

# Paleofish Mitochondries

Objectif : Lancement de Eager sur référence majoritaire

# CR réunion du 27 novembre 2025

- Pour les 5 échantillons not salmo, travailler sur les espèces en génome de référence et non les sous espèces (même pour Thymallus): *Anguilla anguilla*, *Thymallus thymallus*, *Coregonus clupeaformis*, *Brachymystax lenok* --> Eager:

```
smaman@genobioinfo1 /work/project/crucial/PALEOFISH/ARCHIVES/archive_02_eager_44_mito_genomes $ ls
```

```
1-bwa_sur_ref_majoritaire.sh NC_001960.1 NC_006531.1 NC_010007.1 NC_012928.1 NC_018341.1 NC_024032.1 NC_025648.1 OLD REF_MAJORITAIRE scripts-archives
```

--> Joelle: voir multiQC, comparaison consensus double pass Eager sur espèces et sous-espèces.

- Pas besoin d'IGV.js pour la visualisation des consensus.

- Sarah : comparer consensus mpileup avec consensus double pass sur espèces.

Autres points abordés lors d'une précédente réunion :

\* Nettoyer les fichiers undetermined pour les traiter avec Eager (pour fev/mars 26).

\* Point projet : <https://annuel.framapad.org/p/pnqb6r9hsu-ai09?lang=fr>

Comparer consensus mpileup avec consensus double pass sur espèces.

# Visualisation consensus mpileup

```
smaman@genobioinfo1 /work/project/crucial/PALEOFISH/Sarah/1-pipeline_fastpDefault_bwa/CONSENSUS/consensus $ ls *fa
L0_1_EG4_S81_cns.fa      L10_5_BC9_S62_cns.fa    L2_7_TP2_S76_cns.fa    L4_13_MB5_S23_cns.fa    L5_13_BC11_S38_cns.fa    L6_4_SH1979-4-6843-c_S42_cns.fa
L0_2_EG5_S82_cns.fa      L11_1 SCO_98_S66.cns.fa  L2_8_BC1_S77_cns.fa    L4_14_SM2_S24_cns.fa    L5_14_BC12_S39_cns.fa    L6_5_SH1979-4-6843-d_S43_cns.fa
Terminal L11_2 SCO_81_S67.cns.fa  L3_10_MB13_S5_cns.fa  L4_15_BC10_S25_cns.fa  L5_15_BC13_S40_cns.fa    L6_6_SH2001-106_11-b_S44_cns.fa
L0_4_EG9_S84_cns.fa      L11_3 SCO_114_S68_cns.fa  L3_11_MB18_S6_cns.fa  L4_1_MB14_S11_cns.fa    L5_1_MD15_S26_cns.fa    L6_8_SH2000-99-394-e_S45_cns.fa
L0_5_EG10_S85_cns.fa     L11_6 SCO_116_S69_cns.fa  L3_12_MB6_S7_cns.fa    L4_2_MD13_S12_cns.fa    L5_2_MD7_S27_cns.fa    L7_14_22d_S48_cns.fa
L0_6_EG13_S86_cns.fa     L11_7 SCO_1354_S70_cns.fa  L3_13_MB11_S8_cns.fa  L4_3_MB16_S13_cns.fa    L5_3_MD8_S28_cns.fa    L8_15_HTMK99-XXII-0-4287_S49_cns.fa
L0_7_OLG1_S87_cns.fa     L2_10_MPMB2_S78_cns.fa   L3_14_MB10_S9_cns.fa  L4_4_MB9_S14_cns.fa    L5_4_MD4_S29_cns.fa    L9_1_26-52_S50_cns.fa
L0_8_SH2_S88_cns.fa       L2_11_MPMB3_S79_cns.fa   L3_15_MB12_S10_cns.fa  L4_5_MD5_S15_cns.fa    L5_5_MD10_S30_cns.fa    L9_13_SH2000-99-390-b_S55_cns.fa
L10_11 SCO_1338_S63_cns.fa L2_12_MPMB4_S80_cns.fa  L3_1_MB15_S1_cns.fa   L4_6_MB2_S16_cns.fa    L5_6_MB8_S31_cns.fa    L9_14_SH1979-4-6843-a_S56_cns.fa
L10_12 SCO_2213_S64_cns.fa L2_1_EG1_S71_cns.fa    L3_5_SM3_S2_cns.fa   L4_7_MB1_S17_cns.fa    L5_7_MD3_S32_cns.fa    L9_15_SH2000-99-390-c_S57_cns.fa
L10_13 SCO_2764_S65_cns.fa L2_2_EG2_S72_cns.fa   L3_6_MB3_S3_cns.fa   L4_8_MB4_S18_cns.fa    L5_8_MD16_S33_cns.fa    L9_2_2777-17_S51_cns.fa
L10_1_BC5_S58_cns.fa      L2_3_EG3_S73_cns.fa   L3_9_MB17_S4_cns.fa   L4_9_MD12_S19_cns.fa    L5_9_MD6_S34_cns.fa    L9_4_SM1_S52_cns.fa
L10_2_BC6_S59_cns.fa      L2_4_CD2_S74_NC_009263_cns.fa L4_10_BC14_S20_cns.fa  L5_10_MD11_S35_cns.fa    L6_11_SH2000-99-394-g_S46_cns.fa  L9_6_SH1979-4-6843-b_S53_cns.fa
L10_3_BC7_S60_cns.fa       L2_4_CD2_S74_NC_018341_cns.fa L4_11_MD2_S21_cns.fa  L5_11_MD14_S36_cns.fa    L6_13_SH2000-99-394-f_S47_cns.fa  L9_7_SH2000-99-394-c_S54_cns.fa
L10_4_BC8_S61_cns.fa       L2_6_TP1_S75_cns.fa   L4_12_MD1_S22_cns.fa  L5_12_MD9_S37_cns.fa    L6_3_SH1979-4-6843-e_S41_cns.fa
```

Réflexions autour du L11\_7

# Visualisation consensus mpileup

## L11 7 Anguilla

/save/sigenae/public\_html/sarah/mito/consensus\_mpileup/.

```
smaman@genobioinfo1 /save/sigenae/public_html/sarah/mito/consensus_mpileup $ ls ..../consensus_not_salmo  
Anguilla Brachymystax_leinox_tsinlingensis Coregonus Thymallus_ligericus  
smaman@genobioinfo1 /save/sigenae/public_html/sarah/mito/consensus_mpileup $ ls ..../consensus_not_salmo/Anguilla  
CM077320.2.fna      GCA_039654925.2_AngPac_1_genomic.fna      SCO_1354.realign.bai  SCO_1354.unifiedgenotyper.vcf.gz  
CM077320.2.fna.fai  igvjs_SCO_1354_Anguilla_bicolor_pacifica.html  SCO_1354.realign.bam
```

```
smaman@genobioinfo2 /work/project/crucial/PALEOFISH/comparison_mpileup_eager $ ls /work/project/crucial/PALEOFISH/ARCHIVES/archive_02_eager_44_mito_genomes/NC_006531.1/03_eager_only_g  
enotyping/results/consensus_sequence/  
SCO_1354.fasta.gz  SCO_1354.fasta_refmod.fasta.gz  SCO_1354.fasta_uncertainty.fasta.gz
```

# Consensus Eager

## L11 7 Anguilla

```
smaman@genobioinfo2 /work/project/crucial/PALEOFISH/comparison_mpileup_eager $ ls /work/project/crucial/PALEOFISH/ARCHIVES/archive_02_eager_44_mito_genomes/NC_006531.1/03_eager_only_g  
enotyping/results/consensus_sequence/  
SCO_1354.fasta.gz  SCO_1354.fasta_refmod.fasta.gz  SCO_1354.fasta_uncertainty.fasta.gz
```

# Comparaison des consensus Eager/mpileup L11 7 Anguilla

## mpileup / Eager => Comparaison des IGV ?

```
smanan@genobioinfo2:/work/project/crucial/PALEOFISH/comparaison  
4_S70_cns.fa  
>NC_006531.1 Anguilla anguilla mitochondrial, complete genome  
GTTAACGTAGCTTAAACAAAAAGCATGGCACTGAAGATGCCAAGATGAGGCCATAAAAAGC  
TCCGTGACACAAAAGCTGGTCTGACTTAACATCAGTCTGGCTGACTTACACATG  
CAAGTACCCGGCACCGTGAGAATGCCCTATATCCCCTCCGGGAAAGGGCCGGCA  
TCAGGCACACCAACGTAGCCCCAAACACCTTGCTTAACCACACCCCCAAGGGATTCAAC  
AGTGTAGACATTGAGCAATAAGCGCAAGCTTGACTTAGCTCAAGGCCAAAAGAGTTGGTT  
AATCTCTGGCAGCCACCGGGTTACAGAGTAACTCACATTAACTTCACGGCTAAA
```

# Nettoyer les fichiers undetermined pour les traiter avec Eager (pour fev/mars 26).

## CR 11 sept 2025

Voici les points abordés:

Présentation des résultats et des slides d'avancement envoyés par mail ce jour au groupe de travail: multimapping sur 44 références mito, répartition des jeux de données capture\_janv22 par lot, pipelines Eager sur chaque lot, chemin d'accès aux résultats, exclusion des fichiers Sxx issus des fastq undetermined.

Lancement FastQScreen pour "coller" au pipeline Paleotrutta et pour se confronter sur les résultats du multimapping, sans a priori.

Il a été décidé:

Ne pas traiter les données shotgun car elles ne peuvent pas être exploitées par Joelle. Traiter uniquement les données capture. @Joelle : Confirmer le chemin d'accès et la liste des données à traiter : /work/project/crucial/PALEOFISH/DATA/capture\_jan22

Ne pas traiter les jeux de données récupérés avec les barcodes depuis les fichiers Undetermined car les fichiers contiennent des reads non pairés. Ces fichiers \*Sxx\* sont donc archivés dans un nouveau répertoire Undetermined/ : smaman@genobioinfo2 /work/project/crucial/PALEOFISH/DATA/capture\_jan22 \$ ls Undetermined/

L1\_1\_AUD20212-2\_Sxx\_R1\_001\_GACGATT+TCGCAGG.fastq L1\_2\_BC4\_Sxx\_R2\_001\_AACCTGC+CTCTGCA.fastq  
L1\_7\_BC2\_Sxx\_R1\_001\_GTCCGGC+CTCGATG.fastq Undetermined\_S0\_R2\_001.fastq.gz

L1\_1\_AUD20212-2\_Sxx\_R2\_001\_GACGATT+TCGCAGG.fastq L1\_4\_AUD11764-47\_Sxx\_R1\_001\_GCCTACG+GGATCAA.fastq  
L1\_7\_BC2\_Sxx\_R2\_001\_GTCCGGC+CTCGATG.fastq Undetermined\_S0\_R1\_001.fastq.gz

L1\_2\_BC4\_Sxx\_R1\_001\_AACCTGC+CTCTGCA.fastq L1\_4\_AUD11764-47\_Sxx\_R2\_001\_GCCTACG+GGATCAA.fastq Undetermined\_S0\_R1\_001.fastq.gz

Noms de paires de reads qui ne matchent pas => Retrait des \*Sxx\* du lot traité.

Command error:

```
Trimming paired end reads ...
Opening FASTQ file 'repaired_L1_2_BC4_Sxx_R1_001_AACCTGC+CTCTGCA.fastq.pG.fq.gz', line numbers start at 1
Opening FASTQ file 'repaired_L1_2_BC4_Sxx_R2_001_AACCTGC+CTCTGCA.fastq.pG.fq.gz', line numbers start at 1
ERROR: Unhandled exception in thread:
      Pair contains reads with mismatching names:
      - 'NB501044:604:HHNTFAFX3:1:11101:6394:7897'
      - 'NB501044:604:HHNTFAFX3:1:11101:2238:8023'
```

Note that AdapterRemoval by determines the mate numbers as the digit found at the end of the read name, if this is preceded by the character '/'; if these data makes use of a different character to separate the mate number from the read name, then you will need to set the --mate-separator command-line option to the appropriate character.

ERROR: AdapterRemoval did not run to completion;  
do NOT make use of resulting trimmed reads!

# Fastq Sxx reconstruits à partir des fichiers Undetermined

Séquences non reconnues par Eager :

```
/work/project/crucial/PALEOFISH/02_eager/NC_001960.1$ seqtk seq -A L1_1_AUD20212-2_Sxx_R1_001_GACGATT+TCGCAGG.fastq | head  
>NB501044:604:HHNTFAFX3:1:11101:1980:1042 1:N:0:GACGATT+TCGCAGG  
NATTTCTGTAGACAACGACACCTAACACGATTTCGCCTTCCACTTCCATTCCCATTGCTTATTGAGCTGC  
>NB501044:604:HHNTFAFX3:1:11101:20400:1043 1:N:0:GACGATT+TCGCAGG  
NAGGCTCTGGCTTAGTGTCTACCTAACGCCCTGTTATAAGAGATCGGAAGAGCACACGCTGAACCTCAGTCAC  
>NB501044:604:HHNTFAFX3:1:11101:5386:1044 1:N:0:GACGATT+TCGCAGG  
NCTTATAATTCAAGTAGCCCCAACTATCAACTCTTCTACTCATTGGATAGGCCCTTATCAATACTTGTAGGA  
>NB501044:604:HHNTFAFX3:1:11101:19272:1045 1:N:0:GACGATT+TCGCAGG  
AGATCCCCCGCTTCCGCGCGAAACAGATCGGAAGAGCACACGCTGAACCTCAGTCACGACGATTATCGCG  
>NB501044:604:HHNTFAFX3:1:11101:6772:1045 1:N:0:GACGATT+TCGCAGG  
TATAATTAAAGCTCTTCGCTAGTAGATCTCGTGCAGATCGGAAGAGCACACGCTGAACCTCAGTCACGACGA  
/work/project/crucial/PALEOFISH/02_eager/NC_001960.1$ file L1_1_AUD20212-2_Sxx_R1_001_GACGATT+TCGCAGG.fastq  
L1_1_AUD20212-2_Sxx_R1_001_GACGATT+TCGCAGG.fastq: ASCII text  
/work/project/crucial/PALEOFISH/02_eager/NC_001960.1$ wc -l L1_1_AUD20212-2_Sxx_R1_001_GACGATT+TCGCAGG.fastq  
12785788 L1_1_AUD20212-2_Sxx_R1_001_GACGATT+TCGCAGG.fastq  
/work/project/crucial/PALEOFISH/02_eager/NC_001960.1$ cat -A L1_1_AUD20212-2_Sxx_R1_001_GACGATT+TCGCAGG.fastq | head  
@NB501044:604:HHNTFAFX3:1:11101:1980:1042 1:N:0:GACGATT+TCGCAGG$  
NATTTCTGTAGACAACGACACCTAACACGATTTCGCCTTCCACTTCCATTCCCATTGCTTATTGAGCTGC$  
+$  
#AAAAEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEAEAEAE$  
@NB501044:604:HHNTFAFX3:1:11101:20400:1043 1:N:0:GACGATT+TCGCAGG$  
NAGGCTCTGGCTTAGTGTCTACCTAACGCCCTGTTATAAGAGATCGGAAGAGCACACGCTGAACCTCAGTCAC$  
+$  
#AAAAEEEEEEEEEEEEEEEEEE/EEEEEEEAEAEAEAEAEAE/EEEEEEEEEE/E/EEE/E/EE$  
@NB501044:604:HHNTFAFX3:1:11101:5386:1044 1:N:0:GACGATT+TCGCAGG$  
NCTTATAATTCAAGTAGCCCCAACTATCAACTCTTCTACTCATTGGATAGGCCCTTATCAATACTTGTAGGA$
```

# Réparation des fichiers undetermined pour les traiter avec Eager

```
smaman@genobioinfo1 /work/project/crucial/PALEOFISH/DATA_UNDETERMINATED $ ls
L1_1_AUD20212-2_Sxx_R1_001_GACGATT+TCGCAGG.fastq L1_2_BC4_Sxx_R1_cleaned.fastq L1_4_AUD11764-47_Sxx_R2_001_GCCTACG+GGATCAA.fastq L1_7_BC2_Sxx_R2_cleaned.fastq
L1_1_AUD20212-2_Sxx_R1_cleaned.fastq L1_2_BC4_Sxx_R2_001_AACCTGC+CTCTGCA.fastq L1_4_AUD11764-47_Sxx_R2_cleaned.fastq long_reads.sh
L1_1_AUD20212-2_Sxx_R2_001_GACGATT+TCGCAGG.fastq L1_2_BC4_Sxx_R2_cleaned.fastq L1_7_BC2_Sxx_R1_001_GTCCGGC+CTCGATG.fastq README
L1_1_AUD20212-2_Sxx_R2_cleaned.fastq L1_4_AUD11764-47_Sxx_R1_001_GCCTACG+GGATCAA.fastq L1_7_BC2_Sxx_R1_cleaned.fastq
L1_2_BC4_Sxx_R1_001_AACCTGC+CTCTGCA.fastq L1_4_AUD11764-47_Sxx_R1_cleaned.fastq L1_7_BC2_Sxx_R2_001_GTCCGGC+CTCGATG.fastq
```

## Résultats de l'alignement et de la couverture :

```
/work/project/crucial/PALEOFISH/02_eager/UNDETERMINED_ref_majoritaire$ more STATS/*align
::::::::::::::::::
STATS/L1_1_AUD20212-2_Sxx_GACGATT+TCGCAGG.align
::::::::::::::::::
NC_024032.1      154514
NC_010007.1      161407
NC_001960.1      4690381
::::::::::::::::::
STATS/L1_2_BC4_Sxx_AACCTGC+CTCTGCA.align
::::::::::::::::::
NC_024032.1      25409
NC_010007.1      25667
NC_001960.1      1012069
::::::::::::::::::
STATS/L1_4_AUD11764-47_Sxx_GCCTACG+GGATCAA.align
::::::::::::::::::
NC_024032.1      40467
NC_010007.1      42482
NC_001960.1      1200198
::::::::::::::::::
STATS/L1_7_BC2_Sxx_GTCCGGC+CTCGATG.align
::::::::::::::::::
NC_001960.1      100882
NC_010007.1      1727220
NC_024032.1      2029130
```

# Lancement des étapes Eager sur les fichiers undetermined nettoyés

```
smaman@genobioinfo1 ~ $ cd /work/project/crucial/PALEOFISH/DATA_UNDETERMINED/EAGER
smaman@genobioinfo1 /work/project/crucial/PALEOFISH/DATA_UNDETERMINED/EAGER $ ls *
NC_001960.1:
01_script_eager_part1_NC_001960.1.sh    02ter_tsv_modif.sh      crucial_nextflow_VersionPaleofish.config  fastqc          pipeline_trace.txt  work
02bis-bam_tsv.sh                         03_eager_only_genotyping documentation           input_NC_001960.1.tsv   reference_genome
02_rehead_bam_files_v2_NC_001960.1.sh    adapterremoval        FastP                  pipeline_info      slurm-29462790.out

NC_024032.1:
01_script_eager_part1_NC_024032.1.sh    03_eager_only_genotyping documentation           input_NC_024032.1.tsv   reference_genome
02bis-bam_tsv.sh                         adapterremoval        FastP                  pipeline_info      slurm-29462832.out
02_rehead_bam_files_v2_NC_024032.1.sh    crucial_nextflow_VersionPaleofish.config  fastqc          pipeline_trace.txt  work
```

## Pipeline et paramétrages mito

### Phase 1 : Premier lancement d'Eager

tsv\_modif.sh : Génération du fichier d'entrée input.csv listant les FastQ en entrée d'Eager.  
01\_script\_eager\_part1\_NC\_001960.1.sh

### Phase 2 : Filtre des BAM pour ne conserver que les séquences mito

02\_rehead\_bam\_files\_v2\_NC\_001960.1.sh : Enlever les séquences génomiques du BAM pour ne conserver que les séquences mito. Sinon fichier FASTQ trop gros qui fait planter le script 3.

02bis-bam\_tsv.sh : Créer le fichier tabulé listant les BAM, en entrée de la phase 2 d'Eager.

### Phase 3 : Second lancement d'Eager uniquement sur les mitochondries

03\_script\_eager\_only\_genotyping\_v2.sh : Skip des étapes précédentes avec les étapes GATK et VCF.

# Lancement des étapes Eager sur les fichiers undetermined nettoyés

Nom de la librairie (Sxx) insuffisant pour identifier les échantillon lors de l'étape rehead BAM.  
=> Ajout Sxx\_NomSample pour rendre l'information unique.

```
smaman@genobioinfo1 /work/project/crucial/PALEOFISH/DATA_UNDETERMINED/EAGER/NC_001960.1/03_eager_only_genotyping $ ls consensus_sequence/
AUD11764-47.fasta.gz          AUD11764-47.fasta_uncertainty.fasta.gz  AUD20212-2.fasta_refmod.fasta.gz    BC4.fasta.gz           BC4.fasta_uncertainty.fasta.gz
AUD11764-47.fasta_refmod.fasta.gz AUD20212-2.fasta.gz            AUD20212-2.fasta_uncertainty.fasta.gz BC4.fasta_refmod.fasta.gz

smaman@genobioinfo1 /work/project/crucial/PALEOFISH/DATA_UNDETERMINED/EAGER/NC_024032.1/03_eager_only_genotyping $ ls
03_script_eager_only_genotyping_v2_NC_024032.1.sh  damage_rescaling  genotyping      multiqc
bcftools                                         documentation       input_bam.tsv   OLD_03_script_eager_only_genotyping_v2_NC_024032.1.sh reference_genome
consensus_sequence                           endorspy          mapdamage      pipeline_info
                                             samtools          work

smaman@genobioinfo1 /work/project/crucial/PALEOFISH/DATA_UNDETERMINED/EAGER/NC_024032.1/03_eager_only_genotyping $ ls consensus_sequence/
BC2.fasta.gz  BC2.fasta_refmod.fasta.gz  BC2.fasta_uncertainty.fasta.gz

smaman@genobioinfo1 /work/project/crucial/PALEOFISH/DATA_UNDETERMINED/EAGER/NC_024032.1/03_eager_only_genotyping $ [
```

# Point projet

<https://annuel.framapad.org/p/pnqb6r9hsu-ai09?lang=fr>