

lab seminar

Session1

JUN SOUNG KWAK



FASTA files

```
Junsoungs-Mac-mini:~ jun$ head /Volumes/Jun_SSD/genome/Paralichthys_olivaceus-Ensembl/ensembl/genome/Paralichthys_olivaceus-GCA_024713975.2-softmasked.fasta
>1 softmasked:chromosome primary_assembly:ASM2471397v2:1:1:31859138:1
CGAACTAGGTGGATGGGCTTTCTGTAAAGACGGCTCCTGTCGCGGTTTCAAAGCTCTTCG
TGCTTTCTGGCATTGCTGGAGGCCCGTCAAATTGCAATCTATAATGAAAGGGGAATCAG
AAGTTAGTGTAAGGCCATTATCAATGCCTAGTATATGCTGGCCCCTCATTTAGCCAAA
AAAGCCTCAAAGTGTGAAAGTTACACCTTTGCGATCATTAGGCCCCCGATCAACATT
TTTTCCCTCTCATTACCCGTATGTACTCTGTATGGCTATGAAATGGAGCCCTTAGAAAT
GGGCTTAATTGATACCAAAGGTCACAAAGGTCACCTGGGCCCTTCTGCGTTTTTGGTGAA
AATCGGTAATGCAGGCTCAAAGAGCCCCTTGTGACCAGTGGGTCAAATTTGGGGCTCCTA
GGCCCCTGGGAAGGCCCGAAAGGGATCGTAATTTTCCAAAAATCTCCAGAACCAAGAG
GCTGCGAGCCCTGTTTAGGCCCTGACCAGCTAGAGGGCTGAAATCCAAGTGACACCTCA
```

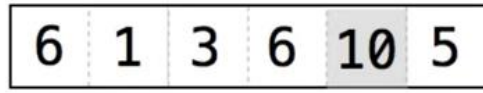
header

sequence

FASTA file is a text file that contains a nucleotide sequence of DNA or RNA or an amino acid sequence of a protein and related information.

Extension	Meaning	Notes
fasta, fa	generic FASTA	Any generic fasta file
fna	FASTA nucleic acid	Used generically to specify nucleic acids
ffn	FASTA nucleotide of gene regions	Contains coding regions for a genome
faa	FASTA amino acid	Contains amino acid sequences
frn	FASTA non-coding RNA	Contains non-coding RNA regions for a genome

Index (Indexing): 색인, 목록



`vec[5]`

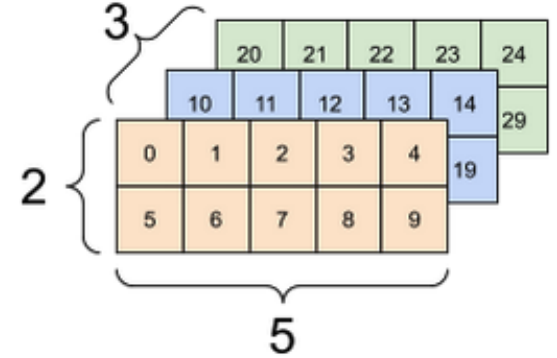
John	1940	guitar
Paul	1941	bass
George	1943	guitar
Ringo	1940	drums

`df[2, c(2,3)]`

```
[[[ 0  1  2  3  4]
   [ 5  6  7  8  9]]

 [[10 11 12 13 14]
   [15 16 17 18 19]]

 [[20 21 22 23 24]
   [25 26 27 28 29]]]
```



3D Array (Tensor)

Index notation is a standard method used to specify or reference individual elements of an array, vector, matrix, or tensor in mathematics, computer science, and related fields.

It provides a systematic way to describe specific elements within these data structures.

1D Array (Vector)

2D Array (Matrix, data frame)

3D Array (Tensor)

FASTA Index

: 유전자 데이터를 효율적으로 저장하고 검색하기 위한 기술.

```
Junsoungs-Mac-mini:~ jun$ head -n 50 /Volumes/Jun_SSD/genome/Paralichthys_olivaceus/Ensembl_rapid/genome/Paralichthys_olivaceus-GCA_024713975.2-softmasked.fa.fai
1      31859138      70      60      61
2      29887412      32390264      60      61
3      29078691      62775870      60      61
4      28192728      92339276      60      61
5      28185977      121001953     60      61
6      28085797      149657767     60      61
7      27885095      178211731     60      61
8      27660206      206561648     60      61
9      27364311      234682928     60      61
10     26501863      262503383     60      61
11     26045262      289447016     60      61
12     25908411      315926438     60      61
13     25560671      342266728     60      61
14     24856284      368253483     60      61
15     24354382      393524111     60      61
16     23785052      418284472     60      61
17     21261208      442466014     60      61
18     20994079      464081648     60      61
19     20610995      485425701     60      61
20     20325903      506380285     60      61
21     19848577      527045026     60      61
22     17781304      547224485     60      61
23     17387852      565302217     60      61
24     14749316      582979939     60      61
MT     17090      597975146     60      61
```

()D array (,)

1st column: 염색체 이름, 또는 시퀀스 이름

2nd column: 염색체 또는 시퀀스의 길이 (염기 수)

3rd column: 해당 염색체 또는 시퀀스가 **FASTA file** 내에서 시작되는 바이트 위치

4th column: **FASTA file**에서 한 줄의 염기 서열에 포함된 문자의 개수 (보통 60으로 설정)

5th column: **FASTA file**에서 한 줄의 염기 서열과 줄 바꿈을 포함한 길이 (보통 61로 설정, 60 (염기서열) + 1 (줄바꿈))

GFF / GTF format

GFF3

```
1  ASM2471397v2  region  1  31859138  .  .  .  ID=region:1;Alias=CM045671.1
###
1  ensembl gene    30565  32060  .  -  .  ID=gene:ENSPOLG00010000718;biotype=protein_coding;gene_id=ENSPOLG00010000718;version=1
1  ensembl mRNA    30565  32060  .  -  .  ID=transcript:ENSPOLT00010001876;Parent=gene:ENSPOLG00010000718;biotype=protein_coding;tag=Ensembl_canonical;transcript_id=ENSPOLT00010001876;version=1
1  ensembl exon    30565  30740  .  -  .  Parent=transcript:ENSPOLT00010001876;constitutive=1;exon_id=ENSPOLE00010009024;rank=2;version=1
1  ensembl three_prime_UTR 30565  30740  .  -  .  Parent=transcript:ENSPOLT00010001876
1  ensembl three_prime_UTR 30829  31048  .  -  .  Parent=transcript:ENSPOLT00010001876
1  ensembl exon    30829  32060  .  -  .  Parent=transcript:ENSPOLT00010001876;constitutive=1;exon_id=ENSPOLE00010009023;rank=1;version=1
1  ensembl CDS     31049  31696  .  -  0  ID=CDS:ENSPOLP00010001795;Parent=transcript:ENSPOLT00010001876;protein_id=ENSPOLP00010001795;version=1
1  ensembl five_prime_UTR 31697  32060  .  -  .  Parent=transcript:ENSPOLT00010001876
```

GTF

```
1  ensembl gene    30565  32060  .  -  .  gene_id "ENSPOLG00010000718"; gene_version "1"; gene_source "ensembl"; gene_biotype "protein_coding";
1  ensembl transcript 30565  32060  .  -  .  gene_id "ENSPOLG00010000718"; gene_version "1"; transcript_id "ENSPOLT00010001876"; transcript_version "1"; gene_source "ensembl"; gene_biotype "protein_coding"; transcript_source
"ensembl"; transcript_biotype "protein_coding"; tag "Ensembl_canonical";
1  ensembl exon    30829  32060  .  -  .  gene_id "ENSPOLG00010000718"; gene_version "1"; transcript_id "ENSPOLT00010001876"; transcript_version "1"; exon_number "1"; gene_source "ensembl"; gene_biotype "protein_coding";
transcript_source "ensembl"; transcript_biotype "protein_coding"; exon_id "ENSPOLE00010009023"; exon_version "1"; tag "Ensembl_canonical";
1  ensembl CDS     31052  31696  .  -  0  gene_id "ENSPOLG00010000718"; gene_version "1"; transcript_id "ENSPOLT00010001876"; transcript_version "1"; exon_number "1"; gene_source "ensembl"; gene_biotype "protein_coding";
transcript_source "ensembl"; transcript_biotype "protein_coding"; protein_id "ENSPOLP00010001795"; protein_version "1"; tag "Ensembl_canonical";
1  ensembl start_codon 31694  31696  .  -  0  gene_id "ENSPOLG00010000718"; gene_version "1"; transcript_id "ENSPOLT00010001876"; transcript_version "1"; exon_number "1"; gene_source "ensembl"; gene_biotype "protein_coding";
transcript_source "ensembl"; transcript_biotype "protein_coding"; tag "Ensembl_canonical";
1  ensembl stop_codon 31049  31051  .  -  0  gene_id "ENSPOLG00010000718"; gene_version "1"; transcript_id "ENSPOLT00010001876"; transcript_version "1"; exon_number "1"; gene_source "ensembl"; gene_biotype "protein_coding";
transcript_source "ensembl"; transcript_biotype "protein_coding"; tag "Ensembl_canonical";
1  ensembl exon    30565  30740  .  -  .  gene_id "ENSPOLG00010000718"; gene_version "1"; transcript_id "ENSPOLT00010001876"; transcript_version "1"; exon_number "2"; gene_source "ensembl"; gene_biotype "protein_coding";
transcript_source "ensembl"; transcript_biotype "protein_coding"; exon_id "ENSPOLE00010009024"; exon_version "1"; tag "Ensembl_canonical";
```

GFF는 일반적인 유전체 annotation 형식인 반면, GTF는 엄격하게 유전자 annotation 형식임.

GFF는 덜 제어됨으로 더 많은 변수를 필요로 함.

GTF → GFF 변환은 쉬움 / 필요한 gene/transcript_id 속성을 GFF안의 ID 속성으로 변환하면 됨.

GFF → GTF 변환은



Sequence alignment and map (SAM) or binary alignment and map (BAM) file

NGS 데이터를 정렬한 결과를 저장하는 파일 형식.

SAM file은 사람이 읽을 수 있지만, BAM은 사람이 이해할 수 없음 - 컴퓨터만 이해함.

SAM file 구조 –
1. Header (헤더)

@HD VN:1.6 SO:coordinate
@SQ SN:chr1 LN:248956422
@SQ SN:chr2 LN:242193529

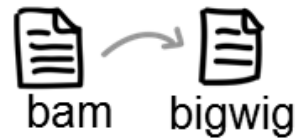
2. Alignment section (정렬 데이터)

1	2	3	4	5	6	7	8	9	10	11	12
리드 이름	FLAG	참조 서열 이름	정렬 시작 위치	매핑 품질 점수	CIGAR 문자열	정렬된 참조 서열 이름	정렬된 위치	길이	염기서열	품질 기호	사전정의 태그
SRR067577.2766	99	chr1	73240003	60	101M	=	73240004	102	GCTA.....	FHG@...	NM:I:0

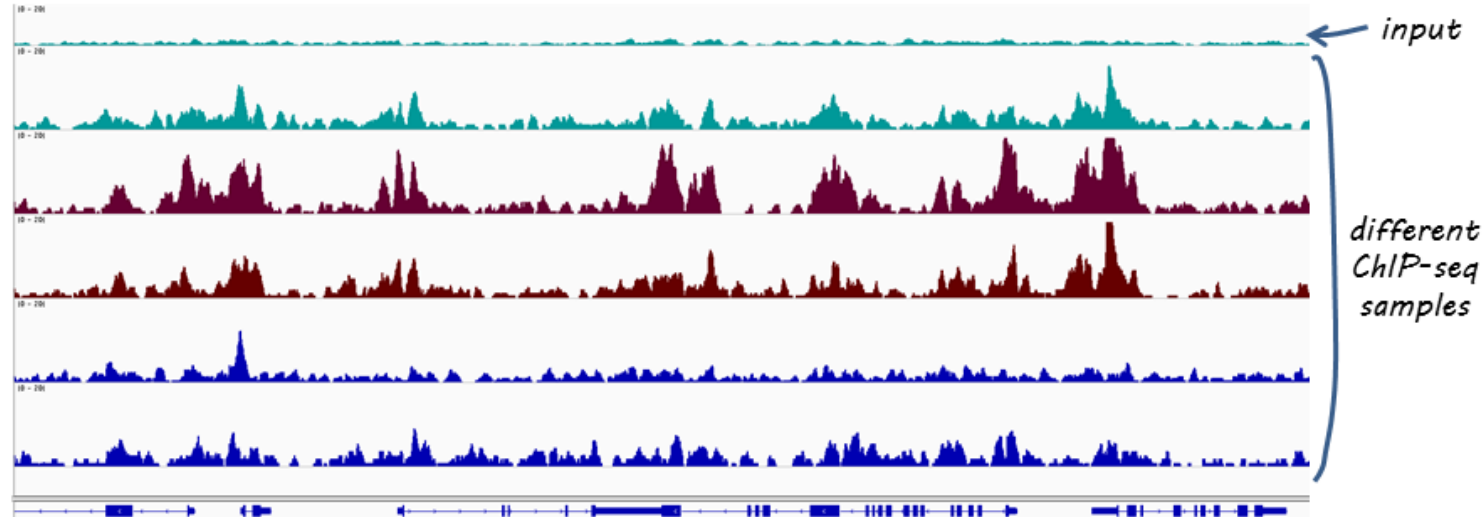
쌍을 이루는 리드


```
Jun$ head /Volumes/Scratch/_nf-core/RNA_seq/2024_SMU_KYJ/2_Olive_floUNDER/OLive_floUNDER_RNA_seq/star_salmon/C1.markdup.sorted.bam
?BCG)?{$$}??{[?]3?H?h&U??qmqm?????9
D?j? ,w???=A?w????
?c?Eq?@?22?J<%?*_ÿ?#U??JI???R????k?????
}}?f{~??;??~????{u[??]}j???1?b'??
?????0?o?n??\????2?8p|?q????rXë?"
,??v!#+?"??U??*'s??E?n*n?d2?(
?^?l"2??*v?$1#????]?K?z??e?RI dXE?R?.e?????UF|w?k%?#BR?Z?"J?"x?v?^xep_?b?
u?v?XNb9???UcTU`d?+???H,
??(pv???W7I*???19?n?,+,?x1??RA????3IÄ"2?-{[?]f???))????\??/w?????OP)nB?`?{Loçyn?撰Vo?g}?f????r?j???
?6:??_[=5_e????????SU?*wmc?畚~????z[?]pl????Fu}?·m??)ve???|k?vV?2q?Cp?u册???+?7??<{vu?T?????
??
?xn~???????V?? ?6???mm?NlIm?4???(\??}?y1}???)B???77?6??0??X?wu??????,?h??22???ķ7?-fáo?U??
?? .n?Y??j?r?g?93q?_?{???????X?Uoo>?v??U{?????uj?ç??;
?s??{??????[?]?\> _?&?????5g?NYot???u?????IW??V?non??B??N?[?]?????
??
t6NF[???|?0?????V?8wf?????X;?????'? ?k/?_7;?x???x??,|?s??s?c?çc.
????k???'&????]?[?czE??x|vJ?L/WO????Z?3?????入δ^>
;s?+??Ru??9xw?v|V?:?
?+???g?????????X7?6?ep???0?A"
t_s?UP??%?-..?,Mk?????Z?-VE[k?bh?aq?? ?}~}?b+S
I00?_u?????UF??V?|?????\{zi???|A?a
1??
?Z\?3wU?M????{
??q)?g??\917??)??_X?[4?%??ğ?V?0??????]qby?2;?rb~Q?]??' +?SS?g?-D?
??~????w?2??X?mvX/-wV?W?dyp??
????~a>pS????c??u?K????q??"tM?WgV?}g????>?b?j?jPw?pl??-?????b?K
3K??>H?jK?Kf?N,??ks??k,?Sp???;??-
?????/?\=?}??Zu?t?
?l???Z???Mh?s??W?w????`???U?pv
M??T;z)?$?m?g?cz???
?"'?????S8[?_??wq??[;s?{??[?m~?X?ujm#??? ?à??
.l????t?r'l'[??P$e?1?"?p??0?i}?;
m5?qxLRh?!Gb?@?qD=?0????????z??????/?n???S?:?????H????D???/XG2u?ð????=c????ó??;3?????IK???k?
!G}????!0?\??}(Pz?a?P?????)CQ???C????u?P?n?s?>T??>θ9?:t?d?T?{?!G)K??C?:~??P?????s??m??a?????of???I?
JD? ???Pè??YM:??9?/?èidSN?$
??p?????F???e????z?G???r?R?^*?{?f??(??~]??5z<N*??C~??i?L-???????S??2N^m??;9?-0?t?v4j???s%"
_??0\????A????^????A??[?]σ?w??+?U~D?'l??? .Rk?_????Y}?r????JQ?k?\{f?\ojZ??F??+5C?5_? ?i?}?6{f?_??jZ=
G"K?51?X
?1?????zj?4?<??zd8X???{/
?W?杖????[?]??{Y? X{?:wLHXt?ek_XwL
??f????t????J
8Vb*öxA?y???Z?[?]Fb?z?^??????q?? .^?{?ZM3{4],?z%?V3?i?"iF?I????~Bo??(???N+???%?XAn?"?
?*?c?l@_?[?].?e?[+?p???f?YF?z?T6k4]3?5?X,?t?W,=?V?B?.$BW9?( ?B??)?<E?? ?%?G?y?^??+????~V?K????
???%tU"???~?GP6G?C"pTUR?g???J%r?J:~1'YAB?BB???1?o?????DH[?r-Q?????B*速??F ???
```

Bigwig file



for visualizing continuous data, e.g. in the UCSC Genome Browser or IGV, bigWig files come in really handy




remember that there are 2 deepTools for bam → bigWig conversion:

- ❖ ***bamCoverage**: for individual files (like those shown here)*
- ❖ ***bamCompare**: to normalize two files to each other*

BAM file 을 시각화 하기 위해 변환한 **data** 형식으로 마찬가지로 **binary** 형태 – 인간이 이해 불가.

JBrowse 2

JBrowse 2: a modular genome browser with views of synteny and structural variation

Colin Diesh, [Garrett J Stevens](#), [Peter Xie](#), [Teresa De Jesus Martinez](#), [Elliot A. Hershberg](#), [Angel Leung](#), [Emma Guo](#), [Shihab Dider](#), [Junjun Zhang](#), [Caroline Bridge](#), [Gregory Hogue](#), [Andrew Duncan](#), [Matthew Morgan](#), [Tia Flores](#), [Benjamin N. Bimber](#), [Robin Haw](#), [Scott Cain](#), [Robert M. Buels](#), [Lincoln D. Stein](#) & [Ian H. Holmes](#) 

Genome Biology 24, Article number: 74 (2023) | [Cite this article](#)

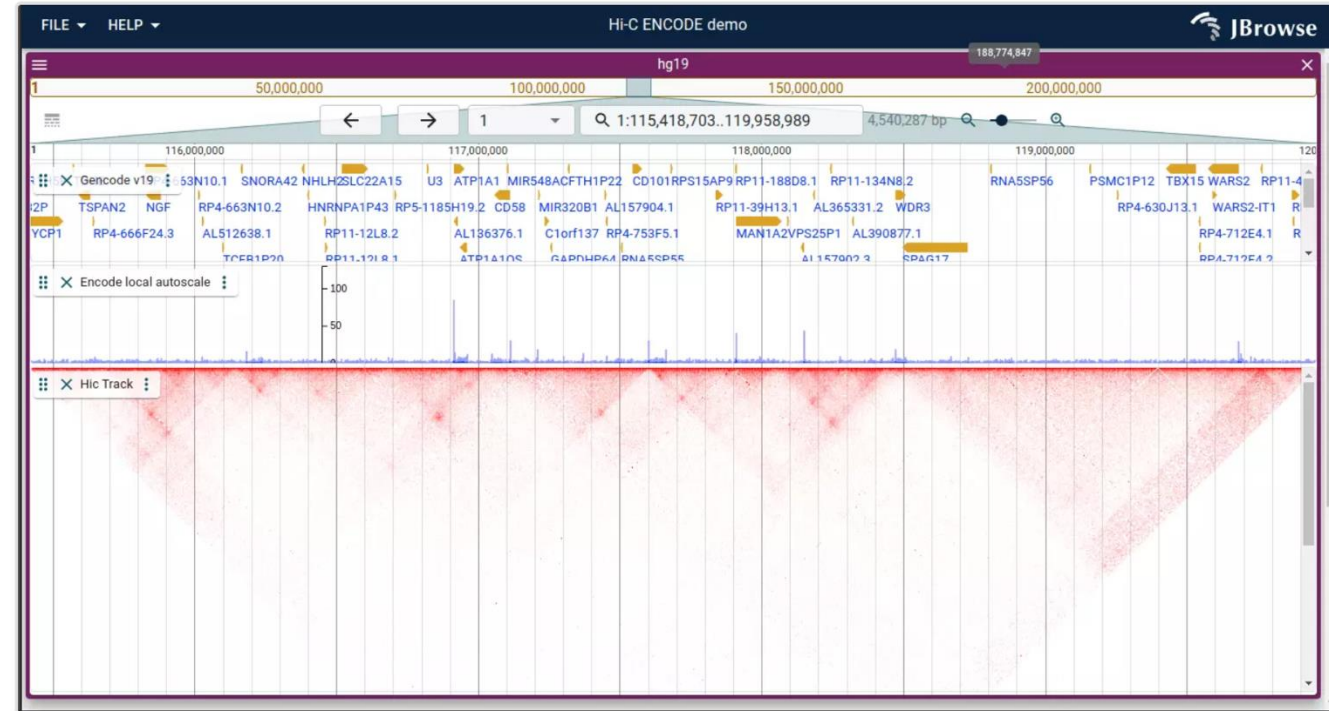
11k Accesses | 393 Citations | 17 Altmetric | [Metrics](#)

Abstract

We present JBrowse 2, a general-purpose genome annotation browser offering enhanced visualization of complex structural variation and evolutionary relationships. It retains core features of JBrowse while adding new views for synteny, dotplots, breakpoints, gene fusions, and whole-genome overviews. It allows users to share sessions, open multiple genomes, and navigate between views. It can be embedded in a web page, used as a standalone application, or run from Jupyter notebooks or R sessions. These improvements are enabled by a ground-up redesign using modern web technology. We describe application functionality, use cases, performance benchmarks, and implementation notes for web administrators and developers.

Download:

<https://jbrowse.org/jb2/>



Load FASTA file on IGV

The screenshot displays the IGV (Integrative Genomics Viewer) interface. The main window shows a genomic track with a scale from 0 to 1,350 kb. A menu is open, showing options to load a genome from a file, URL, or hosted genome. The file explorer on the right shows a directory structure with a file named 'Paralichthys_olivaceus-GCA_024713975.2-softmasked.fa.gz' selected. The file details on the right indicate it is a gzip compressed archive, created on December 22, 2024, and contains a single file named 'Paralichthys_olivaceus-GCA_024713975.2-softmasked.fa.gz'.

Load FASTA file on IGV

IGV 2.16.2

File Genomes View Tracks Regions Tools Amazon Help

Paralichthys_olivaceus-GCA_024713975.2-softmasked.fa

All

Go

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24

contig1:500,355

87M of 268M

genome

star_salmon

Paralichthys_olivaceus-GCA_024713975.2-softmasked.fa

Paralichthys_olivaceus-GCA_024713975.2-softmasked.fa.fai

Paralichthys_olivaceus-GCA_024713975.2-softmasked.fa.gz

Jun SSD

Scratch

2TB SSD

MG-4f4f4

Shared

Blue

Red

Orange

Purple

보라색

보라색_Confirmed vector

주황색

빨간색

옅은색

서식_맞춰야함

All Tags...

Paralichthys_olivaceus-GCA_024713975.2-softmasked.fa

Paralichthys_olivaceus-GCA_024713975.2-softmasked.fa.fai

Paralichthys_olivaceus-GCA_024713975.2-softmasked.fa.gz

Information

Created December 22, 2024

Modified Dec 22, 2024

Where from

https://ftp.ncbi.nlm.nih.gov/pub/rapid-release/Paralichthys_olivaceus/GCA_024713975.2-softmasked.fa.gz

Tags

Add Tags...

.fai file 만들어짐.

Load FASTA File and fai file on JBrowse

JBrowse 2

File

Edit

View

Window

Help

JBrowse

JBrowse

v2.18.0

Launch new session

Select a sequence file e.g. FASTA file

OPEN SEQUENCE FILE(S)

-or-

Quickstart list

Select one or more entries from this quickstart list to start your session. If more than one entry is selected, they will be combined. Q

☐ hg19 ...

☐ hg38 ...

☐ mm10 ...

GO

Recently opened sessions

☐ Show autosaves

OPEN SAVED SESSION (.JBrowse) FILE

No sessions available

Load FASTA File and fai file on JBrowse

JBrowse 2

File

Edit

View

Window

Help

JBrowse

v2.18.0

Launch new session

Select a sequence file e.g. FASTA file

OPEN SEQUENCE FILE(S)

-or-

Quickstart list

Select one or more entries from this quickstart list to start your session. If more than one entry is selected, they will be combined

hg19

...

hg38

...

mm10

...

GO

Recently opened sessions

Show autosaves

OPEN SAVED SESSION (.JBrowse) FILE

No sessions available

Open sequence(s)

×

Use this dialog to open one or more indexed FASTA files (IndexedFastaAdapter), bgzipped+indexed FASTA files (BgzipFastaAdapter), or .2bit files (TwoBitAdapter) of a genome assembly or other sequence. A plain FASTA file can also be supplied (FastaAdapter), which will be indexed on the fly.

Assembly name

The assembly name e.g. hg38

Assembly display name

(optional) A human readable display name for the assembly e.g. "Homo sapiens (hg38)"

Type

IndexedFastaAdapter

Type of adapter to use

fastaLocation

failLocation

FILE

URL

FILE

URL

CHOOSE FILE

No file chosen

CHOOSE FILE

No file chosen

SHOW ADVANCED OPTIONS

ADD ANOTHER ASSEMBLY

CANCEL

SUBMIT

이름 넣어주기

fasta file

fai file

JBrowse 2

File

Edit

View

Window

Help

JBrowse

v2.18.0

Recently opened sessions

Show autosaves

OPEN SAVED SESSION (.JBrowse) FILE

No sessions available

Open sequence(s)

×

Use this dialog to open one or more indexed FASTA files (IndexedFastaAdapter), bgzipped+indexed FASTA files (BgzipFastaAdapter), or .2bit files (TwoBitAdapter) of a genome assembly or other sequence. A plain FASTA file can also be supplied (FastaAdapter), which will be indexed on the fly.

Assembly name

The assembly name e.g. hg38

Assembly display name

(optional) A human readable display name for the assembly e.g. "Homo sapiens (hg38)"

Type

IndexedFastaAdapter

Type of adapter to use

fastaLocation

failLocation

FILE

URL

FILE

URL

CHOOSE FILE

No file chosen

CHOOSE FILE

No file chosen

SHOW ADVANCED OPTIONS

ADD ANOTHER ASSEMBLY

CANCEL

SUBMIT

이름 넣어주기

fasta file

fai file

JBrowse 2

File

Edit

View

Window

Help

JBrowse

v2.18.0

Recently opened sessions

Show autosaves

OPEN SAVED SESSION (.JBrowse) FILE

No sessions available

Open sequence(s)

×

Use this dialog to open one or more indexed FASTA files (IndexedFastaAdapter), bgzipped+indexed FASTA files (BgzipFastaAdapter), or .2bit files (TwoBitAdapter) of a genome assembly or other sequence. A plain FASTA file can also be supplied (FastaAdapter), which will be indexed on the fly.

Assembly name

The assembly name e.g. hg38

Assembly display name

(optional) A human readable display name for the assembly e.g. "Homo sapiens (hg38)"

Type

IndexedFastaAdapter

Type of adapter to use

fastaLocation

failLocation

FILE

URL

FILE

URL

CHOOSE FILE

No file chosen

CHOOSE FILE

No file chosen

SHOW ADVANCED OPTIONS

ADD ANOTHER ASSEMBLY

CANCEL

SUBMIT

이름 넣어주기

fasta file

fai file





Launch new session

Select a sequence file e.g. FASTA file

OPEN SEQUENCE FILE(S)

-or-

Quickstart list

Select one or more entries from this quickstart list to start your session. If more than one entry is selected, they will be combined. Q

☐ hg19 ...

☐ hg38 ...

☐ mm10 ...

GO

Recently opened sessions

Show autosaves

OPEN SAVED SESSION (.JBrowse) FILE

No sessions available

Open sequence(s)

Use this dialog to open one or more indexed FASTA files (IndexedFastaAdapter), bgzipped+indexed FASTA files (BgzipFastaAdapter), or .Zbit files (TwoBitAdapter) of a genome assembly or other sequence. A plain FASTA file can also be supplied (FastaAdapter), which will be indexed on submit.

Assembly name

flounder

The assembly name e.g. hg38

Assembly display name

flounder

(optional) A human readable display name for the assembly e.g. "Homo sapiens (hg38)"

Type

IndexedFastaAdapter

Type of adapter to use

fastaLocation

FILE

URL

CHOOSE

FILE

/Volumes/Jun_SSD/genome/Paralichthys_olivaceus/Ensembl_rapid/genome/P

GCA_024713975.2-softmasked.fa

fastaLocation

FILE

URL

CHOOSE

FILE

/Volumes/Jun_SSD/genome/Paralichthys_olivaceus/Ensembl_rapid/genome/P

GCA_024713975.2-softmasked.fa.fai

SHOW ADVANCED OPTIONS

ADD ANOTHER ASSEMBLY

CANCEL

SUBMIT



Open sequence(s)

Select a view to launch

Linear genome view ▾

Linear genome view

Circular view

Dotplot view

Linear synteny view

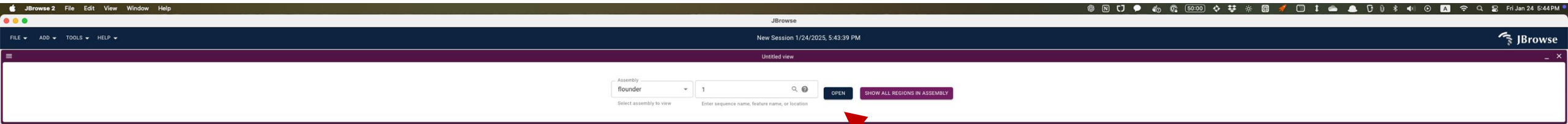
Spreadsheet view

SV Inspector

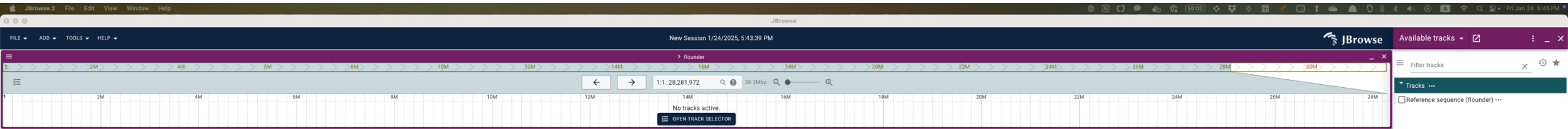
Breakpoint split view



linear genome view

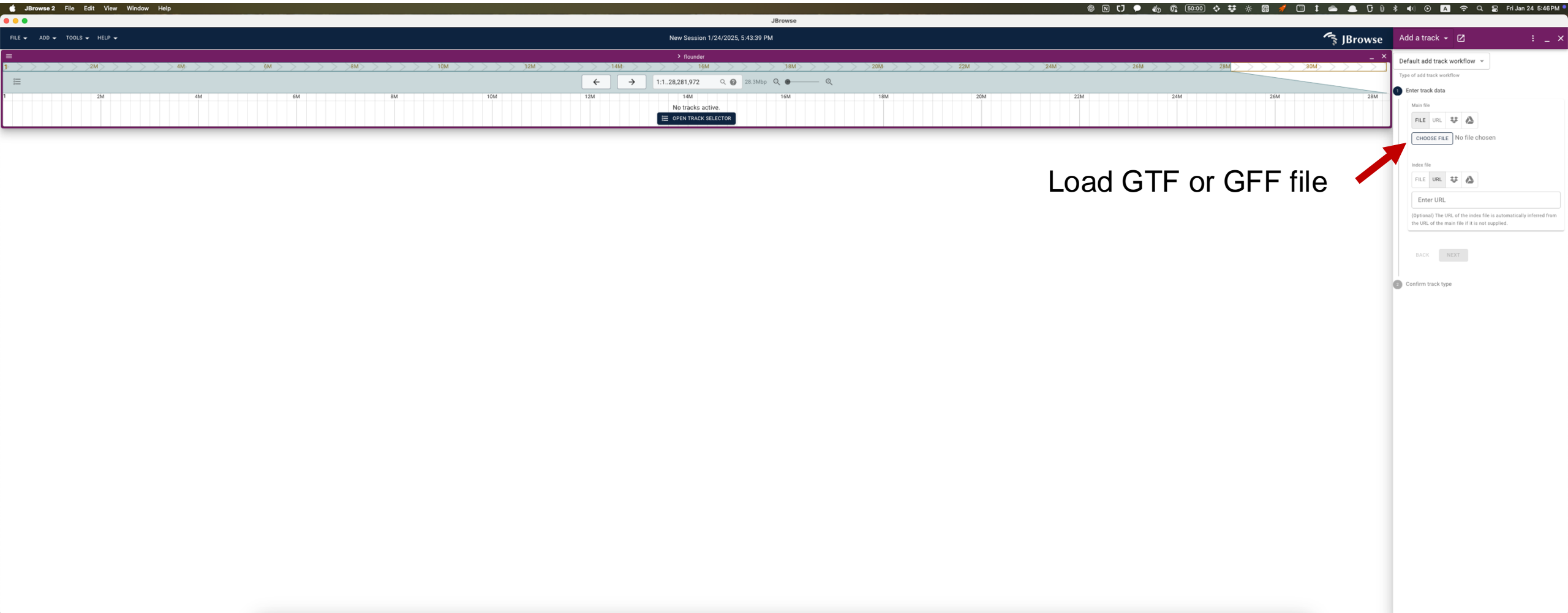


Open

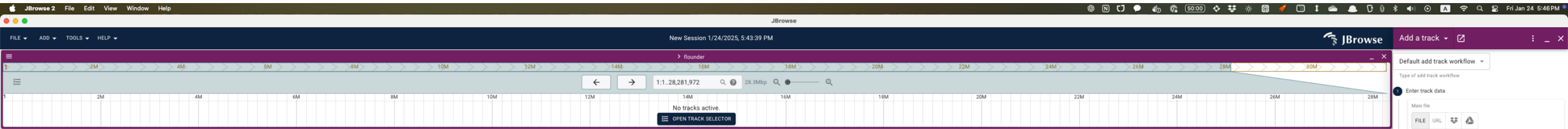


Click



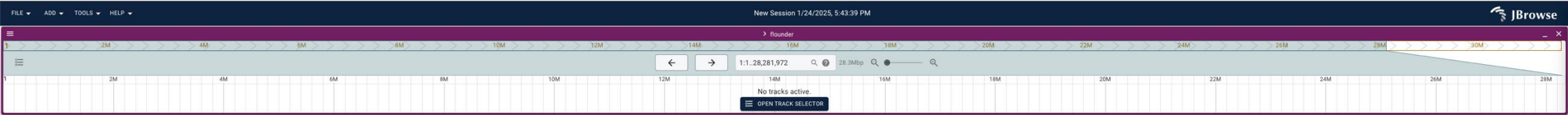


Load GTF or GFF file



Load GTF or GFF file

Indexing is option.
Load 안하면 프로그램이
자동으로 indexing 진행.



Add a track ▾

Default add track workflow ▾

Type of add track workflow

Enter track data

Confirm track type

Using adapter GtfAdapter and guessing track type FeatureTrack. Please enter a track name and, if necessary, update the track type.

trackName

Paralichthys_olivaceus-GCA_024713975.2-2023_02-gen

A name for this track

Adapter type

GTF adapter

Select an adapter type

Track type

Feature track

Select track type

Assembly

flounder

Select an assembly to add track to

BACK ADD

Click

