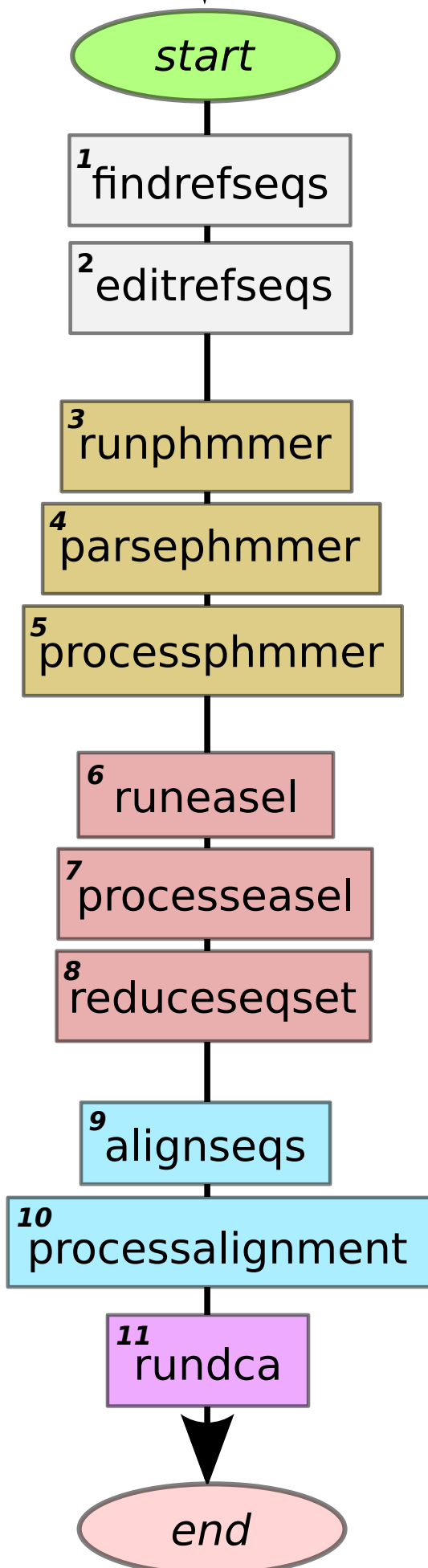


run_workflow.py

file formats

descriptions of steps



xxxx_A_refseq.fasta
xxxx_B_refseq.fasta

Finds paths to reference
fasta sequence files.
Edits refseqs by
removing nonstandard
residues.

xxxx_A_refseq_phmmer.log
xxxx_B_refseq_phmmer.log

Runs PHMMER search
on both reference seq
fastas against UniProt
(release 4.21)

xxxx_A_refseq_phmmer.keyfile
xxxx_B_refseq_phmmer.keyfile

Parses accession IDs
from both phmmer
logfiles into keyfiles

xxxx_A_refseq_phmmer_matched.keyfile
xxxx_B_refseq_phmmer_matched.keyfile

Checks for min/max
nr of accIDs per
keyfile, matches
accIDs based on
organism

xxxx_A_refseq_phmmer_matched.fasta
xxxx_B_refseq_phmmer_matched.fasta

Runs HMMER easel to
extract sequences
from Uniprot based on
accid.

Removes seqs from both
fasta files corresponding
to organisms that could
not be extracted (if any).
Removes seqs from both
fasta files that are above
a certain length (default:
1600). Does this to avoid
Muscle aln memory
errors.

xxxx_A_refseq_phmmer_matched.aln
xxxx_B_refseq_phmmer_matched.aln

Aligns sequences in
both fasta files

Joint_xxxx_A_xxxx_B_aln.fasta

Joins matched sequences
together to make a joint
alignment file

Joint_xxxx_A_xxxx_B_aln_mfdca_scores
.dat,etc.

Runs DCA on joint
alignment, gives out a
scores file