

# Release coordination documentation

- Declaration of Intentions
- Environment variables
- Update your checkouts
- NCBI taxonomy data
- Compara servers
- Configuration file
- Update schema version
- Update table.sql and create patch files
  - Check the patch files
- Run CheckTaxon healthcheck
- Add new genome\_db to compara master database
- Add method\_link\_species\_set entries to compara master database
- Add new species to phylogenetic tree
- Final checks to compara master database
- Create Release Database
- Merge DNA data
- Merge GeneTrees+Families+NCTrees+PatchProjectionsAsHomologies
- Final database checks
- Run the healthchecks
- Test web server
- Run ANALYZE\_TABLE and OPTIMIZE TABLE
- Copy databases to staging servers
- Final handover of databases
- Update documentation and diagrams
- Data dumps
- Final things

## Declaration of Intentions

- ☒ Set up a web page

Once the release coordinator has sent out the email for the declaration of intentions, set up a web page with intentions in the Confluence wiki system to allow easy tracking of the progress.  
Release plans

- ☒ Ask compara team members of their intentions  
Compara has one extra day to declare their intentions because of the need to know what the genebuilders and associated teams (eg wormbase, ensembl genomes) will declare
- ☒ Submit the declaration of intentions on the <http://admin.ensembl.org/index.html> website.

## Environment variables

- ☒ Define \$ENSEMBL\_CVS\_ROOT\_DIR  
This is necessary to run the Hive and is used by many scripts/files in this document. Make sure this is defined in your terminal
- ☒ Define \$ENSADMIN\_PSW  
The password for the mysql 'ensadmin' user also needed for many scripts
- ☐ Define \$COMPARA\_REG variable to simplify connecting to databases via registry  
export COMPARA\_REG="-reg\_conf \$ENSEMBL\_CVS\_ROOT\_DIR/ensembl-compara/scripts/pipeline/production\_reg\_conf.pl -reg\_type compara -reg\_alias"

## Update your checkouts

Ensure you have up-to-date git checkouts of at least the following repositories, pointing at master branch:

- ☐ ensembl-compara
- ☐ ensembl
- ☐ ensembl-hive
- ☐ ensembl-analysis
- ☐ ensj-healthcheck

## NCBI taxonomy data

The production team updates the ncbi\_taxonomy database on livemirror just before the handover to us (please check that this has been done).

We then need to update the tables on our master DB. The current (rel.75) master database is sf5\_ensembl\_compara\_master on compara1

- ☐ Update the ncbi\_taxa\_node and ncbi\_taxa\_name in the master database

▼ [Click here for details](#)

The ncbi\_taxonomy database is located in [mysql://ens-livemirror:3306/ncbi\\_taxonomy](#)

### mysqldump

```
time eval mysqldump --extended-insert --compress --delayed-insert `db_cmd.pl
$COMPARA_REG ncbi_taxonomy -to_params` ncbi_taxa_node ncbi_taxa_name |
db_cmd.pl $COMPARA_REG compara_master
```

### ▼ Times

rel.64: 45 sec

rel.65: 47 sec

rel.66: 47 sec

rel.67: 30 sec

rel.68: 30 sec

rel.69: 35 sec

rel.70 32 sec

rel.71 34 sec

rel.72 36 sec

rel.74 36 sec

rel.75 38 sec

rel.76 44 sec

- ☐ Check new extant taxon names

Each new species must have a 'ensembl alias name' tag in the \$ENSEMBL\_CVS\_ROOT\_DIR/ensembl-compara/scripts/taxonomy/ensembl\_aliases.sql file. It should match the "web\_name" used by the production team in the "species" table of their "ensembl\_production" database on staging1. This may have already been added by Ensembl Production.

### check species/taxa from ensembl\_production

```
db_cmd.pl $COMPARA_REG ensembl_production -sql 'select production_name, web_name,
taxon from species'
```

- ☐ Update the ancestral taxon names

Each new extant species is anchored to the species tree at a certain taxon. This taxon must be described with two fields in the ncbi\_taxa\_name table:

- 'ensembl alias name': a "simple English" description of the taxon.

The script in scripts/taxonomy/place\_species.pl helps in placing the new species in the current compara species tree and discover the new extant species:

### check species/taxa from ensembl\_production

```
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/taxonomy/place_species.pl  
-master_url mysql://ensro@compara1/sf5_ensembl_compara_master -taxon_ids  
9940,9361,7994,7918
```

- 'ensembl timetree mya': the age of the taxon. It can be obtained from the TimeTree database (<http://www.timetree.org>)

For any new species, update the file \$ENSEMBL\_CVS\_ROOT\_DIR/ensembl-compara/scripts/taxonomy/ensembl\_aliases.sql to add the two new tags

- ☐ Load ensembl\_aliases.sql onto the master database

▼ [Click here for details](#)

The script will report any discrepancies that need to be resolved ie any nodes which have been deleted from the ncbi\_taxonomy database but still have entries in the ensembl\_aliases.sql file. Check if these have an entry in the species\_set\_tag table. If not, it is probably safe to delete them. Check with other compara team members.

### load ensembl\_aliases

```
db_cmd.pl $COMPARA_REG compara_master <  
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/taxonomy/ensembl_aliases.sql
```

## Compara servers

Check out the current space on the compara servers and delete the last but one release. Leave the previous release for healthchecks. Check with the other compara team members before deleting.

- ☐ Check space on [http://www.ebi.ac.uk/~mp/compara\\_servers\\_disk\\_load.html](http://www.ebi.ac.uk/~mp/compara_servers_disk_load.html)
- ☐ Ask compara team members to tidy up any unwanted databases (run the command below for all compara servers) and inform them of the intention to delete the last but one release

### how much space do databases take?

```
perl $ENSEMBL_CVS_ROOT_DIR/ensembl/misc-scripts/db/db-space.pl -host compara4  
-port 3306 -user ensadmin -password $ENSADMIN_PSW
```

- ☐ Delete ???\_ensembl\_compara\_xx
- ☐ Delete ???\_ensembl\_compara\_ancestral\_xx

## Configuration file

- ☐ Update production\_reg\_conf.pl and check back into git:

▼ [Click for details](#)

Update the registry configuration file \$ENSEMBL\_CVS\_ROOT\_DIR/ensembl-compara/scripts/pipeline/production\_reg\_conf.pl that will be used throughout the release process.

Make sure to have edited the release numbers, added external core databases and fixed name prefixes.

The convention right now (since rel.66) is to have the release database in compara3.

## Update schema version

- ☐ Update the schema\_version in the master database

```
update meta set meta_value = XX where meta_key = 'schema_version';
```

## Update table.sql and create patch files

Update the \$ENSEMBL\_CVS\_ROOT\_DIR/ensembl-compara/table.sql file and create any patch files.

- ☐ Create a patch file for the schema\_version
- ☐ Check if any other patch files need creating by looking at the Declaration of Intentions and checking with other compara team members
- ☐ Update the schema\_version in table.sql
- ☐ Delete the previous patch INSERT statements from table.sql
- ☐ Add an INSERT statement for the new schema\_version in table.sql and for any other new patches

## Check the patch files

- ☐ Create a new database from the current schema
  - ✓ [Click here for details](#)

### create new database

```
db_cmd.pl -url "mysql://ensadmin:${ENSADMIN_PSW}@compara3/" -sql 'CREATE
DATABASE mp12_current_schema_test'
db_cmd.pl -url
"mysql://ensadmin:${ENSADMIN_PSW}@compara3/mp12_current_schema_test" <
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/sql/table.sql
mysqldump -u ensro -h compara3 -P3306 --no-data --skip-add-drop-table
mp12_current_schema_test | sed 's/AUTO_INCREMENT=[0-9]*\b//' >new_schema.sql
```

- ☐ Generate an empty database from the old schema, populate the meta table, apply the patches, dump it, and check that you get the new schema.

✓ [Click for details](#)

We need to populate the meta table from the previous release to allow the schema\_patcher.pl script to work. The final sdiff will display the peptide\_align\_feature\_xxx tables that are in the previous release. These can be ignored.

### Check patches

```
mysqldump -u ensro -h compara5 -P3306 --no-data --skip-add-drop-table
lg4_ensembl_compara_75 | sed 's/AUTO_INCREMENT=[0-9]*\b//' > old_schema.sql
db_cmd.pl -url
"mysql://ensadmin:${ENSADMIN_PSW}@compara5/lg4_schema_patch_test" -sql 'CREATE
DATABASE'
db_cmd.pl -url
"mysql://ensadmin:${ENSADMIN_PSW}@compara5/lg4_schema_patch_test" <
old_schema.sql
mysqldump --defaults-group-suffix=_compara5 lg4_ensembl_compara_75 meta |
db_cmd.pl -url
"mysql://ensadmin:${ENSADMIN_PSW}@compara5/lg4_schema_patch_test"
$ENSEMBL_CVS_ROOT_DIR/ensembl/misc-scripts/schema_patcher.pl --host compara5
--port 3306 --user ensadmin --pass $ENSADMIN_PSW --database
lg4_schema_patch_test --type compara --from 74 --release 75 --verbose
mysqldump --defaults-group-suffix=_compara5 --no-data --skip-add-drop-table
lg4_schema_patch_test | sed 's/AUTO_INCREMENT=[0-9]*\b//'
>patched_old_schema.sql
sdiff -w 200 -bs patched_old_schema.sql new_schema.sql | less
```

☐ git commit table.sql and any patch files

 After Handover

## Run CheckTaxon healthcheck

☐ Run the CheckTaxon healthcheck early to find any discrepancies between the ncbi\_taxon\_name table and the core databases (information about how to set up the healthchecks can be found [here](#))

### Run healthcheck

```
# make sure you are using the right version of JAVA:
export JAVA_HOME=/software/jdk1.6.0_14

# if you need to recompile (submit to the farm, because you need more memory than is
available on the head) :
bsub -I ant clean jar

# run the healthchecks (submit to the far(submit to the farm, because you need more
memory than is available on the head) :m, because you need more memory than is
available on the head) :
time bsub -I ./run-configurable-testrunner.sh -h compara1 -d
sf5_ensembl_compara_master -t org.ensembl.healthcheck.testcase.compara.CheckTaxon
```

We can use the compara master database as the source, before the creation of the release database.

(see xxxx on how to set run the healthchecks)

## Add new genome\_db to compara master database

The current master database (rel.75) is called sf5\_ensembl\_compara\_master on compara1. You have to create new genome\_dbs and dnafrags when there is a new assembly or a new species. Any new genome\_dbs, dnafrags and method\_link\_species\_set\_ids need to be added before production starts.

☒ Add genome\_db

▼ [Click here for details](#)

This may have already been done if the dna guys started early, please check.

### Add genome\_db

```
perl $ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/update_genome.pl
--reg_conf
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/production_reg_conf.pl
--compara compara_master --species "gadus_morhua" --collection ensembl
```

Add the new genome\_db\_id to the confluence page [Release plans](#). This script may take a while if the species you are adding is new. You can check the progress by counting dnaFrag entries in the master database:

```
select count(*) from dnafrag;
```

NB: The updated update\_genome.pl should do the following automatically. Please check that it does, and remove the italicised text if ok:

*If it is a new genebuild and the assembly hasn't changed, you can just edit the entry in master (genome\_db table) introducing the new genebuild from meta.genebuild.start\_date in the core database.*

```
update genome_db set genebuild = "2011-07-FlyBase" where genome_db_id = 105;
```

☒ Add in extra non-reference patches.

▼ [Click here for details](#)

This is currently done when a new patch for either human or mouse is released. This may have already been done, please ask.

Details about the patches can be found here [ftp://ftp.ncbi.nlm.nih.gov/genbank/genomes/Eukaryotes/vertebrates\\_mammals/Homo\\_sapiens/](ftp://ftp.ncbi.nlm.nih.gov/genbank/genomes/Eukaryotes/vertebrates_mammals/Homo_sapiens/)

eg for patch 11: [ftp://ftp.ncbi.nlm.nih.gov/genbank/genomes/Eukaryotes/vertebrates\\_mammals/Homo\\_sapiens/GRCh37.p11/README](ftp://ftp.ncbi.nlm.nih.gov/genbank/genomes/Eukaryotes/vertebrates_mammals/Homo_sapiens/GRCh37.p11/README)

It is first necessary to find if any patches have been deleted or updated since these need to be deleted from the master before the new and replacement patches are added. This is done by running the find\_assembly\_patches.pl script on the new and previous release of the core database

### Find assembly patches

```
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/find_assembly_patches.pl
l -new_core
"mysql://ensro@ens-staging1:3306/homo_sapiens_core_73_37?group=core&species=homo_sapiens"
-prev_core
"mysql://ensro@ens-livemirror:3306/homo_sapiens_core_72_37?group=core&species=homo_sapiens"
```

The output looks like this (format: name, seq\_region\_id, date)

### Sample output

```
NEW patches
HG29_PATCH 1001061122 2013-02-18 14:16:55
HG1592_PATCH 1001061114 2013-02-18 14:16:55
HG385_PATCH 1001061116 2013-02-18 14:16:55
HSCHR6_2_CTG5 1001061112 2013-02-18 14:16:55
HG1079_PATCH 1001061118 2013-02-18 14:16:55
CHANGED patches
HG1436_HG1432_PATCH new=1001061124 2013-02-18 14:16:55
prev=1000859885 2012-10-08 16:48:36
HG1292_PATCH new=1001061108 2013-02-18 14:16:55      prev=1000759258
2012-04-27 12:12:07
HSCHR22_1_CTG1 new=1001061120 2013-02-18 14:16:55      prev=1000057052
2010-09-07 14:22:38
HG1287_PATCH new=1001061110 2013-02-18 14:16:55      prev=1000859831
2012-10-08 16:48:36
DELETED patches
Patches to delete:
( "HG1436_HG1432_PATCH" , "HG1292_PATCH" , "HSCHR22_1_CTG1" , "HG1287_PATCH" )
```

In this case, there are 5 NEW patches, 4 CHANGED patches and no DELETED patches. Any CHANGED or DELETED patches must be deleted from the master before importing the new patch set.

To add extra non-reference patches to an existing assembly, you need the `-force` option to just add those dnafrags which aren't already in the database.

### Add patches

```
perl $ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/update_genome.pl
--reg_conf
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/production_reg_conf.pl
--compara compara_master --species human --force
```

## Add method\_link\_species\_set entries to compara master database

These are usually added by the people that need them, please check.

The release coordinator (or any team member) should create a new `method_link_species_set` in the master database before starting a new pipeline in order to get a unique `method_link_species_set_id`. Ideally they can be created before starting to build the new database although new `method_link_species_sets` can be added later on.

☒ Add dna method\_link\_species\_set entries

### Pairwise method\_link\_species\_set

```
perl $ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/create_mlss.pl
--method_link_type LASTZ_NET --genome_db_id 90,142 --source "ensembl" --reg_conf
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/production_reg_conf.pl
--compara compara_master
```

☒ Add syntenic method\_link\_species\_set entries

### Synteny method\_link\_species\_set

```
perl $ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/create_mlss.pl
--method_link_type SYNTENY --genome_db_id 90,142 --source "ensembl" --reg_conf
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/production_reg_conf.pl
-compara compara_master
```



Add homology method\_link\_species\_set\_entries

▼ [Click here for details](#)

Choose a temp. directory where the output will be generated:

### Choose temp directory

```
export MLSS_DIR="/tmp/mlss_creation"
mkdir $MLSS_DIR
```

Run the loading script several times:

--pw stands for all pairwised genome\_db\_ids in the list provided

--sg stands for keep genome\_db\_id in the list alone (singleton)



### Protein method\_link\_species\_set

```
# orthologues
echo -e "201\n" | perl
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/create_mlss.pl --f \
--reg_conf
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/production_reg_conf.pl
\
--pw --collection ensembl 1>$MLSS_DIR/create_mlss.ENSEMBL_ORTHOLOGUES.201.out
2>$MLSS_DIR/create_mlss.ENSEMBL_ORTHOLOGUES.201.err

# paralogues wth
echo -e "202\n" | perl
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/create_mlss.pl --f \
--reg_conf
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/production_reg_conf.pl
\
--sg --collection ensembl
1>$MLSS_DIR/create_mlss.ENSEMBL_PARALOGUES.wth.202.out
2>$MLSS_DIR/create_mlss.ENSEMBL_PARALOGUES.wth.202.err

# proteintrees
echo -e "401\n" | perl
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/create_mlss.pl --f \
--reg_conf
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/production_reg_conf.pl
\
--name "protein trees" --collection ensembl
1>$MLSS_DIR/create_mlss.PROTEIN_TREES.401.out
2>$MLSS_DIR/create_mlss.PROTEIN_TREES.401.err

# nctrees
echo -e "402\n" | perl
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/create_mlss.pl --f \
--reg_conf
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/production_reg_conf.pl
\
--name "nc trees" --collection ensembl
1>$MLSS_DIR/create_mlss.NC_TREES.402.out
2>$MLSS_DIR/create_mlss.NC_TREES.402.err

# families
echo -e "301\n" | perl
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/create_mlss.pl --f \
--reg_conf
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/production_reg_conf.pl
\
--name "families" --collection ensembl 1>$MLSS_DIR/create_mlss.FAMILY.301.out
2>$MLSS_DIR/create_mlss.FAMILY.301.err
```

If output/error files are ok, remove them all:

### Remove files

```
rm -rf $MLSS_DIR
```

Otherwise make yourself a nice cup of tea and then \*PANIC\*

## Add new species to phylogenetic tree

- ✓ Add new species to phylogenetic tree

✓ [Click here for details](#)

The easiest way to use this is to use the phylowidget.

From the Ensembl home page:

View full list of all Ensembl species

Species tree (Requires Java)

Select Arrow and select where you want the new species to go (use ncbi taxonomy or wikipedia etc) eg *Gadus morhua*

Then select in the menu "Tree Edit > Add > Sister"

Click on the empty node and edit name (add new name) and branch length

The tree should appear in the Toolbox but if not, then save the tree

Copy the new tree into

```
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/species_tree_blength.nh
```

git commit

```
## rel.65 -- Added Gadus morhua
```

```
## rel.66 -- Added Latimeria chalumnae (Coelacanth)
```

```
## Also included 'ensembl timetree mya' in ncbi_taxa_name for
```

```
## taxon_id: 8287 -- value: 414.9 in the master and final database.
```

```
## rel.69 added Xiphophorus maculatus and Mustela putorius furo
```

## Final checks to compara master database

- ✓ Check if any new species have been postponed  
If a new species is postponed for this release, check that no entries (genome\_db, dnafrags, etc) were added to the master database. If they were, you can simply switch the assembly\_default value in the genome\_db table.  
[1] for species making it / used in the pipeline  
[0] for species not making it / or old assemblies
- ✓ Check that all the new species have been added to the 'ensembl' collection:  
✓ [Click here for details...](#)

```
select name from species_set_tag join species_set using(species_set_id) join
genome_db using(genome_db_id) where value like 'collection-ensembl' order by
genome_db_id;
```

- ☒ Drop method\_link\_species\_set entries for alignments which did not make it.  
    [Click here for details](#)

Check with other members of compara. Remove redundant entries.

### mlss

```
SELECT ml.* FROM method_link ml LEFT JOIN method_link_species_set mlss ON  
ml.method_link_id=mlss.method_link_id WHERE mlss.method_link_id IS NULL;
```

## Create Release Database

Create the new database for the new release and add it to your registry configuration file. Use the \$ENSEMBL\_CVS\_ROOT\_DIR/ensembl-compara/sql/table.sql file to create the tables and populate the database with the relevant primary data and genomic alignments that can be reused from the previous release. This can be done with the \$ENSEMBL\_CVS\_ROOT\_DIR/ensembl-compara/scripts/pipeline/populate\_new\_database.pl script. It requires the master database, the previous released database and the fresh new database with the tables already created. The script will copy relevant data from the master and the old database into the new one.

- ☒ Create new database  
    [Click here for details](#)

### Create database

```
db_cmd.pl $COMPARA_REG compara_curr -sql "CREATE DATABASE"  
db_cmd.pl $COMPARA_REG compara_curr <  
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/sql/table.sql
```

- ☒ Populate the new database  
    [Click here for details](#)

Before you start copying, make a dry run of the populate\_new\_database.pl with -intentions flag to review the list of mlss\_ids to be copied:

### populate\_new\_database intentions

```
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/populate_new_datab  
ase.pl --reg-conf  
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/production_reg_con  
f.pl \  
--master compara_master --old compara_prev --new compara_curr  
--intentions > populate_new_database.intentions
```

This takes about a minute and produces a long list.

If you believe some of the entries should NOT be copied, you should manually add 'skip\_mlss' and 'skip\_ss' entries into master database meta table.

NB There are cases where the mlss does not change but the underlying data does eg the "patch-to-ref" alignment (H.sap-H.sap lastz-patch and M.mus-M.mus lastz-patch). These have a mlss\_id of 556 (H.sap) and 624 (M.mus) and are currently set in the skip\_mlss. If there are no new patches, this needs to be removed to allow the existing data to be copied. If there are new patches, please ensure the 'skip\_mlss' is set in the meta table. However, the entry in the method\_link\_species\_set table will not be copied and will need to be added manually.

Start the copying:

### populate\_new\_database

```
time
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/populate_new_database.pl \
    --reg-conf
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/production_reg_config.pl --master compara_master --old compara_prev --new compara_curr >
populate_new_database.out
```

#### ▼ [Click here for run times](#)

took 3 hours for rel.pre57 (copied from rel.56)

took 3 hours for rel.57 (copied from rel.pre57)

took 2:15 hours for rel.58 (copied from rel.57)

took 2:09 hours for rel.59 (copied from rel.58)

took 3 hours for rel. 60 (copied from rel.59)

rel.64: 2.6h

rel.65: 2.5h

rel.66: 4.8h

rel.67: 2.1h (launched from compara3)

rel.68: 1h40m (run on compara3)

rel.69: 2.5h

rel.70: ~3.5h (compara1 was slow)

rel.71: 4.1h (compara3)

rel.72: 5.1h (compara3)

rel.73: 5.5h (compara2)

rel.74: 2h:3' (compara3)

rel.75: 5.5h (compara5)

rel77: 9.7hr(compara5)

If new method\_link\_species\_sets are added in the master after this, you use this script again to copy the new relevant data. In such case, you will have to:

- skip the old\_database in order to avoid trying to copy the dna-dna alignments and syntenies again
- empty ncbi\_taxa\_name before running

### populate\_new\_database from master only

```
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/populate_new_database.pl \
    --reg-conf
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/production_reg_config.pl --master compara_master --new compara_curr
```



Delete any pairwise alignments on non-reference patches that have been DELETED or UPDATED.

#### ▼ [Click here for details](#)

For the UPDATED patches, you need to use the original dnafrag\_ids (ask Compara members who ran the alignments on patches -

usually Kathryn).

### delete patches

```
SELECT count(*) FROM genomic_align gal, genomic_align ga2, genomic_align_block
gab WHERE gal.dnafrag_id in (13768162,13728785,13768189,12967785) AND
gal.genomic_align_block_id = ga2.genomic_align_block_id AND
gal.genomic_align_id != ga2.genomic_align_id AND gal.genomic_align_block_id =
gab.genomic_align_block_id;
99906
```

```
DELETE gal, ga2, gab FROM genomic_align gal, genomic_align ga2,
genomic_align_block gab WHERE gal.dnafrag_id in
(13768162,13728785,13768189,12967785) AND gal.genomic_align_block_id =
ga2.genomic_align_block_id AND gal.genomic_align_id != ga2.genomic_align_id
AND gal.genomic_align_block_id = gab.genomic_align_block_id;
Query OK, 299718 rows affected (34.34 sec)
99906*3=299718
```

#### ☐ Run healthchecks on the release database

✓ [Click here for details](#)

Run the healthchecks to make sure the the release database is consistent after the initial population of data.

Click [here](#) for how to setup and run the healthchecks

Run the compara\_external\_foreign\_keys healthcheck

### healthcheck

```
cd $ENSEMBL_CVS_ROOT_DIR/ensj-healthcheck

# make sure you are using the right version of JAVA:
export JAVA_HOME=/software/jdk1.6.0_14

# if you need to recompile (submit to the farm, because you need more memory
than is available on the head) :
bsub -I ant clean jar

# some tests need more memory than the farm3's default:
time bsub -q yesterday -M8000 -R"select[mem>8000] rusage[mem=8000]" -I
./run-configurable-testrunner.sh -h compara5 -d sf5_ensembl_compara_77 --host2
ens-staging2 -g ComparaExternalForeignKeys
```

and correct any problems, if any

## Merge DNA data

NOTE: All the runs of copy\_data.pl (except the last one) should have the flag "-re\_enable 0" to avoid constantly recomputing the indices

#### ☒ Pairwise alignments: LASTZ\_NET, BLASTZ\_NET, TRANSLATED\_BLAT\_NET and the special case of LASTZ\_PATCH.

\*NOTE\* : For merging pw alignments involving haplotypes, go to the next point

✓ [Click here for details](#)

These data are usually in separate production databases. You can copy them using the \$ENSEMBL\_CVS\_ROOT\_DIR/ensembl-compara/scripts/pipeline/copy\_data.pl script. This script requires write access to the production database if the dnafrag\_ids need fixing. Use the flag -re\_enable 0 on all calls apart from the last one to avoid recomputing the indices.

Example:

### copy\_data

```
# for each source URL: first plug in the --from_url and add --dry_run to check
that the script has found the right MLSS:
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/copy_data.pl --reg_conf
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/production_reg_conf.pl
--to_reg_name compara_curr --method_link_type LASTZ_NET --re_enable 0
--from_url mysql://ensadmin:${ENSADMIN_PSW}@compara4/sf5_ggal_falb_lastz_73
--dry_run

# if happy, remove the --dry_run flag and run it again, preferably on the
farm:
bsub -q yesterday -R "select[mem>5000] rusage[mem=5000]" -M5000 -I time
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/copy_data.pl --reg_conf
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/production_reg_conf.pl
--to_reg_name compara_curr --method_link_type LASTZ_NET --re_enable 0
--from_url mysql://ensadmin:${ENSADMIN_PSW}@compara4/sf5_ggal_falb_lastz_73

# sometimes you will also need to copy Human_ref_vs_Human_patches (note
the method_link_type is LASTZ_PATCH !)
bsub -q yesterday -R "select[mem>5000] rusage[mem=5000]" -M5000 -I time
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/copy_data.pl --reg_conf
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/production_reg_conf.pl
--to_reg_name compara_curr --method_link_type LASTZ_PATCH --re_enable 0
--from_url
mysql://ensadmin:${ENSADMIN_PSW}@compara3/sf5_compara_human_lastz_patch_and_ha
plotype_73
```

- ☐ Pairwise alignments: non-reference patches for the high coverage LASTZ\_NET alignments. This is to be used when merging pw alignments involving haplotypes.

▼ [Click here for details](#)

This step is now very similar to the previous.

Do not forget --merge and --patch\_merge options.

Also, if it's the last one you might want to switch keys back on

### copy\_data --merge

```
# first plug in the --from_url and add --dry_run to check that the script has
found the right MLSS:
bsub -q yesterday -R "select[mem>5000] rusage[mem=5000]" -M5000 -I time
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/copy_data.pl --reg_conf
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/production_reg_conf.pl
--to_reg_name compara_curr --method_link_type LASTZ_NET --method_link_type
BLASTZ_NET --method_link_type TRANSLATED_BLAT_NET --re_enable 1 --merge
--patch_merge --from_url
mysql://ensadmin:${ENSADMIN_PSW}@compara2/kb3_hsap_lastz_hap_73 --dry_run

# if happy, remove the --dry_run flag and run it again, preferably on the
farm:
bsub -q yesterday -R "select[mem>5000] rusage[mem=5000]" -M5000 -I time
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/copy_data.pl --reg_conf
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/production_reg_conf.pl
--to_reg_name compara_curr --method_link_type LASTZ_NET --method_link_type
BLASTZ_NET --method_link_type TRANSLATED_BLAT_NET --re_enable 1 --merge
--patch_merge --from_url
mysql://ensadmin:${ENSADMIN_PSW}@compara2/kb3_hsap_lastz_hap_73
```

- ✓ Multiple alignments: PECAN, EPO, EPO\_LOW\_COVERAGE, GERP\_CONSTRAINED\_ELEMENT, GERP\_CONSERVATION\_SCORE  
✓ [Click here for details](#)

These data are usually in separate production databases. You can copy them using the \$ENSEMBL\_CVS\_ROOT\_DIR/ensembl-compara/scripts/pipeline/copy\_data.pl script. This script requires write access to the production database if the dnafrag\_ids need fixing or the data must be copied in binary mode (this is required for conservation scores).

Some alignments produce conservation scores and constrained elements (check the [Release plans](#)) and these need to be copied separately.

eg

### copy\_data multiple alignment

```
bsub -q yesterday -R "select[mem>5000] rusage[mem=5000]" -M5000 \
-I time
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/copy_data.pl \
--reg_conf
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/production_reg_conf.pl
--to_reg_name compara_curr \
--method_link_type EPO --method_link_type EPO_LOW_COVERAGE
--method_link_type PECAN \
--method_link_type GERP_CONSTRAINED_ELEMENT --method_link_type
GERP_CONSERVATION_SCORE \
--from_url
mysql://ensadmin:${ENSADMIN_PSW}@compara2/sf5_epo_low_8way_fish_71 -re_enable
0
```

EPO alignments produce ancestral sequences and a separate core database which must also be copied. See below.

- ✓ [Click here for run times](#)

rel 71.

2m: kb3\_hsap\_ggal\_lastz\_71 mlss\_id=632

```

1m kb3_mmus_ggal_lastz_71 mlss_id=633
1m kb3_ggal_mgap_lastz_71 mlss_id=634
1m kb3_ggal_xtro_tblast_71 mlss_id=638
1m kb3_hsap_ggal_tblast_71 mlss_id=637
1m sf5_olat_gmor_lastz_71 mlss_id=625
3m kb3_pecan_20way_71 mlss_id=630
4m kb3_pecan_20way_71 mlss_id=631
35m kb3_pecan_20way_71 mlss_id=50045
3m sf5_compara_epo_6way_71 mlss_id=548
1m sf5_olat_onil_lastz_71 mlss_id=626
1m sf5_olat_xmac_lastz_71 mlss_id=627
1m sf5_epo_low_8way_fish_71 mlss_id=628
2m sf5_epo_low_8way_fish_71 mlss_id=629
9m sf5_epo_low_8way_fish_71 mlss_id=50044
1m kb3_ggal_drer_tblast_71 mlss_id=639
1m kb3_ggal_csav_tblast_71 mlss_id=640
1m sf5_ggal_acar_lastz_71 mlss_id=636
91m sf5_ggal_tgut_lastz_7 mlss_id=635 (re-enable 1)
93m sf5_compara_epo_3way_birds_71 mlss_id=641 (re-enable 1)
14m sf5_compara_epo_3way_birds_71 mlss_id=642 (re-enable 1)
16m sf5_compara_epo_3way_birds_71 mlss_id=50046 (re-enable 1)

```

☒ Check the keys have been re-enabled

✓ [Click here for details](#)

Use mysqlshow to highlight if the table still has disabled keys. The text "disabled" will be shown in the Comment column if the key is disabled. An empty Comment column indicates the keys are enabled.

mysqlshow interprets any underscores in the last argument as a wildcard so to get round this, we need to use % as the last argument.

### mysqlshow

```

eval mysqlshow `db_cmd.pl $COMPARA_REG compara_curr -to_params` --keys
genomic_align_block %
eval mysqlshow `db_cmd.pl $COMPARA_REG compara_curr -to_params` --keys
genomic_align %
eval mysqlshow `db_cmd.pl $COMPARA_REG compara_curr -to_params` --keys
genomic_align_tree %
eval mysqlshow `db_cmd.pl $COMPARA_REG compara_curr -to_params` --keys
conservation_score %
eval mysqlshow `db_cmd.pl $COMPARA_REG compara_curr -to_params` --keys
constrained_element %

```

If there are still tables with keys disabled run the following on them:

```

db_cmd.pl $COMPARA_REG compara_curr -sql "ALTER TABLE <table_name> ENABLE
KEYS" ;

```





## Synteny

✓ [Click here for details](#)

First make sure the entries in \$ENSEMBL\_CVS\_ROOT\_DIR/ensembl-compara/scripts/pipeline/production\_reg\_conf.pl file point at the latest (staging) versions of the core databases.

Load the syntenic data by running the \$ENSEMBL\_CVS\_ROOT\_DIR/ensembl-compara/scripts/syntenic/LoadSyntenicData.pl script. This requires a syntenic file. The location of this should be on the [Release plans](#).

The "-ref" and "-nonref" species should be taken from the name of the method\_link\_species\_set that corresponds to the pairwise alignments containing both species

### Example

#### load syntenic data

```
perl $ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/syntenic/LoadSyntenicData.pl
--reg_conf
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/production_reg_conf.pl
\
    --dbname compara_curr -ref "Homo sapiens" -nonref "Callithrix jacchus"
-mlss_id 10052 \

/lustre/scratch101/ensembl/kb3/scratch/hive/release_64/kb3_hsap_cjac_syntenic_6
4/syntenic/all.100000.100000.BuildSyntenic
```



## Ancestral sequence core database

✓ [Click here for details](#)

Putting together the database of ancestral sequence is now done using a dedicated Hive-Core mini-pipeline.

Check you have the most recent core checkout ie the correct schema and patch files are added to the meta table.

Go to \$ENSEMBL\_CVS\_ROOT\_DIR/ensembl-compara/modules/Bio/EnsEMBL/Compara/PipeConfig and open the PipeConfig file AncestralMerge\_conf.pm .

Make sure you have edited/checked the following:

1) current release number

2) names and locations of current and previous ancestral core databases

3) the table of ancestral sequence sources in the second analysis (some entries might point to the previous release ancestral database, some will be new)

For (3), you can run the following query on your release database and on the previous database: (NB: method\_link\_id=13 is equivalent to method\_link\_type = "EPO")

#### EPO query

```
SELECT * FROM method_link_species_set WHERE method_link_id = 13;
```

The new mlss\_id should be attached to their production database:

'641' => 'mysql://ensadmin:\$ENSADMIN\_PSW@compara3/sf5\_3birds\_ancestral\_sequences\_core\_71'

The mlss\_id that are reused should be linked to the previous database

'505' => \$self->o('prev\_ancestral\_db'),

The current (as of rel.75) list of ancestral alignments are:

5 teleost fish

6 primates

4 sauropsids ("birds")

15 eutherian mammals

Save the changes, exit the editor and run init\_pipeline.pl with this file:

### init\_pipeline

```
init_pipeline.pl AncestralMerge_conf.pm -host compara5
```

Then run both -sync and -loop variations of the beekeeper.pl command suggested by init\_pipeline.pl . This pipeline will merge the separate ancestral core sources into ensembl\_ancestral\_{rel\_number}.

You may want to check the msg table for errors and have a look at the result of the merger:

### Which Ancestral sequences do we have?

```
SELECT left(name,12) na, count(*), min(seq_region_id), max(seq_region_id),  
max(seq_region_id)-min(seq_region_id)+1 FROM seq_region GROUP BY na;
```

If everything is ok, measure the time:

### how much time did running of the pipeline take?

```
call time_analysis('%')
```

✓ [Click here for run times](#)

rel.67: 20min

rel.71: 20min

rel.75: 21min

Then drop hive-specific tables:

### drop hive tables

```
CALL drop_hive_tables;
```

Make sure all tables are myISAM.

```
SHOW TABLE STATUS where engine != 'MyISAM';
```

## Merge GeneTrees+Families+NCTrees+PatchProjectionsAsHomologies

- ✓ Check that CAFE has been run on the ProteinTrees and NCTrees
- ✓ Check that LRGs were included in the Families
- ✓ Check that PatchProjectionsAsHomologies (that includes loading DisplayLabels+Descriptions) pipeline has been run
- ✓ Run the Hive pipeline ( MergeDBsIntoRelease\_conf ) to merge tables from all the four products into the release database

✓ [Click here for details](#)

Go to \$ENSEMBL\_CVS\_ROOT\_DIR/ensembl-compara/modules/Bio/Ensembl/Compara/PipeConfig and open the PipeConfig file MergeDBsIntoRelease\_conf.pm

It has a 'urls' hash where you will have to change the names of the databases and possibly their locations:

master\_db - is the main compara master

prev\_rel\_db - should point to the previous release database

curr\_rel\_db - should point to the current release database being merged into ( not the Hive pipeline database, but purely Compara schema product )

protein\_db - should point to the current GeneTrees pipeline database

family\_db - should point to the current Families pipeline database

ncrna\_db - should point to the current ncRNAtrees pipeline database

projection\_db - should point to the current PatchProjectionsAsHomologies pipeline database.

Also choose the server to run the merging pipeline on ( you don't need a lot of resources or memory, as it is purely Hive book-keeping ) and set the 'host' default\_option.

Save the changes, exit the editor, init and run the merging pipeline with this file:

#### running the merging pipeline

```
init_pipeline.pl MergeDBsIntoRelease_conf.pm
```

```
beekeeper.pl ... -sync
```

```
beekeeper.pl ... -loop
```

This pipeline will merge all the protein-side products into the release database.

#### ▼ [Click here for times](#)

rel.73 was the first experimental run, code had to be fixed, servers had to be reconfigured, so merging took one whole working day.

In the merging database run: call time\_analysis('%');

rel.75 : 5 hours

rel.76 : 5.6 hours

#### ☒ Add the LRG dnafrags

##### ▼ [Click here to expand...](#)

```
mysql.d compara2 lg4_families_77 dnafrag -t --where 'coord_system_name =  
"lrg"' | mysql.1 compara5 sf5_ensembl_compara_77
```

#### ☒ Populate the member\_production\_counts table

##### ▼ [Click here for details](#)

Populate the member\_production\_counts table using the script at

\$ENSEMBL\_CVS\_ROOT\_DIR/ensembl-compara/scripts/production/populate\_member\_production\_counts\_table.sql

time db\_cmd.pl \$COMPARA\_REG compara\_curr <

\$ENSEMBL\_CVS\_ROOT\_DIR/ensembl-compara/scripts/production/populate\_member\_production\_counts\_table.sql

#### ▼ [Click here for times](#)

rel.73: 75min

rel.75: 43min

rel.76: 30min

- ✓ Clean up the merging database

✓ [Click here for details](#)

After you are happy about the result of protein side merging you can drop the "yourname\_pipeline\_dbmerge\_75" database.

- ✓ git commit the changes to the PipeConfig files that you have made.

## Final database checks

- ✓ Remove redundant method\_link entries

✓ [Click here for details](#)

In most cases they can be removed, but check with other members of Compara. Remove redundant method\_link entries

### method\_link entries

```
SELECT ml.* FROM method_link ml LEFT JOIN method_link_species_set mlss
USING(method_link_id) WHERE mlss.method_link_id IS NULL;
DELETE ml FROM method_link ml LEFT JOIN method_link_species_set mlss
USING(method_link_id) WHERE mlss.method_link_id IS NULL;
```

- ✓ Check that all the schema patches have been declared and applied.

✓ [Click here for details](#)

If unsure, recheck the current schema against the previous schema. See Check the patch files for details

## Run the healthchecks

- ✓ Update the code

✓ [Click here for details](#)

The healthchecks are written in java and need to be recompiled after a git pull.

### compile healthchecks

```
cd $ENSEMBL_CVS_ROOT_DIR/ensj-healthcheck
export JAVA_HOME=/software/jdk1.6.0_14
bsub -I ant clean jar
```

We don't need to configure a database.properties any more. Everythin is done from the command line

- ✓ Run the healthchecks for ancestral database

✓ [Click here for details](#)

```
time bsub -M8000 -R"select[mem>8000] rusage[mem=8000]" -I
./run-configurable-testrunner.sh -h compara5 -d sf5_ensembl_ancestral_77 -g
ComparaAncestral
```

✓ [Click here for times](#)

```
# rel.62: 4sec, all successful
# rel.63: 14sec, all successful after analyzing 3 tables
# rel.65: 13sec, all successful after analyzing the tables.
# rel.67: 8sec, all successful
# rel.68: 8sec, all successful
# rel.71: 10sec all successful apart from UTR.java (ignore)
# rel.75: 20sec, all successful after analyzing 3 tables
```

- ✓ Update the max\_alignment\_length IF NECESSARY.

✓ [Click here for details](#)

Check that the max\_alignment\_lengths have been computed.

#### update max\_alignment\_length

```
time bsub -I ./run-configurable-testrunner.sh -h compara5 -d  
sf5_ensembl_compara_77 -t  
org.ensembl.healthcheck.testcase.compara.MLSSTagMaxAlign
```

If not (the healthcheck is failing), you can repair it by adding the --repair flag:

#### update max\_alignment\_length

```
time bsub -I ./run-configurable-testrunner.sh -h compara5 -d  
sf5_ensembl_compara_77 -t  
org.ensembl.healthcheck.testcase.compara.MLSSTagMaxAlign --repair 1
```

✓ Update the alignment mlss\_id of the conservation score IF NECESSARY

✓ [Click here for details](#)

#### update conservation score mlss\_id

```
time bsub -I ./run-configurable-testrunner.sh -h compara5 -d  
sf5_ensembl_compara_77 -t  
org.ensembl.healthcheck.testcase.compara.MLSSTagGERPMSA
```

If the healthcheck is failing, you can repair it by adding the --repair flag:

#### update conservation score mlss\_id

```
time bsub -I ./run-configurable-testrunner.sh -h compara5 -d  
sf5_ensembl_compara_77 -t  
org.ensembl.healthcheck.testcase.compara.MLSSTagGERPMSA --repair 1
```

✓ Run the compara\_external\_foreign\_keys healthcheck

✓ [Click here for details](#)

#### compara\_external\_foreign\_keys

```
time bsub -q yesterday -M8000 -R"select[mem>8000] rusage[mem=8000]" -I  
./run-configurable-testrunner.sh -h compara5 -d sf5_ensembl_compara_77 --host2  
ens-staging2 -g ComparaExternalForeignKeys
```

✓ [Click here for previous results](#)

rel.56 everything passed apart from CheckTaxon - according to Javier in this particular case it was not a problem

rel.pre57 it took 20 minutes (all passed).

rel.57 it took ?? minutes ('genbank common name' for 4 species had to be copied from their 'ensembl common name' in ncbi\_taxa\_name table)

rel.58: 22m

rel.61: 32m, 1 failure

rel.63: 23m, 1 failure ( taeniopygia\_guttata\_core\_63\_1: common\_name::zebra finch is not in lg4\_ensembl\_compara\_63 )

rel.65: 25m, 1 failure ( taeniopygia\_guttata\_core\_65\_1: common\_name::zebra finch is not in mp12\_ensembl\_compara\_65 )

rel.66: 12m, 1 failure ( taeniopygia\_guttata\_core\_66\_1: common\_name::zebra finch is not in mp12\_ensembl\_compara\_66 )  
rel.71: 1 failure (ncbi\_taxa\_name: Lepidion inosimae: common\_name "Morid cod" and "morid cod", this was agreed to be OK in rel.70).  
rel.72 12m  
rel.73 12.5m

- ☒ Run the compara\_genomic healthcheck  
    [Click here for details](#)

### compara\_genomic

```
time bsub -q yesterday -M8000 -R"select[mem>8000] rusage[mem=8000]" -I  
./run-configurable-testrunner.sh -h compara5 -d sf5_ensembl_compara_77 --host2  
ens-staging2 -g ComparaGenomic
```

- [Click here for previous results](#)

rel.58: 47m, 7 failures  
rel.61: 51m, 3 failures  
rel.63: 84m, (2.5 errors that are "ok")  
rel.65: 75m, (5 errors that are "ok")  
rel.66: 23m, (4 errors that are "ok")  
rel.72: 57m  
rel.73 43m

- ☒ Run the compara\_homology healthcheck  
    [Click here for details](#)

### compara\_homology

```
time bsub -q yesterday -M8000 -R"select[mem>8000] rusage[mem=8000]" -I  
./run-configurable-testrunner.sh -h compara5 -d sf5_ensembl_compara_77 --host2  
ens-staging2 -g ComparaHomology
```

- [Click here for previous results](#)

rel.58: 14m, 5 failures  
rel.61: 30m, 3 failures  
rel.62: 3m, 1 failure (CheckSpeciesSetTag, ok?)  
rel.63: 52m, success (after some fixing, of course)  
rel.65: 48m, 3 errors  
rel.66: 18m  
rel.72: 1h6m  
rel.73 45m, 3 errors ( CheckSpeciesSetTag:2 and ForeignKeyMethodLinkSpeciesSetId:1 - detected unused SpeciesSets that were erroneously transferred over from prev. release)

## Test web server

- ☒ Ask ensembl-production to point the test web server to the compara release database  
Upon confirmation from the release coordinator ask other members of Compara to check their data on:  
    <http://staging.ensembl.org/>

## Run ANALYZE\_TABLE and OPTIMIZE TABLE

This is required for the CopyDbOverServer script to work properly.

- ☒ Run ANALYZE\_TABLE on compara and ancestral databases  
    [Click here for details](#)

### analyze table

```
time eval mysqlcheck `db_cmd.pl $COMPARA_REG compara_curr -to_params`  
--analyze --verbose  
time eval mysqlcheck `db_cmd.pl $COMPARA_REG ancestral_curr --reg_type core  
-to_params` --analyze --verbose
```

▼ [Click here for times](#)

rel.56 12min

rel.pre57 30+105min

rel.57 9+4+5min

rel.58 3min

rel.62 6min

rel.63 25m

rel.64 16m

rel.65 21m

rel.66 9m

rel.71 7m

rel.72 10m49.630s

rel.73 6m

rel.75 23+1sec

☒ Run OPTIMIZE TABLE on compara and ancestral databases

▼ [Click here for details](#)

### optimize table

```
time eval mysqlcheck `db_cmd.pl $COMPARA_REG compara_curr -to_params`  
--optimize --verbose  
time eval mysqlcheck `db_cmd.pl $COMPARA_REG ancestral_curr --reg_type core  
-to_params` --optimize --verbose
```

▼ [Click here for times](#)

rel.56 2.5 hours

rel.pre57 : took several iterations (not all tables were MyISAM initially), last one 132min.

rel.57 2+1.6 hours

rel.62 32min

rel.63 61m

rel.64 1.5h

rel.65 3h

rel.66 2h

rel.71 2h

rel.72 104m53.572s

rel.73 4h

rel.75 8.5min + 2sec

## Copy databases to staging servers

- ✓ Logon to ens-staging1 and create a file containing the copying options

✓ [Click here for details](#)

First, ssh into the DESTINATION machine and switch to the bash shell

# NB: ask for the password for mysqlens well in advance - there may be no-one around you at the right moment!

### ssh

```
ssh mysqlens@ens-staging1
bash
```

Create a file that will contain one line for each database with the source/destination parameters, like this:

### copy\_options

```
cat <<EOF >/tmp/lg4_ensembl_compara_75.copy_options
#from_host      from_port  from_dbname      to_host
to_port        to_dbname
#
compara5        3306        lg4_ensembl_compara_75      ens-staging1      3306
ensembl_compara_75
compara5        3306        lg4_ensembl_ancestral_75    ens-staging1      3306
ensembl_ancestral_75
EOF
```

Set the password into the environment variable:

```
export ENSADMIN_PSW='...'
```

- ✓ Copy the databases to ens-staging1

✓ [Click here for details](#)

You should check whether there is enough space on the disk before starting the copy.

### CopyDBoverServer

```
time perl ~lg4/work/ensembl/misc-scripts/CopyDBoverServer.pl -pass
$ENSADMIN_PSW \
    -noflush /tmp/lg4_ensembl_compara_75.copy_options >
/tmp/lg4_ensembl_compara_75.copy.err 2>&1
```

✓ [Click here for times](#)

copying of rel\_56 took 2 hours (SUCCESSFUL for both databases - you should check the output file)

copying of ensembl\_compara\_pre57 took 2 hours (SUCCESSFUL)

copying of ensembl\_compara\_57 took 2 hours (SUCCESSFUL)

copying of ensembl\_ancestral\_57 took 20 minutes (only SUCCESSFUL after analyzing/optimizing)

copying of ensembl\_compara\_58 and ensembl\_ancestral\_58 together took 1:30h (SUCCESSFUL)

copying of ensembl\_compara\_62 and ensembl\_ancestral\_62 took 2h38

copying of ensembl\_compara\_63 and ensembl\_ancestral\_63 took 2h



copying of ensembl\_compara\_64 and ensembl\_ancestral\_64 took 2h15m

rel.65 2h51m

rel.66 1h12m to copy over ensembl\_compara\_66

0h11m to copy over ensembl\_ancestral\_66

rel.68 90m39.996s

rel.70 ens-staging2 (1h29m4s)

rel.71: 1h39m

rel.72 1h40m50s

rel.73 2h

rel.75: 1h30m (both successful)

- ✓ Logon to ens-staging2 and create a file containing the copying options

✓ [Click here for details](#)

First, ssh into the DESTINATION machine and switch to the bash shell

# NB: ask for the password for mysqlens well in advance - there may be no-one around you at the right moment!

### ssh

```
ssh mysqlens@ens-staging2
bash
```

Create a file that will contain one line for each database with the source/destination parameters, like this:

### copy\_options

```
cat <<EOF >/tmp/lg4_ensembl_compara_75.copy_options
#from_host      from_port  from_dbname      to_host
to_port        to_dbname
#
compara5        3306          lg4_ensembl_compara_75      ens-staging2      3306
ensembl_compara_75
compara5        3306          lg4_ensembl_ancestral_75    ens-staging2      3306
ensembl_ancestral_75
EOF
```

Set the password into the environment variable:

```
export ENSADMIN_PSW='...'
```

- ✓ Copy the databases to ens-staging2

✓ [Click here for details](#)

You should check whether there is enough space on the disk before starting the copy.

### CopyDBoverServer

```
time perl ~lg4/work/ensembl/misc-scripts/CopyDBoverServer.pl -pass  
$ENSADMIN_PSW \  
    -noflush /tmp/lg4_ensembl_compara_75.copy_options >  
/tmp/lg4_ensembl_compara_75.copy.err 2>&1
```

✓ [Click here for times](#)

copying of ensembl\_compara\_58 took 1:15h

copying of ensembl\_compara\_62 took 1:34h

copying of ensembl\_compara\_63 took 3:44h

copying of ensembl\_compara\_64 took 3h51m

rel.65: 3h28m

rel.66: ~2h

rel.68: 104m11.046s

rel.72: 1h40m24s

rel.73: 3.75h (started at the same time as ens-staging1, but start was delayed because of locking)

rel.75: 1h25m (both successful)

## Final handover of databases

- ✓ Send an email to ensembl-production to announce the handover the databases
- ✓ Update the Declaration of Intentions <http://admin.ensembl.org/index.html> to indicate what has been handed over and what didn't make it and has been postponed

## Update documentation and diagrams

- ✓ Update pipeline diagrams
  - ✓ [Click here for details](#)
  - Update the pipeline diagrams in the docs directory

### pipeline diagrams

```
export COMPARA_REG="-reg_conf  
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/production_reg_conf.pl  
-reg_type compara -reg_alias"  
  
cd $ENSEMBL_CVS_ROOT_DIR/ensembl-compara/docs/pipelines/diagrams  
generate_graph.pl $COMPARA_REG compara_ptrees -output ProteinTrees.png  
generate_graph.pl $COMPARA_REG compara_nctrees -output ncRNATrees.png  
generate_graph.pl $COMPARA_REG compara_families -output Families.png  
generate_graph.pl -url mysql://ensro@compara3/kb3_pecan_20way_71 -output  
MercatorPecan.png  
generate_graph.pl -url mysql://ensro@compara4/sf5_epo_35way_68 -output  
EpoLowCoverage.png  
generate_graph.pl -url mysql://ensro@compara4/sf5_compara_epo_13way_69 -output  
Epo.png
```

Commit any changed diagrams to git and push.

- ☒ Update the schema documentation and diagrams  
    [Click here for details](#)

#### generate new schema documentation

```
perl $ENSEMBL_CVS_ROOT_DIR/ensembl-production/scripts/sql2html.pl -i
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/sql/table.sql -o
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/docs/schema/compara_schema.html -d
Compara -host compara5 -user ensro -dbname lg4_ensembl_compara_75
-sort_headers 0 -sort_tables 0 -intro
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/docs/schema/intro.html
```

Open the output file `$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/docs/schema/compara_schema.html` in your browser and check that no example errors are reported.

If everything looks fine, copy this file to public-plugins and commit&push both (the compara one and the webcode one) :

#### update schema documentation for web

```
cp $ENSEMBL_CVS_ROOT_DIR/ensembl-compara/docs/schema/compara_schema.html
$ENSEMBL_CVS_ROOT_DIR/public-plugins/docs/htdocs/info/docs/api/compara/
```

If necessary, update schema diagrams by loading the `$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/sql/table.sql` schema file into MySQL Workbench, rearrange/colour the nodes and export into PNG.

The schema diagrams will have to be copied both to `$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/docs/schema/diagrams` and to public-plugins and committed&pushed in both repositories:

#### update schema diagrams for web

```
cp -r $ENSEMBL_CVS_ROOT_DIR/ensembl-compara/docs/schema/diagrams/*.png
$ENSEMBL_CVS_ROOT_DIR/public-plugins/docs/htdocs/info/docs/api/compara/diagrams/
```

- ☒ Update the tutorial documentation  
    [Click here for details](#)

Update the tutorial documentation `compara_tutorial.html` in this directory:

`$ENSEMBL_CVS_ROOT_DIR/ensembl-webcode/htdocs/info/docs/api/compara/`

Be careful that the embedded Perl snippets must use HTML-escaped characters (e.g. `&lt;` and `&gt;`;) and be wrapped in a `<pre class="code sh_perl">`

Open the URL `/info/docs/api/compara/compara_tutorial.html` from a sandbox / test website and export it as a PDF in

`$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/docs/ComparaTutorial.pdf`

- ☒ Update main ensembl species tree if there are any new species

See `$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/docs/species_tree/README`

The end result should go here (check how it should be moved) : [http://www.ensembl.org/info/about/species\\_tree.pdf](http://www.ensembl.org/info/about/species_tree.pdf)

- ☐ Check examples work in `ensembl-compara/scripts/examples/`
- ☐ git commit and push any modified files or added tutorial examples

☒ Update the declared intention with removed / deprecated methods

[Click for details](#)

We need to generate the list of methods exported by the objects / adaptors on the master and release/{ $n$ -1} branches, and compare (diff) them.

### Check deprecated / removed methods

```
# on master/ensembl-compara/modules/Bio/Ensembl/Compara/
grep "^sub " *pm | sort > ~/MASTER
# on the previous release/ensembl-compara/modules/Bio/Ensembl/Compara/
grep "^sub " *pm | sort > ~/RELEASE75
sdiff -w 200 -bs ~/RELEASE75 ~/MASTER | less

# Let's do the same for the adaptors
# on master/ensembl-compara/modules/Bio/Ensembl/Compara/
grep "^sub " DBSQL/*pm | sort > ~/MASTER
# on the previous release/ensembl-compara/modules/Bio/Ensembl/Compara/
grep "^sub " DBSQL/*pm | sort > ~/RELEASE75
sdiff -w 200 -bs ~/RELEASE75 ~/MASTER | less
```

In both cases, make sure the methods are really removed, and not moved to a base / sub class, etc

## Data dumps

These should only be done once ensembl-production has given the go-ahead for this. This is to avoid overloading the databases whilst biomart is being run.

☒ DNA data dumps

Generally the person who ran the pipeline will also do the data dumps. The instructions are in \$ENSEMBL\_CVS\_ROOT\_DIR/ensembl-compara/docs/README.multi\_align.dumps

☒ Gene tree dumps

Generally the person who ran the pipeline will also do the data dumps.

[Click here for details](#)

Go to \$ENSEMBL\_CVS\_ROOT\_DIR/ensembl-compara/modules/Bio/Ensembl/Compara/PipeConfig and open the PipeConfig file DumpTrees\_conf.pm

Check that you are happy about all parameters. In usual cases, they can all be set from the command line and the config file does not need editing.

Make sure you have the XML::Writer module in your PERL5LIB (there is a copy in ~mm14/src/perl/orthoxml/)

Run init\_pipeline.pl with this file:

### init\_pipeline

```
init_pipeline.pl Bio::Ensembl::Compara::PipeConfig::DumpTrees_conf -rel_coord
lg4 -rel_db_host compara5 -pipeline_db -host=compara1 -member_type protein
```



rel.64: testing sqlite mode failed: too many occurrences of "database locked". We should stick to mysql.

Then run both -sync and -loop variations of the beekeeper.pl command suggested by init\_pipeline.pl .

This pipeline takes a couple of hours to run and will produce protein\_tree dumps in the directory pointed at by 'target\_dir' parameter. (/lustre/scratch110/ensembl/.\$self->o('ENV', 'USER'))./.\$self->o('pipeline\_name'))

Create the ncRNA pipeline from the same config file:

### init\_pipeline

```
init_pipeline.pl Bio::Ensembl::Compara::PipeConfig::DumpTrees_conf -rel_coord  
lg4 -rel_db_host compara5 -pipeline_db -host=compara1 -member_type ncRNA
```

Then run both -sync and -loop variations of the beekeeper.pl command suggested by init\_pipeline.pl .

This pipeline will produce ncRNA\_tree dumps in the directory pointed at by 'target\_dir' parameter. This is much faster: less than an hour

Commit the DumpTrees\_conf.pm file into git if you'd like to keep the changes.

☒ Copy the tree content dump for Uniprot

▼ [Click here for details](#)

The file 'target\_dir/ensembl.GeneTree\_content.{release}.txt.gz' needs to be copied to the EBI ftp server.

You can scp the file to login.ebi.ac.uk:/nfs/ftp/pub/databases/ensembl/ensembl\_compara/ and from there, create its MD5 checksum

☒ Report locations of the dumps to ensembl-production

▼ [Click here for details](#)

Dna dumps:

Not all the multiple alignments are run for each release. The data dumps for any multiple alignments not run this release should be copied from the previous release. Currently (e72), the full set of multiple alignments with the corresponding dump files, are:

6 primates EPO (emf)

13 eutherian mammals EPO (emf)

5 teleost fish EPO (emf)

3 neognath birds EPO (emf)

8 teleost fish EPO\_LOW\_COVERAGE (emf, bed)

36 eutherian mammals EPO\_LOW\_COVERAGE (emf, bed)

20 amniota vertebrates Pecan (emf, bed)

Gene tree dumps:

'target\_dir/emf' should go to /emf/ensembl-compara/homologies/ on the ftp

'target\_dir/xml' should go to /xml/ensembl-compara/homologies/ on the ftp

☒ Ancestral alleles (for the Variation team)

▼ [Click here to expand...](#)

```
perl
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/ancestral_sequences/get_ancestra
l_sequence.pl --conf
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/pipeline/production_reg_conf.pl
--compara_url mysql://ensro@compara5/sf5_epo_8primates_77 --species
homo_sapiens
dirname=homo_sapiens*
cd $dirname
perl
$ENSEMBL_CVS_ROOT_DIR/ensembl-compara/scripts/ancestral_sequences/get_stats.pl
> summary.txt
cd ..
tar -jcf ${dirname}.tar.bz2 $dirname
md5sum ${dirname}.tar.bz2 > ${dirname}.tar.bz2.MD5SUM
```

## Final things

☒ Dump the master database and place the copy in a safe place

✓ [Click here for details](#)

### dump master database

```
time eval mysqldump -t `db_cmd.pl $COMPARA_REG compara_master -to_params` | gzip -
> /warehouse/ensembl01/compara/master_dumps/ensembl_compara_master_75.mysql.gz
```

prev (rel65) /warehouse/ensembl01/sf5/Compara\_master\_dumps/sf5\_ensembl\_compara\_master.sql.04\_01\_12.gz

rel.65 /warehouse/ensembl01/compara/master\_dumps/sf5\_ensembl\_compara\_master.sql.04\_01\_12.gz

rel.66 /warehouse/ensembl01/compara/master\_dumps/sf5\_ensembl\_compara\_master.66.gz (3min to dump)

rel.71 /warehouse/ensembl01/compara/master\_dumps/sf5\_ensembl\_compara\_master.71.gz (3min to dump)

rel.75 /warehouse/ensembl01/compara/master\_dumps/ensembl\_compara\_master\_75.mysql.gz (46 sec to dump)

rel.76 /warehouse/ensembl01/compara/master\_dumps/ensembl\_compara\_master\_76.gz

rel.77 /warehouse/ensembl01/compara/master\_dumps/ensembl\_compara\_master\_77.gz (33 sec to dump)

☐ Create a word document and a pdf dump of this document

✓ [Click here for details](#)

In the top-right menu of this Confluence page, choose "Tools -> Export to PDF" and "Tools -> Export to Word".

Put these files into \$ENSEMBL\_CVS\_ROOT\_DIR/ensembl-compara/docs/

git commit and push

```
./run-configurable-testrunner.sh -h compara5 -d sf5_ensembl_compara_77
```