

**14 March 2022**

# **AutoCoEv** (v0.2.7 beta)

## *Manual*

by Petar Petrov

## Structure and settings

The AutoCoEv folder contains the following (Box 1). Configuration is done in a single file, called **settings.conf** (Box 2). Input files, working folder, databases paths, as well as, run-time and post-run options are configured there.

### Box 1. AutoCoEv/

<u>start.sh</u>	→ The main script, that needs to be executed.
settings.conf	→ Configuration file.
proteins/	→ <b>Folder</b> for the list(s) of proteins. Put proteins list here.
placental.tsv	→ Example list of species.
placental.nwk	→ Example species tree
functions/	→ Folder containing functions that <i>start.sh</i> calls.
doc/	→ Folder containing documentation, licensing and credits.
patches/	→ Patches for CAPS2, build instructions and precompiled binary.

### Box 2. settings.conf

```
## INPUT FILES
PROTEIN="protein/"      # **FOLDER** with proteins list(s)
SPECIES="placental.tsv" # FILE list of species
EXTTREE=""              # External species tree file. Specify if PhyML/PRANK need a guide.
PAIRLIST=""             # A list of defined protein pairs (only if PAIRINGMANNER="defined")

## REFERENCE ORGANISM AND ORTHOLOGUES
ORGANISM="10090"         # Taxid of the reference organism (e.g. Mouse)
LEVEL="40674"           # Level at which to search for orthologues (e.g. Mammalia)

## WORKING AND DATABASE DIRS
TMP="/var/tmp/work"      # Working folder
DTB="/var/tmp/DB10v1"    # Folder where databases are

## THREADS UTILIZATION
THREADS="$(nproc)"        # Number of (logical) cores to use

## BLAST OPTIONS
DETBAST="yes"            # Detailed BLAST results ("yes", "no"). Leave to "yes" for now.
PIDENT="35"              # Minimum allowed identity (%) to the reference organism
PGAPS="25"               # Maximum allowed gaps (%) to the reference organism

## ORTHOLOGUES ASSESSMENT BY GUIDANCE
GUIDANCEMSA="muscle"     # MSA method to use (see below for options).
GUIDANCECUT="0.95"       # Cutoff value, to exclude too divergent sequences (0: none)

## MSA OPTIONS
MSAMETHOD="prank"       # MSA method to use ("mafft-linsi", "muscle", "prank", ...)
MUSCLEOPTIONS=""         # Any additional options to pass to MUSCLE
MAFFTOPTIONS=""          # Any additional options to pass to MAFFT
PRANKOPTIONS=""          # Any additional options to pass to PRANK
PRANKGUIDE="noguide"     # Use external guide tree for PRANK ("exguide" or "noguide")?
GBLOCKSOPT="-b5=h"       # Gblocks options, e.g. allowed gaps: "-b5=h" (half)

## PhyML OPTIONS
PHYMLOPTIONS=""          # Any additional options to pass to PhyML
PHYMLGUIDE="noguide"     # Use external guide tree for PhyML ("exguide" or "noguide")?
TREESROOT="rooted"       # Root the generated trees by TreeBeST? ("rooted" or "noroot")

## PAIRING
PAIRINGMANNER="all"      # Pairing manner ("all" or "defined")
MINCOMMONSPCS="20"       # Minimum number of common species per protein pair
TREESCAPS="auto"         # Tree to use with CAPS ("phyml" or "auto")

## CAPS RUN-TIME OPTIONS
ALPHA="0.01"             # Alpha value for threshold cut-off. Do NOT leave blank
BOOT="0.6"              # Bootstrap threshold. Do NOT leave blank
REFER="-H ${ORGANISM}"   # Reference organism sequence for CAPS run
```

## Proteins list

The list of proteins is a simple text file, containing 3 columns (Box 3) that should be placed in the **proteins/** folder (e.g. PROTEIN="proteins/"). Column 1: protein UniProt identifiers; Column 2: protein names; Column 3: protein group. They should be tab separated, with **Unix line endings (LF)**, no headers, no spaces within the columns and no empty rows, except for the bottom one. Make sure you **do not have duplicates** in Column 1! Name of the file itself does not matter, and the file should not exceed 2000 rows. Therefore, if you have, say 5000 input proteins, they can be divided into 5 files (e.g. proteins1.tsv, proteins2.tsv,...) of 1000 proteins each and placed in the **proteins/** folder.

The UniProt identifiers must be from the same organism (e.g. ORGANISM="10090" for mouse), referred later as the *reference organism*. Column 3 is useful for the network analyses, as it makes possible to easily select nodes (proteins) of the same group. In case this is not necessary, just put the same identifier for all proteins (e.g. "NNN") in Column 3.

### Box 3. proteins/proteins.tsv

```
Q9EPQ1  Tlr1    TLR
Q9QUN7  Tlr2    TLR
Q9EPW9  Tlr6    TLR
Q61696  Hspala  HSP
P17879  Hspalb  HSP
```

## Species list

The list of species (e.g. SPECIES="species.tsv") is a simple text file, containing two columns (Box 4). They should be tab separated, with **Unix line endings (LF)**, no headers, no spaces within the columns and no empty rows, except for the bottom one. Column 1: species taxid codes; Column 2: species name. Make sure you **do not have duplicates!** Depending on the species, an appropriate taxonomic level should be specified in **settings.conf** (e.g. LEVEL="32523" for Tetrapoda), at which orthologues will be searched.

### Box 4. species.tsv

```
9031    Gallus_gallus
9595    Gorilla_gorilla
9606    Homo_sapiens
9993    Marmota_marmota
10090   Mus_musculus
```

## External tree

An external tree (e.g. EXTTREE="placental.nwk") should be provided if to be used as a guide by PRANK and/or PhyML. The tree should be in Newick format (nwk). Make sure the species names are exactly the same as in **species.tsv**. The script uses the taxa identifiers and will automatically convert the names to them. A suitable place to obtain an external tree is the TimeTree knowledge-base (<http://www.timetree.org/>).

## Pairs list (optional)

The list of defined protein pairs (e.g. PAIRLIST="pairs.tsv") is a simple text file, containing 2 columns (Box 5). They should be tab separated, with **Unix line endings (LF)**, no headers, no spaces within the columns and no empty rows, except the bottom one. Column 1: protein A UniProt identifiers; Column 2: protein B UniProt identifiers. This file is needed **only** if you want to define specific pairs to be searched for co-evolution (so if you set PAIRINGMANNER="defined"). By default, AutoCoEv creates all possible pairwise combinations (by PAIRINGMANNER="all") between the proteins and this list is **not** required.

### Box 5. pairs.tsv

```
Q9EPQ1  Q9QUN7
Q9EPQ1  Q9EPW9
Q9EPW9  Q9QUN7
Q61696  P17879
P17879  Q9QUN7
```

## Databases

Three databases from OrthoDB (<https://www.orthodb.org/?page=filelist>) are required. These are *all\_fasta*, *gene\_xrefs* and *OG2genes* (Box 6). The script will offer to automatically download them (see next) in the specified folder (e.g. DTB="/var/tmp/DB10v1") and run the necessary preparations. At the moment, the databases are at version 10v1 and require 30GB of disk space when extracted:

### Box 6. /var/tmp/DB10v1

odb10v1_all_fasta.tab.gz	8.8 GB (archive) → 17.1 GB (extracted)
odb10v1_gene_xrefs.tab.gz	1.1 GB (archive) → 7.3 GB (extracted)
odb10v1_OG2genes.tab.gz	1.2 GB (archive) → 5.6 GB (extracted)

## Parallelization

AutoCoEv uses GNU/Parallel for the simultaneous execution of multiple processes. By default, all detected logical cores will be used, but this can be changed if you want to keep some cores free (e.g. THREADS="6" on an 8-core CPU). Run once the following in terminal (outside the script), in order to get familiar with the bibliography information of Parallel and silence its citation notice:

```
$ parallel --citation
```

## Script run

Navigate to the AutoCoEv directory, open terminal there and start the main script:

```
$ bash ./start.sh
```

The script will check if all required executables are in place and will ask you to verify the working directory (e.g. TMP="/var/tmp/work"). The menu of AutoCoEv is simple: it presents the different steps of the pipeline, as an enumerated list of choices. Typing the corresponding number and pressing ENTER will run the respective step. In fact, the whole workflow can be carried out by simply pressing 1, 2, 3 ...

AutoCoEv will offer to run several preparations, such as databases retrieval and processing (Box 7). Once these have been set up, you can skip the preparations menu next time you run the script, by going straight to **step 11**.

### Box 7. Initial preparations menu:

```
Prepare databases or skip [11]:
1) Download odb10v1_gene_xrefs      7) Extract odb10v1_all_fasta
2) Download odb10v1_OG2genes        8) Index odb10v1_all_fasta
3) Download odb10v1_all_fasta       9) Trim odb10v1_gene_xrefs.10090
4) Check MD5sum of databases        10) Trim odb10v1_OG2genes.32523
5) Extract odb10v1_gene_xrefs       11) [DONE AND CONTINUE]
6) Extract odb10v1_OG2genes
```

- **Steps 1-3** download the archived databases from OrthoDB to the specified location.
- **Step 4** checks the MD5SUMs of the databases
- **Steps 5-7** extract the downloaded databases. Make sure you have enough space.
- **Step 8** creates an index file for odb10v1\_all\_fasta.tab.
- **Step 9** extracts from the odb10v1\_gene\_xrefs.tab databases the entries of the specified and reference organism (e.g. ORGANISM="10090"), creating a sub-database.
- **Step 10** extracts from odb10v1\_OG2genes.tab, the entries of the specified level (e.g. LEVEL="32523"), creating a sub-database.
- **Step 11 continues to the Main menu** (Box 8).

The script then verifies the correct names of databases and input files, outputting an excerpt from each to the terminal. A summary of the user-specified settings will be displayed and the main menu of the workflow will be presented (Box 8).

**Box 8. Main menu:**

```

Select a step:
1) Pair UniProt <-> OrthoDB <-> OGUUniqueID
2) Get orthologues (level: 40674)
3) Download sequences from UniProt (organism: 10090)
4) BLAST orthologues against UniProt sequence (10090, detailed: yes)
5) Get FASTA sequences of the best hits (identity: 35.000; gaps: 25)
6) [MSA] Exclude too divergent sequences
7) [MSA] Create MSA with selected method (prank)
8) [TRE] Prepare trees (auto, prank, noguide, rooted)
9) [RUN] Create pairs (all)
10) [RUN] CAPS run (alpha: 0.01, prank, auto)
11) [RUN] Inspect results and re-run CAPS (REV)
12) [RES] Generate columns stats
13) [Exit script]

```

- **Step 1** reads the list of proteins, matches their UniProt identifiers to the ones of OrthoDB and finds the orthologues group (OG) identifier of each. This step will create a folder **Orthologues/** with subfolders for each protein (Box 9). Several report files will be generated in **tsv/** (Box 10).

**Box 9. Orthologues/**

```

Q9EPQ1/      Q9QUN7/      Q9EPW9/      Q61696/      P17879/      ...

```

**Box 10. tsv/**

```

Summary.tsv           → List of matched identifiers between databases
OrthoDB_Missing.tsv   → Uniprot identifiers not found in OrthoDB
duplicates_UniProt.tsv → UniProt identifiers with more than one OrthoDB ID
duplicates_OrthoDB.tsv → OrthoDB IDs that correspond to more than one UniProt ID
duplicates_OrthoGroup.tsv → Proteins belonging to the same OrthoGroup
proteinsFound.tsv     → The entries from proteins.tsv that were found in ODB

```

- **Step 2** prepares a homologues list of each protein for the user provided species. Check the individual protein folders in (Box 9) for details, such as species where homologues were found (**\*.speciesFound.tsv**) or missing (**\*.speciesMissing.tsv**). Then, a subfolder **FASTA/** is created for each protein, where homologue sequences are collected (Box 11), named by species (taxid).

**Box 11. Orthologues/Q9EPQ1/FASTA/**

```

9031.fa      9595.fa      9606.fa      9993.fa      ...

```

- **Step 3** creates a subfolder of the reference organism (e.g. **ORGANISM="10090"**) for each protein, where the protein sequence is downloaded from UniProt. Any failed downloads will be reported in **\$TMP/tsv/UniProt.failed**, so check it to make sure everything is in place!
- **Step 4** creates a mini BLAST database for each of the downloaded sequences from step 4 (Box 12). Then runs BLAST of all sequences collected at Step 3, against the UniProt sequence from the reference organism, downloaded at step 3. The results are in table format, but BLAST will run a second time to generate detailed output (e.g. if **DETBlast="yes"**) with sequence alignments. Hits for each organism are stored in a new subfolder called **BLAST/** (Box 13).

**Box 12. Orthologues/Q9EPQ1/10090/**

```

Q9EPQ1.fa      Q9EPQ1.fa.phr  Q9EPQ1.fa.pog  Q9EPQ1.fa.pot  Q9EPQ1.fa.ptf
Q9EPQ1.fa.pdb  Q9EPQ1.fa.pin  Q9EPQ1.fa.pos  Q9EPQ1.fa.psq  Q9EPQ1.fa.pto

```

**Box 13. Orthologues/Q9EPQ1/BLAST/**

```

9031/      9595/      9606/      9993/      ...

```

- **Step 5** retrieves the sequences of the best BLAST hits from each species, that also pass the identity (e.g. `PIDENT="35"`) and gaps (e.g. `PGAPS="25"`) thresholds, specified by the user. The sequences are stored in a new folder `$TMP/BestBLASTfasta/` (Box 14) and two new report files are placed in `$TMP/tsv/` (box 15).

#### Box 14. BestBLASTfasta/

Q9EPQ1.fa      Q9QUN7.fa      Q9EPW9.fa      Q61696.fa      P17879.fa      ...

#### Box 15. tsv/

blastBestFasta.tsv      → Summary of the sequences that passed the filter  
blastBestExclude.tsv      → Summary of the sequences that did not pass the filter

- **Step 6** calls program Guidance to evaluate the similarity between orthologues for each protein and exclude sequences that are too divergent. Guidance creates a preliminary alignments in `$TMP/Guidance/` by a method of choice (e.g. `GUIDANCEMSA="muscle"`). A similarity threshold is set at runtime (e.g. `GUIDANCECUT="0.95"`) and sequences that pass the threshold are saved in `$TMP/BestGuidedFasta/` for next step. The workfolder of Guidance is named after the MSA method used for the preliminary alignment and the threshold cutoff value, e.g.: `muscle-0.95`.
- **Step 7** creates MSA by the method of choice (e.g. `MSAMETHOD="muscle"`) on the orthologous sequences from the previous step. The generated MSA is processed by Gblocks to report poor quality regions. PRANK can be run with a guide tree (e.g. `PRANKGUIDE="exguide"`). Results are stored in folder `$TMP/MSA/`, in a subfolder named after the Guidance assessment and MSA method used here (Box 16), e.g.: `muscle-0.95-prank`, for MSA created by PRANK from sequences assessed by Guidance as described above. Orthologues with sequences that do not pass the threshold will be listed in `$TMP/tsv/excluded-muscle-0.95.tsv`.

#### Box 16. MSA/muscle-0.95-prank/

Q9EPQ1.fa    Q9EPQ1.fa.10090.ref    Q9EPQ1.fa.gbl    Q9EPQ1.fa.gbl.txt    Q9EPQ1.species

Where:

**Q9EPQ1.fa** → produced alignment

**Q9EPQ1.fa.10090.ref** → alignment quality plotted onto the sequence of reference organism

**Q9EPQ1.fa.gbl** → Gblocks-filtered alignment

**Q9EPQ1.fa.gbl.txt** → Gblocks-filtered alignment in pretty format

**Q9EPQ1.species** → List of species for which an orthologue was found

- **Step 8** is optional. It calculates phylogenetic trees by PhyML from the MSAs created in Step 7 in `$TMP/Trees/`. An external tree can be provided as a guide (e.g. `PHYMLGUIDE="exguide"`). The produced trees will be placed in a subfolder named after the Guidance settings, MSA method and PhyML settings. E.g. a `muscle-0.95-prank-noguide/` folder (Box 17) would contain trees calculated from MSAs produced by PRANK, not using an external tree as a guide.

#### Box 17. Trees/

```
Trees/
├── phylml
│   ├── ext                → external tree copies trimmed for each MSA
│   └── muscle-noguide
│       ├── noroot         → produced trees will be placed here
│       ├── phy            → PhyML working directory
│       └── rooted         → TreeBeST rooted copies of the produced trees
```

**NOTE.** If you wish CAPS to generate trees (e.g. `TREESCAPS="auto"`), you can **skip Step 9**.

- **Step 9** prepares all (e.g. `PAIRINGMANNER="all"`) unique pairwise combinations between proteins, screening in the process for the number of species that are found in both MSAs (have in common). Pairs of proteins where the number of common species is below a

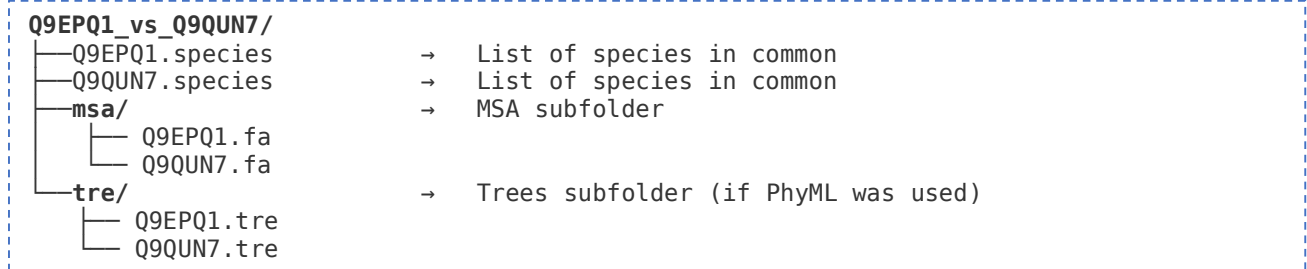
minimum threshold (e.g. `MINCOMMONSPCS="20"`) are excluded. The script calls SeqKit to remove the “unneeded” sequences from the MSA, and TreeBeST to trim the trees accordingly. Pairs that pass the minimum common species requirement are placed in `$TMP/Pairs-all/`, while the ones that do not -- in `$TMP/Pairs-all-excluded/`, under a sub-folder named after the MSA and trees conditions.

E.g. a `muscle-0.95-prank..TREE_auto/` folder contains:

- Orthologues assessed from preliminary MUSCLE alignment with Guidance cutoff 0.95
- MSAs produced by PRANK
- Trees were calculated by CAPS2 (`TRESCAPS="auto"`)

Each protein pair is placed in its individual sub-folder (Box 18).

**Box 18.** Protein pair: Q9EPQ1 vs Q9QUN7



Alternatively, you may wish to search for co-evolution between concrete protein pairs only (e.g. `PAIRINGMANNER="defined"`), in stead of all possible combinations. In this case, a list of protein pairs should be provided (e.g. `PAIRLIST="pairs.tsv"`). Pair folders names will be changed accordingly, with a `"-defined"` suffix.

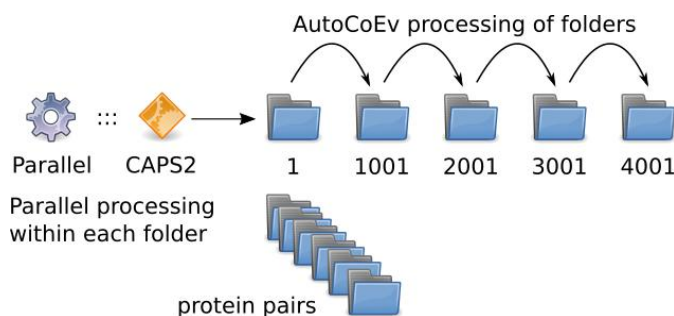
- **Step 10** carries out the actual CAPS2 run in folder `CAPS-all/`, where several folder levels reflect the run settings.

E.g. `muscle-0.95-prank..TREE_auto/Alpha0.01/` contains:

- The MSAs and trees from Step 10
- CAPS run done with run-time Alpha set to 0.01 (`ALPHA="0.01"`)

This ensures that CAPS2 runs under different conditions can be performed without the results being overwritten. GNU/Parallel has a limitation of about 2000 input strings that can be piped in. Therefore, the folders of protein pairs are further divided into groups (Box 19), for example 1000 pairs per folder (e.g. `INCR="1000"`), the maximum number being 2000. The script navigates to the first group, runs CAPS in parallel on all 1000 pairs, then moves to the next group and so on (Figure 1). Progress is logged in file `progress-?.txt`.

**Box 19.** `muscle-0.95-prank..TREE_auto/Alpha0.01/`



**Figure 1.** Sequential processing of protein pairs subfolders by AutoCoEv. The protein pairs within each subfolder are processed in parallel. When work is, e.g. Folder 1 is complete, AutoCoEv moves to the next (1001) and so forth.



- **Step 11** inspects the CAPS results and places them into folder `Results-all/`, creating three subfolders: `coev/`, `fail/` and `nocoev/` (Box 20). For proteins where co-evolution was detected, AutoCoEv runs CAPS2 a second time, this time reversing the order in which the two MSAs are loaded. E.g. A vs B, followed by B vs A. The script creates a second MSA folder, called `msa-rev`, which contains slightly renamed MSAs files, so that they are loaded by CAPS2 the other way around. Same is done for trees if they are used. Then, AutoCoEv runs CAPS2 again and after completion the results from both runs are compared. Protein pairs for which co-evolution was not detected in the second run are moved to `nocoev/`. Only amino acid pairs for which co-evolution was detected “bidirectionally” are considered. These results (files named `bothWays-?.tsv`) are processed in a series of steps of cleaning up and FDR correction by R. Protein pairs for which co-evolution was detected in both runs but not on the same amino acid pairs are moved to `noBothWays/`. In addition, AutoCoEv extracts the number of co-evolving sites for each protein in folder `Chi/`, necessary for Chi<sup>2</sup> test, see Supplementary information for detailed explanation.

**Box 20.** `Results/MSA_muscle_gblocks..PhyML_muscle_gblocks-exguide-rooted/Alpha0.01/`

<code>chi/</code>	→	Data necessary for the Chi <sup>2</sup> test
<code>coev/</code>	→	Pairs for which coevolution was detected
<code>fail/</code>	→	Pairs where CAPS run failed.
<code>noBothWays/</code>	→	Pairs where co-evolution was detected but not on the same residues
<code>nocoev/</code>	→	Pairs for which coevolution was not detected

- **Step 12** creates summary tables in tab separated values format (tsv) that can be readily imported in Cytoscape. Files are:
  - **allResidues.tsv**: results for all co-evolving residues pairs
  - **proteinpairs.tsv**: summarised results per protein pair

Consult with `README.proteinPairs` and `README.allResidues` found in `doc/examples/` for information about columns.

- **Step 13** exits AutoCoEv