

Jingyi Zhou



Yilong Ju

Melissa Suter Spectroscopies with Machine Learning for Detection of Environmental Contaminants

**Ankit Patel** 

Bhagavatula Moorthy N. J. Halas, Rice University

P. Nordlander

**Thomas Senftle** 

ftle Seminar, Stuttgart, Germany: October 17, 2024

Andres Sanchez Lexie Cher

# Gas Chromatography-Mass Spectrometry: large, costly, time-consuming



## **PROJECT VISION:**

#### Can we

-rapidly and inexpensively- detect and identify molecules from complex environmental samples by combining Surface-enhanced Spectroscopy with Machine Learning strategies?

> Baylor College of Medicine

B C M - RICE SUPERFUND RESEARCH P R O G R A M

### Polycyclic Aromatic Hydrocarbons (PAHs)



- 16 Priority pollutants (EPA): toxic, carcinogenic, mutagenic
- Found in atmosphere, in soil and in virtually all water sources: "legacy pollutants"
- Many societal sources: incomplete combustion, diesel fuel, sidestream tobacco smoke, barbecue, burnt toast, drinking water
- Many cancers (skin, liver, testicular, bladder)
- Adult onset diseases (Parkinson's?)
- Poor fetal development: increases in premature birth rates
- Always found in complex mixtures





#### Surface Enhanced Raman Spectroscopy (SERS)

Molecular adsorbates on metal surfaces respond to the local electromagnetic field, emitting inelastically scattered light (Van Duyne, 1974):



Plasmonic substrate requires intense near field at pump, Stokes frequencies:

$$<|\mathsf{E}_{\mathsf{NS}}(\omega_{\mathsf{L}})|^{2} \cdot |\mathsf{E}_{\mathsf{NS}}(\omega_{\mathsf{S}})|^{2} > \approx <|\mathsf{E}|^{4} >$$

#### Surface Enhanced Infrared Absorption Spectroscopy

#### (SEIRA)

• IR absorption, 
$$A \propto \left| \mathbf{E} \cdot \left( \frac{\partial \mu}{\partial Q} \right) \right|^2$$
  $\mathbf{E} = \text{Electric field}$   
 $\mathbf{\mu} = \text{Dipole moment}$   
 $\mathbf{Q} = \text{Normal coordinate}$ 

- SEIRA ~  $|E|^2$ , not  $|E|^4$  like SERS
- spans near-IR to far-IR regime
- •IR cross sections ~10<sup>10</sup> larger than Raman cross sections!

#### SERS and SEIRA combined on the same substrate

H. Wang et al., Angewandte Chemie International Edition 46, 9040-9044 (2007).

F. Le et al., Metallic nanoparticle arrays: a common substrate for both SERS and SEIRA, ACS Nano 2, 707-718 (2008)



#### Combined SERS and SEIRA on the same substrate



## Outline:

- Computational Chromatography: identifying unknown molecules in mixtures without separations
- (M. Bajomo, Y. Ju et al., PNAS 119, e2211406119 (2022))
- Identifying SERS spectra with a Raman library using Machine Learning (Y. Ju et al., ACS Nano 17, 21251-21261 (2023))
- Machine Learning-enhanced SERS + SEIRA Detection of Polycyclic Aromatic Hydrocarbons in Human Placenta

(O. Neumann et al., TBP)

Computational Chromatography: a Machine Learning strategy for demixing individual chemical components in complex mixtures

M. Bajomo, Y. Ju, et al., PNAS 119 (2022) e2211406119

Can we identify the chemical components in a mixture without physically separating them?



Mixture of Chemicals

# Computational Chromatography: a Machine Learning strategy for demixing individual chemical components in complex mixtures

M. Bajomo, Y. Ju, et al., PNAS 119 (2022) e2211406119



Can we identify the chemical components in a mixture without physically separating them?

## The Cocktail Party Problem

(aka "Blind Source Separation")



### Independent Component Analysis: Key Assumptions about the Data-Generating Process

Independent component analysis (ICA) is a ML method for separating a multivariate signal into additive subcomponents.

Signals must be linearly mixed

 $x_1 = a_{11}s_1 + a_{12}s_2$  $x_2 = a_{21}s_1 + a_{22}s_2$ 

x = detected signal s = source signal a = mixing weight Signals must be non-gaussian

Signals must be independent



$$p(\mathbf{s_1}, \mathbf{s_2}) = p(\mathbf{s_1})p(\mathbf{s_2})$$

Chechkin, Aleksei V., et al. "Brownian yet non-Gaussian diffusion: from superstatistics to subordination of diffusing diffusivities." *Physical Review X* 7.2 (2017): 021002.

#### ICA Approximately Demixes to Recover Source Signals (Individual Component Spectra)

Source Signals (SERS of Individual PAHS)



 $x_1 = a_{11}s_1 + a_{12}s_2$  $x_2 = a_{21}s_1 + a_{22}s_2$ 

$$x = As$$
$$s = A^{-1}x$$

$$\hat{s} = Wx$$

**Goal of ICA:** Find <u>demixing matrix W</u> (approximation of  $A^{-1}$ ) so that  $\hat{s} \approx s$ .

## Example: SERS of 2-Component Mixtures of PAHs





Relative intensities of different peaks in the SERS of ANTH+PYR varies as the concentration ratio of ANTH:PYR varies.

## Estimated Independent Components Match PAH Source Spectra





Independent components produced match the SERS of the components of the mixture

## Quantifying Demixing Performance: Evaluation Metrics



## Performance of ICA on 2-Component Mixtures



## How Demixing works "under the hood":

Multiple mixture spectra with varying concentration ratios are required



### Comparing Different ML-based Demixing Algorithms: Informing based on Prior Domain Knowledge is Valuable

Demixing Algos with various assumptions/constraints imposed (aka *inductive biases*):

- Independent Components Analysis (ICA)
- Nonnegative ICA (NICA)
- Sparse ICA (SICA)
- Nonnegative Matrix Factorization (NMF)
- Near-Separable NMF (NSNMF)
- <u>Cha</u>racteristic <u>P</u>eak
  <u>E</u>xtraction-based NMF (CaPE)



Extracting Characteristic Peaks in a shift-tolerant manner significantly improves performance



## Can we extend this to multiple components?

How few samples do we need with varied concentrations?

A synthesized 4component mixture: 8 concentration ratios



- 0.0



# Four Component PAH Demixing



Benzo[a]pyrene )B[a]P)

Anthracene (ANTH)



Pyrene (PYR)

Benz[a]anthracene (B[a]A)

Despite more peak overlap, there is acceptable agreement between demixed spectra and PAH component spectra



B C M - R I C E SUPERFUND RESEARCH P R O G R A M

# Can we use Facial Recognition Approaches to identify SERS spectra of molecules using a Raman Database?

Yilong Ju et al., ACS Nano 17, 21251-21261 (2023).



Medicine

PROGRAM



#### Present and Future of Surface-Enhanced Raman Scattering

Judith Langer,<sup>†</sup> Dorleta Jimenez de Aberasturi,<sup>†</sup><sup>®</sup> Javier Aizpurua,<sup>‡</sup><sup>®</sup> Ramon A. Alvarez-Puebla,<sup>§,||</sup><sup>®</sup> Baptiste Auguié,<sup>L,#,7</sup> Jeremy J. Baumberg,<sup>8</sup><sup>®</sup> Guillermo C. Bazan,<sup>9</sup><sup>®</sup> Steven E. J. Bell,<sup>10</sup><sup>®</sup> Anja Boisen,<sup>11</sup> Alexandre G. Brolo,<sup>12,13</sup><sup>®</sup> Jaebum Choo,<sup>14</sup> Dana Cialla-May,<sup>15,16</sup> Volker Deckert,<sup>15,16</sup><sup>®</sup> Laura Fabris,<sup>17</sup> Karen Faulds,<sup>18</sup><sup>®</sup> F. Javier García de Abajo,<sup>||,19</sup> Royston Goodacre,<sup>20</sup><sup>®</sup> Duncan Graham,<sup>18</sup> Amanda J. Haes,<sup>21</sup><sup>®</sup> Christy L. Haynes,<sup>22</sup><sup>®</sup> Christian Huck,<sup>23</sup><sup>®</sup> Tamitake Itoh,<sup>24</sup> Mikael Käll,<sup>25</sup><sup>®</sup> Janina Kneipp,<sup>26</sup><sup>®</sup> Nicholas A. Kotov,<sup>27</sup><sup>®</sup> Hua Kuang,<sup>28,29</sup> Eric C. Le Ru,<sup>⊥,#,7</sup><sup>®</sup> Hiang Kwee Lee,<sup>30,31</sup> Jian-Feng Li,<sup>32</sup><sup>®</sup> Xing Yi Ling,<sup>30</sup><sup>®</sup> Stefan A. Maier,<sup>33</sup> Thomas Mayerhöfer,<sup>15,16</sup> Martin Moskovits,<sup>34</sup><sup>®</sup> Kei Murakoshi,<sup>35</sup><sup>®</sup> Jwa-Min Nam,<sup>36</sup><sup>®</sup> Shuming Nie,<sup>37</sup> Yukihiro Ozaki,<sup>38</sup><sup>®</sup> Isabel Pastoriza-Santos,<sup>39</sup><sup>®</sup> Jorge Perez-Juste,<sup>39</sup><sup>®</sup> Juergen Popp,<sup>15,16</sup><sup>®</sup> Annemarie Pucci,<sup>23</sup><sup>®</sup> Stephanie Reich,<sup>40</sup> Bin Ren,<sup>32</sup><sup>®</sup> George C. Schatz,<sup>41</sup><sup>®</sup> Timur Shegai,<sup>25</sup><sup>®</sup> Sebastian Schlücker,<sup>42</sup><sup>®</sup> Li-Lin Tay,<sup>43</sup> K. George Thomas,<sup>44</sup><sup>®</sup> Zhong-Qun Tian,<sup>32</sup><sup>®</sup> Richard P. Van Duyne,<sup>41</sup> Tuan Vo-Dinh,<sup>45</sup><sup>®</sup> Yue Wang,<sup>46</sup> Katherine A. Willets,<sup>47</sup><sup>®</sup> Chuanlai Xu,<sup>28,29</sup><sup>®</sup> Hongxing Xu,<sup>48</sup> Yikai Xu,<sup>10</sup><sup>®</sup> Yuko S. Yamamoto,<sup>49</sup> Bing Zhao,<sup>50</sup><sup>®</sup> and Luis M. Liz-Marzán<sup>#,†,51</sup><sup>®</sup>



## The BIG Problem with SERS:

- SERS spectra are different than Raman spectra (in database)!
- SERS spectra are different on each SERS substrate!

#### **Solution Proposed by others:**

create a SERS database for each type of SERS substrate...(not realistic!)

GOAL: Can we develop ML algorithms so we can use a Raman library (database)to identify chemicals through their SERS spectrum?

SUPERFUND

RESEARCH

PROGRAM

College of

Medicine



#### **ML-Based Spectral Recognition: CaPE and CaPSim**



**Yilong Ju** 





Machine Learning-enhanced SERS + SEIRA Detection of Polycyclic Aromatic Hydrocarbons in Human Placenta

Smoking during pregnancy is associated with increased risk of:

- miscarriage,
- prematurity,
- stillbirth,
- low birth weight,
- perinatal morbidity,
- sudden infant death syndrome (SIDS),
- adverse neurodevelopmental disorders such as Attention-Deficit Hyperactivity (ADHD),
- anxiety, and depression

PAHs can cross the placenta from the mother to the fetus, exposing the fetus to these chemicals



PROGRAM

Machine Learning-enhanced SERS + SEIRA Detection of Polycyclic Aromatic Hydrocarbons in Human Placenta



## SERS and SEIRA of Human Placenta



B C M - R I C E SUPERFUND RESEARCH P R O G R A M

#### CaPE Analysis of Smoker vs. Non-Smoker SERS data:





## CONCLUSIONS:

- ML can enable the detection of specific PAHs in complex mixtures without separations
- Characteristic Peak Extraction (CaPE) simplifies SERS spectra so a Raman spectral library can be used for identification with CaPSim
- CaPE and CaPSim can differentiate PAHs in the placenta of smokers versus non-smokers



ROGRAM