

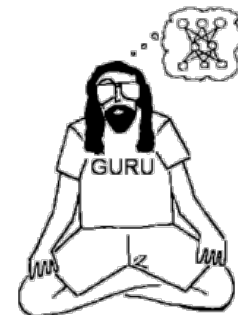
# SPECTRE: A Spectral Transformer for Molecule Identification

*A collaboration between the Gerwick & GURU labs*

Gary Cottrell, CSE

With Wangdong Xu (CSE), Byeol Ryu (SIO), Huanru Henry Mao (Calclavia),  
Hyunwoo Kim (Dongguk U.), James Zhao (CSE), Chen Zhang (SIO),  
Anthony Tong (CSE), Yiran Xu (CSE).  
and Bill Gerwick (SIO)





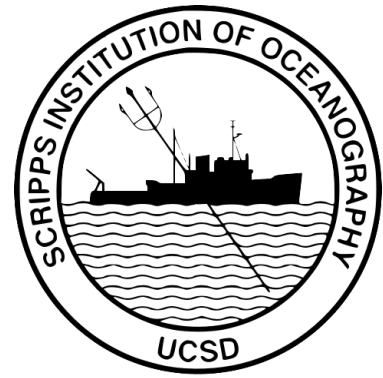
# SPECTRE: A Spectral Transformer for Molecule Identification

*A collaboration between the Gerwick & GURU labs*

Gary Cottrell, CSE

With **Wangdong Xu** (CSE), **Byeol Ryu** (SIO), Huanru Henry Mao (Calclavia),  
Hyunwoo Kim (Dongguk U.), James Zhao (CSE), Chen Zhang (SIO),  
Anthony Tong (CSE), Yiran Xu (CSE).  
and **Bill Gerwick** (SIO)





Bill Gerwick collecting a natural product

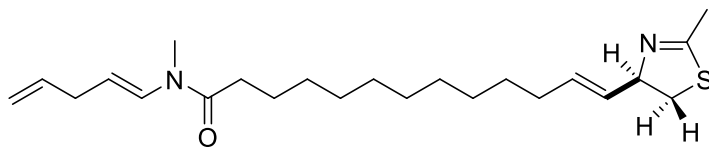
# Outline

- **Motivation**
- SMART/DeepSAT
- SPECTRE
  - Methods
  - Results
- Conclusion

# Motivation

- Natural products (NPs) comprise somewhere between 30-50% of drugs on the market
- The pipeline for novel NPs (after collection) starts with purification and structure determination
- 2D HSQC NMR is the preferred method for initial molecule structure elucidation
- But this interpretation step requires a great deal of human expertise, i.e., Bill Gerwick!

# 2D NMR is an Indicator of Compound Structure

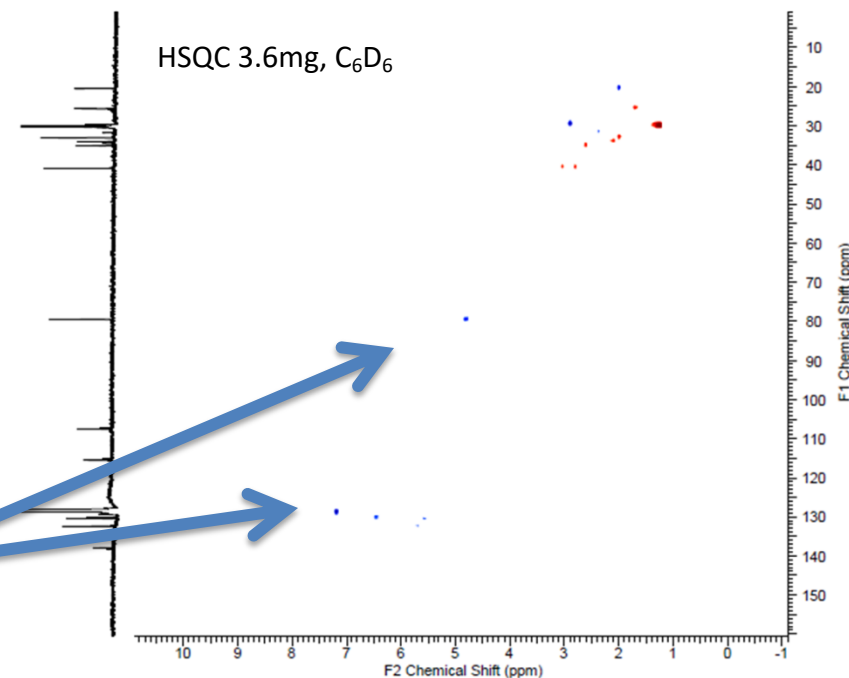
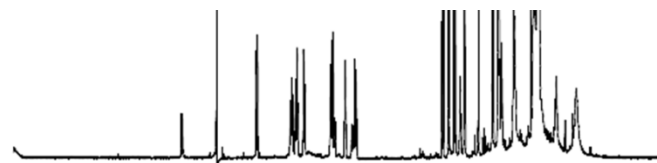


laucysteinamide A

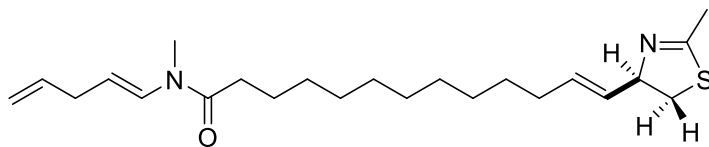
*Caldora penicillata*



Each “dot” here  
corresponds to a bond  
between a hydrogen  
and a carbon

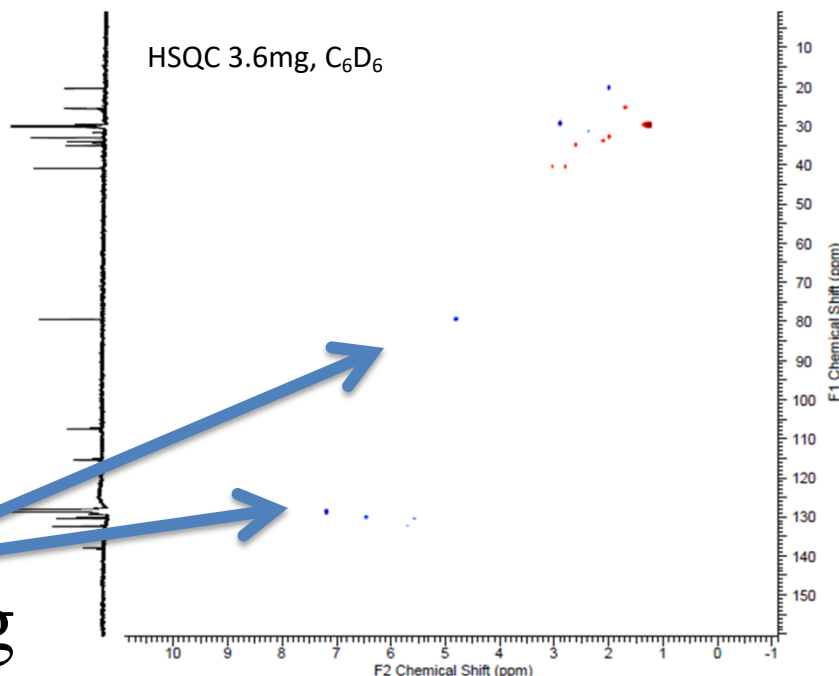
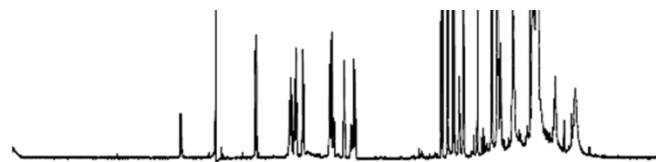


# 2D NMR is an Indicator of Compound Structure



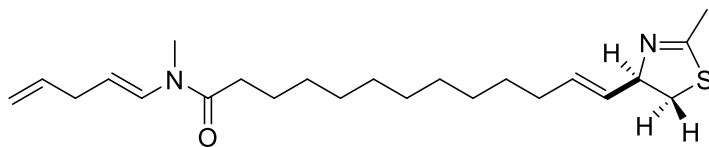
laucysteinamide A

*Caldora penicillata*



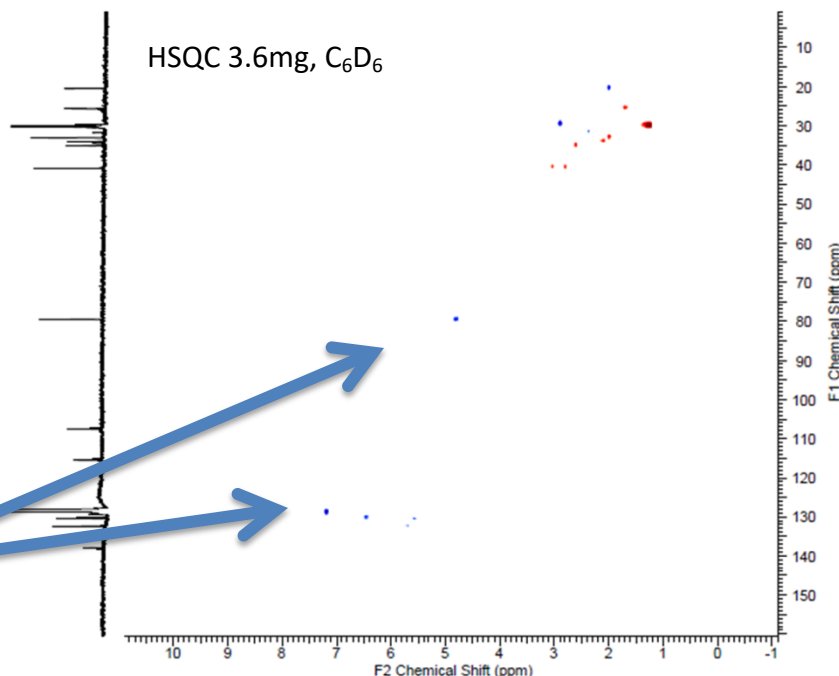
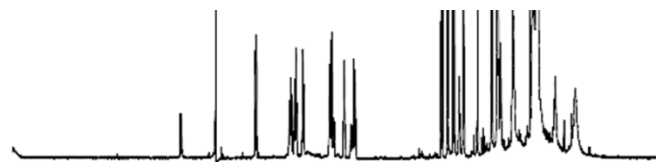
Where the dots appear  
depend on the neighboring  
atoms – this is called the  
“chemical shift”

# 2D NMR is an Indicator of Compound Structure



laucysteinamide A

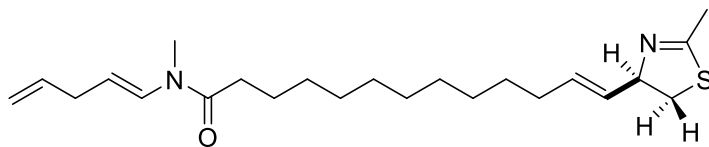
*Caldora penicillata*



And now I've told you  
everything I know about  
NMR!!! This is why it's  
good to have collaborators!



# 2D NMR is an Indicator of Compound Structure

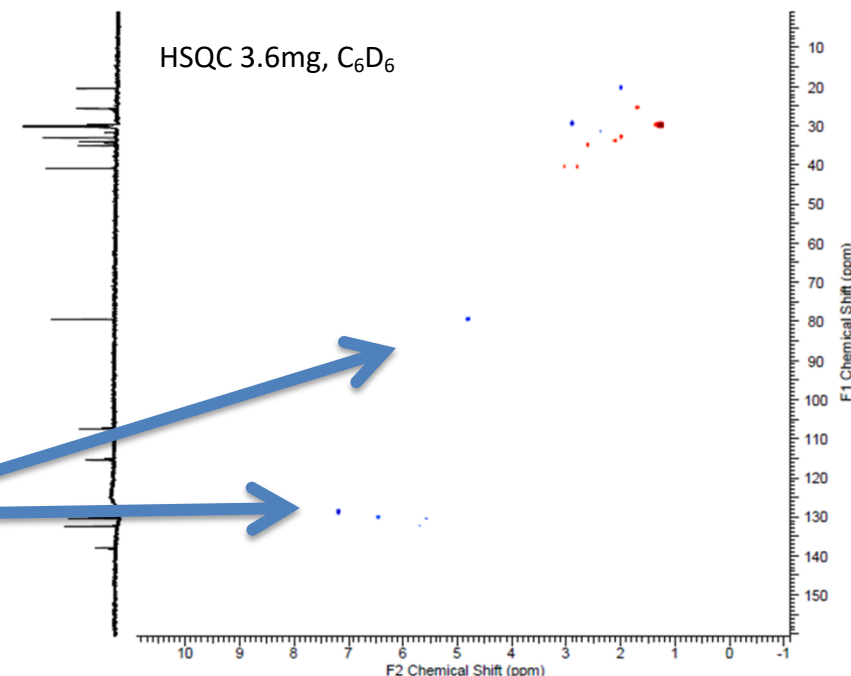
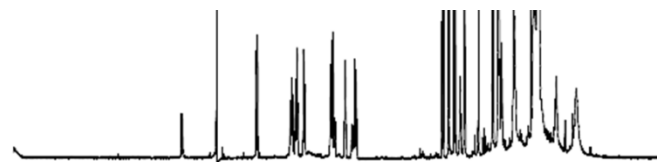


laucysteinamide A

*Caldora penicillata*



An expert (e.g., Bill Gerwick) can look at this and say “Ok, looks like we have a methyl group”



We think of this as “the face of the molecule”

# Outline

- Motivation
- **SMART/DeepSAT**
- SPECTRE
  - Methods
  - Results
- Conclusion

# Enter Deep Learning

- Our student Chen Zhang had the idea that what Bill was doing was like face recognition – so he came to me.
- We started the SMART project – treating the 2D NMR as an image, and using Convolutional Neural Networks to map that image to a cluster space where similar compounds had similar locations in the space.
- Given a new compound, nearby points in the space suggested possible structures

We created a sequence of better and better models...

# SCIENTIFIC REPORTS

JOURNAL OF  
NATURAL  
PRODUCTS

Cite This: *J. Nat. Prod.* 2020, 83, 617–625

Article

[pubs.acs.org/jnp](https://pubs.acs.org/jnp)

Received: 2

Accepted: 2

Published: 2

Pa  
Ch  
Ne

I | A | C | S

Kim et al. *Journal of Cheminformatics* (2023) 15:71  
<https://doi.org/10.1186/s13321-023-00738-4>

Journal of Cheminformatics

Yuey  
Che  
Tho

SOFTWARE

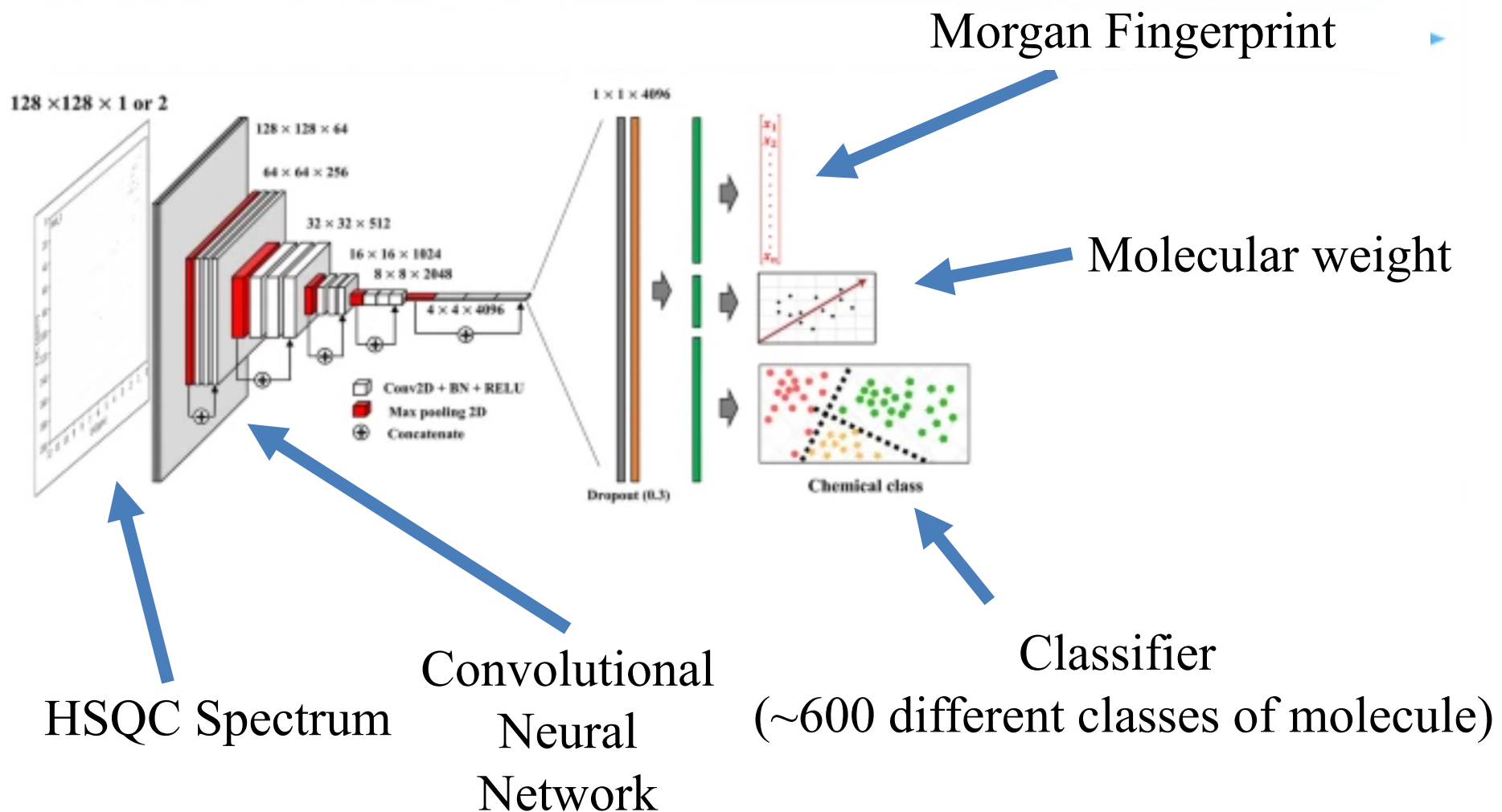
Open Access

## DeepSAT: Learning Molecular Structures from Nuclear Magnetic Resonance Data

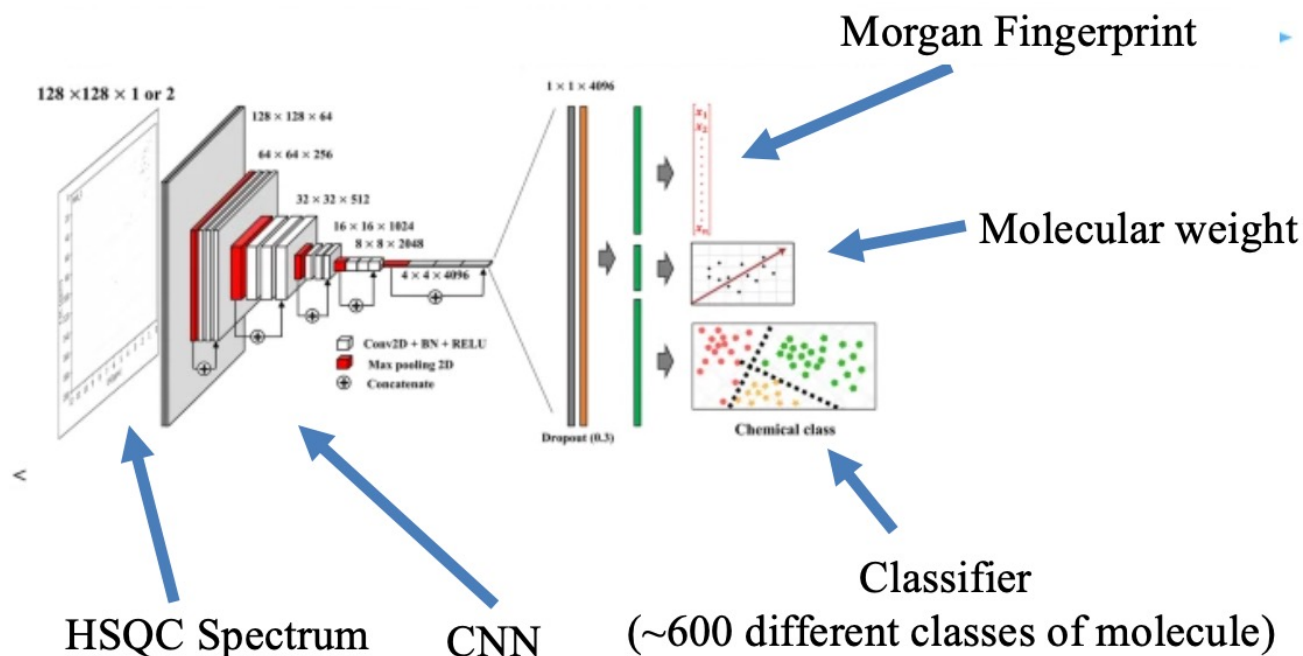


Hyun Woo Kim<sup>1,2</sup>, Chen Zhang<sup>1,3</sup>, Raphael Reher<sup>1,4</sup>, Mingxun Wang<sup>5,6,7</sup>, Kelsey L. Alexander<sup>1,8</sup>, Louis-Félix Nothias<sup>9</sup>, Yoo Kyong Han<sup>10</sup>, Hyeji Shin<sup>10</sup>, Ki Yong Lee<sup>1,10</sup>, Kyu Hyeong Lee<sup>2</sup>, Myeong Ji Kim<sup>2</sup>, Pieter C. Dorrestein<sup>5</sup>, William H. Gerwick<sup>1,5\*</sup> and Garrison W. Cottrell<sup>3\*</sup>

# DeepSAT Network architecture: Supervised multi-task CNN



# DeepSAT Network architecture: Supervised multi-task CNN

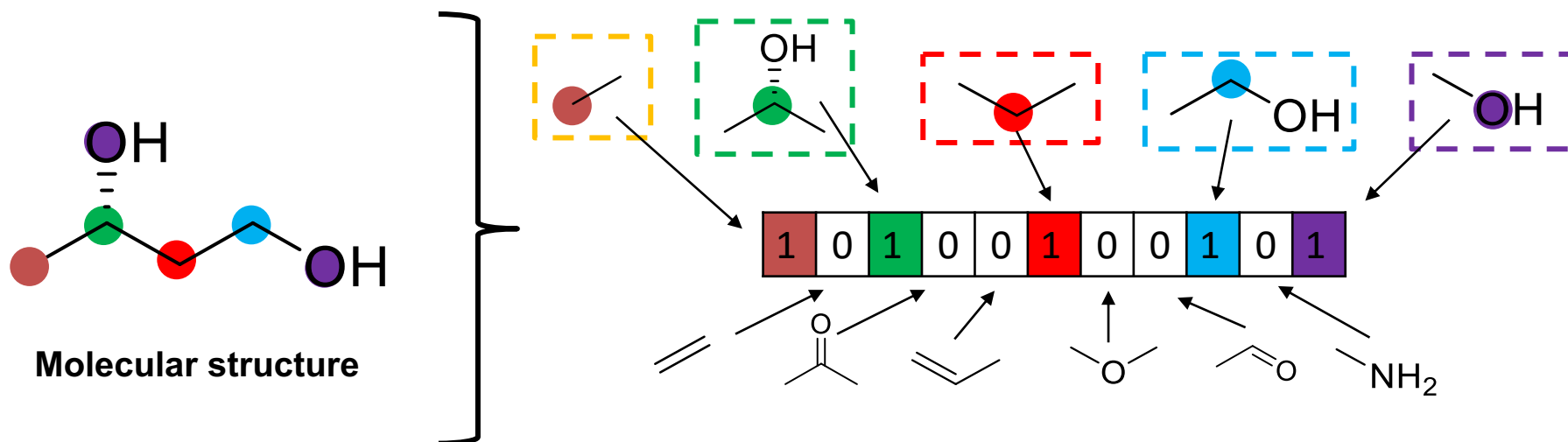


Main output: Morgan Fingerprints, a vector-based representation of molecular structure

These are compared to a database of MFs, and a list of similar molecules are returned

# Morgan fingerprints

- Morgan Fingerprints are a vector-based representation of molecular structure used in many computational tools for cheminformatics

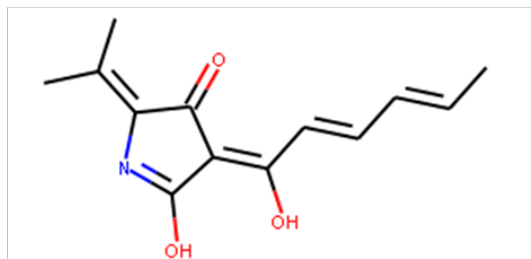


- DeepSAT uses a 6,144 bit vector, with each position associated with a specific partial structure
- A “1” means that this substructure is in the molecule,
- A “0” means it’s absent

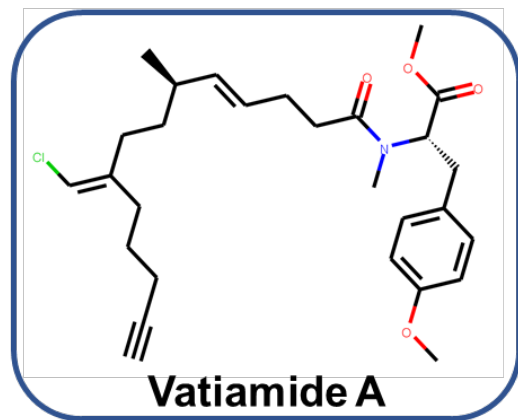
# DeepSAT – Using Morgan-type Fingerprints to Determine Chemical Similarity

## Predicted Fingerprint from DeepSAT

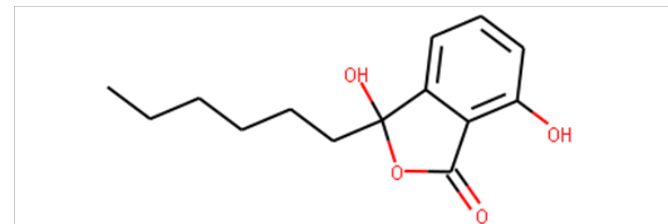
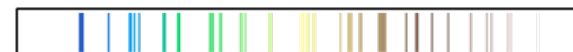
Similarity: 0.29



**Similarity: 0.74**



Similarity: 0.36



Hyunwoo Kim

This result gives the researcher clues to the chemical structure of a novel compound – speeding structure identification



# Outline

- Motivation
- SMART/DeepSAT
- **SPECTRE**
  - Methods
  - Results
- Conclusion

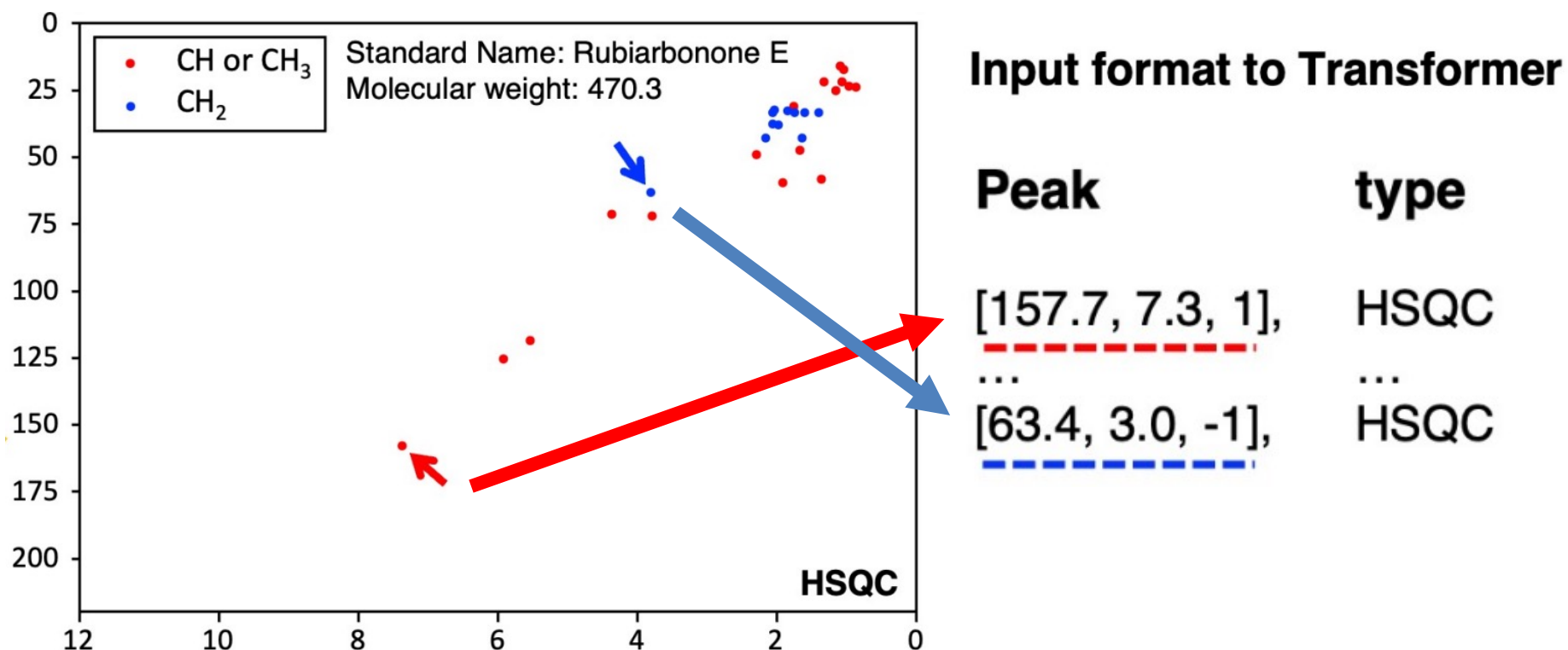
# What's wrong with this picture?

- These models have over 8k users and over 450k queries.
- But they're limited and inefficient for three reasons:
  - They only take 2D HSQC NMR as inputs
    - Inflexible
  - Over 99% of the pixels are zero!
    - A lot of wasted compute.
  - The target, Morgan Fingerprints, are a hash table
    - Collisions: Locations in the table are *ambiguous*

# The First Main Idea of this talk: Transformers

- A transformer is what underlies ChatGPT:
  - It takes words as input
  - It processes those words in the context of other words to extract meaning
- Instead of words, we give it the only the (x,y) locations of the peaks
  - much more efficient (no wasted computation)
  - It processes the peaks in the context of the other peaks
- Just like DeepSAT, we train it to produce Morgan Fingerprints from thousands of examples

# The First Main Idea of this talk: Transformers

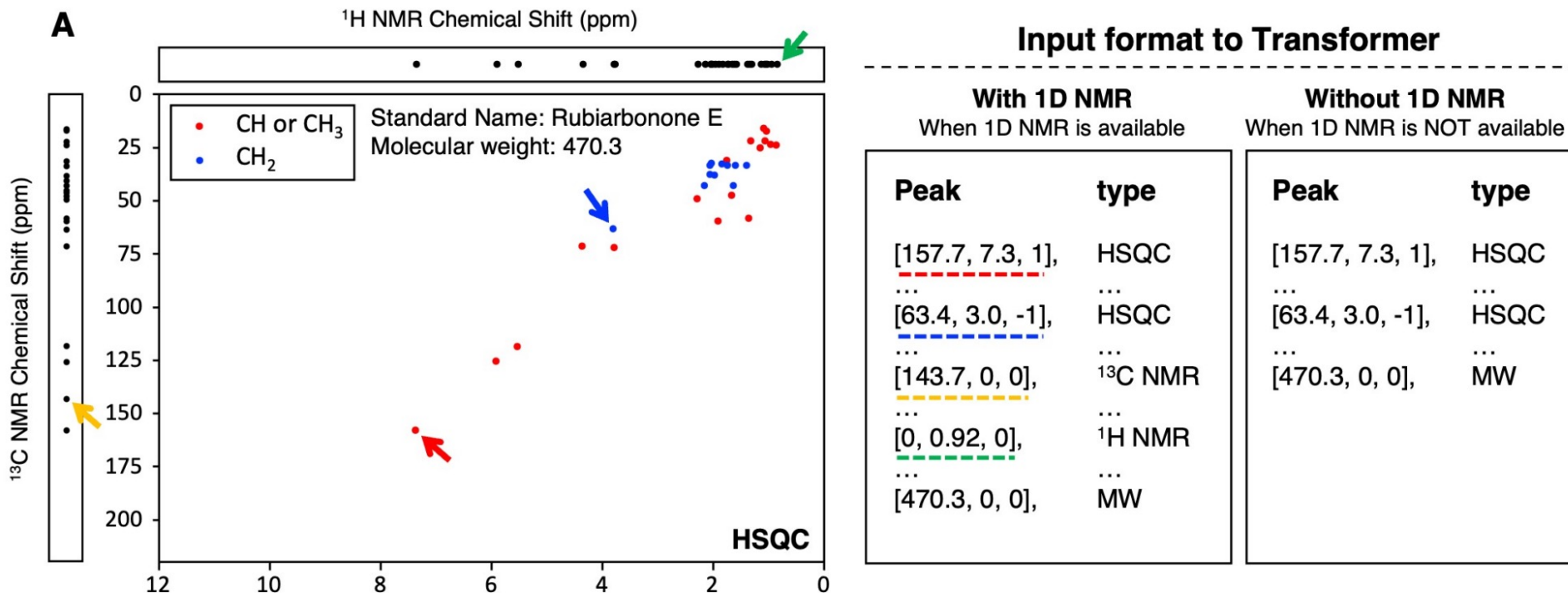


Peaks are encoded by their chemical shift as well as their multiplicity - odd or even # of bonds

# The Second Main Idea of this talk: Flexible Inputs

- A transformer is what underlies ChatGPT:
  - It takes words as input
  - It processes those words in the context of other words to extract meaning
- Instead of words, we give it the (x,y) locations of the peaks
  - much more efficient!
- We can *also* give it other data, tagged by its type:
  - 1D  $^{13}\text{C}$  NMR peaks:  $(\text{C}, x_1), (\text{C}, x_2), (\text{C}, x_3), \dots (\text{C}, x_{\text{NC}})$
  - 1D  $^1\text{H}$  NMR peaks:  $(\text{H}, x_1), (\text{H}, x_2), (\text{H}, x_3), \dots, (\text{H}, x_{\text{NH}})$
  - Molecular weight:  $(\text{MW}, 470.3)$
- We train it by *randomly choosing what data types to give it as input* - over many training trials, *it learns to use whatever data is available*

# The Second Main Idea of this talk: Flexible Inputs



- Now we can give it *multiple data types* – 2D NMR, 1D NMR, Molecular weight
- We train it by randomly giving it different types of data for each example – one time, it might just have 2D HSQC and 1D Carbon NMR, other times, just 1D Carbon, etc.

# The Third Main Idea of this talk:

## Better Morgan Fingerprints

- We created Morgan Fingerprints up to radius 9 over a large dataset of molecules
- We sort them by their entropy, keeping the bits with the most information and label them with the substructure
- We keep 16,384 of these bits in the vector
- In our hands, these are collision-free: every bit in the vector corresponds to a unique substructure
  - So we can *label* which parts of the retrieved structures match and which don't

# Network Architecture

Input Data

```
<HSQC_start>  
[5.48, 81.28, 1]  
[2.37, 38.00, -1]  
.....  
[2.27, 44.60, -1]  
.....  
<HSQC_end>  
<H-NMR_start>  
[5.21, 0, 0]  
[4.37, 0, 0]  
.....  
<H-NMR_end>  
<C-NMR_start>  
[138.48, 0, 0]  
[81.52, 0, 0]  
.....  
<C-NMR_end>  
<Mol-Weight_start>  
[502.726, 0, 0]  
<Mol-Weight_end>
```

Transformer  
Encoder



Transformer  
Encoder



CLS Token



Multi-layer  
Perceptron

Final Layer  
Representation

```
encoded token 1  
encoded token 2  
encoded token 3  
encoded token 4  
.  
.  
.  
.  
.  
.  
encoded token n
```

16,384-D Morgan Fingerprint  
[1 0 0 1 1 1 0 ... 0 1 1]

Instead of predicting the next word, like ChatGPT, it is trained to predict our super-duper Morgan Fingerprint

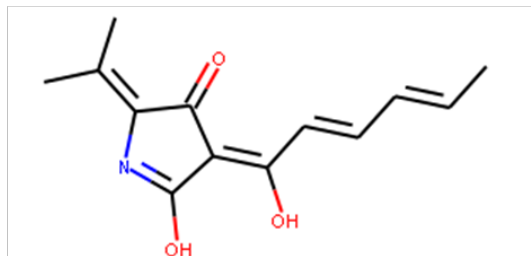


Like DeepSAT, SPECTRE uses the super-duper Morgan Fingerprint to find similar molecules

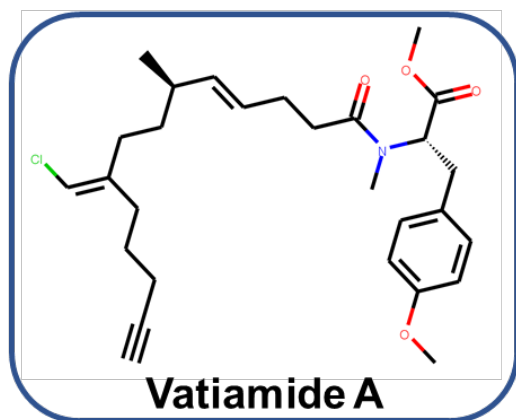
Our retrieval set is very large - over **520,000 NP candidates**

### Predicted Fingerprint from SPECTRE

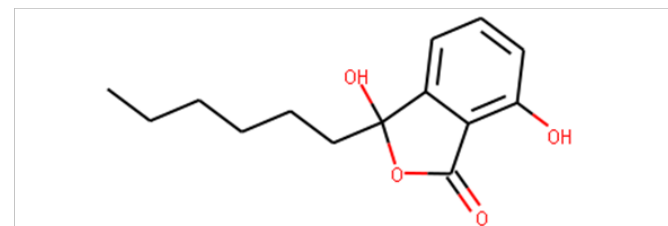
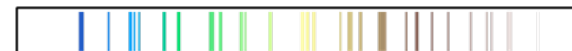
Similarity: 0.29



Similarity: **0.74**



Similarity: 0.36



Hyunwoo Kim

The result is a list of molecules ordered by their similarity to the predicted fingerprint – speeding structure identification

# Data Collection & Training, Test, and Validation Set:

We collected a LOT of data from publicly available datasets, as well as *predicting* data to train on

## Data Collection



### CH-NMR-NP

Literature NMR spectra  
29,500 NPs + 6,000 organic  
compounds



### ACD/Labs

Computed NMR spectra  
ACD/Labs 113,967



1D NMR spectra ( $^1\text{H}$  &  $^{13}\text{C}$ )  
Natural Products 155,815

**1D NMR (n = 155,815)**  
chemical names, SMILES strings, and molecular weight



**Retrieval DB**  
526,316 NPs

## Training Set

Model Input	SPECTRE (with DTD) occurrence rate during training
All 3 NMR spectra types	7.6%
$^{13}\text{C}$ NMR and $^1\text{H}$ NMR	12.2%
HSQC and $^{13}\text{C}$ NMR	7.6%
HSQC and $^1\text{H}$ NMR	7.6%
Only $^{13}\text{C}$ NMR	12.2%
Only $^1\text{H}$ NMR	12.2%
Only HSQC	40.4%

Model Input	Specialized Model (w/o DTD) training set size
All 3 NMR spectra types	33,203
$^{13}\text{C}$ NMR and $^1\text{H}$ NMR	93,370
HSQC and $^{13}\text{C}$ NMR	39,685
HSQC and $^1\text{H}$ NMR	33,203
Only $^{13}\text{C}$ NMR	103,409
Only $^1\text{H}$ NMR	93,374
Only HSQC	109,694

### Test Set

n = 4,096

### Validation Set

n = 4,056

Chosen when all three NMR spectra  
( $^1\text{H}$ ,  $^{13}\text{C}$ , and HSQC) were available

# Results, quantitative

We compared SPECTRE against specialized models trained only on one type of data; here, we're measuring when the correct molecule is the top hit.

(in practice, we can apply the best model in each case)

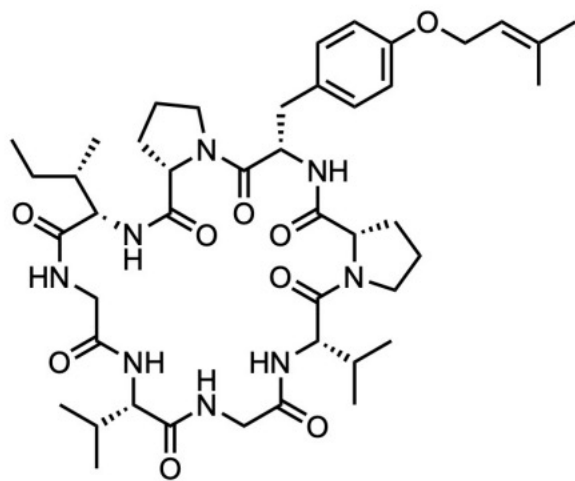
Model Input Type	Top-1	
	SPECTRE	Specialized
ME-HSQC, $^{13}\text{C}$ NMR, and $^1\text{H}$ NMR	<b>79.78%<math>\pm</math>0.82%</b>	72.32 % $\pm$ 0.19%
Standard HSQC, $^{13}\text{C}$ NMR, and $^1\text{H}$ NMR	<b>78.41%<math>\pm</math>0.77%</b>	69.61 % $\pm$ 0.36%
ME-HSQC and $^{13}\text{C}$ NMR	<b>79.81%<math>\pm</math>0.47%</b>	73.14 % $\pm$ 0.33%
Standard HSQC and $^{13}\text{C}$ NMR	<b>77.98%<math>\pm</math>0.76%</b>	69.57 % $\pm$ 0.42%
ME-HSQC and $^1\text{H}$ NMR	<b>75.65%<math>\pm</math>0.82%</b>	65.50 % $\pm$ 1.95%
Standard HSQC and $^1\text{H}$ NMR	<b>72.83%<math>\pm</math>0.98%</b>	61.88 % $\pm$ 2.01%
ME-HSQC	74.25% $\pm$ 0.79%	<b>76.52 % <math>\pm</math>0.40%</b>
Standard HSQC	70.20% $\pm$ 0.90%	<b>72.83 % <math>\pm</math>0.23%</b>
$^{13}\text{C}$ NMR and $^1\text{H}$ NMR	59.36% $\pm$ 0.90%	<b>68.91 %<math>\pm</math>0.65%</b>
$^{13}\text{C}$ NMR	51.96% $\pm$ 0.73%	<b>57.59 %<math>\pm</math>0.26%</b>
$^1\text{H}$ NMR	15.38% $\pm$ 0.25%	<b>19.63 %<math>\pm</math>0.46%</b>

# Results, qualitative

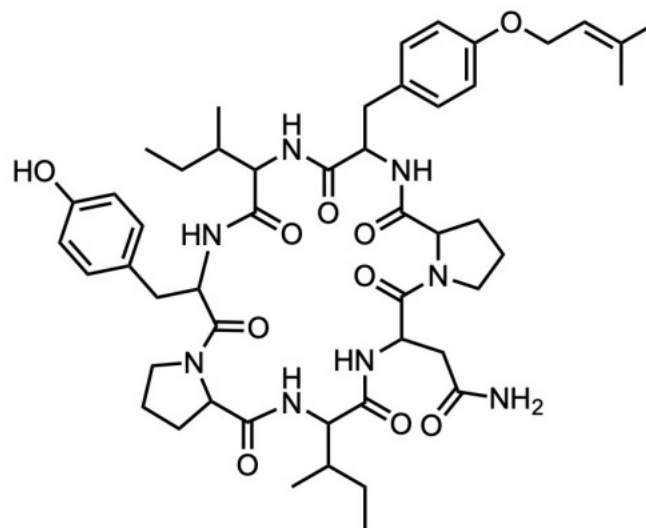
Here, using Multiplicity-Edited HSQC as input

## A. Multiplicity-Edited HSQC

---



Monchicamide I  
(input)



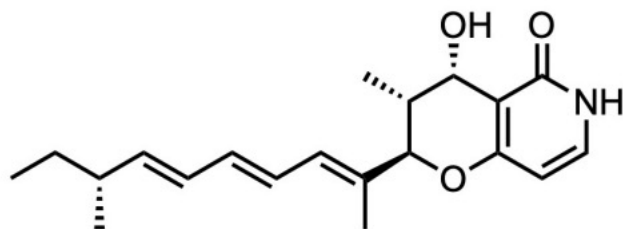
Closest retrieved molecule

# Results, qualitative

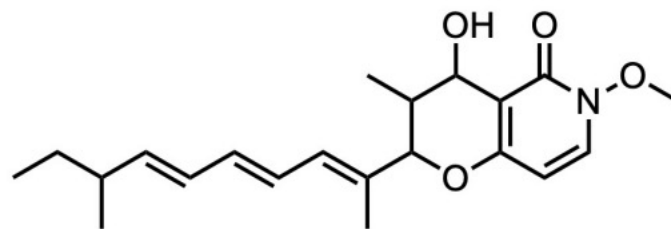
Using only Standard HSQC as input

## B. Standard HSQC

---



Aculeapuridone A  
(input)



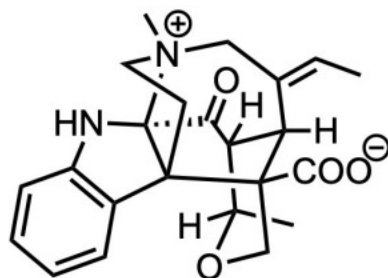
Closest retrieved molecule

# Results, qualitative

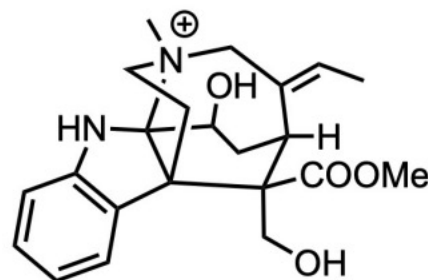
Here, using only  $^{13}\text{C}$  NMR as input

## C. $^{13}\text{C}$ NMR

---



Alstolarsine A  
(input)



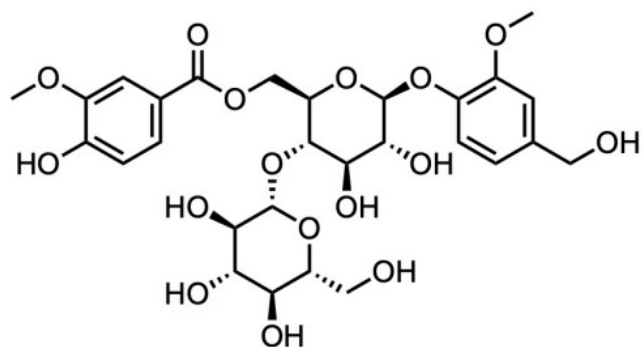
Closest retrieved molecule

# Results, qualitative

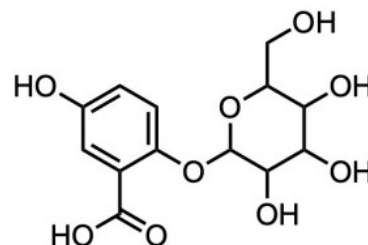
Here, using only  $^1\text{H}$  NMR as input  
(doesn't work well)

## D. $^1\text{H}$ NMR

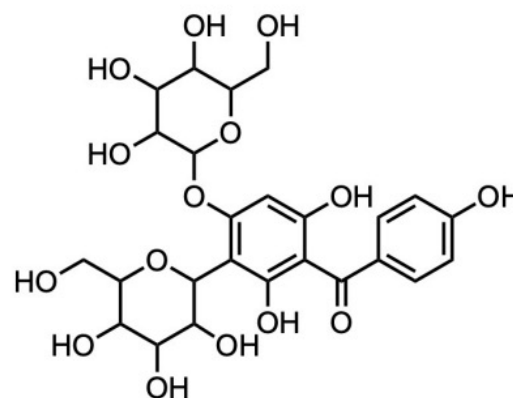
---



Wrightioside A  
(input)



1<sup>st</sup> Closest  
retrieved  
molecule

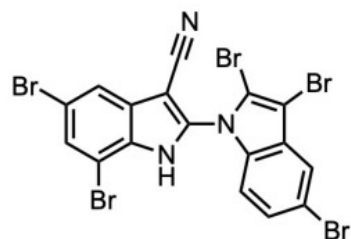


2<sup>nd</sup> Closest  
retrieved  
molecule

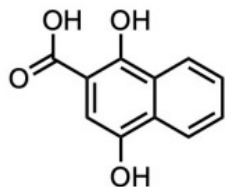
# Results, qualitative

Proton-deficient compounds:

**The power of multiple input data types**



Aetokthonotoxin  
(input)



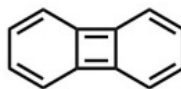
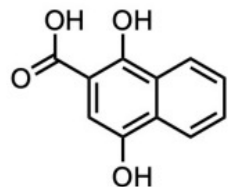
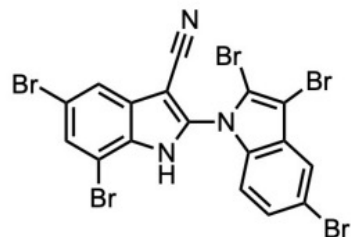
Using  
only  
ME-  
HSQC



# Results, qualitative

## Proton-deficient compounds:

# The power of multiple input data types



Aetokthonotoxin  
(input)

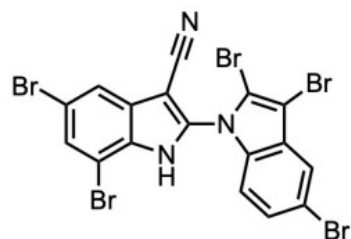
# Using only ME- HSQC

# Using only $^{13}\text{C}$

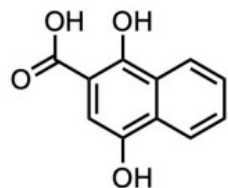
# Results, qualitative

Proton-deficient compounds:

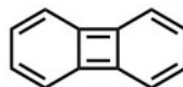
**The power of multiple input data types**



Aetokthonotoxin  
(input)



Using  
only  
ME-  
HSQC



Using  
only  
 $^{13}\text{C}$

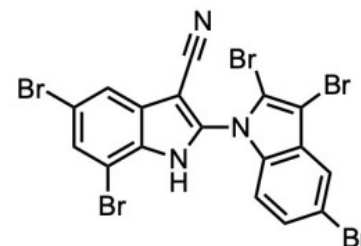
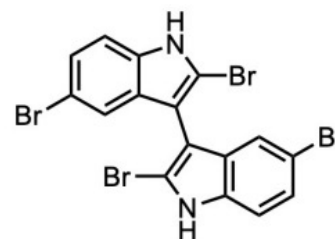
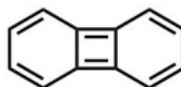
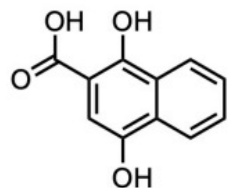
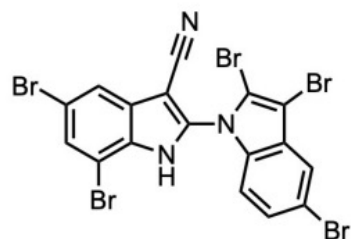


Using  
ME-  
HSQC  
+  $^{13}\text{C}$

# Results, qualitative

Proton-deficient compounds:

**The power of multiple input data types**



Aetokthonotoxin  
(input)

Using  
only  
ME-  
HSQC

Using  
only  
 $^{13}\text{C}$

Using  
ME-  
HSQC  
+  $^{13}\text{C}$

Using  
ME-  
HSQC  
+  $^{13}\text{C}$   
+ MW

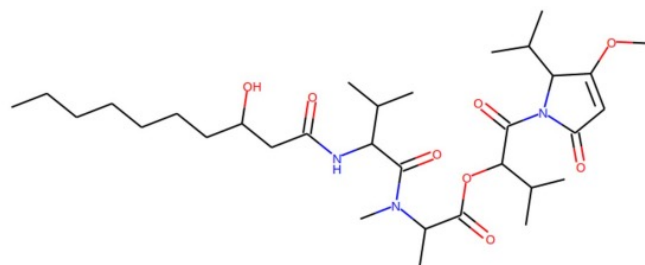
A perfect match – the molecule has been “dereplicated”

# The advantage of unique bits in our entropy-optimized Morgan Fingerprints

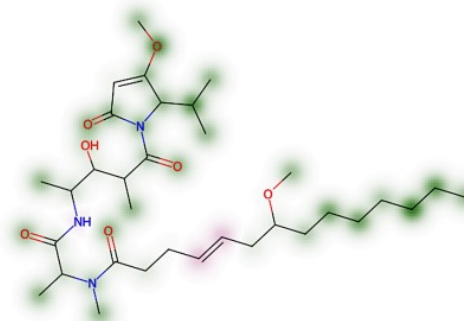
## We can mark what matches and what doesn't

### A. Kavaratamide A - Standard HSQC

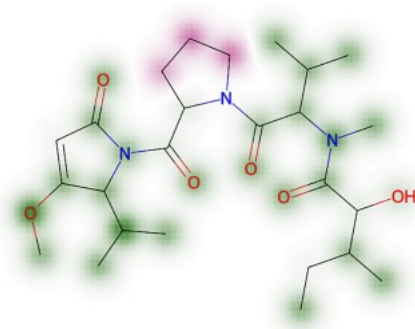
---



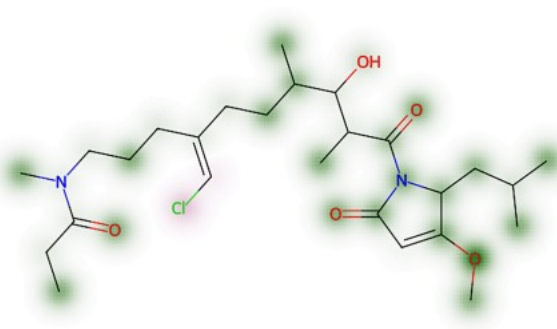
**Kavaratamide A**



**Retrieval 1**



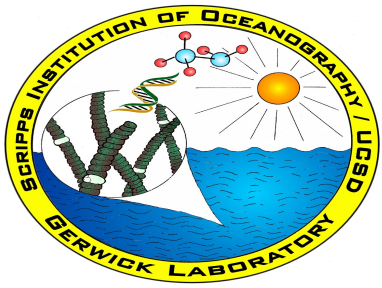
**Retrieval 2**



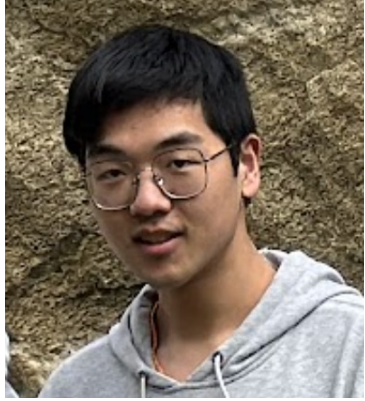
**Retrieval 4**

# Conclusions

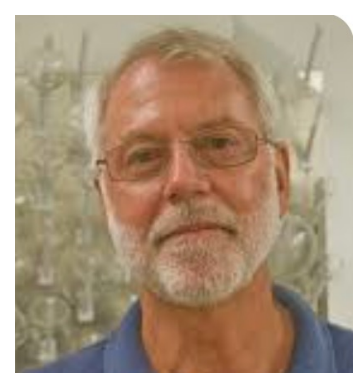
- SPECTRE is a major advance over our previous work
- It automagically combines data from multiple sources to obtain the best result, given the amount of data provided
- It uses an advanced form of Morgan Fingerprint – we call “entropy-optimized Morgan Fingerprints”
- These allow highlighting of matching substructures – providing more information than just – hey, this is similar!
- It beats our previous state of the art model by a lot!



**Thanks!!**  
**To these very smart folks!**



Wangdong Xu    Byeol Ryu    Henry Mao    Hyunwoo Kim    James Zhao



Chen Zhang    Anthony Tong    Yiran Xu    Ming Wang    Bill Gerwick

And thank YOU for listening!

Questions?