

SOFTWARE

Open Access



# Hitac: a hierarchical taxonomic classifier for fungal ITS sequences compatible with QIIME2

Fábio M. Miranda<sup>1,2</sup>, Vasco C. Azevedo<sup>3</sup>, Rommel J. Ramos<sup>3,4,5</sup>, Bernhard Y. Renard<sup>1†</sup> and Vitor C. Piro<sup>1,2\*†</sup>

<sup>†</sup>Bernhard Y. Renard and Vitor C. Piro shared senior co-authorship.

\*Correspondence:  
Vitor.Piro@fu-berlin.de

<sup>1</sup> Data Analytics and Computational Statistics, Hasso Plattner Institute, Digital Engineering Faculty, University of Potsdam, Potsdam, Germany

<sup>2</sup> Department of Mathematics and Computer Science, Freie Universität Berlin, Berlin, Germany

<sup>3</sup> Institute of Biological Sciences, Universidade Federal de Minas Gerais, Belo Horizonte, Brazil

<sup>4</sup> Institute of Biological Sciences, Federal University of Pará, Belém, Brazil

<sup>5</sup> Centro de Computação de Alto Desempenho, Universidade Federal do Pará, Belém, Brazil

## Abstract

**Background:** Fungi play a key role in several important ecological functions, ranging from organic matter decomposition to symbiotic associations with plants. Moreover, fungi naturally inhabit the human body and can be beneficial when administered as probiotics. In mycology, the internal transcribed spacer (ITS) region was adopted as the universal marker for classifying fungi. Hence, an accurate and robust method for ITS classification is not only desired for the purpose of better diversity estimation, but it can also help us gain a deeper insight into the dynamics of environmental communities and ultimately comprehend whether the abundance of certain species correlate with health and disease. Although many methods have been proposed for taxonomic classification, to the best of our knowledge, none of them fully explore the taxonomic tree hierarchy when building their models. This in turn, leads to lower generalization power and higher risk of committing classification errors.

**Results:** Here we introduce HiTaC, a robust hierarchical machine learning model for accurate ITS classification, which requires a small amount of data for training and can handle imbalanced datasets. HiTaC was thoroughly evaluated with the established TAXXI benchmark and could correctly classify fungal ITS sequences of varying lengths and a range of identity differences between the training and test data. HiTaC outperforms state-of-the-art methods when trained over noisy data, consistently achieving higher F1-score and sensitivity across different taxonomic ranks, improving sensitivity by 6.9 percentage points over top methods in the most noisy dataset available on TAXXI.

**Conclusions:** HiTaC is publicly available at the Python package index, BIOCONDA and Docker Hub. It is released under the new BSD license, allowing free use in academia and industry. Source code and documentation, which includes installation and usage instructions, are available at <https://gitlab.com/dacs-hpi/hitac>.

**Keywords:** Python, Taxonomy, ITS, Local Hierarchical Classification



## Background

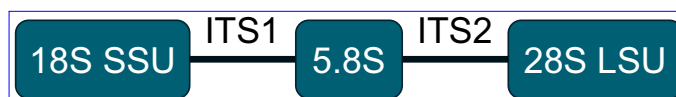
In the last decades, there has been a surge of interest in characterizing the mycoflora of communities. This is due to the fact that fungi are key organisms in several important ecological functions, ranging from nutrient recycling to mycorrhizal associations and decomposition of wood and litter [1]. Moreover, fungi naturally inhabit the human body and may benefit the host when administered as probiotics [2], yet little is known about the fungal microbiota of healthy individuals when compared to the bacterial microbiome [3]. Furthermore, identifying the presence of fungi can be challenging, as they rarely form structures visible to the naked eye, and similar structures are frequently composed of several distinct species [4]. Hence, it became common practice to use DNA sequencing in addition to morphological studies in contemporary mycology [5].

One of the approaches that can be employed in environmental studies is whole-genome shotgun sequencing (WGS), where the purpose is to obtain the genetic material of all organisms present in a given microbiome simultaneously [6]. Nonetheless, in studies where the end goal is only to estimate the diversity or discover the evolutionary distance of organisms in a microbiome, it is cheaper and more convenient to sequence only genetic markers such as 16 S rRNA, which is well conserved in all bacteria [7] and also widely used to classify and identify archaea [8]. In mycology, the internal transcribed spacer (ITS) region was adopted as the universal marker for classifying fungi [9].

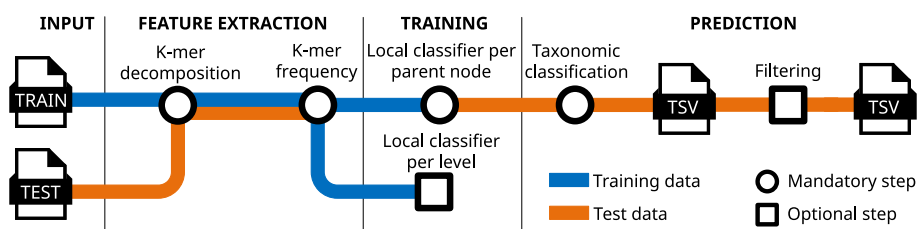
The main characteristic of 16 S rRNA that makes it a good genetic marker is the presence of nine hypervariable regions V1–V9, which are flanked by conserved regions that can be used to amplify target sequences using universal primers [10]. However, 16 S rRNA has fewer hypervariable domains in fungi, and consequently, it is more appropriate to use the ITS region (Fig. 1), which is flanked by the 18 S and 28 S rRNA genes and interrupted by the conserved 5.8S rRNA gene [11]. Additionally, the ease with which ITS is amplified makes it an appealing target for sequencing samples from environments as challenging as soil and wood, where the initial quantity and quality of DNA is low [4].

From a computational perspective, the taxonomic classification of ITS fungal sequences can be performed by using either similarity-based or machine learning techniques. Similarity-based approaches, such as BLAST [12], require the alignment of a query sequence against all sequences available in a reference database, e.g., Silva [13] and UNITE [14]. However, a major limitation of similarity-based methods is the dependence on homologs in databases, which are not always available, while machine learning methods can learn only the relevant features and criteria for classification and can, therefore, be applied to any new sequence [15].

One of the first machine learning software proposed to perform taxonomic classification was the RDP-Classifier, which uses 8-mer frequencies to train a Naive Bayes classification algorithm [16]. Improving upon the RDP-Classifier, a novel Naive Bayes classifier



**Fig. 1** Visual representation of the ITS region. The ITS region consists of the ITS1 and ITS2 regions, which are flanked by the 18 S and 28 S rRNA genes and separated by the conserved 5.8S rRNA gene [11]



**Fig. 2** Sketch of the key algorithmic components of HiTaC

with multinomial models was developed to increase the accuracy [17]. This multinomial classifier was later reimplemented and optimized on Microclass [18].

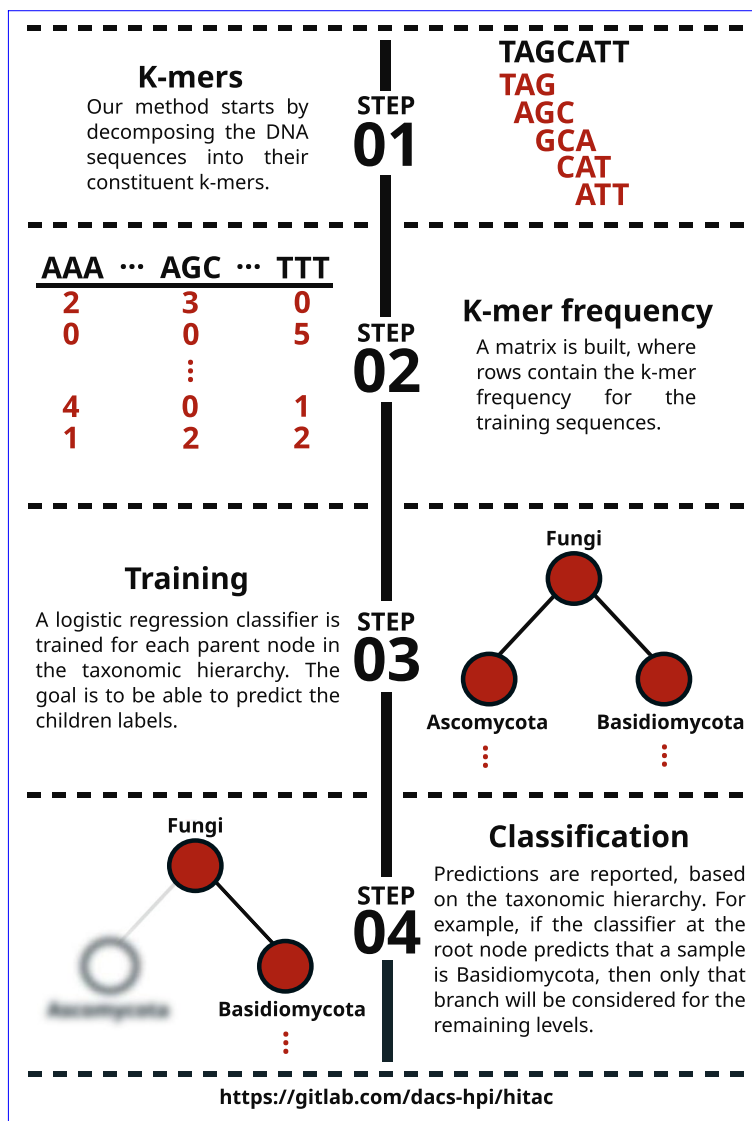
Other machine learning approaches have also been proposed in the literature, such as the k-Nearest Neighbor (KNN) algorithm [19]. Given a query sequence, the KNN algorithm identifies the k-most similar sequences in a database and uses the taxonomic information from each of those sequences to determine the consensus taxonomy. Q2\_SK [20] implements several supervised learning classifiers from scikit-learn [21], e.g., Random Forest, Support Vector Machines (SVM) and Gradient Boosting. Q2\_SK is flexible and allows the selection of features and classifiers.

Despite the advances made in taxonomic classification, there are still opportunities for enhancement due to the complexity of the data. Furthermore, to the best of our knowledge, the machine learning software available may lack in predictive performance since they perform flat classification, which means that they are only trained on leaf node labels and do not fully explore the hierarchical structure of the taxonomic tree when building their models (see Fig. S1 for a depiction of the hierarchical structure of the taxonomic tree). The main disadvantage of this approach is that the information about parent–child class relationships in the hierarchical structure is not entirely considered, which may cause flat classifiers to commit more errors than hierarchical approaches [22]. Therefore, we developed a new hierarchical taxonomic classifier for fungal ITS sequences, denominated HiTaC, which explores the hierarchical structure of the taxonomic tree during the training stage and improves the F1-score and sensitivity of fungal ITS predictions. Furthermore, HiTaC can be easily installed and is compatible with QIIME2, facilitating its adoption in existing ITS sequence analysis pipelines.

## Implementation

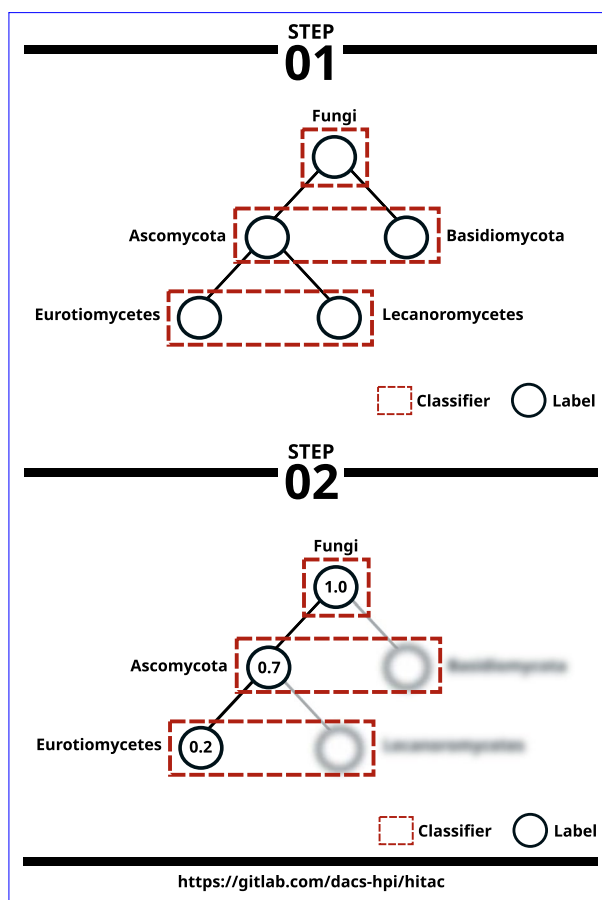
### The local hierarchical taxonomic classifier

HiTaC performs feature extraction, building of the hierarchical model, taxonomic classification, and reporting of the predictions. Figure 2 illustrates the algorithmic steps of HiTaC, where it starts by decomposing the DNA sequences from both reference and query files into their constituents k-mers. Next, two k-mer frequency matrices are built for the reference and query sequences, respectively, where each row contains the k-mer frequency for a sequence. Logistic regression classifiers implemented in the library scikit-learn [21] are trained for each parent node in the taxonomic hierarchy, using the local classifier per parent node implementation from a library called HiClass that we developed and published previously [23]. The k-mer frequency matrix and taxonomic labels of the reference sequences are used as training



**Fig. 3** HiTaC algorithm overview. To compute the k-mer frequency of DNA sequences, we first decompose them into their constituent k-mers. In the example shown in step 01,  $k = 3$ , resulting in 5 sub-strings. Afterwards, HiTaC computes a matrix containing the k-mer frequency for all training sequences, where each row contains the k-mer frequency for a given sequence and the columns are all possible k-mers for a given  $k$ , i.e., all  $k$  nucleotide permutations of A, C, G and T. Next, HiTaC trains a logistic regression classifier for each parent node in the taxonomic tree, where the features are the k-mer frequencies computed in step 2 and the labels are the taxonomic annotations of the training sequences. In step 03, only the first 2 levels are shown for brevity. Lastly, HiTaC makes predictions in a top-down approach, based on the taxonomic hierarchy

data to predict children labels (see Fig. 3 for a visual representation). The mycological nomenclature used as labels for the internal nodes originates from the taxonomy provided by the user, which allows the user to decide whether to use Index Fungorum [24], Faces of Fungi [25], Mycobank [26], or The NCBI taxonomy database [27], for example. Predictions are performed using the k-mer frequency matrix computed for the query sequences and are reported in a top-down approach, based on the taxonomic hierarchy. For example, if the classifier at the root node predicts that a query



**Fig. 4** HiTaC filter overview. To build the filter, we first train a logistic regression classifier for each level of the hierarchy, using the same k-mer frequency algorithm introduced previously on Fig. 3. In the example shown in step 01, there are only 3 levels in the hierarchy for simplicity. Afterwards, HiTaC computes the confidence score for all taxonomic ranks assigned to a sample by the local classifier per parent node, then on a bottom-up approach it removes labels that are below a certain threshold. In step 02, supposing that the threshold is 0.7 and that the given sample has a confidence score of 0.2 assigned for the lowest rank, then the final prediction for this sample will only contain the labels Fungi and Ascomycota

sequence belongs to the phylum Ascomycota, then only that branch of the taxonomic tree will be considered for the remaining taxonomic levels.

#### Classification uncertainty

To reduce classification errors when faced with uncertainties, we also construct a filter by training a logistic regression classifier for each level of the hierarchy. This filter computes confidence scores for all taxonomic ranks assigned by the local classifier per parent node and excludes labels that are below a user-defined threshold in a bottom-up approach (Fig. 4). The confidence scores are computed using the `predict_proba` method implemented in scikit-learn’s logistic regression classifier, which in turn uses a one-vs-rest approach; that is, it calculates the probability of each class using the logistic function and assuming it to be positive, then it normalizes these values across all the classes [28].

### Code quality assurance

The code base adheres to the PEP 8 code style [29], which is enforced by `flake8` and the uncompromising code formatter `black` to ensure high code quality. Versioning complies with `SemVer` to increase reproducibility and facilitate dependency management by end users. The code is accompanied by unit tests that cover 98% of the code base and are automatically executed by our *continuous integration workflow* upon commits.

### Installation and usage

HiTaC is hosted on GitLab<sup>1</sup>, where there is also documentation with installation and usage instructions. HiTaC is compatible with Python 3.8+ and can be installed on GNU/Linux, Windows and macOS. It can be easily obtained with `pip install hitac`, `conda install -c bioconda hitac` or `docker pull mirand863/hitac`. Listing 1 shows a basic example of how to fit a hierarchical model and predict the taxonomy using HiTaC. More elaborate examples can be found in the documentation.

```
hitac-fit \  
  --reference reference.fasta \  
  --kmer 6 \  
  --threads 8 \  
  --classifier classifier.pickle  
  
hitac-classify \  
  --reads query.fasta \  
  --classifier classifier.pickle \  
  --kmer 6 \  
  --threads 8 \  
  --classification predictions.tsv
```

**Listing 1** Example of how to use HiTaC to train a hierarchical classifier and predict taxonomy

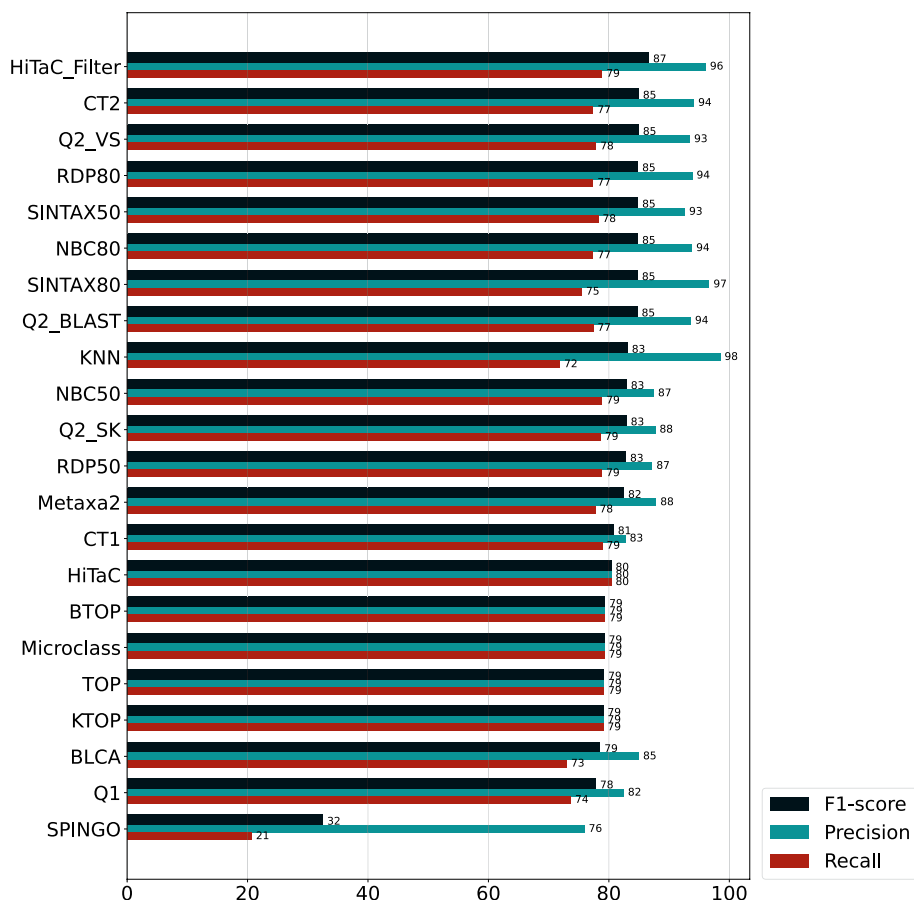
### Evaluation

We evaluated HiTaC in five datasets with varying sequence lengths (Tables S2–S3) and compared its performance against seventeen taxonomic classifiers using the TAXXI benchmark [30]. All methods evaluated were trained and tested with the same data in order to avoid bias. The commands and parameters executed are summarized in Tables S4–S22.

To measure the performance of our model, we adopted all metrics available in the TAXXI benchmark and extended them by adding standard machine learning metrics. Furthermore, we also computed hierarchical metrics, which can give better insights into which algorithm is better at classifying hierarchical data [22]. The definitions of all these metrics are listed in the supporting information text.

The datasets in the TAXXI benchmark were created using real environmental data, publicly available on NCBI [31], RDP [32], the Warcup fungal ITS training set v2 (WITS) [33], UNITE [34] and other *in vivo* samples [30]. These datasets were created with a

<sup>1</sup> <https://gitlab.com/dacs-hpi/hitac>.

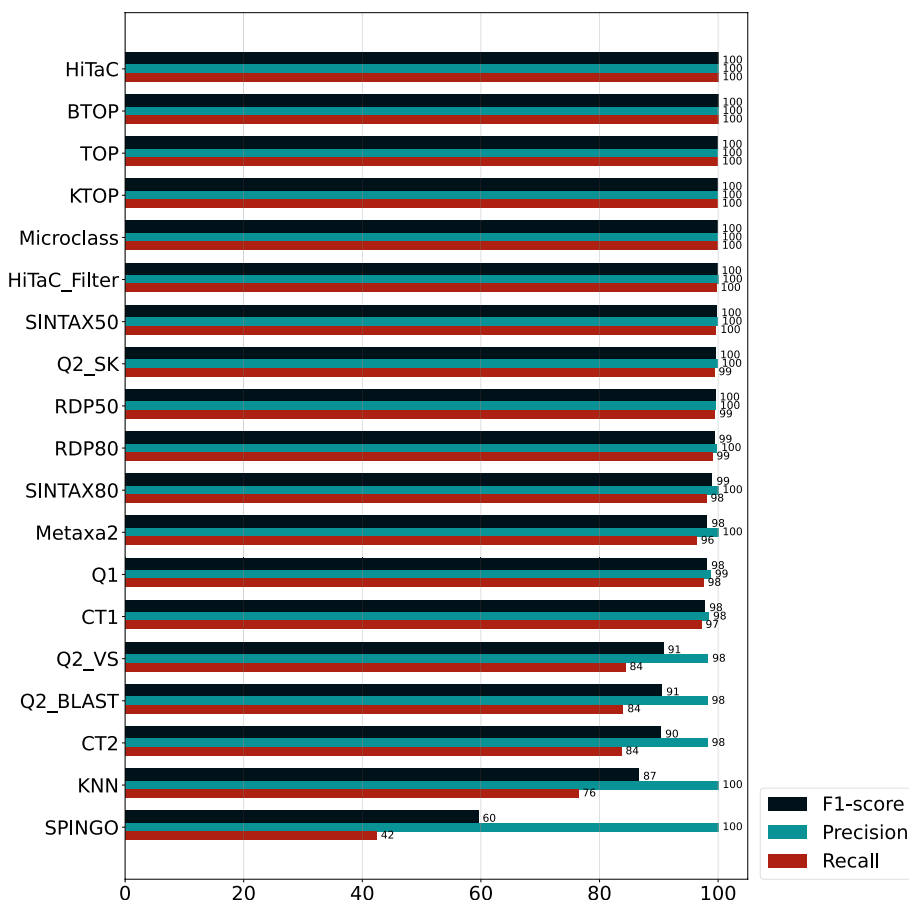


**Fig. 5** Hierarchical f1-score, precision and recall computed for all taxonomic ranks for the dataset with 90% identity, sorted by f1-score. HiTaC\_Filter was one of the best-performing methods when dealing with uncertainty, achieving high f1-score, precision and recall, while the filter-less version obtained the highest recall of all methods

strategy known as cross-validation by identity, where varying distances between query sequences and the reference database are accounted for. In practice, this means that in the dataset with 90% identity all species are novel and there is a mix of novel and known taxa at the genus level, making it the most difficult dataset in the benchmark. As the identity increases, so does the amount of known labels at all taxonomic ranks; hence, the dataset with 100% identity can be considered the easiest in the benchmark.

**Results and discussion**

First, we evaluated HiTaC on the dataset with 90% identity to assess its behavior under uncertainty. In order to do this, we computed the hierarchical precision to measure the proportion of correct predictions, the hierarchical recall to measure the proportion of detected positive samples and the hierarchical F1-score, which is the harmonic mean of the precision and recall. As shown in Fig. 5, HiTaC\_Filter was one of the top-performing methods in terms of precision, scoring 96.05 while maintaining a recall of 78.86, equivalent to other high-ranking methods. Although other software obtained similar or slightly better precision, it cost them a steeper decline in recall, which is also evidenced



**Fig. 6** Hierarchical f1-score, precision and recall computed for all taxonomic ranks for the dataset with 100% identity. Both HiTaC and HiTaC\_Filter achieved almost perfect scores when all organisms were known

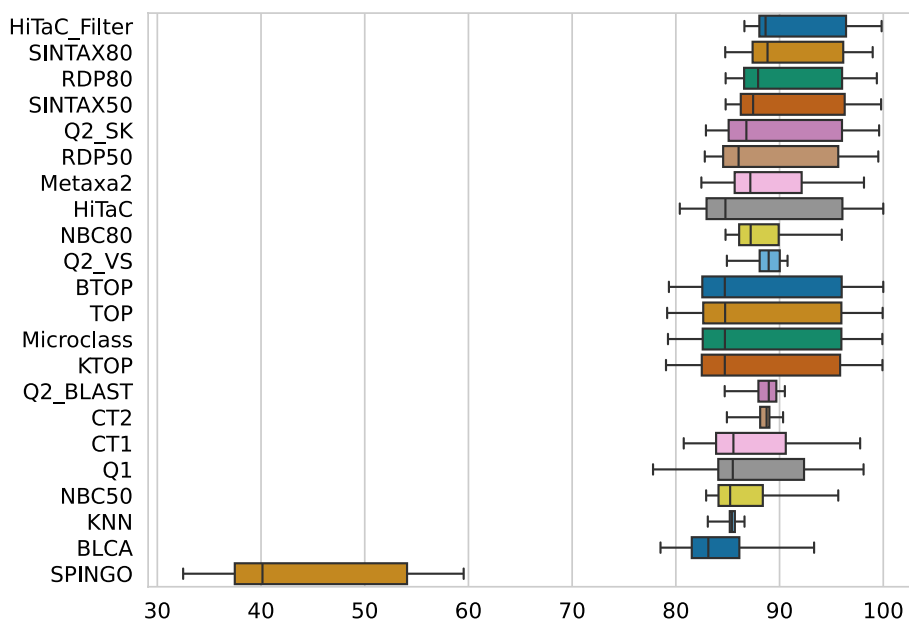
by the lower F1 scores. This result shows that HiTaC performed well in the presence of uncertainty.

To appraise HiTaC’s behavior when classifying known organisms, we evaluated it on the dataset with 100% identity. As shown in Fig. 6, HiTaC\_Filter achieved a perfect precision score of 100 when the taxonomy was fully annotated in the reference database. Furthermore, the recall persisted at 99.68; that is, there was no prominent decrease in recall when compared with the filter-less version. This result shows the potential of HiTaC in accurately classifying known organisms, given that query sequences have perfect matches in the reference database.

In Fig. 7, we summarize the results achieved for the hierarchical metrics by computing the F1-score for all five datasets and sorting by the average of F1 scores. HiTaC\_Filter obtained an F1-score equivalent or higher to the top-performing approaches for all datasets, corroborating the results presented in the previous paragraphs. This result demonstrates that HiTaC\_Filter can positively identify organisms annotated in databases while being able to leave most new species unannotated.

Similar conclusions can be drawn from Fig. 8, in which we evaluated the percentage of known sequences correctly predicted, i.e., the true positive rate (TPR) for all methods and datasets available in the benchmark. The trend is that HiTaC achieved a TPR higher



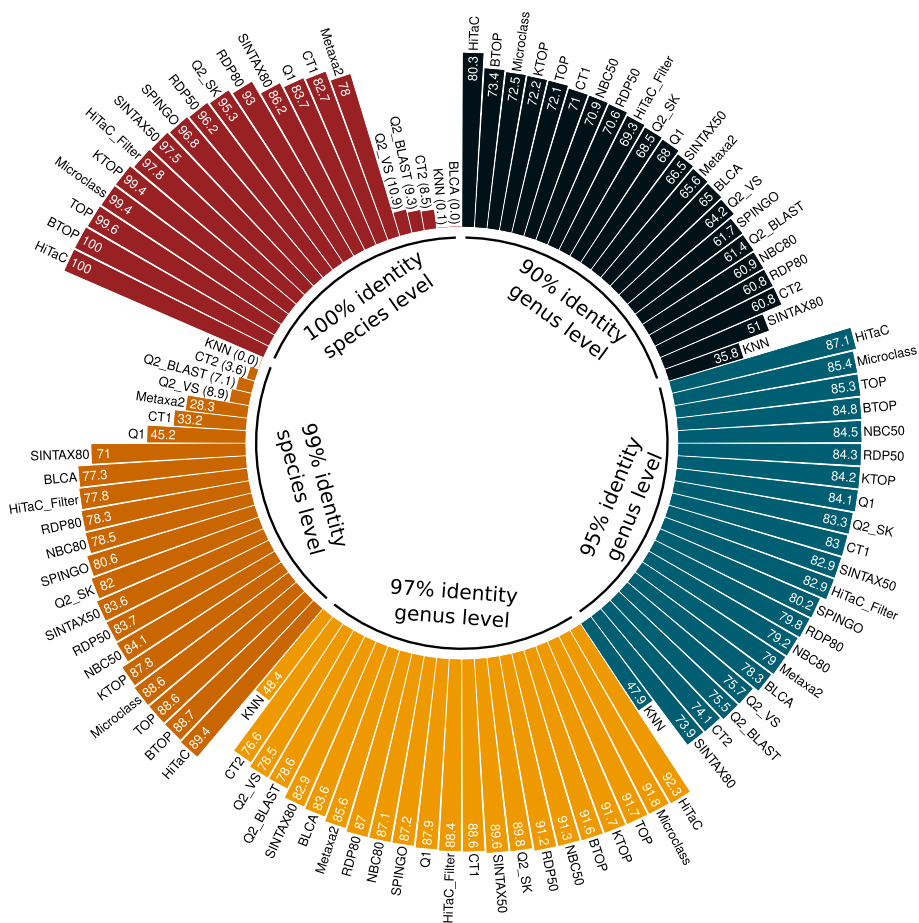


**Fig. 7** Hierarchical F1-score computed for all taxonomic ranks for the datasets with 90, 95, 97, 99 and 100% identity, sorted by average. HiTaC\_Filter achieved F1-scores equal to or above the best methods available in the literature

or equal to top methods. For instance, the best true positive rate at the genus level for the dataset with 90% identity was achieved by HiTaC, which sharply increased the TPR by 6.9 percentage points when compared with BTOP. Regarding the dataset with 95% identity at the genus level, HiTaC obtained a 1.7 percentage points improvement over Microclass, while for the dataset with 97% identity at the genus level, HiTaC surpassed Microclass with a 0.5 increase in percentage points. At the species level, HiTaC achieved a 0.7 percentage points gain over BTOP for the dataset with 99% identity, while for the dataset with 100% identity, HiTaC tied with BTOP, obtaining a perfect score of 100% TPR. These results suggest that the local hierarchical classification approach implemented in HiTaC commits fewer errors than flat classifiers as long as the sequences are previously known.

By default, HiTaC uses 6-mer frequency as a feature extraction method, but the user can change this parameter. Nevertheless, increasing the k-mer size provides small gains in predictive performance and comes with higher computational costs. The k-mer frequency computation method discards substrings with ambiguous nucleotides since they do not often occur in short reads, and accounting for them would require a slower algorithm. For example, a V means that there is either an A, C, or G in that position, and computing all these possibilities would make the method slower. Hence, we opted to disregard ambiguous nucleotides to keep the k-mer counting process as fast as possible. In future releases, we intend to implement a flexible user interface to allow users to use third-party feature extraction methods. Furthermore, we believe allowing for third-party feature extraction methods could also enable accurate taxonomic classification for 16 S rRNA.

Training the filter is slower than most methods evaluated in the benchmark (Tables S23–S27). However, it must only be performed once for the reference database,



**Fig. 8** True positive rates for all methods for the datasets with 90%, 95% and 97% identity at the genus level, and 99% and 100% identity at the species level. The trend is that HiTaC achieved a sensitivity higher or equal to top methods. For instance, HiTaC tied with the best method for the dataset with 100% identity, obtaining a perfect score. Moreover, for the datasets with 90–99% identity, HiTaC improved the true positive rates upon top methods

and after training, the filter can be used to estimate the uncertainty quickly. Moreover, some public databases, such as UNITE and SILVA, are not updated frequently. Nevertheless, we provide models pre-trained on the public database UNITE<sup>2</sup>, speeding up the process for users. Some of these pre-trained models contain all eukaryotic ITS sequences available on UNITE, which enable detection and removal of nonfungal sequences mistakenly amplified by polymerase chain reaction (PCR) [35]. At the time of writing, the current version of UNITE uses the mycological nomenclature published in Fungal Diversity [36]; hence, this is the nomenclature HiTaC learns when trained with the UNITE database and is reported to the user with the pre-trained models. Furthermore, HiTaC uses the unique species hypotheses identifiers provided in the database UNITE as the last taxonomic level in the hierarchy during training and reports them to the user, which increases taxonomic reproducibility.

<sup>2</sup> <https://gitlab.com/dacs-hpi/hitac#pre-trained-models>

On average, the entire ITS region is approximately 600 base pairs (bp) in length [37]. This trend can also be observed in the UNITE database, with averages ranging from 578 bp to 651 bp. However, in the UNITE database, the sequence lengths can vary from 250 bp to 16,761 bp, while the ITS sequences in the TAXXI benchmark have a maximum length of 1373 bp.

In summary, HiTaC is a Python package optimized to produce accurate taxonomic classification of fungal ITS sequences. Due to the extensive use of the taxonomic hierarchy during training, HiTaC produces fewer errors than other methods compared in the benchmark. Furthermore, its filter is a reliable classifier under uncertainty. HiTaC provides standalone Python scripts and a QIIME 2 plugin that can be quickly adopted into existing analysis pipelines. HiTaC and its dependencies can be easily installed via `pip`, `conda` or `docker`, which is not the case for most of the taxonomic classifiers available in the literature. HiTaC is released under the new BSD license, allowing free use in academia and industry and free copy, modification, and redistribution of the code as long as a duplicate of the original license is kept and proper acknowledgments are given.

## Conclusion

HiTaC is a Python package for the taxonomic classification of fungal ITS sequences, which applies local hierarchical classifiers for accurate predictions. It provides scripts to train classifiers on specialized datasets and to classify new sequences with the trained models. The standard version has high sensitivity and is ideal for exploratory analysis. HiTaC also implements a filter that can express uncertainties in classifications, indicating if the input sequences are complex to recognize. Thanks to its compatibility with QIIME 2, users can quickly adopt it in existing mycobiome analysis pipelines. HiTaC is an open-source software available at <https://gitlab.com/dacs-hpi/hitac>.

## Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s12859-024-05839-x>.

Supporting information text, Figures S1–S4 and tables S1–S92

## Acknowledgements

FMM gratefully acknowledges Marcos Augusto dos Santos (UFMG) for initial inspiration, Robert C. Edgar (independent researcher) for providing the scripts and datasets in the TAXXI benchmark, Sven Giese (HPI) for code improvement suggestions and Elizabeth Yuu (HPI) for proofreading.

## Author contributions

FMM designed and implemented HiTaC's algorithm, performed analysis, benchmarked the experiments and wrote the manuscript. VAA and RJR assisted with analysis and revised the manuscript. BYR assisted with analysis, visualizations and reviewed the manuscript. VCP oversaw experimental and analysis aspects of the project and reviewed the manuscript. All authors edited and approved the final manuscript.

## Funding

Open Access funding enabled and organized by Projekt DEAL. VCP was financed by the Deutsche Forschungsgemeinschaft (DFG) under grant number 458163427.

## Availability of data materials

HiTaC is freely available on the Python package index, BIOCONDA and docker hub. It is easier to install it by executing `pip install hitac` in the terminal. All datasets and evaluation scripts used in this paper are publicly available in the TAXXI benchmark [30], and a reproducible evaluation pipeline is available on GitLab <https://gitlab.com/dacs-hpi/hitac/-/tree/main/benchmark>.

## Declarations

### Ethics approval and consent to participate

Not applicable.

### Consent for publication

Not applicable.

### Competing interest

The authors declare that they have no competing interest.

Received: 15 March 2024 Accepted: 11 June 2024

Published online: 02 July 2024

## References

- Hawksworth DL. The magnitude of fungal diversity: the 1.5 million species estimate revisited. *Mycol Res.* 2001;105(12):1422–32.
- Banik A, Halder SK, Ghosh C, Mondal KC. Fungal probiotics: opportunity, challenge, and prospects. In: *Recent advancement in white biotechnology through fungi: volume 2: perspective for value-added products and environments*; 2019. pp. 101–117.
- Huffnagle GB, Noverr MC. The emerging world of the fungal microbiome. *Trends Microbiol.* 2013;21(7):334–41.
- Nilsson RH, Ryberg M, Abarenkov K, Sjökvist E, Kristiansson E. The ITS region as a target for characterization of fungal communities using emerging sequencing technologies. *FEMS Microbiol Lett.* 2009;296(1):97–101.
- Hibbett DS. After the gold rush, or before the flood? Evolutionary morphology of mushroom-forming fungi (agari-comycetes) in the early 21st century. *Mycol Res.* 2007;111(9):1001–18.
- Venter JC, Remington K, Heidelberg JF, Halpern AL, Rusch D, Eisen JA, Wu D, Paulsen I, Nelson KE, Nelson W. Environmental genome shotgun sequencing of the Sargasso sea. *Science.* 2004;304(5667):66–74.
- Fuhrman JA. Metagenomics and its connection to microbial community organization. *F1000 Biol Rep.* 2012;4:12.
- Kim M, Chun J. 16S rRNA gene-based identification of bacteria and archaea using the EzTaxon server. In: *Methods in microbiology*, vol 41; Elsevier. pp. 61–74.
- Schoch C.L., Seifert, K.A., Huhndorf, S., Robert, V., Spouge, J.L., Levesque, C.A., Chen, W., Consortium, F.B., Nuclear ribosomal internal transcribed spacer (ITS) region as a universal DNA barcode marker for fungi. *Proc Natl Acad Sci.* 2012;109(16):6241–6.
- Chakravorty S, Helb D, Burday M, Connell N, Alland D. A detailed analysis of 16S ribosomal rRNA gene segments for the diagnosis of pathogenic bacteria. *J Microbiol Methods.* 2007;69(2):330–9.
- Dos Reis JBA, Lorenzi AS, Vale HMM. Methods used for the study of endophytic fungi: a review on methodologies and challenges, and associated tips. *Arch Microbiol.* 2022;204(11):675.
- Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ. Gapped blast and psi-blast: a new generation of protein database search programs. *Nucleic Acids Res.* 1997;25(17):3389–402.
- Pruesse E, Quast C, Knittel K, Fuchs BM, Ludwig W, Peplies J, Glöckner FO. SILVA: a comprehensive online resource for quality checked and aligned ribosomal RNA sequence data compatible with ARB. *Nucleic Acids Res.* 2007;35(21):7188–96.
- Nilsson RH, Larsson K-H, Taylor AFS, Bengtsson-Palme J, Jeppesen TS, Schigel D, Kennedy P, Picard K, Glöckner FO, Tedersoo L. The UNITE database for molecular identification of fungi: handling dark taxa and parallel taxonomic classifications. *Nucleic Acids Res.* 2018;47(D1):259–64.
- Bzhalava Z, Tampuu A, Bała P, Vicente R, Dillner J. Machine learning for detection of viral sequences in human metagenomic datasets. *BMC Bioinform.* 2018;19:1–11.
- Wang Q, Garrity GM, Tiedje JM, Cole JR. Naive Bayesian classifier for rapid assignment of rRNA sequences into the new bacterial taxonomy. *Appl Environ Microbiol.* 2007;73(16):5261–7.
- Liu K-L, Wong T-T. Naive Bayesian classifiers with multinomial models for rRNA taxonomic assignment. *IEEE/ACM Trans Comput Biol Bioinform.* 2013;10(5):1–1.
- Liland KH, Vinje H, Snipen L. microclass: an R-package for 16S taxonomy classification. *BMC Bioinform.* 2017;18(1):172.
- Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, Lesniewski RA, Oakley BB, Parks DH, Robinson CJ. Introducing mothur: open-source, platform-independent, community-supported software for describing and comparing microbial communities. *Appl Environ Microbiol.* 2009;75(23):7537–41.
- Bokulich NA, Dillon MR, Bolyen E, Kaehler BD, Huttley GA, Caporaso JG. q2-sample-classifier: machine-learning tools for microbiome classification and regression. *J Open Res Softw.* 2018;3:30.
- Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V. Scikit-learn: machine learning in python. *J Mach Learn Res.* 2011;12:2825–30.
- Silla CN, Freitas AA. A survey of hierarchical classification across different application domains. *Data Min Knowl Disc.* 2011;22(1–2):31–72.
- Miranda FM, Köhnecke N, Renard BY. Hiclass: a python library for local hierarchical classification compatible with scikit-learn. *J Mach Learn Res.* 2023;24(29):1–17.
- Index Fungorum. <https://www.indexfungorum.org/>. Accessed 13 May 2024.
- Jayasiri SC, Hyde KD, Ariyawansa HA, Bhat J, Buyck B, Cai L, Dai Y-C, Abd-El Salam KA, Ertz D, Hidayat I. The faces of fungi database: fungal names linked with morphology, phylogeny and human impacts. *Fungal Divers.* 2015;74:3–18.
- Robert V, Vu D, Amor ABH, Wiele N, Brouwer C, Jabas B, Szoke S, Dridi A, Triki M, Daoud SB. Mycobank gearing up for new horizons. *IMA Fungus.* 2013;4:371–9.
- Federhen S. The NCBI taxonomy database. *Nucleic Acids Res.* 2012;40(D1):136–43.

28. Scikit-learn: logistic regression probability estimates. 2023. [https://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.LogisticRegression.html](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html). Accessed 26 Oct 2023.
29. Rossum G, Warsaw B, Coghlan N. PEP 8-style guide for Python code. python.org 2001.
30. Edgar RC. Accuracy of taxonomy prediction for 16s rRNA and fungal ITS sequences. *PeerJ*. 2018;6:4652.
31. Sayers EW, Barrett T, Benson DA, Bolton E, Bryant SH, Canese K, Chetverin V, Church DM, DiCuccio M, Federhen S. Database resources of the national center for biotechnology information. *Nucleic Acids Res*. 2012;40(D1):13–25.
32. Cole JR, Wang Q, Fish JA, Chai B, McGarrell DM, Sun Y, Brown CT, Porras-Alfaro A, Kuske CR, Tiedje JM. Ribosomal database project: data and tools for high throughput rRNA analysis. *Nucleic Acids Res*. 2014;42(D1):633–42.
33. Deshpande V, Wang Q, Greenfield P, Charleston M, Porras-Alfaro A, Kuske CR, Cole JR, Midgley DJ, Tran-Dinh N. Fungal identification using a bayesian classifier and the warcup training set of internal transcribed spacer sequences. *Mycologia*. 2016;108(1):1–5.
34. Kõljalg U, Larsson K-H, Abarenkov K, Nilsson RH, Alexander IJ, Eberhardt U, Erland S, Høiland K, Kjøller R, Larsson E. Unite: a database providing web-based methods for the molecular identification of ectomycorrhizal fungi. *New Phytol*. 2005;166(3):1063–8.
35. Rawson C, Zahn G. Inclusion of database outgroups reduces false positives in fungal metabarcoding taxonomic assignments. *Mycologia*. 2023;8:1–7.
36. Tedersoo L, Sánchez-Ramírez S, Kõljalg U, Bahram M, Döring M, Schigel D, May T, Ryberg M, Abarenkov K. High-level classification of the fungi and a tool for evolutionary ecological analyses. *Fungal Divers*. 2018;90:135–59.
37. Porras-Alfaro A, Liu K-L, Kuske CR, Xie G. From genus to phylum: large-subunit and internal transcribed spacer rRNA operon regions show similar classification accuracies influenced by database composition. *Appl Environ Microbiol*. 2014;80(3):829–40. <https://doi.org/10.1128/AEM.02894-13>.

### **Publisher's Note**

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.