

The University of South Bohemia in České Budějovice
Faculty of Science

**Searching for the Polyenes in *Streptomyces* Associated with
Bark Beetles**

Master's Thesis

Bc. Sara Khaled Youssef Sayed

Supervisor: RNDr. Alica Chronáková, PhD.

Co-supervisor: Ana Catalina Lara Rodriguez, PhD.

České Budějovice

2023

Genome Screening of antifungal polyenic secondary metabolites

1. Research Objective:

Data mining from polyene databases was carried out to search for a specific gene marker that is essential for the synthesis of polyenes. This was accomplished by a deep understanding of the biosynthetic gene clusters (BGCs) of the polyenes produced by *Streptomyces sp.* and evolutionary related actinomycetes. The first hypothesis was that Cytochrome P450 could be used as a specific gene marker for polyene production. The main aim of this work was to evaluate the matching of the primer pair that was designed in a previous study with CYP 450 sequences of a more robust dataset of polyenes. Furthermore, an alternative hypothesis was formulated after conducting an extensive data mining to search for a specific gene marker. The hypothesis was that thioesterase type I is a good gene marker for the identification of polyene-producing actinomycetes, specifically streptomycetes.

2. Methodology:

2.1 Data mining from polyene databases

Since polyenes have a common structure essential for their functions, they have common genes which may act as gene markers (Figure 1). Data mining from polyene databases such as the National Center for Biotechnology Information (NCBI), Minimum Information about a Biosynthetic Gene cluster (MIBiG) and Database of Biosynthesis clusters Curated and Integrated (DOBISCUIT), was carried out to identify a specific gene marker required for the production of Polyenes.

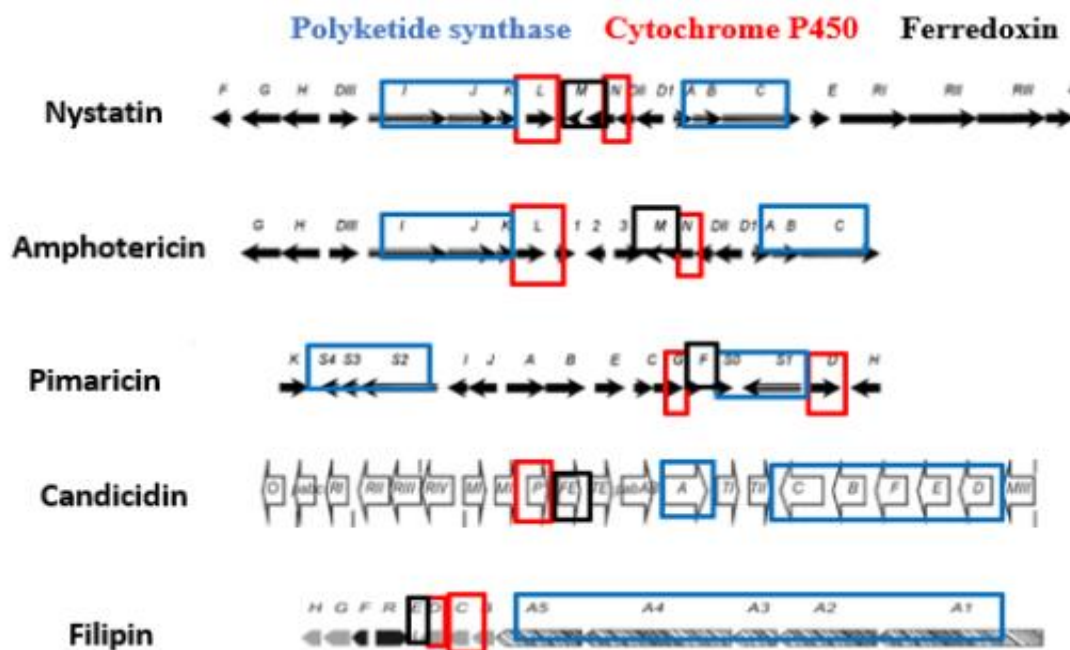


Figure 1: Illustration for the biosynthetic gene cluster of the polyenes with common genes highlighted.

• Cytochrome P450 as a gene marker

Based on a previous study, where a primer pair was specifically designed for cytochrome P450

and utilized for PCR-based genome screening to isolate potential polyene-producing actinomycetes strains.

Firstly, nucleotide and amino acid sequences of polyene BGCs were downloaded from publicly available databases such as DOBISCUIT (Accessed in October 2022) and MIBiG (Accessed in October 2022). Then a FASTA file was created for the amino acid sequence of CYP 450 of 6 polyenes. Multiple sequence alignment of CYP 450 using the graphical interface of Geneious version 2022.0.2 (Biomatters Ltd., Auckland, New Zealand) was done using Muscle (Version 3.5). The aligned sequences were edited manually by removing the gaps and used for phylogenetic tree construction using Molecular Evolutionary Genetics Analysis (MEGA) software version 11.0.13. Maximum likelihood statistical method and 100 bootstrap replications were selected.

● **Thioesterase type I as a gene marker**

After conducting an extensive data mining process to search for another gene marker, the alternative hypothesis was that thioesterase type I could serve as an effective gene marker for identifying polyene-producing actinomycetes, particularly streptomycetes. It was based on a previous literature that highlighted the presence of a distinct group of PKS-TE domains responsible for producing polyene macrolides. These domains exhibited relatively low sequence similarity compared to TEs of other non-polyene macrolides such as Erythromycin and Pikromycin.

Firstly, nucleotide sequences of polyene biosynthetic gene clusters were downloaded from JDB, NCBI and MIBiG. Then FASTA file was created for nucleotide sequences of the last module of type I polyketide synthase that associated with thioesterase type I, which included 26 polyenes and 2 non-polyenes (Pikromycin and Erythromycin). The sequences were checked that they are starting with start codon and have stop codon at the end. Then the sequences of the last module were uploaded on AliView software (Alignment Viewer and Editor). The sequences were translated into amino acid sequences using bacterial genetic code and translation frame 1. Then multiple sequence alignment of the last module of polyketide synthase using AliView was performed using Muscle and selecting the option "Realign everything as Translated Amino Acids". The aligned amino acid sequences were checked, edited manually and the gaps were removed. Then the amino acid sequences of ACP-TE were annotated using antiSMASH bacterial version and extracted from the whole module. The results of the alignment of amino acid sequences of the thioesterase (TE) gene were compared to the results provided by Zhou, Yucong, et al.. The Conserved region of TE for the polyenes and non-polyenes was extracted and the large gaps that occupied only by non-polyenes were removed. Then the Amino acid sequences were translated back to nucleotide sequences and used for designing the forward and reverse Primers.

2.2 Designing of primer pairs

Ten primer pairs were designed on the consensus sequences using Geneious software version R8.1.9 (Biomatters Ltd., Auckland, New Zealand) by Clicking the Primers button and selecting Design New Primers. The parameters were adjusted as shown in Figure 2 for optimization of the primer design. The optimum product size was adjusted to 500 bp. Then the best option for the primer pairs was selected based on the % GC content, primer melting temperature (T_m), Hairpin T_m and Self-Dimer T_m and used for designing degenerate primer pairs. Moreover, it is important that the sequence of the primers have guanine (G) or cytosine (C) within the last 5 nucleotides at the 3' end as they form 3 hydrogen bonds thus better stability. This is known as GC clamp.

Characteristics

Primer DNA Probe

Size Min: 12 Optimal: 18 Max: 27

Tm Min: 50 Optimal: 60 Max: 65

%GC Min: 40 Optimal: 60 Max: 80

Product Tm Min: 0 Optimal: 0 Max: 0

Max Tm Difference: 5

GC Clamp: 1

Max Dimer Tm: 55

Max Poly-X: 3

Max 3' Stability: 9

☐ Allow primers inside target with penalty: 0

Figure 2: Optimization of the primer's properties.

Moreover, the degenerate primer pair was designed manually using Geneious software version 8.0.4 (Biomatters Ltd., Auckland, New Zealand) by following the procedure that mentioned on Geneious website. Firstly, the file with the aligned sequences was selected. Then the consensus bottom was selected, and the threshold was set to 50 % at which the bases matching was at least 50% of the sequences. The highlighting of the sequences was adjusted as depicted in Figure 3 for easy observation of the differences.

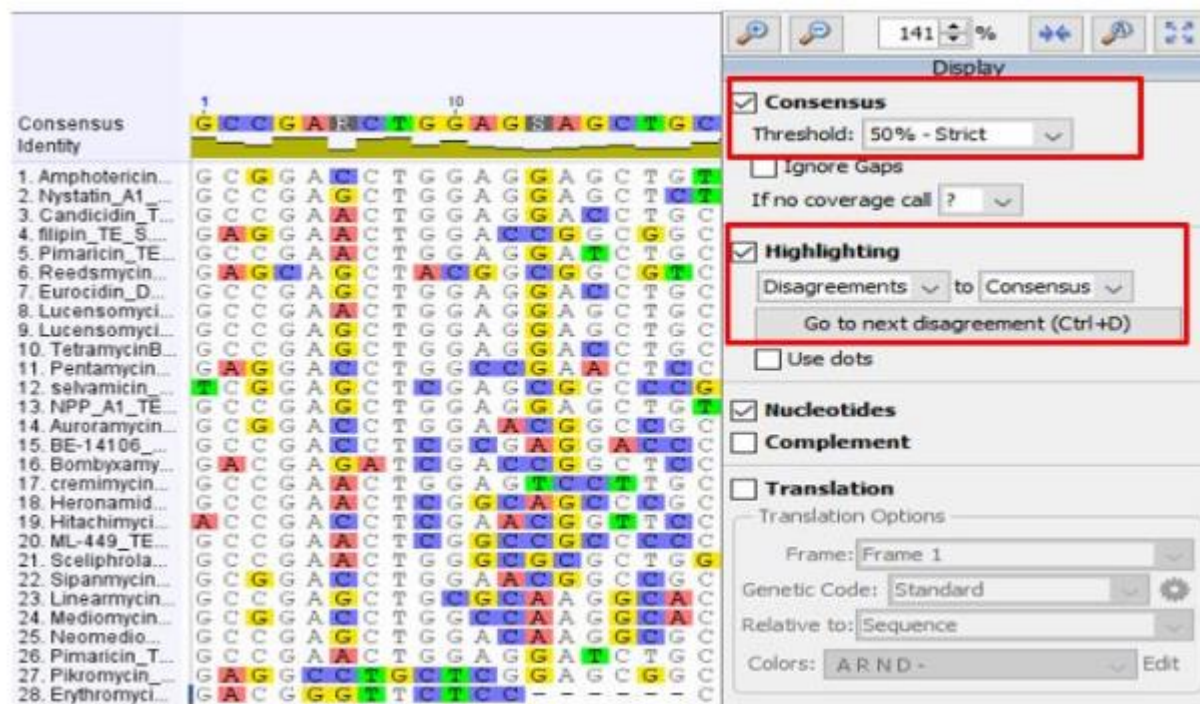


Figure 3: Screenshot illustrating the settings for designing degenerate primers.

The regions of the best option out of the 10 primer pairs that were designed before were selected on the consensus sequence in the alignment. For the forward primer, the region was selected from left to right while the region of the reverse primer was selected from right to left. The Add Annotation button was selected and the calculated characteristics of the primer, including the T_m range based on the degeneracy were checked. Then the name of the primer was entered, and the primer was added to the consensus as an annotation by clicking Ok.

2.3 In silico PCR

In silico PCR, also known as virtual PCR, is a computational method used to simulate the polymerase chain reaction (PCR) process. An online tool called “In silico PCR amplification” (<http://insilico.ehu.es/PCR/index.php?mo=Streptomyces>, accessed in October 2023) was used for testing the efficiency and specificity of the designed primer pair. The primer pair was uploaded and tested against all available *Streptomyces* strains allowing up to 2 mismatches, but in 1 nucleotide in 3' end. The maximum lengths of the bands were adjusted to 600 nucleotides (Figure 4).

In silico PCR amplification

[Input primers in fasta format](#)

Primer 1¹ 5'- -3' [C](#)

Primer 2¹ 5'- -3' [C](#)

Microorganism
 ▼

☒ Include plasmids (if available)

Allow mismatches, but in nucleotides in 3' end

Maximum length of bands
 nucleotides

¹ Degenerated nucleotides are allowed; A+T+G+C must be 10 or more.

[Info](#)

Figure 4: Screenshot illustrating the settings for testing the primer pair.

3. Results:

3.1 Data mining from polyene databases

Robust dataset of 26 linear and cyclic polyenes (macrolactone and macrolactam compounds) was created and used for the analysis of cytochrome P450 (CYP 450) and thioesterase type I (TE I) genes. They are listed in groups based on their structures together with bacterial strain, classification, database, and biological activity (Table 1).

Table 1: List of polyene secondary metabolites produced by *Streptomyces* sp. and phylogenetically close actinomycetes that were used in this study.

Compound	Bacterial strain	Classification	Biological activity	Database	Used for analysis of CYP 450	Used for analysis of TE
Macrolactone						
Amphotericin	<i>Streptomyces nodosus</i> -ATCC 14899	Heptaenes	Antifungal	JDB	✓	✓
Candididin	<i>Streptomyces</i> sp. FR-008	Heptaenes	Antifungal	Mibig	✓	✓
Filipin	<i>Streptomyces avermitilis</i> MA-4680	Pentaenes	Antifungal and low anti-	Mibig	✓	✓

			bacterial activity			
Nystatin	<i>Streptomyces noursei</i> ATCC 11455	Tetraenes	Antifungal	JDB	✓	✓
Pimaricin	<i>Streptomyces Natalensis</i> ATCC-27448	Tetraenes	Antifungal and anti-protozoal	JDB	✓	✓
Rimocidin	<i>Streptomyces diastaticus</i> 108	Tetraenes	Antifungal	JDB	✓	
Reedsmycins	<i>Streptomyces</i> Sp. OUC6819	Pentaenes	Antifungal	Mibig		✓
Eurocidin-D	<i>Streptomyces eurocidicus</i> ATCC 27428	Pentaenes	Antifungal	Mibig		✓
Lucensomycin	<i>Streptomyces achromogenes</i> _NBRC-14001 <i>Streptomyces cyanogenus</i>	Tetraenes	Antifungal and antiviral activity.	Mibig		✓
Tetramycin	<i>Streptomyces hygrospinosus</i> CGMCC 4.1123	Tetraenes	Antifungal	Mibig		✓
NPP	<i>Pseudonocardia</i> -sp. AL041005-10	Tetraenes	Antifungal	NCBI		✓
pentamycin	<i>Streptomyces</i> sp. S816	Pentaenes	Antifungal	NCBI		✓
Macrolactam						
Auroramycin	<i>Streptomyces filamentosus</i> strain NRRL 15998	Pentaenes	Antifungal, Anti-MRSA	Mibig, NCBI		✓
ML-449	<i>Streptomyces</i> sp. MP39-85	Tetraenes	Anti-bacterial, Cytotoxicity	Mibig		✓
BE-14106	<i>Streptomyces</i> sp. DSM-21069	Tetraenes	Anti-bacterial, Cytotoxicity	Mibig		✓
Sipanmycin	<i>Streptomyces</i> sp. CS149	Pentaenes	Cytotoxicity	Mibig		✓

Hitachimycin	<i>Streptomyces scabrisporus</i> CMB-0406	Triene	Anti- protozoal antifungal	Mibig		✓
Heronamide_C	<i>Streptomyces</i> <i>sp.</i> CMB-0406	Tetraenes	No anti- bacterial activity and cytotoxicity	Mibig		✓
Sceliphrolactam	<i>Streptomyces</i> <i>sp.</i> SD85	Pentaenes	Antifungal	NCBI		✓
Bombyxamycin A	<i>Streptomyces</i> <i>sp.</i> SD53	Pentaenes	Anti- bacterial, cytotoxicity	NCBI		✓
cremimycin	<i>Streptomyces</i> <i>sp.</i> MJ635-86F5	Triene	Anti- bacterial (G+), cytotoxicity	NCBI		✓
Linear						
Linearmycin A	<i>Streptomyces</i> <i>sp.</i> Mg1		Antifungal	Mibig		✓
Mediomycin A	<i>Streptomyces</i> <i>blastmyceticus</i>		Antifungal	Mibig		✓
Neomediomycin	<i>Streptomyces</i> <i>sp.</i> RK95-74		Antifungal	Mibig		✓

3.2 Analysis of Cytochrome P450

Amino acid sequences of cytochrome P450 from 6 polyene biosynthetic gene clusters (BGCs) were used for multiple sequence alignment as depicted in Figure 5. The majority of those polyene BGCs contain two cytochrome P450 genes, with the exception of Candicidin and Rimocidin BGCs.

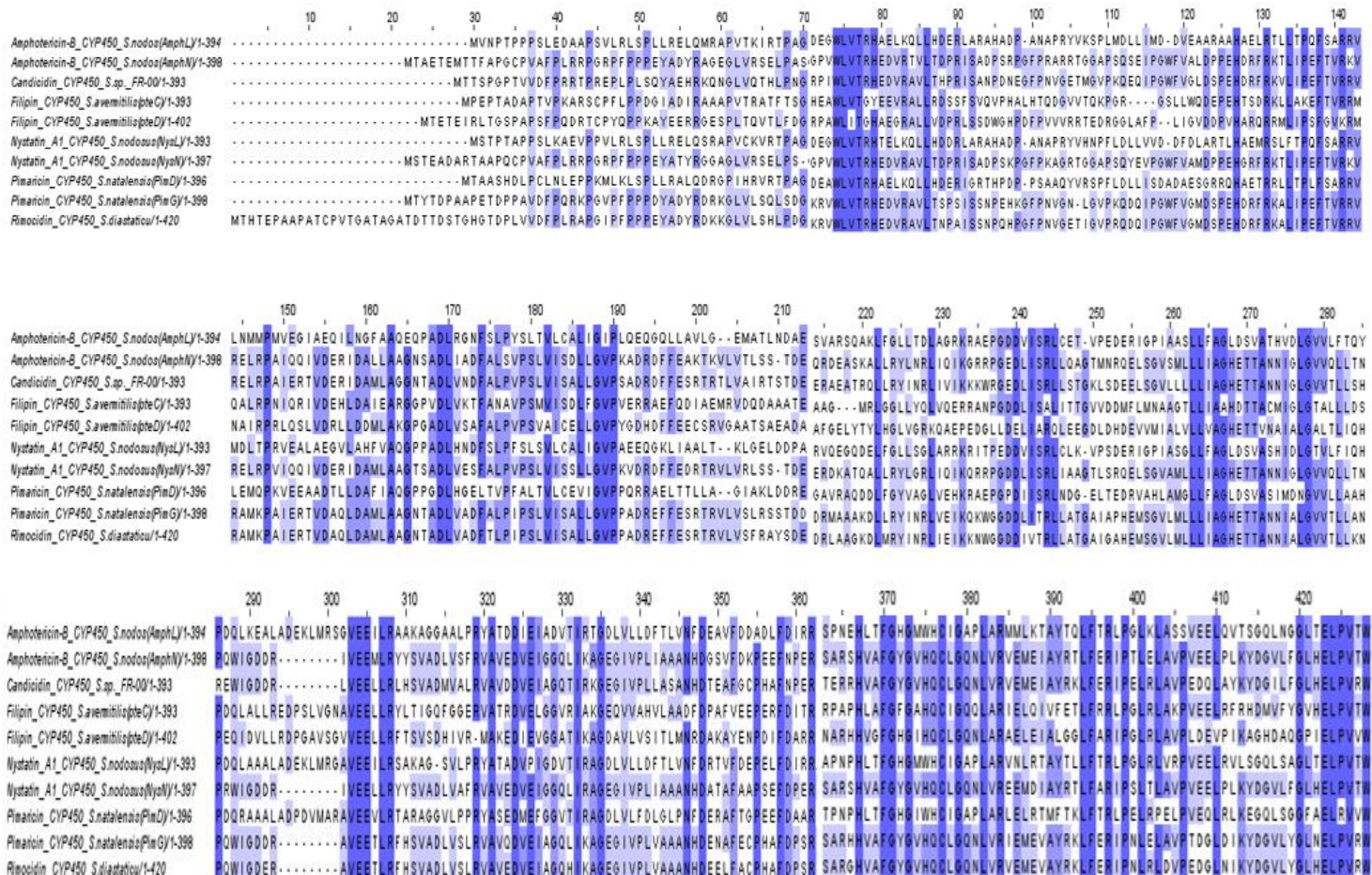


Figure 5: Multiple sequence alignment of the amino acid sequences of cytochrome P450. The highly conserved amino acid residues were highlighted in dark blue.

The results demonstrate the presence of the regions where the amino acid residues are highly conserved across the CYP sequences of the six-polyene BGCs. Based on the alignment, a phylogenetic tree was constructed to determine the similarity of the sequences among the CYPs from different BGCs (Figure 6).

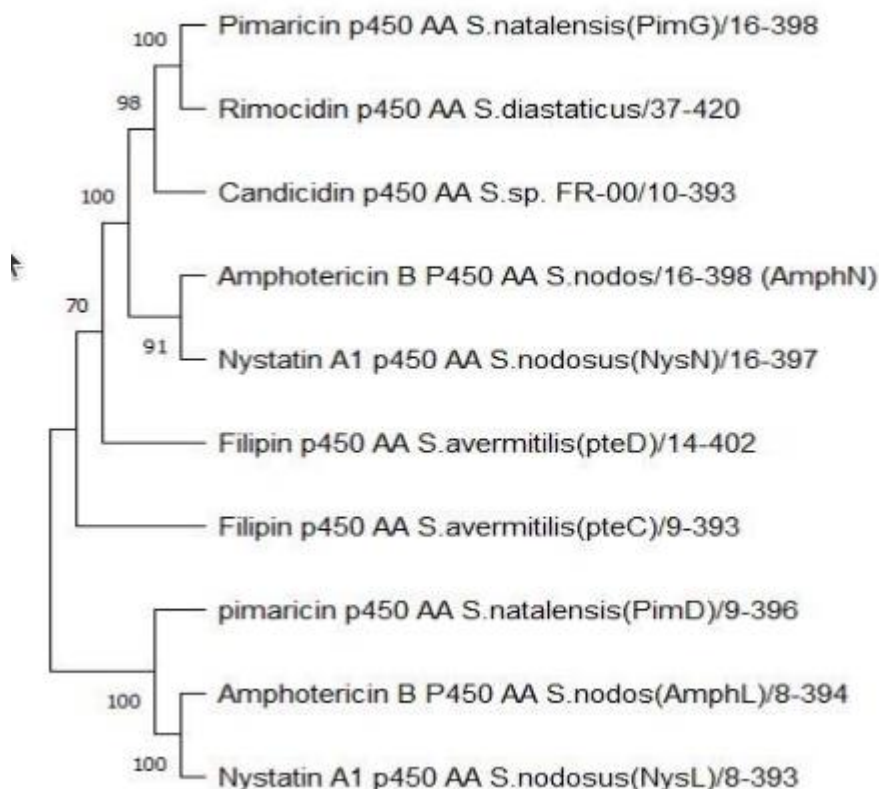


Figure 6: Phylogenetic tree of CYP450 gene extracted from BGCs coding for the 6 different Polyenes.

This analysis showed that cytochrome P450 genes encoded by the polyene clusters can be classified into two distinct similarity groups within the constructed tree. As shown in Figure 6, the first group consists of PimG, AmphN, NysN, PteD, RimG and fscP whereas the second group includes PimD, AmphL, PteC and NysL.

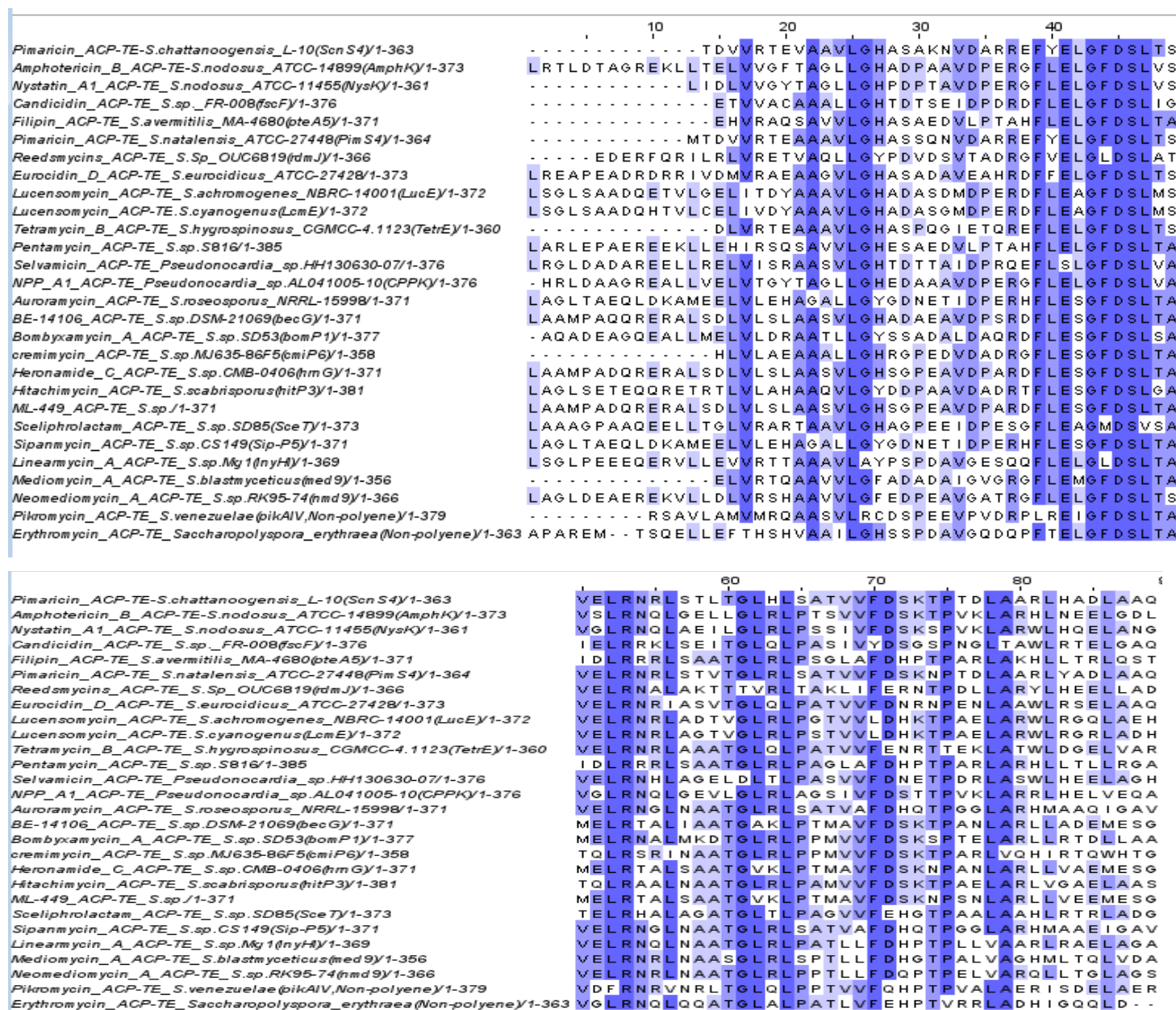
Upon searching for more polyene biosynthetic gene clusters (BGCs) for further analysis of cytochrome P450 (CYP) sequences, it has been discovered that certain polyenes, such as Reedsmycin (a cyclic polyene), and linear polyenes like Linearmycin-A, Mediomycin-A, and Neomediomycin-B, do not possess the genes encoding for cytochrome P450 as well as ferredoxin.

Therefore, the alternative hypothesis was formulated, and the dataset was expanded by including an additional 20 polyene and 2 non-polyene biosynthetic gene clusters (BGCs).

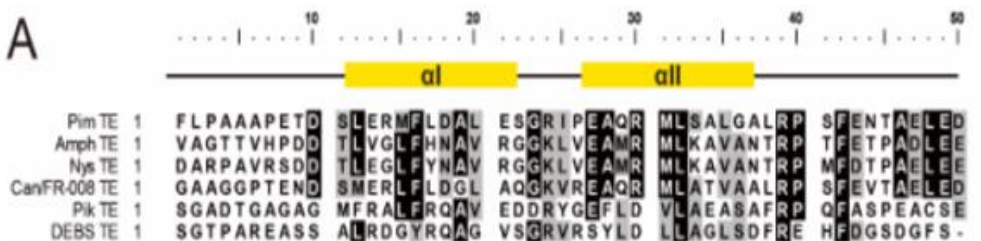
3.3 Analysis of Thioesterase type I

Since thioesterase type I is associated with the acyl carrier protein (ACP) of the last module of type I polyketide synthase, the nucleotide sequence of the whole module was downloaded. The analysis of the nucleotide sequences of the last module of type I polyketide synthase that associated with thioesterase type I revealed that there is only one starting codon and one stop codon, which means that the genes within the module cannot be synthesized individually but

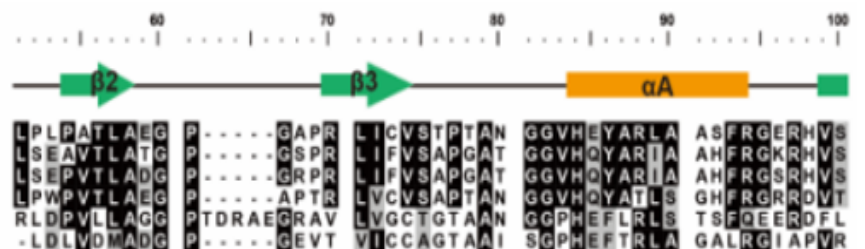
Then the hypothesis was tested by analyzing the sequence similarity of the gene responsible for encoding TE I within the dataset. The results of the multiple sequence alignment of thioesterase gene were compared to the results provided by Zhou, Yucong, et al (Figure 7).



A

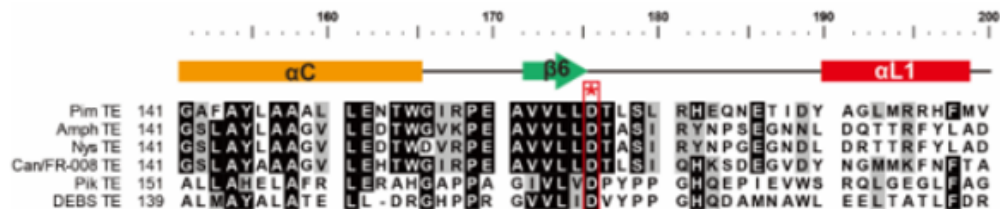


Pimaricin_ACP-TE-S.chattanoogaensis_L-10(Scn S4Y1-363
Amphotericin_B_ACP-TE-S.nodosus_ATCC-14899(AmphKY1-373
Nystatin_A1_ACP-TE-S.nodosus_ATCC-11455(NysKY1-361
Candididin_ACP-TE_S.sp_FR-008(fscFY1-376
Filipin_ACP-TE_S.avermitilis_MA-4680(pteA5Y1-371
Pimaricin_ACP-TE_S.natalensis_ATCC-27448(Pim S4Y1-364
Reedomyacin_ACP-TE_S.Sp_OUC6819(dmJY1-366
Eurocidin_D_ACP-TE_S.eurocidicus_ATCC-27428/1-373
Lucensomycin_ACP-TE_S.achromogenes_NBRC-14001(LucEY1-372
Lucensomycin_ACP-TE_S.cyanogenus(LcmEY1-372
Tetramycin_B_ACP-TE_S.hygrosiposus_CGMCC-4.1123(TetrEY1-360
Pentamycin_ACP-TE_S.sp.SB16/1-385
Selvamicin_ACP-TE_Pseudonocardia.sp.HH130630-07/1-376
NPP_A1_ACP-TE_Pseudonocardia.sp.AL041005-10(CPPKY1-376
Auroramyacin_ACP-TE_S.roseosporus_NRR-15998/1-371
BE-14106_ACP-TE_S.sp.DSM-21069(tecGY1-371
Bombyxamycin_A_ACP-TE_S.sp.SD53(bomP1Y1-377
cremimycin_ACP-TE_S.sp.MJ635-86F5(emiP6Y1-358
Heronamide_C_ACP-TE_S.sp.CMB-0406(hmGY1-371
Hitachimycin_ACP-TE_S.scabrisporus(hitP3Y1-381
ML-449_ACP-TE_S.sp./1-371
Sceliphrolactam_ACP-TE_S.sp.SD85(SeeTY1-373
Sipanmycin_ACP-TE_S.sp.CS149(Sip-P5Y1-371
Lineamycin_A_ACP-TE_S.sp.Myl1(nyHY1-369
Mediomycin_A_ACP-TE_S.sp.blastmyceticus(med9Y1-356
Neomediomycin_A_ACP-TE_S.sp.RK95-74(hmd9Y1-366
Pikromycin_ACP-TE_S.venezuelae(pikAIV,Non-polyeneY1-379
Erythromycin_ACP-TE_Saccharopolyspora_erythraea(Non-polyeneY1-363

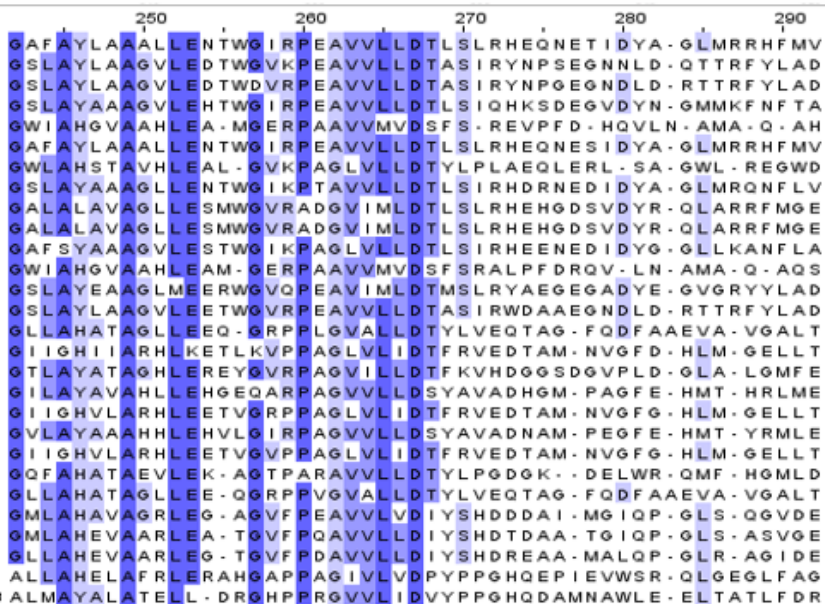


Pimaricin_ACP-TE-S.chattanoogaensis_L-10(Scn S4Y1-363
Amphotericin_B_ACP-TE-S.nodosus_ATCC-14899(AmphKY1-373
Nystatin_A1_ACP-TE_S.nodosus_ATCC-11455(NysKY1-361
Candididin_ACP-TE_S.sp_FR-008(fscFY1-376
Filipin_ACP-TE_S.avermitilis_MA-4680(pteA5Y1-371
Pimaricin_ACP-TE_S.natalensis_ATCC-27448(Pim S4Y1-364
Reedomyacin_ACP-TE_S.Sp_OUC6819(dmJY1-366
Eurocidin_D_ACP-TE_S.eurocidicus_ATCC-27428/1-373
Lucensomycin_ACP-TE_S.achromogenes_NBRC-14001(LucEY1-372
Lucensomycin_ACP-TE_S.cyanogenus(LcmEY1-372
Tetramycin_B_ACP-TE_S.hygrosiposus_CGMCC-4.1123(TetrEY1-360
Pentamycin_ACP-TE_S.sp.SB16/1-385
Selvamicin_ACP-TE_Pseudonocardia.sp.HH130630-07/1-376
NPP_A1_ACP-TE_Pseudonocardia.sp.AL041005-10(CPPKY1-376
Auroramyacin_ACP-TE_S.roseosporus_NRR-15998/1-371
BE-14106_ACP-TE_S.sp.DSM-21069(tecGY1-371
Bombyxamycin_A_ACP-TE_S.sp.SD53(bomP1Y1-377
cremimycin_ACP-TE_S.sp.MJ635-86F5(emiP6Y1-358
Heronamide_C_ACP-TE_S.sp.CMB-0406(hmGY1-371
Hitachimycin_ACP-TE_S.scabrisporus(hitP3Y1-381
ML-449_ACP-TE_S.sp./1-371
Sceliphrolactam_ACP-TE_S.sp.SD85(SeeTY1-373
Sipanmycin_ACP-TE_S.sp.CS149(Sip-P5Y1-371
Lineamycin_A_ACP-TE_S.sp.Myl1(nyHY1-369
Mediomycin_A_ACP-TE_S.sp.blastmyceticus(med9Y1-356
Neomediomycin_A_ACP-TE_S.sp.RK95-74(hmd9Y1-366
Pikromycin_ACP-TE_S.venezuelae(pikAIV,Non-polyeneY1-379
Erythromycin_ACP-TE_Saccharopolyspora_erythraea(Non-polyeneY1-363

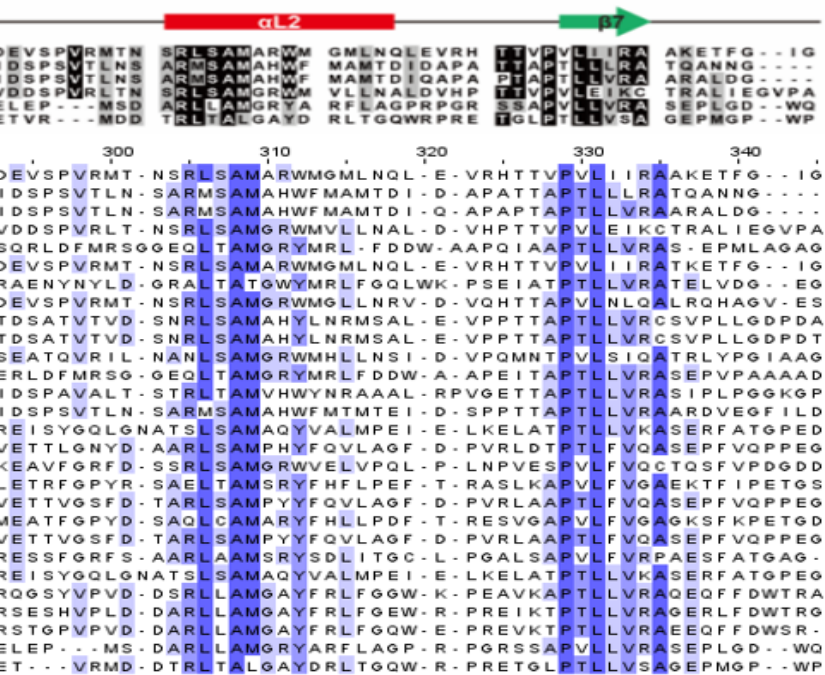
LPLPATLAEG P - - - - GAP R L I C V S T P T A N G G V H E Y A R L A A S F R G E R H V S
 LSEAVTLATG P - - - - GSP R L I F V S A P G A T G G V H Q Y A R I A A H F R G K R H V S
 LSEPVTLAGG P - - - - GRP R L I F V S A P G A T G G V H Q Y A R I A A H F R G S R H V S
 LPWPVTLAEG P - - - - APT R L V C V S A P T A N G G V H Q Y A T L S G H F R G R D V T
 RLDPVLLAGG PTDRAEGRAV LVGCTGTAAI SGPHEFTRLA GALRGIAPVR
 -LDLVDMAAG P - - - - GEV T V I C C A G T A A I S G P H E F T R L A G A L R G I A P V R
 LPLPATLAEG P - - - - GAP R L I C V S T P T A N G G V H E Y A R L A A S F R G E R H V S
 LSEAVTLATG P - - - - GSP R L I F V S A P G A T G G V H Q Y A R I A A H F R G K R H V S
 LSEPVTLAGG P - - - - GRP R L I F V S A P G A T G G V H Q Y A R I A A H F R G S R H V S
 LPWPVTLAEG P - - - - APT R L V C V S A P T A N G G V H Q Y A T L S G H F R G R D V T
 RPAPVRLSQD - - - - DHPLALMCFSPYVVPAGAHQYARFAAPFRDRLDVW
 LPLPATLAEG P - - - - GTP R L I C V S T P T A N G G V H E Y A R L A A S F R G E R H V S
 RPDVVPLVKGES - - - - DAA P L V C F A P P N A F A G P H Q Y Q F A R A F P G H R E V A
 LPLPTTLAEGS - - - - GTP R L I C I S T P T A N G G V H E Y V R F A A H F R G E R H V C
 LPEPVTLAGGPA - - - - G P R L I C I S S P V S V G G A H Q Y A R I A A H F R G D R G V Q
 LQPVTLAGGPA - - - - E P R L I C I S S P V S V G G A H Q Y A R L A A H F R G D R G V Q
 LPLPTTLAEGPA - - - - S P R L I C I S T P T A N G G V H E Y A R A A R F R G R N V S
 LPAPVRLSQDD - - - - L P Y T V M C F S P Y V V P A G A H Q Y A R F A A P F R D R V D M W
 PASPVVLADGPT - - - - T P K L I F V S A P G A T G G V H Q Y A R L A A H F R G R R R V L
 LSEPVTLAEGPA - - - - A P R L I F V S A P G A T G G V H Q Y A R I A A H F R G S R H V S
 PPVSVMADGDE - - - - A P H L V C L C T P A A M G G A Y Q Y A K L V S A F K G A R T V T
 TPKTVRLADGP - - - - GRP R L I C L A T P M A G G V H Q H A R L G S E F R D V R H V S
 LPAEPTTLADGPA - - - - G P H L I C L S T P M A G G V H Q H A R L V S H F G R H K I S
 LPEPTTLSEOPR - - - - G P R L I C V S T P M A G G V H Q H A R L A A H F R G R K S I T
 PPKTVQLADGAR - - - - R P R L I C L S T P M A G G V H Q H A R L G S E F R D I R P V S
 FPAPAGLADGTS - - - - G P R L I C V S T P M A G G V H Q H A R L A A H F R G V R P V S
 PPKTVQLADGAR - - - - R P R L I C L S T P M A G G V H Q H A R L G S E F R D R R P V S
 LEQPVRLADGGE - - - - G P K L I C F S P M A L G G A Q Q Y A R F A A R F R G R E V V
 PPVSVMADGDE - - - - A P H L V C L C T P A A M G G A Y Q Y A K L I S A F K G A R T V T
 APSLVRLSRGT - - - - R P G L V C F S S I L S I S P H Q Y A R F A S A F R G R D V H
 APALVRLSGAG AEPVP P A L V C F S S I L P I S P H Q Y A R F A A G F R G R D V W
 APTLVRLSRGET - - - - G P A L V C F S S I L S I G G P H Q Y A R F A A G F R G R D V W
 RLDPVLLAGGPTDRAEGRA - - - - V L V G C T G T A A N G G P H E F L R L S T S F Q E E R D F L
 -LDLVDMAAG P - - - - GEV T V I C C A G T A A I S G P H E F T R L A G A L R G I A P V R



Pimaricin_ACP-TE-*S.chattanoogaensis*_L-10(ScnS4Y1-363
Amphotericin_B_ACP-TE-*S.nodosus*_ATCC-14899(AmphKY1-373
Nystatin_A1_ACP-TE-*S.nodosus*_ATCC-11455(NysKY1-361
Candididin_ACP-TE-*S.sp.*_FR-008(fscFY1-376
Filipin_ACP-TE-*S.avermitilis*_MA-4680(pteA5Y1-371
Pimaricin_ACP-TE-*S.natalensis*_ATCC-27448(PimS4Y1-364
Reedsmycin_ACP-TE-*S.sp.*_OUC6819(dmJY1-366
Eurocidin_D_ACP-TE-*S.eurocidicus*_ATCC-27428/1-373
Lucensomycin_ACP-TE-*S.achromogenes*_NBRC-14001(LucEY1-372
Lucensomycin_ACP-TE-*S.cyanogenus*(LmEY1-372
Tetramycin_B_ACP-TE-*S.hydrospinosus*_CGMCC-4.1123(TetrEY1-360
Pentamycin_ACP-TE-*S.sp.*_S816/1-385
Selvamycin_ACP-TE-*Pseudonocardia*_sp.HH130630-07/1-376
NPP_A1_ACP-TE-*Pseudonocardia*_sp.AL041005-10(CPPKY1-376
Auroramycin_ACP-TE-*S.roseosporus*_NRRL-15998/1-371
BE-14106_ACP-TE-*S.sp.*_DSM-21069(hecGY1-371
Bombyxamycin_A_ACP-TE-*S.sp.*_SD53(bomP1Y1-377
cremimycin_ACP-TE-*S.sp.*_MJ635-86F5(emiP6Y1-358
Heronamide_C_ACP-TE-*S.sp.*_CMB-0406(hmGY1-371
Hitachimycin_ACP-TE-*S.scabrisporus*(hitP3Y1-381
ML-449_ACP-TE-*S.sp.*_1-371
Sceliphrolactam_ACP-TE-*S.sp.*_SD85(SeeTY1-373
Sipanmycin_ACP-TE-*S.sp.*_CS149(Sip-P5Y1-371
Linearmycin_A_ACP-TE-*S.sp.*_Mg1(nyHY1-369
Mediomycin_A_ACP-TE-*S.blastmyceticus*(med9Y1-356
Neomediomycin_A_ACP-TE-*S.sp.*_RK95-74(hmd9Y1-366
Pikromycin_ACP-TE-*S.venezuelae*(pikAIV,Non-polyeneY1-379
Erythromycin_ACP-TE-*Saccharopolyspora_erythraea*(Non-polyeneY1-363



Pimaricin_ACP-TE-*S.chattanoogaensis*_L-10(ScnS4Y1-363
Amphotericin_B_ACP-TE-*S.nodosus*_ATCC-14899(AmphKY1-373
Nystatin_A1_ACP-TE-*S.nodosus*_ATCC-11455(NysKY1-361
Candididin_ACP-TE-*S.sp.*_FR-008(fscFY1-376
Filipin_ACP-TE-*S.avermitilis*_MA-4680(pteA5Y1-371
Pimaricin_ACP-TE-*S.natalensis*_ATCC-27448(PimS4Y1-364
Reedsmycin_ACP-TE-*S.sp.*_OUC6819(dmJY1-366
Eurocidin_D_ACP-TE-*S.eurocidicus*_ATCC-27428/1-373
Lucensomycin_ACP-TE-*S.achromogenes*_NBRC-14001(LucEY1-372
Lucensomycin_ACP-TE-*S.cyanogenus*(LmEY1-372
Tetramycin_B_ACP-TE-*S.hydrospinosus*_CGMCC-4.1123(TetrEY1-360
Pentamycin_ACP-TE-*S.sp.*_S816/1-385
Selvamycin_ACP-TE-*Pseudonocardia*_sp.HH130630-07/1-376
NPP_A1_ACP-TE-*Pseudonocardia*_sp.AL041005-10(CPPKY1-376
Auroramycin_ACP-TE-*S.roseosporus*_NRRL-15998/1-371
BE-14106_ACP-TE-*S.sp.*_DSM-21069(hecGY1-371
Bombyxamycin_A_ACP-TE-*S.sp.*_SD53(bomP1Y1-377
cremimycin_ACP-TE-*S.sp.*_MJ635-86F5(emiP6Y1-358
Heronamide_C_ACP-TE-*S.sp.*_CMB-0406(hmGY1-371
Hitachimycin_ACP-TE-*S.scabrisporus*(hitP3Y1-381
ML-449_ACP-TE-*S.sp.*_1-371
Sceliphrolactam_ACP-TE-*S.sp.*_SD85(SeeTY1-373
Sipanmycin_ACP-TE-*S.sp.*_CS149(Sip-P5Y1-371
Linearmycin_A_ACP-TE-*S.sp.*_Mg1(nyHY1-369
Mediomycin_A_ACP-TE-*S.blastmyceticus*(med9Y1-356
Neomediomycin_A_ACP-TE-*S.sp.*_RK95-74(hmd9Y1-366
Pikromycin_ACP-TE-*S.venezuelae*(pikAIV,Non-polyeneY1-379
Erythromycin_ACP-TE-*Saccharopolyspora_erythraea*(Non-polyeneY1-363



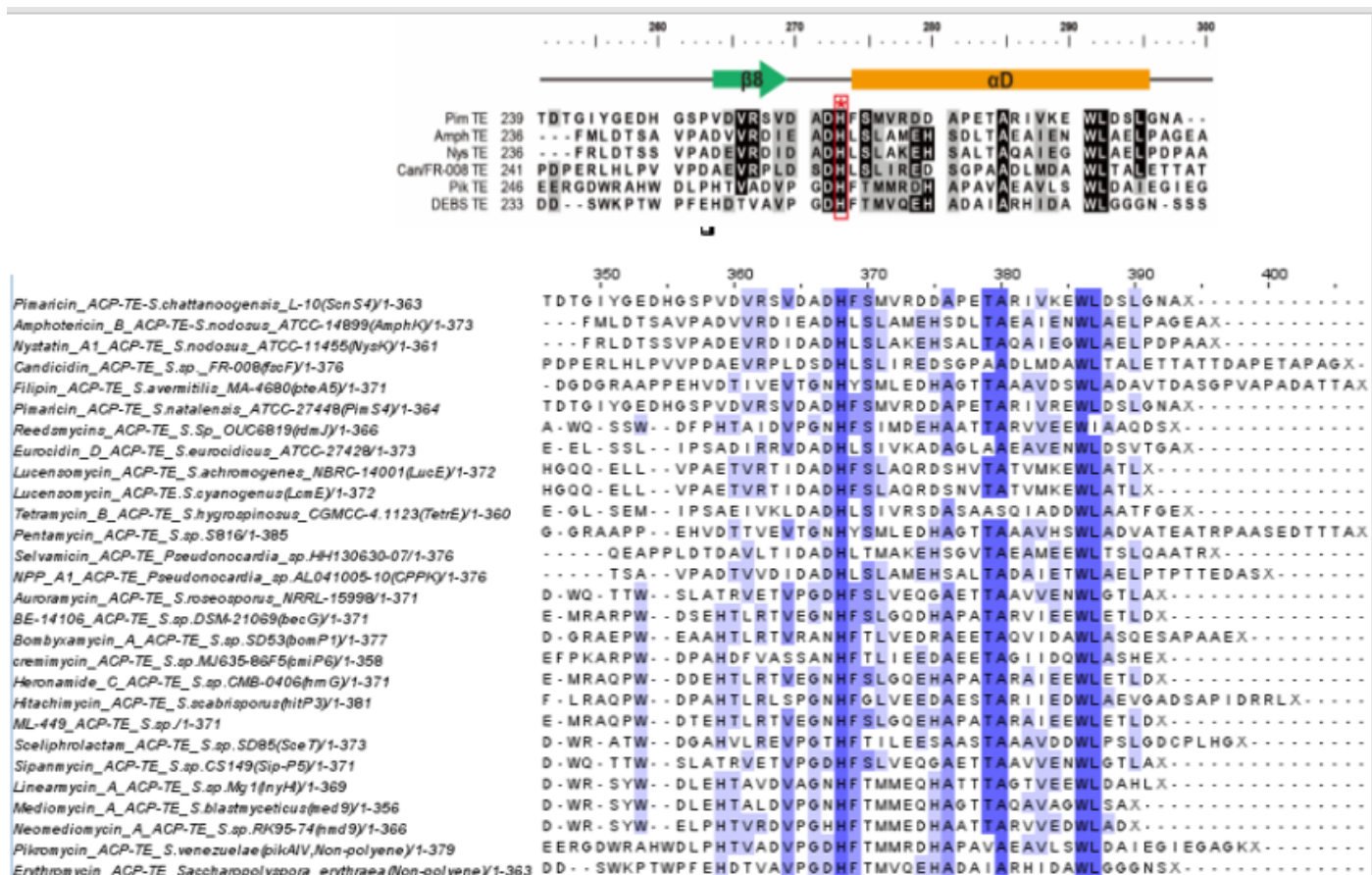


Figure 7: Multiple sequence alignment of amino acid residues of TE type I gene associated with ACP of the last module of polyketide synthase (PKS). The result of the alignment of thioesterase gene was compared to the results provided by Zhou, Yucong, et al.

The results demonstrated that the ACP gene is highly conserved among both polyene and non-polyene polyketide synthases, since it codes for important functional domains that are crucial for its activity. Moreover, the analysis of the multiple sequence alignment of the comprehensive dataset and its comparison to the aligned sequence in previous literature indicated that the cyclic (macrolactone and macrolactam) and linear polyenes share conserved domains with some variability in this region. Although different amino acids may be present within the conserved regions, these variations do not alter the electric charges and the overall geometry of the conserved domain.

Then the highly conserved region of the thioesterase for the polyenes from amino acid 45 to 387 was extracted and translated back into nucleotide sequences for further analysis at the nucleotide level and for the purpose of designing primers.

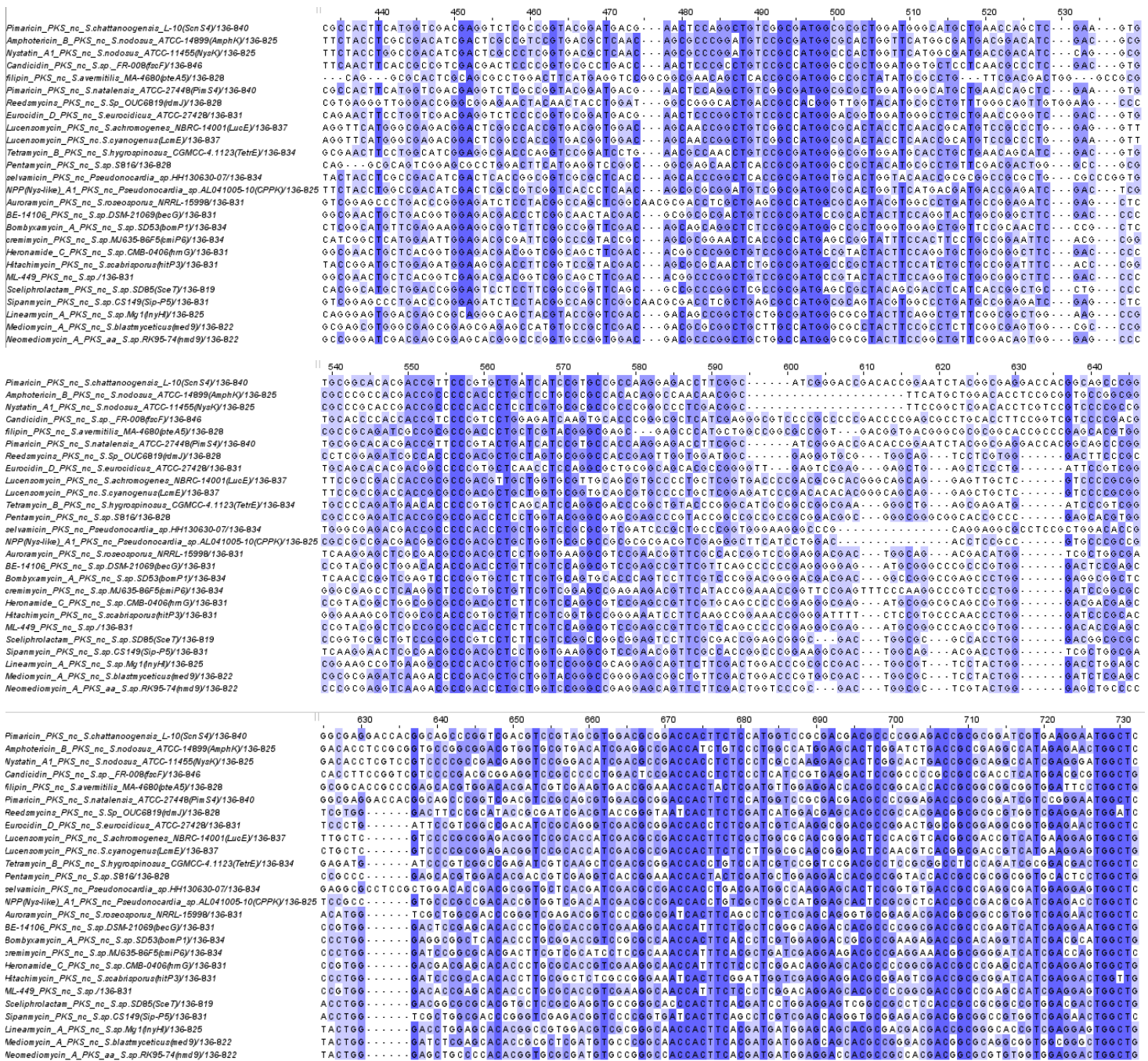


Figure 8: Multiple sequence alignment of the conserved region of TE type I nucleotide sequences. The regions of interest that could be used for designing the primers are indicated with black boxes.

The analysis of the conserved region of the thioesterase gene showed that there is a sequence variability at the level of nucleotides so that the degenerative positions needed to be used.

3.4.4 Designing of primers

Factors such as G+C content, length, melting temperature (T_m) and degeneracy should be considered during designing of primers as they provide optimal specificity and efficiency in PCR amplification.

The results of the designed primers are listed in table 2, analyzed and compared. Then the best option for the primer pairs was selected based on the factors mentioned above.

Table 2: List of the ten primers that were designed by Geneious software version R8.1.9. The best primer pair is highlighted in green.

Primer pairs	Sequences	length	Interval	%GC	Tm	Hairpin Tm	Dimer Tm	Product size
1	TE-Forward GCCGAGCTGGAGGAGC	16	1-16	75.0	59.5	49.6	6.4	501
	TE-Reverse GCCCATCGCGGACAG	15	501-484	72.0	62.9	None	None	
2	TE-Forward GCCGAGCTGGAGGAGC	16	1-16	75.0	59.5	49.6	6.4	499
	TE-Reverse CAGCAGGTGGAAGTAG CGC	19	499-481	63.2	61.1	None	None	
3	TE-Forward GCCGAGCTGGAGGAGC	16	1-16	75.0	59.5	49.6	6.4	499
	TE-Reverse CAGCAGGTGGAAGTAG CG	18	499-482	61.1	57.8	None	None	
4	TE-Forward GGCCCGCGGCTGATC	15	51-65	80.0	60.3	49.6	19.3	500
	TE-Reverse GGCCCGGACGAGCAG	15	550-536	80.0	59.4	None	2.8	
5	TE-Forward GGCCCGCGGCTGATC	15	51-65	80.0	60.3	49.6	19.3	499
	TE-Reverse GCCCCGACGAGCAGC	15	549-535	80.0	60.2	None	None	
6	TE-Forward CCCGCGGCTGATCTGC	16	53-68	75.0	60.7	43.7	22.9	500
	TE-Reverse GTGGCCCGGACGAGC	15	552-538	80.0	59.7	None	None	
7	TE-Forward TGATCTGCGTCAGCACC C	18	61-78	61.1	59.7	50.4	2.7	502
	TE-Reverse GAACGGCTCGGTGGCC	16	562-547	75.0	60.5	47.9	5.1	
8	TE-Forward TGATCTGCGTCAGCACC C	18	61-78	61.1	59.7	50.4	2.7	501
	TE-Reverse ACGGCTCGGTGGCCC	15	560-546	80.0	61.6	48.8	8.6	
9	TE-Forward TGATCTGCGTCAGCACC C	18	61-78	61.1	59.7	50.4	2.7	500
	TE-Reverse		561-547	73.3	59.2			

	AACGGCTCGGTGGCC	15				47.9	5.1	
10	TE-Forward TGATCTGCGTCAGCACC C	18	61-78	61.1	59.7	50.4	2.7	501
	TE-Reverse AACGGCTCGGTGGCCC	16	561-546	75.0	62.2	48.8	8.6	

The analysis of the results indicated that primer pair 1 is the optimal choice. This is because not only does it fulfill the required criteria, but upon examining the binding regions for this primer pair, it was observed that these regions were highly conserved. Consequently, designing a less degenerate pair was possible.

Then the sequences of the best primer pair were utilized to design degenerate primers with the objective of targeting a minimum of 50% of the polyenes present in the dataset.

Table 3: Degenerate PCR primer pair properties.

Sequence	Length	Tm	%GC	Hairpin Tm	Self-Dimer Tm	Degeneracy
TE-Forward 5`GCCGARCTGGAGSAGC3`	16	52.5 - 56.6	66.7- 73.3	45.2	6.4	4
TE-Reverse 5`SCCCATCGCGGACAG3`	15	56.3 - 57.1	73.3	36.6	None	2

3.5. *In silico* PCR

The specificity of the primer pair against the *Streptomyces* genomes was assessed by conducting *in silico* PCR amplification.

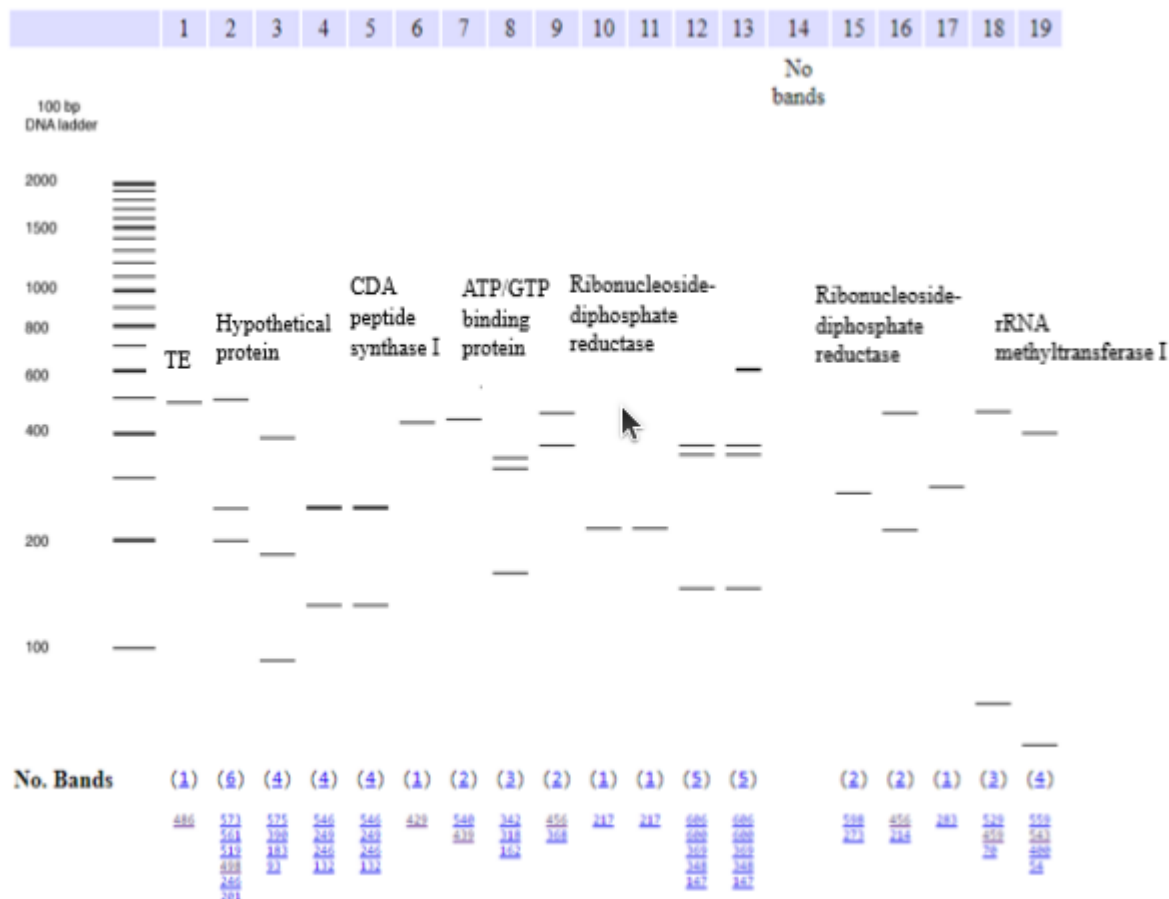


Figure 9: *In-silico* PCR gel illustrating the amplified genes.

The results of the analysis showed that the primer pair, which was designed to target the thioesterase gene, exhibited some degree of non-specific binding. This non-specific binding resulted in the amplification of the intended target as well as other genes with the same band size (around 500 bp).

5. Discussion:

Two hypotheses were formulated to identify a specific gene marker crucial for the biosynthesis of polyenes in *Streptomyces* and evolutionarily related actinomycetes. The first hypothesis was that Cytochrome P450 could serve as a reliable gene marker for polyene production. The multiple sequence alignment results revealed that the CYP sequences of the six-polyene BGCs exhibited shared regions where the amino acid residues are highly conserved. This suggests that these conserved regions may correspond to important catalytic sites or binding domains involved in the biosynthesis of polyene antifungal compounds. Furthermore, apart from the Candicidin and Rimocidin BGCs, the other 4 polyene BGCs possess two cytochrome P450 genes. The phylogenetic tree analysis revealed the presence of two distinct groups within the constructed tree. The possible explanation is that each group is involved in regiospecific oxidation. The first group in figure 6, including PimG, AmphN, NysN, PteD, RimG and

fscP are catalyzing the oxidation of the exocyclic methyl branch to the carboxyl group, whereas the second group including PimD, AmphL, PteC and NysL, are involved in oxidative modifications of the polyol segment. Example for the regiospecific oxidation by PimG and PimD is illustrated in Figure 10.

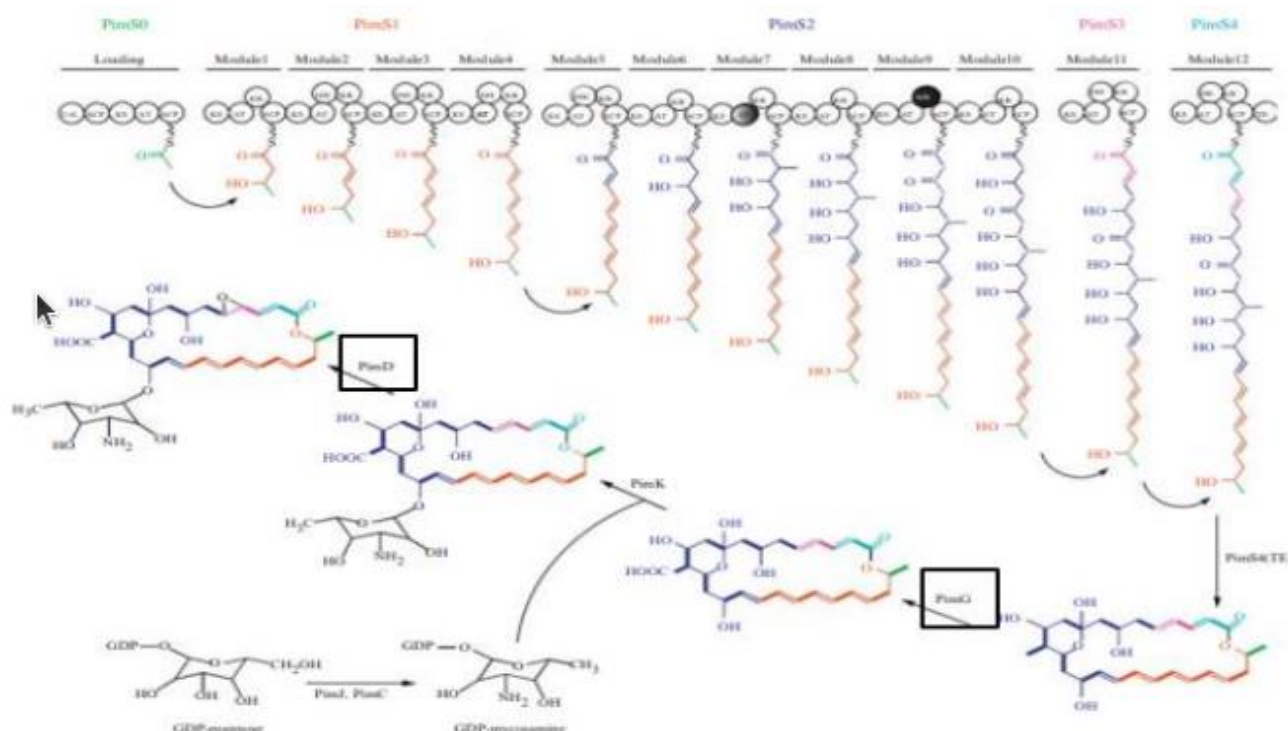


Figure 10: The biosynthesis of Pimaricin. The genes involved in regiospecific oxidation are marked with black boxes.

However, the results of searching for more polyene BGCs indicated that certain polyenes such as Reedsmycin, Linearmycin-A, Mediomycin-A, and Neomediomycin-B are devoid of cytochrome P450 genes. The biosynthetic gene clusters of those polyenes are listed in the appendix. This means that cytochrome P450 cannot serve as a reliable gene marker for polyene production as it is not a gene marker for all polyene biosynthetic gene clusters. Therefore, the hypothesis was rejected, and we have focused on alternative gene. The second hypothesis was that thioesterase type I can be used as an effective gene marker for identifying polyene-producing actinomycetes, particularly streptomycetes. Then the BGCs of 20 polyenes and 2 non-polyenes were included to the dataset.

The results of multiple sequence alignment of thioesterase type I within the biosynthetic gene clusters (BGCs) of 26 polyenes and 2 non-polyenes were compared to the aligned sequences of TE type I of the cyclic (macrolactone) polyenes in previous literature. The analysis revealed that the results in our study are similar to the results in the literature. Despite including additional

sequences of cyclic (macrolactone and macrolactam) and linear polyenes, they still exhibited conserved domains with some degree of variability in these regions. However, these variations in the amino acids that present within the conserved regions do not significantly impact the electric charges or overall geometry of the conserved domain for example, at position 47 (Figure 7), both aspartate (D) and glutamate (E) amino acids have negative charge. These domains share low sequence similarity with non-polyenes at the amino acid level, suggesting that these domains are specific for cyclization and offloading of the polyhydroxylated intermediates with continuous conjugated double bonds. As a result, the amino acid sequences of these domains are translated back to nucleotide sequences and used for designing the primers. The analysis of the domains at the nucleotide level indicated that those different amino acids, that have the same properties, are represented by different codons, and even identical amino acids have multiple possible codons. This poses a challenge when designing primers for thioesterase type I.

Based on the circumstances, 10 primer pairs could be designed. The results showed that primer pair 1 is the best option as both forward and reverse primers have the optimum GC %. This indicates the stability and specificity of primer binding to the target DNA sequence. Furthermore, the melting temperatures (T_m) of the primers are close to each other. Having closely matched T_m values for the forward and reverse primers is important so that they are likely to anneal to the target DNA sequence at a similar temperature during PCR amplification.

Moreover, formation of dimers and hairpins are possible only for the forward primer, unlike primer pairs 7-10. This difference in behaviour lowers the probability of the primers annealing effectively. The most important advantage of this primer pair over the others is that the analysis of the binding regions for this primer pair revealed a high degree of conservation for polyenes thus, designing less degenerate primer pair.

Then the sequences of the best primer pair were used to design the degenerate primers. Degenerate primers provide flexibility in primer design, enabling them to bind to target regions with slight variations. However, degeneracy of the primers was kept low to ensure specific binding. *In silico* PCR was used to evaluate the specificity of the degenerate primer pairs by testing these primers against 19 *Streptomyces* strains. The results show that not only the thioesterase type I gene was amplified, but also other genes. The lack of specificity observed with this primer pair suggests that the sequences of the primer pair may have similarity with other genes within the genome. Moreover, the relatively short length of the primer pair (15 bp) may increase the probability of non-specific binding to other regions of the genome as shorter primers have a higher likelihood of encountering sequences with partial complementarity elsewhere in the genome, leading to unintended amplification during PCR. Therefore, the hypothesis was rejected. To enhance primer specificity and minimize non-specific amplification, it is recommended to design longer primers (18-30 bases) that are more unique to the target gene sequence or to use nested PCR to ensure specific binding. Due to time constraints, the primer pair was only tested using *in silico* PCR. However, it is preferable to conduct experimental testing in order to facilitate further analysis, such as sequencing the amplicons.

6. Conclusion:

The analysis of cytochrome P450 sequences in the polyene biosynthetic gene clusters revealed two different groups. Each group leads to regiospecific oxidation. However, some polyene biosynthetic gene clusters did not have cytochrome P450 enzymes, making it an unreliable gene marker. Furthermore, the thioesterase type I gene was found in all polyene BGCs, but was not effective as a gene marker due to the lack of specificity of the designed primer pair. Therefore, testing the primer pair in the lab and sequencing the products can provide additional information about the amplified regions.