

ARTICLE

Artificial intelligence in microbial natural product drug discovery: current and emerging role

Vinodh J Sahayasheela,^a Manendra B Lankadasari,^b Vipin Mohan Dan,^c Syed G Dastager,^{*d}

Ganesh N Pandian,^{*e} Hiroshi Sugiyama^{*a,e}

Received 00th January 20xx,
Accepted 00th January 20xx

DOI: 10.1039/x0xx00000x

Microorganisms are exceptional sources of a wide array of unique natural products and play a significant role in drug discovery. During the golden era, several life-saving antibiotics and anticancer agents were isolated from microbes; moreover, they are still widely used. However, difficulties in the isolation methods and repeated discoveries of the same molecules have caused a setback in the past. Artificial intelligence (AI) has had a profound impact on various research fields, and its application allows the effective performance of data analyses and predictions. With the advances in omics, it is possible to obtain a wealth of information for the identification, isolation, and prediction of the targets of secondary metabolites. In this review, we discuss drug discovery based on natural products from microorganisms with the help of AI and machine learning.

Introduction

Microorganisms are well known to produce structurally diverse secondary metabolites that are widely used in clinical settings for treating various clinical conditions, such as cancer, infectious disease, and inflammation.¹ Conversely, they are also used in various other sectors, such as agriculture (as herbicides and insecticides), the food sector (as nutraceuticals), enzyme inhibitors, and for bioremediation, which uses natural products (NPs) directly or develops molecules derived from their scaffolds.^{2,3} Compared with synthetic molecules, NPs offer specific features in terms of structural complexity and scaffold diversity.⁴ The discovery of NPs has also revealed previously unknown targets in cells. For instance, rapamycin, which was isolated from a strain of *Streptomyces hygroscopicus*, has resulted in the identification of the mechanistic target of the rapamycin (mTOR) cell signaling pathway.⁵

Artificial intelligence (AI) uses computers to perform complex functions, analyse large datasets, and interpret them based on algorithms.⁶ AI has been used widely in various research fields and industries for decision-making and processing tasks because it provides efficient analysis and faster results with reduced human error and at times uncovers data

structures difficult to obtain from other sources.⁷ Recently, AI has received increased attention and is being used by chemists to perform various tasks in drug discovery, as well as to identify molecular properties, process automation, plan synthetic routes, and predict the bioactivity of molecules.^{8–10} Based on the recent prolific growth in machine learning (ML) and the wealth of information in cloud computing in the form of databases and repositories, researchers can now gain access to big data and integrate AI/ML approaches into their tasks.

Despite the unparalleled role of NPs in drug discovery, this approach has various challenges, such as the isolation, screening, purification, and structural characterization of the NPs derived from microbial sources.¹¹ However, in the past two decades, the repetitive identification of existing and already known NPs, the demand for resources, and the time-consuming nature of the tasks have curbed interest in NPs among researchers and industries.¹² With the advancement of genomics, proteomics, metabolomics, and other omics technologies recently, it is now possible to obtain a wealth of information to identify the biosynthetic dark matter.^{13,14} AI/ML in the field of NPs has been growing, to analyse the extensive amount of data stemming from the omics techniques (Figure 1) and open the microbial Pandora's box for the discovery of bioactive molecules.

This review features the existing and emerging AI- and ML-based tools in various stages of the investigation of NPs from microorganisms. (Figure 2) We will highlight the techniques available to identify the microbes and prioritize them based on their genome and metabolite potentials. Subsequently, we will discuss fast dereplication, which is one of the major challenges in NP discovery, together with the tools available for this type

^a Department of Chemistry, Graduate School of Science, Kyoto University, Kitashirakawa-Oiwakecho, Sakyo-Ku, Kyoto 606-8502, Japan.

^b Thoracic Service, Department of Surgery, Memorial Sloan Kettering Cancer Center, New York, NY 10065, USA

^c Microbiology Division, Jawaharlal Nehru Tropical Botanic Garden and Research Institute, Thiruvananthapuram, Kerala, India

^d NCIM Resource Centre, Division of Biochemical Sciences, CSIR - National Chemical Laboratory, Pune, Maharashtra, India

^e Institute for Integrated Cell-Material Sciences (WPI-iCeMS), Kyoto University, Yoshida-Ushinomaecho, Sakyo-Ku, Kyoto 606-8501, Japan

of analysis. Furthermore, we will address the expedited elucidation of the structure of compounds and the identification of their targets with the aid of AI/ML. Finally, we will discuss the development of new powerful tools and the integration of multiple techniques that will speed up NP discovery, thus leading to a boom in the identification of potent drug candidates in the future.

2. Application of AI/ML in natural product discovery

2.1 Selection of organism and Taxonomic Identification.

The selection of organisms is the preliminary step in NP discovery. Certain species, such as actinomycetes, have been among the most prolific sources of pharmaceutical candidates in the past.¹² However, the overmining of this limited resource has led to the repeated rediscovery of known compounds and has exhausted the identification of novel molecules in this setting.¹⁵ Although the isolation of NPs is very laborious and challenging, careful selection of underexplored microorganisms¹⁶ from untapped environments, such as marine sources¹⁷ and symbiotic sponges,¹⁸ increases the chance of identifying molecules with different scaffolds. In addition to cultured microorganisms, nearly 99% of microbial species are uncultured in the lab and hold promise in the search for new NPs. This has led to the identification of potent antibiotics, such as teixobactin¹⁹ and lassomycin,²⁰ using specialized culture techniques.

The classical approaches in bacterial identification according to taxonomy are time-consuming and misleading; however, with the advent of the omics and ML techniques, it is possible to predict microbes efficiently.²¹ Although Gram staining is the gold-standard technique for the initial classification of bacteria, it is a highly time-intensive and manual-dependent activity. In contrast, using convolutional neural networks (CNNs), researchers were able to classify different shapes of Gram-positive and Gram-negative bacteria via imaging with high confidence.²² This technique can be further extended to various microorganisms, for their identification and classification using ML tools. DNA-based identification is the most accurate method of classification of various microorganisms, as in the identification of DNA from bacteria, which can also be distinguished based on the specialized metabolites they produce. In the past, the ability to correlate microbial identity with signature metabolites was limited, even with access to the vast amount of data generated by mass spectrometry. However, recently, researchers developed a technique termed IDBac with the help of ML to classify microbes based on their proteins and specialized metabolites using matrix-assisted laser desorption/ionization time-of-flight mass spectrometry (MALDI-TOF MS).²³ Using this approach, those authors could discriminate *Bacillus subtilis* at the strain level based on its ability to produce cyclic peptide antibiotics and a group of *Micromonospora* with 99% sequence similarity with high confidence. Another algorithm called SPeDE also facilitates the identification of microbes at taxonomic resolution from a mass spectral dataset of both culture-dependent and -independent samples.²⁴ MALDI-TOF is a powerful tool that is known for its versatility and is used in various fields with the advantage of being relatively easy to operate, fast, and accurate. The high-

throughput capacity of MALDI-TOF combined with ML tools allows the rapid identification of microbial communities compared with traditional biochemical or molecular biology techniques.²⁵ Hence, in the future, rare and underexplored microbes can be identified directly from samples with the help of ML-assisted MALDI-TOF, which will accelerate the process of candidate selection for NP screening and isolation. Another interesting application of MALDI is imaging MS (IMS), which has been used to map the spatial distribution of various secondary metabolites.^{26–29}

2.2 Genome mining with the aid of AI/ML

The use of genome mining for secondary metabolite identification has been rapidly increasing in recent years with the advent of next-generation sequencing techniques, followed by bioinformatics pipelines.³⁰ Although NPs are highly diverse in structure, their biosynthetic machinery, which is known as biosynthetic gene clusters (BGCs), is highly conserved in the microbes that fall under the class of polyketide synthases (PKSs),³¹ nonribosomally synthesized peptides (NRPs),³² ribosomally synthesized and post-translationally modified peptides, alkaloids,³³ and terpenes.³⁴ The technique begins with the identification of existing and novel BGCs from the genome sequences and further characterization of novel gene clusters, to complete the analysis. To perform this type of complex analysis using big data, ML algorithms are widely used to predict the BGC assembly lines and predict the putative encoded structure from the sequence.³⁵ With the help of BGC databases^{36–41} and computational tools,^{42–49} NPs can be predicted based on previously characterized pathways (Table 1). Using one such tool, antiSMASH, which employs profile hidden Markov models (pHMMs) to identify the BGC, a novel polyketide named formicamycin (Figure 3) has been isolated.⁵⁰ In another study, a potent antituberculous compound, gladiolin (Figure 3), was isolated with the help of genome mining from *Burkholderia gladioli*, which is a previously unknown source of NPs, in a patient with cystic fibrosis.⁵² More recently, a new class of previously unknown cryptic BGCs, i.e., lanthipeptides,⁵¹ was identified with the help of ML and deep learning (DL) strategies.

Conventionally, the process of NP isolation uses a “grind and find” approach, which involves culturing the microorganism followed by purification and structure elucidation; however, with the advent of genome mining and ML/DL-based approaches, novel metabolites have been isolated from uncultured organisms.⁵² For instance, the combination of the two strategies has led to the discovery of the antibiotic malacidin from the global microbiome using heterologous expression without culturing the organism.⁵³ A computational algorithm based on hidden Markov models (HMMs) is available for BGC identification from metagenomic samples, which allows the identification of interesting molecules from the human microbiome.^{54,55} In many cases, most of the BGCs remain silent, without expression, which hinders the production of secondary metabolites; nevertheless, using elicitors (e.g., small molecules and coculture), it is possible to predict the biosynthetic genes and express them with the help of ML tools.⁵⁶ One of the major

obstacles to NP discovery is the identification of secondary metabolites from unconventional sources because of the lack of

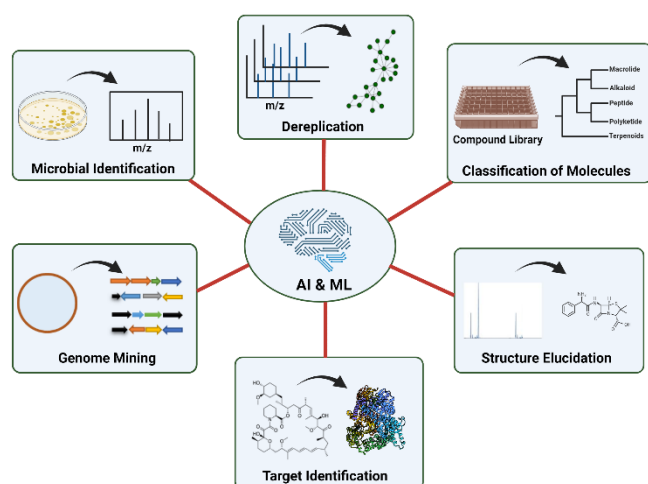


Figure 1: Application of AI/ML to various areas of microbial natural product drug discovery.

cultivation of the microbes. However, with the emergence of metagenomics and ML, it is now possible to predict NPs in environmental or biological niches using specialized ML tools.^{54,57}

2.3 AI/ML tools for Metabolite production and expression

Many microorganisms, such as those in the genera *Streptomyces* and *Myxococcus*, have been predicted to have large secondary metabolite BGCs with the advent of genome sequencing and bioinformatics. However, they usually do not code for NPs and remain as silent gene clusters.⁵⁸ Therefore, various genome engineering techniques have been applied to activate those silent gene clusters, such as cocultivation,⁵⁹ one strain many compounds,⁶⁰ elicitors,⁶¹ ribosome engineering,⁶² chemical epigenetics,⁶³ epigenetic modification,⁶⁴ overexpression of transcription factors,⁶⁵ and heterologous expression,⁶⁶ which have had huge success in identifying new compounds. Despite the success in the control of parameters such as growth and strain engineering, media optimization remains challenging.⁶⁷ To overcome this hurdle, various AI/ML techniques have been developed to control and monitor the production of metabolites. A study reported by Neythen et al. has used deep reinforcement learning, an approach from AI, for the control of cocultures in a continuous bioreactor.⁶⁸ Using this approach, those authors were able to optimize the output of the coculture bioprocess by controlling various parameters. This type of study can be considered for controlling various factors in the production of NPs. Another study reported by Fei et al. used a high-throughput method to activate the silent BGCs in various organisms.⁶⁹ The authors screened elicitors to induce secondary metabolite production with the help of IMS in nearly 500 conditions. Using this approach, they identified a new glycopeptide from *Amycolatopsis keratiniphila*, NRRL B24117, with the help of laser-ablation-coupled electrospray ionization

MS. Although this approach can perform HTS to overcome the drawback of IMS and to analyze complex datasets, Brett et al. have developed a work tool for Metabolomics Explorer (MetEx) that enables users quickly and intuitively to analyze complex liquid chromatography (LC)-MS and metabolomics datasets.⁷⁰

2.4 Dereplication of NPs with AI/ML techniques

During the golden era of NP development, several drug candidates were identified, most of which are still widely used for treating various diseases and infections.⁷¹ However, in the late 20th century, NP discovery started slowing down because of the repeated isolation of known compounds.⁷² To overcome this issue, fast identification of the known secondary metabolites is necessary, to reduce the analytical time and resources.⁷³ Dereplication is a key process in the quick identification of previously known compounds in microbial extracts.⁷⁴ Microbial extracts contain various compounds; therefore, the use of dereplication techniques helps eliminate redundancy and provides knowledge regarding novel compounds. To perform this highly efficient and robust task, ML tools with high accuracy are required. Previously, the dereplication techniques were carried out using high-performance liquid chromatography connected with a UV or photodiode array (PDA) detector with an automated bioassay and inbuilt library databases.⁷⁵ However, structural information is lacking when using UV/PDA-based detection, and a more powerful instrument is required to capture additional spectral properties of the compounds.

2.4.1 Mass spectrometry-based dereplication using AI/ML

MS is a technique that has been widely used recently for dereplication in NPs because of its sensitivity, accuracy, and rapidity. Another major advantage of MS is its ability to gain a large amount of structural information from a trace amount of sample using an untargeted approach.¹⁴ The combination of mass information with UV/PDA can readily identify compounds with the help of databases such as Dictionary of Natural Products⁷⁶ (<http://dnp.chemnetbase.com/intro/>), MarInLit⁷⁷ (<https://marinlit.rsc.org/>), StreptomeDB⁷⁸ (<http://www.pharmbioinf.uni-freiburg.de/streptomedb>), NPedia⁷⁹ (<http://www.cbgr.riken.jp/npedia/>), and The Natural Products Atlas⁸⁰ (<https://www.npatlas.org/>). Using this approach, secondary metabolites from various actinomycetes have been dereplicated.⁸¹ LC coupled with MS can achieve high-throughput screening of metabolites; however, the analysis of the data in an efficient way remains challenging. Moreover, this requires researchers manually to search various datasets, such as UV signatures, mass spectra, and microorganisms in different databases, which are far from comprehensive.¹⁴ ML-based approaches could be a good solution for the in-line identification of NPs using spectral information without manual support against the available databases.

Although MS plays an important role in the identification and dereplication of NPs, it has several drawbacks and major problems arise regarding the overlapping parent molecular masses of various metabolites based on MS spectra alone.^{82,83} Therefore, a more efficient MS-based dereplication technique, such as tandem MS, is required and can increase the sensitivity

of the detection of compounds based on MS/MS fragmentation.⁸⁴ However, the analysis of MS/MS data is a

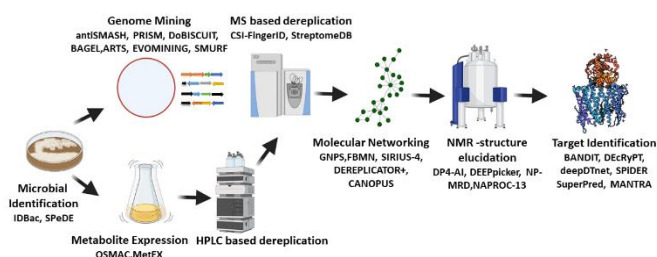


Figure 2: Various stages of natural product drug discovery with the corresponding available AI/ML tools.

cumbersome and intensive manual task, and an automated untargeted metabolomics pipeline is thus warranted to identify the metabolites efficiently. Recently, using various ML tools and algorithms, it was possible to interpret high-resolution mass spectra with reduced noise.⁸⁵ Several AI/ML-based tool has been developed for mass spectral data processing and analysis such as MZmine⁸⁶, Metaboanlyst⁸⁷, MS-Dial⁸⁸, Decon2LS⁸⁹, XCMS⁹⁰, THRASH⁹¹ and some are available as part of commercial vendor packages such as XCalibur (Thermo Fisher), MassHunter (Agilent), and using those metabolites has been predicted with high confidence manually.⁹² Metabolomics databases that are available based on MS/MS patterns are Massbank⁹³, Metlin⁹⁴, LMSD⁹⁵, MoNA (<https://mona.fiehnlab.ucdavis.edu/>), Massbank (<https://massbank.eu/MassBank/>) and GNPS⁹⁶. But in terms of microbial NPs identification, these are not widely used due to the scarcity of spectral data of natural products with the exception of GNPS⁹⁶.

Recently, molecular networking (MN) has received widespread attention in the NP community for the dereplication and delineation of novel secondary metabolites from various sources with minimal manual interference. This approach was first introduced in 2012 for metabolite analysis from a set of living microbial colonies,⁹⁷ yielding results that were comparable to the DNA sequencing of environmental samples to study microbial communities.⁹⁸ MN is a computational technique that interprets the complex dataset that arises from MS analysis and visualizes it in the form of a network.⁹⁹ To enable the analysis of MN, a crowdsourced library of reference spectra from a large number of compounds has been deposited from various communities and is available for analysis in GNPS⁹⁶ (Global Natural Products Social Molecular Networking (<http://gnps.ucsd.edu>)). MN can identify compounds based on MS/MS spectral similarities and can also link the unknown molecules with related ones by exploiting similar fragmentation patterns. MN has been recognized for its high success rate and is becoming a routine tool for dereplication. For example, using MN and indexing 260 strains with ecologically diverse origins, the *Pseudomonas*-specialized metabolome led to the discovery of poeamide B and bananamides (Figure 3).¹⁰⁰ In another study, two novel chlorinated metabolites, isoconulothiazole B and

conulothiazole C, were isolated from cyanobacteria using the MN strategy.¹⁰¹

Moreover, based on MN, further developments have been made to render the road toward the identification of NPs more straightforward. Using classical MN, various features have been incorporated with MS/MS, and feature-based MN (FBMN) has been introduced.¹⁰² It can efficiently distinguish isomers based on chromatography and ion mobility, while also facilitating spectral annotations and quantifications, thereby enabling robust analyses. Further, during ionization molecules form different adduct which limits the library annotation in MN to overcome this bottleneck Ion Identity Molecular Networking (IIMN) was developed.¹⁰³ This feature improved the network connectivity for structurally related molecule and can be used to reveal unknown ion-ligand complexes. Very recently to identify bioactive compounds a scalable native metabolomics approach integrating non-targeted liquid chromatography tandem mass spectrometry, and simultaneous detection of protein binding via native mass spectrometry was developed.¹⁰⁴ Using this integrated technique, rivulariapeptolides a family of serine protease inhibitors with nanomolar potency was identified and such approach could be central importance for drug discovery in future.

Hosein et al. have developed DEREPLICATOR+, an algorithm that can aid in the identification of NP classes such as NRPs, polyketides, terpenes, benzenoids, alkaloids, and flavonoids.¹⁰⁵ A common problem in NP identification is the isolation of active compounds during bioassay-guided purification from the extract. To overcome this hurdle, bioactivity-based MN, which integrates bioinformatics workflow to map the bioactive score using MN, was developed.¹⁰⁶ Using this approach, antiviral compounds were isolated from extracts of *Euphorbia dendroides*, for which a classical bioassay-guided fractionation procedure had previously failed.¹⁰⁶ Further, a versatile, open-access platform NP Analyst was developed as a user friendly web-based infrastructure enabling NP community to analyze without the need for intense data processing.¹⁰⁷ Although in the past MN could only be done via the web with GNPS, now many off-line tools such as MZmine3.0⁸⁶, MS-DIAL⁸⁸, Metaboseek¹⁰⁸, NetID¹⁰⁹ and commercial software like Compound Discoverer (Thermo Scientific) have the ability to perform MN without the online platform making it easier.

Although mass spectral analytical tools are available for the identification of known compounds from databases, predicting the structure of unknown metabolites is a very challenging task. However, with the advent of ML, it is improving fast. Bocker et al. developed a tool (SIRIUS 4) that can identify the structure based on MS/MS datasets using a support vector machine.¹¹⁰ Further, advancing SIRIUS 4, ZODIAC, a network-based algorithm for the de novo annotation of database-independent molecular formulas was developed by the same group.¹¹¹ Employing Bayesian statistics and Gibbs sampling it ensures fast processing in practice and is found to be better than SIRIUS by 16.5 fold. Using such ML tools novel molecular formula can be annotated. In another study that used a Deep Neural Network (DNN), a computational tool (class assignment and ontology prediction using MS, CANOPUS) was developed that could

predict unknown metabolites for which spectral and structural reference data were not available.¹¹² Similar to CANOPUS, a

Task	Tool	Features	Ref
Microbial Identification with AI/ML tools			
MALDI-TOF analysis	IDBac	Bioinformatics pipeline that integrates both intact protein and metabolite for detection	23
	SpeDE	Identification based on unique features instead of global similarity	24
Genome Mining AI/ML tools			
BGC databases	antiSMASH database	Popular and comprehensive resource on secondary metabolite BGC	36
	DoBISCUIT	Curated and literature-based collection of PKS and NRPS biosynthetic gene clusters	37
	IMG-ABC	Largest database of curated BGC from microbial genomes and metagenomes	38
	MIBiG	Collection of large curated BGC	39
	ClusterMine360	Curated database of BGCs including produced compound(s), taxonomic information	40
	Bactibase	Integrated open-access database of bacterial antimicrobial peptides/bacteriocins	41
BGC Identification from Genomes BGC databases	antiSMASH	Most widely used tool for BGC detection based on profile Hidden Markov Models (pHMMs)	42
	PRISM	BGC identification along with cheminformatic dereplication and biological activity	43
	BAGEL	Mining tool for ribosomally synthesized and post-translationally modified peptides (RIPPs)	44
	ARTS	Prioritization of the most promising BGCs encoding antibiotics with novel modes of action	45
	EvoMining	Identify secondary metabolite biosynthetic gene clusters (BGCs) based on phylogenomics	46
	SMURF	HMM-based BGC identification tool from fungi	47
	MIPS-CG	Identify completely novel BGCs using genome data in fungus alone	48
	DeepBGC	Deep learning genome-mining strategy for BGC cluster prediction	49
BGC identification from Metagenome	MetaBGC	A read-based algorithm for the detection of BGCs directly in metagenomic sequencing data	54
	eSNaPDA	Surveying and Mining BGCs from Metagenomes also take into account metadata	57
Metabolite production and expression			
Elicitor screening	MetEx	UPLC-MS-guided high-throughput elicitor screening	70
Natural product dereplication and structure elucidation with help of AI/ML			
Databases	DNP	Structure database containing over 226,000 NPs with physical and chemical properties	76
	MarinLit	Database of the marine natural products (Not open access)	77
	StreptomeDB	Database of NP isolated from streptomycetes with chemical and biological information	78
	NPEdia	Database for Natural Products	79
	NPATlas	Online database of microbial-derived natural products with structures and features	80
MS based dereplication/Identification	GNPS	Online repository for untargeted MS/MS data with sample information	96
	FBMN	Incorporates isotope patterns and retention time along with MN	102
	DEREPLICATOR+	Molecular Network combined with dereplication workflow	105
	Bioactive-MN	MN guided bioassay-guided fractionation of bioactive compound(s)	106
	SIRIUS-4	Molecular structure identification from MS/MS	110
	CANOPUS	Predict structure exclusively for which neither spectral nor reference data are available	112
	MetGem	Molecular Networks Based on the t-SNE Algorithm	113
	MESSAR	Automated prediction of metabolite substructures from tandem mass spectra	114
	Moldiscovery	Molecule identification by probabilistic model with their mass spectra	115
	FALCON	Density-based clustering of MS/MS spectra for unsupervised structure prediction	116
	SIMILE	Significant Interrelation of MS/MS Ions via Laplacian Embedding to predict the structural relationships of compounds	117
	MolNetEnhancer	Enhanced Molecular Networks by Integrating Metabolome Mining and Annotation Tools	118
	COSMIC	High-confidence structural annotation of metabolites absent from spectral libraries	119

NMR based structure elucidation/dereplication	NAPROC-13	A database for dereplication of NPs in mixtures based on C13 NMR	120
	NP-MRD	A huge structural and NMR database of nearly 41,000 NP	121
	DP4-AI	Automated NMR data analysis for structure prediction	122
	DEEP picker	Deconvolution of complex two-dimensional NMR spectra based on DNN	123
	MixONat	Software for the Dereplication of Mixtures Based on 13C NMR Spectroscopy	124
	ELINA	1H NMR based identification of bioactive compounds in a mixture prior to purification	125
	SMART-Miner	A convolutional neural network-based metabolite identification from NMR spectra	126
	SMART 2.0	NMR-based machine learning tool for annotation of molecularly diverse Natural Products	127
Integrated approach	NPClassifier	A Deep Neural Network-Based Structural Classification Tool for Natural Products	128
	GNP	Identify polyketides and NRP using genome and LC-MS/MS	129
	NRPminer	NRP identification by Integrating genomics and metabolomics dataset	130
	NRPquest	Integrates genomics and metabolomics for NRP discovery	131
	DEREP-NP	Database for dereplication from MS and NMR Experiments	132
ML-based target identification	deepDTnet	Target identification by deep learning from heterogeneous networks	133
	BANDIT	Bayesian ML approach that integrates multiple data types to predict drug binding targets	134
	SPiDER	ML tool using self-organizing maps built from various features for target prediction	135
	DEcRyPT	Machine Intelligence workflow-based target prediction	136
	SuperPred	Drug classification and target prediction using 2D, Fragment, and 3D similarity	137
	MANTRA2.0	Mechanism of action prediction using transcriptional profiles.	138
	Openchem	A Deep Learning Toolkit for Computational Chemistry and Drug Design	139
	DeepTox	Toxicity Prediction using Deep Learning	140

Table 1. List of the AI/ML tools available for various phases of natural product identification and drug leads

high-confidence structural annotation tool COSMIC based on SVM was developed.¹¹⁹ MS2DeepScore, which is an ML-supported mass spectral similarity-predicting algorithm was developed that allowed clustering, to identify metabolites similar to GNPS.^{96,141} Further, FALCON¹¹⁶ a density-based clustering of MS/MS spectra¹¹⁶, MS2LDA combined with Mass2Motif¹⁴² an unsupervised substructure discovery platform¹⁴³ and Significant Interrelation of MS/MS Ions via Laplacian Embedding (SIMILE)¹¹⁷ are also available to predict the structural relationships of compounds. MN-based approaches for dereplication can be carried out with high success and can be further employed for the structural elucidation of novel compounds in the future with the support of the ML approaches developed recently.^{112–115,119}

2.4.2 AI for the NMR-based structure elucidation/dereplication of NPs

The structural elucidation of molecules is a challenging problem in NP research. Although X-ray crystallography provides unambiguous structural information, it is often impeded by the requirement of a single crystal, and the limited amount of the

isolated molecule restricts its wide application.¹⁴⁴ Nuclear magnetic resonance (NMR) is a universally employed spectroscopy method that allows NP chemists to deduce molecular structures from spectra.¹⁴⁵ Computer-aided structural elucidation (CASE) still plays a marginal role in this setting, although it was one of the earlier applications of AI.¹⁴⁶ Although databases for NMR are available (NAPROC-13,¹²⁰ CH-NMR-NP (<https://www.j-resonance.com/en/nmrdb/>), BMRB,¹⁴⁷ and Spektraris NMR),¹⁴⁸ they have several drawbacks and, thus, do not truly satisfy the requirements of NP communities.¹⁴⁹ To overcome this issue, NP-MRD,¹²¹ which is an NMR database including over 41,000 NP compounds from >7400 different living species with various features, was introduced very recently.¹²¹ This database is still under development; however in the future, it will allow automated dereplication and CASE to be performed much more efficiently.

To assist the structure elucidation and perform dereplication, ML tools and software, such as logic for structure elucidation,¹⁵⁰ ACD/Structure elucidator,¹⁵¹ Mestrelab

ARTICLE

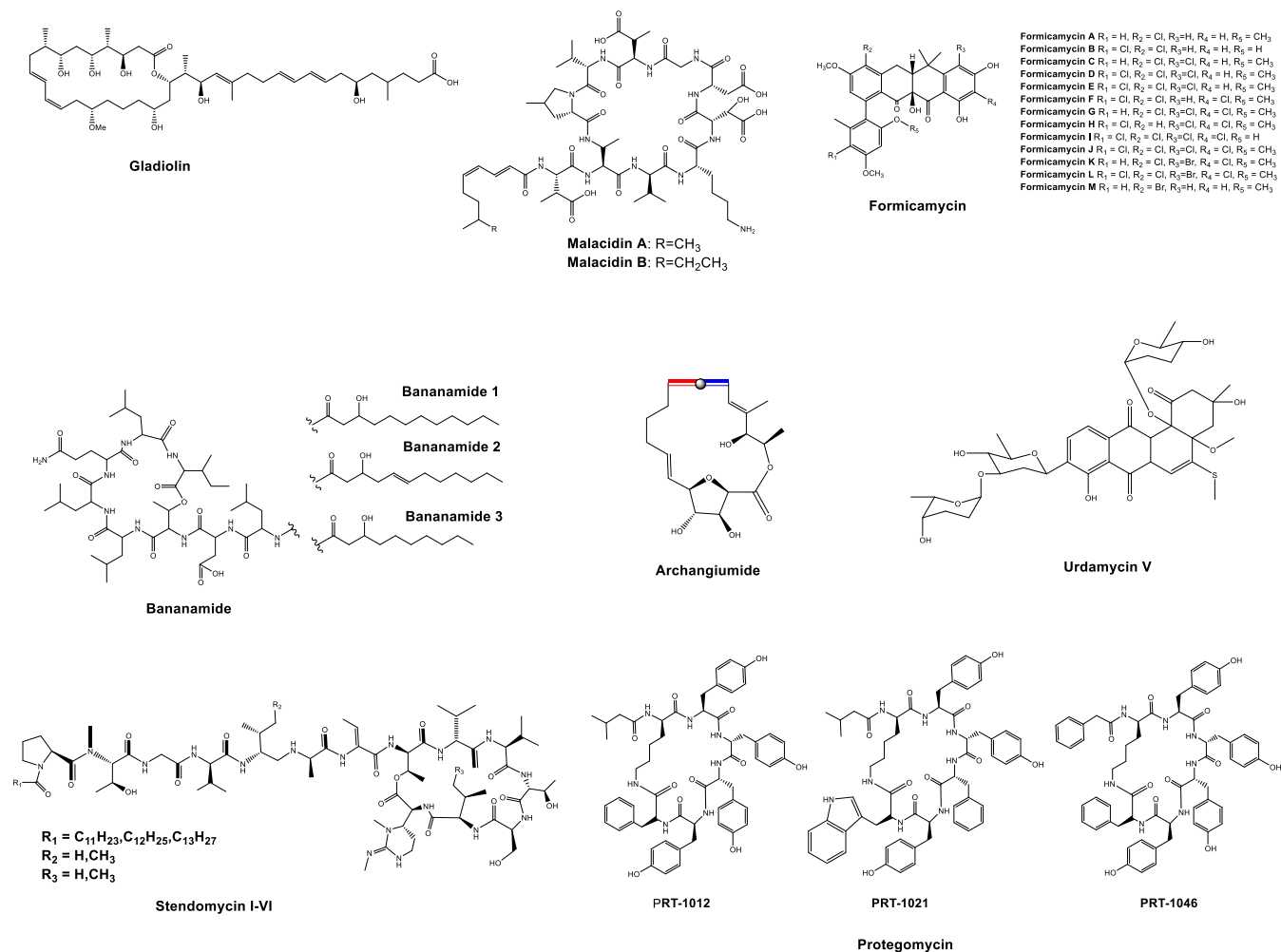


Figure 3: Novel natural products predicted with the support of AI/ML tools.

Mnova,¹⁵² and Computer-aided Spectral Assignment,¹⁵³ were developed and have aided NP identification.^{154,155} Recently, a robust AI-powered structure-prediction tool (DP4-AI)¹²² was developed and allowed the successful assignment of the structure of complex NPs.^{156–158} Using the CNN-based approach NMR-based machine learning tool “Small Molecule Accurate Recognition Technology” (SMART 2.0) for mixture analysis and characterization of new natural products were developed.¹²⁷ This led to the identification of a new chimeric swinholide-like macrolide, symplocolide A, as well as the annotation of swinholide A, samholides A–I, and several new derivatives. In another study, SMART-Miner a metabolite identification tool from the ¹H-¹³C HSQC NMR spectra with the support of CNN was developed. The model was trained on 657 chemical entities collected from HMDB and BMRB to

subsequently identify those molecules in complex mixtures with an accuracy of 88%.

To analyze the two-dimensional NMR spectra, a DNN-based approach for peak picking and spectral deconvolution (DEEP Picker) was developed very recently.¹²³ In another study, various classes of NPs were predicted using ML from ¹³C-NMR spectroscopic data.¹⁵⁹ NMR is relatively less explored for dereplication compared with HRMS because of its sensitivity; nevertheless, it can offer high accuracy in terms of the prediction of stereoisomers and the detection of all organic compounds in a mixture.¹⁶⁰ Recently, using ¹³C-NMR, a dereplication software (MixONat)¹²⁴ was developed that allowed the distinction of structurally close NPs, including stereoisomers, and aided the identification of xanthenes in *Calophyllum brasiliense*.^{124,161} In another study that used ¹H-

NMR, Grienke et al. developed a workflow ELINA (Eliciting Nature's Activities) based on a heterocovariance analysis, which can detect chemical features that correlate with bioactivity before isolation; using this approach, the authors discovered lanostane triterpenes from the extract of the fungus *Fomitopsis pinicola*.¹²⁵

2.5 Integrated approach for NP discovery using AI/ML

Multiple strategies have been developed over the years for NP prioritization, and a combination of various approaches (e.g., genomic, metabolomic, taxonomic, spectral information, and bioactivity) can be used as a factor for ranking before the downstream process of purification and structure elucidation of NPs.¹⁶² More recently, Kim et al. developed NPClassifier,¹²⁸ which is a tool that can classify NPs using a DL approach. They have been categorized into three hierarchical levels based on the pathway and chemical properties; moreover, structural details can be classified using this NP, which indicates its applicability for drug discovery and the elucidation of biological interactions. In another study, an automated genome-guided NP discovery tool, with the support of an LC-MS/MS dataset, was developed that could automatically predict, combinatorialize, and identify polyketides and NRPs from crude extracts.¹²⁹ Hosein et al. developed NRPquest, which is an ML tool that integrates MS and Genome Mining for Nonribosomal Peptide (NRP) discovery.¹³⁰ Similarly, another tool (NRPminer) was developed very recently that combined both genomics and metabolomics to identify novel NRPs; using this approach, four unknown NRP families were identified from microbes and human microbiota and shown to exhibit antiparasitic activity.¹³¹ By integrating genomics and metabolomics focusing on NRPs, several novel protegomyacin derivatives from a previously unknown NP source (*X. doucetiae* and *X. poinarii*) were identified (Figure 3).¹³¹ A study reported by Kleigrew et al. integrated metabolomics and genome analysis to discover NPs from cyanobacteria; using this innovative approach, the authors discovered a new class of di- and trichlorinated acyl amide columbamides with cannabinomimetic activity.¹⁶³ Previously, we combined genome mining with MN and identified urdamycin E and a novel derivative, urdamycin V (Figure 3), from *Streptomyces* spp., which induce cell death by inhibiting mTOR in cancerous cells.^{164,165} Carlos et al. developed a database (DEREP-NP) to dereplicate metabolites efficiently by integrating MS and NMR spectra.¹³² Another study that combined NMR-based profiling with genome mining led to the discovery of the allenic macrolide Archangiumide (Figure 3) from *Myxobacterium*.¹⁶⁶ Using MS-guided genome mining, which connects the chemotypes of peptide NPs with their BGCs by iteratively matching *de novo* tandem MS, a new NP peptidogenomics approach was developed.¹⁶⁷ Using this combined approach, five new stendomycin analogues were identified that differed in the acyl chain and in valine or isoleucine substitutions at positions 5 and 13 from *S. hygrosopicus* ATCC 53653 (Figure 3).

3. Bioactivity and Target Identification of NPs with AI/ML techniques

One of the challenges in the development of NP-based drug candidates is the identification of their mechanism of action and side effects, which is a costly and lengthy process.^{168,169} Because of the enormous structural diversity and broad chemical spaces, the bioactivity of NPs is discovered based on phenotypic effects or via high-throughput phenotypic screening.^{170,171} To identify the targets experimentally, chemical genomics^{172,173} and chemical proteomics¹⁷⁴ approaches are generally used; however, although they can validate the targets they are often laborious and time-consuming processes.¹³³ To overcome this, computational approaches can narrow down the large search space of the targets.¹⁷⁵ There are three computational approaches and, in addition to the traditional structure-based¹⁷⁶ and ligand-based target identification methods,¹⁷⁷ ML-based approaches have numerous advantages and can be promising strategies for NP target identification.¹⁷⁸ To identify drug targets, Madhukar et al. developed BANDIT,¹³⁴ a Bayesian machine-learning approach that integrates multiple data types to predict drug binding targets.¹³⁴ Using this approach, the authors predicted the targets of nearly 4,000 compounds with 90% accuracy and further validated 14 novel microtubule inhibitors. In another study aimed at identifying drug–target interactions (DTIs), a CNN-based tool, NeoDTI, was developed.¹⁷⁹ NeoDTI mines large-scale graph data and automatically learns the topology-preserving representations of drugs and targets, to facilitate DTI prediction with compound–protein binding affinity. Using such approaches, the drug targets of NPs can be identified, which can accelerate the drug-discovery platform. In another study, a DL toolkit, “Openchem,” which is based on the PyTorch framework, was developed for drug design and computational chemistry.¹³⁹ It can enable drug discovery and molecular modeling applications using DL algorithms. This DL-based approach can help in various tasks in NP discovery, such as their physical properties and structure–activity relation. A recent study reported by Walker and Clardy described an ML-based approach to predict the biological activity of NPs using genome mining without isolation.¹⁸⁰ The authors used ML classifiers to predict antibacterial or antifungal activity based on known NP BGCs with an accuracy of 80%.

The SPiDER ML tool merges the concept of self-organizing maps, consensus scoring, and statistical analysis to successfully identify targets for both known drugs and computer-generated molecular scaffolds; moreover, using this method, off-target fenofibrate-related compounds were identified.¹³⁵ Furthermore, to increase the confidence, the Drug–Target Relationship Predictor (DECRyPT) machine intelligence workflow, which uses regression random forest technology as an orthogonal learning approach to self-organizing maps, was developed.¹³⁶ Using this ML tool, the targets of β -lapachone were identified and validated as potent modulators of 5-lipoxygenase.¹³⁶ SuperPred¹³⁷ provides drug classification and target prediction considering features such as 2-D, Fragment, and 3-D similarity and adapting concepts of the basic local

alignment search tool (BLAST) algorithm.^{137,181} These ML approaches can innovate the drug target identification process and serve as an alternative powerful strategy to chemoproteomics. Another study reported by Carrella et al. developed MANTRA 2.0, which is a transcriptional profile-based drug target identification that uses a microarray dataset.¹³⁸ By uploading the gene expression profile of the compound in cell lines, an ML-based automated pipeline revealed its mechanism of action based on the transcriptional signature of existing drugs.¹³⁸ Despite the advantages of the ML tools, they can sometimes be inaccurate and only the previously studied targets can be predicted with further target validation.^{182,183} In the drug-discovery process, one of the key criteria for candidate molecules is that they have fewer adverse effects; however, numerous time- and cost-intensive *in vitro* and *in vivo* studies are required to assess toxicity.¹⁸⁴ Computational toxicology can be effectively used to screen a large number of compounds without the use of time-consuming animal studies; nevertheless, this approach has severe drawbacks in terms of accuracy.¹⁸⁵ To overcome this issue, a recent study reported a DL pipeline, “DeepTox,” which exhibited a high accuracy of toxicity prediction.¹⁴⁰ Such a DL-based approach can be utilized in the future effectively to predict the toxicity of NPs and to tweak molecules with less adverse effects.

4. Conclusions and future perspectives

NPs from microorganisms and their molecular frameworks have a long tradition for many drug leads and are still widely used for treating various diseases and infections.^{182,186,187} The bioprospecting of the NP leads is challenging because of the amount of data generated and technical barriers, such as screening, isolation, characterization, and target identification. AI approaches can be used to address these problems and uncover hidden patterns by employing algorithms and decreasing the analytical time, resources, and costs required to identify NPs.¹⁸⁸ As proof of concept, recently, a highly effective antibiotic (halicin) with an entirely new mechanism of action was identified from the ZINC15 database using a DL approach.¹⁰ AI can help prioritize the microbes for screening based on their taxonomic novelty and genomes regarding the ability to produce novel NPs. Furthermore, it can help rapid dereplication and assist in the identification of active compounds using LC-HRMS and NMR.

Several NPs were isolated during the golden age of NPs, but most of them have been neglected or are limited by specific bioactivity with the discovery of various lead compounds at similar times.^{1,189} However, the surge of antimicrobial resistance and technological advancements have rekindled the interest in NPs as drug leads and repurposing is being assessed.¹⁹⁰ The cyclic peptide griselimycin was identified in 1960 from *Streptomyces*¹⁹¹ and exhibited potent antituberculous activity, but was neglected; however, very recently, it was modified and introduced into the drug-development pipeline.¹⁹² Similarly, another NP, chrysomycin A, which is a rare C-aryl glycoside, was first discovered over 60

years ago and has anticancer activity^{193,194} with no further studies; however, recently, it was reported as inhibiting multidrug-resistant tuberculosis effectively (MDR-TB).^{195,196} Drug repurposing and alternate bioactivity prediction are cost-effective processes compared with drug discovery; nevertheless, they are quite challenging. To overcome this drawback, AI/ML can be used for candidate selection.¹⁹⁷ Furthermore, AI can also assist in macromolecular target identification in a fast and effective manner.

A big obstacle in the full-fledged implementation of AI in NP research is the lack of integrated and curated databases.¹⁹⁸ Most of the data, such as taxonomic, structural, genomic, and metabolomic data, for the specific compounds are not available compiled in the form of databases and presented in the form of scientific literature, which is very difficult to access and analyze manually.^{198,199} Hence, an integrated approach is required for the effective analysis of NPs, as is a single algorithm for the management of the entire process of NP discovery alone. By addressing these issues, the common problems associated with AI, such as errors and repeatability, can be controlled in the learning process from reliable datasets.^{200–202} With the worsening drug-resistance scenario and the increase in the number of new infections, the search for novel NPs is essential. Nature is extremely generous to mankind by providing diverse compounds over the centuries to cure diseases. With the advent of technological advancements and AI, can we expect a new golden era of NP drug discovery?

5. Author contributions

VJS conceptualized the review process, wrote the manuscript, and drew the figures with support from MBL and VMD. SGD edited the manuscript. GNP: funding acquisition, HS conceptualized, supervised and edited the review process as well as funding acquisition.

Conflicts of interest

There are no conflicts to declare.

Acknowledgments

VJS acknowledges MEXT for fellowship support. We express sincere thanks for a by JSPS KAKENHI (grant numbers 20H05936 and 21H04705) to H.S. This work was also supported by NIH award R01CA236350 to H. S.

Notes and references

- 1 A. G. Atanasov, S. B. Zotchev, V. M. Dirsch, International Natural Product Sciences Taskforce and C. T. Supuran, *Nat. Rev. Drug Discov.*, 2021, **20**, 200–216.
- 2 L. Katz and R. H. Baltz, *J. Ind. Microbiol. Biotechnol.*, 2016, **43**, 155–176.
- 3 B. O. Bachmann, S. G. Van Lanen and R. H. Baltz, *J. Ind. Microbiol. Biotechnol.*, 2014, **41**, 175–184.
- 4 K. Grabowski, K.-H. Baringhaus and G. Schneider, *Nat. Prod.*

- Rep., 2008, **25**, 892–904.
- 5 D. M. Sabatini, *Proc. Natl. Acad. Sci. U. S. A.*, 2017, **114**, 11818–11825.
- 6 M. Skinnider, F. Wang, D. Pasin, R. Greiner, L. Foster, P. Dalsgaard and D. S. Wishart, *ChemRxiv*, DOI:10.26434/chemrxiv.14644854.v1.
- 7 J. Jiménez-Luna, F. Grisoni and G. Schneider, *Nature Machine Intelligence*, 2020, **2**, 573–584.
- 8 Z. J. Baum, X. Yu, P. Y. Ayala, Y. Zhao, S. P. Watkins and Q. Zhou, *J. Chem. Inf. Model.*, 2021, **61**, 3197–3212.
- 9 N. Choudhary, R. Bharti and R. Sharma, *Materials Today: Proceedings*, DOI:10.1016/j.matpr.2021.09.428.
- 10 J. M. Stokes, K. Yang, K. Swanson, W. Jin, A. Cubillos-Ruiz, N. M. Donghia, C. R. MacNair, S. French, L. A. Carfrae, Z. Bloom-Ackermann, V. M. Tran, A. Chiappino-Pepe, A. H. Badran, I. W. Andrews, E. J. Chory, G. M. Church, E. D. Brown, T. S. Jaakkola, R. Barzilay and J. J. Collins, *Cell*, 2020, **181**, 475–483.
- 11 Silver Lynn L., *Clin. Microbiol. Rev.*, 2011, **24**, 71–109.
- 12 D. Lyddiard, G. L. Jones and B. W. Greatrex, *FEMS Microbiol. Lett.*, DOI:10.1093/femsle/fnw084.
- 13 T. Hautbergue, E. L. Jamin, L. Debrauwer, O. Puel and I. P. Oswald, *Nat. Prod. Rep.*, 2018, **35**, 147–173.
- 14 A. Bouslimani, L. M. Sanchez, N. Garg and P. C. Dorrestein, *Nat. Prod. Rep.*, 2014, **31**, 718–729.
- 15 O. Genilloud, *Nat. Prod. Rep.*, 2017, **34**, 1203–1232.
- 16 K. Gerth, N. Bedorf, G. Höfle, H. Irschik and H. Reichenbach, *J. Antibiot.*, 1996, **49**, 560–563.
- 17 R. H. Felting, G. O. Buchanan, T. J. Mincer, C. A. Kauffman, P. R. Jensen and W. Fenical, *Angew. Chem. Int. Ed Engl.*, 2003, **42**, 355–357.
- 18 M. Rust, E. J. N. Helfrich, M. F. Freeman, P. Nanudorn, C. M. Field, C. Rückert, T. Kündig, M. J. Page, V. L. Webb, J. Kalinowski, S. Sunagawa and J. Piel, *Proc. Natl. Acad. Sci. U. S. A.*, 2020, **117**, 9508–9518.
- 19 L. L. Ling, T. Schneider, A. J. Peoples, A. L. Spoering, I. Engels, B. P. Conlon, A. Mueller, T. F. Schäberle, D. E. Hughes, S. Epstein, M. Jones, L. Lazarides, V. A. Steadman, D. R. Cohen, C. R. Felix, K. A. Fetterman, W. P. Millett, A. G. Nitti, A. M. Zullo, C. Chen and K. Lewis, *Nature*, 2015, **517**, 455–459.
- 20 E. Gavrish, C. S. Sit, S. Cao, O. Kandror, A. Spoering, A. Peoples, L. Ling, A. Fetterman, D. Hughes, A. Bissell, H. Torrey, T. Akopian, A. Mueller, S. Epstein, A. Goldberg, J. Clardy and K. Lewis, *Chem. Biol.*, 2014, **21**, 509–518.
- 21 P. Hugenholtz, M. Chuvochina, A. Oren, D. H. Parks and R. M. Soo, *ISME J.*, 2021, **15**, 1879–1892.
- 22 K. P. Smith, A. D. Kang and J. E. Kirby, *J. Clin. Microbiol.*, DOI:10.1128/JCM.01521-17.
- 23 C. M. Clark, M. S. Costa, L. M. Sanchez and B. T. Murphy, *Proc. Natl. Acad. Sci. U. S. A.*, 2018, **115**, 4981–4986.
- 24 C. Dumolin, M. Aerts, B. Verheyde, S. Schellaert, T. Vandamme, F. Van der Jeugt, E. De Canck, M. Cnockaert, A. D. Wieme, I. Cleenwerck, J. Peiren, P. Dawyndt, P. Vandamme and A. Carlier, *mSystems*, 2019, **4**.
- 25 M. Sauguet, B. Valot, X. Bertrand and D. Hocquet, *Trends Microbiol.*, 2017, **25**, 447–455.
- 26 E. Esquenazi, Y.-L. Yang, J. Watrous, W. H. Gerwick and P. C. Dorrestein, *Nat. Prod. Rep.*, 2009, **26**, 1521–1534.
- 27 Y.-L. Yang, Y. Xu, P. Straight and P. C. Dorrestein, *Nat. Chem. Biol.*, 2009, **5**, 885–887.
- 28 D. J. Gonzalez, N. M. Haste, A. Hollands, T. C. Fleming, M. Hamby, K. Pogliano, V. Nizet and P. C. Dorrestein, *Microbiology*, 2011, **157**, 2485–2492.
- 29 E. Esquenazi, C. Coates, L. Simmons, D. Gonzalez, W. H. Gerwick and P. C. Dorrestein, *Mol. Biosyst.*, 2008, **4**, 562–570.
- 30 R. H. Baltz, *J. Ind. Microbiol. Biotechnol.*, DOI:10.1093/jimb/kuab044.
- 31 A. Nivina, K. P. Yuet, J. Hsu and C. Khosla, *Chem. Rev.*, 2019, **119**, 12524–12547.
- 32 Y. Hai, A. Huang and Y. Tang, *J. Nat. Prod.*, 2020, **83**, 593–600.
- 33 M. W. Mullowney, R. A. McClure, M. T. Robey, N. L. Kelleher and R. J. Thomson, *Nat. Prod. Rep.*, 2018, **35**, 847–878.
- 34 M. Baunach, J. Franke and C. Hertweck, *Angew. Chem. Int. Ed Engl.*, 2015, **54**, 2604–2626.
- 35 K. Scherlach and C. Hertweck, *Nat. Commun.*, 2021, **12**, 3864.
- 36 K. Blin, V. Pascal Andreu, E. L. C. de Los Santos, F. Del Carratore, S. Y. Lee, M. H. Medema and T. Weber, *Nucleic Acids Res.*, 2019, **47**, D625–D630.
- 37 N. Ichikawa, M. Sasagawa, M. Yamamoto, H. Komaki, Y. Yoshida, S. Yamazaki and N. Fujita, *Nucleic Acids Res.*, 2013, **41**, D408–14.
- 38 K. Palaniappan, I.-M. A. Chen, K. Chu, A. Ratner, R. Seshadri, N. C. Kypides, N. N. Ivanova and N. J. Mouncey, *Nucleic Acids Research*, 2019.
- 39 S. A. Kautsar, K. Blin, S. Shaw, J. C. Navarro-Muñoz, B. R. Terlouw, J. J. van der Hooft, J. A. van Santen, V. Tracanna, H. G. Suarez Duran, V. Pascal Andreu, N. Selem-Mojica, M. Alanjary, S. L. Robinson, G. Lund, S. C. Epstein, A. C. Sisto, L. K. Charkoudian, J. Collemare, R. G. Linington, T. Weber and M. H. Medema, *Nucleic Acids Res.*, 2020, **48**, D454–D458.
- 40 K. R. Conway and C. N. Boddy, *Nucleic Acids Res.*, 2013, **41**, D402–7.
- 41 R. Hammami, A. Zouhir, C. Le Lay, J. Ben Hamida and I. Fliss, *BMC Microbiol.*, 2010, **10**, 22.
- 42 K. Blin, S. Shaw, K. Steinke, R. Villebro, N. Ziemert, S. Y. Lee, M. H. Medema and T. Weber, *Nucleic Acids Res.*, 2019, **47**, W81–W87.
- 43 M. A. Skinnider, C. A. Dejong, P. N. Rees, C. W. Johnston, H. Li, A. L. H. Webster, M. A. Wyatt and N. A. Magarvey, *Nucleic Acids Res.*, 2015, **43**, 9645–9662.
- 44 A. de Jong, S. A. F. T. van Hijum, J. J. E. Bijlsma, J. Kok and O. P. Kuipers, *Nucleic Acids Res.*, 2006, **34**, W273–9.
- 45 M. D. Mungan, M. Alanjary, K. Blin, T. Weber, M. H. Medema and N. Ziemert, *Nucleic Acids Res.*, 2020, **48**, W546–W552.
- 46 N. Sélem-Mojica, C. Aguilar, K. Gutiérrez-García, C. E. Martínez-Guerrero and F. Barona-Gómez, *Microb Genom.*, DOI:10.1099/mgen.0.000260.
- 47 N. Khaldi, F. T. Seifuddin, G. Turner, D. Haft, W. C. Nierman, K. H. Wolfe and N. D. Fedorova, *Fungal Genet. Biol.*, 2010, **47**, 736–741.
- 48 I. Takeda, M. Umemura, H. Koike, K. Asai and M. Machida, *DNA Res.*, 2014, **21**, 447–457.
- 49 G. D. Hannigan, D. Prihoda, A. Palicka, J. Soukup, O. Klempir, L. Rampula, J. Durcak, M. Wurst, J. Kotowski, D. Chang, R. Wang, G. Piizzi, G. Temesi, D. J. Hazuda, C. H. Woelk and D. A. Bitton, *Nucleic Acids Research*, 2019, **47**, e110–e110.
- 50 Z. Qin, J. T. Munnoch, R. Devine, N. A. Holmes, R. F. Seipke, K. A. Wilkinson, B. Wilkinson and M. I. Hutchings, *Chemical Science*, 2017, **8**, 3218–3227.
- 51 A. M. Kloosterman, P. Cimermanic, S. S. Elsayed, C. Du, M. Hadjithomas, M. S. Donia, M. A. Fischbach, G. P. van Wezel and M. H. Medema, *PLoS Biol.*, 2020, **18**, e3001026.
- 52 S. J. Miller and J. Clardy, *Nat. Chem.*, 2009, **1**, 261–263.
- 53 B. M. Hover, S.-H. Kim, M. Katz, Z. Charlop-Powers, J. G. Owen, M. A. Ternei, J. Maniko, A. B. Estrela, H. Molina, S. Park, D. S.

- Perlin and S. F. Brady, *Nat Microbiol*, 2018, **3**, 415–422.
- 54 Y. Sugimoto, F. R. Camacho, S. Wang, P. Chankhamjon, A. Odabas, A. Biswas, P. D. Jeffrey and M. S. Donia, *Science*, 2019, 366.
- 55 M. S. Donia and M. A. Fischbach, *Science*, 2015, **349**, 1254766.
- 56 M. Banf, K. Zhao and S. Y. Rhee, *Bioinformatics*, 2019, **35**, 3178–3180.
- 57 B. V. B. Reddy, A. Milshteyn, Z. Charlop-Powers and S. F. Brady, *Chem. Biol.*, 2014, **21**, 1023–1033.
- 58 C. D. Bader, F. Panter and R. Müller, *Biotechnol. Adv.*, 2020, **39**, 107480.
- 59 M. C. Stroe, T. Netzker, K. Scherlach, T. Krüger, C. Hertweck, V. Valiante and A. A. Brakhage, *Elife*, DOI:10.7554/eLife.52541.
- 60 H. B. Bode, B. Bethe, R. Höfs and A. Zeeck, *Chembiochem*, 2002, **3**, 619–627.
- 61 M. R. Seyedsayamdost, *Proc. Natl. Acad. Sci. U. S. A.*, 2014, **111**, 7266–7271.
- 62 J. Shima, A. Hesketh, S. Okamoto, S. Kawamoto and K. Ochi, *J. Bacteriol.*, 1996, **178**, 7276–7284.
- 63 T. Asai, T. Yamamoto, N. Shirata, T. Taniguchi, K. Monde, I. Fujii, K. Gomi and Y. Oshima, *Org. Lett.*, 2013, **15**, 3346–3349.
- 64 X.-M. Mao, W. Xu, D. Li, W.-B. Yin, Y.-H. Chooi, Y.-Q. Li, Y. Tang and Y. Hu, *Angew. Chem. Int. Ed Engl.*, 2015, **54**, 7592–7596.
- 65 T. Zhang, J. Wan, Z. Zhan, J. Bai, B. Liu and Y. Hu, *Acta Pharm Sin B*, 2018, **8**, 478–487.
- 66 J. B. Biggins, X. Liu, Z. Feng and S. F. Brady, *J. Am. Chem. Soc.*, 2011, **133**, 1638–1641.
- 67 A. B. Ramzi, S. N. Baharum, H. Bunawan and N. S. Scrutton, *Front. Bioeng. Biotechnol.*, 2020, **8**, 608918.
- 68 N. J. Treloar, A. J. H. Fedorec, B. Ingalls and C. P. Barnes, *PLoS Comput. Biol.*, 2020, **16**, e1007783.
- 69 F. Xu, Y. Wu, C. Zhang, K. M. Davis, K. Moon, L. B. Bushin and M. R. Seyedsayamdost, *Nat. Chem. Biol.*, 2019, **15**, 161–168.
- 70 B. C. Covington and M. R. Seyedsayamdost, *ACS Chem. Biol.*, DOI:10.1021/acscchembio.1c00737.
- 71 B. Shen, *Cell*, 2015, **163**, 1297–1300.
- 72 P. R. Jensen, K. L. Chavarria, W. Fenical, B. S. Moore and N. Ziemert, *Journal of Industrial Microbiology and Biotechnology*, 2014, **41**, 203–209.
- 73 T. Ito and M. Masubuchi, *The Journal of Antibiotics*, 2014, **67**, 353–360.
- 74 S. P. Gaudêncio and F. Pereira, *Natural Product Reports*, 2015, **32**, 779–810.
- 75 D. J. Hook, C. F. More, J. J. Yacobucci, G. Dubay and S. O'Connor, *J. Chromatogr.*, 1987, **385**, 99–108.
- 76 J. Buckingham, *Dictionary of natural products, supplement 3: Third supplement*, CRC Press, London, England, 1996.
- 77 J. W. Blunt, A. R. Carroll, B. R. Copp, R. A. Davis, R. A. Keyzers and M. R. Prinsep, *Nat. Prod. Rep.*, 2018, **35**, 8–53.
- 78 A. F. A. Moumbock, M. Gao, A. Qaseem, J. Li, P. A. Kirchner, B. Ndingkokhar, B. D. Bekono, C. V. Simoben, S. B. Babiaka, Y. I. Malange, F. Sauter, P. Zierrep, F. Ntie-Kang and S. Günther, *Nucleic Acids Res.*, 2021, **49**, D600–D604.
- 79 T. Tomiki, T. Saito, M. Ueki, H. Konno, T. Asaoka, R. Suzuki, M. Uramoto, H. Kakeya and H. Osada, *J Comput Aid Chem*, 2006, **7**, 157–162.
- 80 J. A. van Santen, G. Jacob, A. L. Singh, V. Aniebok, M. J. Balunas, D. Bunsko, F. C. Neto, L. Castañó-Espriu, C. Chang, T. N. Clark, J. L. Cleary Little, D. A. Delgadillo, P. C. Dorrestein, K. R. Duncan, J. M. Egan, M. M. Galey, F. P. J. Haeckl, A. Hua, A. H. Hughes, D. Iskakova, A. Khadiikar, J.-H. Lee, S. Lee, N. LeGrow, D. Y. Liu, J. M. Macho, C. S. McCaughey, M. H. Medema, R. P. Neupane, T. J. O'Donnell, J. S. Paula, L. M. Sanchez, A. F. Shaikh, S. Soldatou, B. R. Terlouw, T. A. Tran, M. Valentine, J. J. J. van der Hooft, D. A. Vo, M. Wang, D. Wilson, K. E. Zink and R. G. Lington, *ACS Cent Sci*, 2019, **5**, 1824–1833.
- 81 G. T. Mehetre, J. S. Vinodh, B. B. Burkul, D. Desai, B. Santhakumari, M. S. Dharne and S. G. Dastager, *RSC Advances*, 2019, **9**, 9850–9859.
- 82 L. K. Caesar, J. J. Kellogg, O. M. Kvalheim and N. B. Cech, *J. Nat. Prod.*, 2019, **82**, 469–484.
- 83 J. Hubert, J.-M. Nuzillard and J.-H. Renault, *Phytochemistry Reviews*, 2017, **16**, 55–95.
- 84 T. Hoffmann, D. Krug, S. Hüttel and R. Müller, *Anal. Chem.*, 2014, **86**, 10780–10788.
- 85 P. Kaur and P. B. O'Connor, *Journal of the American Society for Mass Spectrometry*, 2006, **17**, 459–468.
- 86 T. Pluskal, S. Castillo, A. Villar-Briones and M. Oresic, *BMC Bioinformatics*, 2010, **11**, 395.
- 87 Z. Pang, G. Zhou, J. Ewald, L. Chang, O. Hacariz, N. Basu and J. Xia, *Nat. Protoc.*, DOI:10.1038/s41596-022-00710-w.
- 88 H. Tsugawa, T. Cajka, T. Kind, Y. Ma, B. Higgins, K. Ikeda, M. Kanazawa, J. VanderGheynst, O. Fiehn and M. Arita, *Nat. Methods*, 2015, **12**, 523–526.
- 89 N. Jaitly, A. Mayampurath, K. Littlefield, J. N. Adkins, G. A. Anderson and R. D. Smith, *BMC Bioinformatics*, 2009, **10**, 87.
- 90 C. A. Smith, E. J. Want, G. O'Maille, R. Abagyan and G. Siuzdak, *Anal. Chem.*, 2006, **78**, 779–787.
- 91 D. M. Horn, R. A. Zubarev and F. W. McLafferty, *J. Am. Soc. Mass Spectrom.*, 2000, **11**, 320–332.
- 92 V. Kumar, A. A. Kumar, V. Joseph, V. M. Dan, A. Jaleel, T. R. S. Kumar and C. C. Kartha, *Mol. Cell. Biochem.*, 2020, **463**, 147–160.
- 93 H. Horai, M. Arita, S. Kanaya, Y. Nihei, T. Ikeda, K. Suwa, Y. Ojima, K. Tanaka, S. Tanaka, K. Aoshima, Y. Oda, Y. Kakazu, M. Kusano, T. Tohge, F. Matsuda, Y. Sawada, M. Y. Hirai, H. Nakanishi, K. Ikeda, N. Akimoto, T. Maoka, H. Takahashi, T. Ara, N. Sakurai, H. Suzuki, D. Shibata, S. Neumann, T. Iida, K. Tanaka, K. Funatsu, F. Matsuura, T. Soga, R. Taguchi, K. Saito and T. Nishioka, *J. Mass Spectrom.*, 2010, **45**, 703–714.
- 94 C. A. Smith, G. O'Maille, E. J. Want, C. Qin, S. A. Trauger, T. R. Brandon, D. E. Custodio, R. Abagyan and G. Siuzdak, *Ther. Drug Monit.*, 2005, **27**, 747–751.
- 95 M. Sud, E. Fahy, D. Cotter, A. Brown, E. A. Dennis, C. K. Glass, A. H. Merrill Jr, R. C. Murphy, C. R. H. Raetz, D. W. Russell and S. Subramaniam, *Nucleic Acids Res.*, 2007, **35**, D527–32.
- 96 M. Wang, J. J. Carver, V. V. Phelan, L. M. Sanchez, N. Garg, Y. Peng, D. D. Nguyen, J. Watrous, C. A. Kapono, T. Luzzatto-Knaan, C. Porto, A. Bouslimani, A. V. Melnik, M. J. Meehan, W.-T. Liu, M. Crüsemann, P. D. Boudreau, E. Esquenazi, M. Sandoval-Calderón, R. D. Kersten, L. A. Pace, R. A. Quinn, K. R. Duncan, C.-C. Hsu, D. J. Floros, R. G. Gavilan, K. Kleigrew, T. Northen, R. J. Dutton, D. Parrot, E. E. Carlson, B. Aigle, C. F. Michelsen, L. Jelsbak, C. Sohlenkamp, P. Pevzner, A. Edlund, J. McLean, J. Piel, B. T. Murphy, L. Gerwick, C.-C. Liaw, Y.-L. Yang, H.-U. Humpf, M. Maansson, R. A. Keyzers, A. C. Sims, A. R. Johnson, A. M. Sidebottom, B. E. Sedio, A. Klitgaard, C. B. Larson, C. A. B. P, D. Torres-Mendoza, D. J. Gonzalez, D. B. Silva, L. M. Marques, D. P. Demarque, E. Pociute, E. C. O'Neill, E. Briand, E. J. N. Helfrich, E. A. Granatosky, E. Glukhov, F. Ryffel, H. Houson, H. Mohimani, J. J. Kharbush, Y. Zeng, J. A. Vorholt, K. L. Kurita, P. Charusanti, K. L. McPhail, K. F. Nielsen, L. Vuong, M. Elfeki, M. F. Traxler, N. Engene, N. Koyama, O. B. Vining, R. Baric, R. R. Silva, S. J. Mascuch, S. Tomasi, S. Jenkins, V.

- Macherla, T. Hoffman, V. Agarwal, P. G. Williams, J. Dai, R. Neupane, J. Gurr, A. M. C. Rodríguez, A. Lamsa, C. Zhang, K. Dorrestein, B. M. Duggan, J. Almaliti, P.-M. Allard, P. Phapale, L.-F. Nothias, T. Alexandrov, M. Litaudon, J.-L. Wolfender, J. E. Kyle, T. O. Metz, T. Peryea, D.-T. Nguyen, D. VanLeer, P. Shinn, A. Jadhav, R. Müller, K. M. Waters, W. Shi, X. Liu, L. Zhang, R. Knight, P. R. Jensen, B. O. Palsson, K. Pogliano, R. G. Linington, M. Gutiérrez, N. P. Lopes, W. H. Gerwick, B. S. Moore, P. C. Dorrestein and N. Bandeira, *Nat. Biotechnol.*, 2016, **34**, 828–837.
- 97 J. Watrous, P. Roach, T. Alexandrov, B. S. Heath, J. Y. Yang, R. D. Kersten, M. van der Voort, K. Pogliano, H. Gross, J. M. Raaijmakers, B. S. Moore, J. Laskin, N. Bandeira and P. C. Dorrestein, *Proc. Natl. Acad. Sci. U. S. A.*, 2012, **109**, E1743–52.
- 98 A. T. Aron, E. C. Gentry, K. L. McPhail, L.-F. Nothias, M. Nothias-Esposito, A. Bouslimani, D. Petras, J. M. Gauglitz, N. Sikora, F. Vargas, J. J. J. van der Hooft, M. Ernst, K. B. Kang, C. M. Aceves, A. M. Caraballo-Rodríguez, I. Koester, K. C. Weldon, S. Bertrand, C. Roullier, K. Sun, R. M. Tehan, C. A. Boya P, M. H. Christian, M. Gutiérrez, A. M. Ulloa, J. A. Tejada Mora, R. Mojica-Flores, J. Lakey-Beitia, V. Vásquez-Chaves, Y. Zhang, A. I. Calderón, N. Tayler, R. A. Keyzers, F. Tugizimana, N. Ndlovu, A. A. Aksenov, A. K. Jarmusch, R. Schmid, A. W. Truman, N. Bandeira, M. Wang and P. C. Dorrestein, *Nat. Protoc.*, 2020, **15**, 1954–1991.
- 99 F. Vincenti, C. Montesano, F. Di Ottavio, A. Gregori, D. Compagnone, M. Sergi and P. Dorrestein, *Front. Chem.*, 2020, **8**, 572952.
- 100 D. D. Nguyen, A. V. Melnik, N. Koyama, X. Lu, M. Schorn, J. Fang, K. Aguinaldo, T. L. Lincecum Jr, M. G. K. Ghequire, V. J. Carrion, T. L. Cheng, B. M. Duggan, J. G. Malone, T. H. Mauchline, L. M. Sanchez, A. M. Kilpatrick, J. M. Raaijmakers, R. De Mot, B. S. Moore, M. H. Medema and P. C. Dorrestein, *Nat Microbiol*, 2016, **2**, 16197.
- 101 R. Teta, G. D. Sala, G. Esposito, C. W. Via, C. Mazzoccoli, C. Piccoli, M. J. Bertin, V. Costantino and A. Mangoni, *Org Chem Front*, 2019, **6**, 1762–1774.
- 102 L.-F. Nothias, D. Petras, R. Schmid, K. Dührkop, J. Rainer, A. Sarvepalli, I. Protsyuk, M. Ernst, H. Tsugawa, M. Fleischauer, F. Aicheler, A. A. Aksenov, O. Alka, P.-M. Allard, A. Barsch, X. Cachet, A. M. Caraballo-Rodríguez, R. R. Da Silva, T. Dang, N. Garg, J. M. Gauglitz, A. Gurevich, G. Isaac, A. K. Jarmusch, Z. Kameník, K. B. Kang, N. Kessler, I. Koester, A. Korf, A. Le Gouellec, M. Ludwig, C. Martin H, L.-I. McCall, J. McSayles, S. W. Meyer, H. Mohimani, M. Morsy, O. Moyne, S. Neumann, H. Neuweger, N. H. Nguyen, M. Nothias-Esposito, J. Paolini, V. V. Phelan, T. Pluskal, R. A. Quinn, S. Rogers, B. Shrestha, A. Tripathi, J. J. J. van der Hooft, F. Vargas, K. C. Weldon, M. Witting, H. Yang, Z. Zhang, F. Zubeil, O. Kohlbacher, S. Böcker, T. Alexandrov, N. Bandeira, M. Wang and P. C. Dorrestein, *Nat. Methods*, 2020, **17**, 905–908.
- 103 R. Schmid, D. Petras, L.-F. Nothias, M. Wang, A. T. Aron, A. Jagels, H. Tsugawa, J. Rainer, M. Garcia-Aloy, K. Dührkop, A. Korf, T. Pluskal, Z. Kameník, A. K. Jarmusch, A. M. Caraballo-Rodríguez, K. C. Weldon, M. Nothias-Esposito, A. A. Aksenov, A. Bauermeister, A. Albarracin Orío, C. O. Grundmann, F. Vargas, I. Koester, J. M. Gauglitz, E. C. Gentry, Y. Hövelmann, S. A. Kalinina, M. A. Pendergraft, M. Panitchpakdi, R. Tehan, A. Le Gouellec, G. Aleti, H. Mannocho Russo, B. Arndt, F. Hübner, H. Hayen, H. Zhi, M. Raffatellu, K. A. Prather, L. I. Aluwihare, S. Böcker, K. L. McPhail, H.-U. Humpf, U. Karst and P. C. Dorrestein, *Nat. Commun.*, 2021, **12**, 3832.
- 104 R. Reher, A. T. Aron, P. Fajtová, P. Stincone, C. Liu, I. Y. Ben Shalom, W. Bittremieux, M. Wang, M. L. Matos-Hernandez, K. L. Alexander, E. J. Caro-Díaz, C. B. Naman, C. C. Hughes, P. C. Dorrestein, A. J. O'Donoghue, W. H. Gerwick and D. Petras, *bioRxiv*, 2021.
- 105 H. Mohimani, A. Gurevich, A. Shlemov, A. Mikheenko, A. Korobeynikov, L. Cao, E. Shcherbin, L.-F. Nothias, P. C. Dorrestein and P. A. Pevzner, *Nat. Commun.*, 2018, **9**, 4035.
- 106 L.-F. Nothias, M. Nothias-Esposito, R. da Silva, M. Wang, I. Protsyuk, Z. Zhang, A. Sarvepalli, P. Leyssen, D. Touboul, J. Costa, J. Paolini, T. Alexandrov, M. Litaudon and P. C. Dorrestein, *J. Nat. Prod.*, 2018, **81**, 758–767.
- 107 S. Lee, J. A. van Santen, N. Farzaneh, D. Y. Liu, C. R. Pye, T. U. H. Baumeister, W. R. Wong and R. G. Linington, *ACS Cent. Sci.*, 2022, **8**, 223–234.
- 108 M. J. Helf, B. W. Fox, A. B. Artyukhin, Y. K. Zhang and F. C. Schroeder, *Nat. Commun.*, 2022, **13**, 782.
- 109 L. Chen, W. Lu, L. Wang, X. Xing, Z. Chen, X. Teng, X. Zeng, A. D. Muscarella, Y. Shen, A. Cowan, M. R. McReynolds, B. J. Kennedy, A. M. Lato, S. R. Campagna, M. Singh and J. D. Rabinowitz, *Nat. Methods*, 2021, **18**, 1377–1385.
- 110 K. Dührkop, M. Fleischauer, M. Ludwig, A. A. Aksenov, A. V. Melnik, M. Meusel, P. C. Dorrestein, J. Rousu and S. Böcker, *Nat. Methods*, 2019, **16**, 299–302.
- 111 M. Ludwig, L.-F. Nothias, K. Dührkop, I. Koester, M. Fleischauer, M. A. Hoffmann, D. Petras, F. Vargas, M. Morsy, L. Aluwihare, P. C. Dorrestein and S. Böcker, *Nat Mach Intell*, 2020, **2**, 629–641.
- 112 K. Dührkop, L.-F. Nothias, M. Fleischauer, R. Reher, M. Ludwig, M. A. Hoffmann, D. Petras, W. H. Gerwick, J. Rousu, P. C. Dorrestein and S. Böcker, *Nat. Biotechnol.*, 2021, **39**, 462–471.
- 113 F. Olivon, N. Elie, G. Grelier, F. Roussi, M. Litaudon and D. Touboul, *Anal. Chem.*, 2018, **90**, 13900–13908.
- 114 Y. Liu, A. Mrzic, P. Meysman, T. De Vijlder, E. P. Romijn, D. Valkenburg, W. Bittremieux and K. Laukens, *PLoS One*, 2020, **15**, e0226770.
- 115 L. Cao, M. Guler, A. Tagirdzhanov, Y.-Y. Lee, A. Gurevich and H. Mohimani, *Nat. Commun.*, 2021, **12**, 3718.
- 116 W. Bittremieux, K. Laukens, W. S. Noble and P. C. Dorrestein, *Rapid Commun. Mass Spectrom.*, 2021, e9153.
- 117 D. G. C. Treen, M. Wang, S. Xing, K. B. Louie, T. Huan, P. C. Dorrestein, T. R. Northen and B. P. Bowen, *Nat. Commun.*, 2022, **13**, 2510.
- 118 M. Ernst, K. B. Kang, A. M. Caraballo-Rodríguez, L.-F. Nothias, J. Wandy, C. Chen, M. Wang, S. Rogers, M. H. Medema, P. C. Dorrestein and J. J. J. van der Hooft, *Metabolites*, 2019, **9**, 144.
- 119 M. A. Hoffmann, L.-F. Nothias, M. Ludwig, M. Fleischauer, E. C. Gentry, M. Witting, P. C. Dorrestein, K. Dührkop and S. Böcker, *Nat. Biotechnol.*, , DOI:10.1038/s41587-021-01045-9.
- 120 J. L. López-Pérez, R. Therón, E. del Olmo and D. Díaz, *Bioinformatics*, 2007, **23**, 3256–3257.
- 121 D. S. Wishart, Z. Sayeeda, Z. Budinski, A. Guo, B. L. Lee, M. Berjanskii, M. Rout, H. Peters, R. Dizon, R. Mah, C. Torres-Calzada, M. Hiebert-Giesbrecht, D. Varshavi, D. Varshavi, E. Oler, D. Allen, X. Cao, V. Gautam, A. Maras, E. F. Poynton, P. Tavangar, V. Yang, J. A. van Santen, R. Ghosh, S. Sarma, E. Knutson, V. Sullivan, A. M. Jystad, R. Renslow, L. W. Sumner, R. G. Linington and J. R. Cort, *Nucleic Acids Res.*, , DOI:10.1093/nar/gkab1052.
- 122 A. Howarth, K. Ermanis and J. M. Goodman, *Chem. Sci.*, 2020, **11**, 4351–4359.
- 123 D.-W. Li, A. L. Hansen, C. Yuan, L. Bruschweiler-Li and R.

- Brüschweiler, *Nat. Commun.*, 2021, **12**, 5229.
- 124 A. Bruguère, S. Derbré, J. Dietsch, J. Leguy, V. Rahier, Q. Pottier, D. Bréard, S. Suor-Cherer, G. Vialat, A.-M. Le Ray, F. Saubion and P. Richomme, *Anal. Chem.*, 2020, **92**, 8793–8801.
- 125 U. Grienke, P. A. Foster, J. Zwirchmayr, A. Tahir, J. M. Rollinger and E. Mikros, *Sci. Rep.*, 2019, **9**, 11113.
- 126 H. W. Kim, C. Zhang, G. W. Cottrell and W. H. Gerwick, *Magn. Reson. Chem.*, DOI:10.1002/mrc.5240.
- 127 R. Reher, H. W. Kim, C. Zhang, H. H. Mao, M. Wang, L.-F. Nothias, A. M. Caraballo-Rodríguez, E. Glukhov, B. Teke, T. Leao, K. L. Alexander, B. M. Duggan, E. L. Van Everbroeck, P. C. Dorrestein, G. W. Cottrell and W. H. Gerwick, *J. Am. Chem. Soc.*, 2020, **142**, 4114–4120.
- 128 H. W. Kim, M. Wang, C. A. Leber, L.-F. Nothias, R. Reher, K. B. Kang, J. J. van der Hooft, P. C. Dorrestein, W. H. Gerwick and G. W. Cottrell, *J. Nat. Prod.*, 2021, **84**, 2795–2807.
- 129 C. W. Johnston, M. A. Skinnider, M. A. Wyatt, X. Li, M. R. M. Ranieri, L. Yang, D. L. Zechel, B. Ma and N. A. Magarvey, *Nat. Commun.*, 2015, **6**, 8421.
- 130 H. Mohimani, W.-T. Liu, R. D. Kersten, B. S. Moore, P. C. Dorrestein and P. A. Pevzner, *J. Nat. Prod.*, 2014, **77**, 1902–1909.
- 131 B. Behsaz, E. Bode, A. Gurevich, Y.-N. Shi, F. Grundmann, D. Acharya, A. M. Caraballo-Rodríguez, A. Bouslimani, M. Panitchpakdi, A. Linck, C. Guan, J. Oh, P. C. Dorrestein, H. B. Bode, P. A. Pevzner and H. Mohimani, *Nat. Commun.*, 2021, **12**, 3225.
- 132 C. L. Zani and A. R. Carroll, *J. Nat. Prod.*, 2017, **80**, 1758–1766.
- 133 X. Zeng, S. Zhu, W. Lu, Z. Liu, J. Huang, Y. Zhou, J. Fang, Y. Huang, H. Guo, L. Li, B. D. Trapp, R. Nussinov, C. Eng, J. Loscalzo and F. Cheng, *Chem. Sci.*, 2020, **11**, 1775–1797.
- 134 N. S. Madhukar, P. K. Khade, L. Huang, K. Gayvert, G. Galletti, M. Stogniew, J. E. Allen, P. Giannakakou and O. Elemento, *Nat. Commun.*, 2019, **10**, 5221.
- 135 D. Reker, T. Rodrigues, P. Schneider and G. Schneider, *Proc. Natl. Acad. Sci. U. S. A.*, 2014, **111**, 4067–4072.
- 136 T. Rodrigues, M. Werner, J. Roth, E. H. G. da Cruz, M. C. Marques, P. Akkapeddi, S. A. Lobo, A. Koeberle, F. Corzana, E. N. da Silva Júnior, O. Werz and G. J. L. Bernardes, *Chem. Sci.*, 2018, **9**, 6899–6903.
- 137 J. Nickel, B.-O. Gohlke, J. Erehman, P. Banerjee, W. W. Rong, A. Goede, M. Dunkel and R. Preissner, *Nucleic Acids Res.*, 2014, **42**, W26–31.
- 138 D. Carrella, F. Napolitano, R. Rispoli, M. Miglietta, A. Carissimo, L. Cutillo, F. Sirci, F. Gregoretti and D. Di Bernardo, *Bioinformatics*, 2014, **30**, 1787–1788.
- 139 M. Korshunova, B. Ginsburg, A. Tropsha and O. Isayev, *J. Chem. Inf. Model.*, 2021, **61**, 7–13.
- 140 G. Klambauer, T. Unterthiner, A. Mayr and S. Hochreiter, *Toxicol. Lett.*, 2017, **280**, S69.
- 141 F. Huber, S. van der Burg, J. J. van der Hooft and L. Ridder, *J. Cheminform.*, 2021, **13**, 84.
- 142 S. Rogers, C. W. Ong, J. Wandy, M. Ernst, L. Ridder and J. J. van der Hooft, *Faraday Discuss.*, 2019, **218**, 284–302.
- 143 J. J. van der Hooft, J. Wandy, F. Young, S. Padmanabhan, K. Gerasimidis, K. E. V. Burgess, M. P. Barrett and S. Rogers, *Anal. Chem.*, 2017, **89**, 7569–7577.
- 144 A. V. Buevich and M. E. Elyashberg, *J. Nat. Prod.*, 2016, **79**, 3105–3116.
- 145 W. F. Reynolds, *Pharmacognosy*, 2017, 567–596.
- 146 R. K. Lindsay, *Applications of Artificial Intelligence for Organic Chemistry: The DENDRAL Project*, McGraw-Hill Book Company, 1980.
- 147 E. L. Ulrich, H. Akutsu, J. F. Doreleijers, Y. Harano, Y. E. Ioannidis, J. Lin, M. Livny, S. Mading, D. Maziuk, Z. Miller, E. Nakatani, C. F. Schulte, D. E. Tolmie, R. Kent Wenger, H. Yao and J. L. Markley, *Nucleic Acids Res.*, 2008, **36**, D402–8.
- 148 J. T. Fishedick, S. R. Johnson, R. E. B. Ketchum, R. B. Croteau and B. M. Lange, *Phytochemistry*, 2015, **113**, 87–95.
- 149 J. B. McAlpine, S.-N. Chen, A. Kutateladze, J. B. MacMillan, G. Appendino, A. Barison, M. A. Beniddir, M. W. Biavatti, S. Bluml, A. Boufridi, M. S. Butler, R. J. Capon, Y. H. Choi, D. Coppage, P. Crews, M. T. Crimmins, M. Csete, P. Dewapriya, J. M. Egan, M. J. Garson, G. Genta-Jouve, W. H. Gerwick, H. Gross, M. K. Harper, P. Hermanto, J. M. Hook, L. Hunter, D. Jeannerat, N.-Y. Ji, T. A. Johnson, D. G. I. Kingston, H. Koshino, H.-W. Lee, G. Lewin, J. Li, R. G. Linington, M. Liu, K. L. McPhail, T. F. Molinski, B. S. Moore, J.-W. Nam, R. P. Neupane, M. Niemitz, J.-M. Nuzillard, N. H. Oberlies, F. M. M. Ocampos, G. Pan, R. J. Quinn, D. S. Reddy, J.-H. Renault, J. Rivera-Chávez, W. Robien, C. M. Saunders, T. J. Schmidt, C. Seger, B. Shen, C. Steinbeck, H. Stuppner, S. Sturm, O. Tagliatalata-Scafati, D. J. Tantillo, R. Verpoorte, B.-G. Wang, C. M. Williams, P. G. Williams, J. Wist, J.-M. Yue, C. Zhang, Z. Xu, C. Simmler, D. C. Lankin, J. Bisson and G. F. Pauli, *Nat. Prod. Rep.*, 2019, **36**, 35–107.
- 150 J.-M. Nuzillard and B. Plainchont, *Magn. Reson. Chem.*, 2018, **56**, 458–468.
- 151 M. E. Elyashberg, K. A. Blinov, A. J. Williams, S. G. Molodtsov, G. E. Martin and E. R. Martirosian, *J. Chem. Inf. Comput. Sci.*, 2004, **44**, 771–792.
- 152 C. Cobas, F. Seoane, E. Vaz, M. A. Bernstein, S. Dominguez, M. Pérez and S. Sýkora, *Magn. Reson. Chem.*, 2013, **51**, 649–654.
- 153 J.-M. Nuzillard and G. Massiot, *Anal. Chim. Acta*, 1991, **242**, 37–41.
- 154 R. B. Williams, M. O’Neil-Johnson, A. J. Williams, P. Wheeler, R. Pol and A. Moser, *Org. Biomol. Chem.*, 2015, **13**, 9957–9962.
- 155 A. Bakiri, B. Plainchont, V. de Paulo Emerenciano, R. Reynaud, J. Hubert, J.-H. Renault and J.-M. Nuzillard, *Mol. Inform.*, DOI:10.1002/minf.201700027.
- 156 J. K. Cooper, K. Li, J. Aubé, D. A. Coppage and J. P. Konopelski, *Org. Lett.*, 2018, **20**, 4314–4317.
- 157 C. I. MacGregor, B. Y. Han, J. M. Goodman and I. Paterson, *Chemical Communications*, 2016, 52, 4632–4635.
- 158 K. M. Snyder, J. Sikorska, T. Ye, L. Fang, W. Su, R. G. Carter, K. L. McPhail and P. H.-Y. Cheong, *Org. Biomol. Chem.*, 2016, **14**, 5826–5831.
- 159 S. H. Martínez-Treviño, V. Uc-Cetina, M. A. Fernández-Herrera and G. Merino, *J. Chem. Inf. Model.*, 2020, **60**, 3376–3386.
- 160 A. Vignoli, V. Ghini, G. Meoni, C. Licari, P. G. Takis, L. Tenori, P. Turano and C. Luchinat, *Angew. Chem. Int. Ed Engl.*, 2019, **58**, 968–994.
- 161 L. F. Silva-Castro, S. Derbré, A. M. Le Ray, P. Richomme, K. García-Sosa and L. M. Peña-Rodríguez, *Phytochem. Anal.*, 2021, **32**, 1102–1109.
- 162 J.-L. Wolfender, M. Litaudon, D. Touboul and E. F. Queiroz, *Nat. Prod. Rep.*, 2019, **36**, 855–868.
- 163 K. Kleigrew, J. Almaliti, I. Y. Tian, R. B. Kinnel, A. Korobeynikov, E. A. Monroe, B. M. Duggan, V. Di Marzo, D. H. Sherman, P. C. Dorrestein, L. Gerwick and W. H. Gerwick, *J. Nat. Prod.*, 2015, **78**, 1671–1682.
- 164 V. M. Dan, V. J S, S. C J, R. Sanawar, A. Lekshmi, R. A. Kumar, T. R. Santhosh Kumar, U. K. Marelli, S. G. Dastager and M. R. Pillai, *ACS Chem. Biol.*, 2020, **15**, 780–788.
- 165 V. M. Dan, B. Muralikrishnan, R. Sanawar, V. J S, B. B. Burkul, K.

- P. Srinivas, A. Lekshmi, N. S. Pradeep, S. G. Dastager, B. Santhakumari, T. R. Santhoshkumar, R. A. Kumar and M. R. Pillai, *Sci. Rep.*, 2018, **8**, 2810.
- 166 J.-Q. Hu, J.-J. Wang, Y.-L. Li, L. Zhuo, A. Zhang, H.-Y. Sui, X.-J. Li, T. Shen, Y. Yin, Z.-H. Wu, W. Hu, Y.-Z. Li and C. Wu, *Org. Lett.*, 2021, **23**, 2114–2119.
- 167 R. D. Kersten, Y.-L. Yang, Y. Xu, P. Cimermancic, S.-J. Nam, W. Fenical, M. A. Fischbach, B. S. Moore and P. C. Dorrestein, *Nat. Chem. Biol.*, 2011, **7**, 794–802.
- 168 X. Chen, Y. Wang, N. Ma, J. Tian, Y. Shao, B. Zhu, Y. K. Wong, Z. Liang, C. Zou and J. Wang, *Signal Transduct Target Ther*, 2020, **5**, 72.
- 169 L. Dai, Z. Li, D. Chen, L. Jia, J. Guo, T. Zhao and P. Nordlund, *Pharmacol. Ther.*, 2020, **216**, 107690.
- 170 M. Feher and J. M. Schmidt, *ChemInform*, 2003, 34.
- 171 J. G. Moffat, F. Vincent, J. A. Lee, J. Eder and M. Prunotto, *Nat. Rev. Drug Discov.*, 2017, **16**, 531–543.
- 172 L. I. Zon and R. T. Peterson, *Nature Reviews Drug Discovery*, 2005, **4**, 35–44.
- 173 J. Lamb, E. D. Crawford, D. Peck, J. W. Modell, I. C. Blat, M. J. Wrobel, J. Lerner, J.-P. Brunet, A. Subramanian, K. N. Ross, M. Reich, H. Hieronymus, G. Wei, S. A. Armstrong, S. J. Haggarty, P. A. Clemons, R. Wei, S. A. Carr, E. S. Lander and T. R. Golub, *Science*, 2006, **313**, 1929–1935.
- 174 K. M. Lum, Y. Sato, B. A. Beyer, W. C. Plaisted, J. L. Anglin, L. L. Lairson and B. F. Cravatt, *ACS Chem. Biol.*, 2017, **12**, 2671–2681.
- 175 G. R. Langley, I. M. Adcock, F. Busquet, K. M. Crofton, E. Csernok, C. Giese, T. Heinonen, K. Herrmann, M. Hofmann-Apitius, B. Landesmann, L. J. Marshall, E. Mclvor, A. R. Muotri, F. Noor, K. Schutte, T. Seidle, A. van de Stolpe, H. Van Esch, C. Willett and G. Wozczek, *Drug Discov. Today*, 2017, **22**, 327–339.
- 176 G. M. Morris, R. Huey, W. Lindstrom, M. F. Sanner, R. K. Belew, D. S. Goodsell and A. J. Olson, *J. Comput. Chem.*, 2009, **30**, 2785–2791.
- 177 M. J. Keiser, B. L. Roth, B. N. Armbruster, P. Ernsberger, J. J. Irwin and B. K. Shoichet, *Nature Biotechnology*, 2007, **25**, 197–206.
- 178 Y. Luo, X. Zhao, J. Zhou, J. Yang, Y. Zhang, W. Kuang, J. Peng, L. Chen and J. Zeng, *Nat. Commun.*, 2017, **8**, 573.
- 179 F. Wan, L. Hong, A. Xiao, T. Jiang and J. Zeng, *Bioinformatics*, 2019, **35**, 104–111.
- 180 A. S. Walker and J. Clardy, *J. Chem. Inf. Model.*, 2021, **61**, 2560–2571.
- 181 S. F. Altschul, W. Gish, W. Miller, E. W. Myers and D. J. Lipman, *J. Mol. Biol.*, 1990, **215**, 403–410.
- 182 T. Rodrigues, D. Reker, P. Schneider and G. Schneider, *Nat. Chem.*, 2016, **8**, 531–541.
- 183 J. Jeon, S. Nim, J. Teyra, A. Datti, J. L. Wrana, S. S. Sidhu, J. Moffat and P. M. Kim, *Genome Med.*, 2014, **6**, 57.
- 184 F. F. Hefti, *BMC Neurosci.*, 2008, **9 Suppl 3**, S7.
- 185 I. Rusyn and G. P. Daston, *Environ. Health Perspect.*, 2010, **118**, 1047–1050.
- 186 A. L. Demain, *J. Ind. Microbiol. Biotechnol.*, 2014, **41**, 185–201.
- 187 V. J. Sahayasheela, Z. Yu, Y. Hirose, G. N. Pandian, T. Bando and H. Sugiyama, *Bull. Chem. Soc. Jpn.*, 2022, **95**, 693–699.
- 188 G. Li, P. Lin, K. Wang, C.-C. Gu and S. Kusari, *Trends in Cancer*, 2021.
- 189 G. Brahmachari, *Discovery and Development of Therapeutics from Natural Products Against Neglected Tropical Diseases*, Elsevier, 2019.
- 190 U. Theuretzbacher, K. Outtersson, A. Engel and A. Karlén, *Nat. Rev. Microbiol.*, 2020, **18**, 275–285.
- 191 B. Terlain and T. Jean-Pierre, *Bull. Soc. Chim. Fr.*, 1971, **1971**, 2349–2356.
- 192 A. Kling, P. Lukat, D. V. Almeida, A. Bauer, E. Fontaine, S. Sordello, N. Zaburanyi, J. Herrmann, S. C. Wenzel, C. König, N. C. Ammerman, M. B. Barrio, K. Borchers, F. Bordon-Pallier, M. Brönstrup, G. Courtemanche, M. Gerlitz, M. Geslin, P. Hammann, D. W. Heinz, H. Hoffmann, S. Klieber, M. Kohlmann, M. Kurz, C. Lair, H. Matter, E. Nuernberger, S. Tyagi, L. Fraisse, J. H. Grosset, S. Lagrange and R. Müller, *Science*, 2015, **348**, 1106–1112.
- 193 F. Strelitz, H. Flon and I. N. Asheshov, *J. Bacteriol.*, 1955, **69**, 280–283.
- 194 S. K. Jain, A. S. Pathania, R. Parshad, C. Raina, A. Ali, A. P. Gupta, M. Kushwaha, S. Aravinda, S. Bhushan, S. B. Bharate and R. A. Vishwakarma, *RSC Adv.*, 2013, **3**, 21046–21053.
- 195 B. Muralikrishnan, V. M. Dan, J. S. Vinodh, V. Jamsheena, R. Ramachandran, S. Thomas, S. G. Dastager, K. Santhosh Kumar, R. S. Lankalapalli and R. A. Kumar, *RSC Adv.*, 2017, **7**, 36335–36339.
- 196 F. Wu, J. Zhang, F. Song, S. Wang, H. Guo, Q. Wei, H. Dai, X. Chen, X. Xia, X. Liu, L. Zhang, J.-Q. Yu and X. Lei, *ACS Cent Sci*, 2020, **6**, 928–938.
- 197 Z. Tanoli, M. Vähä-Koskela and T. Aittokallio, *Expert Opin. Drug Discov.*, 2021, **16**, 977–989.
- 198 N. B. Cech, M. H. Medema and J. Clardy, *Nat. Prod. Rep.*, 2021, **38**, 1947–1953.
- 199 J. A. van Santen, S. A. Kautsar, M. H. Medema and R. G. Linington, *Nat. Prod. Rep.*, 2021, **38**, 264–278.
- 200 A. Bender and I. Cortés-Ciriano, *Drug Discovery Today*, 2021, **26**, 511–524.
- 201 B. Y. Anom, *Ethics, Medicine and Public Health*, 2020, **15**, 100568.
- 202 S. Vatanserver, A. Schlessinger, D. Wacker, H. Ü. Kaniskan, J. Jin, M.-M. Zhou and B. Zhang, *Med. Res. Rev.*, 2021, **41**, 1427–1473.