

# Accurate MS/MS Spectral Prediction with CFM-ID 4.0



Fei Wang<sup>1,3</sup>, Siyang Tian<sup>1,2</sup>, Jaanus Liggand<sup>2,4</sup>, David Arndt<sup>2</sup>, Russell Greiner<sup>1,3</sup>, David S. Wishart<sup>1,2</sup>

<sup>1</sup>Department of Computing Science, <sup>2</sup>Department of Biological Science, University of Alberta, Edmonton, Alberta, Canada

<sup>3</sup>Alberta Machine Intelligence Institute (AMII), Edmonton, Alberta, Canada

<sup>4</sup>Institute of Chemistry, University of Tartu, Tartu, Estonia

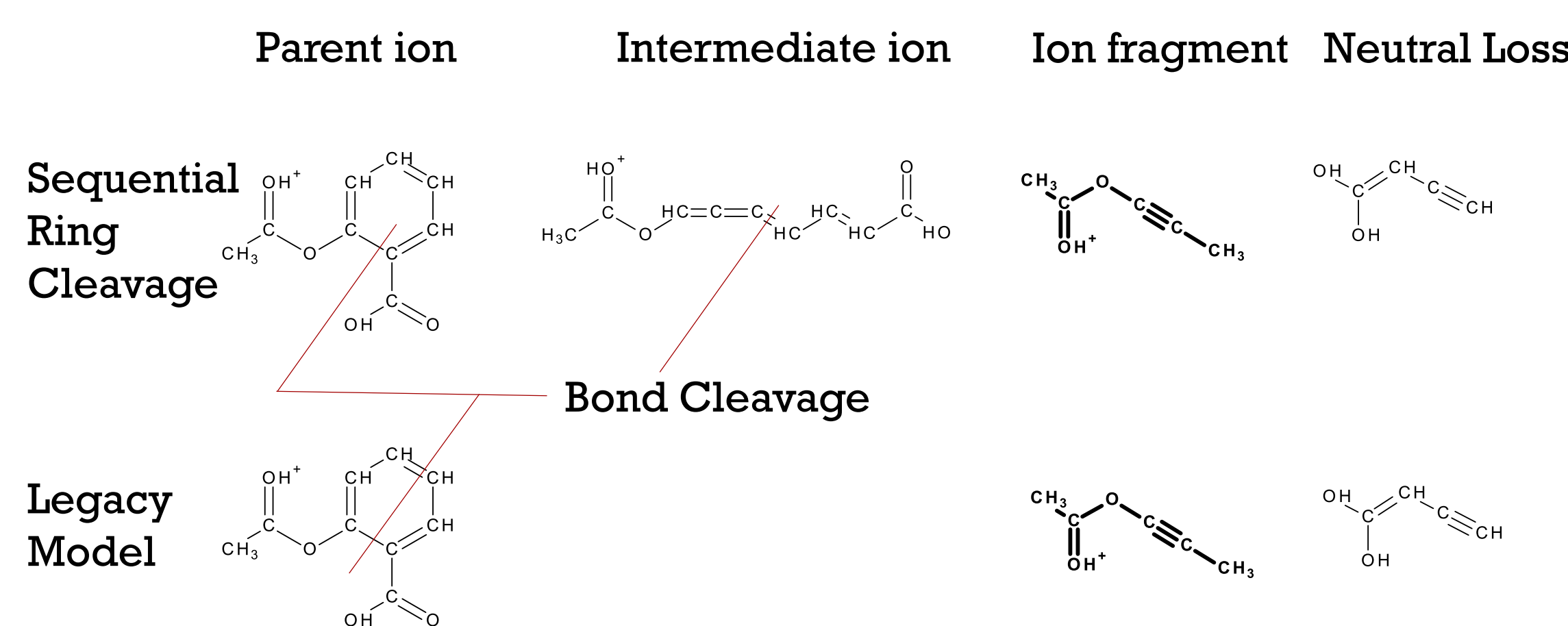
## Introduction

CFM-ID is a software package that uses machine learning to accurately predict the Electrospray Ionization Tandem Mass Spectrometry (ESI-MS/MS) spectra for organic compounds. It can perform the following three tasks: (1) **Predict the spectra** for a given chemical structure, (2) **annotate the peaks** in a set of given spectra of a known chemical structure, and (3) **classify the structure** for a target spectrum. CFM-ID 3.0 is a state-of-the-art program for ESI-MS/MS spectral prediction, as well as spectrum-to-compound classification. In this work, we are introducing the significantly improved version, CFM-ID 4.0. It is freely available as both, a web version (cfmid4.wishartlab.com), and a downloadable software package.

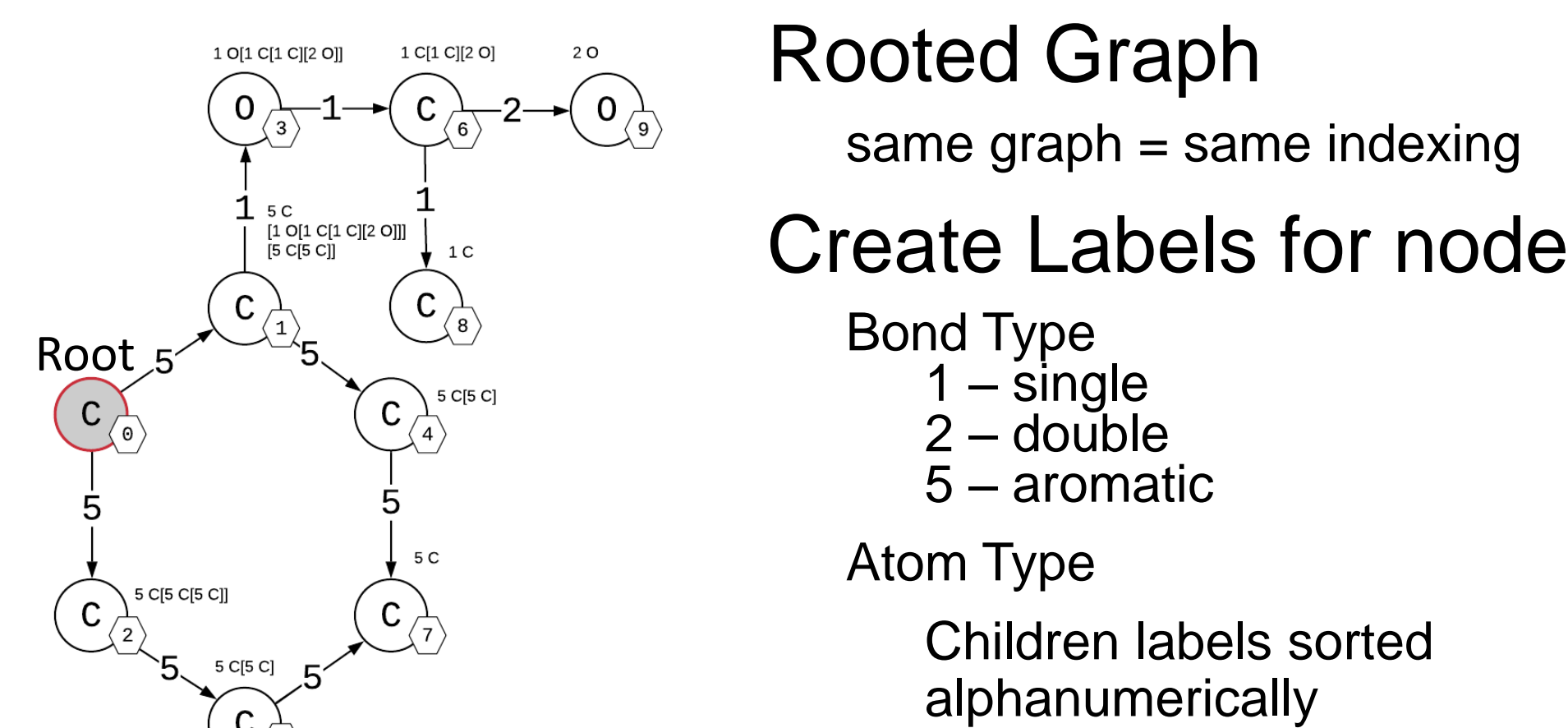
## CFM-ID 4.0 Improvements

- Learned parameters via Deep Neural Networks
- Learned parameters from molecular topology
- Improved ring cleavage model
- Expanded training set by 3x (QToF)
- Added rule-based schema for
  - Acylcarnitines
  - Acylcholines
  - Polyphenols
  - Flavonols

## Ring Cleavage Modeling



## Graph Based Feature Representation



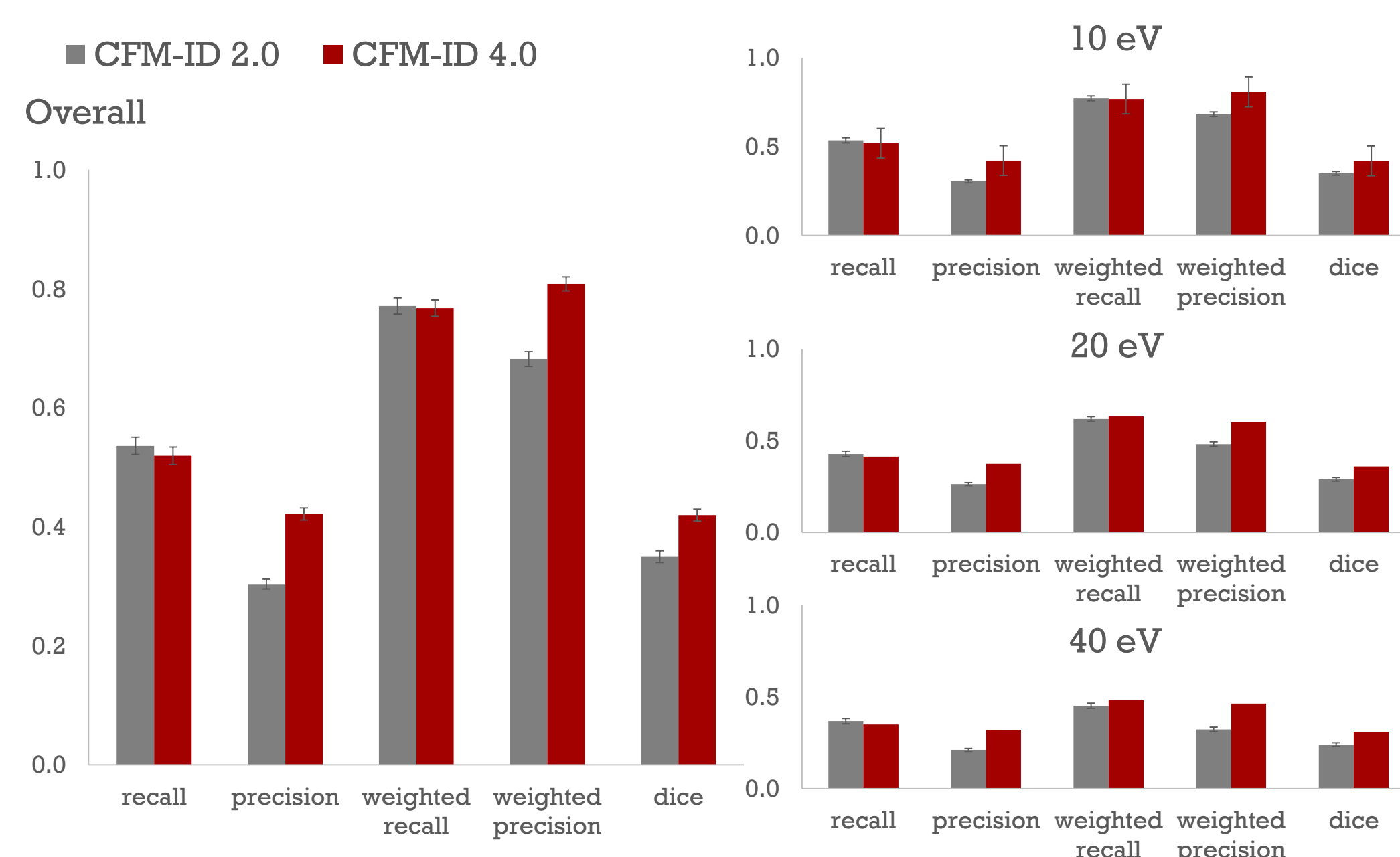
## Expanding Rule Based Predictors

Compound class	Number of Covered Rules	Covered Adduct Types
acylcarnitines	9	[M] <sup>+</sup>
acylcholines	11	[M] <sup>+</sup>
flavonols	47	[M+H] <sup>+</sup> , [M-H] <sup>-</sup>
flavones	20	[M+H] <sup>+</sup> , [M-H] <sup>-</sup>
flavanones	15	[M+H] <sup>+</sup> , [M-H] <sup>-</sup>
flavonoid-3-O-glycosides	36	[M+H] <sup>+</sup> , [M-H] <sup>-</sup>
flavonoid-7-O-glycosides	63	[M+H] <sup>+</sup> , [M-H] <sup>-</sup>
flavonoid-7-O-glucuronides	12	[M+H] <sup>+</sup> , [M-H] <sup>-</sup>
4'-O-methylated flavonoids	53	[M+H] <sup>+</sup> , [M-H] <sup>-</sup>
7-O-methylated flavonoids	36	[M+H] <sup>+</sup> , [M-H] <sup>-</sup>
3'-O-methylated flavonoids	11	[M-H] <sup>-</sup>

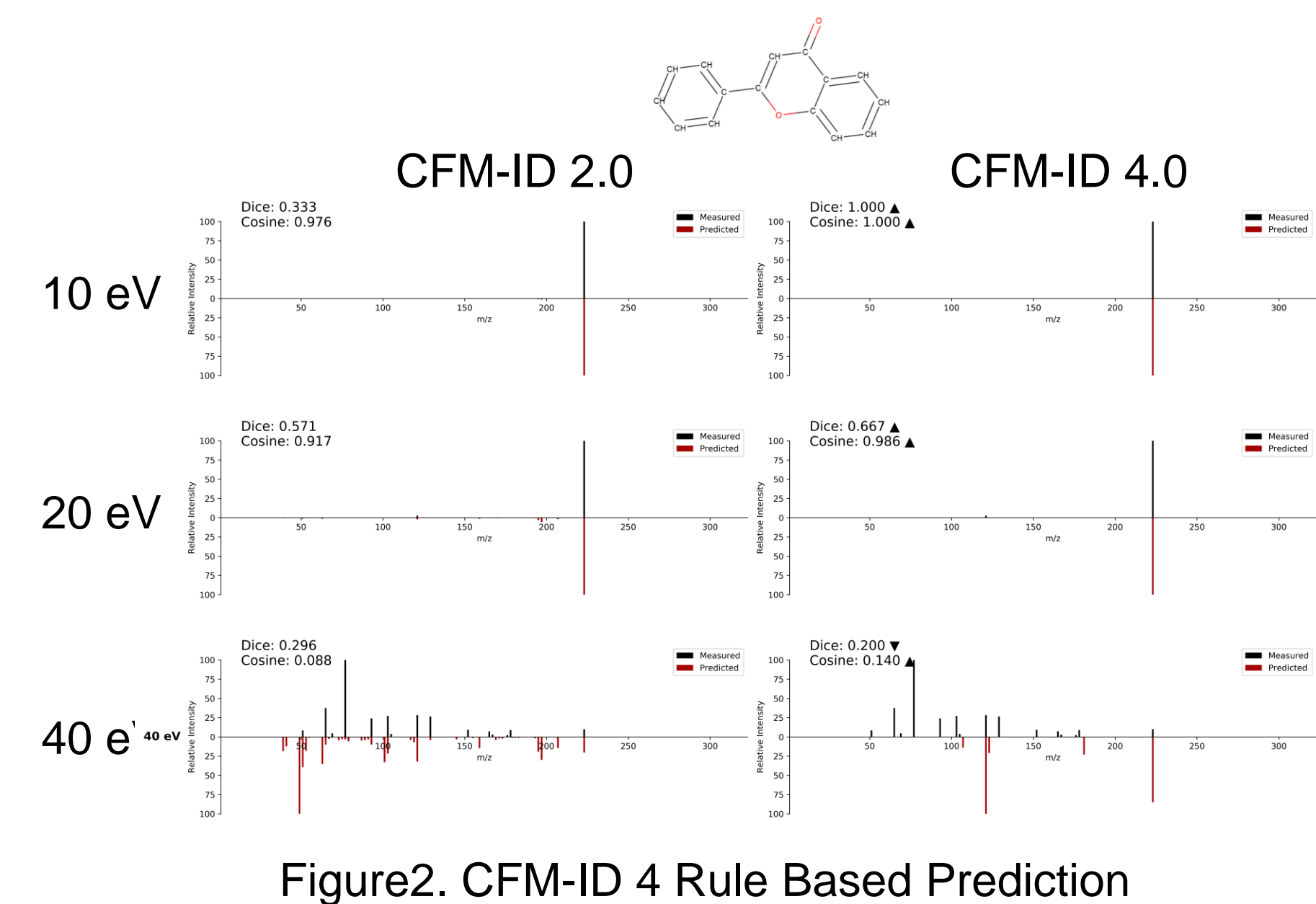
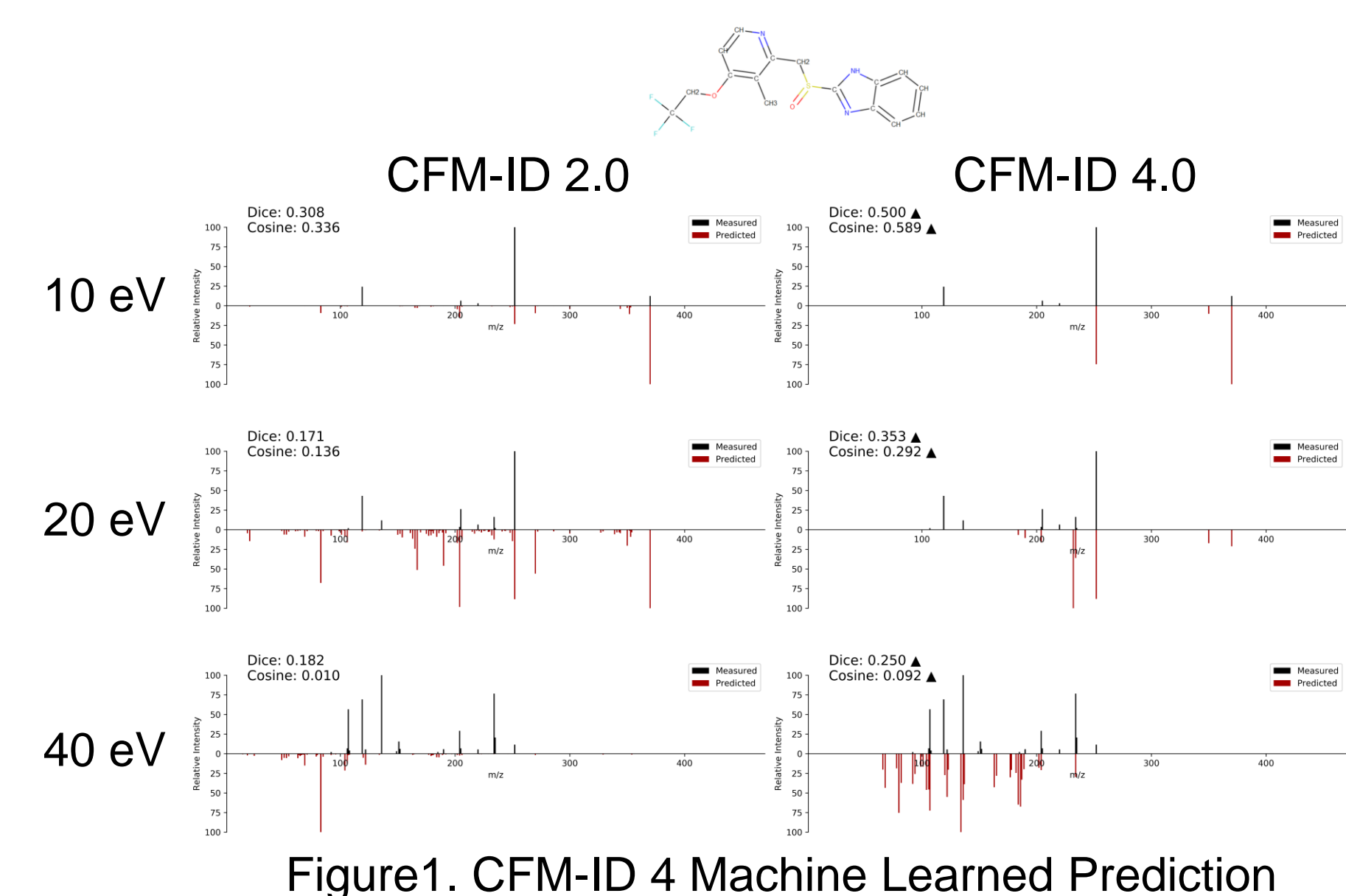
## Results

### Spectral Prediction Performance

10-Fold Cross Validation Results on Metlin 2015 [M+H]<sup>+</sup> Set



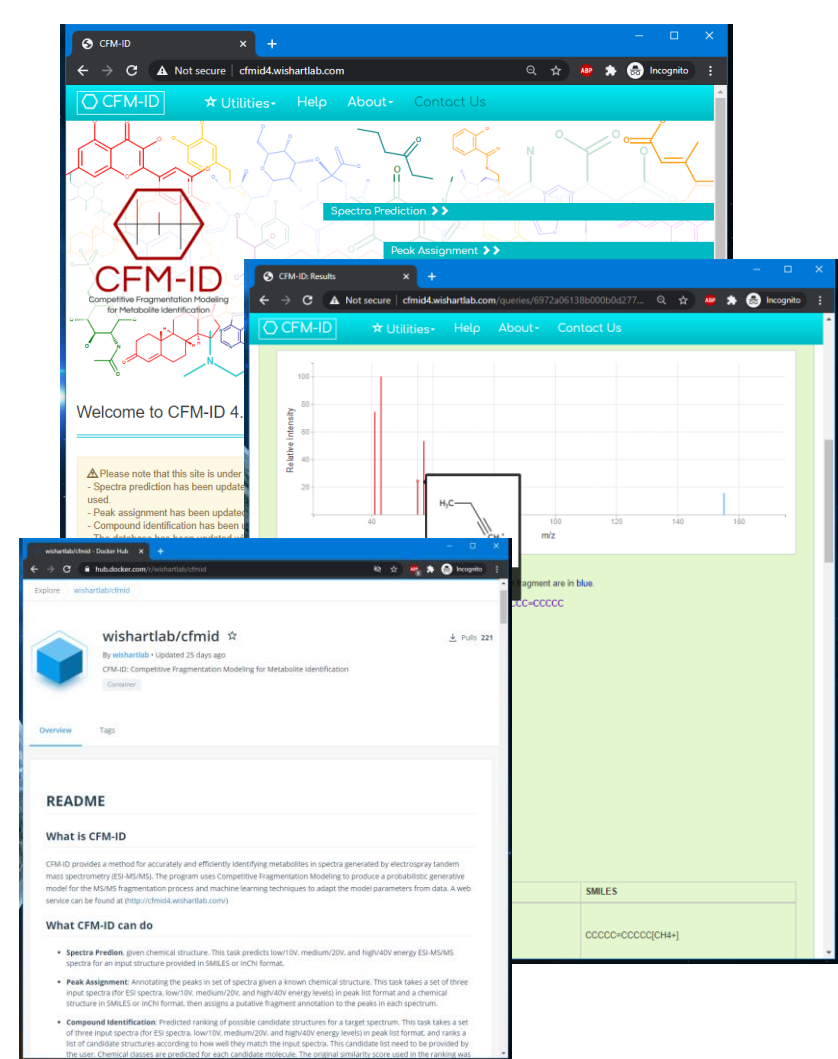
## Predicted Spectra



## Conclusions

In this work, we introduced a novel tensor representation for describing chemical structures and used it to extend the capabilities of existing CFM-ID machine learning methods in ESI-MS/MS spectral prediction tasks. Alongside machine learning based improvements, we proposed new rule-based methods to further enhance the CFM-ID 4.0 system's predictive ability for specific classes of chemicals, where the machine learning model suffers. The *in-silico* spectra prediction performance of these novel methods was examined against empirical results on multiple ESI-MS/MS data sets, encompassing a wide range of chemical classes, in both positive and negative ionization modes. While still imperfect, our proposed method outperformed the legacy CFM-ID model by a significant margin across all data sets. In addition, we demonstrated CFM-ID 4.0's *in-silico* compound identification ability via the CASMI 2016 competition (category 3), where CFM-ID 4.0 achieved better identification results than all existing approaches.

## CFM-ID Webservice



Web Service:  
<http://cfmid4.wishartlab.com>  
Docker Hub Access:  
<https://hub.docker.com/r/wishartlab/cfmid>  
Source Code:  
<https://bitbucket.org/wishartlab/cfmid-code/>  
<https://bitbucket.org/wishartlab/m-srb-fragmenter>

## Acknowledgements

