

Improving Understanding of Dissolved Organic Matter by Using Machine Learning to Predict Stable Carbon Isotope based on Molecular Abundances

DR. YUANBI YI¹, TONGCUN LIU², JULIAN MERDER³,
CHEN HE⁴, HONGYAN BAO⁵, PENGHUI LI⁶, SI-LIANG LI⁷,
QUAN SHI⁴ AND DING HE¹

¹Department of Ocean Science and the Southern Marine Science and Engineering Guangdong Laboratory (Guangzhou), The Hong Kong University of Science and Technology, Hong Kong, China

²Zhejiang A&F University

³Department of Global Ecology, Carnegie Institution for Science, Stanford

⁴State Key Laboratory of Heavy Oil Processing, China University of Petroleum, Beijing, China

⁵Xiamen University

⁶Sun Yat-Sen University

⁷Institute of Surface-Earth System Science, School of Earth System Science, Tianjin University, Tianjin, China

Presenting Author: yuanbiyi@ust.hk

Dissolved organic matter (DOM) consists of lots of molecular formulas (MFs), playing a unique function in aquatic systems. However, it remains unclear how MFs determine bulk DOM properties. Here, we use Fourier-transform ion cyclotron resonance mass spectrometry (FT-ICR MS) to characterize DOM molecular composition from 510 samples, distributed in the China Coastal Sea, of which more than 320 samples contain accompanying stable carbon isotope ($\delta^{13}\text{C}$) measurements. A machine learning model (ML) was created based on $\delta^{13}\text{C}$ contained samples, resulting in a mean absolute error (MAE) of 0.30‰ on the test data. The developed model was used to predict in remaining samples without $\delta^{13}\text{C}$ and other published datasets. The result shows that ML selected 5199 MFs to predict $\delta^{13}\text{C}$, and the observed $\delta^{13}\text{C}$ values were more effectively tracked than traditional approaches. Besides, the ML revealed that the elemental composition would be primarily impacted when DOM was transferred from land to ocean. This work suggests that with increasing molecular research and larger learning datasets, ML could be considered a new paradigm for studying the non-linear and non-monotonic connections behind the complex DOM molecular composition.

