2591  
# The schema allows multiple restrictions.  
2592 if ( !in( 'protect', $editrestriction ) ) {  
2593     $interface_exists = true;  
2594     $intval = true;  
2595     $cascader = array();  
2596     $ip2long = true;  
2597     $iptcembed = true;  
2598 }  
2599  
2600 if ( $cascader ) {  
    $cascader_description = 'The following cascader settings apply to this page:  
    ' . implode( ', ', $cascader );  
}
```

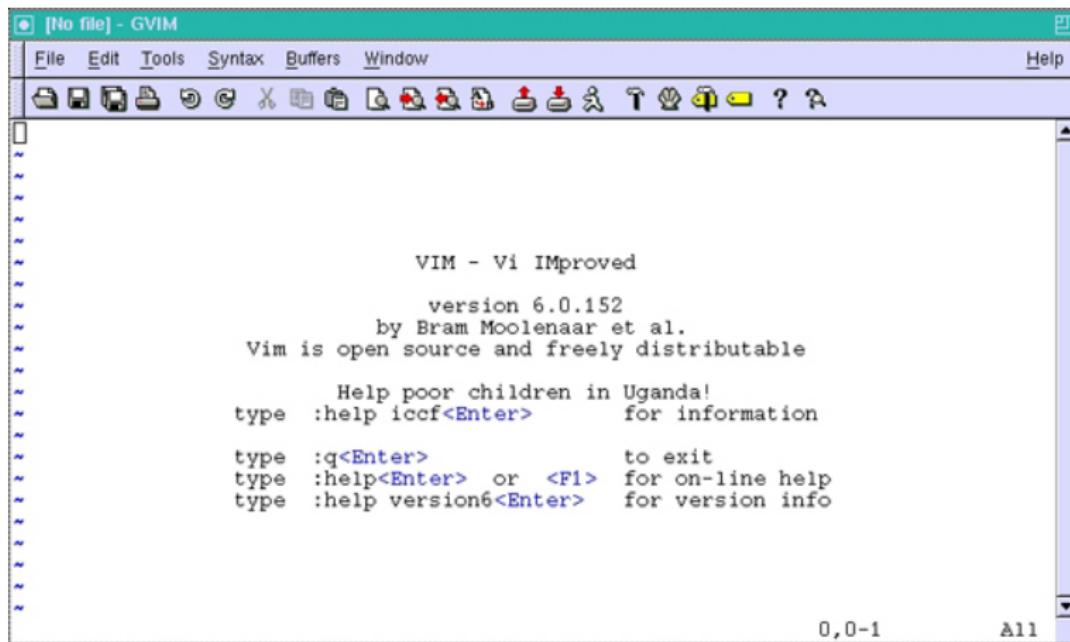
load.php

```
75  
76  
77  
78 if( $wgUseFileCache && isset( $wgFileCacheProfileIn( 'main-try-filecache' ) ) ) {  
79 // Raw pages should handle caching.  
80 // even when using file cache.  
81 if( $action != 'raw' && !HTMLForm::isLowLevelFileCache() ) {  
82     $cache = new HTMLFileCache();  
83     if( $cache->isFileCacheGood() ) {  
84         /* Check incoming headers */  
85         if( !$wgOut->checkLastModified() ) {  
86             $cache->loadFromFile();  
87         }  
88     }  
89 }  
90  
91  
92  
93  
94  
95  
96  
97  
98  
99  
100  
101  
102  
103  
104  
105  
106  
107  
108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161  
162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431  
432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1079  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1129  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187  
1188  
1189  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1239  
1240  
1241  
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1289  
1290  
1291  
1292  
1293  
1294  
1295  
1296  
1297  
1298  
1299  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349  
1349  
1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1399  
1400  
1401  
1402  
1403  
1404  
1405  
1406  
1407  
1408  
1409  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457  
1458  
1459  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1509  
1510  
1511  
1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1559  
1560  
1561  
1562  
1563  
1564  
1565  
1566  
1567  
1568  
1569  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619  
1619  
1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1669  
1670  
1671  
1672  
1673  
1674  
1675  
1676  
1677  
1678  
1679  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727  
1728  
1729  
1729  
1730  
1731  
1732  
1733  
1734  
1735  
1736  
1737  
1738  
1739  
1739  
1740  
1741  
1742  
1743  
1744  
1745  
1746  
1747  
1748  
1749  
1749  
1750  
1751  
1752  
1753  
1754  
1755  
1756  
1757  
1758  
1759  
1759  
1760  
1761  
1762  
1763  
1764  
1765  
1766  
1767  
1768  
1769  
1769  
1770  
1771  
1772  
1773  
1774  
1775  
1776  
1777  
1778  
1779  
1779  
1780  
1781  
1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1799  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1809  
1809  
1810  
1811  
1812  
1813  
1814  
1815  
1816  
1817  
1818  
1819  
1819  
1820  
1821  
1822  
1823  
1824  
1825  
1826  
1827  
1828  
1829  
1829  
1830  
1831  
1832  
1833  
1834  
1835  
1836  
1837  
1838  
1839  
1839  
1840  
1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1849  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1859  
1860  
1861  
1862  
1863  
1864  
1865  
1866  
1867  
1868  
1869  
1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889  
1889  
1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1939  
1940  
1941  
1942  
1943  
1944  
1945  
1946  
1947  
1948  
1949  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978  
1979  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997  
1998  
1999  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  

```



文本 | 编辑器 | Vim



[No file] - GVIM

File Edit Tools Syntax Buffers Window Help

VIM - Vi IMproved

version 6.0.152
by Bram Moolenaar et al.
Vim is open source and freely distributable

Help poor children in Uganda!
type :help icccf<Enter> for information

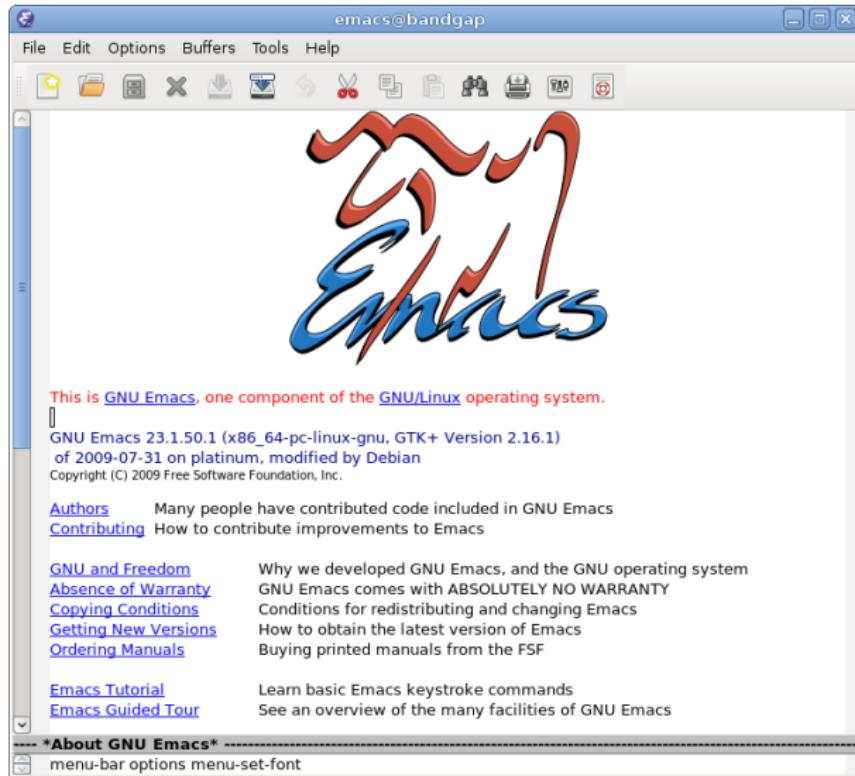
type :q<Enter> to exit
type :help<Enter> or <F1> for on-line help
type :help version6<Enter> for version info

0,0-1 All

This screenshot shows the GVIM interface with the Vim startup screen displayed. The title bar reads '[No file] - GVIM'. The menu bar includes File, Edit, Tools, Syntax, Buffers, Window, and Help. A toolbar with various icons is located above the main editor area. The main window displays the Vim startup message, which includes the version number (6.0.152), author (Bram Moolenaar et al.), and the fact that Vim is open source and freely distributable. It also provides instructions for seeking help, such as using ':help icccf<Enter>' for information and ':q<Enter>' to exit. The status bar at the bottom shows '0,0-1' and 'All'.



文本 | 编辑器 | Emacs



文本 | 编辑器 | Sublime Text

Soda Light.sublime-theme — Theme – Soda

OPEN FILES

- Soda Light.sublime-theme
- Soda Dark.sublime-theme

FOLDERS

- ▼ Theme – Soda
 - Soda Dark
 - Soda Light
- README.md
- Soda Dark.sublime-theme
- Soda Light.sublime-theme

Soda Light.sublime-theme

1 soda light

2

3 658 Soda Light.sublime-theme

4 Soda Light.sublime-theme

5 532 Widget – Soda Light.stTheme

6 Soda Light/Widget – Soda Light.stTheme

7

8 532 Widget – Soda Light.sublime-settings

9 Soda Light/Widget – Soda Light.sublime-settings

10 "layer0.texture": "Theme – Soda/Soda Light/tabsheet-
background.png",

11 "layer0.inner_margin": [1, 7],

12 "layer0.opacity": 1.0,

13 "content_margin": [-4, 0, -4, 3],

14 "tab_overlap": 5,

15 "tab_width": 180,

16 "tab_min_width": 45,

17 "tab_height": 25,

18 "mouse_wheel_switch": false

19 },

20 {

Find What: dark

Replace With: light

Find Replace All

Line 15, Column 1

Spaces: 4

JSON

```
"layer0.texture": "Theme – Soda/Soda Light/tabsheet-background.png",
"layer0.inner_margin": [1, 7],
"layer0.opacity": 1.0,
"content_margin": [-4, 0, -4, 3],
"tab_overlap": 5,
"tab_width": 180,
"tab_min_width": 45,
"tab_height": 25,
"mouse_wheel_switch": false}
```



教学提纲

1

引言

2

基因组组装版本

3

基因组坐标系统

4

基因组注释常用格式

5

文本文件与文本编辑器

6

基因组坐标的逻辑运算

7

总结与答疑

8

引言

9

变异位点的注释

10

基因集富集分析

11

序列标识

12

box plot

13

解析图表

14

总结与答疑

15

引言

16

Galaxy 分析平台

17

Galaxy 使用演示

18

数据处理三段论

19

总结与答疑

20

复习思考题



逻辑运算 | 常见问题

数据

gene1	chr1	10	20	+
gene2	chr1	40	60	+
gene3	chr1	50	100	+
snp1	chr1	15		+
snp2	chr1	55		-
exon3.1	chr1	50	60	+
exon3.2	chr1	90	100	+

问题

- ① 找到 gene1 和 gene2 之间的基因间区域。
- ② snp1 在 gene1 上吗? snp2 在 gene1 上吗 (, 在 gene2 上吗) ?
- ③ 找到与 gene3 重叠和不重叠的基因?
- ④ 找到 gene3 的内含子区域。

逻辑运算 | 常见问题

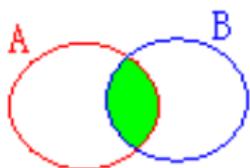
数据

gene1	chr1	10	20	+
gene2	chr1	40	60	+
gene3	chr1	50	100	+
snp1	chr1	15		+
snp2	chr1	55		-
exon3.1	chr1	50	60	+
exon3.2	chr1	90	100	+

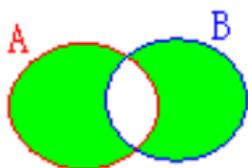
问题

- ① 找到 gene1 和 gene2 之间的基因间区域。
- ② snp1 在 gene1 上吗? snp2 在 gene1 上吗 (, 在 gene2 上吗) ?
- ③ 找到与 gene3 重叠和不重叠的基因?
- ④ 找到 gene3 的内含子区域。

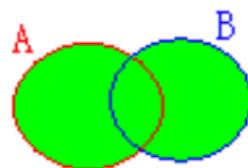
逻辑运算 | 集合运算



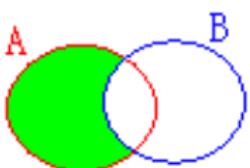
求同
交集



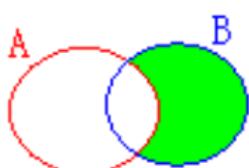
求异



相加
并集



相减 A-B

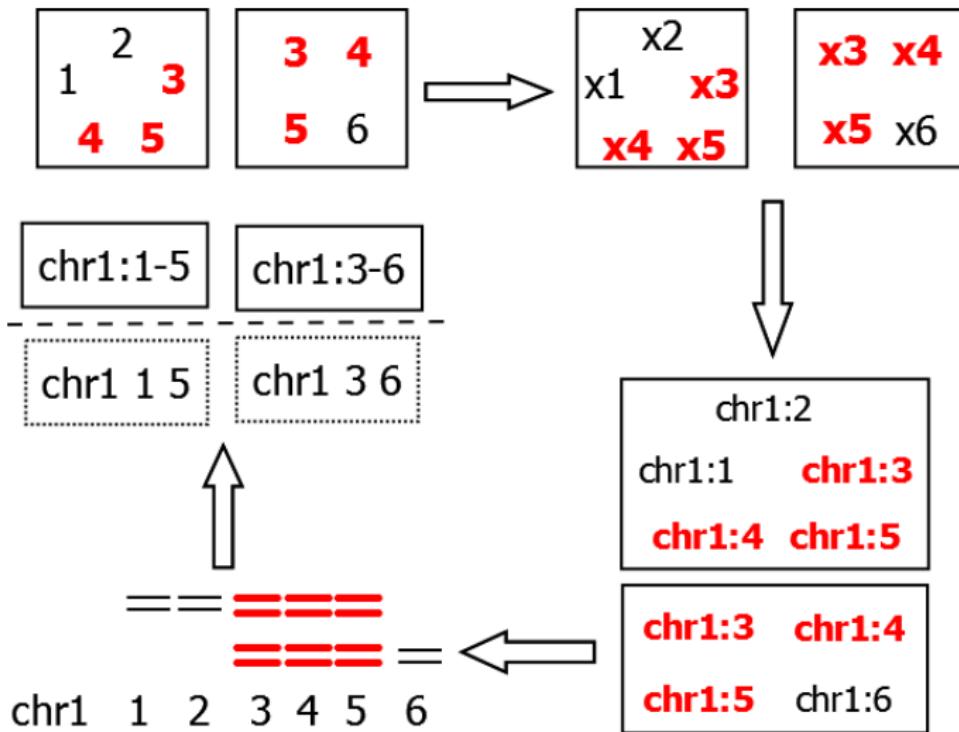


相减 B-A



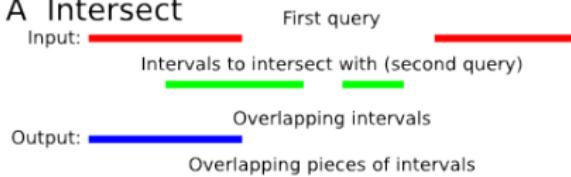
补集

逻辑运算 | 集合 \Rightarrow 基因组

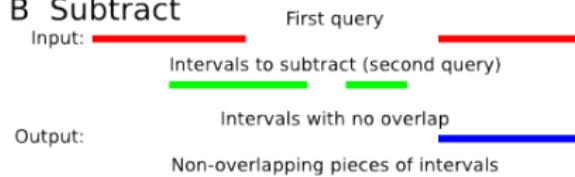


逻辑运算 | 运算模式

A Intersect



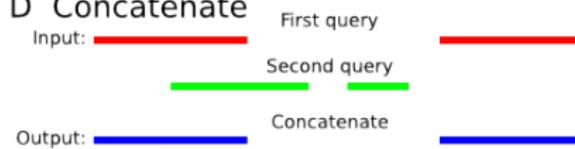
B Subtract



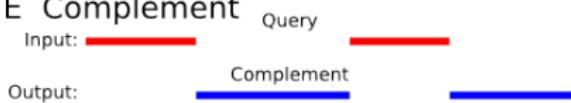
C Merge



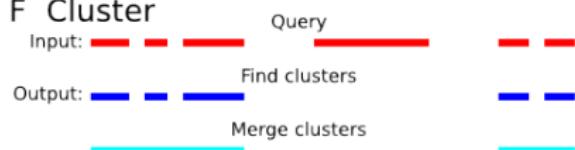
D Concatenate



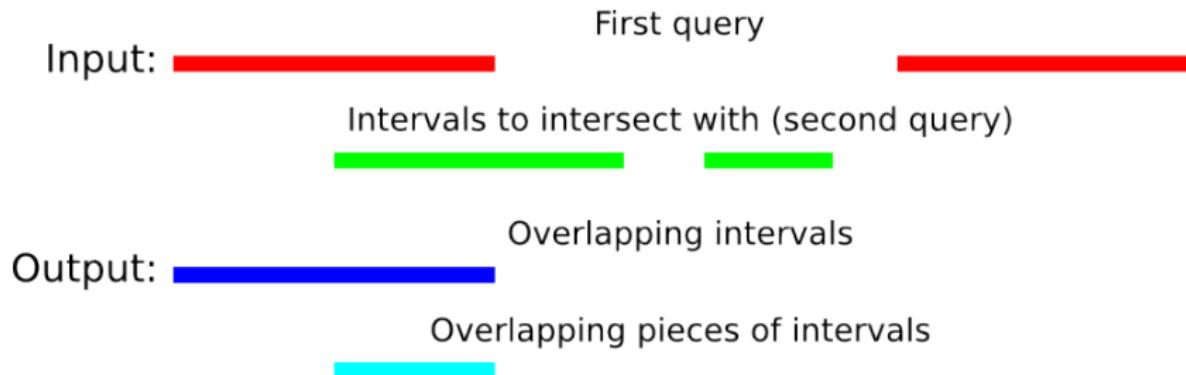
E Complement



F Cluster



逻辑运算 | intersect



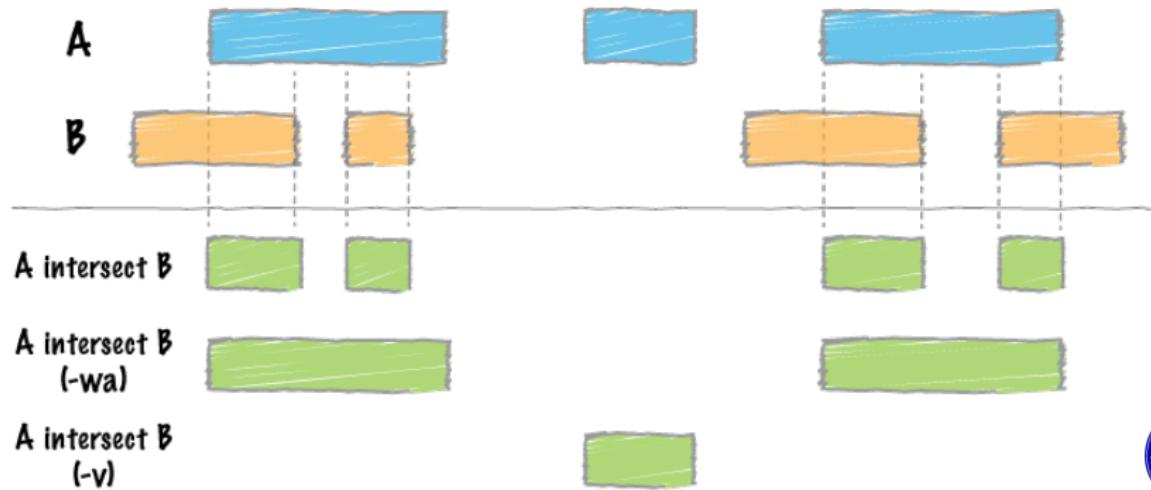
逻辑运算 | intersect

Chromosome ======

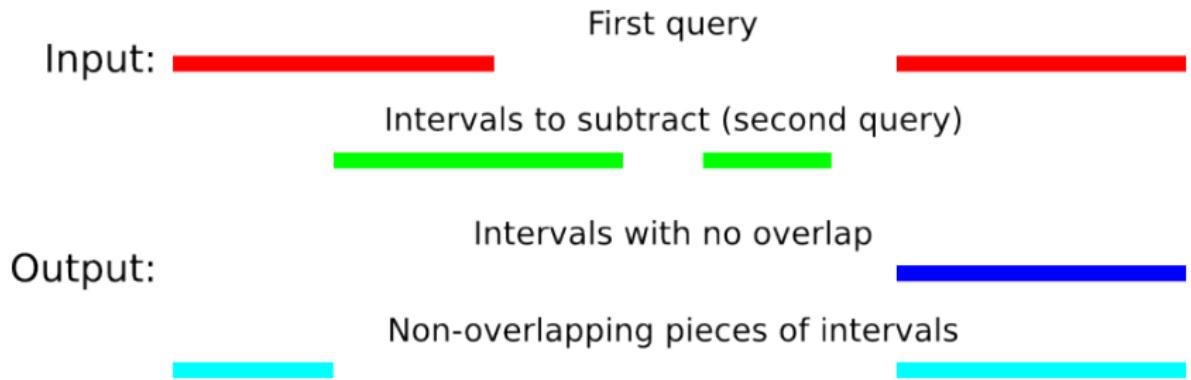
BED/BAM A ====== ======

BED File B ======

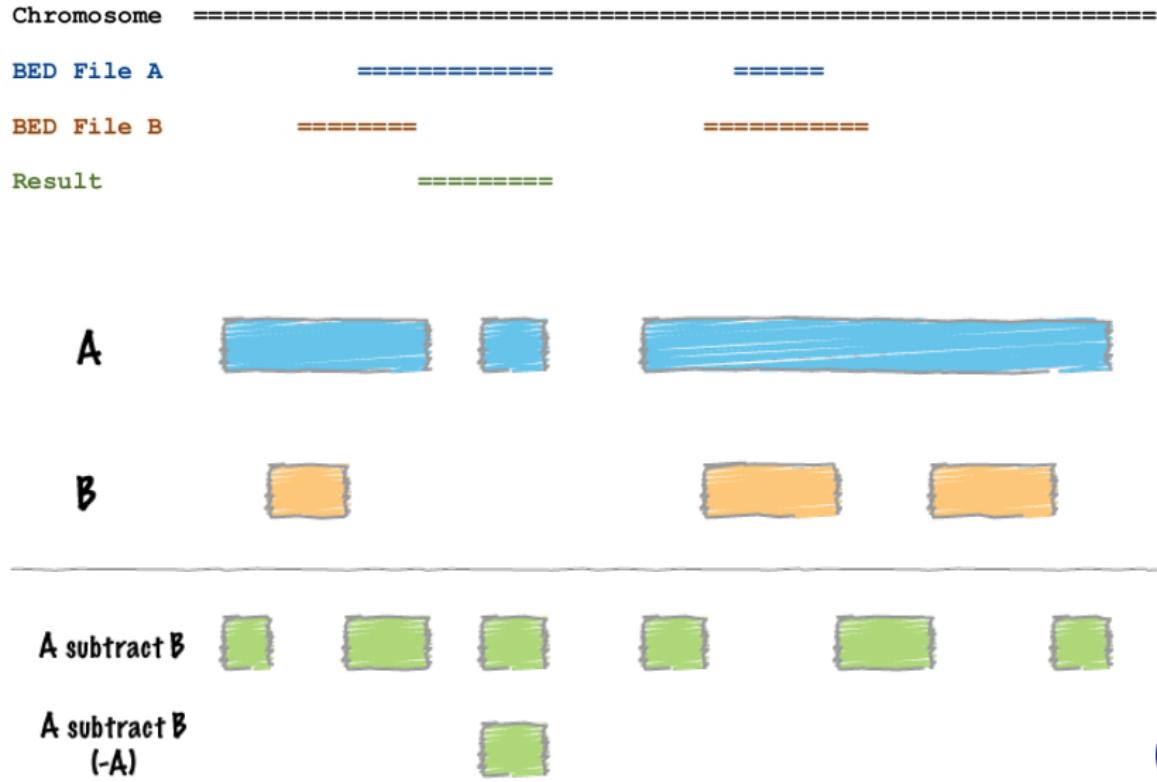
Result ======



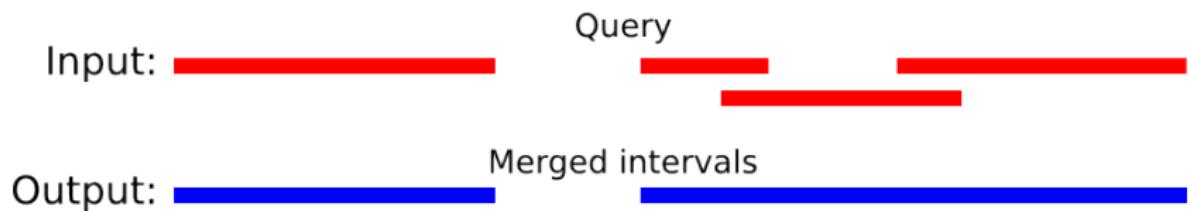
逻辑运算 | subtract



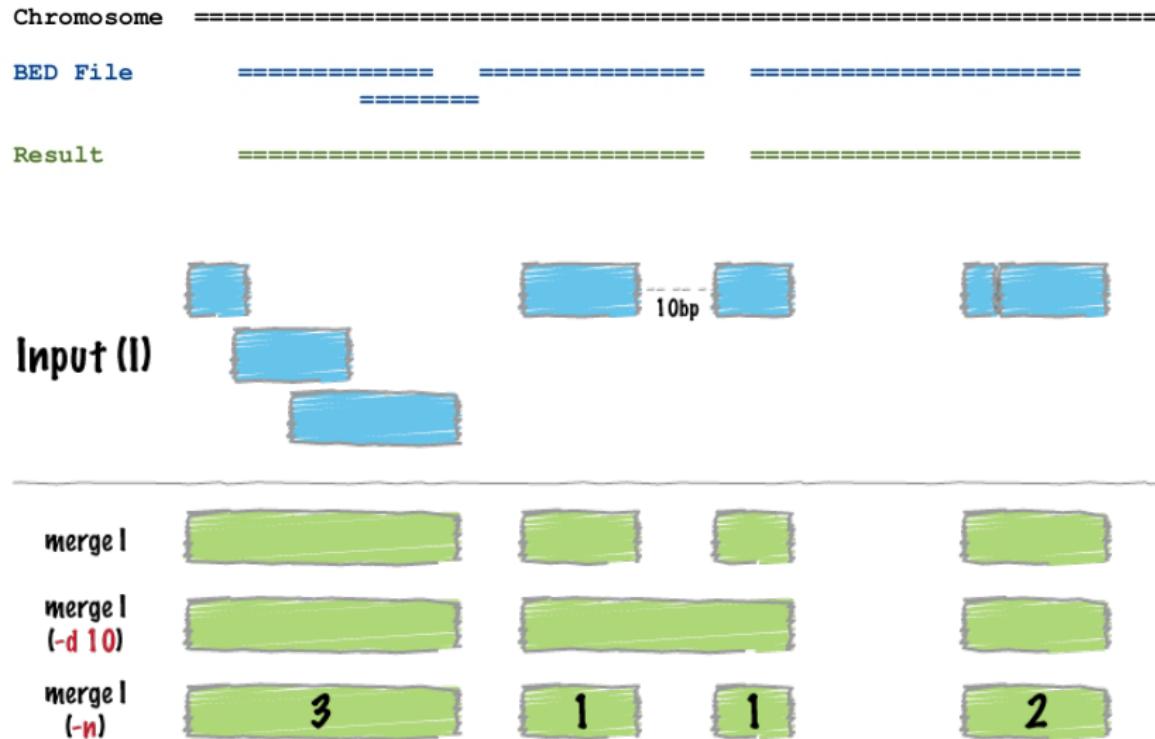
逻辑运算 | subtract



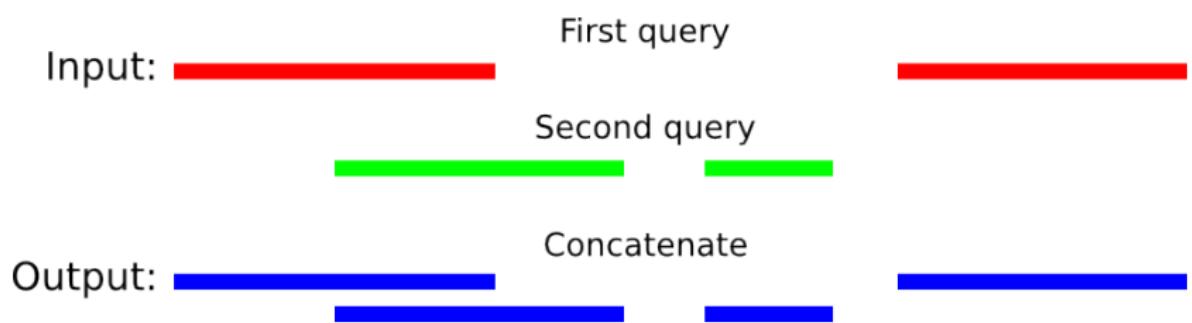
逻辑运算 | merge



逻辑运算 | merge



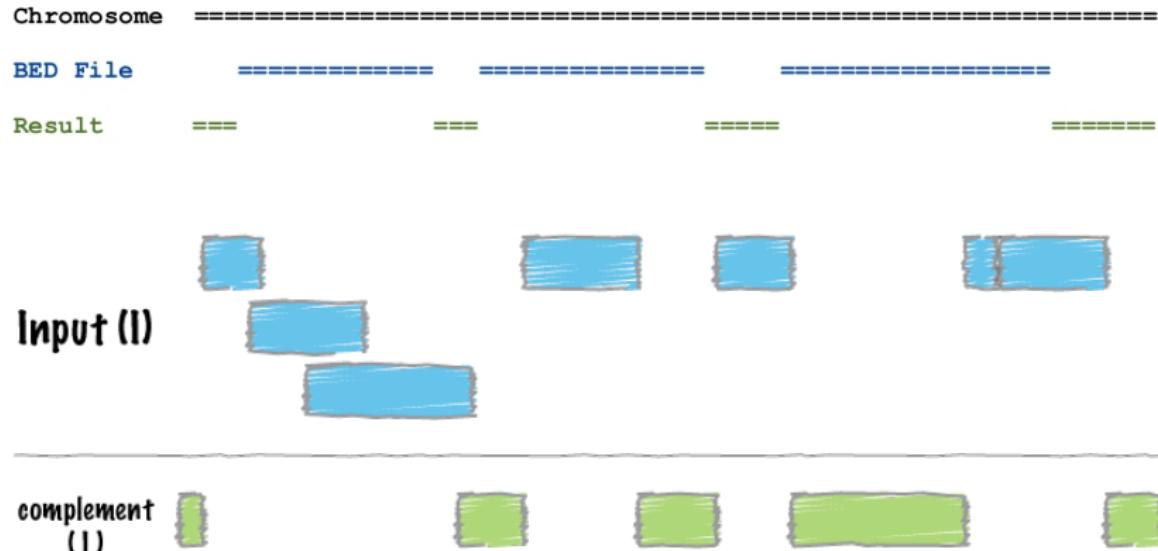
逻辑运算 | concatenate



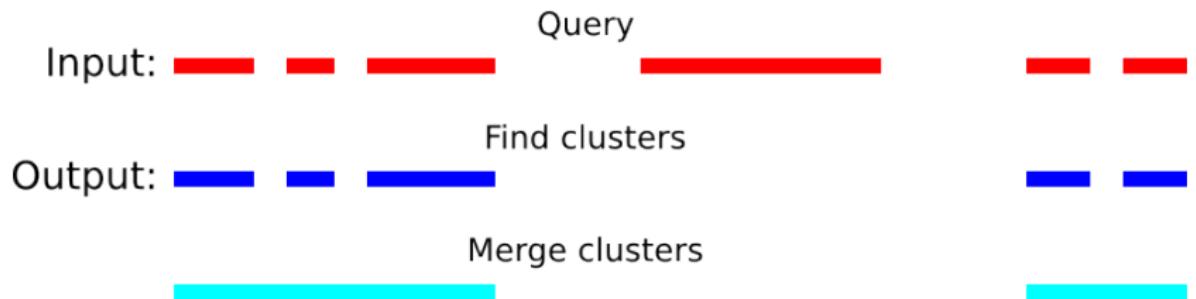
逻辑运算 | complement



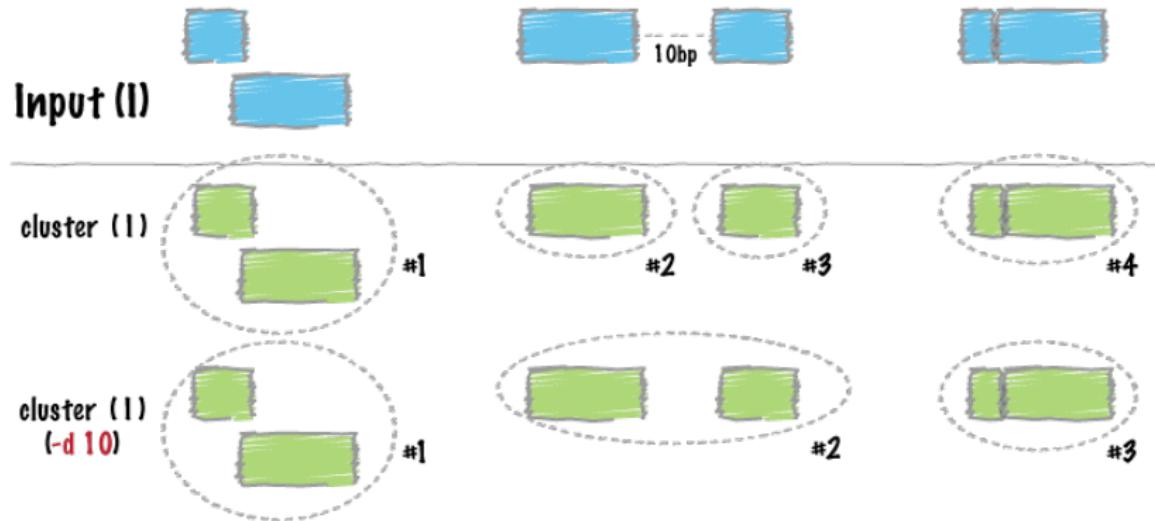
逻辑运算 | complement



逻辑运算 | cluster



逻辑运算 | cluster



逻辑运算 | join

Input						
Query 1:						
chr1 10 100 Query1.1						
chr1 500 1000 Query1.2						
chr1 1100 1250 Query1.3						
Query 2:						
chr1 20 80 Query2.1						
chr1 2000 2204 Query2.2						
chr1 2500 3000 Query2.3						
Output						
(Return only records that are joined)						
chr1 10 100 Query1.1 chr1 20 80 Query2.1						
Return only records that are joined (INNER JOIN) Return all records of first query (fill null with ".") Return all records of second query (fill null with ".") Return all records of both queries (fill nulls with ".")						



逻辑运算 | join

Input																											
Query 1:	<table><tbody><tr><td>chr1</td><td>10</td><td>100</td><td>Query1..1</td><td></td><td></td><td></td></tr><tr><td>chr1</td><td>500</td><td>1000</td><td>Query1..2</td><td></td><td></td><td></td></tr><tr><td>chr1</td><td>1100</td><td>1250</td><td>Query1..3</td><td></td><td></td><td></td></tr></tbody></table>						chr1	10	100	Query1..1				chr1	500	1000	Query1..2				chr1	1100	1250	Query1..3			
chr1	10	100	Query1..1																								
chr1	500	1000	Query1..2																								
chr1	1100	1250	Query1..3																								
Query 2:	<table><tbody><tr><td></td><td></td><td></td><td>chr1</td><td>20</td><td>80</td><td>Query2..1</td></tr><tr><td></td><td></td><td></td><td>chr1</td><td>2000</td><td>2204</td><td>Query2..2</td></tr><tr><td></td><td></td><td></td><td>chr1</td><td>2500</td><td>3000</td><td>Query2..3</td></tr></tbody></table>									chr1	20	80	Query2..1				chr1	2000	2204	Query2..2				chr1	2500	3000	Query2..3
			chr1	20	80	Query2..1																					
			chr1	2000	2204	Query2..2																					
			chr1	2500	3000	Query2..3																					
(Return all records of first query)	<table><tbody><tr><td>chr1</td><td>10</td><td>100</td><td>Query1..1</td><td>chr1</td><td>20</td><td>80</td></tr><tr><td>chr1</td><td>500</td><td>1000</td><td>Query1..2</td><td>.</td><td>.</td><td>.</td></tr><tr><td>chr1</td><td>1100</td><td>1250</td><td>Query1..3</td><td>.</td><td>.</td><td>.</td></tr></tbody></table>						chr1	10	100	Query1..1	chr1	20	80	chr1	500	1000	Query1..2	.	.	.	chr1	1100	1250	Query1..3	.	.	.
chr1	10	100	Query1..1	chr1	20	80																					
chr1	500	1000	Query1..2	.	.	.																					
chr1	1100	1250	Query1..3	.	.	.																					
	<p>Return only records that are joined (INNER JOIN) Return all records of first query (fill null with ".") Return all records of second query (fill null with ".") Return all records of both queries (fill nulls with ".")</p>																										



逻辑运算 | join

Input																														
Query 1:	<table><tbody><tr><td>chr1</td><td>10</td><td>100</td><td>Query1.1</td><td></td><td></td><td></td></tr><tr><td>chr1</td><td>500</td><td>1000</td><td>Query1.2</td><td></td><td></td><td></td></tr><tr><td>chr1</td><td>1100</td><td>1250</td><td>Query1.3</td><td></td><td></td><td></td></tr></tbody></table>						chr1	10	100	Query1.1				chr1	500	1000	Query1.2				chr1	1100	1250	Query1.3						
chr1	10	100	Query1.1																											
chr1	500	1000	Query1.2																											
chr1	1100	1250	Query1.3																											
Query 2:	<table><tbody><tr><td></td><td></td><td></td><td>chr1</td><td>20</td><td>80</td><td>Query2.1</td></tr><tr><td></td><td></td><td></td><td>chr1</td><td>2000</td><td>2204</td><td>Query2.2</td></tr><tr><td></td><td></td><td></td><td>chr1</td><td>2500</td><td>3000</td><td>Query2.3</td></tr></tbody></table>									chr1	20	80	Query2.1				chr1	2000	2204	Query2.2				chr1	2500	3000	Query2.3			
			chr1	20	80	Query2.1																								
			chr1	2000	2204	Query2.2																								
			chr1	2500	3000	Query2.3																								
(Return all records of second query)	<table><tbody><tr><td>chr1</td><td>10</td><td>100</td><td>Query1.1</td><td>chr1</td><td>20</td><td>80</td><td>Query2.1</td></tr><tr><td>.</td><td>.</td><td>.</td><td>.</td><td>chr1</td><td>2000</td><td>2204</td><td>Query2.2</td></tr><tr><td>.</td><td>.</td><td>.</td><td>.</td><td>chr1</td><td>500</td><td>3000</td><td>Query2.3</td></tr></tbody></table>						chr1	10	100	Query1.1	chr1	20	80	Query2.1	chr1	2000	2204	Query2.2	chr1	500	3000	Query2.3
chr1	10	100	Query1.1	chr1	20	80	Query2.1																							
.	.	.	.	chr1	2000	2204	Query2.2																							
.	.	.	.	chr1	500	3000	Query2.3																							
	<p>Return only records that are joined (INNER JOIN) Return all records of first query (fill null with ".") Return all records of second query (fill null with ".") Return all records of both queries (fill nulls with ".")</p>																													



逻辑运算 | join

Input						
Query 1:						
chr1 10 100 Query1..1						
chr1 500 1000 Query1..2						
chr1 1100 1250 Query1..3						
Query 2:						
chr1 20 80 Query2..1						
chr1 2000 2204 Query2..2						
chr1 2500 3000 Query2..3						
(Return all records of both queries)						
chr1 10 100 Query1..1	chr1 20 80 Query2..1	chr1 . . .	chr1 2000 2200 Query2..2	chr1 2500 3000 Query2..3	Return only records that are joined (INNER JOIN) Return all records of first query (fill null with ".") Return all records of second query (fill null with ".") Return all records of both queries (fill nulls with ".")	
chr1 500 1000 Query1..2						
chr1 1100 1250 Query1..3						
.		
.		



- coverage** Finds the number of bases each interval in the first dataset covers of the second dataset.
- flank** Finds the upstream and/or downstream flanking region(s).
- closest** Find the closest, potentially non-overlapping upstream and/or downstream features.
- slop** Adjust the size of intervals.
- window** Find overlapping intervals within a window around an interval.



逻辑运算 | 实例

Dataset 1

chr1	10	49	Feature1.1
chr1	70	119	Feature1.2
chr1	170	209	Feature1.3
chr1	180	229	Feature1.4

Dataset 2

chr1	80	109	Feature2.1
chr1	150	199	Feature2.2
chr1	250	289	Feature2.3
chr1	270	309	Feature2.4



逻辑运算 | 实例

Dataset 1

chr1	10	49	Feature1.1
chr1	70	119	Feature1.2
chr1	170	209	Feature1.3
chr1	180	229	Feature1.4

Dataset 2

chr1	80	109	Feature2.1
chr1	150	199	Feature2.2
chr1	250	289	Feature2.3
chr1	270	309	Feature2.4

intersect

chr1	80	109	Feature3.1
chr1	170	199	Feature3.2
chr1	180	199	Feature3.3



逻辑运算 | 实例

Dataset 1

chr1	10	49	Feature1.1
chr1	70	119	Feature1.2
chr1	170	209	Feature1.3
chr1	180	229	Feature1.4

Dataset 2

chr1	80	109	Feature2.1
chr1	150	199	Feature2.2
chr1	250	289	Feature2.3
chr1	270	309	Feature2.4

subtract (1-2)

chr1	10	49	Feature4.1
chr1	70	80	Feature4.2
chr1	109	119	Feature4.3
chr1	199	209	Feature4.4
chr1	199	229	Feature4.5

subtract (2-1)

chr1	150	170	Feature5.1
chr1	250	289	Feature5.2
chr1	270	309	Feature5.3



逻辑运算 | 实例

Dataset 1

chr1	10	49	Feature1.1
chr1	70	119	Feature1.2
chr1	170	209	Feature1.3
chr1	180	229	Feature1.4

Dataset 2

chr1	80	109	Feature2.1
chr1	150	199	Feature2.2
chr1	250	289	Feature2.3
chr1	270	309	Feature2.4

join

chr1	70	119	Feature1.2
chr1	170	209	Feature1.3
chr1	180	229	Feature1.4

chr1	80	109	Feature2.1
chr1	150	199	Feature2.2
chr1	150	199	Feature2.2



实际问题

- ① Find genes that overlap LINEs.
- ② Remove introns from gene features. Exons will (should) be reported.
- ③ Merge overlapping repetitive elements into a single entry.
- ④ Report all intervals in the human genome that are not covered by repetitive elements.

解决策略

- ① intersect
- ② subtract
- ③ merge
- ④ complement

实际问题

- ① Find genes that overlap LINEs.
- ② Remove introns from gene features. Exons will (should) be reported.
- ③ Merge overlapping repetitive elements into a single entry.
- ④ Report all intervals in the human genome that are not covered by repetitive elements.

解决策略

- ① intersect
- ② subtract
- ③ merge
- ④ complement

- Galaxy 中的 “Operate on Genomic Intervals” 工具集
- bedtools: a powerful toolset for genome arithmetic
- BEDOPS: the fast, highly scalable and easily-parallelizable genome analysis toolkit



教学提纲

1

引言

2

基因组组装版本

3

基因组坐标系统

4

基因组注释常用格式

5

文本文件与文本编辑器

6

基因组坐标的逻辑运算

7

总结与答疑

8

引言

9

变异位点的注释

10

基因集富集分析

11

序列标识

12

box plot

13

解析图表

14

总结与答疑

15

引言

16

Galaxy 分析平台

17

Galaxy 使用演示

18

数据处理三段论

19

总结与答疑

20

复习思考题



知识点——基因组注释基础

- 基因组组装版本——对应关系
- 两种坐标系统——0-based 和 1-based
- 四种常用格式——FASTA, BED, GFF, VCF
- 坐标逻辑运算——常见模式及其适用范围
- 坐标转换、格式转换、逻辑运算的工具

技能——纯文本与文本编辑器

- 纯文本与格式化文本
- 不同操作系统中的换行符
- 文本编辑器——Notepad++, Vim, Emacs



教学提纲

- 1 引言
- 2 基因组组装版本
- 3 基因组坐标系统
- 4 基因组注释常用格式
- 5 文本文件与文本编辑器
- 6 基因组坐标的逻辑运算
- 7 总结与答疑
- 8 引言
- 9 变异位点的注释
- 10 基因集富集分析

- 11 序列标识
- 12 box plot
- 13 解析图表
- 14 总结与答疑
- 15 引言
- 16 Galaxy 分析平台
- 17 Galaxy 使用演示
- 18 数据处理三段论
- 19 总结与答疑
- 20 复习思考题



前期准备工作

- 组装版本
- 坐标系统
- 常用格式
- 逻辑运算

后续功能注释

- 变异位点的注释
- 基因集富集分析
- 制作序列标识
- ...



前期准备工作

- 组装版本
- 坐标系统
- 常用格式
- 逻辑运算

后续功能注释

- 变异位点的注释
- 基因集富集分析
- 制作序列标识
- ...



教学提纲

- 1 引言
- 2 基因组组装版本
- 3 基因组坐标系统
- 4 基因组注释常用格式
- 5 文本文件与文本编辑器
- 6 基因组坐标的逻辑运算
- 7 总结与答疑
- 8 引言
- 9 变异位点的注释
- 10 基因集富集分析

- 11 序列标识
- 12 box plot
- 13 解析图表
- 14 总结与答疑
- 15 引言
- 16 Galaxy 分析平台
- 17 Galaxy 使用演示
- 18 数据处理三段论
- 19 总结与答疑
- 20 复习思考题

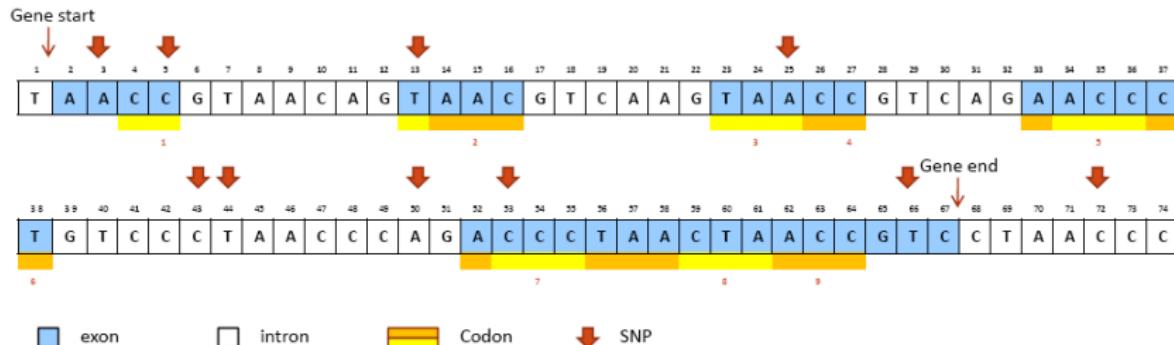


变异位点的注释 | SNP

ID	Chromosome	Position	Reference	Mutation	TotalHit
SNP00000001	15	20833654	C	A	38
SNP00000002	15	23501058	C	A	31
SNP00000003	15	45564496	C	A	20
SNP00000004	15	45564498	A	T	20
SNP00000005	15	45564501	G	T	20
SNP00000006	15	45564504	C	A	20
SNP00000007	15	45564505	C	T	20
SNP00000008	15	45564506	T	C	20
SNP00000009	15	45564508	A	T	20
SNP00000010	15	50212324	G	C	21
SNP00000011	15	50212325	A	T	21
SNP00000012	15	50212326	A	C	21
SNP00000013	15	50212328	G	A	21
SNP00000014	15	50212329	A	G	21
SNP00000015	15	50212330	G	A	21
SNP00000016	15	50212342	G	T	21
SNP00000017	15	52098626	A	G	20
SNP00000018	15	52098627	G	A	20



变异位点的注释 | SNP 注释



Pos	Alt SNP	Ref SNP	Alt SNP Codon	Ref SNP Codon	Alt SNP AA	Ref SNP AA	Anno Type
3	G	A	--	--	--	--	5'UTR
5	A	C	CAT	CCT	His	Pro	Non_Synonymous
13	G	T	CCG	CCT	Pro	Pro	Synonymous
25	C	A	TAC	TAA	Tyr	Stop	Stop Loss
43	A	C	--	--	--	--	Splice Site
44	G	T	--	--	--	--	Intronic
50	C	A	--	--	--	--	Essential Splice Site
53	A	C	ACC	CCC	Thr	Pro	Non_Synonymous
66	C	T	--	--	--	--	3'UTR
72	A	C	--	--	--	--	Downstream

- SNVs 的注释：SeattleSeq Annotation、variant tools、SnpEff
- 非同义多态性的功能注释：SIFT、PolyPhen-2、SNPs3D
- indels 的功能注释：PROVEAN



变异位点的注释 | 结果解析 | SeattleSeq Annotation

File:
/data/jboss-as-
7.1.1.Final/gvsBatchOutput/SeattleSeqAnnotation137.1.individual.294000040650.txt

Title:
1individual

Counts:
HapMapFreqType HapMapFreqMinor
polyPhenType polyPhenScore

Count missense SNPs = 8
Count stop SNPs = 0
Count SNP's in splice sites = 0
Count SNP's in coding synonymous = 8
Count SNP's in coding (not mod 3) = 0
Count SNP's in a UTR = 0
Count SNP's near a gene = 0
Count SNP's in introns = 0
Count intergenic SNPs = 0

number SNPs in microRNAs = 0

number accessions coding-synonymous NCBI = 19
number accessions missense NCBI = 15
number accessions stop NCBI = 0
number accessions splice-site NCBI = 0
number SNPs in dbSNP = 16
number SNPs not in dbSNP = 0
number SNPs total = 16

Add or Remove Columns:

- Sample Alleles
- Alleles in dbSNP
- GVS Function
- dbSNP Function
- Chimp Allele
- Copy Number Variations
- HapMap Rare-Allele Frequencies
- dbSNP Validation
- RepeatMasker
- Tandem Repeats
- microRNAs
- Grantham Score
- cDNA Position
- PolyPhen Prediction
- Clinical Association
- Distance to Nearest Splice Site
- NHLBI ESP Allele Counts

Sort by Column Value:

- Original Order
- dbSNP Function
- GVS Function
- Conservation Score phastCons
- Conservation Score GERP
- In dbSNP

Sort Direction:

- Forward
- Reverse

Filter:

- Only missense, nonsense, splice, frameshift (GVS)
- Only synonymous SNP's or coding (not frameshift) indels (GVS)
- Only intron (GVS)
- Only variations not in dbSNP
- Only variations with clinical association

Table 

reset 

16 SNP locations 36 accession lines page 1 of 1

inDBSNPOrNot	chromosome	position	referenceBase	sampleGenotype	sampleAlleles	allelesDBSNP	accession	functionGVS	functionDBSNP	rsID	aminoAcids	proteinPosition
dbSNP_130	10	1126383	A	R	A/G	A/G	NM_014023.3	coding-synonymous	synonymous-codon	73578536	none	121/495
dbSNP_86	10	3150973	C	Y	C/T	C/T	NM_001242339.1	coding-synonymous	synonymous-codon	1132173	none	309/777
dbSNP_86	10	3150973	C	Y	C/T	C/T	NM_002627.4	coding-synonymous	synonymous-codon	1132173	none	317/785



变异位点的注释 | 结果解析 | SeattleSeq Annotation

inDBSNPOrNot	chromosome	position	referenceBase	sampleGenotype	sampleAlleles	allelesDBSNP
dbSNP_130	10	1126383	A	R	A/G	A/G
dbSNP_86	10	3150973	C	Y	C/T	C/T
dbSNP_86	10	3150973	C	Y	C/T	C/T

accession	functionGVS	functionDBSNP	rsID	aminoAcids	proteinPosition
NM_014023.3	coding-synonymous	synonymous-codon	73578536	none	121/495
NM_001242339.1	coding-synonymous	synonymous-codon	1132173	none	309/777
NM_002627.4	coding-synonymous	synonymous-codon	1132173	none	317/785



变异位点的注释 | 结果解析 | SIFT

Transcript ID	Protein ID	Substitution	Region	dbSNP ID	SNP Type	Prediction	SIFT Score
ENST00000294724	ENSP00000294724	R1487G	EXON CDS	rs12118058:G	Nonsynonymous	TOLERATED	0.46
ENST00000294724	ENSP00000294724	E1405G	EXON CDS	rs28730708:G	Nonsynonymous	DAMAGING	0.01
ENST00000294724	ENSP00000294724	R1487R	EXON CDS	rs12118058:G	Synonymous	TOLERATED	0.64
ENST00000330029	ENSP00000332887	E49A	EXON CDS	novel	Nonsynonymous	DAMAGING	0.02
ENST00000371564	ENSP00000360619	T612N	EXON CDS	rs6067785:T	Nonsynonymous	DAMAGING	0
ENST00000283943	ENSP00000283943	Q1910*	EXON CDS	rs1803846:A	Nonsynonymous	N/A	N/A
ENST00000341772	ENSP00000345229	P433L	EXON CDS	rs17853365:A	Nonsynonymous	DAMAGING	0.02



教学提纲

- 1 引言
 - 2 基因组组装版本
 - 3 基因组坐标系统
 - 4 基因组注释常用格式
 - 5 文本文件与文本编辑器
 - 6 基因组坐标的逻辑运算
 - 7 总结与答疑
 - 8 引言
 - 9 变异位点的注释
 - 10 基因集富集分析



富集分析 | 基因集

Table 7 The minimum gene set selected in PRI dataset (gene scores rank from high to low)

Probe ID	Gene symbol	Gene name	Chromosomal regions
219868_s_at	ANKFY1	Ankyrin repeat and FYVE domain containing1	17p13.3
213613_s_at	NADK	NAD kinase	1p36.33-p36.21
208002_s_at	ACOT7	Acyl-coa thioesterase 7	1p36
222133_s_at	PHF20L1	PHD finger protein 20-like 1	8q24.22
203858_s_at	COX10	COX10 homolog, cytochrome c oxidase assembly protein, heme A: farnesyltransferase (yeast)	17p12
204051_s_at	SFRP4	Secreted frizzled-related protein 4	7p14.1
207567_at	SLC13A2	Solute carrier family 13 (sodium-dependent dicarboxylate transporter), member 2	17p13.2
225803_at	FBXO32	F-box protein 32	8q24.13
205527_s_at	GEMIN4	Gem (nuclear organelle) associated protein 4	17p13
207017_at	RAB27B	RAB27B, member RAS oncogene family	18q21.2
206746_at	BFSP1	Beaded filament structural protein 1, filensin	20p12.1
217099_s_at	GEMIN4	Gem (nuclear organelle) associated protein 4	17p13
233638_s_at	POMGNT1	Protein O-linked mannose beta1,2-N-acetylglucosaminyltransferase	1p34.1
217381_s_at	TRGV5	T cell receptor gamma variable 5	7p14



数据库

GO Gene Ontology

KEGG Kyoto Encyclopedia of Genes and Genomes

分析工具

DAVID Database for Annotation, Visualization and Integrated Discovery



数据库

GO Gene Ontology

KEGG Kyoto Encyclopedia of Genes and Genomes

分析工具

DAVID Database for Annotation, Visualization and Integrated Discovery



三个方面

- biological process, BP, 生物学过程
- molecular function, MF, 分子功能
- cellular component, CC, 细胞组份

两大关系

- is_a: for simple, hierarchical connections between terms
- part_of: for describing how the components of a living system fit together



三个方面

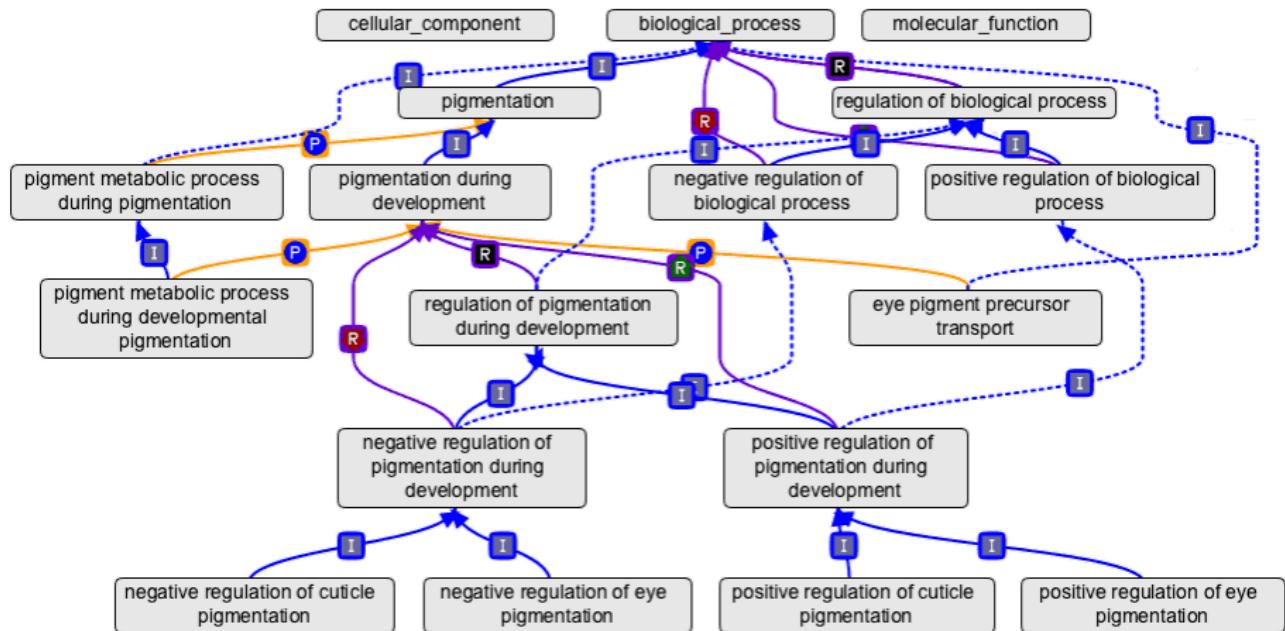
- biological process, BP, 生物学过程
- molecular function, MF, 分子功能
- cellular component, CC, 细胞组份

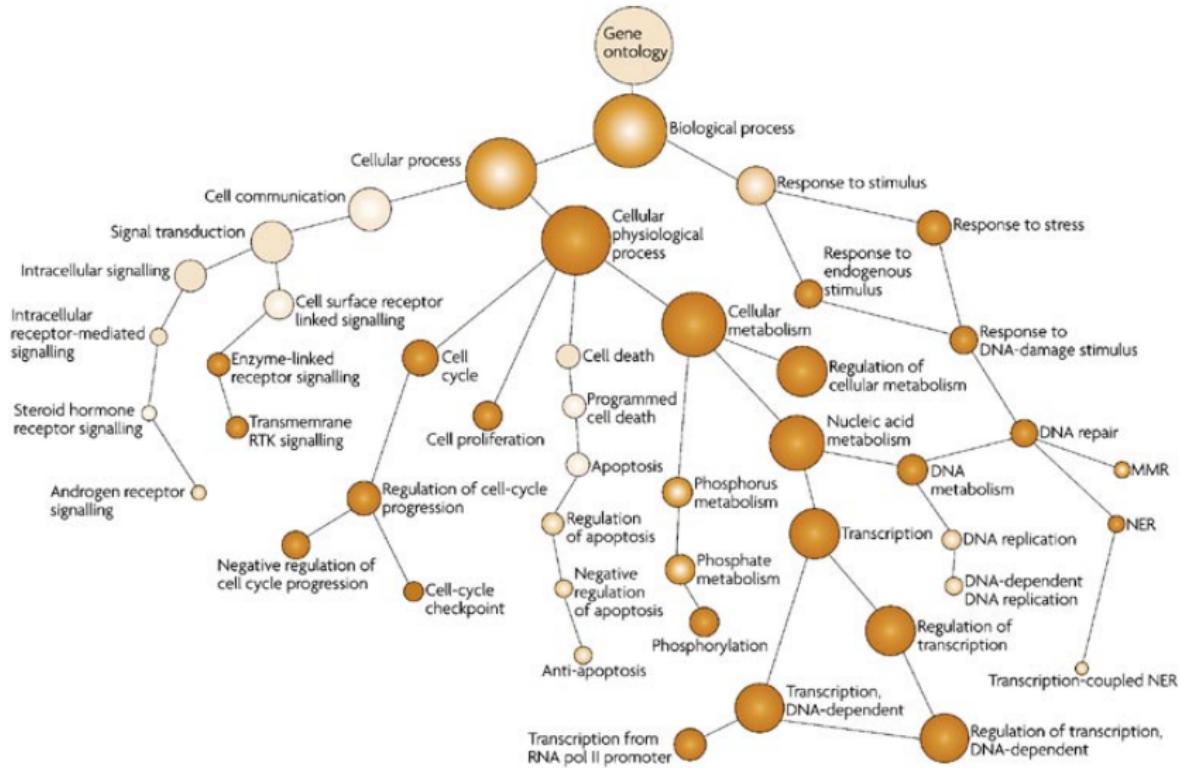
两大关系

- is_a: for simple, hierarchical connections between terms
- part_of: for describing how the components of a living system fit together



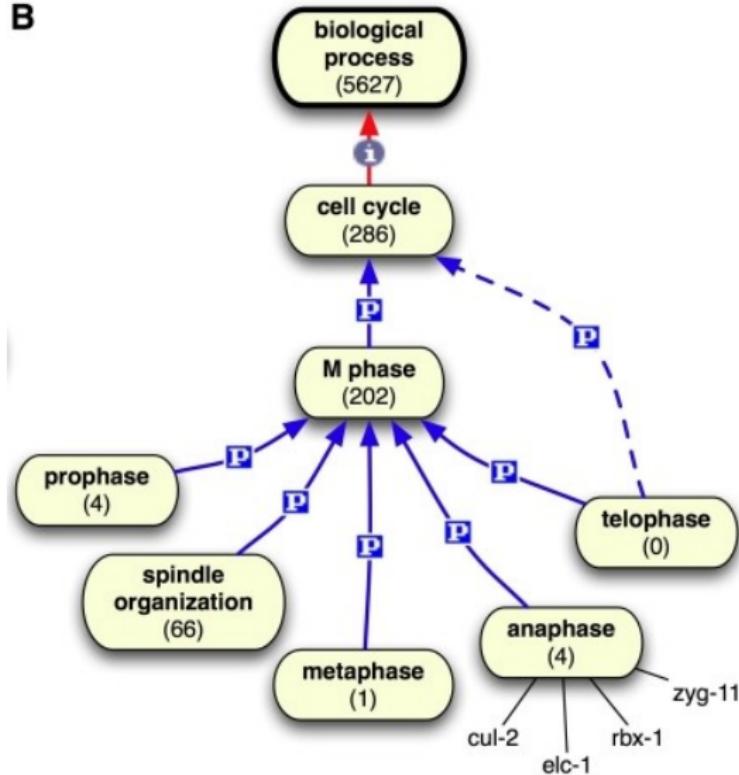
富集分析 | GO



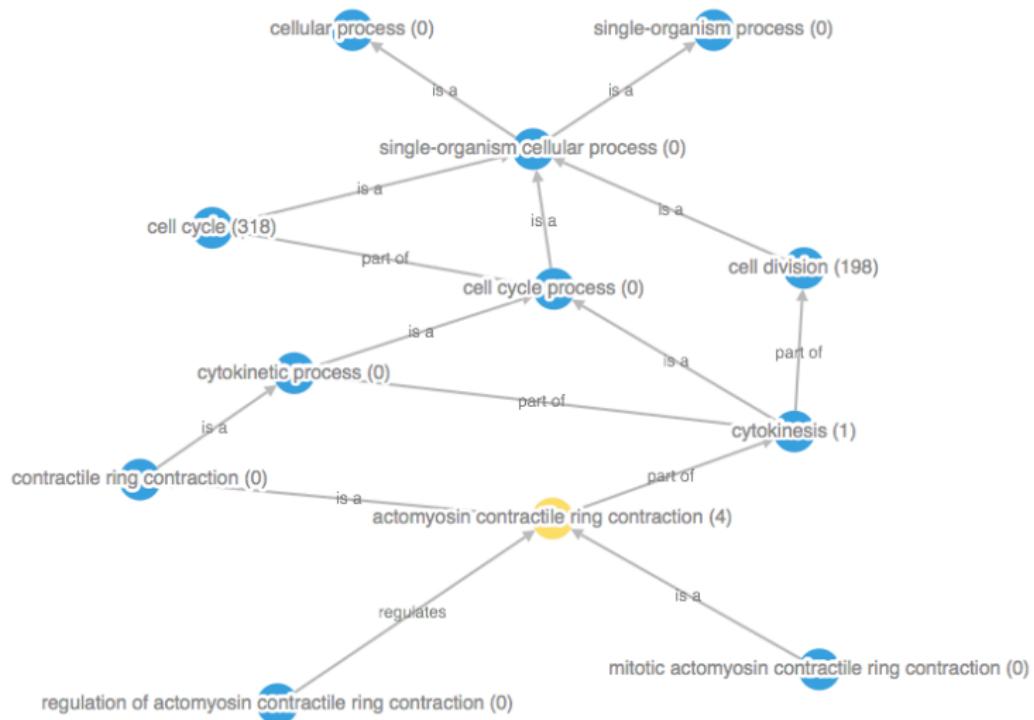


富集分析 | GO

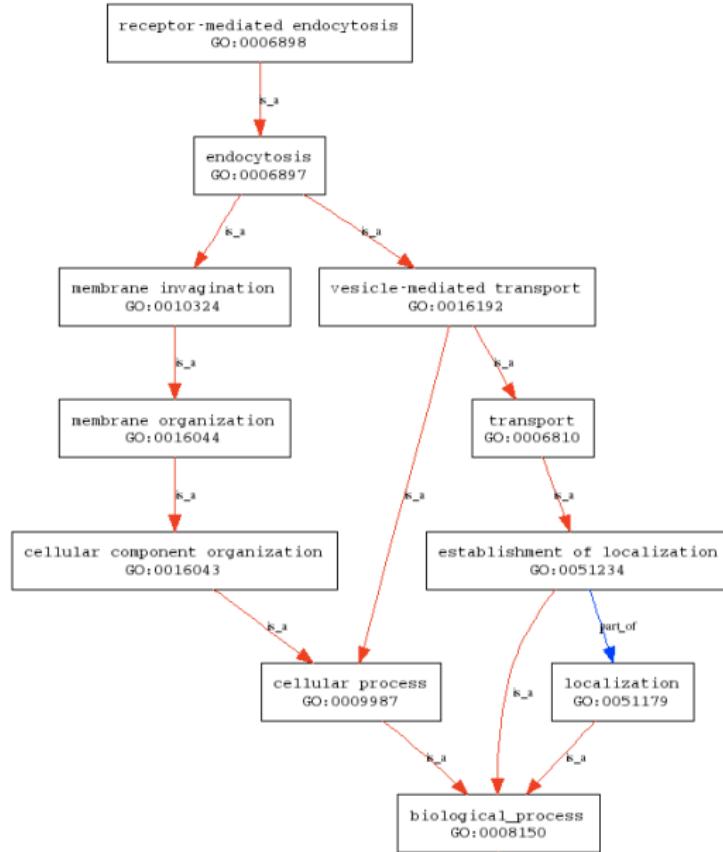
B



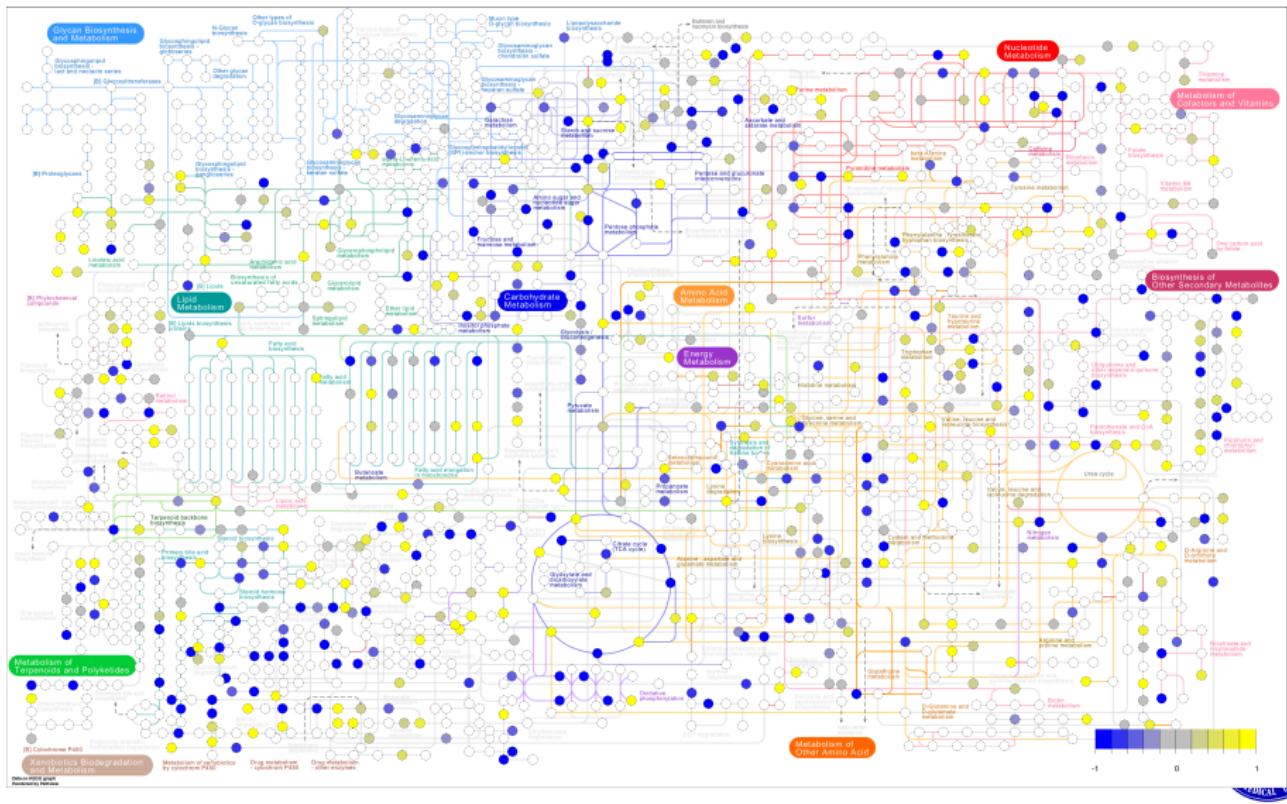
Relationships	Nodes
-P→ part_of	ontology term
-i→ is_a	root node
-d→ develops_from	
-→ inferred	
— annotation	gene name



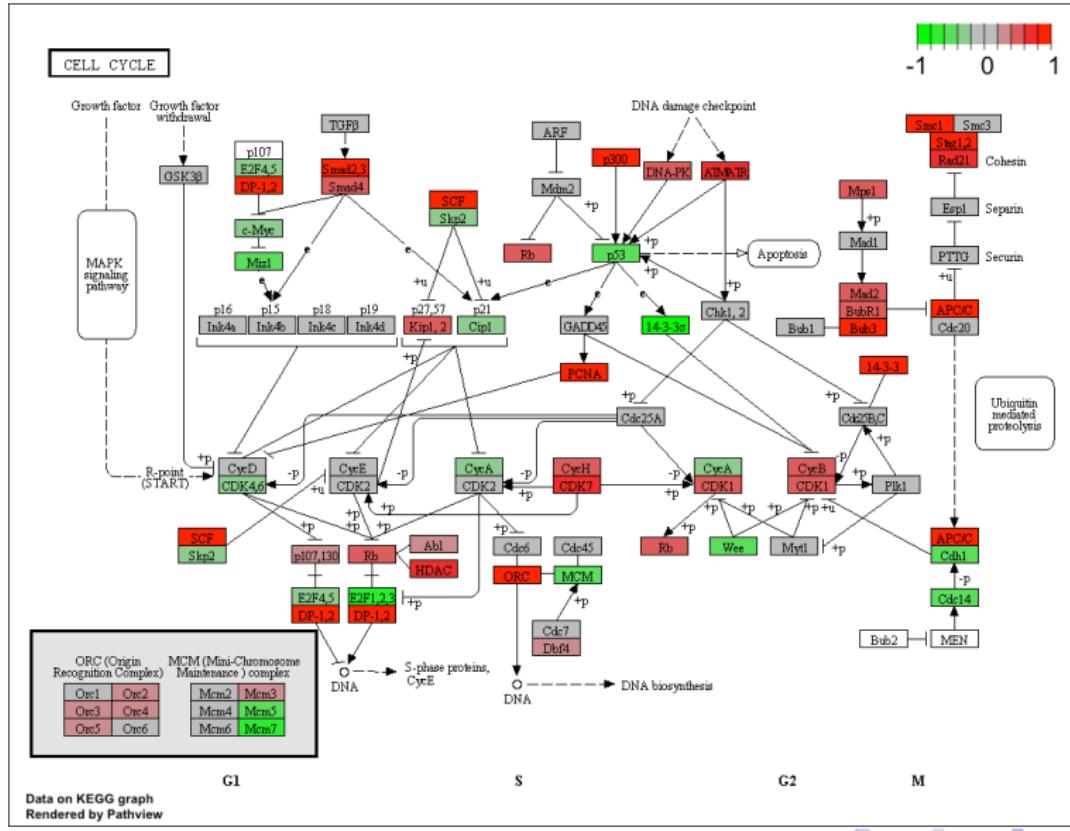
富集分析 | GO



富集分析 | KEGG



富集分析 | KEGG



- Gene Name Batch Viewer
 - Gene ID Conversion Tool
 - Gene Functional Classification Tool
 - Functional Annotation Tool
 - Functional Annotation Clustering
 - **Functional Annotation Chart:** 富集分析
 - Functional Annotation Table



富集分析 | DAVID | 结果解析

Functional Annotation Chart

[Help and Manual](#)**Current Gene List: demolist1****Current Background: Homo sapiens****155 DAVID IDs** Options[Rerun Using Options](#)[Create Sublist](#)**105 chart records** [Download File](#) [1, 2, 3, 4, 5, 6](#)

Category	Term	RT	Genes	Count	%	P-Value	Benjamini
GOTERM_CC_FAT	extracellular_region	RT		40	25.8	6.9E-6	1.5E-3
GOTERM_CC_FAT	extracellular_region_part	RT		24	15.5	3.8E-5	4.0E-3
GOTERM_MF_FAT	oxygen_binding	RT		6	3.9	3.8E-5	1.4E-2
GOTERM_CC_FAT	extracellular_space	RT		19	12.3	9.4E-5	6.5E-3
GOTERM_MF_FAT	heme_binding	RT		8	5.2	1.0E-4	1.9E-2
GOTERM_BP_FAT	defense_response	RT		18	11.6	1.3E-4	1.7E-1
GOTERM_BP_FAT	response_to_bacterium	RT		10	6.5	1.4E-4	9.1E-2
GOTERM_MF_FAT	tetrapyrrole_binding	RT		8	5.2	1.5E-4	1.9E-2
GOTERM_MF_FAT	iron_ion_binding	RT		11	7.1	4.3E-4	3.9E-2
GOTERM_BP_FAT	defense_response_to_bacterium	RT		7	4.5	8.9E-4	3.4E-1
GOTERM_BP_FAT	response_to_drug	RT		9	5.8	1.5E-3	4.0E-1
GOTERM_BP_FAT	regulation_of_response_to_external_stimulus	RT		7	4.5	5.2E-3	7.7E-1
GOTERM_BP_FAT	taxis	RT		7	4.5	5.4E-3	7.2E-1
GOTERM_BP_FAT	chemotaxis	RT		7	4.5	5.4E-3	7.2E-1
GOTERM_CC_FAT	hemoglobin_complex	RT		3	1.9	5.7E-3	2.6E-1
GOTERM_MF_FAT	oxygen_transporter_activity	RT		3	1.9	5.8E-3	3.5E-1

富集分析 | DAVID | 工具选择

- Highly recommended
- Recommended

	Gene ID conversion tool	Gene name batch viewer	Gene functional classification	Functional annotation chart	Functional annotation clustering	Functional annotation table
Convert gene IDs from one type to another	■					
Diagnose and fix problems of gene IDs		■				■
Explore gene names in batch		■	■			■
Discover enriched functionally related gene groups			■	■		
Display relationship of many-genes-to-many-terms on 2D view.				■	■	■
Initial glance of major biological functions associated with gene list	■		■	■		
Identify enriched (overrepresented) annotation terms				■	■	
Visualize genes on BioCarta and KEGG pathway maps				■		
Link gene–disease associations				■		
Highlight protein functional domains and motifs			■			
Redirect to related literatures				■		■
List interacting proteins				■	■	■
Cluster redundant and heterozygous annotation terms						
Search other functionally similar genes in genome, but not in list	■	■		1	1	
Search other annotations functionally similar to one of my interests			■			
Read all annotation contents associated with a gene						■



教学提纲

1

引言

2

基因组组装版本

3

基因组坐标系统

4

基因组注释常用格式

5

文本文件与文本编辑器

6

基因组坐标的逻辑运算

7

总结与答疑

8

引言

9

变异位点的注释

10

基因集富集分析

11

序列标识

12

box plot

13

解析图表

14

总结与答疑

15

引言

16

Galaxy 分析平台

17

Galaxy 使用演示

18

数据处理三段论

19

总结与答疑

20

复习思考题



序列标识 | 徽标



序列标识

定义

序列标识图是显示序列保守区域的共识序列、每个位置上各个氨基酸或核苷酸出现的频率以及各个位点上的序列信息量的一种可视化方法。

含义

根据序列保守区域的多序列比对来绘制序列标识图。

在一个标识图像里，由大小不一的字符形成的一个堆栈代表序列保守区域的一个位点。每个核苷酸或氨基酸的高度和它在对应位点上出现的频率成比例。堆栈的总高度代表对应位点上的序列信息，以比特（bit）为单位。在每个堆栈里，字符按其出现的频率大小自上而下排列。所以，位于各个堆栈最上方的字符组成保守区域的共识序列。



序列标识

定义

序列标识图是显示序列保守区域的共识序列、每个位置上各个氨基酸或核苷酸出现的频率以及各个位点上的序列信息量的一种可视化方法。

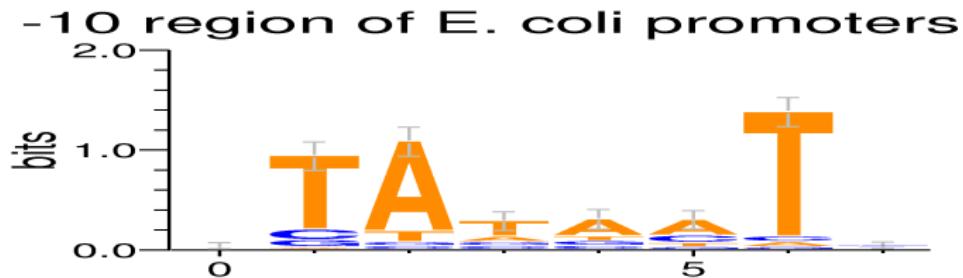
含义

根据序列保守区域的多序列比对来绘制序列标识图。

在一个标识图像里，由大小不一的字符形成的一个堆栈代表序列保守区域的一个位点。每个核苷酸或氨基酸的高度和它在对应位点上出现的频率成比例。堆栈的总高度代表对应位点上的序列信息，以比特（bit）为单位。在每个堆栈里，字符按其出现的频率大小自上而下排列。所以，位于各个堆栈最上方的字符组成保守区域的共识序列。



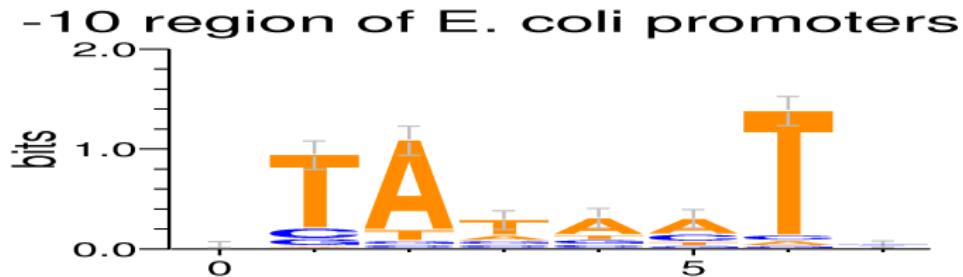
序列标识



序列标识 (sequence logo)

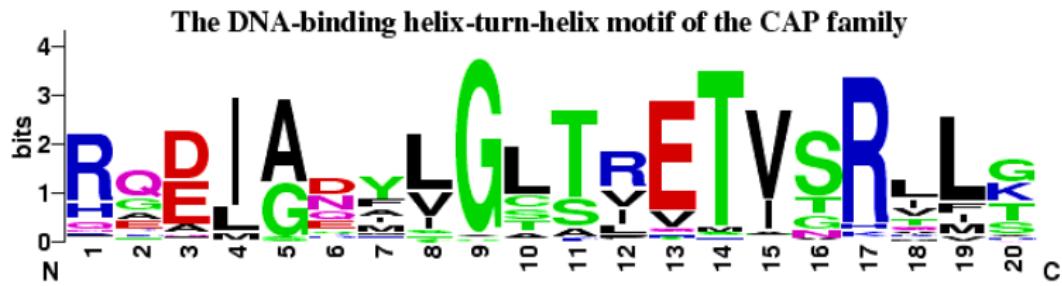
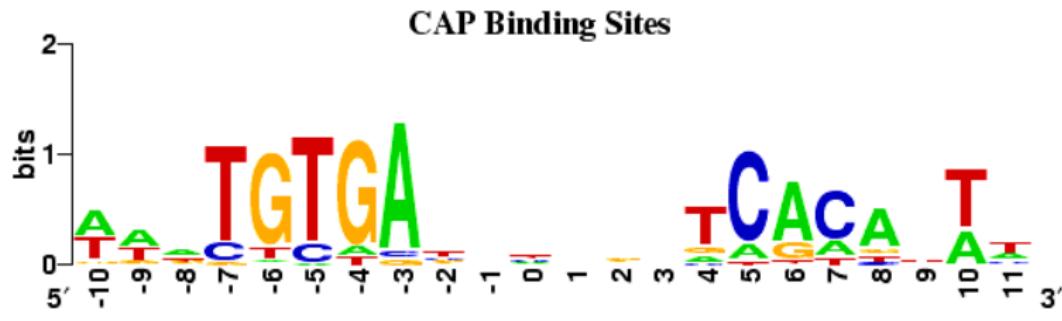
- 数据：多序列比对信息
- 横轴：序列坐标位置
- 纵轴：比特，计量单位
- 总高度：信息量/保守性
- 相对高度：相对频率
- 位置自上而下：频率由大到小
- 制作工具：WebLogo, enoLOGOS, Skylign

序列标识



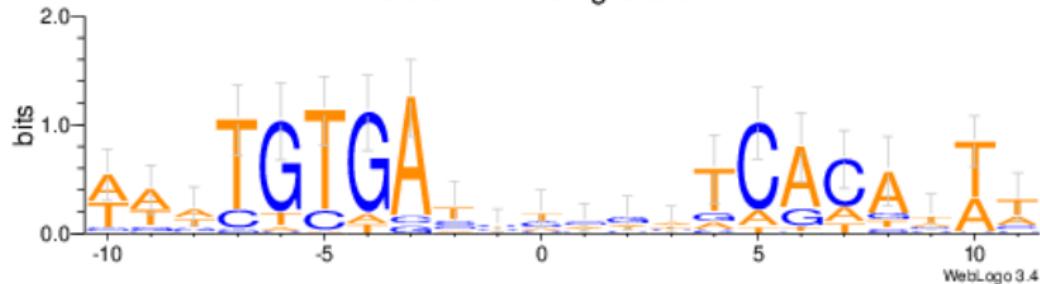
序列标识 (sequence logo)

- 数据：多序列比对信息
- 横轴：序列坐标位置
- 纵轴：比特，计量单位
- 总高度：信息量/保守性
- 相对高度：相对频率
- 位置自上而下：频率由大到小
- 制作工具：WebLogo, enoLOGOS, Skylign

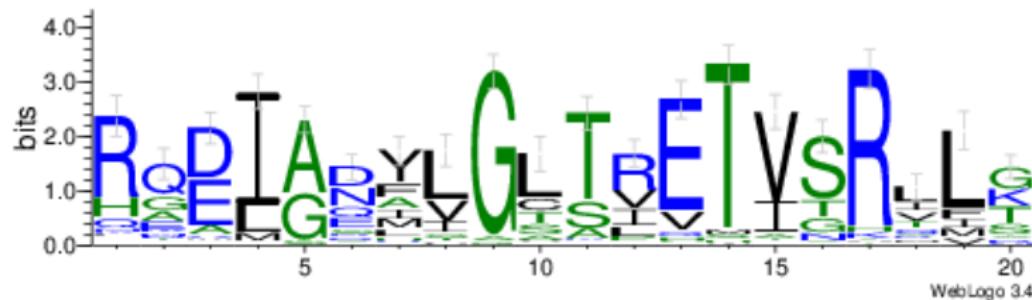


序列标识 | 着色 | WebLogo3

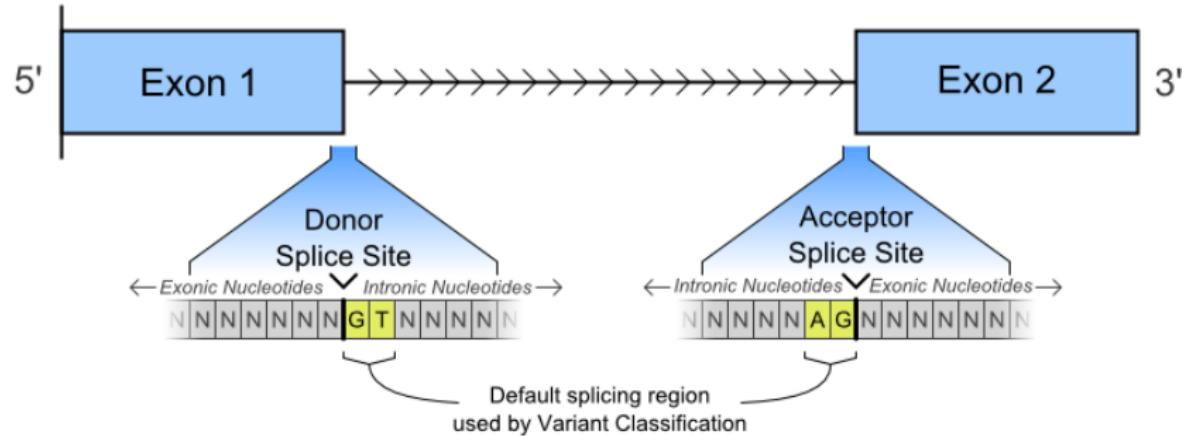
58 CAP Binding Sites



The DNA-binding helix-turn-helix motif of the CAP family

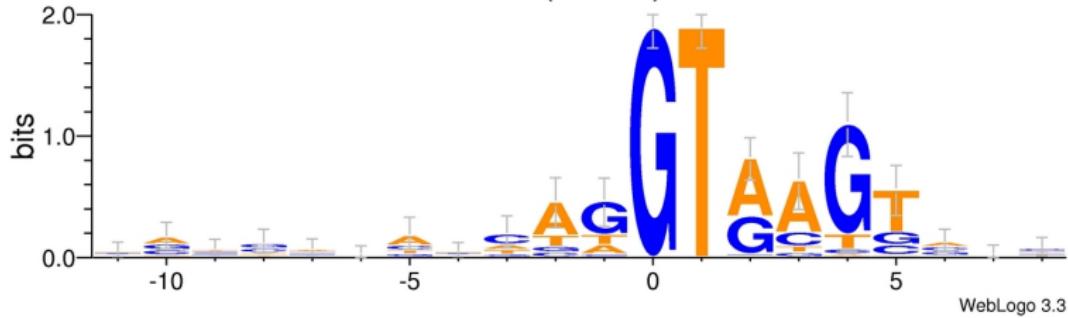


序列标识 | 剪接

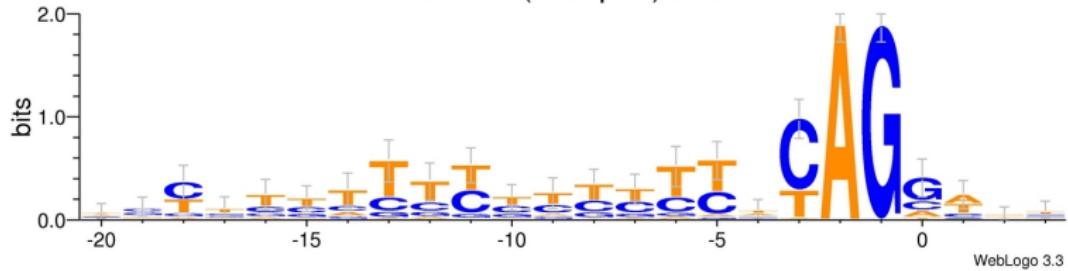


序列标识 | 实例

Exon-Intron (Donor) Sites



Intron-Exon (Acceptor) Sites



序列标识 | 实例 | 真实数据

Donor Sites(%)		Acceptor Site(%)	
GT	98.797	AG	99.714
GC	0.920	AC	0.120
AT	0.143	TG	0.032
GA	0.028	AT	0.024
GG	0.025	GG	0.022
CT	0.018	AA	0.019
TT	0.016	CG	0.010
CC	0.011	CC	0.010
TG	0.007	TT	0.009
AG	0.007	CT	0.008
TA	0.006	CA	0.008
AC	0.006	GC	0.007
CA	0.006	TA	0.006
TC	0.004	TC	0.004
AA	0.004	GT	0.004
CG	0.002	GA	0.003

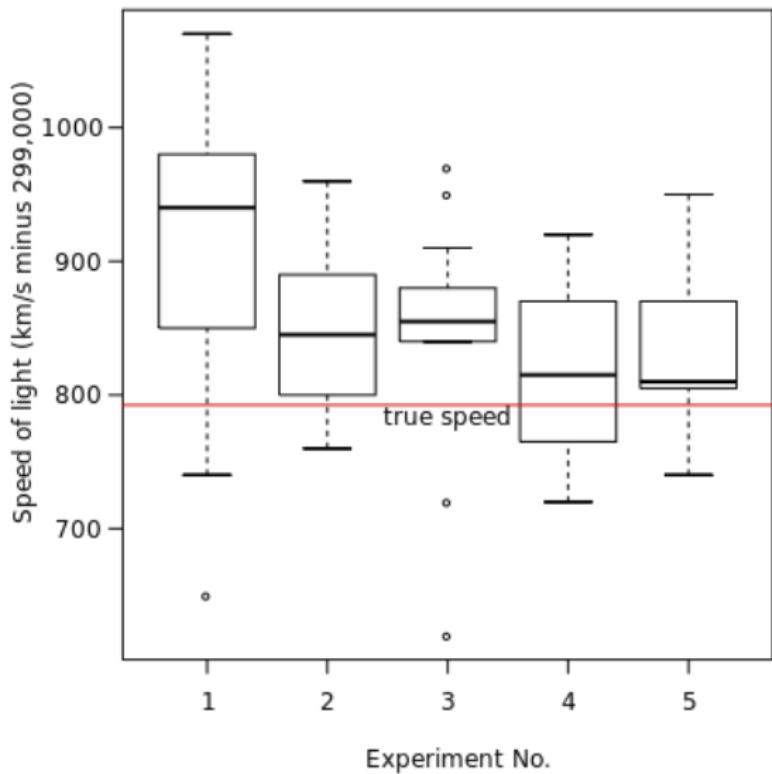
教学提纲

- 1 引言
- 2 基因组组装版本
- 3 基因组坐标系统
- 4 基因组注释常用格式
- 5 文本文件与文本编辑器
- 6 基因组坐标的逻辑运算
- 7 总结与答疑
- 8 引言
- 9 变异位点的注释
- 10 基因集富集分析

- 11 序列标识
- 12 box plot
- 13 解析图表
- 14 总结与答疑
- 15 引言
- 16 Galaxy 分析平台
- 17 Galaxy 使用演示
- 18 数据处理三段论
- 19 总结与答疑
- 20 复习思考题



box plot | 实例



历史

- box plot, boxplot, Box-whisker Plot
- 箱线图、箱须图、盒须图、盒式图、盒状图，因形状如箱子而得名
- 1977 年由美国著名统计学家约翰·图基（John Tukey）发明

简介

- 显示一组数据分散情况的统计图
- 显示最大值、最小值、中位数、下四分位数和上四分位数

优缺点

- 可以粗略地看出数据是否具有对称性、分布的离散程度
- 适合用于几个样本的比较
- 不能提供关于数据分布偏态和尾重程度的精确度量

box plot | 简介

历史

- box plot, boxplot, Box-whisker Plot
- 箱线图、箱须图、盒须图、盒式图、盒状图，因形状如箱子而得名
- 1977 年由美国著名统计学家约翰·图基（John Tukey）发明

简介

- 显示一组数据分散情况的统计图
- 显示最大值、最小值、中位数、下四分位数和上四分位数

优缺点

- 可以粗略地看出数据是否具有对称性、分布的离散程度
- 适合用于几个样本的比较
- 不能提供关于数据分布偏态和尾重程度的精确度量

历史

- box plot, boxplot, Box-whisker Plot
- 箱线图、箱须图、盒须图、盒式图、盒状图，因形状如箱子而得名
- 1977 年由美国著名统计学家约翰·图基（John Tukey）发明

简介

- 显示一组数据分散情况的统计图
- 显示最大值、最小值、中位数、下四分位数和上四分位数

优缺点

- 可以粗略地看出数据是否具有对称性、分布的离散程度
- 适合用于几个样本的比较
- 不能提供关于数据分布偏态和尾重程度的精确度量

- 最小值 min, 最大值 max
- 中位数 median
- 下四分位数 Q1, 上四分位数 Q3
- 四分位数差 IQR (interquartile range) , $IQR = Q3 - Q1$
- 内限: $Q3 + 1.5IQR, Q1 - 1.5IQR$
- 外限: $Q3 + 3IQR, Q1 - 3IQR$
- 异常值 (outliers) : 处于内限以外的数据
- 温和的异常值 (mild outliers) : 在内限与外限之间的异常值
- 极端的异常值 (extreme outliers) : 在外限以外的异常值



- 最小值 min, 最大值 max
- 中位数 median
- 下四分位数 Q1, 上四分位数 Q3
- 四分位数差 IQR (interquartile range) , $IQR = Q3 - Q1$
- 内限: $Q3 + 1.5IQR, Q1 - 1.5IQR$
- 外限: $Q3 + 3IQR, Q1 - 3IQR$
- 异常值 (outliers) : 处于内限以外的数据
- 温和的异常值 (mild outliers) : 在内限与外限之间的异常值
- 极端的异常值 (extreme outliers) : 在外限以外的异常值



- 最小值 min, 最大值 max
- 中位数 median
- 下四分位数 Q1, 上四分位数 Q3
- 四分位数差 IQR (interquartile range) , $IQR = Q3 - Q1$
- 内限: $Q3 + 1.5IQR, Q1 - 1.5IQR$
- 外限: $Q3 + 3IQR, Q1 - 3IQR$
- 异常值 (outliers) : 处于内限以外的数据
- 温和的异常值 (mild outliers) : 在内限与外限之间的异常值
- 极端的异常值 (extreme outliers) : 在外限以外的异常值



- 最小值 min, 最大值 max
- 中位数 median
- 下四分位数 Q1, 上四分位数 Q3
- 四分位数差 IQR (interquartile range) , $IQR = Q3 - Q1$
- 内限: $Q3 + 1.5IQR, Q1 - 1.5IQR$
- 外限: $Q3 + 3IQR, Q1 - 3IQR$
- 异常值 (outliers) : 处于内限以外的数据
- 温和的异常值 (mild outliers) : 在内限与外限之间的异常值
- 极端的异常值 (extreme outliers) : 在外限以外的异常值



- 最小值 min, 最大值 max
- 中位数 median
- 下四分位数 Q1, 上四分位数 Q3
- 四分位数差 IQR (interquartile range) , $IQR = Q3 - Q1$
- 内限: $Q3 + 1.5IQR, Q1 - 1.5IQR$
- 外限: $Q3 + 3IQR, Q1 - 3IQR$
- 异常值 (outliers) : 处于内限以外的数据
- 温和的异常值 (mild outliers) : 在内限与外限之间的异常值
- 极端的异常值 (extreme outliers) : 在外限以外的异常值



- 最小值 min, 最大值 max
- 中位数 median
- 下四分位数 Q1, 上四分位数 Q3
- 四分位数差 IQR (interquartile range) , $IQR = Q3 - Q1$
- 内限: $Q3 + 1.5IQR$, $Q1 - 1.5IQR$
- 外限: $Q3 + 3IQR$, $Q1 - 3IQR$
- 异常值 (outliers) : 处于内限以外的数据
- 温和的异常值 (mild outliers) : 在内限与外限之间的异常值
- 极端的异常值 (extreme outliers) : 在外限以外的异常值



- 最小值 min, 最大值 max
- 中位数 median
- 下四分位数 Q1, 上四分位数 Q3
- 四分位数差 IQR (interquartile range) , $IQR = Q3 - Q1$
- 内限: $Q3 + 1.5/IQR$, $Q1 - 1.5/IQR$
- 外限: $Q3 + 3/IQR$, $Q1 - 3/IQR$
- 异常值 (outliers) : 处于内限以外的数据
 - 温和的异常值 (mild outliers) : 在内限与外限之间的异常值
 - 极端的异常值 (extreme outliers) : 在外限以外的异常值



- 最小值 min, 最大值 max
- 中位数 median
- 下四分位数 Q1, 上四分位数 Q3
- 四分位数差 IQR (interquartile range) , $IQR = Q3 - Q1$
- 内限: $Q3 + 1.5/IQR$, $Q1 - 1.5/IQR$
- 外限: $Q3 + 3/IQR$, $Q1 - 3/IQR$
- 异常值 (outliers) : 处于内限以外的数据
- 温和的异常值 (mild outliers) : 在内限与外限之间的异常值
- 极端的异常值 (extreme outliers) : 在外限以外的异常值



- 最小值 min, 最大值 max
- 中位数 median
- 下四分位数 Q1, 上四分位数 Q3
- 四分位数差 IQR (interquartile range) , $IQR = Q3 - Q1$
- 内限: $Q3 + 1.5/IQR$, $Q1 - 1.5/IQR$
- 外限: $Q3 + 3/IQR$, $Q1 - 3/IQR$
- 异常值 (outliers) : 处于内限以外的数据
- 温和的异常值 (mild outliers) : 在内限与外限之间的异常值
- 极端的异常值 (extreme outliers) : 在外限以外的异常值



box plot | 绘图步骤

● 绘制数轴。

- 计算上四分位数 (Q3) , 中位数, 下四分位数 (Q1) 。
- 计算四分位数差 (IQR) 。
- 绘制箱线图的矩形, 上限为 Q3, 下限为 Q1。在矩形内部中位数的位置画一条横线 (中位线) 。
- 在 $Q3 + 1.5 \times IQR$ 和 $Q1 - 1.5 \times IQR$ 处画两条与中位线一样的线段, 这两条线段为异常值截断点, 称为内限; 在 $Q3 + 3 \times IQR$ 和 $Q1 - 3 \times IQR$ 处画两条线段, 称为外限。(注意: 统计软件绘制的箱线图一般都没有标出内限和外限。)
- 在非异常值的数据中, 最靠近上边缘和下边缘 (即内限) 的两个数值处画横线, 作为箱线图的触须。
- 从矩形的两端向外各画一条线段直到不是异常值的最远点 (即上一步的触须), 表示该批数据正常值的分布区间。
- 温和的异常值用空心圆表示; 极端的异常值用实心点 (一说用 *) 表示。



- 绘制数轴。
- 计算上四分位数 (Q3) , 中位数, 下四分位数 (Q1) 。
- 计算四分位数差 (IQR) 。
- 绘制箱线图的矩形, 上限为 Q3, 下限为 Q1。在矩形内部中位数的位置画一条横线 (中位线) 。
- 在 $Q3 + 1.5/IQR$ 和 $Q1 - 1.5/IQR$ 处画两条与中位线一样的线段, 这两条线段为异常值截断点, 称为内限; 在 $Q3 + 3/IQR$ 和 $Q1 - 3/IQR$ 处画两条线段, 称为外限。 (注意: 统计软件绘制的箱线图一般都没有标出内限和外限。)
- 在非异常值的数据中, 最靠近上边缘和下边缘 (即内限) 的两个数值处画横线, 作为箱线图的触须。
- 从矩形的两端向外各画一条线段直到不是异常值的最远点 (即上一步的触须), 表示该批数据正常值的分布区间。
- 温和的异常值用空心圆表示; 极端的异常值用实心点 (一说用 *) 表示。



- 绘制数轴。
- 计算上四分位数 (Q3) , 中位数, 下四分位数 (Q1) 。
- **计算四分位数差 (IQR) 。**
- 绘制箱线图的矩形, 上限为 Q3, 下限为 Q1。在矩形内部中位数的位置画一条横线 (中位线) 。
- 在 $Q3 + 1.5/IQR$ 和 $Q1 - 1.5/IQR$ 处画两条与中位线一样的线段, 这两条线段为异常值截断点, 称为内限; 在 $Q3 + 3/IQR$ 和 $Q1 - 3/IQR$ 处画两条线段, 称为外限。**(注意: 统计软件绘制的箱线图一般都没有标出内限和外限。)**
- 在非异常值的数据中, 最靠近上边缘和下边缘 (即内限) 的两个数值处画横线, 作为箱线图的触须。
- 从矩形的两端向外各画一条线段直到不是异常值的最远点 (即上一步的触须), 表示该批数据正常值的分布区间。
- 温和的异常值用空心圆表示; 极端的异常值用实心点 (一说用星号 *) 表示。



box plot | 绘图步骤

- 绘制数轴。
- 计算上四分位数 (Q3) , 中位数, 下四分位数 (Q1) 。
- 计算四分位数差 (IQR) 。
- 绘制箱线图的矩形, 上限为 Q3, 下限为 Q1。在矩形内部中位数的位置画一条横线 (中位线) 。
- 在 $Q3 + 1.5/IQR$ 和 $Q1 - 1.5/IQR$ 处画两条与中位线一样的线段, 这两条线段为异常值截断点, 称为内限; 在 $Q3 + 3/IQR$ 和 $Q1 - 3/IQR$ 处画两条线段, 称为外限。**(注意: 统计软件绘制的箱线图一般都没有标出内限和外限。)**
- 在非异常值的数据中, 最靠近上边缘和下边缘 (即内限) 的两个数值处画横线, 作为箱线图的触须。
- 从矩形的两端向外各画一条线段直到不是异常值的最远点 (即上一步的触须), 表示该批数据正常值的分布区间。
- 温和的异常值用空心圆表示; 极端的异常值用实心点 (一说用星号 *) 表示。



- 绘制数轴。
- 计算上四分位数 (Q3) , 中位数, 下四分位数 (Q1) 。
- 计算四分位数差 (IQR) 。
- 绘制箱线图的矩形, 上限为 Q3, 下限为 Q1。在矩形内部中位数的位置画一条横线 (中位线) 。
- 在 $Q3 + 1.5/IQR$ 和 $Q1 - 1.5/IQR$ 处画两条与中位线一样的线段, 这两条线段为异常值截断点, 称为内限; 在 $Q3 + 3/IQR$ 和 $Q1 - 3/IQR$ 处画两条线段, 称为外限。 (注意: 统计软件绘制的箱线图一般都没有标出内限和外限。)
- 在非异常值的数据中, 最靠近上边缘和下边缘 (即内限) 的两个数值处画横线, 作为箱线图的触须。
- 从矩形的两端向外各画一条线段直到不是异常值的最远点 (即上一步的触须), 表示该批数据正常值的分布区间。
- 温和的异常值用空心圆表示; 极端的异常值用实心点 (一说用星 *) 表示。



- 绘制数轴。
- 计算上四分位数 (Q3) , 中位数, 下四分位数 (Q1) 。
- 计算四分位数差 (IQR) 。
- 绘制箱线图的矩形, 上限为 Q3, 下限为 Q1。在矩形内部中位数的位置画一条横线 (中位线) 。
- 在 $Q3 + 1.5/IQR$ 和 $Q1 - 1.5/IQR$ 处画两条与中位线一样的线段, 这两条线段为异常值截断点, 称为内限; 在 $Q3 + 3/IQR$ 和 $Q1 - 3/IQR$ 处画两条线段, 称为外限。**(注意: 统计软件绘制的箱线图一般都没有标出内限和外限。)**
- 在非异常值的数据中, 最靠近上边缘和下边缘 (即内限) 的两个数值处画横线, 作为箱线图的触须。
- 从矩形的两端向外各画一条线段直到不是异常值的最远点 (即上一步的触须), 表示该批数据正常值的分布区间。
- 温和的异常值用空心圆表示; 极端的异常值用实心点 (一说用星 *) 表示。



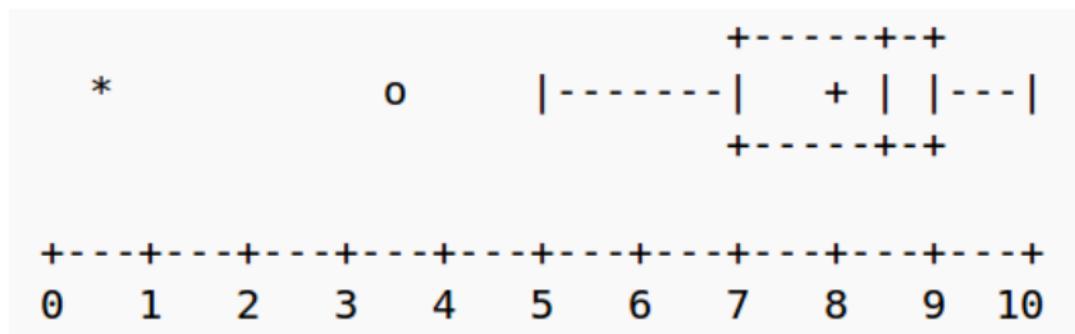
- 绘制数轴。
- 计算上四分位数 (Q3) , 中位数, 下四分位数 (Q1) 。
- 计算四分位数差 (IQR) 。
- 绘制箱线图的矩形, 上限为 Q3, 下限为 Q1。在矩形内部中位数的位置画一条横线 (中位线) 。
- 在 $Q3 + 1.5/IQR$ 和 $Q1 - 1.5/IQR$ 处画两条与中位线一样的线段, 这两条线段为异常值截断点, 称为内限; 在 $Q3 + 3/IQR$ 和 $Q1 - 3/IQR$ 处画两条线段, 称为外限。**(注意: 统计软件绘制的箱线图一般都没有标出内限和外限。)**
- 在非异常值的数据中, 最靠近上边缘和下边缘 (即内限) 的两个数值处画横线, 作为箱线图的触须。
- 从矩形的两端向外各画一条线段直到不是异常值的最远点 (即上一步的触须), 表示该批数据正常值的分布区间。
- 温和的异常值用空心圆表示; 极端的异常值用实心点 (一说用星*) 表示。



- 绘制数轴。
- 计算上四分位数 (Q3) , 中位数, 下四分位数 (Q1) 。
- 计算四分位数差 (IQR) 。
- 绘制箱线图的矩形, 上限为 Q3, 下限为 Q1。在矩形内部中位数的位置画一条横线 (中位线) 。
- 在 $Q3 + 1.5/IQR$ 和 $Q1 - 1.5/IQR$ 处画两条与中位线一样的线段, 这两条线段为异常值截断点, 称为内限; 在 $Q3 + 3/IQR$ 和 $Q1 - 3/IQR$ 处画两条线段, 称为外限。**(注意: 统计软件绘制的箱线图一般都没有标出内限和外限。)**
- 在非异常值的数据中, 最靠近上边缘和下边缘 (即内限) 的两个数值处画横线, 作为箱线图的触须。
- 从矩形的两端向外各画一条线段直到不是异常值的最远点 (即上一步的触须), 表示该批数据正常值的分布区间。
- 温和的异常值用空心圆表示; 极端的异常值用实心点 (一说用星号*) 表示。



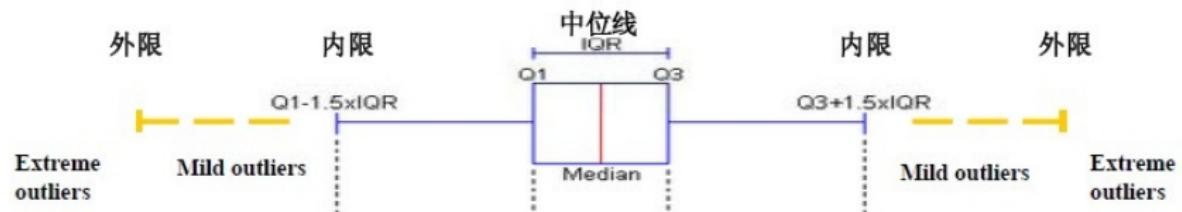
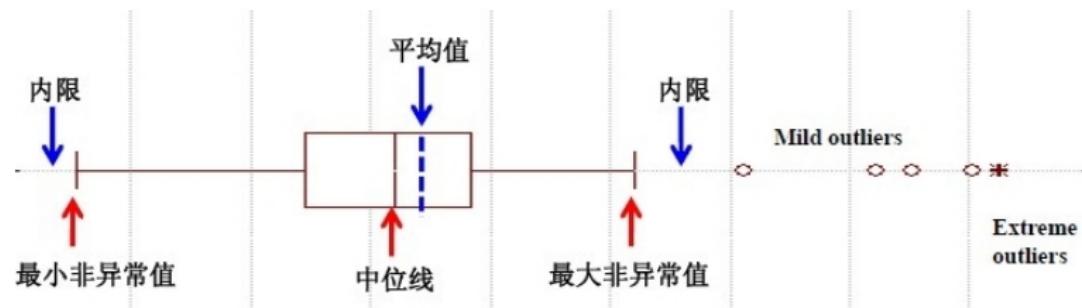
box plot | 图解



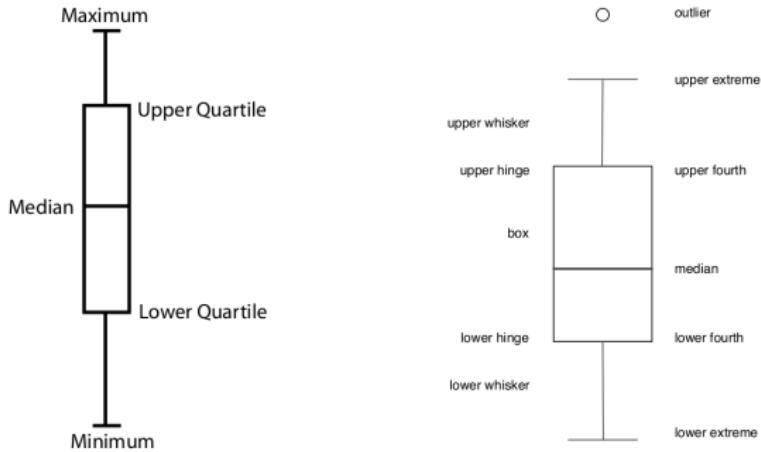
最小值 (min) =0.5; 下四分位数 (Q1) =7; 中位数 (Med) =8.5;
上四分位数 (Q3) =9; 最大值 (max) =10; 平均值 =8;
四分位数差 (interquartile range, 四分位间距) =Q3-Q1=2。



box plot | 图解



box plot | 图解



Sample, $n = 20$ ○ ●○○ ●●○ ○○○○ ○ ○○○

$1.5 \times IQF$

IQR

$$1.5 \times IQR$$

Q

四

Q

whiskers

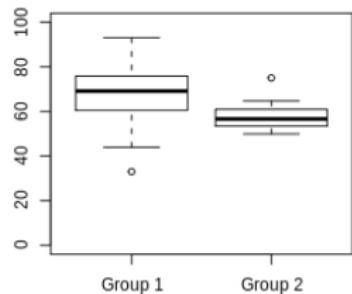
Outliers

6

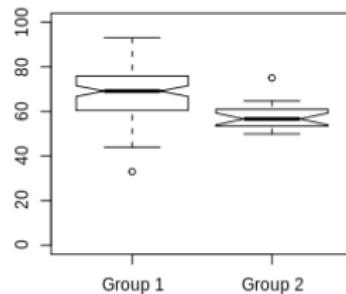


box plot | 变体

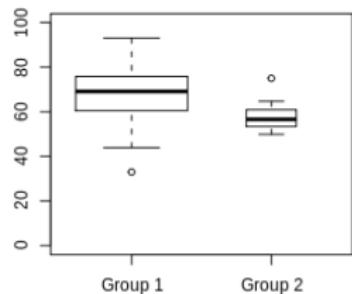
Traditional Box Plot



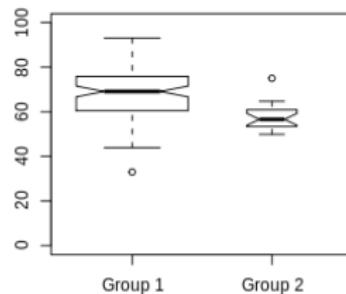
Notched Box Plot



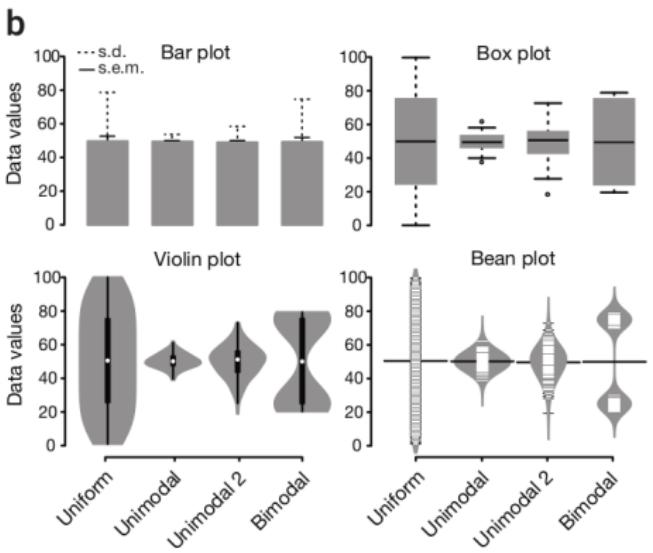
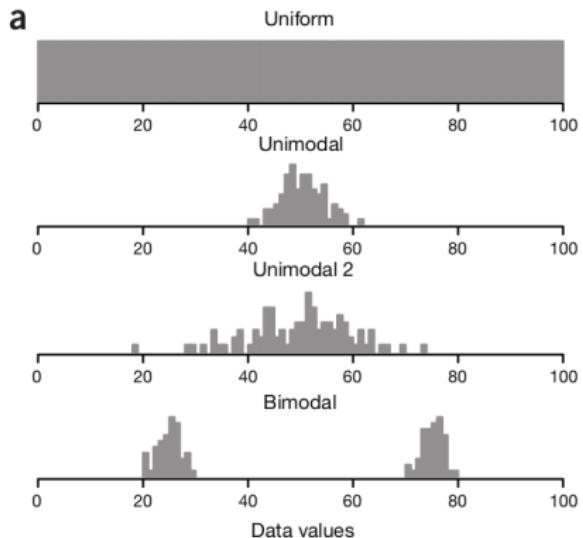
Variable Width Box Plot



Variable Width Notched Box Plot



box plot | 变体



- BoxPlotR
- Plotly
- ECplot
- Galaxy (“Graph/Display Data” 工具集中的 Boxplot)
- R
- ...



教学提纲

- 1 引言
- 2 基因组组装版本
- 3 基因组坐标系统
- 4 基因组注释常用格式
- 5 文本文件与文本编辑器
- 6 基因组坐标的逻辑运算
- 7 总结与答疑
- 8 引言
- 9 变异位点的注释
- 10 基因集富集分析

- 11 序列标识
- 12 box plot
- 13 解析图表
- 14 总结与答疑
- 15 引言
- 16 Galaxy 分析平台
- 17 Galaxy 使用演示
- 18 数据处理三段论
- 19 总结与答疑
- 20 复习思考题



表格

- 行、列的含义
- 缩写的含义
- 数值的含义

图片

- 生成图片的数据
- 横、纵轴的含义
- 图片包含的元素
- 图片元素属性的含义



表格

- 行、列的含义
- 缩写的含义
- 数值的含义

图片

- 生成图片的数据
- 横、纵轴的含义
- 图片包含的元素
- 图片元素属性的含义

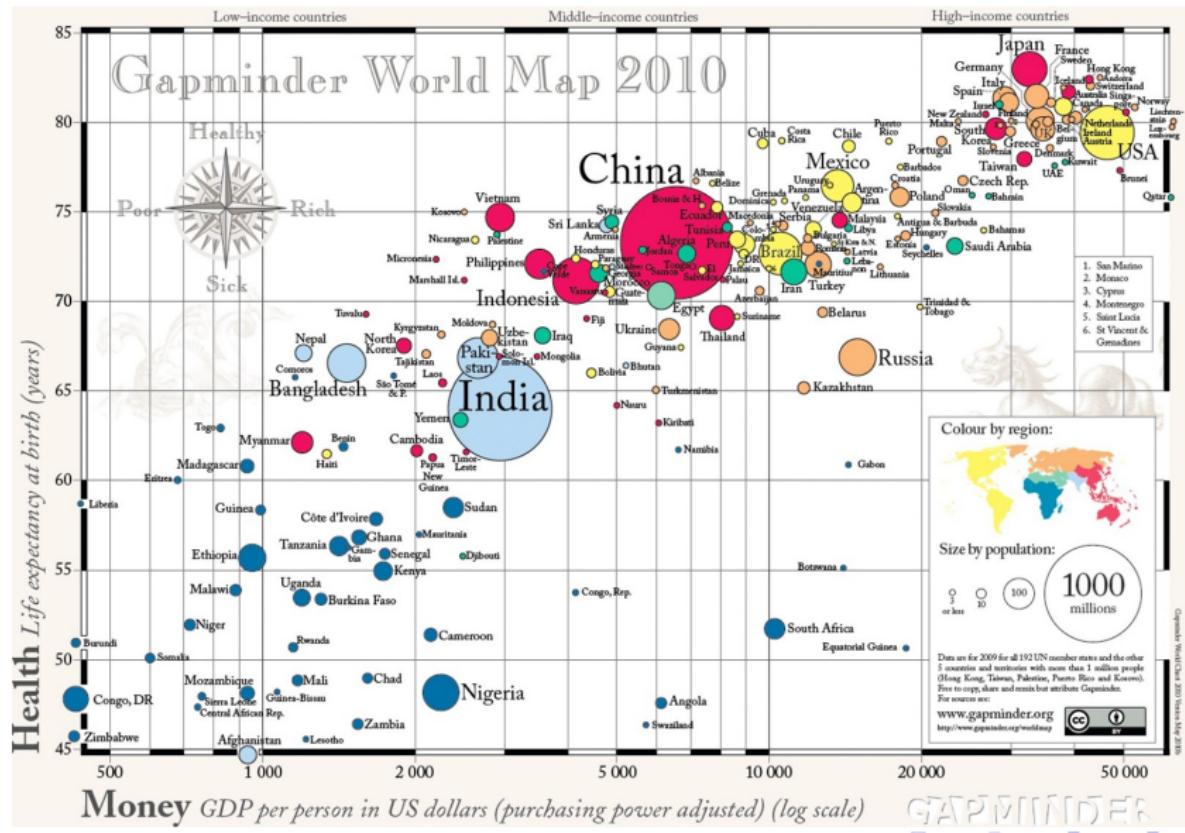


解析图表 | 图形的语法

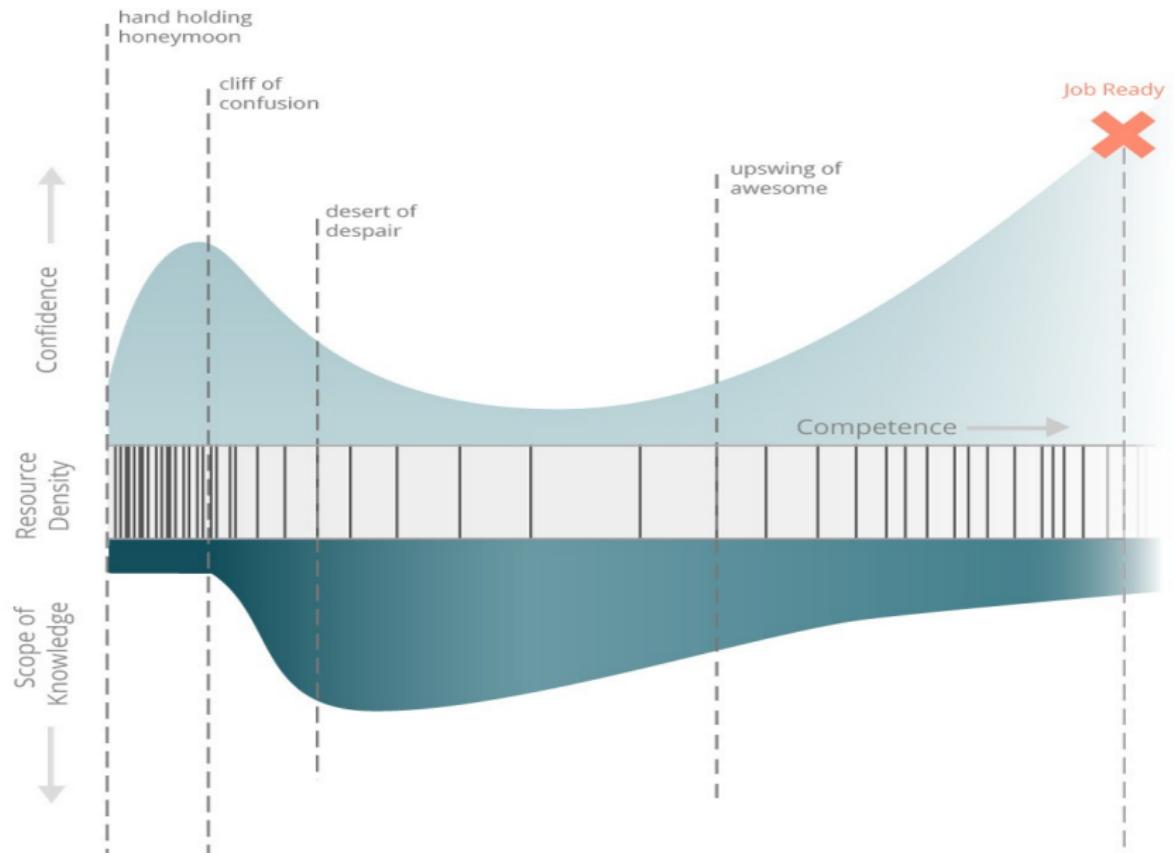
- 数据 (data) : 我们想要可视化的对象。
- 几何对象 (geom) : 用以呈现数据、在图中实际看到的图形元素 (点、线、条形、多边形等)。
- 图形属性 (aes) : 几何对象的视觉属性 (位置、颜色、形状、大小和线条类型)。
- 映射 (mapping) : 从数据中的变量对应到图形属性。
- 统计变换 (stats) : 对数据进行的某种汇总 (将数据分组计数以创建直方图)。
- 标度 (scale) : 控制着数据空间的值到图形属性空间的值的映射 (用颜色、大小或形状来表示不同的取值)。
- 坐标系 (coord) : 描述了数据是如何映射到图形所在的平面的，它同时提供了看图所需的坐标轴和网格线。
- 引导元素 (guide) : 向看图者展示了如何将视觉属性映射回数据空间 (坐标轴上的刻度线和标签、图例)。
- 分面 (facet) : 描述了如何将数据分解为各个子集，以及如何对子集作图并联合进行展示。分面也叫做条件作图或网格作图。
- 位置调整：控制着图形对象的重叠。



解析图表 | 图形的语法 | ggplot2



解析图表 | 其他



教学提纲

- 1 引言
- 2 基因组组装版本
- 3 基因组坐标系统
- 4 基因组注释常用格式
- 5 文本文件与文本编辑器
- 6 基因组坐标的逻辑运算
- 7 总结与答疑
- 8 引言
- 9 变异位点的注释
- 10 基因集富集分析

- 11 序列标识
- 12 box plot
- 13 解析图表
- 14 总结与答疑
- 15 引言
- 16 Galaxy 分析平台
- 17 Galaxy 使用演示
- 18 数据处理三段论
- 19 总结与答疑
- 20 复习思考题



知识点——基因组功能的高级注释

- 变异位点的注释——用途，注释内容，注释工具
- 基因集富集分析——功能，分析工具
- 序列标识——含义，制作工具
- box plot——理解，绘制

技能——解析图表

- 表——行列，缩写，数值
- 图——数据，横纵轴，图元素，元素属性



教学提纲

- 1 引言
- 2 基因组组装版本
- 3 基因组坐标系统
- 4 基因组注释常用格式
- 5 文本文件与文本编辑器
- 6 基因组坐标的逻辑运算
- 7 总结与答疑
- 8 引言
- 9 变异位点的注释
- 10 基因集富集分析

- 11 序列标识
- 12 box plot
- 13 解析图表
- 14 总结与答疑
- 15 引言
- 16 Galaxy 分析平台
- 17 Galaxy 使用演示
- 18 数据处理三段论
- 19 总结与答疑
- 20 复习思考题



基础知识

- 组装版本和坐标系统
- 常用格式
- 坐标的逻辑运算

高级注释

- 变异位点的注释
- 基因集富集分析
- 制作序列标识

分析平台

- Galaxy
- GenePattern

基础知识

- 组装版本和坐标系统
- 常用格式
- 坐标的逻辑运算

高级注释

- 变异位点的注释
- 基因集富集分析
- 制作序列标识

分析平台

- Galaxy
- GenePattern

基础知识

- 组装版本和坐标系统
- 常用格式
- 坐标的逻辑运算

高级注释

- 变异位点的注释
- 基因集富集分析
- 制作序列标识

分析平台

- Galaxy
- GenePattern

教学提纲

- 1 引言
- 2 基因组组装版本
- 3 基因组坐标系统
- 4 基因组注释常用格式
- 5 文本文件与文本编辑器
- 6 基因组坐标的逻辑运算
- 7 总结与答疑
- 8 引言
- 9 变异位点的注释
- 10 基因集富集分析

- 11 序列标识
- 12 box plot
- 13 解析图表
- 14 总结与答疑
- 15 引言
- 16 Galaxy 分析平台
- 17 Galaxy 使用演示
- 18 数据处理三段论
- 19 总结与答疑
- 20 复习思考题



- Get Data
- Text Manipulation
- Convert Formats
- Operate on Genomic Intervals
- Phenotype Association
- Statistics
- Graph/Display Data
- NGS Toolbox
- ...



The screenshot shows the Galaxy web interface with the following components:

- Top Navigation Bar:** Analyze Data, Workflow, Shared Data, Visualization, Cloud, Help, User, and a grid icon.
- Left Panel (Tools):** A sidebar titled "工具" containing a search bar and a list of links:
 - search tools
 - [Get Data](#)
 - [Lift-Over](#)
 - [Text Manipulation](#)
 - [Convert Formats](#)
 - [FASTA manipulation](#)
 - [Filter and Sort](#)
 - [Join, Subtract and Group](#)
 - [Extract Features](#)
 - [Fetch Sequences](#)
 - [Fetch Alignments](#)
 - [Get Genomic Scores](#)
 - [Operate on Genomic Intervals](#)
 - [Statistics](#)
 - [Graph/Display Data](#)
 - [Regional Variation](#)
 - [Multiple regression](#)
 - [Multivariate Analysis](#)
- Main Content Area:** A large central area with the title "Galaxy 101" and subtitle "Start small". Below it is the text "The very first tutorial you need". At the bottom of this area is a horizontal navigation bar with several small circular icons.
- Right Panel (History):** A sidebar titled "历史" showing an "Unnamed history" entry (1.5 MB). It includes a search icon and a message in Chinese: "历史已空，请单击左边窗格中'获取数据'" (The history is empty, please click on the left panel to get data).



教学提纲

- 1 引言
- 2 基因组组装版本
- 3 基因组坐标系统
- 4 基因组注释常用格式
- 5 文本文件与文本编辑器
- 6 基因组坐标的逻辑运算
- 7 总结与答疑
- 8 引言
- 9 变异位点的注释
- 10 基因集富集分析

- 11 序列标识
- 12 box plot
- 13 解析图表
- 14 总结与答疑
- 15 引言
- 16 Galaxy 分析平台
- 17 Galaxy 使用演示
- 18 数据处理三段论
- 19 总结与答疑
- 20 复习思考题



- UCSC liftOver tool：支持 BED 和 “chrN:start-end” 格式的输入
- Galaxy (基于 UCSC liftOver tool)：支持 BED、GFF 和 GTF 格式的输入
- CrossMap：支持 SAM/BAM、Wiggle/BigWig、BED、GFF/GTF 和 VCF 格式的输入，输出对应格式
- NCBI Remap：支持 BED、GFF、GTF 和 VCF 等格式的输入
- Ensembl assembly converter (2015 年退休, CrossMap 继位)：支持 BED、GFF、GTF 和 PSL 格式的输入，但输出都是 GFF 格式的
- pyliftover：仅支持点坐标 (point coordinates) 的转换，无法对区段 (ranges) 坐标进行转换



Galaxy 演示 | 坐标转换 | liftOver

hg19 ⇒ hg18

获取输入

输出文件



hg19 ⇒ hg18

① 获取输入

- 输入文件：hg19 坐标

② 数据处理

- 使用 Galaxy 中的“坐标转换”工具

③ 保存输出



hg19 ⇒ hg18

① 获取输入

- 输入文件：hg19 坐标

② 数据处理

→ 选择工具：hg19 ⇒ hg18

③ 保存输出



hg19 \Rightarrow hg18

① 获取输入

- 输入文件：hg19 坐标

② 数据处理

- 设置参数：hg19 \Rightarrow hg18

③ 保存输出



hg19 \Rightarrow hg18

① 获取输入

- 输入文件：hg19 坐标

② 数据处理

- 设置参数：hg19 \Rightarrow hg18

③ 保存输出

正在处理... MAPPED



hg19 \Rightarrow hg18

- ① 获取输入
 - 输入文件：hg19 坐标
- ② 数据处理
 - 设置参数：hg19 \Rightarrow hg18
- ③ 保存输出
 - 过滤结果：MAPPED VS. UNMAPPED



hg19 \Rightarrow hg18

- ① 获取输入
 - 输入文件：hg19 坐标
- ② 数据处理
 - 设置参数：hg19 \Rightarrow hg18
- ③ 保存输出
 - 过滤结果：MAPPED VS. UNMAPPED



Galaxy 演示 | 格式转换 | BED ⇌ GFF

BED ⇌ GFF

① 获取输入

选择文件



BED ⇌ GFF

① 获取输入

- 输入文件：BED

② 数据处理

③ 保存输出



BED ⇌ GFF

① 获取输入

- 输入文件：BED

② 数据处理

BED → GFF

选择下方的“GFF”

③ 保存输出



BED ⇌ GFF

① 获取输入

- 输入文件：BED

② 数据处理

- ① BED ⇒ GFF
- ② GFF ⇒ BED

③ 保存输出



BED ⇌ GFF

① 获取输入

- 输入文件：BED

② 数据处理

① BED ⇒ GFF

② GFF ⇒ BED

③ 保存输出



BED ⇌ GFF

① 获取输入

- 输入文件：BED

② 数据处理

- ① BED ⇒ GFF
- ② GFF ⇒ BED

③ 保存输出

- 运行结果：互相比较



BED ⇌ GFF

① 获取输入

- 输入文件：BED

② 数据处理

- ① BED ⇒ GFF
- ② GFF ⇒ BED

③ 保存输出

- 查看结果：互相对比



BED ⇌ GFF

① 获取输入

- 输入文件：BED

② 数据处理

- ① BED ⇒ GFF
- ② GFF ⇒ BED

③ 保存输出

- 查看结果：互相比较



- Galaxy 中的 “Operate on Genomic Intervals” 工具集
- BEDTools: a powerful toolset for genome arithmetic
- BEDOPS: the fast, highly scalable and easily-parallelizable genome analysis toolkit



外显子 vs. SNP

① 获取输入



外显子 vs. SNP

① 获取输入

- exon
- SNP

② 数据处理

③ 保存输出



外显子 vs. SNP

① 获取输入

- exon
- SNP

② 数据处理

③ 保存输出



外显子 vs. SNP

① 获取输入

- exon
- SNP

② 数据处理

- subtract
- join

③ 保存输出



外显子 vs. SNP

① 获取输入

- exon
- SNP

② 数据处理

- subtract
- join

③ 保存输出



外显子 vs. SNP

① 获取输入

- exon
- SNP

② 数据处理

- **subtract**
- join

③ 保存输出



外显子 vs. SNP

① 获取输入

- exon
- SNP

② 数据处理

- subtract
- join

③ 保存输出

→ 分析结果



外显子 vs. SNP

① 获取输入

- exon
- SNP

② 数据处理

- subtract
- join

③ 保存输出

- 解析结果



外显子 vs. SNP

① 获取输入

- exon
- SNP

② 数据处理

- subtract
- join

③ 保存输出

- 解析结果



问题

- 找到含有至少 N (2) 个 SNP 的外显子。

输入

- 外显子数据 (BED 格式)
- SNP 数据 (BED 格式)

输出

- 满足要求的外显子信息 (BED 格式)



问题

- 找到含有至少 N (2) 个 SNP 的外显子。

输入

- 外显子数据 (BED 格式)
- SNP 数据 (BED 格式)

输出

- 满足要求的外显子信息 (BED 格式)



问题

- 找到含有至少 N (2) 个 SNP 的外显子。

输入

- 外显子数据 (BED 格式)
- SNP 数据 (BED 格式)

输出

- 满足要求的外显子信息 (BED 格式)



外显子数据

①	chr1	10	20	exon1	0	+
②	chr1	30	40	exon2	0	+
③	chr1	50	60	exon3	0	-
④	chr1	65	75	exon4	0	+
⑤	chr1	85	95	exon5	0	-



SNP 数据

①	chr1	11	12	snp1	0	+
②	chr1	15	16	snp2	0	+
③	chr1	17	18	snp3	0	+
④	chr1	24	25	snp4	0	+
⑤	chr1	33	34	snp5	0	+
⑥	chr1	37	38	snp6	0	+
⑦	chr1	44	45	snp7	0	+
⑧	chr1	54	55	snp8	0	+
⑨	chr1	57	58	snp9	0	-



Finding exons with the highest number (≥ 10) of SNPs

Join-Group-Filter-Compare-Sort

- Input: Get exons, SNPs; UCSC Table Browser
- Join[Operate on Genomic Intervals]: Join exons with SNPs
- Group: Count the number of SNPs per exon
- Filter: Filter exons that have ten or more SNPs
- Compare two Datasets: Recover exon information
- Sort: Sort the start and end coordinates
- Visualize: Display data in genome browser



Finding exons with the highest number (≥ 10) of SNPs

Join-Group-Filter-Compare-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Group: Count the number of SNPs per exon
- ④ Filter: Filter SNPs with count >= 10
- ⑤ Sort: Sort by count (descending)



Finding exons with the highest number (≥ 10) of SNPs

Join-Group-Filter-Compare-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Group: Count the number of SNPs per exon
- ④ Filter: Filter exons that have ten or more SNPs
- ⑤ Compare two Datasets: Recover exon information
- ⑥ Sort: Sort the start and end coordinates
- ⑦ Visualize: Display data in genome browser



Finding exons with the highest number (≥ 10) of SNPs

Join-Group-Filter-Compare-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Group: Count the number of SNPs per exon
- ④ Filter: Filter exons that have ten or more SNPs
- ⑤ Compare two Datasets: Recover exon information
- ⑥ Sort: Sort the start and end coordinates
- ⑦ Visualize: Display data in genome browser



Finding exons with the highest number (≥ 10) of SNPs

Join-Group-Filter-Compare-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Group: Count the number of SNPs per exon
- ④ Filter: Filter exons that have ten or more SNPs
- ⑤ Compare two Datasets: Recover exon information
- ⑥ Sort: Sort the start and end coordinates
- ⑦ Visualize: Display data in genome browser



Finding exons with the highest number (≥ 10) of SNPs

Join-Group-Filter-Compare-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Group: Count the number of SNPs per exon
- ④ Filter: Filter exons that have ten or more SNPs
- ⑤ Compare two Datasets: Recover exon information
- ⑥ Sort: Sort the start and end coordinates
- ⑦ Visualize: Display data in genome browser



Finding exons with the highest number (≥ 10) of SNPs

Join-Group-Filter-Compare-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Group: Count the number of SNPs per exon
- ④ Filter: Filter exons that have ten or more SNPs
- ⑤ Compare two Datasets: Recover exon information
- ⑥ Sort: Sort the start and end coordinates
- ⑦ Visualize: Display data in genome browser



Finding exons with the highest number (≥ 10) of SNPs

Join-Group-Filter-Compare-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Group: Count the number of SNPs per exon
- ④ Filter: Filter exons that have ten or more SNPs
- ⑤ Compare two Datasets: Recover exon information
- ⑥ Sort: Sort the start and end coordinates
- ⑦ Visualize: Display data in genome browser



Finding exons with the highest number (≥ 10) of SNPs

Join-Group-Filter-Compare-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Group: Count the number of SNPs per exon
- ④ Filter: Filter exons that have ten or more SNPs
- ⑤ Compare two Datasets: Recover exon information
- ⑥ Sort: Sort the start and end coordinates
- ⑦ Visualize: Display data in genome browser



Finding exons with the highest number (≥ 10) of SNPs

Join-Count-Filter-Cut-Sort

- Input: Get exons, SNPs; UCSC Table Browser
- Join[Operate on Genomic Intervals]: Join exons with SNPs
- Count: Count the number of SNPs per exon
- Filter: Filter exons that have ten or more SNPs
- Cut: Cut columns to recover BED format
- Sort: Sort the start and end coordinates
- Visualize: Display data in genome browser



Finding exons with the highest number (≥ 10) of SNPs

Join-Count-Filter-Cut-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Count: Count the number of SNPs per exon
- ④ Sort: Sort by the number of SNPs per exon
- ⑤ Filter: Filter exons with at least 10 SNPs
- ⑥ Cut: Cut the top 10 exons
- ⑦ Sort: Sort the exons by their genomic position



Finding exons with the highest number (≥ 10) of SNPs

Join-Count-Filter-Cut-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Count: Count the number of SNPs per exon
- ④ Filter: Filter exons that have ten or more SNPs
- ⑤ Cut: Cut columns to recover BED format
- ⑥ Sort: Sort the start and end coordinates
- ⑦ Visualize: Display data in genome browser



Finding exons with the highest number (≥ 10) of SNPs

Join-Count-Filter-Cut-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Count: Count the number of SNPs per exon
- ④ Filter: Filter exons that have ten or more SNPs
- ⑤ Cut: Cut columns to recover BED format
- ⑥ Sort: Sort the start and end coordinates
- ⑦ Visualize: Display data in genome browser



Finding exons with the highest number (≥ 10) of SNPs

Join-Count-Filter-Cut-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Count: Count the number of SNPs per exon
- ④ Filter: Filter exons that have ten or more SNPs
- ⑤ Cut: Cut columns to recover BED format
- ⑥ Sort: Sort the start and end coordinates
- ⑦ Visualize: Display data in genome browser



Finding exons with the highest number (≥ 10) of SNPs

Join-Count-Filter-Cut-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Count: Count the number of SNPs per exon
- ④ Filter: Filter exons that have ten or more SNPs
- ⑤ Cut: Cut columns to recover BED format
- ⑥ Sort: Sort the start and end coordinates
- ⑦ Visualize: Display data in genome browser



Finding exons with the highest number (≥ 10) of SNPs

Join-Count-Filter-Cut-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Count: Count the number of SNPs per exon
- ④ Filter: Filter exons that have ten or more SNPs
- ⑤ Cut: **Cut columns to recover BED format**
- ⑥ Sort: Sort the start and end coordinates
- ⑦ Visualize: Display data in genome browser



Finding exons with the highest number (≥ 10) of SNPs

Join-Count-Filter-Cut-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Count: Count the number of SNPs per exon
- ④ Filter: Filter exons that have ten or more SNPs
- ⑤ Cut: Cut columns to recover BED format
- ⑥ Sort: Sort the start and end coordinates
- ⑦ Visualize: Display data in genome browser



Finding exons with the highest number (≥ 10) of SNPs

Join-Count-Filter-Cut-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Count: Count the number of SNPs per exon
- ④ Filter: Filter exons that have ten or more SNPs
- ⑤ Cut: Cut columns to recover BED format
- ⑥ Sort: Sort the start and end coordinates
- ⑦ Visualize: Display data in genome browser



Finding 10 exons with the highest number of SNPs

Join-Group-Sort-SelectFirst-Join-Cut-Sort

- Input: Get exons, SNPs; UCSC Table Browser
- Join[Operate on Genomic Intervals]: Join exons with SNPs
- Group: Count the number of SNPs per exon
- Sort: Sort exons by SNPs count
- Select first: Select top ten
- Join[Join two Datasets]: Recover exon information
- Cut: Cut columns to recover BED format
- Sort: Sort the start and end coordinates
- Visualize: Display data in genome browser

Finding 10 exons with the highest number of SNPs

Join-Group-Sort-SelectFirst-Join-Cut-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Group: Count the number of SNPs per exon



Finding 10 exons with the highest number of SNPs

Join-Group-Sort-SelectFirst-Join-Cut-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Group: Count the number of SNPs per exon
- ④ Sort: Sort exons by SNPs count
- ⑤ Select first: Select top ten
- ⑥ Join[Join two Datasets]: Recover exon information
- ⑦ Cut: Cut columns to recover BED format
- ⑧ Sort: Sort the start and end coordinates
- ⑨ Visualize: Display data in genome browser

Finding 10 exons with the highest number of SNPs

Join-Group-Sort-SelectFirst-Join-Cut-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Group: Count the number of SNPs per exon
- ④ Sort: Sort exons by SNPs count
- ⑤ Select first: Select top ten
- ⑥ Join[Join two Datasets]: Recover exon information
- ⑦ Cut: Cut columns to recover BED format
- ⑧ Sort: Sort the start and end coordinates
- ⑨ Visualize: Display data in genome browser

Finding 10 exons with the highest number of SNPs

Join-Group-Sort-SelectFirst-Join-Cut-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Group: Count the number of SNPs per exon
- ④ Sort: Sort exons by SNPs count
- ⑤ Select first: Select top ten
- ⑥ Join[Join two Datasets]: Recover exon information
- ⑦ Cut: Cut columns to recover BED format
- ⑧ Sort: Sort the start and end coordinates
- ⑨ Visualize: Display data in genome browser

Finding 10 exons with the highest number of SNPs

Join-Group-Sort-SelectFirst-Join-Cut-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Group: Count the number of SNPs per exon
- ④ Sort: Sort exons by SNPs count
- ⑤ Select first: Select top ten
- ⑥ Join[Join two Datasets]: Recover exon information
- ⑦ Cut: Cut columns to recover BED format
- ⑧ Sort: Sort the start and end coordinates
- ⑨ Visualize: Display data in genome browser

Finding 10 exons with the highest number of SNPs

Join-Group-Sort-SelectFirst-Join-Cut-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Group: Count the number of SNPs per exon
- ④ Sort: Sort exons by SNPs count
- ⑤ **Select first: Select top ten**
- ⑥ Join[Join two Datasets]: Recover exon information
- ⑦ Cut: Cut columns to recover BED format
- ⑧ Sort: Sort the start and end coordinates
- ⑨ Visualize: Display data in genome browser

Finding 10 exons with the highest number of SNPs

Join-Group-Sort-SelectFirst-Join-Cut-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Group: Count the number of SNPs per exon
- ④ Sort: Sort exons by SNPs count
- ⑤ Select first: Select top ten
- ⑥ Join[Join two Datasets]: Recover exon information
- ⑦ Cut: Cut columns to recover BED format
- ⑧ Sort: Sort the start and end coordinates
- ⑨ Visualize: Display data in genome browser

Finding 10 exons with the highest number of SNPs

Join-Group-Sort-SelectFirst-Join-Cut-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Group: Count the number of SNPs per exon
- ④ Sort: Sort exons by SNPs count
- ⑤ Select first: Select top ten
- ⑥ Join[Join two Datasets]: Recover exon information
- ⑦ Cut: Cut columns to recover BED format
- ⑧ Sort: Sort the start and end coordinates
- ⑨ Visualize: Display data in genome browser

Finding 10 exons with the highest number of SNPs

Join-Group-Sort-SelectFirst-Join-Cut-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Group: Count the number of SNPs per exon
- ④ Sort: Sort exons by SNPs count
- ⑤ Select first: Select top ten
- ⑥ Join[Join two Datasets]: Recover exon information
- ⑦ Cut: Cut columns to recover BED format
- ⑧ Sort: Sort the start and end coordinates
- ⑨ Visualize: Display data in genome browser

Finding 10 exons with the highest number of SNPs

Join-Group-Sort-SelectFirst-Join-Cut-Sort

- ① Input: Get exons, SNPs; UCSC Table Browser
- ② Join[Operate on Genomic Intervals]: Join exons with SNPs
- ③ Group: Count the number of SNPs per exon
- ④ Sort: Sort exons by SNPs count
- ⑤ Select first: Select top ten
- ⑥ Join[Join two Datasets]: Recover exon information
- ⑦ Cut: Cut columns to recover BED format
- ⑧ Sort: Sort the start and end coordinates
- ⑨ Visualize: Display data in genome browser

Finding exons with the highest number (≥ 10) of SNPs

BEDTools-shell-BEDTools

考虑链性

```
bedtools intersect -a exon.bed -b snp.bed -c -s |  
awk '{if($7>=10) print;}' | cut -f1-6 |  
bedtools sort -i stdin
```

不考虑链性

```
bedtools intersect -a exon.bed -b snp.bed -c |  
awk '{if($7>=10) print;}' | cut -f1-6 |  
bedtools sort -i stdin
```

Finding exons with the highest number (≥ 10) of SNPs

BEDTools-shell-BEDTools

考虑链性

```
bedtools intersect -a exon.bed -b snp.bed -c -s |  
awk '{if($7>=10) print;}' | cut -f1-6 |  
bedtools sort -i stdin
```

不考虑链性

```
bedtools intersect -a exon.bed -b snp.bed -c |  
awk '{if($7>=10) print;}' | cut -f1-6 |  
bedtools sort -i stdin
```

Finding exons with the highest number (≥ 10) of SNPs

BEDTools-shell-BEDTools

考虑链性

```
bedtools intersect -a exon.bed -b snp.bed -c -s |  
awk '{if($7>=10) print;}' | cut -f1-6 |  
bedtools sort -i stdin
```

不考虑链性

```
bedtools intersect -a exon.bed -b snp.bed -c |  
awk '{if($7>=10) print;}' | cut -f1-6 |  
bedtools sort -i stdin
```

Finding exons with the highest number (≥ 10) of SNPs

BEDTools-shell-BEDTools

考虑链性

```
bedtools intersect -a exon.bed -b snp.bed -c -s |  
awk '{if($7>=10) print;}' | cut -f1-6 |  
bedtools sort -i stdin
```

不考虑链性

```
bedtools intersect -a exon.bed -b snp.bed -c |  
awk '{if($7>=10) print;}' | cut -f1-6 |  
bedtools sort -i stdin
```

Workflow: create, modify, rerun, share

Save: Rename the history as "Exons and SNPs"

Workflow: Extract workflow from history

Workflow: Create a new workflow from scratch



Workflow: create, modify, rerun, share

- ① Save: Rename the history as “Exons and SNPs”
- ② Workflow: Extract workflow from history
- ③ Modify: Open workflow editor and modify the parameter
- ④ Rerun: Run workflow on whole genome data
- ⑤ Share: Share or publish workflow
- ⑥ Create: Create workflows from scratch (e.g. Find the 50 longest exons)



Workflow: create, modify, rerun, share

- ① Save: Rename the history as “Exons and SNPs”
- ② Workflow: Extract workflow from history
- ③ Modify: Open workflow editor and modify the parameter
- ④ Rerun: Run workflow on whole genome data
- ⑤ Share: Share or publish workflow
- ⑥ Create: Create workflows from scratch (e.g. Find the 50 longest exons)



Workflow: create, modify, rerun, share

- ① Save: Rename the history as “Exons and SNPs”
- ② Workflow: Extract workflow from history
- ③ **Modify: Open workflow editor and modify the parameter**
- ④ Rerun: Run workflow on whole genome data
- ⑤ Share: Share or publish workflow
- ⑥ Create: Create workflows from scratch (e.g. Find the 50 longest exons)



Workflow: create, modify, rerun, share

- ① Save: Rename the history as “Exons and SNPs”
- ② Workflow: Extract workflow from history
- ③ Modify: Open workflow editor and modify the parameter
- ④ **Rerun: Run workflow on whole genome data**
- ⑤ Share: Share or publish workflow
- ⑥ Create: Create workflows from scratch (e.g. Find the 50 longest exons)



Workflow: create, modify, rerun, share

- ① Save: Rename the history as “Exons and SNPs”
- ② Workflow: Extract workflow from history
- ③ Modify: Open workflow editor and modify the parameter
- ④ Rerun: Run workflow on whole genome data
- ⑤ Share: Share or publish workflow
- ⑥ Create: Create workflows from scratch (e.g. Find the 50 longest exons)



Workflow: create, modify, rerun, share

- ① Save: Rename the history as “Exons and SNPs”
- ② Workflow: Extract workflow from history
- ③ Modify: Open workflow editor and modify the parameter
- ④ Rerun: Run workflow on whole genome data
- ⑤ Share: Share or publish workflow
- ⑥ Create: Create workflows from scratch (e.g. Find the 50 longest exons)



教学提纲

1

引言

2

基因组组装版本

3

基因组坐标系统

4

基因组注释常用格式

5

文本文件与文本编辑器

6

基因组坐标的逻辑运算

7

总结与答疑

8

引言

9

变异位点的注释

10

基因集富集分析

11

序列标识

12

box plot

13

解析图表

14

总结与答疑

15

引言

16

Galaxy 分析平台

17

Galaxy 使用演示

18

数据处理三段论

19

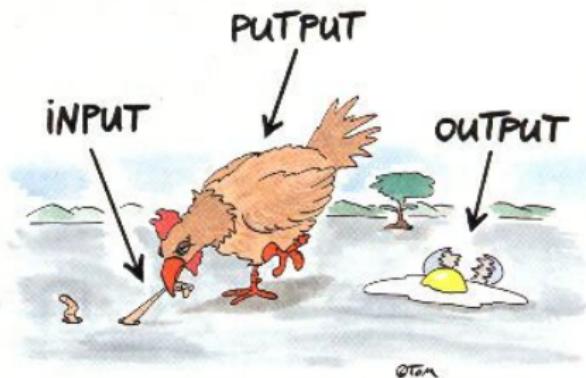
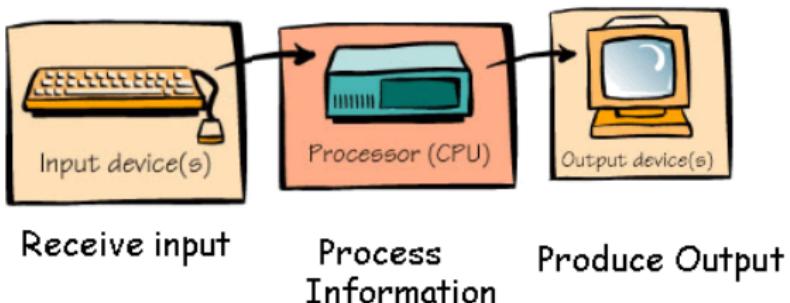
总结与答疑

20

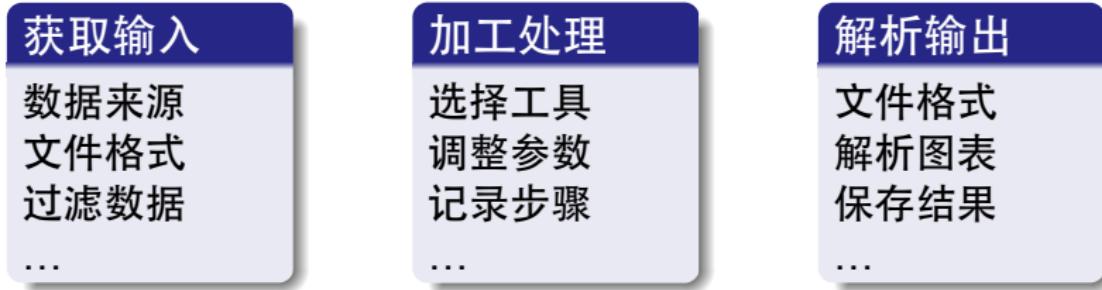
复习思考题



三段论：“输入-加工-输出”



三段论：“输入-加工-输出”



教学提纲

1 引言

基因组组装版本

基因组坐标系统

基因组注释常用格式

文本文件与文本编辑器

基因组坐标的逻辑运算

总结与答疑

引言

变异位点的注释

基因集富集分析

11 序列标识

12 box plot

13 解析图表

14 总结与答疑

15 引言

16 Galaxy 分析平台

17 Galaxy 使用演示

18 数据处理三段论

19 总结与答疑

20 复习思考题



知识点——Galaxy 分析平台

- Galaxy——界面、学习、使用

技能——“输入-加工-输出”三段论

- 获取输入——格式、来源、过滤
- 数据处理——工具、版本、参数
- 解析输出——格式、注释、解析



教学提纲

- 11 序列标识
 - 12 box plot
 - 13 解析图表
 - 14 总结与答疑
 - 15 引言
 - 16 Galaxy 分析平台
 - 17 Galaxy 使用演示
 - 18 数据处理三段论
 - 19 总结与答疑
 - 20 复习思考题



复习思考题

知识点

- ① hg19 和 mm10 分别代表什么含义？hg19 是和 GRCh37 相对应，还是和 GRCm38 相对应？
- ② 常见的基因组坐标系统是哪两种，举例进行说明。
- ③ 简述 BED 格式前 6 列的含义，能解释实际的 BED 记录。
- ④ 基于基因组坐标的常见逻辑运算模式有哪些，画图进行解释。
- ⑤ 简述序列标识的含义，能解释实际的序列标识图。

技能

- ① 不同操作系统的换行符有何区别？
- ② 以 SNP 的注释结果为例，论述如何解析一张表。
- ③ 以 box plot 为例，论述如何解析一张图。
- ④ 以坐标转换为例，论述“输入-加工-输出”的工作流程。

Powered by



T_EX L^AT_EX X_ET_EX Beamer