

Chapter 5
Results:
**In silico analyses of secondary metabolite
biosynthetic gene clusters in different host-specific
lineages of the blast fungus**

5.1 Identification of biosynthetic gene clusters (BGCs) in host-specific strains of *M. oryzae*

In this investigation, we employed a dataset comprising of 68 *M. oryzae* genomes, which were sourced from six distinct host plants, namely rice (*Oryza sativa*), finger millet (*Eleusine coracana*), foxtail millet (*Setaria sp.*), wheat (*Triticum aestivum*), perennial ryegrass (*Lolium sp.*) and weeping lovegrass (*Eragrostis curvula*) (**Fig. 5.1A, Table 5.1**). This dataset encompasses an earlier sequenced 15 field isolates, collected from different parts of India (**Fig. 5.1A, Table 5.1**). The remaining 53 genomes were obtained from publicly available genome sequences. Furthermore, we incorporated sequence data from three *M. grisea* strains, previously isolated from crabgrass (*Digitaria sp.*), to be used as control (**Table 5.1**). Sixteen out of the publicly available assemblies were derived from long-read sequencing technology, resulting in highly contiguous datasets. Subsequently, gene predictions were performed on all 71 assemblies and analysis of BUSCO genes indicated robust gene prediction quality and assembly completeness, with a BUSCO score exceeding 90% (**Table 5.1**).

Utilizing a set of 2655 BUSCO genes found within the 68 *M. oryzae* genomes, we constructed a phylogenomic tree. The resulting topology of the species tree confirms the existence of multiple genetic lineages within *M. oryzae*, each specialized on different host plants, including *Oryza*, *Setaria*, *Eleusine*, *Eragrostis*, *Triticum* and *Lolium* (**Fig. 5.1B**). This phylogenetic arrangement aligns with the previously reported population structure (Gladieux, Condon, et al., 2018). A similar phylogenomic assessment was conducted using an additional three genomes from *M. grisea* (*Digitaria* isolates), reaffirming the divergence of *M. grisea* as a distinct species from *M. oryzae*, as in accordance with prior findings (**Fig. 5.2**; Gladieux et al., 2018).

Our methodology involved the utilization of a comprehensive pipeline to discern the biosynthetic diversity within various lineages of *M. oryzae* (**Fig. 5.3**). This approach encompassed the initial identification of biosynthetic gene clusters (BGCs) responsible for secondary metabolite (SM) production within genomic regions, followed by a subsequent analysis employing similarity networks to explore the genetic variations among the predicted BGCs.

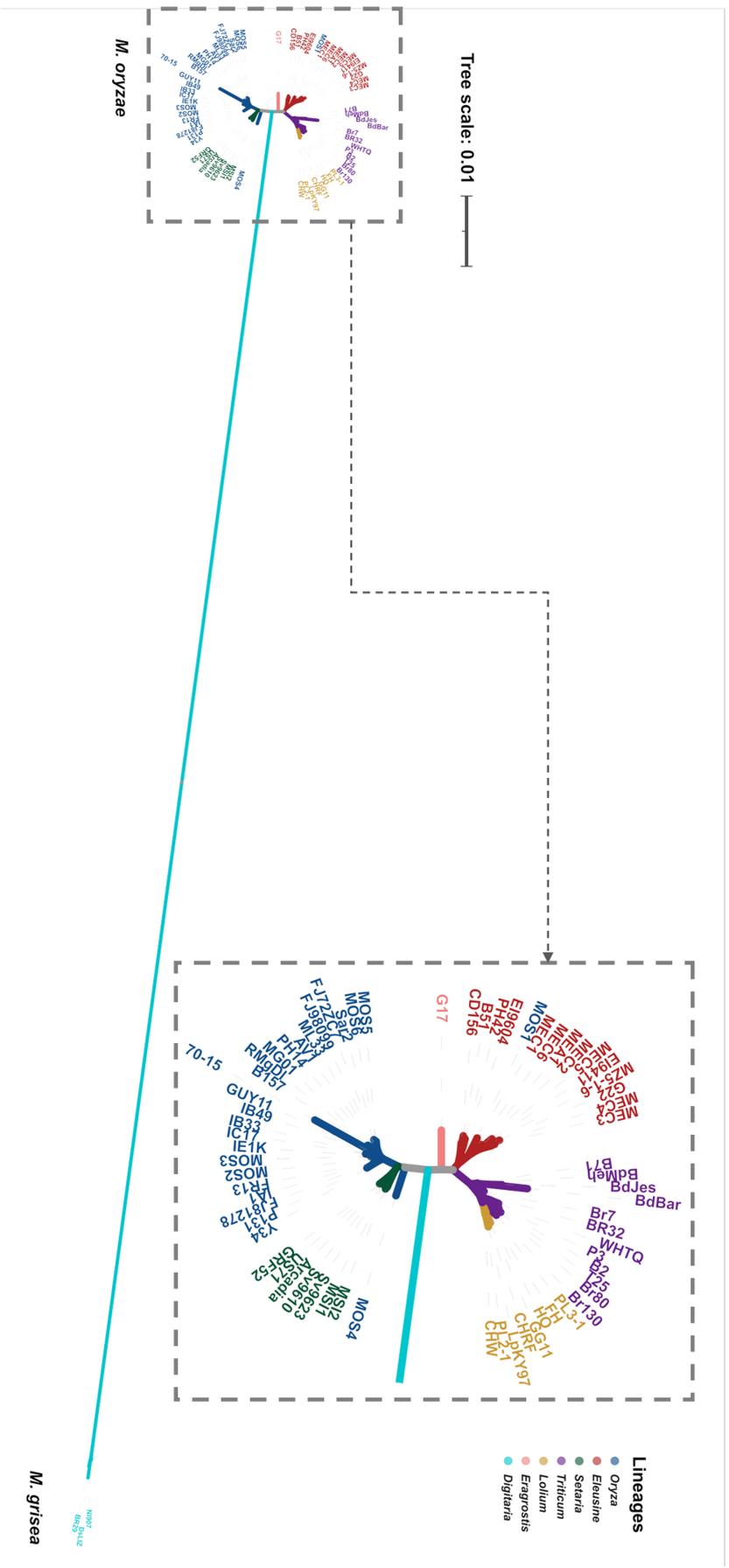


Figure 5.2: *M. oryzae* and *M. grisea* are evolutionarily distinct species adapted to different hosts. Maximum likelihood tree constructed based on concatenation of a total 2557 BUSCOs present in all 71 genomes of *M. oryzae* and *M. grisea* strains used in the study. Colored branches depict different host-specific genetic lineages of *M. oryzae*.

Table 5.1: *M. oryzae* and *M. grisea* genomes used in this study.

Strains	Synonyms	Host of Isolation	Country/Region	Year of Collection	Phylogenetic lineage	Assembly Accession IDs	References	BUSCO (%)
70-15		<i>Oryza sativa</i>	n/a	-	<i>Oryza</i>	GCF_000002495.2	Dean et al. 2005	98.6
Arcadia		<i>Setaria viridis</i>	Lexington, Kentucky, USA	1998	<i>Setaria</i>	GCA_002925445.1	Rahnama et al. 2021	95.7
AV1-1-1		<i>Oryza sativa</i>	Ghana: Aveyime	2015	<i>Oryza</i>	GCA_011799965.1 *	Zhong et al. 2020	97.1
B157		<i>Oryza sativa</i>	India: Marutera, Andhra Pradesh	1989	<i>Oryza</i>	GCA_000832285.1	Gowda et al. 2015	95.8
B2		<i>Triticum aestivum</i>	Bolivia: Okinawa Uno	2011	<i>Triticum</i>	GCA_002218465.1	Rahnama et al. 2021	96.5
B51		<i>Eleusine indica</i>	Bolivia: Quirusillas	2012	<i>Eleusine</i>	GCA_002925415.1	Farman et al. 2017	96.7
B71		<i>Triticum aestivum</i>	Bolivia: Okinawa Uno	2012	<i>Triticum</i>	GCA_004785725.2 *	Peng et al. 2019	97.1
BdBar	BdBar16-1	<i>Triticum aestivum</i>	Bangladesh: Barisal	2016	<i>Triticum</i>	GCA_001675615.1	Rahnama et al. 2021	91.5
BdJes	BdJes16-1	<i>Triticum aestivum</i>	Bangladesh: Jessore district	2016	<i>Triticum</i>	GCA_001675595.1	Rahnama et al. 2021	93.7
BdMeh	BdMeh16-1	<i>Triticum aestivum</i>	Bangladesh: Mehepur district	2016	<i>Triticum</i>	GCA_001675605.1	Rahnama et al. 2021	96.9
BR0032	BR32	<i>Triticum aestivum</i>	Brazil	1991	<i>Triticum</i>	GCA_900474545.3 *	Langner et al. 2021	97
Br130		<i>Triticum aestivum</i>	Brazil: Mato Grosso do Sul	1990	<i>Triticum</i>	GCA_002925325.1	Rahnama et al. 2021	94.6
BR29		<i>Digitaria</i>	Brazil	-	<i>Digitaria</i>	BR29	Chiapello et al. 2015	97.2
Br7		<i>Triticum aestivum</i>	Brazil: Parana	1990	<i>Triticum</i>	GCA_002925335.1	Rahnama et al. 2021	97
Br80		<i>Triticum aestivum</i>	Brazil	1991	<i>Triticum</i>	GCA_002925345.1	Rahnama et al. 2021	96.9
CD156	CD0156	<i>Eleusine indica</i>	Ivory Coast, Ferkessedougou	1989	<i>Eleusine</i>	GCA_900474475.3 *	Langner et al. 2021	97.2
CHRF		<i>Lolium perenne</i>	Siler Springs, MD, USA	1996	<i>Lolium</i>	GCA_002925295.1	Rahnama et al. 2021	97.1

CHW	<i>Lotium perenne</i>	USA: Annapolis, MD	1996	<i>Lotium</i>	GCA_002925285.1	Rahmana et al. 2021	96.8
DSLIZ	<i>Digitaria</i>	Lexington, Kentucky, USA	2000	<i>Digitaria</i>	GCA_002925245.1 *	Rahmana et al. 2021	96
EI9411	<i>Eleusine indica</i>	China, Fujian	1994	<i>Eleusine</i>	GCA_001548775.1	Zhong et al. 2016	96.8
EI9604	<i>Eleusine indica</i>	China, Zhejiang	1996	<i>Eleusine</i>	GCA_001548785.1	Zhong et al. 2016	96.9
FH	<i>Lotium perenne</i>	USA: Hagerstown, MD	1997	<i>Lotium</i>	GCA_002925225.1	Rahmana et al. 2021	95.7
FJ72ZC7-77	<i>Oryza sativa</i>	China: Fujian	1992	<i>Oryza</i>	GCA_011799905.1 *	Zhong et al. 2020	96.9
FJ81278	<i>Oryza sativa</i>	China: Fujian	1981	<i>Oryza</i>	GCA_002368475.1 *	Bao et al. 2017	97
FJ98099	<i>Oryza sativa</i>	China: Fujian	1998	<i>Oryza</i>	GCA_011799925.1 *	Zhong et al. 2020	97.1
FR13	<i>Oryza sativa</i>	France	1990	<i>Oryza</i>	GCA_900474655.3 *	Langner et al. 2021	97.1
G17	<i>Eragrostis curvula</i>	Japan	1976	<i>Eragrostis</i>	GCA_002925205.1	Rahmana et al. 2021	96.9
G22	<i>Eleusine coracana</i>	Japan	1976	<i>Eleusine</i>	GCA_002925165.1	Glaideux et al. 2018	96.9
GG11	<i>Lotium perenne</i>	USA: Lexington, KY	1997	<i>Lotium</i>	GCA_002925155.1	Rahmana et al. 2021	96.9
GRF52	<i>Setaria viridis</i>	Lexington, Kentucky, USA	2001	<i>Setaria</i>	GCA_002925145.1	Rahmana et al. 2021	96.9
GY11	<i>Oryza sativa</i>	French - Guyana	1988	<i>Oryza</i>	GCA_002368485.1 *	Bao et al. 2017	97
HO	<i>Lotium perenne</i>	USA: Richmond, PA	1996	<i>Lotium</i>	GCA_002925105.1	Rahmana et al. 2021	97
IA1	<i>Oryza sativa</i>	Arkansas, USA	2009	<i>Oryza</i>	GCA_002925085.1	Rahmana et al. 2021	97
IB33	<i>Oryza sativa</i>	Texas, USA	-	<i>Oryza</i>	GCA_002925065.1	Rahmana et al. 2021	97.1
IB49	<i>Oryza sativa</i>	AR, USA	1992	<i>Oryza</i>	GCA_002925045.1	Rahmana et al. 2021	97.1
IC17	<i>Oryza sativa</i>	AR, USA	1992	<i>Oryza</i>	GCA_002925025.1	Rahmana et al. 2021	97.1
IEIK	<i>Oryza sativa</i>	AR, USA	2003	<i>Oryza</i>	GCA_002924985.1	Rahmana et al. 2021	97.1
LpKY97	<i>Lotium perenne</i>	USA	1997	<i>Lotium</i>	GCA_012272995.1 *	Rahmana et al. 2021	96.6
MEA1	<i>Eleusine africana</i>	GKVK, Bengaluru, KA, India	2015	<i>Eleusine</i>		This study	96.9
MEC1	<i>Eleusine coracana</i>	Waghai, Dangs, GJ, India	2015	<i>Eleusine</i>		This study	96.9

MEC2	<i>Eleusine coracana</i>	Waghai, Dangs, GJ, India	2015	<i>Eleusine</i>	This study	96.9
MEC3	<i>Eleusine coracana</i>	GKVK, Bengaluru, KA, India	2015	<i>Eleusine</i>	This study	96.8
MEC4	<i>Eleusine coracana</i>	GKVK, Bengaluru, KA, India	2015	<i>Eleusine</i>	This study	96.8
MEC5	<i>Eleusine coracana</i>	Dentam, west sikkim, SK, India	2015	<i>Eleusine</i>	This study	96.4
MEC6	<i>Eleusine coracana</i>	South Sikkim, SK, India	2015	<i>Eleusine</i>	This study	96.8
MG01	<i>Oryza sativa</i>	India: Mandya	2011	<i>Oryza</i>	Gowda et al. 2015	95.4
ML33	<i>Oryza sativa</i>	Mali	1995	<i>Oryza</i>	Rahnama et al. 2021	97
MOS1	<i>Oryza sativa</i>	Waghai, Dangs, GJ, India	2013	<i>Eleusine</i>	This study	97
MOS2	<i>Oryza sativa</i>	Waghai, Dangs, GJ, India	2013	<i>Oryza</i>	This study	97
MOS3	<i>Oryza sativa</i>	Hazaribag, JH, India	2010	<i>Oryza</i>	This study	96.9
MOS4	<i>Oryza sativa</i>	Kanyakumari, TN, India	2015	<i>Oryza</i>	This study	97.1
MOSS	<i>Oryza sativa</i>	Upper Lingthem, North Sikkim, SK, India	2015	<i>Oryza</i>	This study	97
MOS6	<i>Oryza sativa</i>	Geysing, West Sikkim, SK, India	2015	<i>Oryza</i>	This study	96.9
MSI1	<i>Setaria Italica</i>	GKVK, Bengaluru, KA, India	2015	<i>Setaria</i>	This study	96.8
MSI2	<i>Setaria Italica</i>	GKVK, Bengaluru, KA, India	2015	<i>Setaria</i>	This study	96.7
MZ5-1-6	<i>Eleusine coracana</i>	Japan, Miyazaki	1976	<i>Eleusine</i>	Luciano et al 2019	97.1
NI907	<i>Digitaria</i>	Japan:Tochigi	1974	<i>Digitaria</i>	Gomez Luciano et al. 2019	97.3
P131	<i>Oryza sativa</i>	Japan	1976	<i>Oryza</i>	Xue et al. 2012	96.7
P3	<i>Triticum durum</i>	Paraguay: Canindeyu	2012	<i>Triticum</i>	Rahnama et al. 2021	97
PH0014-rn	<i>Oryza sativa</i>	Philippines	-	<i>Oryza</i>	Chiapello et al. 2015	96.6

PH42	<i>Eleusine coracana</i>	Philippines	1983	<i>Eleusine</i>	GCA_002924865.1	Pieck et al. 2017	95.5
PL2-1	<i>Lolium multiflorum</i>	Pulaski Co, KY, USA	2002	<i>Lolium</i>	GCA_002924835.1	Rahmana et al. 2021	97
PL3-1	<i>Lolium multiflorum</i>	Pulaski Co, KY, USA	2002	<i>Lolium</i>	GCA_002924825.1	Rahmana et al. 2021	96.7
RMg-DI	<i>Oryza sativa</i>	India: Madhubani, Bihar	2012	<i>Oryza</i>	GCA_001853415.2 *	Reddy et al. 2021	93.7
Sar-2-20-1	<i>Oryza sativa</i>	Suriname: Saramacca	2013	<i>Oryza</i>	GCA_011799915.1 *	Zhong et al. 2020	97
Sv9610	<i>Setaria viridis</i>	China: Zhejiang	1996	<i>Setaria</i>	GCA_001548845.1	Zhong et al. 2016	97
Sv9623	<i>Setaria viridis</i>	China: Zhejiang	1996	<i>Setaria</i>	GCA_001548855.1	Zhong et al. 2016	96.8
T25	<i>Triticum aestivum</i>	Brazil: Parana	1988	<i>Triticum</i>	GCA_002924745.1	Rahmana et al. 2021	96.8
US0071	<i>Setaria spp.</i>	USA	1998	<i>Setaria</i>	GCA_900474175.3 *	Langner et al. 2021	97.1
WHTQ	<i>Triticum aestivum</i>	Brazil	-	<i>Triticum</i>	GCA_002924665.1	Rahmana et al. 2021	93.8
Y34	<i>Oryza sativa</i>	China: Yunnan	1982	<i>Oryza</i>	GCA_000292585.1	Xue et al. 2012	96.7

* Denotes assemblies obtained using long-read next generation sequencing technologies

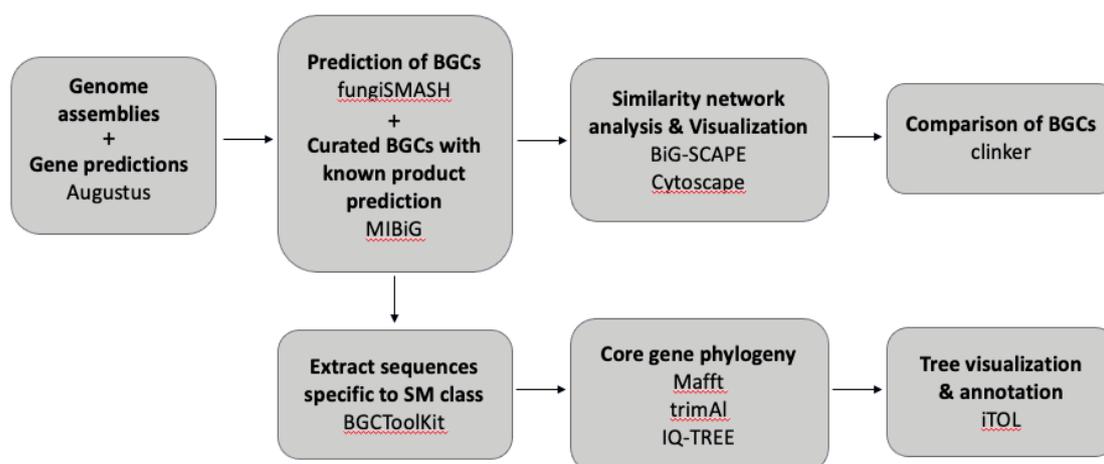


Figure 5.3: Workflow for exploring fungal biosynthetic diversity.

The process of pinpointing genomic regions harboring SM BGCs was executed across all 71 genomes using fungiSMASH (Blin et al., 2019). This investigation yielded a total of 4224 BGCs predictions, with an average of approximately 59 BGCs per strain. These projected BGCs were categorized based on their core biosynthetic genes, such as genes encoding polyketide synthases (PKSs), non-ribosomal peptide synthetases (NRPSs) or terpene cyclases (TCs). It is noteworthy that all lineages specific to particular host plants exhibited a similar number of BGCs across various BGC classes, with type I PKSs being the most prevalent (**Fig. 5.4**). Overall, these analyses underscore the substantial potential of *M. oryzae* to synthesize SMs, some of which might play pivotal roles in virulence and/or host specialization.

5.2 Similarity network analyses to identify biosynthetic diversity in host-specific lineages

To ascertain whether any potential biosynthetic gene cluster (BGC) is linked to the ability to infect specific host plants, we performed a network similarity analysis using BiG-SCAPE (Navarro-Muñoz et al., 2019). This analysis encompassed a dataset of the 4224 predicted

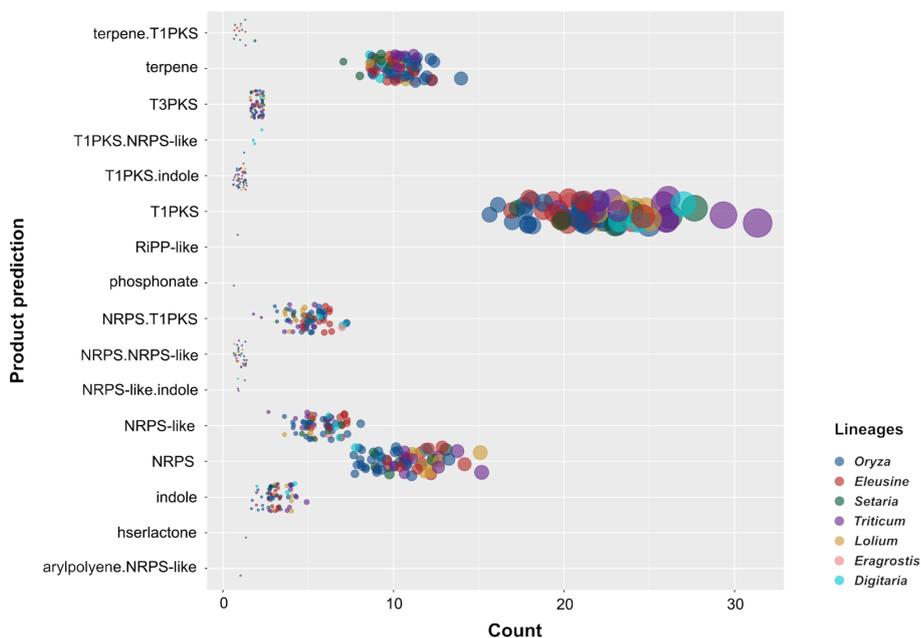
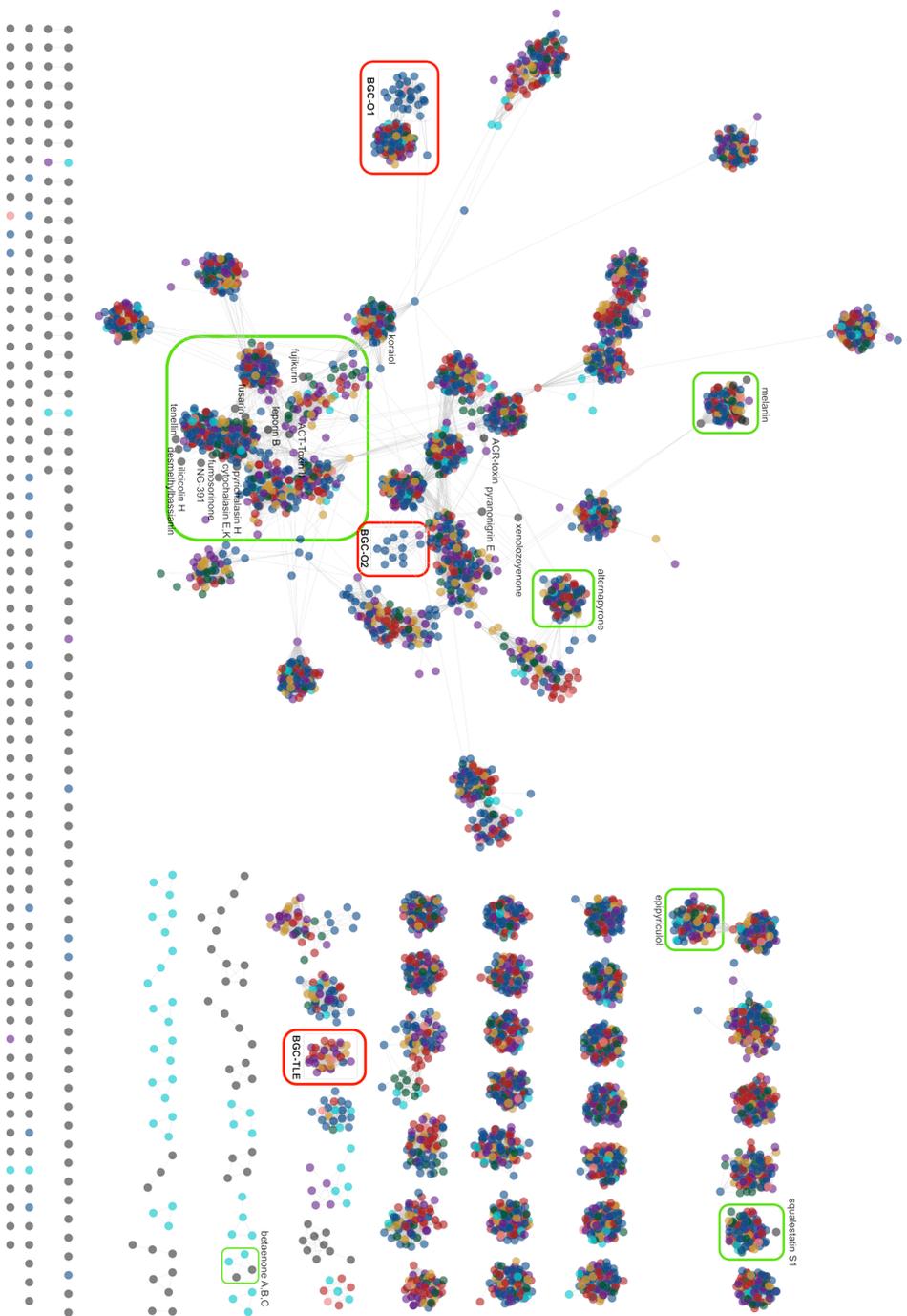


Figure 5.4: Occurrence of BGCs associated with specific classes of SM in individual genomes of *M. oryzae* and *M. grisea*. Area of a given circle is directly proportional to the total number of BGC associated with a class of specific product in a particular genome. Color of a given circle denotes the host-specific lineage it belongs to.

BGCs, complemented by 277 characterized BGCs sourced from the MIBiG database (Kautsar et al., 2019) for reference purposes. This resulted in a total of 4501 BGCs, which were subsequently clustered into 283 gene cluster families (GCFs) or subnetworks, of which 180 represents singletons. Among these, 160 belonged to the reference characterized BGCs (**Fig. 5.5**). Our BiG-SCAPE analysis unveiled that, while the majority of the BGCs, regardless of their SM classification, are distributed across various lineages specific to different host plants (multi-colored closed circles grouped together; **Fig. 5.5**), only a limited number of GCFs or subnetworks exhibited similarities with reference BGCs derived from the MIBiG database. This observation implies that a significant portion of these BGCs remains uncharacterized.

The likely products associated with a specific BGC can be inferred through the examination of homology with reference BGCs known to encode pathways for SMs, particularly those identified in different fungi. As a result of our analysis, we have identified a subset of BGCs



Lineages: ● *Oryza* ● *Elymus* ● *Setaria* ● *Triticum* ● *Lolium* ● *Eragrostis* ● *Digitaria* ● Known fungal BGCs

Figure 5.5: Similarity network analysis of biosynthetic gene clusters (BGCs) from *M. oryzae* and *M. grisea*. A BiG-SCAPE analysis with a cutoff c0.5 depicts similarity of 4224 BGCs from *M. oryzae* or *M. grisea* with 277 reference BGCs from MIBiG database. Each dot represents a BGC and is color-coded according to the lineage. Gene cluster families (GCF; subnetworks) marked with green boxes share significant homology with reference MIBiG BGCs (grey-colored circle). GCFs marked with red boxes are found to be unique to host-specific lineages. The length of the gray lines is proportional to the genetic distance between BGCs. Singletons are shown as individual dots at the bottom.

that are likely associated with the synthesis of SMs such as melanin, cytochalasans, epipyriculol, squalestatin, Fusarin, fujikurin, alternapyrone, cercosporin, ACT-ToxinII, and pyranonigrin, in *M. grisea* and/or different lineages of *M. oryzae* (green boxes; **Fig. 5.5**).

Within the 91 GCFs present in *Magnaporthe* strains, twelve GCFs were found to comprise characterized BGCs, primarily associated with the production of DHN melanin, epipyriculol, alternapyrone, squalestatin and cytochalasans in *M. oryzae*, and betaneone in *M. grisea* (**Fig. 5.5-5.10**). Nonetheless, the majority of the remaining BGCs studied here, did not display any homology with the known/reference BGCs in other fungi. Forty-four GCFs are found in several *M. oryzae* lineages as well as in *M. grisea*, most of them remain uncharacterized. Interestingly, 13 GCFs were specific to *M. grisea*, while 14 GCFs appeared to be exclusive to *M. oryzae*, hinting at their potential roles in pathogenesis on *Digitaria* and other relevant host plants (**Fig. 5.5**). Taken together, our comprehensive analyses strongly indicate that certain SM BGCs exhibit significant diversity among various *M. oryzae* lineages adapted to specific host plants and/or geographical locations.

5.3 Gene Cluster Families unique to host-specific lineages

Within the set of GCFs specific to *M. oryzae*, three GCFs could potentially be associated with host specialization. Specifically, BGC-O1 and BGC-O2 exhibited a predominant presence within the *Oryza* lineage, while BGC-TLE was unique to the *Triticum*, *Lolium*, *Eleusine* and *Eragrostis* lineages (**Fig. 5.5**). These lineages shared a common ancestry and diverged from the *Oryza* and *Setaria* lineages (**Fig. 5.1**). BGC-O1 encompasses genes encoding a Type 1 reducing polyketide synthase (rPKS) and tailoring enzymes (**Fig. 5.14**), whereas BGC-O2 is restricted to a solitary T1 PKS gene (**Fig. 5.11**).

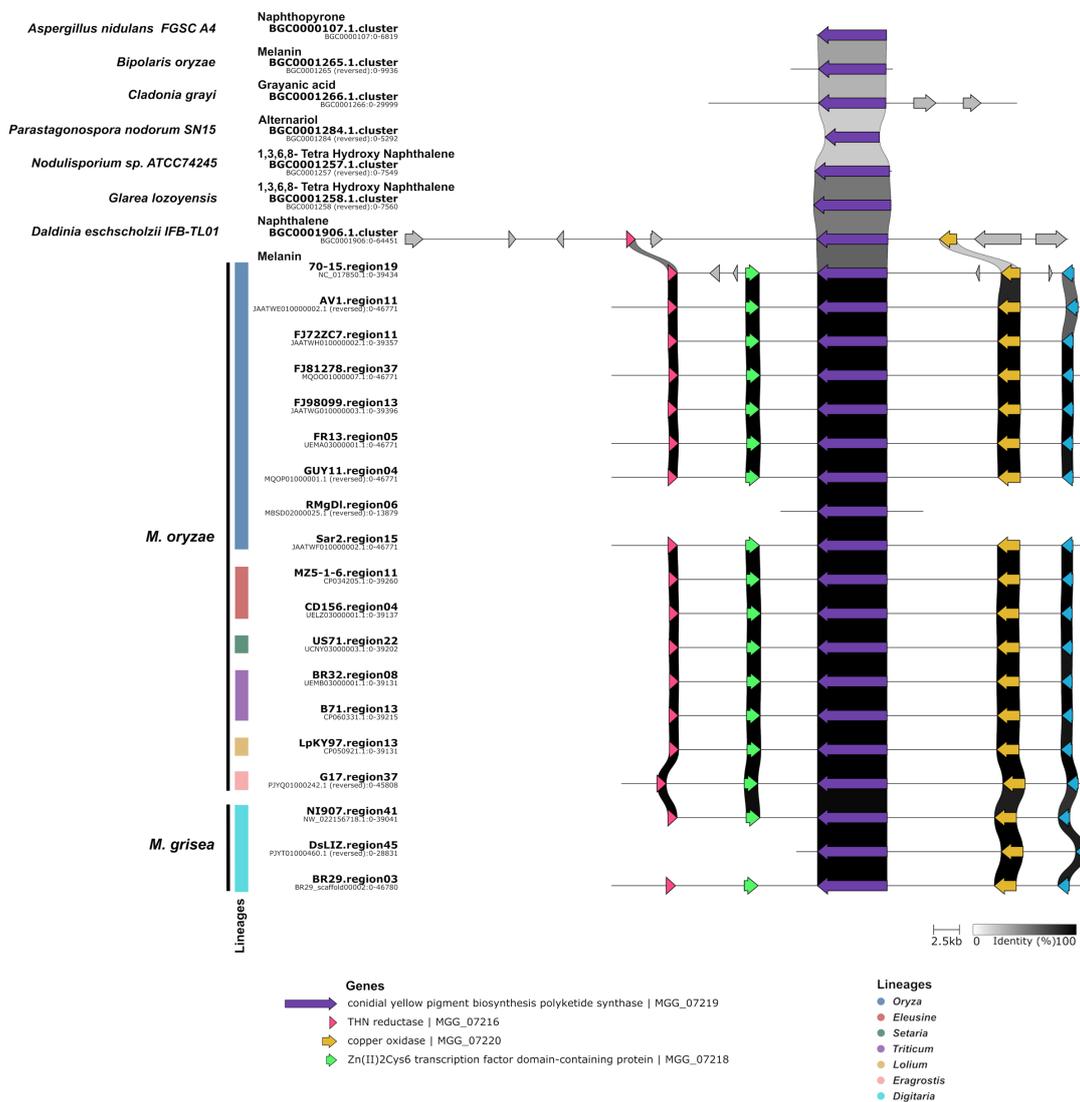


Figure 5.6: Homology of melanin-associated *M. oryzae* and *M. grisea* gene cluster family with reference BGCs in MIBiG database. The map depicts comparison of BGC loci with reference BGCs associated with melanin or melanin-like SM product in MIBiG. The shaded area between any two arrows denotes degree of homology (0 to 100%; white to black, respectively) between the two sequences.

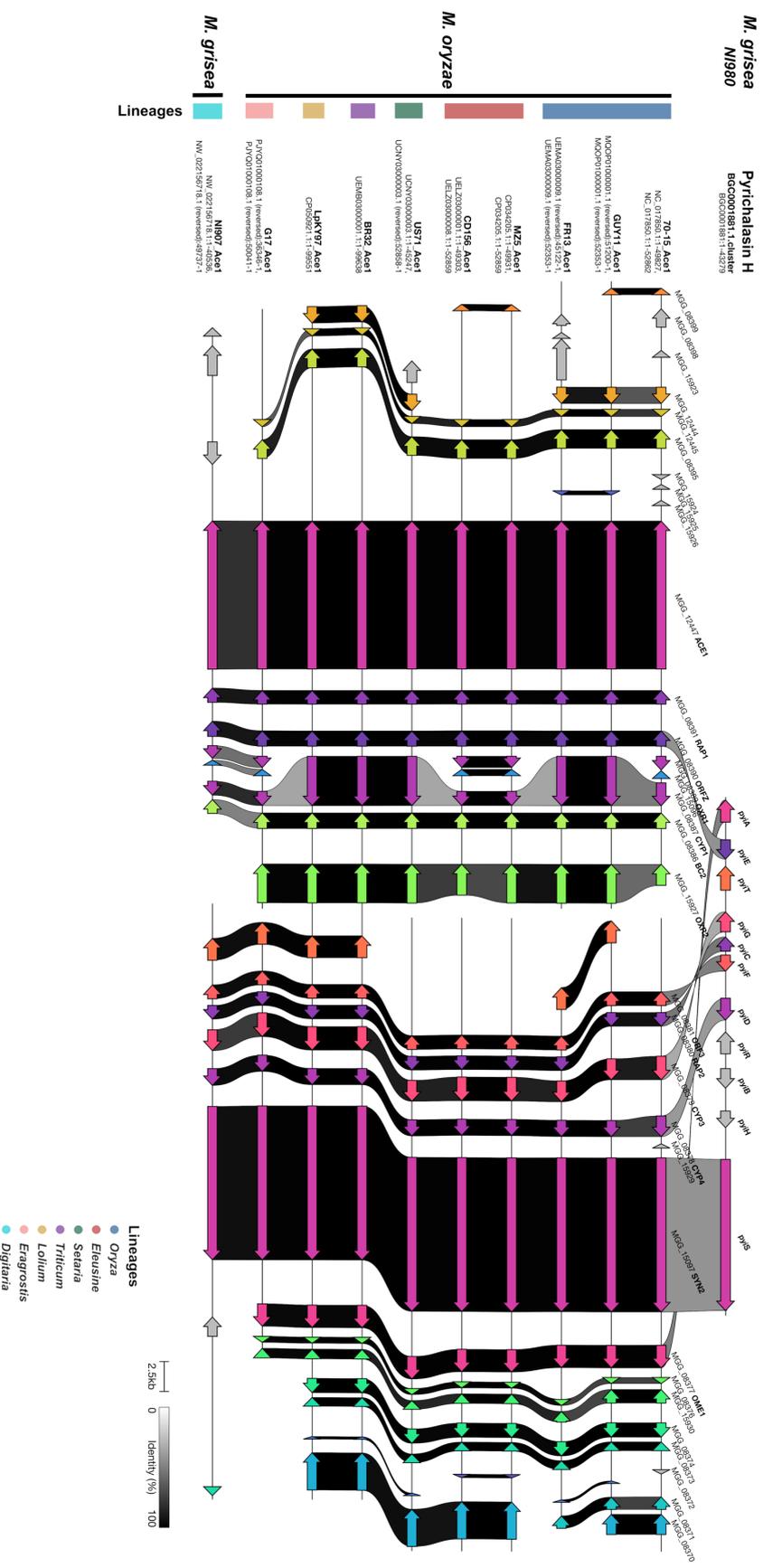


Figure 5.7: Analysis of homology between ACEI gene cluster family in *M. oryzae* and *M. grisea* and reference BGCs in MIBiG database. The map depicts comparison of BGC loci with reference BGCs associated with pyrichalasin H in MIBiG. The shaded area between any two arrows denotes degree of homology (0 to 100%; white to black, respectively) between the two sequences.

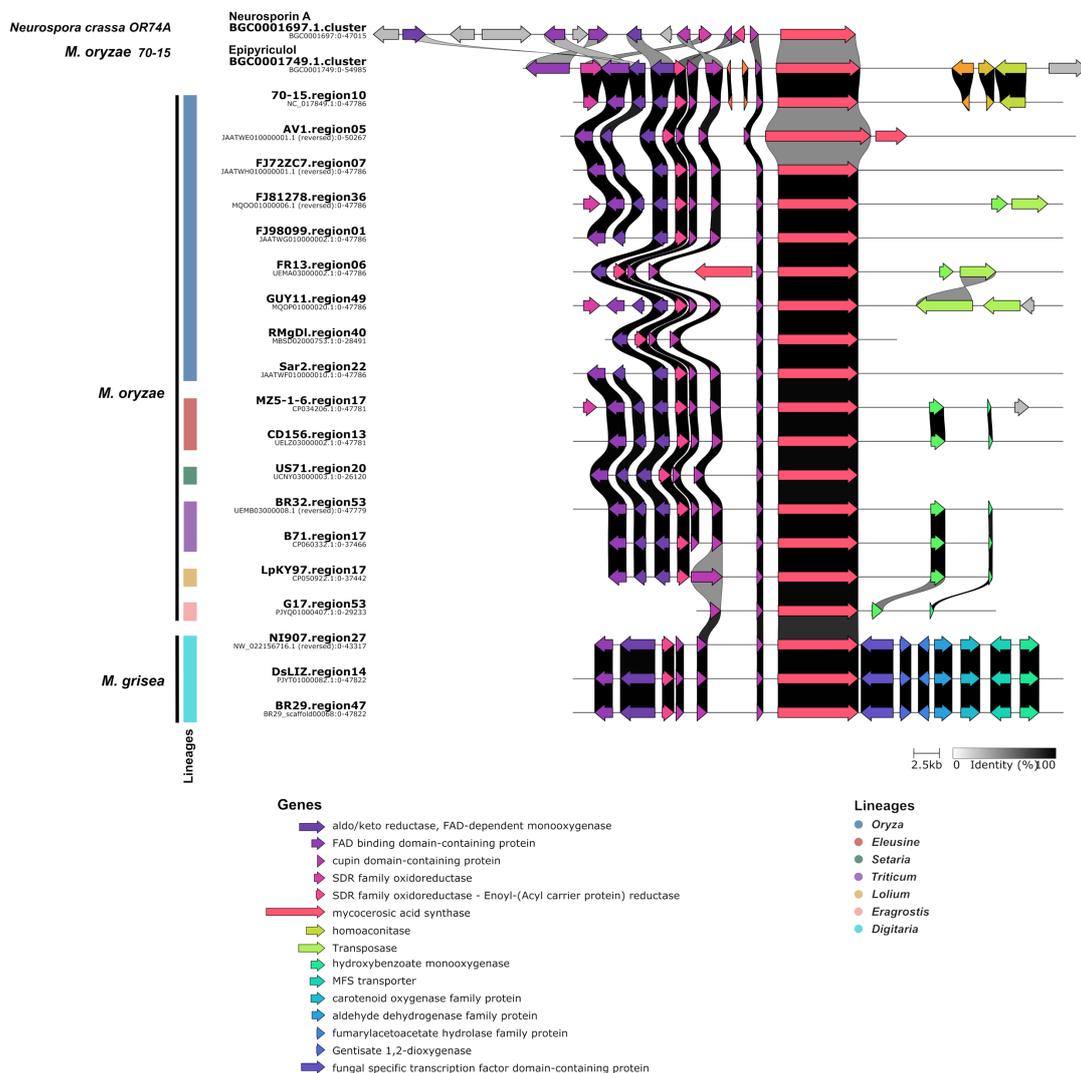


Figure 5.8: Analysis of homology between putative epipyrliculol-associated gene cluster family in *M. oryzae* and *M. grisea* and reference BGCs in MIBiG database. The map depicts comparison of BGC loci with reference BGCs associated with epipyrliculol in MIBiG. The shaded area between any two arrows denotes degree of homology (0 to 100%; white to black, respectively) between the two sequences.

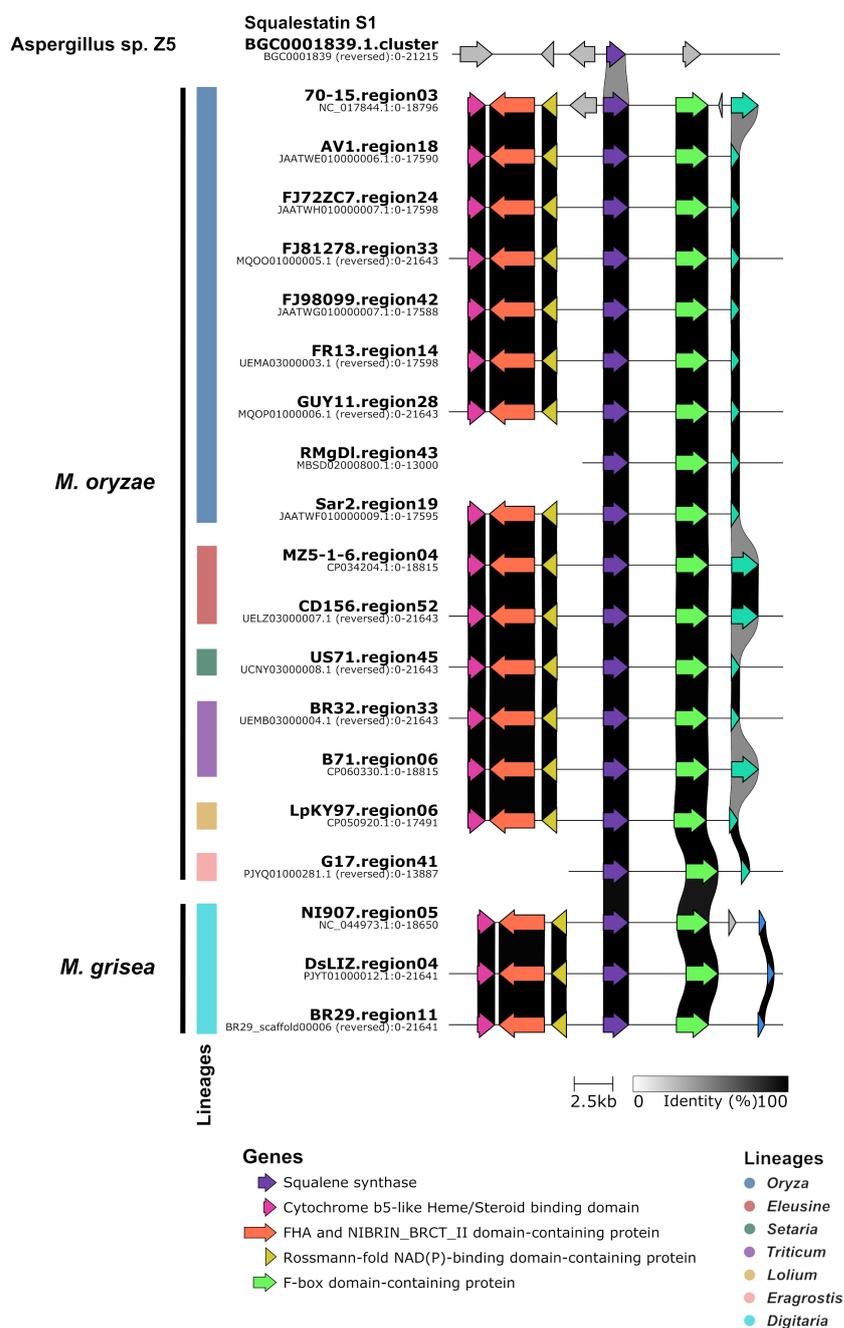


Figure 5.9: Homology of putative Squalestatin S1-associated *M. oryzae* and *M. grisea* gene cluster family with reference BGCs in MIBiG database. The map depicts comparison of BGC loci with reference BGCs associated with Squalestatin S1 in MIBiG. The shaded area between any two arrows denotes degree of homology (0 to 100%; white to black, respectively) between the two sequences.

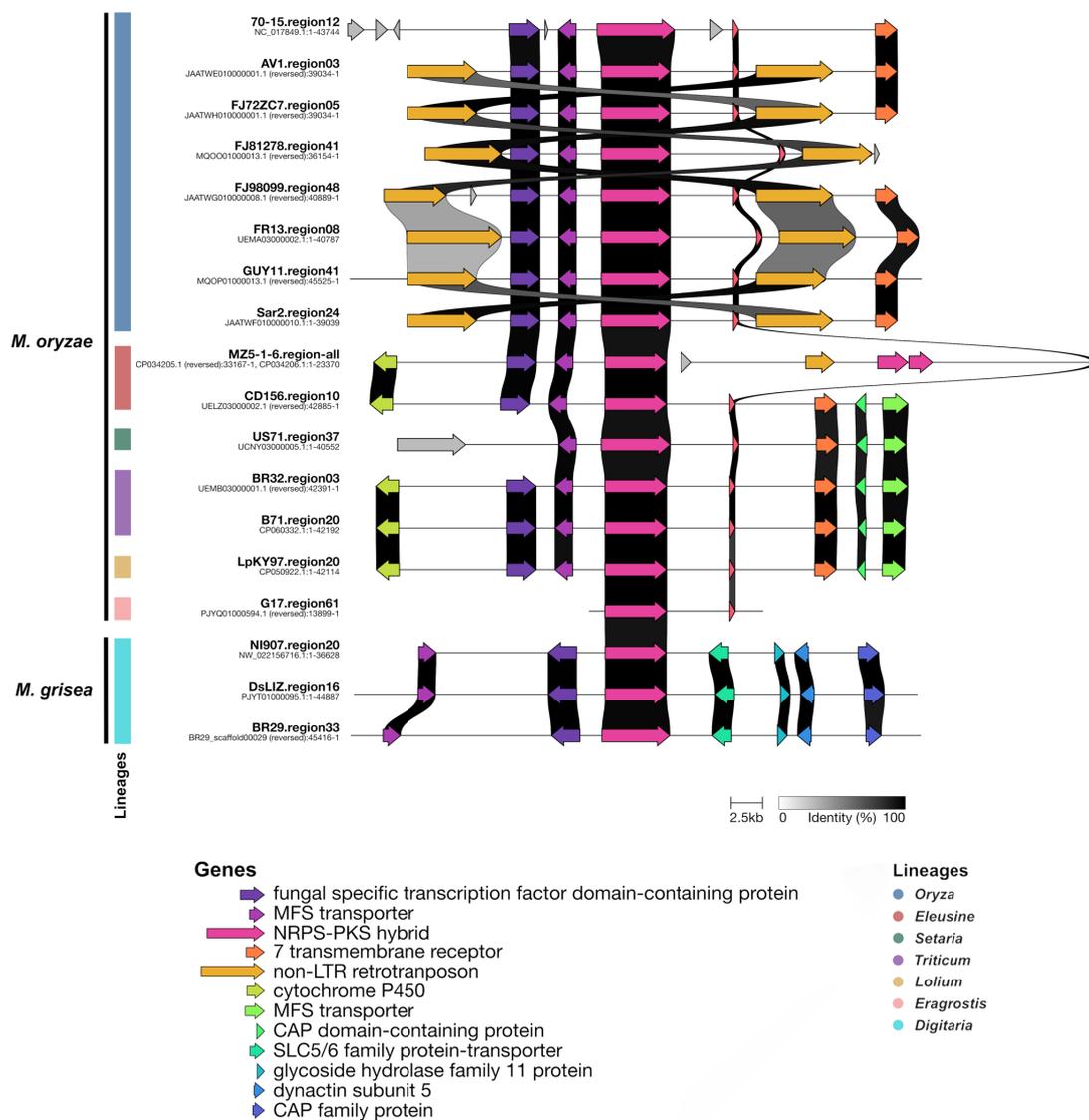


Figure 5.10: Analysis of homology between Tenuazonic acid-associated gene cluster families in host-adapted *M. oryzae* and *M. grisea* isolates. The shaded area between any two arrows denotes degree of homology (0 to 100%; white to black, respectively) between the two sequences.

In the case of BGC-TLE, although the standalone core biosynthetic NRPS gene therein is present in all the strains, there were notable differences in the gene structure (**Fig. 5.12**). Through manual curation of the NRPS gene, it was discerned that the *Triticum* and *Lolium* lineages had a deletion of 557 bps, corresponding to second exon observed in the *Eleusine* lineage (CD156 strain). Furthermore, the presence of a stop codon within the Amp-binding domain of the NRPS gene in the *Triticum* and *Lolium* lineages suggested pseudogenization of the NRPS gene and, consequently, the possible non-functionality of the corresponding BGC in those strains (**Fig. 5.13**). It is conceivable that BGC-TLE may either play a pivotal role in virulence exclusively on *Eleusine* host plants, or its product could serve an avirulence effector-like function in *Triticum* and *Lolium* hosts. In response to an evolutionary arms-race, pathogen from the *Triticum* and *Lolium* lineages have likely undergone adaptations leading to the loss of a functional NRPS gene, allowing them to continue infecting these hosts. Therefore, among the three candidate BGCs, our further exploration has focused on BGC-O1, as it holds the potential for involvement in specialization on the rice host.

5.4 Identification of a novel reducing polyketide synthase BGC unique to *Oryza* lineage

BGC-O1 was detected in 23 out of the 24 strains belonging to *Oryza* lineage that were employed in this study, as well as in a single strain from the *Eragrostis* lineage (**Fig. 5.14**). This GCF seems to encompass two distinct networks, and a comparative analysis of the genomic loci indeed clearly distinguishes between two potential BGCs - one conserved across all the lineages, while the other, BGC-O1 is specifically found in the *Oryza* and *Eragrostis* lineages (**Fig. 5.14**). BGC-O1 comprises a novel core biosynthetic gene, a reducing type I polyketide synthase (rPKS) gene, encoded by MGG_08236, along with adjacent tailoring genes. These tailoring genes include one methyl transferase (MGG_15107), one Co-A transferase (MGG_15108) and two cytochrome P450 monooxygenases (MGG_12496 and MGG_12497).

The predominance of BGC-O1 in the *Oryza*-specific lineage strongly suggests that the product generated by the core PKS might play a role in specialization or adaptation to the rice host. Crucially, it is noteworthy that BGC-O1 did not exhibit any similarity with any of the reference BGCs present in the MIBiG database.

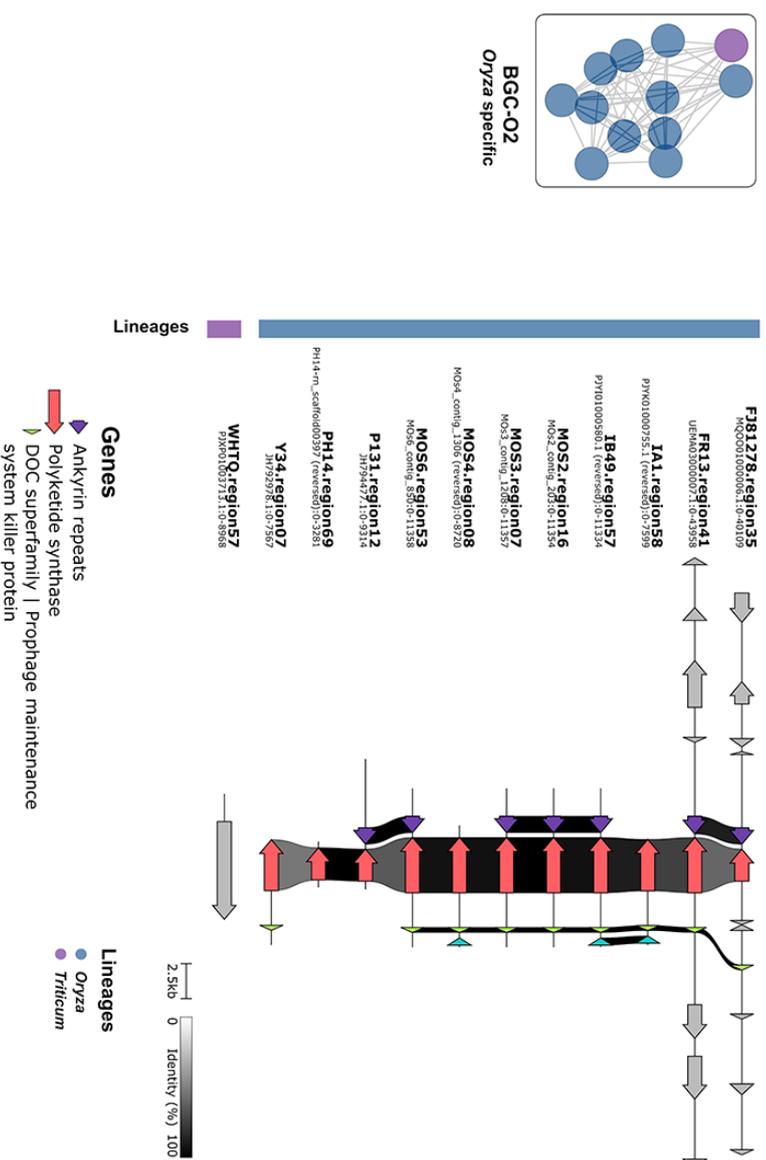


Figure 5.11: BGC-O2 is predominantly present in *Oryza*-specific lineage of *M. oryzae*. BIG-SCAPE analysis showing BGC-O2 GCF with all the BGCs specifically present in 11 genomes from *Oryza* lineage and one genome from *Triticum* lineage. The shaded area between any two arrows denotes degree of homology (0 to 100%; white to black, respectively) between the two sequences.

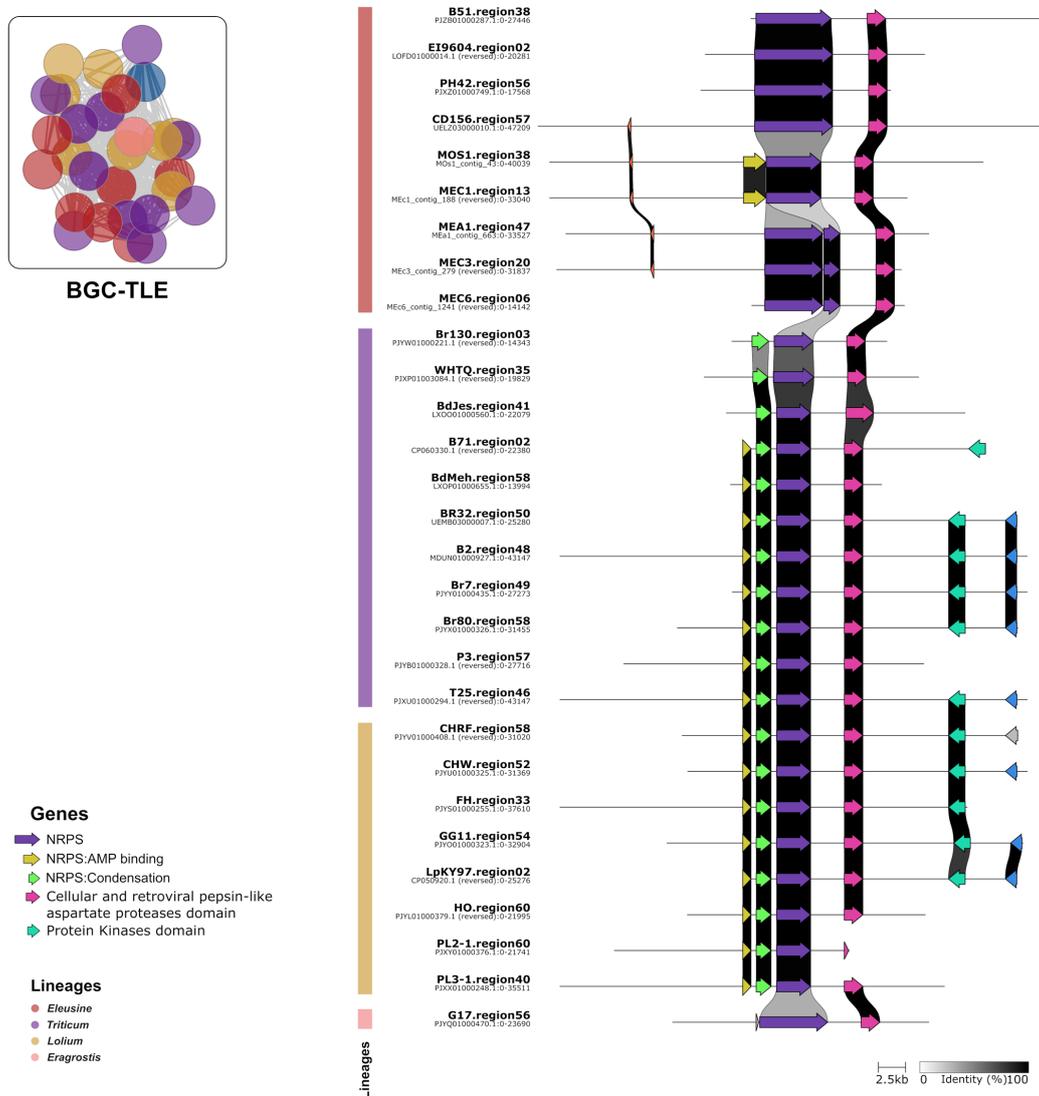


Figure 5.12: BGC-TLE is predominantly present in *Triticum*, *Lolium* and *Eleusine*-specific lineages of *M. oryzae*. BiG-SCAPE analysis showing BGC-TLE GCF with all the BGCs specifically present in 12, 8, and 9 genomes from *Triticum*-, *Lolium*- and *Eleusine*-specific lineages, respectively. The shaded area between any two arrows denotes degree of homology (0 to 100%; white to black, respectively) between the two sequences.

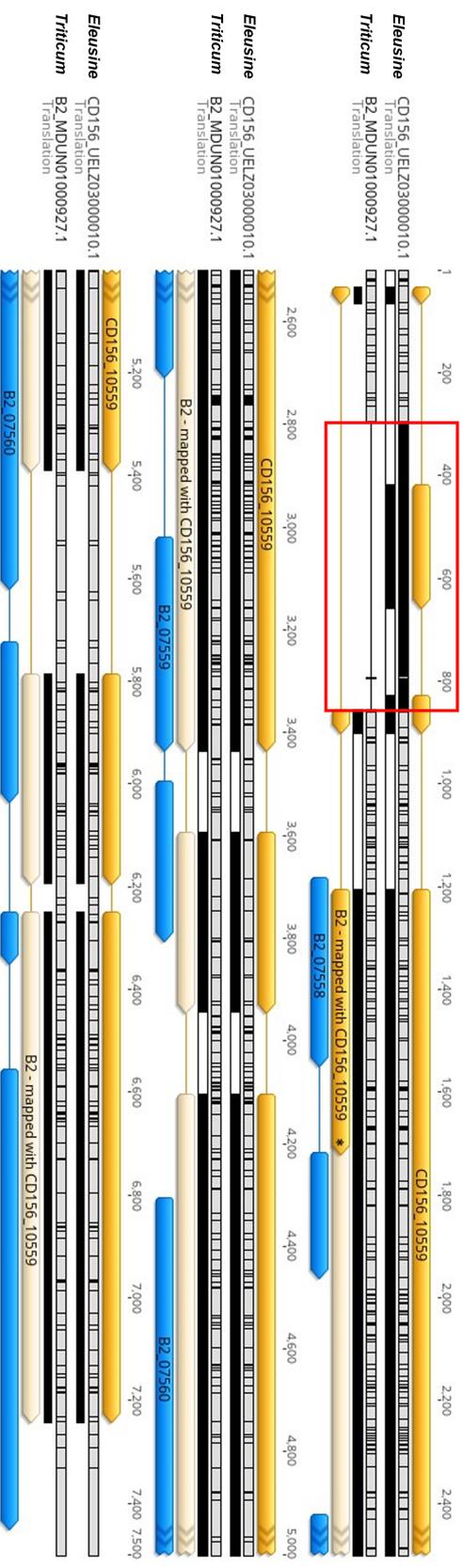


Figure 5.13: Comparison of core biosynthetic gene region of BGC-TLE shows the pseudogenization in NRPS gene belonging to *Triticum* lineage. Alignments between representative strains CD156 (*Elusine*) and B2 (*Triticum*) shows the presence of premature stop codon depicted as * in AMP-binding domain of NRPS gene. Yellow arrows represent the gene models in B2 in accordance with the one in CD156, whereas blue arrows in B2 strain represents the earlier predicted gene model by Augustus. Red box displays the deletion of 557 bps in B2 strains, which corresponds with the second exon of CD156_10559 NRPS gene.

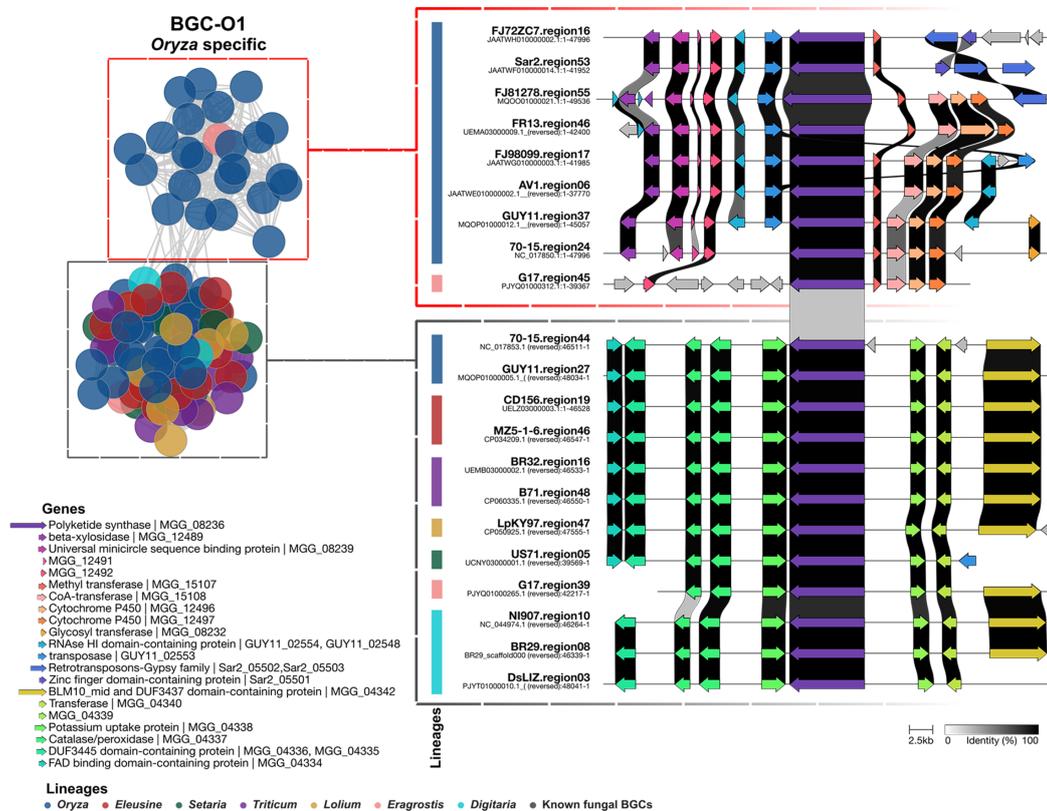


Figure 5.14: BGC-O1 is predominantly present in *Oryza* lineage of *M. oryzae*. BiG-SCAPE network showing BGC-O1 present in 23 genomes from *Oryza* and one genome from *Eragrostis* lineages. Conservation of the BGC in selected strains is depicted using Clinker. Vertical bars besides the name of strain are colored according to the lineages. Red and gray dashed boxes separate BGC-O1 and BGC-O1-like clusters, respectively.

This observation suggests that BGC-O1 represents a distinctive gene cluster likely involved in the biosynthesis of a novel secondary metabolite, associated with the specialization on the rice host.

We were curious to determine whether the specificity of the BGC-O1 cluster region was exclusive to the *Oryza* lineage, and thus consequently we assessed the conservation or variability of the flanking genomic regions surrounding BGC-O1 in different lineages. To achieve this, we aligned representative high-quality genome assemblies from each lineage against the reference genome assembly of the 70-15 strain, which belongs to the *Oryza* lineage. This allowed us to evaluate synteny at a global level. Our analysis revealed that BGC-O1 is situated in the sub-telomeric region of chromosome 2 (NC_017850.1) and is

located approximately 528 kb downstream of the *ACE1* gene cluster in the reference strain 70-15. When comparing the macrosynteny between 70-15 and GUY11, both strains from the *Oryza* lineage, we observed conservation of ~134 kb upstream and ~117 kb downstream flanking regions, encompassing the BGC-O1 cluster (**Fig. 5.15A**). Notably, the immediate flanking regions, which span ~10 kb, exhibited a translocation between the GUY11 and 70-15 genomes (blue lines crossing over within the orange band; black arrowhead; **Fig. 5.15A**). A similar comparison between 70-15 and FR13 strains displayed partial conservation with ~14 kilobases of the BGC-O1 region retained in the FR13 strain, while an approximately 9.5 kilobase segment from the BGC-O1 region in the 70-15 strain was relocated further downstream on the same contig in the FR13 strain (red arrowhead; **Fig. 5.15A**). Furthermore, although ~102 kilobases of the upstream flanking region remained syntenic between 70-15 and FR13, the downstream flanking region exhibited several structural rearrangements (**Fig. 5.15A**).

In contrast, when comparing with non-*Oryza* lineages, such as *Eleusine* and *Setaria*, a significant loss of synteny was observed in the sub-telomeric region of chromosome 2. While the upstream flanking region located ~97 kb upstream of BGC-O1 in 70-15 strain, exhibited synteny with a corresponding genomic locus in MZ5-1-6, CD156 and US71, the BGC-O1 cluster and downstream sequences were either missing or exhibited limited conservation in the genomes of these lineages (**Fig. 5.15B**).

Further, in order to understand the genetic diversity in these SM-BGCs, we looked for overall synteny between the two genomes (70-15 and MZ5-1-6) with chromosome-level assemblies. The orthologous gene-pairs, among all the genes present in SM-BGCs, were investigated using a Bi-directional Best Blast Hit approach, where we identified 529 orthologous SM gene pairs between the two isolates (**Table 5.2**). Whereas, we identified 638 orthologous gene pairs, when SM genes of MZ5-1-6 were compared with total genes of 70-15. This suggests that the orthologs of 109 SM-BGC genes of MZ5-1-6 are likely located in the regions outside of the predicted SM-BGCs in the 70-15 genome, likely due to differential rearrangements in these two genomes.

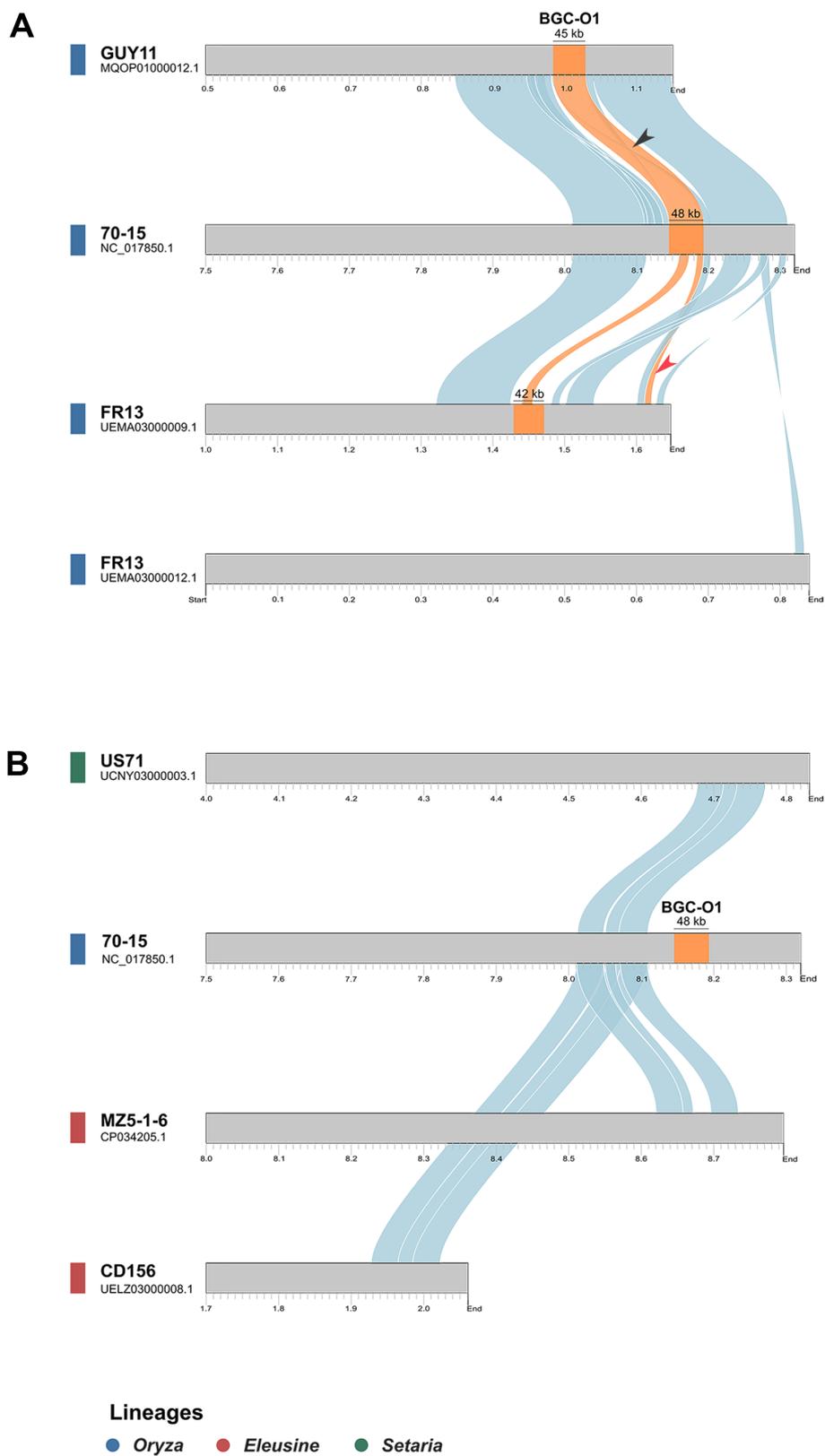


Figure 5.15: Genomic localization of BGC-O1 and synteny between the *Oryza* and non-*Oryza* lineages. A) Synteny analysis using pairwise comparison within *Oryza* lineage. Genomes of GUY11 and FR13 aligned individually with reference genome 70-15. The syntenic BGC-O1 locus (orange) and the flanking regions (blue) on chromosome 2 of 70-15 are depicted. The differential lengths of BGC-O1 in different isolates are marked with a bar and corresponding length in kilobases. The black arrowhead depicts genomic rearrangement (swapping) in the flanking ~10 Kb region. The red arrowhead marks the rearrangement of ~9 Kb region of the BGC-O1 in FR13. B) Synteny analyses of representative genome assemblies belonging to *Setaria* (US71), *Eleusine* (MZ5-1-6 and CD156) and *Oryza* (70-15) lineages. The BGC-O1 locus (orange) and flanking region (blue) sequences from US71, MZ5-1-6 and CD156 aligned with that of 70-15 are depicted.

The Circos plots were constructed to determine the synteny of these orthologous SM-genes between the genomes of the rice isolate (70-15) and finger millet isolate (MZ5-1-6) (**Fig. 5.16**). While, most of the genes were highly syntenic, a total of 36 SM-gene pairs located on to different chromosomes in the two genomes. Notably, 30 out of 36 SM-genes, which located on the chromosome 6 of MZ5-1-6, were rather found on chromosome-1 of 70-15. Such large chromosomal translocation events have been reported earlier for MZ5-1-6 and 70-15 (Luciano et al. 2019). Thus, the differential genomic rearrangement, especially with respect to the genes involved in secondary metabolism, in the rice and millet isolates might have some evolutionary role in adaptation to a specific host plant.

Table 5.2: Summary of ortholog pairs between the genomes of millet isolate (MZ5-1-6) and rice isolate (70-15).

Pairs	Total Genes MZ5 vs 70-15	MZ5 SM-genes vs 70-15 total genes	MZ5 SM-genes vs 70-15 SM- genes
Total number of orthologs	10813	638	529
Total number of orthologs rearranged	671	84	36
Total number of orthologs rearranged between MZ5_Chr6 and 70-15_Chr1	549	50	30

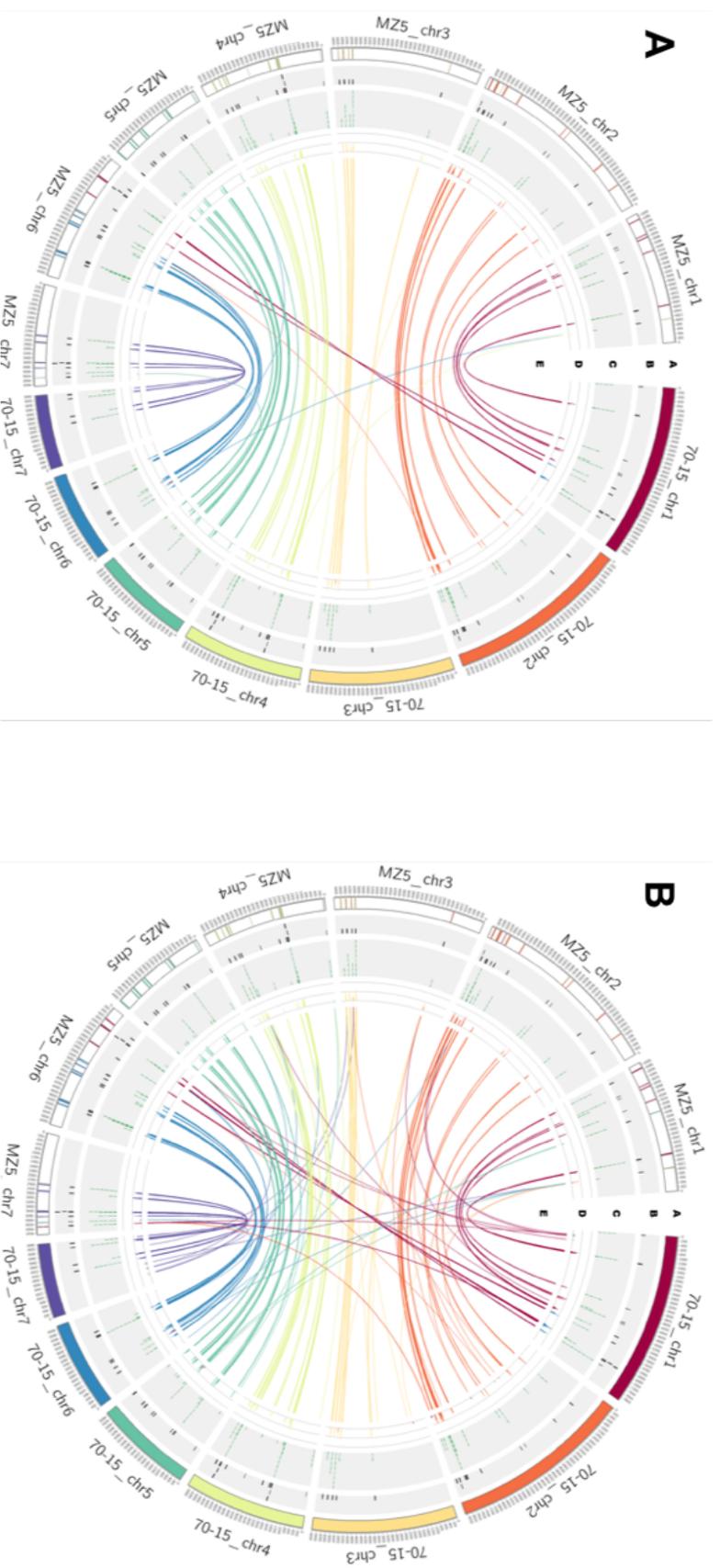


Figure 5.16: Circos plots showing synteny among SM related genes between representative strains from *Oryza* and *Eleusine* lineages. The orthologous gene pairs represented by connecting links between (A) SM-genes of both MZ5-1-6 and 70-15 and (B) SM-genes of MZ5-1-6 and total genes of 70-15. In each plot, individual tracks (A to E) denote different characteristics – Track A, chromosomes scaled according to their length; Track B, candidate clusters positioned accordingly on the chromosomes; Track C, core biosynthetic genes with respect to their positions

on the chromosomes; Track D, histogram showing orthologous gene-pairs between the two genomes, with color codes according to those of the chromosomes of 70-15; and Track E, links representing orthologous gene-pairs between the two genomes, with color codes same as those for the 70-15 chromosomes.

5.5 Evolutionary history of the novel polyketide synthase gene

In our quest to unravel the evolutionary history of BGC-O1, we embarked on a search for orthologous and closely related counterparts of the MGG_08236 rPKS protein within the *Pezizomycotina* group from the MycoCosm repository (Grigoriev et al., 2014;). Interestingly, only two closely resembling homologues were identified from the fungus *Colletotrichum eremochloae*. Employing a phylogenetic analysis, it became evident that the rPKS 670826 from *C. eremochloae* is a true orthologue of MGG_08236 rPKS, being in the adjacent clade to the *M. oryzae* lineage (**Fig 5.17A**). Meanwhile, the other paralogue in *C. eremochloae*, 679399, falls within a clade restricted to the *Colletotrichum* genus and constitutes a sister clade to the MGG_08236 rPKS clade. Upon conducting a comparative analysis of the genomic loci encompassing both homologues in *C. eremochloae* and BGC-O1 in *M. oryzae*, it was revealed that these two clusters shared the rPKS and two cytochrome P450 genes with BGC-O1. However, the methyltransferase and Co-A transferase genes were found as single copies elsewhere within the *C. eremochloae* genome (**Fig. 5.17B**). Notably, the orthologous BGC in *C. eremochloae* shares an average nucleotide identity of 76% when compared to the *M. oryzae* BGC. Intriguingly, both BGCs in *C. eremochloae* share only 50% nucleotide identity with each other.

5.6 BGC-O1 genes are expressed specifically during host invasion

In our study, we conducted PCR targeting the open reading frame (ORF) region of the rPKS gene MGG_08236, using genomic DNA extracted from *M. oryzae* strains from rice and finger millet host plants (**Fig. 5.18A**). The MGG_08236 gene was found to be present in all the examined *Oryza* strains studied, except for MOS3, which was stood out as the sole *Oryza* lineage strain lacking the BGC-O1 as investigated by our in-silico analysis.

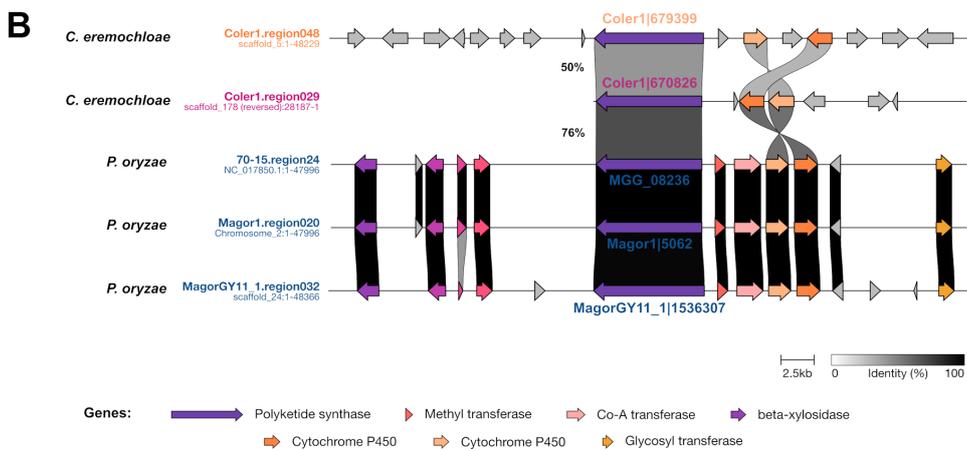


Figure 5.17: Evolutionary origin of the reducing polyketide synthase (PKS) MGG_08236 gene. A) Maximum-likelihood phylogenetic tree using protein sequences of reducing PKS genes from BGC-O1 and related gene cluster families, as well as homologues retrieved from MycoCosm repository. Tip labels depicting sequences from *Oryza* and *Eragrostis*-specific clade are marked with blue and pink background, respectively. The *Colletotrichum eremocloae* ortholog (jgi.p_Coler1_670826), closer to MGG_08236, is denoted with pink label; whereas the more distant paralogue (jgi.p_Coler1_679399) is shown in orange. Gray shaded triangles denote collapsed clades with distant sequences. Branches were supported by > 95% Bootstrap values indicated with gray circles at the nodes. The tree is rooted at midpoint. B) Comparative analysis, using Clinker tool, of the BGC-O1 or BGC-O1-like loci from three *Oryza*-specific and *C. eremochloae* genomes.

The MGG_08236 ORF was notably absent in all the *Eleusine* strains investigated, as well as in the MOS1 and MOS4 strains, both of which are placed outside of *Oryza* lineage according to our phylogenetic analysis (**Fig. 5.1B and 5.18A**).

Subsequently, we delved into the investigation of the expression profiles of the rPKS genes (MGG_08236) and three associated tailoring genes, namely methyl transferase (MGG_15107), CoA-transferase (MGG_15108) and cytochrome P450 (MGG_12496), all of which are members of the BGC-O1 locus. This was achieved through semi-quantitative RT-PCR analysis conducted at different stages of infection. Total RNA was isolated from fungal vegetative mycelia cultured in complete medium, barley leaves inoculated with *M. oryzae* strain and incubated for different time intervals, and uninoculated barley leaves as mock samples. These assays were conducted to assess the transcript levels of the aforementioned genes relative to those of β -Tubulin, serving as an endogeneous control. Our finding indicated that the expression of the core rPKS gene was either undetected or barely seen during vegetative growth – mycelium and pre-invasive stage - 12 hours post inoculation (hpi), respectively. However, transcript accumulation commenced during the progression of pathogenesis, with a substantial increase in expression at 24 hpi, followed by a consistent, albeit gradual, rise in expression levels until 72 hpi (**Fig. 5.18B and 5.18C**). The expression patterns of MGG_15107 and MGG_15108 tailoring genes displayed partial correlation with that of rPKS, albeit with a lower expression level (**Fig. 5.18B and 5.18C**). In contrast, MGG_12496 exhibited a different expression pattern with its highest expression levels observed in mycelium and at 48 hpi during infection (**Fig. 5.18B and 5.18C**).

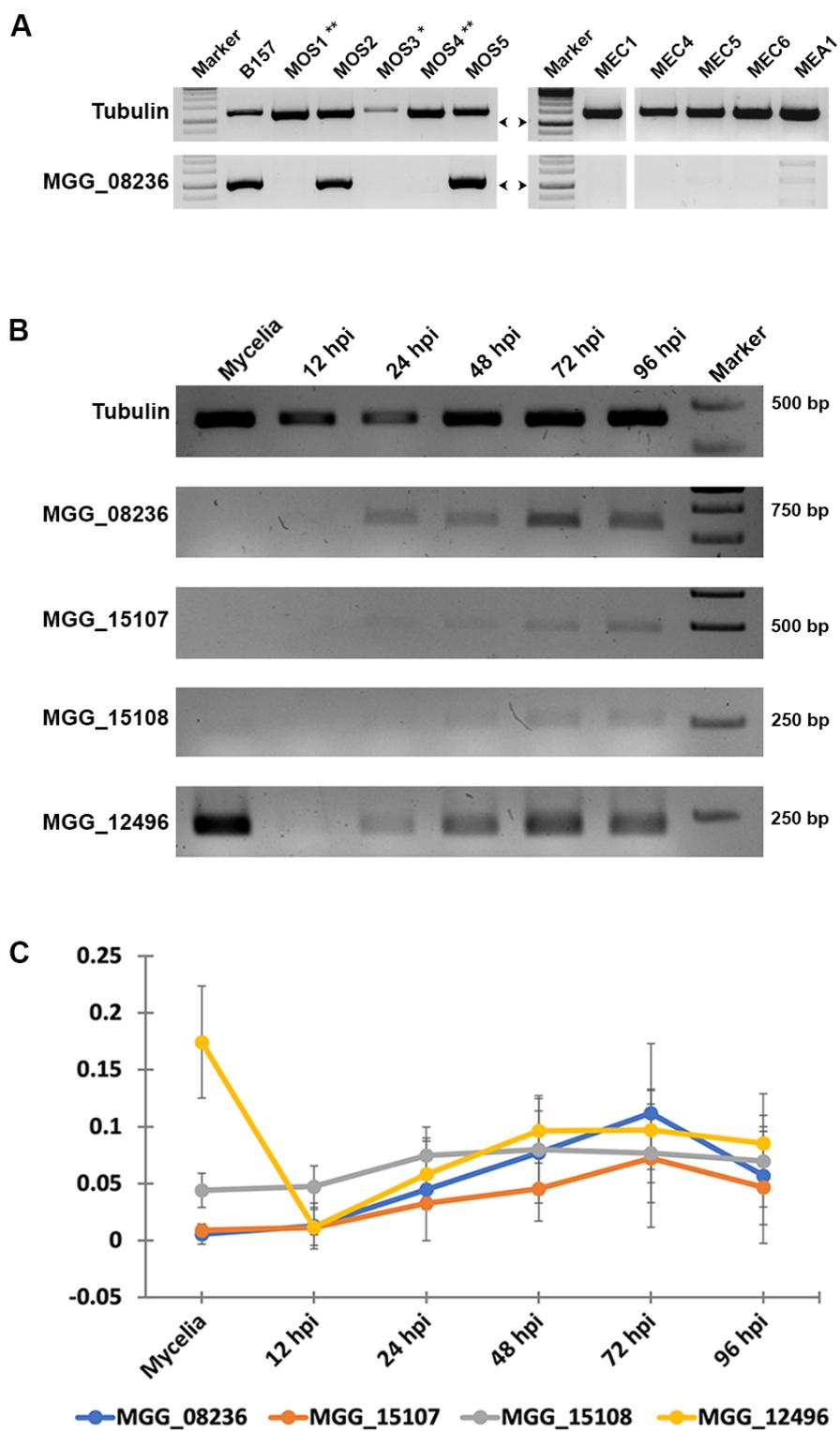


Figure 5.18: BGC-O1 genes are specifically expressed during pathogenic development.

A) Gel electrophoresis of PCR products displaying presence or absence of ORF region of rPKS MGG_08236 from genomic DNA of the indicated *M. oryzae* strains. B157, MOS2, MOS3 and MOS5 belong to the *Oryza* lineage; whereas, the MOS1, MEC1, MEC4, MEC5, MEC6 and MEA1 belong to the *Eleusine* lineage. * - the only *Oryza* strain that lacked the BGC-O1 in in-silico analysis. ** - the *Oryza* strains placed outside *Oryza* lineage. Arrowheads corresponds to the 1 Kb size of band from Marker. B) RT-PCR gel depicting expression of genes MGG_08236 (Polyketide synthase), MGG_15107 (Methyl transferase), MGG_15108 (Co-A transferase) and MGG_12496 (Cytochrome P450) in *Oryza*-specific strain B157 at different stages of barley infection (12 to 96 hpi) and during vegetative growth (mycelium) in complete medium. C) Quantification of the expression of BGC-O1 genes. The intensity of each band was measured using ImageJ, and relative gene expression was calculated relative to that of β -tubulin (MGG_00604) as an endogenous control. The data on expression of MGG_08236 represents mean \pm standard deviation of mean (SDM) from three independent biological experiments. Data on expression of the tailoring genes (MGG_15107, MGG_15108 and MGG_12496) represent observation from a single experiment.

These results indicate that the MGG_08236 PKS gene is specifically expressed during pathogenesis and potentially has a key role to play during host colonization. Additionally, our findings suggest that the BGC-O1 likely comprises only two co-regulated tailoring genes.

Altogether, our in-silico analyses identified a novel PKS gene cluster in the *Oryza*-specific lineage, which likely played a key role in shaping specialization of the blast fungus to rice host.