

# Supporting supervised learning in fungal Biosynthetic Gene Cluster discovery: new benchmark datasets

Hayda Almeida<sup>1,2</sup>

Adrian Tsang<sup>2</sup>

Abdoulaye Baniré Diallo<sup>1</sup>

<sup>1</sup>University of Quebec in Montreal, Montreal, Canada

<sup>2</sup>Concordia University, Montreal, Canada

Corresponding author: diallo.abdoulaye@uqam.ca

**Abstract**—Fungal Biosynthetic Gene Clusters (BGCs) of secondary metabolites are clusters of genes capable of producing natural products, compounds that play an important role in the production of a wide variety of bioactive compounds, including antibiotics and pharmaceuticals. Identifying BGCs can lead to the discovery of novel natural products to benefit human health. Previous work has been focused on developing automatic tools to support BGC discovery in plants, fungi, and bacteria. Data-driven methods, as well as probabilistic and supervised learning methods have been explored in identifying BGCs. Most methods applied to identify fungal BGCs were data-driven and presented limited scope. Supervised learning methods have been shown to perform well at identifying BGCs in bacteria, and could be well suited to perform the same task in fungi. But labeled data instances are needed to perform supervised learning. Openly accessible BGC databases contain only a very small portion of previously curated fungal BGCs. Making new fungal BGC datasets available could motivate the development of supervised learning methods for fungal BGCs and potentially improve prediction performance compared to data-driven methods. In this work we propose new publicly available fungal BGC datasets to support the BGC discovery task using supervised learning. These datasets are prepared to perform binary classification and predict candidate BGC regions in fungal genomes. In addition we analyse the performance of a well supported supervised learning tool developed to predict BGCs.

**Index Terms**—biosynthetic gene clusters, secondary metabolites, supervised learning, BGC, fungi, dataset

## I. INTRODUCTION

Natural products (NPs) are specialized bioactive compounds primarily produced by plants, fungi and bacteria. NPs are a vital source for drugs: from anti-cancer, anti-virus, and cholesterol-lowering medications to antibiotics, and immunosuppressants [1]–[3]. Unlike those in plants, genes involved in the biosynthesis of many NPs in bacteria and fungi are co-localized in the genome of organisms and usually organized as clusters of genes [4]. Gene clusters capable of producing NPs are known as Biosynthetic Gene Clusters (BGC).

The task of identifying new BGCs could potentially lead to the discovery of novel NPs to benefit human health. However this task involves complex and costly processes, as well as the analysis of large amounts of biological data. Development of

automatic tools that can support the identification of BGCs is therefore highly relevant. Various approaches have been used to develop such tools, such as data-driven methods, probabilistic methods, and supervised learning methods. In supervised learning the BGC discovery task can be represented as binary classification task. The goal in a binary classification task is to classify data instances as belonging to one out of two different categories. A binary classification BGC dataset would therefore be composed of positive and negative BGC instances.

Supervised learning has been previously used to predicting bacterial BGCs [5], [6] and shown to perform well. Supervised learning methods however are developed primarily based on annotated datasets, for which all instances are labeled as belonging to a specific class. Unlike for bacteria, the number of known fungal BGC data previously validated by curators is rather limited. The Minimum Information about a Biosynthetic Gene cluster (MIBiG) [7]<sup>1</sup> repository is one of the largest freely available BGC databases.

As an example of the disparity between known available BGC from bacteria versus fungi that has been annotated by curators, MIBiG holds over 1,196 bacteria BGCs, while only 206 are fungal BGCs<sup>2</sup>.

Generating fungal BGC datasets for supervised learning approaches imposes a few challenges. For instance, negative samples are needed for binary classification, and they are not directly provided by BGC databases just as annotated BGC data. To be able to support a robust classification approach, fungal BGC datasets used as input should include a variety of organisms and BGC types to properly represent fungal genomic profiles.

The availability of fungal BGC datasets could leverage the development of new supervised learning approaches to tackle BGC discovery in fungi. This work presents new datasets prepared to tackle fungal BGC discovery as a binary classification task. These datasets are constructed in such way that they include most variety of BGC types from different organisms, attempting to represent fungal genomic profiles to better suit

<sup>1</sup><http://mibig.secondarymetabolites.org/>

<sup>2</sup>As of July 2019.

the fungal BGC classification task. Finally we also analyse the usage of fungal BGC datasets with one of the state-of-the-art supervised learning methods developed for BGC discovery, DeepBGC [6].

## II. PREVIOUS WORK

In this section we present previous work on the availability of BGC data previously predicted or annotated by curators that can support BGC discovery, and previous work conducted towards developing automatic approaches to identify fungal BGCs. BGC databases and some of their characteristics are discussed in Section II-A. Previous work on predicting BGCs in fungi is presented in Section II-B.

### A. BGC Databases

Only a small number of open access BGC databases is currently available to support research on automatic tools to identify BGCs. The majority of entries in these databases corresponds to bacteria data, while only a small portion are fungal BGCs.<sup>3</sup> MIBiG is a BGC repository in which curated entries are submitted by curators, and added to the database in a format compliant with the Minimum Information about any Sequence (MIxS) framework data standard. It holds 206 fungi BGCs and 1,196 for bacteria. Clustermine360 [8] contains microbial polyketide synthases (PKS) and non-ribosomal peptide synthetases (NRPS) biosynthesis. It holds a total of 29 fungal BGCs, while over 900 are from bacteria. Clustermine360 entries are curated and submitted by curators, enriched with information from the National Center for Biotechnology Information (NCBI)<sup>4</sup>, and analysed with the antiSMASH [9] tool. The antiSMASH database [10] has 24,773 microbial BGCs predicted based on its homonymous tool. Unlike its bacteria version, the fungal version of antiSMASH does not provide a database of fungal BGCs to the best of our knowledge.

The Integrated Microbial Genomes: Atlas of Biosynthetic Gene Clusters [11] (IMG/ABC) database contains BGCs predicted using the ClusterFinder algorithm [12]. IMG/ABC holds 127 fungal BGCs and 1,025 from bacteria.

These databases are not connected. Since it is likely that there are overlaps among the different databases, the number of unique fungal BGCs could be even smaller. The small proportion of fungal BGCs across databases is an example of the challenges in developing automatic tools to tackle BGC discovery in fungi. This work proposes new publicly available datasets to be an input of supervised learning tools to predict fungal BGCs, based on MIBiG and orthologous genes. The details on our datasets and their analysis are discussed in Section III.

### B. BGC discovery in Fungi

Significant effort has been put towards developing approaches to discover BGCs [2], [3]. The majority of approaches focused on processing bacterial data, while some of them are specially focused on fungi. Identifying BGCs remains

a challenging task specially in fungal genomes, due to the diversity of clusters [13].

Previous work on fungal BGC discovery made use mostly of data-driven methods, which are heavily based on the analysis of the input or output data and require fine parameter-tuning. These methods required as input the genome sequence combined with transcription data [14], [15], or gene functional annotations [16], as well as both nucleotide and amino acid sequences [17]. [14] and [15] focused on analysing similar gene expression levels, while [15] used virtual clusters. [14] looked at motif co-occurrence in promoters around anchor genes, and [17] analysed homologous genes through a comparative genomics approach.

Such data-driven methods are less dependent on curated BGC data, which are time consuming to obtain, but they all present limitations. [16] requires gene functional annotations, which may not be available, and [14] relies heavily on manual curation of output to achieve the expected results. A very limited BGC prediction scope is considered in [18] and [17]. Both approaches are developed based on biological sequences from a single species, and they also require fine parameter-tuning. Such limitations of data-driven methods can restrict their ability to generalize to new data, and as a consequence compromise the discovery of novel BGCs.

Likely due to the larger availability of curated BGC data, probabilistic [9], [12], [19] and machine learning approaches [5], [6] have been more explored in bacteria compared to fungi, and shown to perform well. Probabilistic and machine learning approaches could be beneficial for BGC discovery, since by nature they are more capable of generalizing given new data, and will likely perform better at identifying data patterns and discovering novel BGCs, when compared to data-driven methods. In this study we also analyse the performance of a supervised learning approach developed to tackle BGC discovery using the fungal BGC datasets proposed by our work. The details on our experimental setup are further discussed in Section III.

## III. METHODOLOGY

Some of the challenges in generating fungal BGC datasets for binary classification are the need of negative instances, which are not directly provided in BGC databases; and accounting for a variety of organisms, BGC types, and also fungal genomic profiles. The availability of new fungal BGC datasets however could potentially motivate the development of supervised learning approaches to tackle fungal BGC discovery.

In this work we propose new publicly available fungal BGC datasets to support supervised learning approaches tackling BGC discovery as a binary classification task. We present here the methodology adopted to prepare fungal BGC datasets and their analysis using a supervised learning method, with the goal of analysing the method performance in fungal BGC data.

Details on our proposed fungal BGC datasets are presented in Section III-A. Section III-B presents the test datasets with which we analysed the performance of classification models

<sup>3</sup>Number of entries for databases are reported as of July 2019.

<sup>4</sup><https://www.ncbi.nlm.nih.gov/>

built on fungal BGC datasets. In Section III-C we provide details on the parameters considered in our analysis based on a supervised learning method, as well as the classification models considered.

### A. Proposed Datasets

Supervised learning was shown to perform well at BGC discovery in previous work that focused on handling bacteria data [5], [6]. Given that annotated data are needed to perform a supervised learning approach, we propose here fungal BGC datasets to support the development of this approach for fungi.

As mentioned in Section I, positive and negative instances are needed to perform fungal BGC discovery as a binary classification task using supervised learning. To create our fungal BGC datasets, we extracted and filtered positive instances from the MIBiG [7] repository, previously presented in Section II-A. MIBiG has the highest number of unique fungal BGCs among the BGC databases previously presented. Additionally, MIBiG BGCs were annotated and submitted by the research community, unlike BGCs in other databases that were automatically predicted.

From all MIBiG instances, we have selected only the fungal BGC subset, excluding BGCs belonging to *Aspergillus niger* (*A. niger*) to avoid overlaps during the test phase, resulting in a total of 200 positive instances.

We generated synthetic negative instances collecting and integrating orthologous genes from OrthoDB<sup>5</sup> [20]. Orthologs are homologous genes descendants from a single gene of a last common ancestor. The OrthoDB database contains protein-coding genes that represent the last common ancestors given a specific phylogeny radiation of a species, and are therefore known to retain ancestral function [20]. Orthologs represent regions conserved across species. They can correspond to a relevant negative instances for BGC discovery. this is due to the fact that fungal BGCs are known to have opposite characteristics and show large genomic diversity even in otherwise closely-related or same genus species [13]. Genes belonging to fungal BGCs have been previously referred to as “species-specific” [21], unlike orthologs.

Orthologous genes have been previously used to discover BGCs in fungi. In [17], the authors presented an alignment-based approach focused on identifying syntenic block regions, which are more likely to contain orthologs and less likely to contain BGCs. Non-syntenic blocks were then used to search for candidate BGCs and to better define candidate cluster boundaries. The approach in [17] was explored in small set of 10 filamentous fungi. The results showed good performance, predicting correctly 21 out of 24 fungal BGCs.

In this study we selected the fungal OrthoDB subset to construct the synthetic negative BGC instances. The OrthoDB fungal subset contains a total of 5,083,652 non-redundant orthologs. To avoid potential overlaps, we performed a BLAST analysis between the fungal subsets of both OrthoDB and MIBiG. We discarded 11,000 ortholog matches found using the BLAST parameter *eval*ue (expected value) set to  $1e - 60$ .

To generate synthetic negative instances, we then concatenated the amino acid sequence of fungal orthologs using a fixed length of 7,000 amino acids to create synthetic gene clusters. The 7,000 amino acid length is chosen since it corresponds to the average length of fungal BGC amino acid sequences in MIBiG. Figure 1 shows an example of positive instances in our datasets and negative instances being generated from OrthoDB orthologs. After processing OrthoDB fungal orthologs a total of 693,195 synthetic negative clusters were generated.

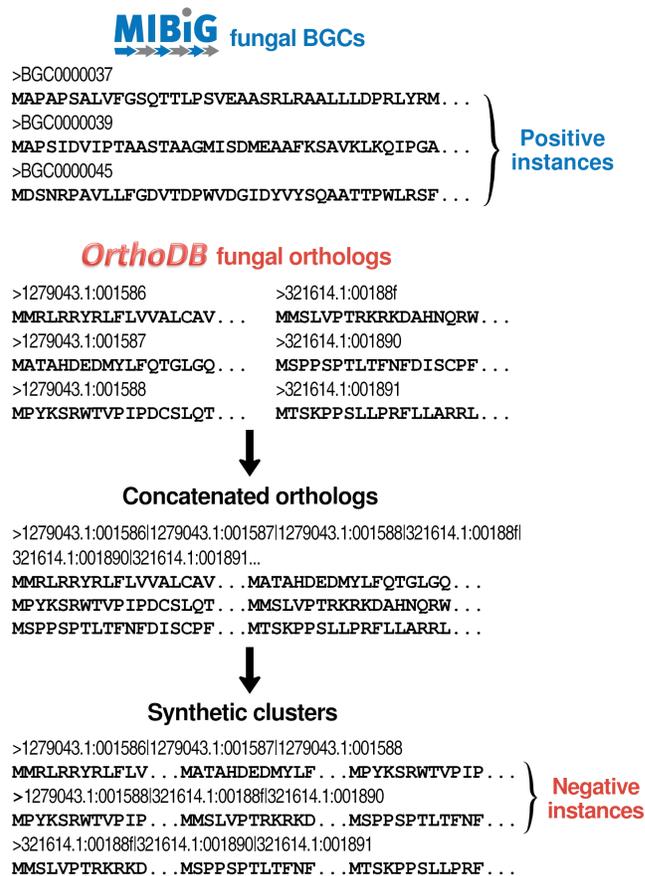


Fig. 1. Example of positive instances and the process to generate synthetic negative instances from orthologs

The MIBiG fungal subset and the pool of OrthoDB synthetic negative clusters were then considered to generate fungal BGC datasets with different distributions of positive and negative instances. Among the MIBiG fungal subset the annotated BGC regions corresponded in average to  $\approx 1\%$  of the total genome length of an organism, which provides a hint on the imbalance in class distribution that can be seen in a real test case scenario. Due to the natural imbalance of BGC regions versus non-BGC regions in a genome, we are interested in analysing the performance of a supervised learning approach based on datasets with various distributions of positive and negative instances. To analyse this aspect, we generated fungal BGC datasets with varying distributions by increasing the number of synthetic negative instances randomly selected from the

<sup>5</sup><http://orthodb.org/>

OrthoDB synthetic negative clusters pool. Table I shows the positive vs. negative distributions in each dataset.

TABLE I  
DISTRIBUTION OF INSTANCES ACROSS FUNGAL BGC DATASETS

Dataset	Train		Validation	
	Pos	Neg	Pos	Neg
50%-50%	160	160	40	40
40%-60%	160	240	40	60
30%-70%	160	373	40	93
20%-80%	160	640	40	160
10%-90%	160	1,440	40	360
05%-95%	160	3,040	40	760
01%-99%	160	15,840	40	3,960

To generate classification models based on a supervised learning method, we extracted Pfam [22]<sup>6</sup> IDs from the positive and negative instances. All datasets were converted into `pfamtsv` format [6], which is required as input in the supervised learning approach applied in this work. For each dataset, 80% were randomly selected for the training phase, while 20% were held out for the validation phase, as shown in Table I.

### B. Test Datasets

To analyse the performance of classification models built based on fungal BGC datasets, we selected a fungal genome from the *Aspergillus* genus to represent a real test case scenario. *Aspergillus* is the most frequent genus among fungal species in MIBiG, together with *Penicillium*. For this evaluation we focused specifically on the *A. niger* species. *A. niger* is a genome of interest due to its biological diversity and major relevance to industrial processes [23]. In [24] the authors present manual annotation of BGCs in *Aspergilli*, among which a total of 79 BGCs are found in *A. niger*.

To generate candidate clusters for the test phase, we collected a manually curated *A. niger* genome sequence made publicly available through the Genozymes project<sup>7</sup>. We generated test candidate clusters by considering a sliding window of 30,000 nucleotides in the *A. niger* genome. The 30,000 sliding window length is defined based on the average length of the nucleotide sequence of MIBiG fungal BGCs. A similar approach was previously applied in fungal BGC discovery to generate virtual clusters [15].

The 30,000 sliding window was shifted along the genome using either a 50% or a 30% overlap. The overlaps in a sliding window mean that each test candidate cluster will contain the last 15,000 nucleotides (if a 50% overlap) or the last 9,000 nucleotides (if a 30% overlap) of the immediate previous candidate cluster. With the strategy of generating candidate clusters using overlaps, we are more likely to cover regions in between two or more genes. Figure 2 shows an example of candidate clusters being generated from *A. niger* genes using overlaps. The test datasets based on a 50% overlap contains a total of 1,184 candidate clusters, while the test datasets based on a 30% overlap contains a total of 846 candidate clusters.

<sup>6</sup><http://pfam.xfam.org>

<sup>7</sup><https://gb.fungalgenomics.ca/portal/>

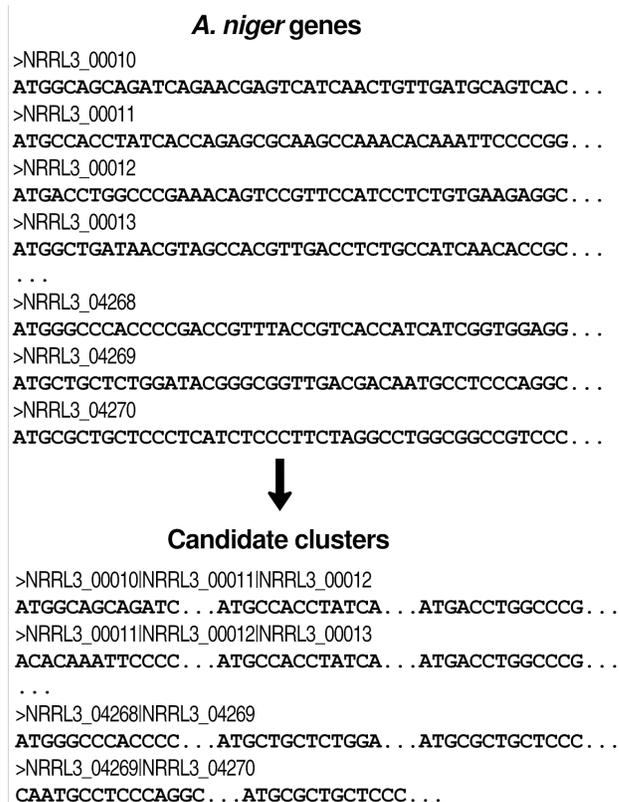


Fig. 2. Example of *A. niger* candidate clusters generated for test phase

### C. Classification Models

In this section we describe the methods applied to analyse the performance of a supervised learning approach using the fungal BGC datasets presented in Section III-A and the test data presented in Section III-B. To generate classification models with our fungal BGC datasets, we utilized the DeepBGC system [6]. DeepBGC executable, source code and other resources are openly available<sup>8</sup>. Among these resources, there are pre-built BGC classification models and word2vec-based embeddings built using Pfam IDs, referred to as `pfam2vec` embeddings. In [6] the authors explained that `pfam2vec` embeddings were trained based in a skipgram architecture with 100 dimensions and over 15,686 unique Pfam IDs. DeepBGC classification is based on a Bidirectional Long Short Term Memory (BiLSTM) neural network, for which the input are `pfam2vec` embeddings. In [6] DeepBGC hyperparameters are described as a BiLSTM layer size of 128, dropout of 0.2, sigmoid activation, batch size of 64, 256 timestamps over 328 epochs, using Adam optimizer at a learning rate of 1e-4, with weighted binary cross-entropy loss. To generate classification models using fungal BGC datasets on the DeepBGC system we adopted the same hyperparameters described in [6], as well as the `pfam2vec` embeddings as input for training. For each fungal BGC dataset, we have generated a different

<sup>8</sup><https://github.com/Merck/deepbgc>

classification model using DeepBGC. Fungal BGC models are named by their positive instance percentage:

- pos50 (50%-50%)
- pos40 (40%-60%)
- pos30 (30%-70%)
- pos20 (20%-80%)
- pos10 (10%-90%)
- pos05 (05%-95%)
- pos01 (01%-99%)

To complement our analysis, we also analysed the performance of our test datasets using the four bacteria-based models made available at the DeepBGC repository:

- deepbgc
- cf\_o (clusterfinder\_original)
- cf\_r (clusterfinder\_retrained)
- cf\_g (clusterfinder\_geneborder)

According to the models description at the DeepBGC releases page<sup>9</sup> and [6], the deepbgc model is based on the BiLSTM DeepBGC architecture and trained on a MIBiG dataset. The other models are built based on ClusterFinder [12], which is a Hidden Markov Model (HMM). cf\_o is a ClusterFinder HMM using original parameters; cf\_r is also a ClusterFinder HMM but trained on a MIBiG dataset; and cf\_g is a ClusterFinder HMM that switches stages only on gene borders, and trained on a MIBiG dataset.

#### IV. RESULTS AND DISCUSSION

We present here statistics and further details on the publicly available fungal BGC datasets proposed in this study. We also present results of validation and test phase obtained with classification models based on fungal BGC datasets and built using DeepBGC. Section IV-A has further information and statistics on the fungal BGC datasets proposed in our work. In Section IV-B we present results obtained at validation of training DeepBGC using the models pos50, pos40, pos30, pos20, pos10, pos05, and pos01. In Section IV-C we present results obtained at test phase. For the sake of comparison, we also report the results on test data using BGC classification models provided by DeepBGC and built based on bacteria data, as listed in Section III-C. All performance metrics are reported on the positive class only.

##### A. Fungal BGC datasets

The fungal BGC datasets proposed in this work are composed of positive and negative instances, as mentioned in Section III-A. These datasets are suitable for performing binary classification to predict fungal BGCs, and are made publicly available at <https://github.com/bioinfoUQAM/fungalbgcdata>. The availability of such resource can potentially motivate the development of supervised learning approaches to tackle BGC discovery in fungi.

Positive instances in our datasets represent fungal BGCs from 52 different fungal genera. The variety of fungal genus

is relevant to provide a large representation of BGC occurrence through different organisms. Additionally, the positive instances contain samples of over 10 different BGC types. Table II shows the different BGC types and a summary of fungal genera in our datasets. As the table shows, the most common BGC type is Polyketide synthase (PKS), followed by Non-ribosomal peptide synthase (NRP) and Terpene synthase (TC). The presence of different fungal genus and BGC types in the datasets are important for representing a wide variety of BGC occurrences, and therefore contribute to building more robust supervised learning approaches.

BGC types		BGC fungi genus	
	#		#
Alkaloid	3	Acremonium	1
Alkaloid/NRP	3	Alternaria	5
Alkaloid/TC	1	Armillaria	1
Alkaloid/NRP/TC	1	Aspergillus	9
NRP	41	Aureobasidium	1
NRP/PKS	19	Beauveria	1
PKS	90	Bipolaris	3
PKS/TC	5	Botrytis	1
RiPP	3	Byssochlamys	1
Saccharide	1	Cercospora	1
TC	23	Chaetomium	1
Other	10	Cladonia	2
Total	200	Claviceps	2
		Diaporthe	1
		Elsinoe	1
		Epichloe	2
		Fusarium	8
		Glarea	1
		Glycomyces	1
		Hypholoma	1
		Hypomyces	1
		Isaria	1
		Lasioidiplodia	1
		Lecanicillium	1
		Leptosphaeria	1
		Malbranchea	1
		Metacordyceps	1
		Metarhizium	1
		Monascus	3
		Mycosphaerella	1
		Myrothecium	1
		Neosartorya	1
		Neotyphodium	2
		Nodulisporium	1
		Paecilomyces	1
		Parastagonospora	1
		Penicillium	13
		Pestalotiopsis	1
		Phoma	2
		Phomopsis	1
		Purpureocillium	1
		Sarocladium	1
		Shiraia	1
		Sordaria	1
		Sphaceloma	1
		Stachybotrys	1
		Starmerella	1
		Talaromyces	3
		Tapinella	1
		Topolycladium	2
		Trichophyton	1
		Ustilago	1

TABLE II  
FUNGAL GENERA AND BGC TYPES IN POSITIVE INSTANCES OF DATASETS

Negative instances in our datasets represent synthetic gene clusters composed of fungal orthologs. By using fungal orthologs as source for the negative instances, we can generate synthetic gene clusters that depict the genomic profile of fungi. A total of 549 fungal species are present in orthologs composing our negative instances. The main fungal groups to which the orthologs belong to are shown in Table III, according to their taxonomy level. In this table we show the number of species clustered under different taxonomy levels (genus, family, order, or class), and the corresponding total of non-redundant orthologous genes for each group.

The 52 fungal genera in positive instances together with the 549 fungal species in negative instance orthologs contribute to represent the genomic diversity in fungi, and therefore support the development of more robust classification models.

##### B. Validation performance

Table IV shows validation metrics obtained with fungal BGC datasets. During training phase, all models using fungal BGC datasets had early stopping before completing the total 328 epochs. This could be a sign that the models were overfitting, a possible consequence due to the size of the

<sup>9</sup><https://github.com/Merck/deepbgc/releases>

Group	Taxonomy level	# Species	# Genes
Aspergillus	Genus	30	309,629
Cryptococcus	Genus	7	44,028
Exophiala	Genus	7	67,291
Metarhizium	Genus	5	45,563
Penicillium	Genus	21	208,580
Phytophthora	Genus	6	89,378
Hypocreaceae	Family	7	66,815
Pleosporaceae	Family	9	94, 817
Polyporaceae	Family	6	61,584
Saprotlegniaceae	Family	6	81,114
Trichocomaceae	Family	6	52,941
Agaricales	Order	25	293,149
Eurotiales	Order	60	608,401
Helotiales	Order	14	162,251
Hypocreales	Order	50	512,282
Mucorales	Order	15	164,081
Polyporales	Order	17	169,368
Sordariales	Order	8	66,549
Agaricomycetes	Class	77	912,187
Eurotiomycetes	Class	103	1,002,099
Microbotryomycetes	Class	9	59,326
Pucciniomycetes	Class	6	64,018
Saccharomycetes	Class	76	390,808
Tremellomycetes	Class	18	121,702
Ustilaginomycetes	Class	9	55,465

TABLE III

MAIN FUNGAL GROUPS PRESENT IN NEGATIVE INSTANCES OF DATASETS

datasets and the imbalanced distribution between the two classes.

The best performing model `pos50` is the one with the most balanced distribution of positive and negative instances. It yield Precision (P) of 0.598, Recall (R) of 0.995, and F-measure (F) of 0.747. Models `pos10`, `pos05`, and `pos01`, the ones with the most imbalanced distributions, had the lowest validation loss but also the lowest P, R and F.

TABLE IV

VALIDATION PERFORMANCE USING MODELS BUILT ON PROPOSED DATASETS

Model	Epochs	Loss	P	R	F
<code>pos50</code>	91	0.683	0.598	0.995	0.747
<code>pos40</code>	52	0.719	0.407	1	0.578
<code>pos30</code>	108	0.667	0.536	0.743	0.623
<code>pos20</code>	97	0.758	0.230	0.991	0.373
<code>pos10</code>	70	0.389	0	0	0
<code>pos05</code>	73	0.240	0	0	0
<code>pos01</code>	57	0.062	0	0	0

### C. Test performance

The test phase show how the models would perform in a real case scenario, when a complete genome is being processed to predict candidate BGC regions. The dataset inputted in the test phase is composed of candidate clusters from the *A. niger* genome sequence, as described in Section III-B. The performance on the test data is presented in two ways: gene metrics and cluster metrics. Gene metrics show P, R, and F for genes that belong to knownBGCs. Cluster metrics show P, R, and F for knownBGCs where a minimum of one cluster gene must be correctly classified for the cluster to be predicted as positive. Tables V and VI show the results on *A. niger* test

datasets, with overlaps of respectively 50% and 30%. These results were obtained using classification models built with the fungal BGC datasets described in Section III-A.

TABLE V

PERFORMANCE FOR *A. niger* TEST DATA USING MODELS BUILT ON FUNGAL BGC DATASETS USING 50% OVERLAP

Model	Gene metrics			Cluster metrics		
	P	R	F	P	R	F
<code>pos50</code>	0.049	1.0	0.094	0.072	0.988	0.134
<code>pos40</code>	0.048	0.962	0.091	0.073	0.988	0.136
<code>pos30</code>	0.044	0.867	0.083	0.073	0.977	0.136
<code>pos20</code>	0.039	0.694	0.074	0.079	0.93	0.146
<code>pos10</code>	0	0	0	0	0	0
<code>pos05</code>	0	0	0	0	0	0
<code>pos01</code>	0	0	0	0	0	0

TABLE VI

PERFORMANCE FOR *A. niger* TEST DATA USING MODELS BUILT ON FUNGAL BGC DATASETS USING 30% OVERLAP

Model	Gene metrics			Cluster metrics		
	P	R	F	P	R	F
<code>pos50</code>	0.05	1.0	0.096	0.1	0.988	0.182
<code>pos40</code>	0.048	0.951	0.092	0.099	0.953	0.179
<code>pos30</code>	0.045	0.865	0.085	0.1	0.942	0.18
<code>pos20</code>	0.039	0.669	0.073	0.105	0.884	0.188
<code>pos10</code>	0	0	0	0	0	0
<code>pos05</code>	0	0	0	0	0	0
<code>pos01</code>	0	0	0	0	0	0

Results in the test phase show an important decrease in performance compared to the validation phase metrics. However the behaviors observed at the validation step also appear in test. Similarly to the validation phase, the more imbalanced models `pos10`, `pos05`, `pos01` did not predict any candidate cluster as positive. This behavior happened with both test datasets of 50% or 30% overlap, and it could indicate that the model is sensitive to an imbalanced distribution of classes.

Also similarly to the validation phase the more balanced models `pos50`, `pos40`, `pos30`, `pos20` tended to predict most of candidate clusters as positives, leading to high recall but very low precision. Table VI shows slightly better performance for P, R, and F compared to table V. This behavior could indicate that using a 30% overlap in the test data is better suited for the task.

Following the results obtained with models based on fungal BGC datasets, we would like to also analyse the performance of DeepBGC models built using bacteria data on *A. niger* test datasets. Tables VII and VIII show the results obtained on *A. niger* data with respectively 50% and 30% overlap using DeepBGC bacteria models.

Among all DeepBGC bacteria models, `deepbgc` performed best at both gene and cluster metrics, either using 30% or 50% overlap, with 0.273 F. The model `cf_o` showed the lowest performance, with 0.138 F. Models `cf_r` and `cf_g` showed in both cases better performance than `cf_o`. The results using DeepBGC trained models yield a similar tendency than that of the fungal BGC models: high

TABLE VII  
PERFORMANCE FOR *A. niger* TEST DATA WITH 50% OVERLAP USING  
MODELS PROVIDED BY DEEPBGC

Model	Gene metrics			Cluster metrics		
	P	R	F	P	R	F
deepbgc	0.074	0.972	0.138	0.114	0.988	0.205
cf_o	0.05	1.0	0.096	0.074	0.988	0.138
cf_r	0.056	0.997	0.106	0.083	0.988	0.153
cf_g	0.06	0.989	0.113	0.09	0.988	0.166

TABLE VIII  
PERFORMANCE FOR *A. niger* TEST DATA WITH 30% OVERLAP USING  
MODELS PROVIDED BY DEEPBGC

Model	Gene metrics			Cluster metrics		
	P	R	F	P	R	F
deepbgc	0.074	0.954	0.138	0.159	0.988	0.273
cf_o	0.051	0.984	0.096	0.103	0.988	0.187
cf_r	0.058	0.994	0.109	0.118	0.988	0.211
cf_g	0.061	0.992	0.116	0.126	0.988	0.223

recall but very low precision.

A loss in performance between validation and test results is evident, either when using fungal BGC based models or DeepBGC bacteria models.

As mentioned in Section III-A, fungal BGCs seem to show larger genomic diversity, which possibly makes it more complex to perform BGC discovery in fungi if compared to bacteria. Therefore, performance is expected to be somewhat affected by performing fungal BGC classification using bacteria-based models.

The dataset size at training time could also have had an impact on training pos50, pos40, pos30, pos20, pos10, pos05 models, given DeepBGC classification approach. As the authors in [25] explained, the suitability of deep learning approaches varies according to the problem at hand; and in cases when available data is limited conventional approaches could be relevant and more advantageous. As discussed in Section III-A the number of known fungal BGC data previously validated by curators is rather limited, which as a consequence will limit the size of fungal BGC datasets. It is possible and worth investigating that different classification methods, apart from a BiLSTM neural network as adopted in DeepBGC, will be better suited for handling fungal BGC discovery.

## V. CONCLUSION

NPs are bioactive compounds that play a vital role in the production of a large variety of drugs, and the discovery of novel NPs can potentially benefit human health. Great effort has been put on identifying BGCs that are capable of producing NPs in plants, bacteria and fungi. Identifying BGCs is a challenging task, specially in fungi given the clusters genomic diversity.

Previous work on identifying BGCs in bacteria have resulted in a large variety of approaches and annotated data available. In fungi most previous approaches are based on data-driven

methods and present a limited scope, such as covering only certain types of BGCs, or have been developed based on a single species data. The availability of new fungal BGC datasets could potentially motivate the development of new methods to identify BGCs in fungi. One example is supervised learning, a method that have shown to perform well in bacteria data.

In this work, we present new fungal BGC datasets to leverage supervised learning in the fungal BGC discovery task. These datasets are made publicly available at <https://github.com/bioinfoUQAM/fungalbgcdata>. The availability of such fungal BGC datasets can potentially motivate the development of binary classification approaches to tackle the BGC discovery task. We have shown results obtained on these fungal BGC datasets using a supervised learning approach developed for bacteria BGCs. We also analysed the performance of bacteria-based classification models applied on a fungal genome. The test performance on both fungal-based generated models or bacteria-based models was similar given precision (low) and recall (high) metrics using the same supervised learning method. This points to an opportunity to explore different supervised learning approaches than the one adopted by the DeepBGC system, that might be more suitable to handle fungal BGC datasets.

## ACKNOWLEDGMENT

This work was supported by a fellowship of the Natural Sciences and Engineering Research Council of Canada (NSERC) to H.A., and NSERC Discovery Grant to A.B.D. and A.T.

## REFERENCES

- [1] A. K. Chaudhary, D. Dhakal, and J. K. Sohng, "An insight into the -omics based engineering of streptomycetes for secondary metabolite overproduction," *BioMed Research International*, vol. 2013, 2013.
- [2] M. H. Medema and M. A. Fischbach, "Computational approaches to natural product discovery," *Nature Chemical Biology*, vol. 11, no. 9, pp. 639–648, 2015.
- [3] A. K. Chavali and S. Y. Rhee, "Bioinformatics tools for the identification of gene clusters that biosynthesize specialized metabolites," *Briefings in Bioinformatics*, vol. 19, pp. 1022–1034, 04 2017.
- [4] A. Osbourn, "Secondary metabolic gene clusters: evolutionary toolkits for chemical innovation," *Trends in Genetics*, vol. 26, no. 10, pp. 449–457, 2010.
- [5] P. Agrawal, S. Khater, M. Gupta, N. Sain, and D. Mohanty, "RiPP-Miner: a bioinformatics resource for deciphering chemical structures of RiPPs based on prediction of cleavage and cross-links," *Nucleic Acids Research*, vol. 45, no. W1, pp. W80–W88, 2017.
- [6] G. D. Hannigan, D. Prihoda, A. Palicka, J. Soukup, O. Klempir, L. Rampula, J. Durcak, M. Wurst, J. Kotowski, D. Chang, *et al.*, "A deep learning genome-mining strategy for biosynthetic gene cluster prediction," *Nucleic Acids Research*, 08 2019.
- [7] M. H. Medema, R. Kottmann, P. Yilmaz, M. Cummings, J. B. Biggins, K. Blin, I. De Bruijn, Y. H. Chooi, J. Claesen, R. C. Coates, *et al.*, "Minimum information about a biosynthetic gene cluster," *Nature Chemical Biology*, vol. 11, no. 9, p. 625, 2015.
- [8] K. R. Conway and C. N. Boddy, "ClusterMine360: a database of microbial PKS/NRPS biosynthesis," *Nucleic Acids Research*, vol. 41, no. D1, pp. D402–D407, 2012.
- [9] K. Blin, T. Wolf, M. G. Chevrette, X. Lu, C. J. Schwalen, S. A. Kautsar, H. G. Suarez Duran, E. L. De Los Santos, H. U. Kim, M. Nave, *et al.*, "antiSMASH 4.0 improvements in chemistry prediction and gene cluster boundary identification," *Nucleic Acids Research*, vol. 45, no. W1, pp. W36–W41, 2017.

- [10] K. Blin, M. H. Medema, R. Kottmann, S. Y. Lee, and T. Weber, "The antiSMASH database, a comprehensive database of microbial secondary metabolite biosynthetic gene clusters," *Nucleic Acids Research*, p. gkw960, 2016.
- [11] M. Hadjithomas, I.-M. A. Chen, K. Chu, J. Huang, A. Ratner, K. Palaniappan, E. Andersen, V. Markowitz, N. C. Kyrpides, and N. N. Ivanova, "IMG-ABC: new features for bacterial secondary metabolite analysis and targeted biosynthetic gene cluster discovery in thousands of microbial genomes," *Nucleic Acids Research*, vol. 45, no. D1, pp. D560–D565, 2016.
- [12] P. Cimermancic, M. H. Medema, J. Claesen, K. Kurita, L. C. W. Brown, K. Mavrommatis, A. Pati, P. A. Godfrey, M. Koehrsen, J. Clardy, *et al.*, "Insights into secondary metabolism from a global analysis of prokaryotic biosynthetic gene clusters," *Cell*, vol. 158, no. 2, pp. 412–421, 2014.
- [13] I. Kjærboelling, T. C. Vesth, J. C. Frisvad, J. L. Nybo, S. Theobald, A. Kuo, P. Bowyer, Y. Matsuda, S. Mondo, E. K. Lyhne, *et al.*, "Linking secondary metabolites to gene clusters through genome sequencing of six diverse *Aspergillus* species," *Proceedings of the National Academy of Sciences*, vol. 115, no. 4, pp. E753–E761, 2018.
- [14] T. C. Vesth, J. Brandl, and M. R. Andersen, "FunGeneClusterS: predicting fungal gene clusters from genome and transcriptome data," *Synthetic and Systems Biotechnology*, vol. 1, no. 2, pp. 122–129, 2016.
- [15] M. Umemura, H. Koike, N. Nagano, T. Ishii, J. Kawano, N. Yamane, I. Kozono, K. Horimoto, K. Shin-ya, K. Asai, J. Yu, J. W. Bennett, and M. Machida, "MIDDAS-M: motif-independent de novo detection of secondary metabolite gene clusters through the integration of genome sequencing and transcriptome data," *PLOS ONE*, vol. 8, no. 12, p. e84028, 2013.
- [16] T. Wolf, V. Shelest, N. Nath, and E. Shelest, "CASSIS and SMIPS: promoter-based prediction of secondary metabolite gene clusters in eukaryotic genomes," *Bioinformatics*, vol. 32, no. 8, pp. 1138–1143, 2016.
- [17] I. Takeda, M. Umemura, H. Koike, K. Asai, and M. Machida, "Motif-independent prediction of a secondary metabolism gene cluster using comparative genomics: application to sequenced genomes of *Aspergillus* and ten other filamentous fungal species," *DNA Research*, vol. 21, no. 4, pp. 447–457, 2014.
- [18] N. Khaldi, F. T. Seifuddin, G. Turner, D. Haft, W. C. Nierman, K. H. Wolfe, and N. D. Fedorova, "SMURF: genomic mapping of fungal secondary metabolite clusters," *Fungal Genetics and Biology*, vol. 47, no. 9, pp. 736–741, 2010.
- [19] M. A. Skinnider, C. A. Dejong, P. N. Rees, C. W. Johnston, H. Li, A. L. Webster, M. A. Wyatt, and N. A. Magarvey, "Genomes to natural products prediction informatics for secondary metabolomes (PRISM)," *Nucleic Acids Research*, vol. 43, no. 20, pp. 9645–9662, 2015.
- [20] E. V. Kriventseva, D. Kuznetsov, F. Tegenfeldt, M. Manni, R. Dias, F. A. Simão, and E. M. Zdobnov, "OrthoDB v10: sampling the diversity of animal, plant, fungal, protist, bacterial and viral genomes for evolutionary and functional annotations of orthologs," *Nucleic Acids Research*, vol. 47, no. D1, pp. D807–D811, 2018.
- [21] T. C. Vesth, J. L. Nybo, S. Theobald, J. C. Frisvad, T. O. Larsen, K. F. Nielsen, J. B. Hoof, J. Brandl, A. Salamov, R. Riley, *et al.*, "Investigation of inter-and intraspecies variation through genome sequencing of *Aspergillus* section *Nigri*," *Nature Genetics*, vol. 50, no. 12, p. 1688, 2018.
- [22] S. El-Gebali, J. Mistry, A. Bateman, S. R. Eddy, A. Luciani, S. C. Potter, M. Qureshi, L. J. Richardson, G. A. Salazar, A. Smart, *et al.*, "The Pfam protein families database in 2019," *Nucleic Acids Research*, vol. 47, no. D1, pp. D427–D432, 2018.
- [23] R. P. de Vries, R. Riley, A. Wiebenga, G. Aguilar-Osorio, S. Amillis, C. A. Uchima, G. Anderluh, M. Asadollahi, M. Askin, K. Barry, *et al.*, "Comparative genomics reveals high biological diversity and specific adaptations in the industrially and medically important fungal genus *Aspergillus*," *Genome biology*, vol. 18, no. 1, p. 28, 2017.
- [24] D. O. Inglis, J. Binkley, M. S. Skrzypek, M. B. Arnaud, G. C. Cerqueira, P. Shah, F. Wymore, J. R. Wortman, and G. Sherlock, "Comprehensive annotation of secondary metabolite biosynthetic genes and gene clusters of *Aspergillus nidulans*, *A. fumigatus*, *A. niger* and *A. oryzae*," *BMC Microbiology*, vol. 13, no. 1, p. 91, 2013.
- [25] C. Angermueller, T. Pärnamaa, L. Parts, and O. Stegle, "Deep learning for computational biology," *Molecular Systems Biology*, vol. 12, no. 7, p. 878, 2016.