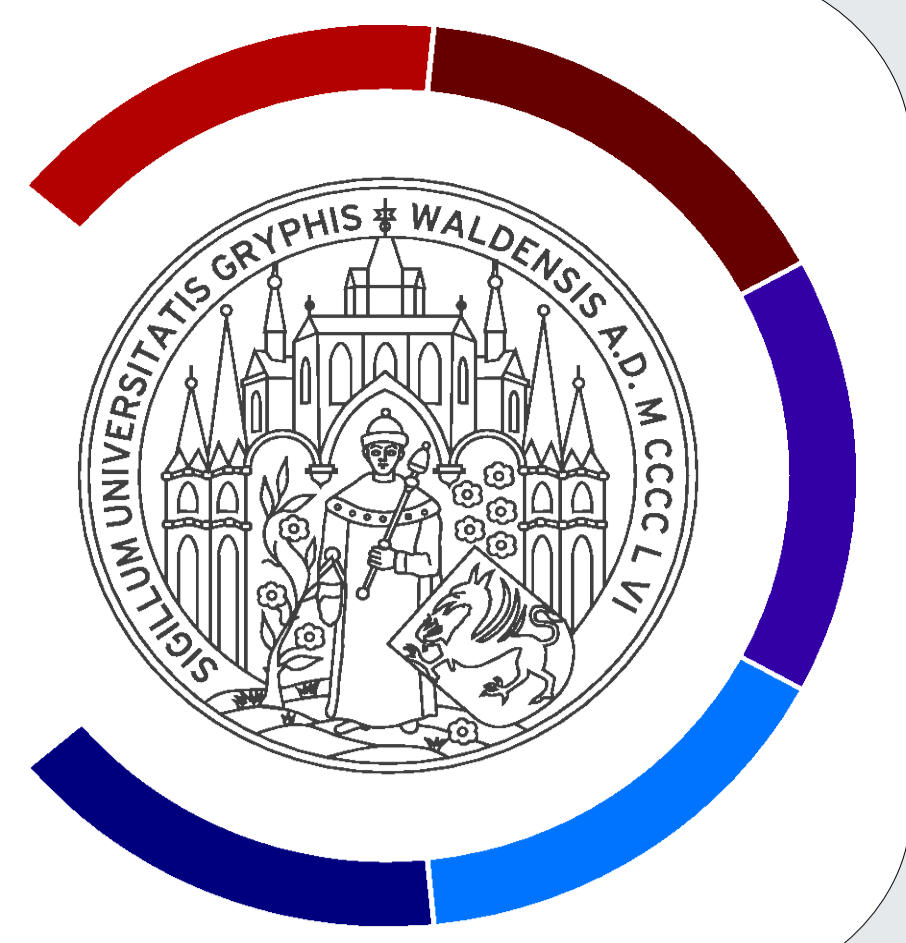# TrainAUGUSTUS – a Web Server Application for Parameter Training and Gene Prediction in Eukaryotes

Katharina J. Hoff and Mario Stanke, University of Greifswald, GERMANY. Contact: katharina.hoff@uni-greifswald.de

## Abstract

AUGUSTUS is a tool for predicting genes in eukaryotic genomic sequences [1, 2]. For achieving accurate gene predictions, a species-specific set of parameters is needed. Due to the rapidly growing number of newly sequenced genomes, an automated and easy-to-use procedure is needed in order to make gene prediction parameters for new species availabe.

Gene prediction parameters are optimized using annotated genes from the species of interest. Such initial gene sets may be generated automatically, e.g. from aligning expressed sequence tags (ESTs) to genomic sequences, or by mapping protein coding genes from other species to the genome.

We present a web server application for creating high quality training gene sets from ESTs or protein sequences. Subsequent to finding training genes, the web server application optimizes AUGUSTUS parameters and makes predictions in the supplied genomic sequence using the newly trained parameters and the supplied ESTs or protein sequences as external supporting evidence ("hints"). It is also possible to supply hints that were created externally, e.g. through manual editing, or from RNAseq data alignments.

The web server application is available at http://bioinf.uni-greifswald.de/trainaugustus

## AUGUSTUS Training Web Interface



## Training: cDNA and Genome File



## Training: Protein and Genome File



## AUGUSTUS Prediction Web Interface



## Prediction Pipeline



A project identifier is assigned to each AUGUSTUS training run. This ID may be used to call pre-trained parameters. The upload of external parameters is also possible.

*) Hints can be provided in a file, and hints are generated from provided cDNA sequence data by the prediction pipeline. Externally provided hints are treated by AUGUSTUS as "manually created", i.e. they have a stronger influence on predictions than hints that are generated by the web server application pipeline.

## Training: Gene Structure and Genome File (gff option)



## Training: Gene Structure and Genome File (gb option)



## Training: cDNA, Protein and Genome File

**Training gene generation** see "Protein and Genome File". **Hints file generation** see "cDNA and Genome File".
**Results:** *ab initio* gene predictions and gene predictions with hints.

## Training: cDNA, Gene Structure and Genome File
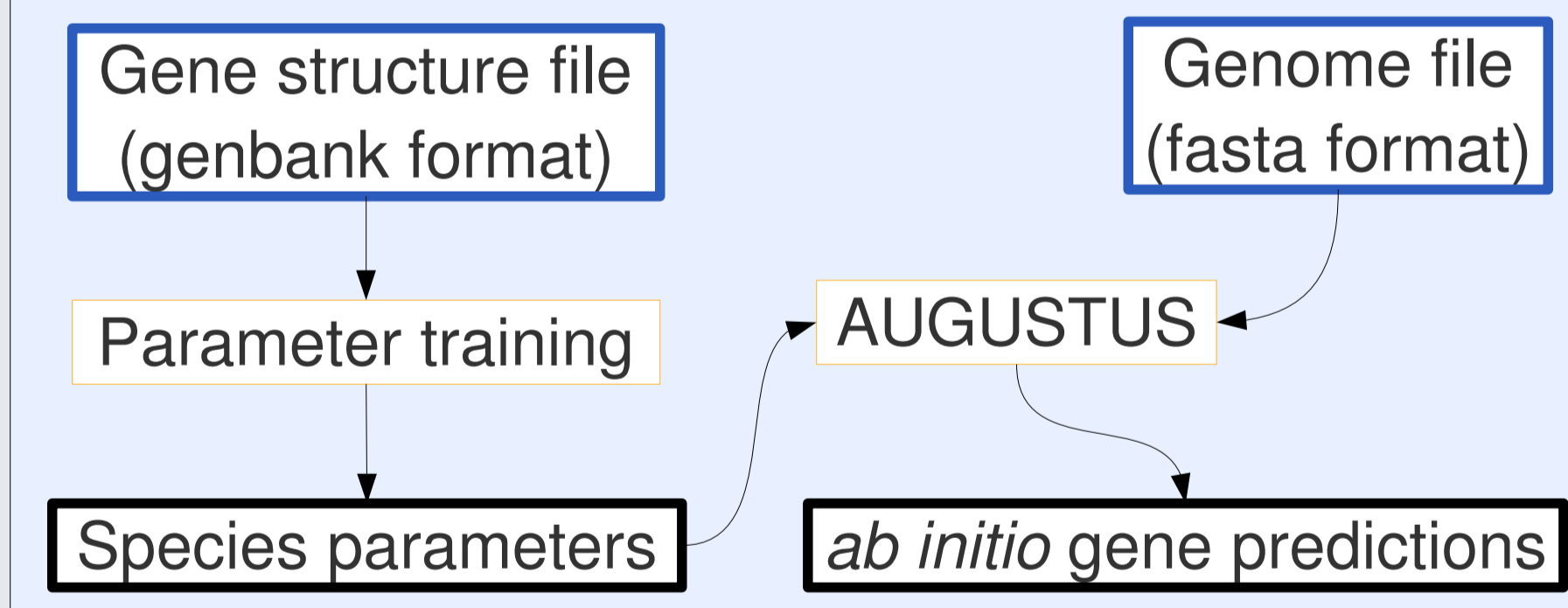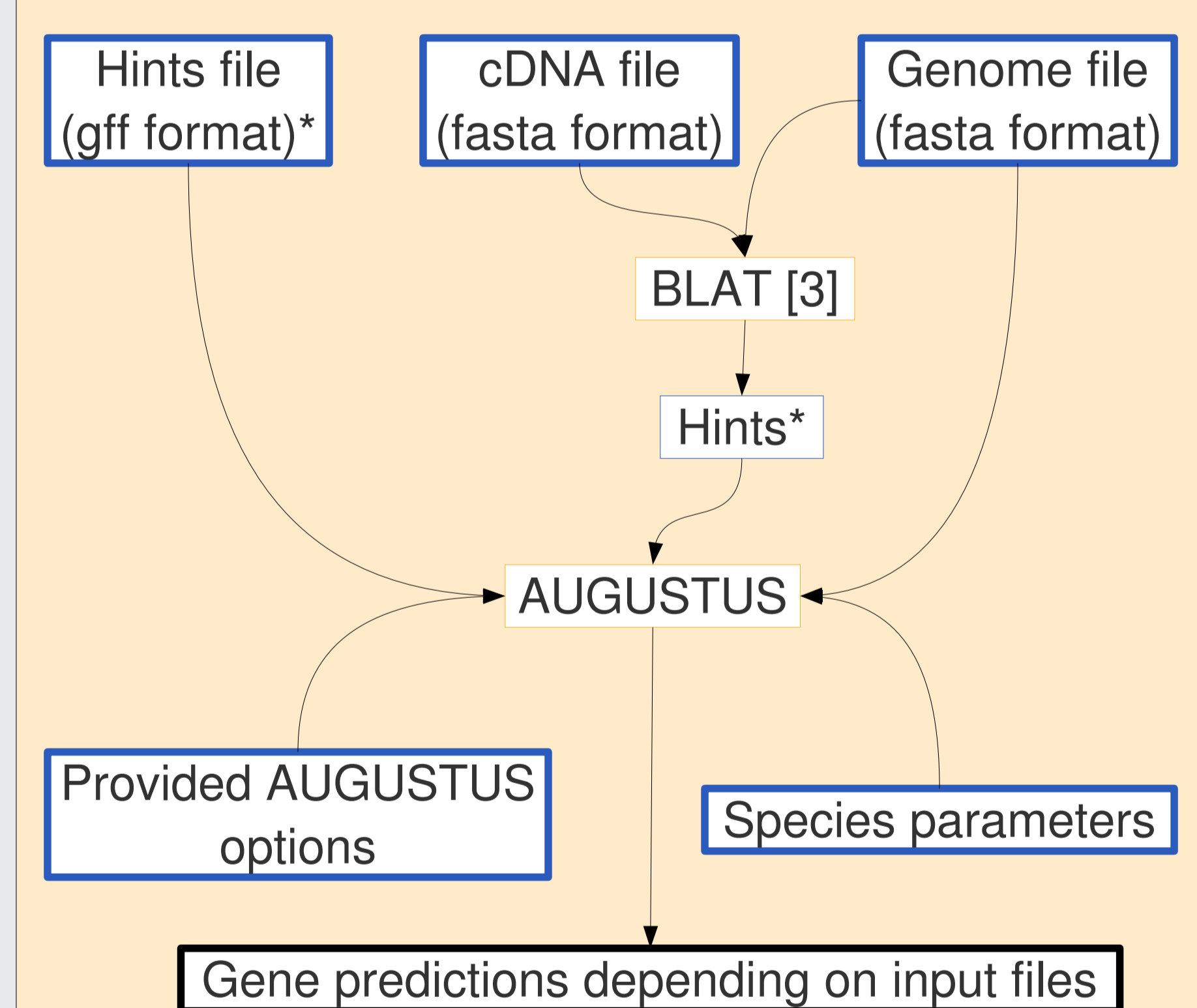
**Training gene generation** see "Gene Structure and Genome File". **Hints file generation** see "cDNA and Genome File".
**Results:** *ab initio* gene predictions and gene predictions with hints.

## Results of TrainAUGUSTUS

**Each training run produces at least the following files:**
• AutoAug.log → an event log file
• AutoAug.err → an error log file
**In addition, the following files may be produced depending on the input file combination and depending on training success:**
• parameters.tar.gz → archive with AUGUSTUS species parameters
• training.gb.gz → gzipped genbank file with produced training genes
• ab_initio.tar.gz → archive with ab initio gene predictions
• hints_pred.tar.gz → archive that contains gene predictions with hints
**Gene prediction archives** contain at least a file that contains gene predictions in gff format.
Additionally, a gtf-file and fasta files with amino acid sequences, exon sequences, coding sequences and mRNA sequences may be included.
Also a gbrowse-file may be produced.

## References

[1] M. Stanke and S. Waack (2003) "Gene prediction with a hidden Markov model and a new intron submodel", Bioinformatics, Vol. 19, Suppl. 2, pages ii215-ii225
[2] M. Stanke, M. Diekhans, R. Baertsch, D. Haussler (2008) "Using native and syntenically mapped cDNA alignments to improve de novo gene finding", Bioinformatics, 24(5):637
[3] Kent, W.J. (2002) BLAT—The BLAST-Like Alignment Tool. Genome Res, 12, 656-664.
[4] Keller, O., Odronitz, F., Stanke, M., Kollmar, M., Waack, S. (2008) Scipio: Using protein sequences to determine the precise exon/intron structures of genes and their orthologs in closely related species. BMC Bioinformatics 9, 278.
[5] Haas, B.J., Delcher, A.L., Mount, S.M., Wortman, J.R., Smith Jr, R.K., Jr., Hannick, L.I., Maiti, R., Ronning, C.M., Rusch, D.B., Town, C.D. et al. (2003) Improving the Arabidopsis genome annotation using maximal transcript alignment assemblies. Nucleic Acids Res, 31, 5654-5666.

## Funding