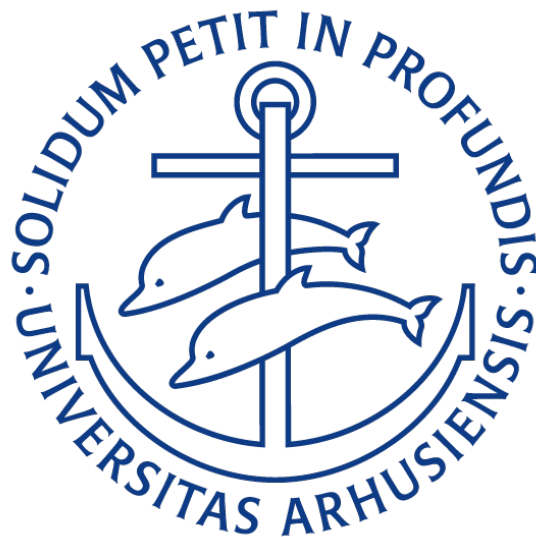# Antibiotic potential of soil bacteria targeting ESKAPE pathogens

**Laura Marie Hagen Fur**

MSc Student in Bioinformatics

Bioinformatics Research Center, Aarhus University

**Spring 2025**

## Supervisors

**Thomas Bataillon**

Bioinformatics Research Center

**Thomas Tørring**

Biological and Chemical Engineering

**Aarhus University**

## Abstract

Antibiotic resistance is a growing global concern, with the emergence and spread of antibiotic-resistant bacteria continuing to increase. Antibiotics are small metabolites naturally produced by microorganisms to ensure their survival in competitive environments by inhibiting or killing other competitive microorganisms. Conventional treatment options are getting exhausted, calling for an urgent need to discover novel antibiotics.

Here, I apply an *in silico* approach to discover natural products with potential antibacterial effects against ESKAPE pathogens by utilizing the bioinformatics pipeline, antiSMASH, to mine genomes of bacterial soil isolates for biosynthetic gene clusters (BGCs). In the study, unsupervised and supervised learning methods along with phylogenetic analysis, were integrated to find associations between observed bioactivity and the presence of specific BGCs. A total of 558 BGCs were predicted by antiSMASH, of which several BGCs were flagged in the analyses after correlating these with the results from the bioactivity assays. Some of these BGCs have reported antimicrobial activities with proposed broad- and narrow-spectrum activity, such as epilancin 15x, paenilamicin, and tridecaptin, which highlights the promising potential of combining supervised learning methods with laboratory work to the discovery of novel antibiotics.

# Contents

# Acronyms