# SeMPI– a genome-based Secondary Metabolite Prediction and Identification web server
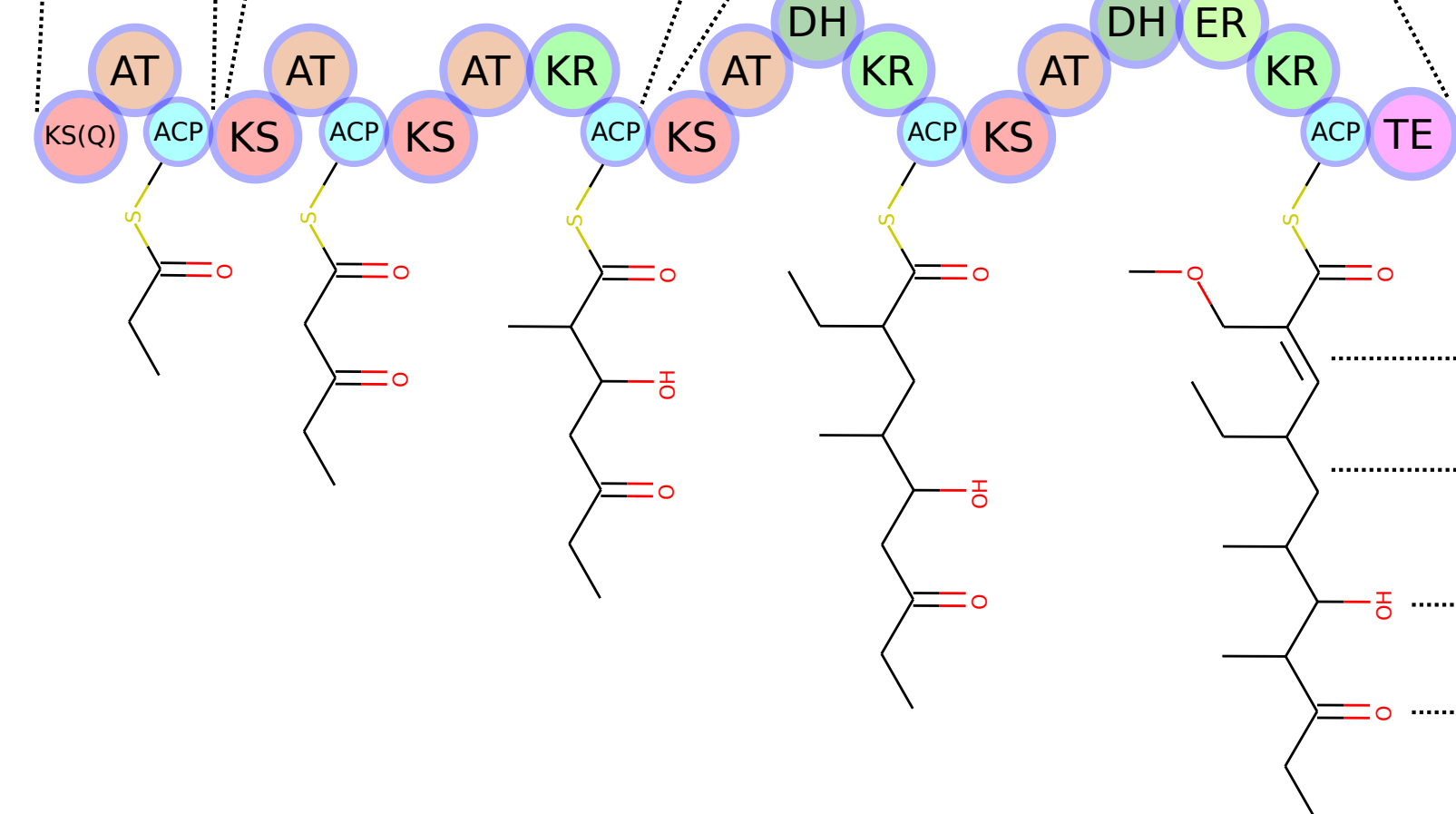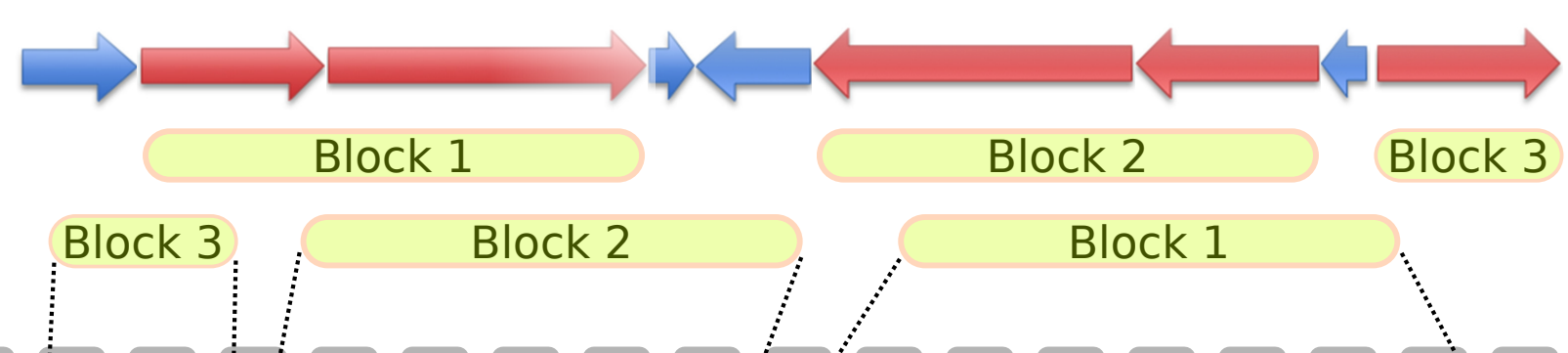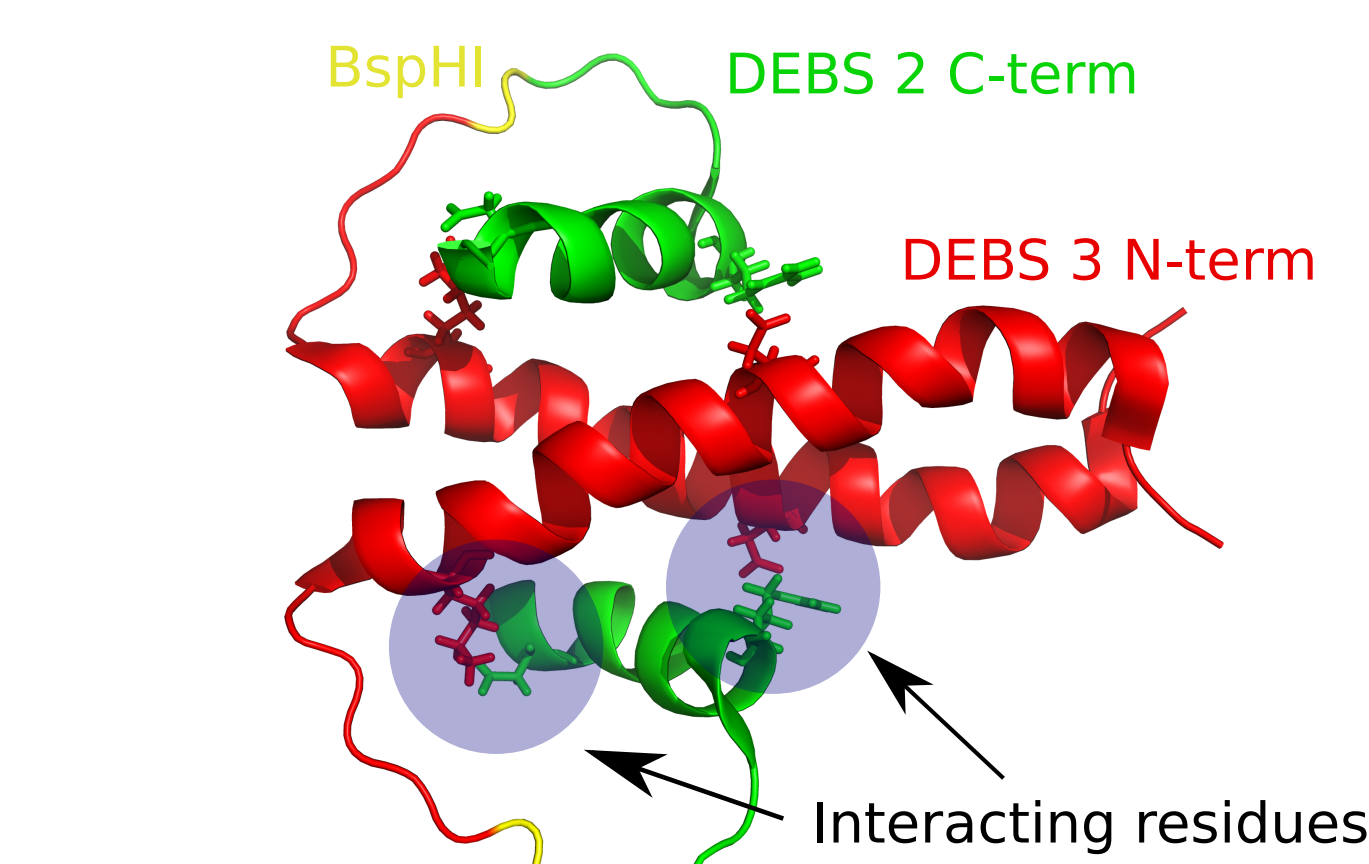
Paul Zierep, Natàlia Padilla, Dimitar Yonchev, Kiran Telukunta, Dennis Klementz and Stefan Günther

Pharmaceutical Bioinformatics, Institute of Pharmaceutical Science, Germany

paul.zierep@pharmazie.uni-freiburg.de

The secondary metabolism of bacteria, fungi and plants yields a vast number of bio-active substances. The constantly increasing amount of published genomic data provides the opportunity for an efficient identification of gene clusters by genome mining. Conversely, for many natural products with resolved structures, the encoding gene clusters have not been identified yet. Structural elucidation of the actual secondary metabolite is still challenging, especially due to as yet unpredictable post-modifications. Here we introduce SeMPI, a web server providing a prediction and identification pipeline for natural products synthesized by polyketide synthases (PKS) of type I modular.
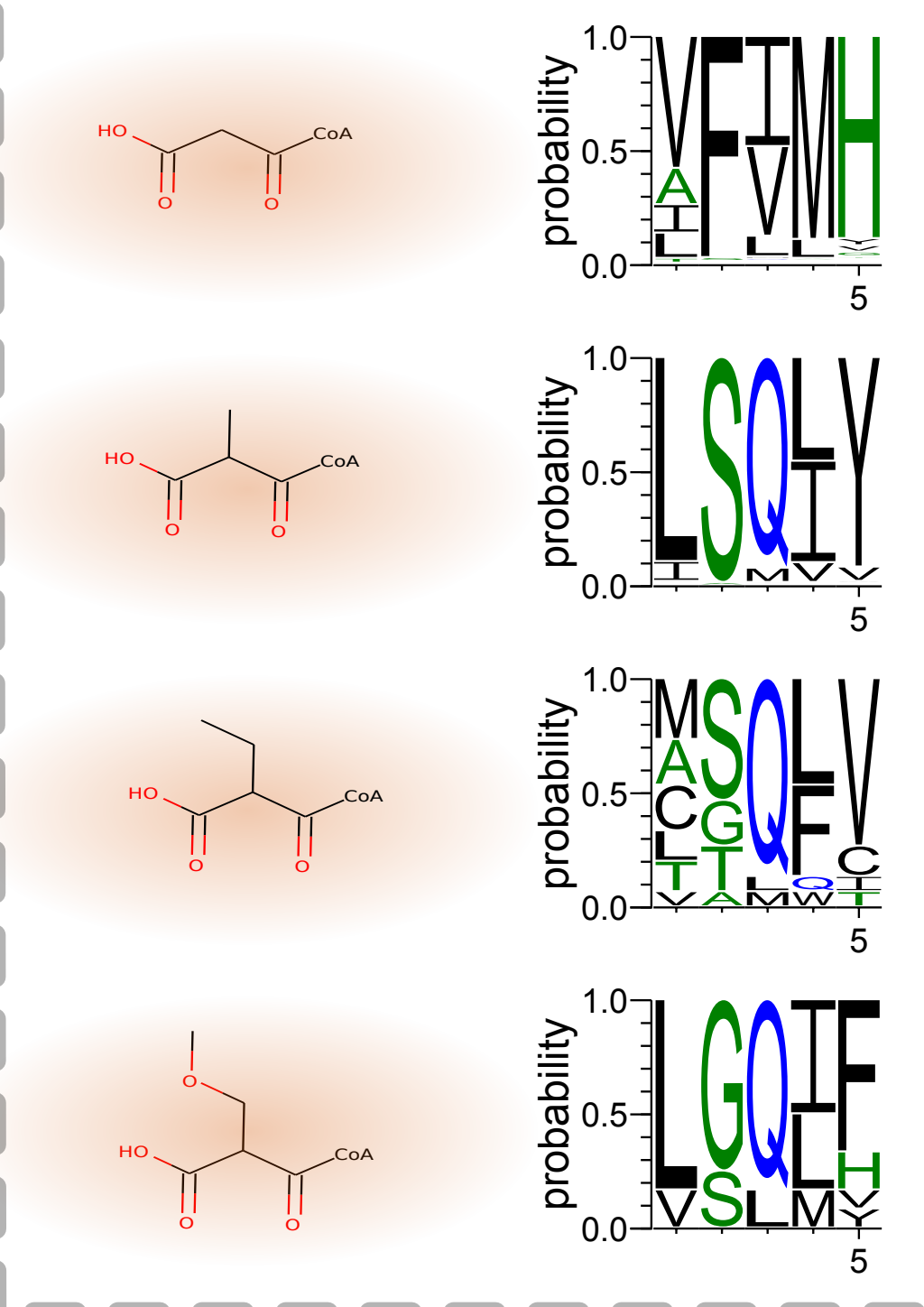
## Chain prediction

>New Gene Sequence
VAMVFPGQGAQWQGMARDLLRESQVFAD ...



a) Domain recognition and ordering
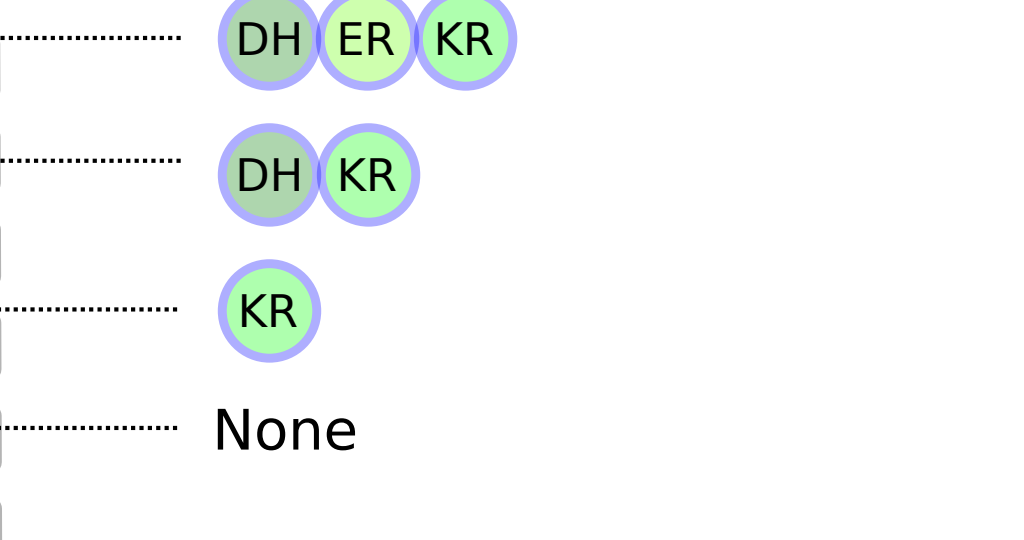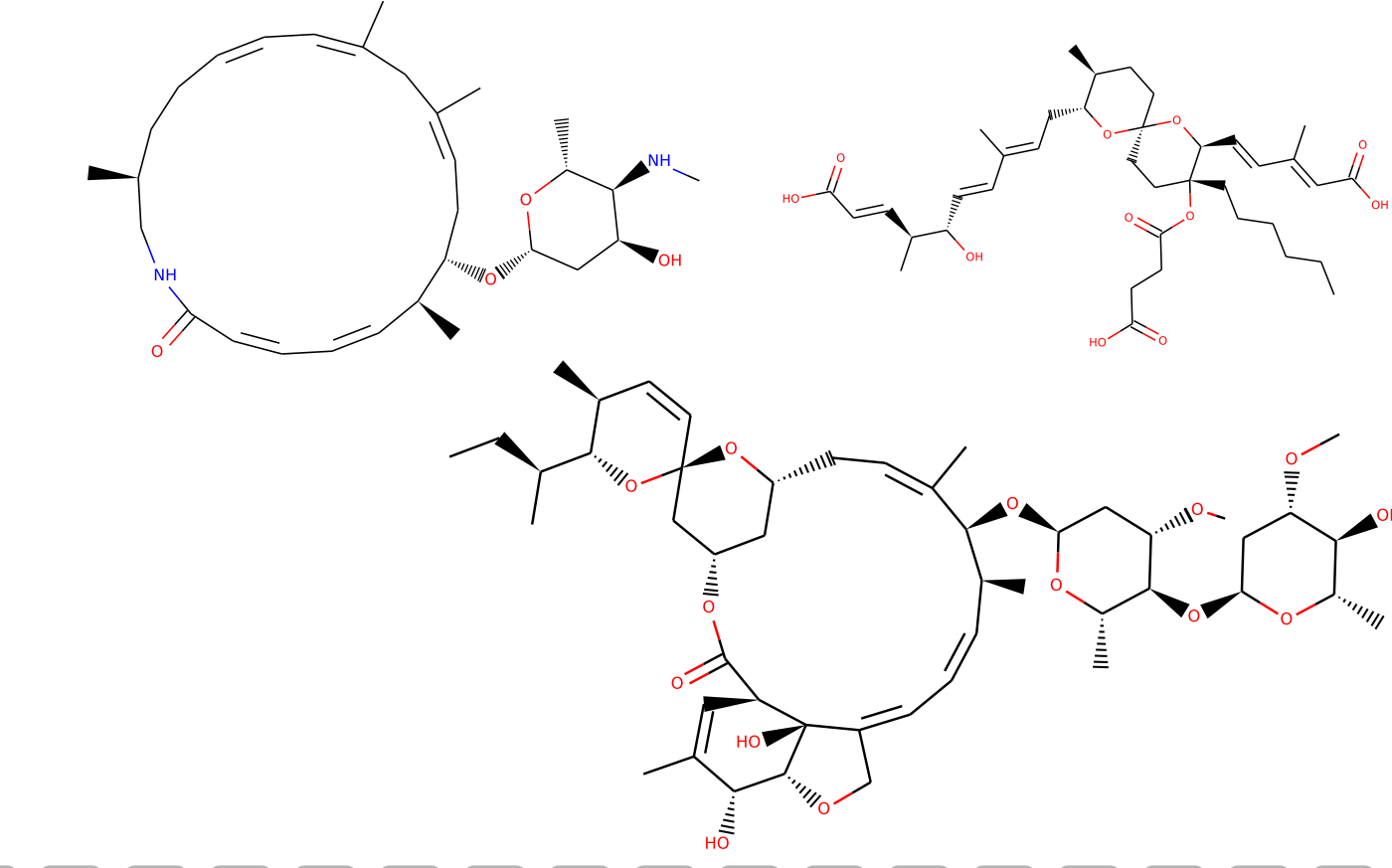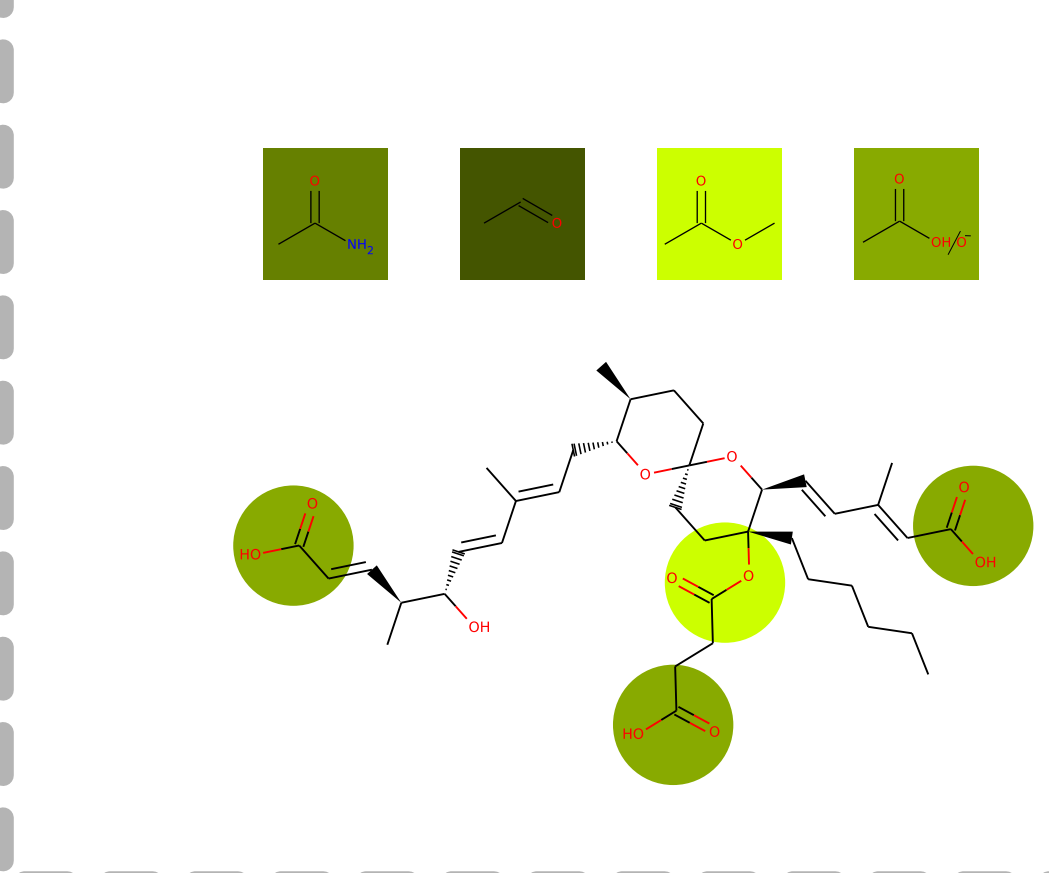b) AT substrate specificity
c) β-keto reduction profile

**Figure 1.** The gene cluster boundaries are based on the antiSMASH 3.0 [1] annotation. Domain signatures are identified using profiles Hidden Markov Models and ordered based on interacting residues, determined by alignment of the docking domains with 6-deoxyerythronolide B synthase (DEBS) (a). Acetyltransferase (AT) specificity is predicted using position probability matrices (b) and β-keto modification is based on presence of reductive domains (c). This information is used to predict the basic PK-chain.
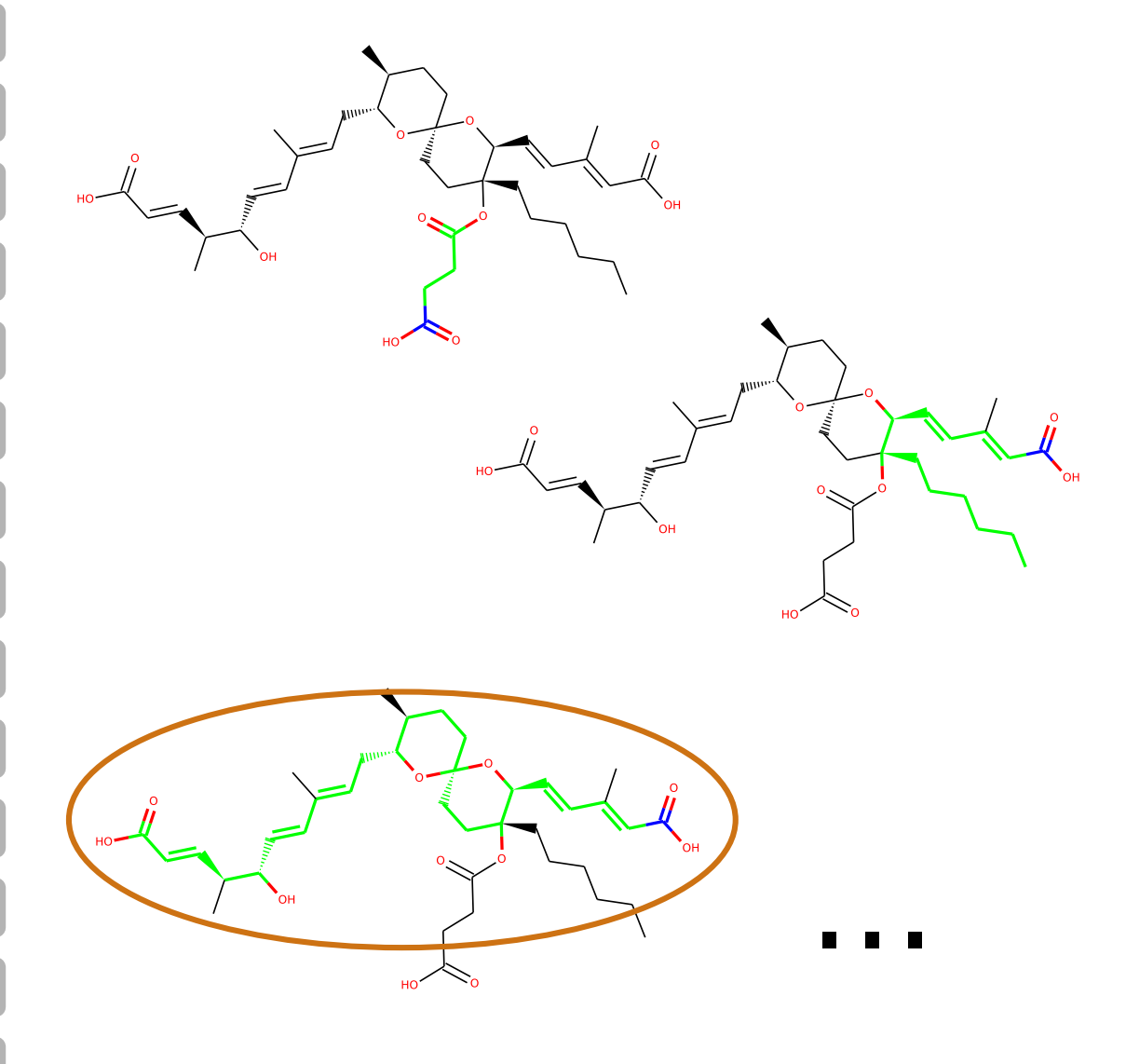
## Database screening



d) Metabolite database
e) Starting unit identification
f) Path generation
g) Matrix annotation
h) Matrix scoring and metabolite ranking

**Figure 2.** A database [2] with more then 4000 diverse molecules is used for chain comparison (d). Based on conserved starting units (e) a set of paths, representing the putative initial biosynthesized carbon chain of the metabolite is computed for each molecule (f). All paths a stored in a specific format, referred as matrix (g) and compared with the predicted initial chain (h).

## Benchmark



h) ROC curve
i) Individual ranking

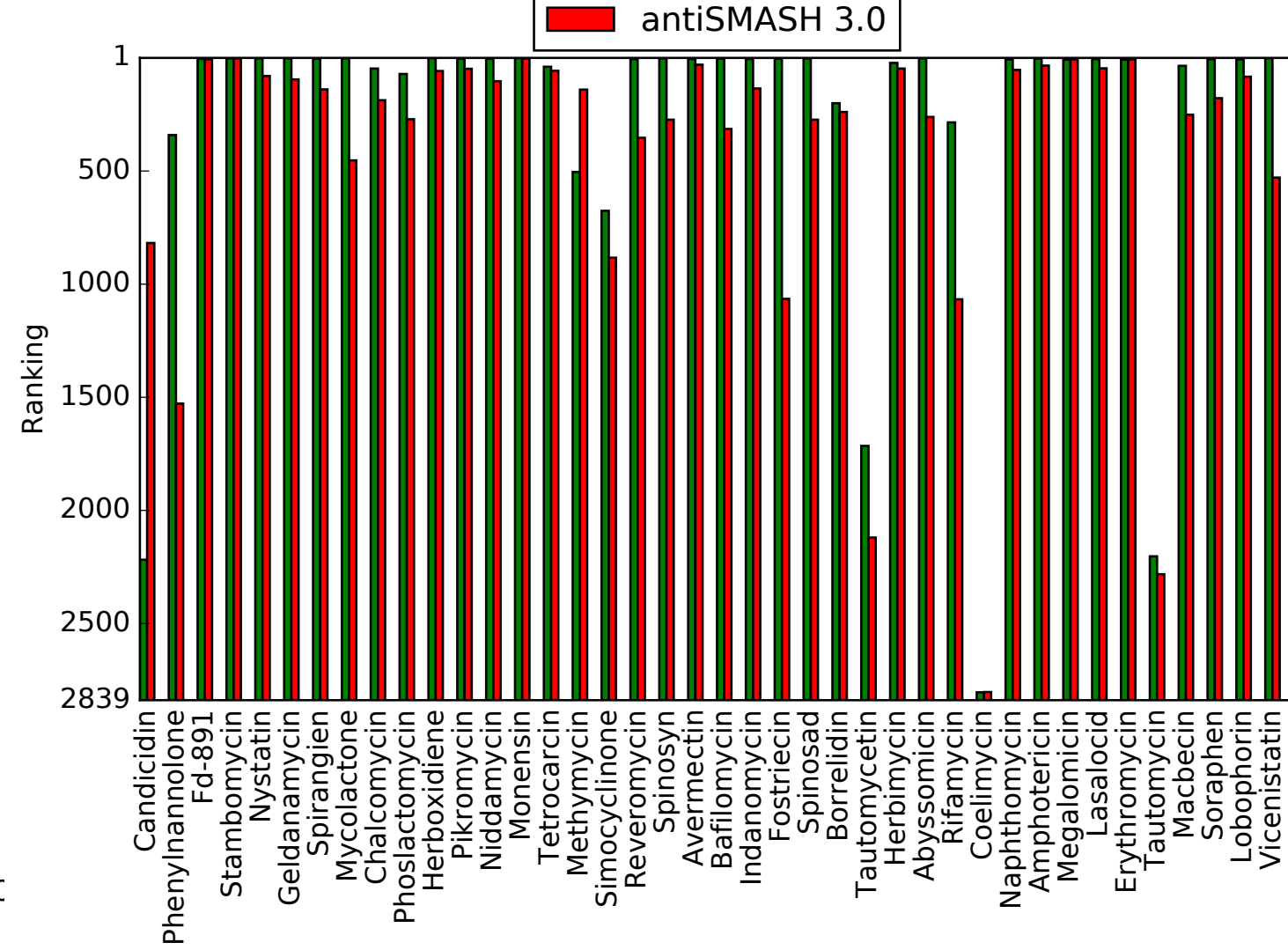**Figure 3.** Based on the 40 gene clusters from the test data-set, SeMPI reached an AUC-value of 0.85, compared to antiSMASH with an AUC-value of 0.69 (l). In the individual ranking of the gene cluster products among natural products from the StreptomeDB, SeMPI ranked 27 correct gene cluster products (antiSMASH: 5) within the first ten of 2839 possible ranks (m).

## Website output



Results

**Figure 4.** Besides comprehensive information about the predicted chain, SeMPI also lists matching molecules form the database, based on a ranking of their best scoring paths. This provides researchers with extensive data about the investigated gene cluster and might help to identify the actual secondary metabolite or to indicate the possibility for a not-yet annotated compound.

## References

[1] Weber, T., Blin, K., Duddela, S., Krug, D., Kim, H.U., Bruccoleri, R., Lee, S.Y., Fischbach, M.A., Muller, R., Wohlleben, W. et al. (2015) antiSMASH 3.0-a comprehensive resource for the genome mining of biosynthetic gene clusters. Nucleic Acids Res., 43, W237-243.

[2] Klementz, D., Doring, K., Lucas, X., Telukunta, K.K., Erxleben, A., Deubel, D., Erber, A., Santillana, I., Thomas, O.S., Bechthold, A. et al. (2016) StreptomeDB 2.0--an extended resource of natural products produced by streptomycetes. Nucleic Acids Res., 44, D509-514.

**Pharmazeutische Bioinformatik**

www.pharmaceutical-bioinformatics.de

DFG

C-factors RTG1976

UNI FREIBURG