# UCLA
## UCLA Electronic Theses and Dissertations

**Title**
Alzheimer's Disease Diagnosis Using SERS and Advanced Data Analysis

**Permalink**
https://escholarship.org/uc/item/7vg9b178

**Author**
Yu, Xinke

**Publication Date**
2019

Peer reviewed|Thesis/dissertation

UNIVERSITY OF CALIFORNIA

Los Angeles

Alzheimer's Disease Diagnosis Using SERS

and Advanced Data Analysis

A dissertation submitted in partial satisfaction of the

requirements for the degree Doctor of Philosophy

in Materials Science and Engineering

by

Xinke Yu

2019

ABSTRACT OF THE DISSERTATION

Alzheimer's Disease Diagnosis Using SERS

and Advanced Data Analysis

by

Xinke Yu

Doctor of Philosophy in Materials Science and Engineering

University of California, Los Angeles, 2019

Professor Ya-Hong Xie, Chair

Early diagnosis of Alzheimer's disease (AD) is critical for disease prevention and cure, but no method to do so have yet been developed that had the required sensitivity and specificity. Computational methods are increasingly being applied in these efforts. One such method is "deep learning." We propose here a convolutional neural network-based AD diagnosis approach using SERS fingerprints of human cerebrospinal fluid to have a preliminary test of the feasibility of early diagnosis using the hybrid platform. To realize the testing, we further prove the reliability of a SERS hybrid platform from its quantification capability and orientation dependence. Analysis using Amyloid beta (Aβ) to prove the biological feasibility and test the

specificity of the platform is also done.

We report results demonstrating the reproducibility and accuracy of this novel SERS data analysis platform. We have achieved 100% reproducibility in double blind experiments and 92% accuracy in disease diagnosis. Comparison of the SERS-neural network approach with single biomarker tests shows it is more accurate, thus it may have substantial value in the differential diagnosis of AD as well as other neurodegenerative disorders.

We also show here that surface-enhanced Raman spectroscopy (SERS) coupled with principal component analysis (PCA) readily distinguishes small biological differences: Aβ40 and Aβ42. We show further, through comparison of assembly-dependent changes in secondary structure and morphology, that the SERS/PCA approach readily and unambiguously differentiates closely related assembly stages not readily differentiable by circular dichroism spectroscopy, electron microscopy, or other techniques.

To test the substrate feature, we demonstrate, using a biologically relevant test analyte, the amyloid β-protein (Aβ), a seminal pathologic agent of Alzheimer's disease (AD), that linear relationships exist between (a) peak intensity and concentration at a single plasmonic hot spot smaller than 100 nm, and (b) frequency of hot spots with observable protein signals, i.e. the co-location of an A    protein and a hot spot. We demonstrate the detection of Aβ at a concentration as low as 10-18 M after a single 20 μl aliquot of the analyte onto the hybrid platform.

Orientation dependence is also proved by analyzing the standard deviation of spectral feature. The standard deviation in the intensity of individual Raman peaks diminishes

for protein size larger than 13 amino acids. Secondary structure of protein (such as protein-protein interaction) remains unchanged regardless of protein orientation. Numerical simulation studies corroborate the experimental observation in that the SERS spectral features of biomedically relevant protein (of larger than 13 amino acids in size, which represent all human protein types) are not affected by the orientation of amino acids randomly dispersed on SERS-active surfaces.

The dissertation of Xinke Yu is approved.

Yang Yang

Gerald Chee Lai Wong

Ya-Hong Xie, Committee Chair

University of California, Los Angeles

2019

**To my parents, my boyfriend and my dog, who supported me throughout my PhD journey, and without whom it would not have been possible.**

# Table of Contents

# List of Tables

# List of Figures

# Acknowledgement

First and foremost, I acknowledge my great thank to my advisor Dr. Ya-hong Xie, without whom none of the work presented in this dissertation would have been accomplished. I came to work with Prof. Xie six years ago for a summer research program and did not regret a single day for coming back to the same lab as a PhD student. Working with Prof. Xie was an amazing experience. I gradually understand what real science is and begin to understand the beauty of physics. Through his mentorship, I gained theoretical and experimental skills, and most importantly, learned to be a rigorous researcher and an integrate person. Being Prof. Xie's student is beyond lucky and I will always remember what I have learned from him through the five years.

I would like to thank Dr. David Teplow and Eric Hyden for your help in guiding me to the world of neuroscience. All the discussions are insightful and working with them was an amazing experience.

One of the best things of my PhD is knowing all the fantastic people in Xie's Group. I would like to thank Dr. Wei Zhang, Dr Pu Wang, Dr Ming Xia, Dr Zhongbo Yan, Owen Liang, Peiyi Ye, Zirui Liu, Shan Huang, Ming Dong, Jun Liu and Tieyi Li for not only being great co-workers but good friends as well. I felt gratitude for working with all these smart minds.

I would like to thank my boyfriend, Shilei, for the unconditional love, trust and support throughout the five whole years. Long distance relationship is hard, but he

makes it much easier. He has guided me to become a better and stronger person. He is always there for me no matter what happens. I will not finish my PhD so smoothly without his help.

I have adopted a dog, Stoney, during my PhD and it is one of the best decisions I made. He is always friendly and loving, soothing my mind and making me relax after long days' work.

At the end, I would like to show my great gratitude to my parents, for their support and understanding in my pursing my dreams. I am honored to be their daughter.

VITA

| 2014 | B.S. in Electrical Engineering and Computer Science |
| | Peking University, Beijing, China |
| 2016 | M.S. in Materials Science and Engineering |
| | University of California at Los Angeles |
| 2019 | Ph.D. Candidate in Materials Science and Engineering |
| | University of California at Los Angeles |

# Chapter 1

# Introduction

## 1.1   Motivation and Innovation of Thesis

Alzheimer's disease (AD), being one of the most fatal diseases in the world, has attracted more attention and the myth of the disease still remains unknown to researcher around the world. Many of the hypotheses regarding the disease contain Amyloid beta (Aβ) peptide, which can become toxic with its secondary structure changes. No solid prediction method exists, not to mention a cure for this disease.

Surface enhanced Raman spectroscopy (SERS) has attracted a growing attention in the area of chemistry, biology, medicine, pharmacology and environmental sciences due to its fast speed, single molecule detection sensitivity and molecular specificity. However, while the direct SERS measurement of small molecules is simple and the results are easy to validate, when it comes to molecule with larger molecular weight (such as protein), questions on quantification, specificity and reproducibility of SERS testing arise. To better apply this powerful analytical technique into real medical applications,

a detailed analysis on the quantification capability and specificity <u>at protein level</u> is needed.

With the protein level analysis capability, a huge number of diseases can be diagnosed using SERS as they are closely related to secondary folding of protein. Huntington's, Parkinson's and "mad cow" disease are all protein misfolding diseases. Hungtington, alpha-synuclein, and prion are proteins that may misfold and accumulate leading to these diseases, respectively.

A combination of SERS-active metal nanostructure (gold pyramid structure) and bio-compatible material (graphene) create multiple synergies. A proven extra high enhancement factor and ultra-high sensitivity has been shown in previous work and by utilizing this platform, I present here several mechanisms to quantify protein concentration; validate the molecular specificity on protein level especially regarding orientation; monitor the Aβ peptide and mutation; distinguish AD patients from normal controls. Several novel and significant approaches and results will be presented in this work:

**<u>Aβ differentiation using SERS:</u>** Aβ is one of the key biomarkers of AD and understanding the structure difference between the two types of Aβ peptides (Aβ40 and Aβ42) is of vital importance. A detailed analysis between the spectral feature of the two peptides is done. Time dependence structural change of Aβ42 is also demonstrated.

**<u>SERS Quantification Mechanisms:</u>** The platform has reproducible quantification

capability with 2 quantification mechanisms and 7 orders of magnitude dynamic range. The outstanding quantitative characteristic is demonstrated both theoretically and experimentally.

**Orientation Dependent Specificity:** Single molecule level orientation analysis via bio-analyte method is done. Specificity on protein level is proved using theoretical analysis, simulation and experimental data. Secondary structure stability is further proved by analyzing Amide peaks.

**Cerebrospinal Fluid (CSF) based AD diagnostic:** CSF is used as a key body fluid to diagnose neuro degenerative diseases. Spectral features of 26 different individuals are collected and tested using machine learning algorithms and deep learning methods to classify AD patients vs. normal individuals at a high accuracy.

The above highlighted points establish a solid foundation for protein level SERS detection and further apply the novel platform into medical fields for accurate disease diagnostic.

This chapter presents an introduction on the background of the whole research. To better understand Alzheimer's disease, Chapter 1.2 shows basic knowledge on the disease. In Chapter 1.3, information related to SERS is presented and we expand the technology to our special platform in Chapter 1.4. Chapter 1.5 gives a concise summary for all the data analysis methods used in our spectral analysis and Chapter 1.6 presents the outline of the dissertation.

## 1.2 Alzheimer's disease

Alzheimer's disease (AD) is a disease of aging. It is characterized in part by progressive loss of memory and executive functions and primarily affects older adults and is the most common cause of dementia. The loss of functions is attributed to synaptic damage and neuronal loss in the hippocampus, cerebral cortex and other brain regions. The disease worsens as it progresses and eventually leads to death. While treatments to ameliorate some symptoms exist, there is currently no cure for AD.

The disease was first described by German neuropathologist Alois Alzheimer in 1906, and the disease was later named after him. Early-onset AD (onset of symptom before 65 years old) is rare and usually gene related (present before their 50s). The age of people with AD in 2018 is shown in Fig 1.1. As is shown in the figure, only 4% of the AD patients have onset of symptom before 65 and with the aged tendency of population, the percentage of AD patient over 85 may further increase to be larger than 75-84 regime.



- 85+ years, 37%
- 75-84 years, 44%
- 65-74 years, 16%
- <65 years, 4%

**Fig 1.1.** Ages of people with Alzheimer's dementia in the United States, 2018

For most cases, the prevalence rate of AD increases exponentially after the age of 65. In 2010, there were an estimated 454,000 new cases of AD and the number is projected to be 959,000 in 2015 (110% growth). And in 2050, the number of people with AD will grow from 5.5 M to 13.8M as is shown in the Fig 1.2.



**Fig 1.2.** Projected number of people age 65 and older (total and by age) in the U.S. population with Alzheimer's disease, 2010 to 2050.

The cause of AD remains unknown and one of the most popular theories is amyloid hypothesis. It postulates that extracellular beta-amyloid (Aβ) deposits are the fundamental cause of AD. Facts are presented to prove the hypothesis: location of the gene for the amyloid precursor protein (APP) is on chromosome 21 and people with trisomy 21 almost universally exhibit AD by the age of 40 years of age. Aβ peptide plays an important role in this process as it is the main component of senile plaques. The ratio between the two types of Aβ peptides (Aβ40 and Aβ42) can possibly lead to

AD.

Aβ belongs to a class of protein that is "natively unfolded" and preferentially form amyloid fibrils rather than protein crystals and the pathway of fibril assembly is illustrated in Fig 1.3. Aβ structures aggregate from monomer to paranucleus and protofibril (identified more than a decade ago) and eventually mature to fibril. The length of these structures is <150nm and ~5nm in diameter.



**Fig 1.3.** The pathway of Aβ fibril assembly.

Even though the Aβ peptides and Aβ fibril are the same protein, Aβ fibrils have a fiber-like structure and the peptide has a β-sheet arrangement. Whether a β-sheet structure transforms into a fiber structure and what the transformation speed is depend on whether the protein folds or misfold.

For Aβ, protein misfolding happens when the hydrophobic section of the protein fails to fold into the interior and thus bonding with other water-repelling portions of other unfolded proteins, which eventually lead to protein aggregation and form fibril structure. To monitor the folding process and to better understand the pathology of AD, detecting the changes in the senconday folding state of Aβ peptide leads to a more

complete image and is vital for the process.

## 1.3　Surface enhanced Raman Spectroscopy (SERS)

Raman spectroscopy has been widely used in the area of non-destructive testing and molecular recognition technology. It can help to provide fingerprint information of chemical and biomolecular structure. However, the cross section of conventional Raman scattering were only $10^{-6}$ to $10^{-14}$ of that of IR and florescence process. This inherent low sensitivity has tremendous restriction in its application in the field of trace detection and surface science.

In 1974, Fleischmann et al. roughened the smooth surface of the silver electrode, and for the first time, they found the high quality Raman spectra of the mono molecular layer pyridine molecules adsorbed on the surface of the silver electrode. In 1977, Van Duyne and Creighton independently discovered that the Raman signal of the pyridine molecules is about 106 times higher than that of single pyridine in solution. They pointed out that this is a surface enhancement related to surface roughness and is named as SERS effect (as is illustrated in Fig 1.4). And SERS effectively solve the critical problem of low sensitivity in Raman.

**Fig 1.4.** Conceptual illustration of SERS.

SERS can be simply described as amplified Raman scattering by the presence of a plasmonic structure in the close vicinity of the target analyte. It is a valuable tool in multiple research fields such as biology, pharmaceutical, chemical, etc. The high sensitivity and specificity make SERS an incomparable technology. Moreover, the label free detection, non-destructive nature and minimal waste feature make SERS a unique analytical tool in bio-sensing. As SERS contain unique lattice vibration information of the analyte, more detailed and specific information can be acquired with high sensitivity.

Most disease states start with a small change in cellular processes and the change becomes amplified along the disease progression. As a result, in bio-medical disciplines, sensitive enough to measure small concentration and outstanding specificity to distinguish minute differences are vital in diagnostic. Besides, non-destructive and minimal invasive make a technique even more attractive. With all these features, SERS has attracted a growing attention and researchers have been

8

working for over 30 years to apply the technology into bio medical applications.

Nanoparticle has been widely used as SERS probe: Ahmed *et al.* used Ag nanoparticles (NP) and AuNP in rats' brains to understand neurological diseases; Masson *et al.* applied functionalized AuNP for *IgG* protein detection. One popular substrate is called SERS-active substrate with self-assembled monolayer (SAM) and have been widely used: Zhang et al. reported the use of SAM substrate in glucose detection; Zou *et al.* have distinguished diabetes from normal individuals. Other substrates such as colloidal SERS substrates (Pucetaite *et al.* have used the substrate for cardiovascular disease), SERS nanowire sensor (Eom et al. have used the substrate for breast cancer diagnostic) and some other nano structured substrate including nano pyramids, nanocrystal, nanorods .etc.

SERS detection is popular among disease diagnostic, however, many problems appear in this process, such as the ambiguous enhancement factor definition, the difficulties in quantification and the toxicity of the metal used. Even with these problems, SERS continue to remain at the forefront of disease diagnostic in both in vitro and in vivo applications.

## 1.4   Graphene Hybrid Platform

An ideal SERS platform for bio sensing needs to have a high enhancement factor (to reach single molecule sensitivity) and bio compatibility (to have minimal influence on

the features of the analytes). The hybrid platform is a combination of periodic Au pyramid structure and single layer graphene, which is a perfect match for these two requirements.

The Au nano structure boosts strong plasmonic enhancement to provide ultra-high enhancement factor. The most optimized structure has a pyramid tip of 200nm, which provides the highest sensitivity. Single layer graphene provides bio-compatible environment for the analytes with the proof that cell can directly survive on graphene without any glial layer. Besides, graphene can increase the resistance of metal to oxidation and electrochemical degradation and make the substrate more stable.

The hybrid platform is fabricated by transferring monolayer graphene on the Au pyramid substrates, as is shown in Fig 1.5.



**Fig 1.5.** Schematic process showing the synthesis of the hybrid platform. The CVD monolayer graphene is transferred onto the Au tip substrate.

The hybrid platform enables single molecule detection and provides a reproducible and uniform response. It allows a sub-$10^{-18}$ M detection of Aβ42 with the high electromagnetic enhancement(EM) of the nanostructure and extra chemical

enhancement (CM) of graphene. It also provides a built-in hotspot marker, helping to define the hotness of the hotspot and locate the exact position of the hotspot.

This platform overcomes the limitation of conventional nanoparticle SERS systems and makes accurate analysis of protein concentration possible. Further analysis on disease diagnostic is also made possible with the features of the substrate.

## 1.5 Data Analysis

The complexity of biological-analyses (such as protein, exosome and cell) has led to complicated and diverse Raman spectrum. To ensure the inclusion of the ever present statistical variation due to factors such as biological and individual variability as well as the many co-factors such as a patient suspected of cancer could also be suffering from high blood pressure or diabetes, large number (on the order of hundreds) of spectra from each sample must be collected with their spectral features categorized and latter compared with through comparison with the diagnosis rendered by the current laboratory practice. And thus, it is not difficult to appreciate the massive volume of data for the platform to render clinical usefulness. To make the mass data analysis possible, statistical methods needs to be applied in the analysis.

**Fig 1.6.** A diagram of standard SERS analysis work flow

Besides, the nature of SERS spectrum itself makes the analysis complicated. To start with, peak intensity does not solely depend on the analyte, difference in the surface structure (hotness of the hotspot) will lead to different peak intensity. Furthermore, one of the methods to interpret Raman spectra is to take the intensity of each wavenumber as a dimension, each of the spectra can be taken as a ~1500 dimension data. It is relative inefficient to operate on a high dimension data. All these complexity makes statistical methods a must in SERS analysis.

Multiple data analysis methods can be used in the spectral analysis and the performance of those methods varies over different applications because of their difference in factors

such as size of dataset and analyte composition. A detailed analysis is needed for each type of sample we analyze to optimize the analysis results. A detailed description of those methods and their application in SERS analysis will be detailed described in Chapter 4 and 5.

## 1.6   Outline of Thesis

The remainder of this thesis is divided into 5 chapters, which go from the theoretical analysis of SERS hybrid platform to the application in its disease diagnostic. As concisely summarized above, Chapter 1 discusses the complex field of SERS platform and data analysis. These two segments work together to build a diagnostic system for protein misfolding diseases.

Chapter 2 shows the quantification capability of the platform. 7 orders of magnitude dynamic range with 2 different quantification mechanisms is presented. Such capability enables detailed concentration analysis using SERS with high reproducibility.

Chapter 3 elaborates the specificity of SERS when doing protein analysis, especially from the orientation dependence perspective. The study simulates protein orientation and the change of orientation dependence as molecular weight of protein increases. The stability of secondary structure features is further proved, providing a solid basis for SERS diagnostic using protein as biomarker.

Chapter 4 discusses the application of SERS hybrid platform protein detection. By

applying principal component analysis we are able to differentiate the 2 amino acid difference between the two types of Aβ peptides. Similar method is applied to trace the structural change in Aβ42 peptide as time changes. The successful application shows the potential of SERS in clinical applications.

Chapter 5 details the diagnostic of Alzheimer's disease using SERS. Multiple machine learning algorithms are applied to do the patient classification. Patient specificity of 100% was reached doing double blind experiment and the prediction accuracy has reached over 80%. The high accuracy proves that SERS can be used in human fluid application and have huge potential in disease diagnostic.

# Chapter 2

# Protein differentiation using SERS

## 2.1 Introduction

Amyloid β-protein (Aβ) assembly into neurotoxic structures appears to be a seminal pathogenetic event in Alzheimer's disease (AD)[1]. For this reason, intense efforts have been devoted to understand the physiology of Aβ production, how the peptide assembles into neurotoxic structures, and the mechanism of amyloid plaque formation. Each of these efforts represents a potential therapeutic approach to prevent or treat AD. Equally important has been the search for biomarkers that would enable accurate diagnosis of disease state and provide metrics for evaluation of clinical trial efficacy. Unfortunately, thus far, no effective therapeutic agents or reliable biomarkers are available for clinical use. One reasonable approach to address these unmet needs is to develop new methods for dissecting the Aβ assembly process, methods that could

reveal structural details of assembly at heretofore unsurpassed resolution and sensitivity and enable identification of novel structural biomarkers, e.g., Aβ structural states that correlate with disease status. Surface enhanced Raman spectroscopy (SERS) is one such method, which is known to have single molecule sensitivity[2].

SERS has been applied to the Aβ system in the past, but no systematic studies of the system have been published, to our knowledge. Beier *et al.* used SERS to detect Congo Red bound to Aβ, reporting a linear detection regime of $10^{-12}$-$10^{-8}$ M[3]. Benford used a nanofluidic device containing trapped gold colloid particles (60 nm) to physically restrict Aβ in an illuminated volume[4]. This allowed them to detect Aβ40 at concentrations as low as 11.5 pM. Bhowmik *et al.* used lipid bilayer-coated silver nanoparticles to bind Aβ40 and determine its secondary structure[5]. A number of nanofluidic devices have been fabricated in which Aβ can be concentrated and its spectra acquired at concentrations as low as 10 fM[4,6]. Voiciuk adsorbed oligomeric forms of Aβ42 onto self-assembled monolayers terminated by heptanethiol, octadecanethiol, or N-(6-mercapto)pyridinium groups in an effort to detect unique spectral features[7]. Changes in carboxylate stretching modes were observed upon binding. Most recently, Nabers *et al.* suggested that spectral shifts in the amide I band of Aβ in cerebrospinal fluid might discriminate dementia of the Alzheimer's disease type from other types of dementia[8]. This study used FT-IR, not Raman, spectroscopy, but these results illustrate the potential usefulness of SERS because both FT-IR and

SERS probe the vibrational modes of molecules[9].

Recently, Wang *et al.* developed a graphene-gold hybrid plasmonic SERS platform with intrinsic electromagnetic field enhancement normalization capabilities and extremely high sensitivity. The platform is capable of single molecule detection sensitivity[10] and has been shown to detect Aβ42 at attomolar ($10^{-18}$) concentration[11]. We present here results of studies demonstrating that this novel platform readily distinguishes Aβ42 from Aβ40 and reveals distinct spectral signatures for different conformational and assembly states. These capabilities allow monitoring of time-dependent conformational and assembly changes as well as the potential of defining new disease state biomarkers based on specific spectral signatures.

## 2.2  Experimental procedures

### 2.2.1  *Principal component analysis*

Principal component analysis (PCA) is a statistical procedure that is primarily used for dimensionality reduction. It uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called principal components (PCs).

**Fig 2.1.** A 2D example of how PCA works.

As is shown in Fig 2.1, PCA is useful for eliminating dimensions, the two graphs show the exact same data, but the right panel reflects the original data transformed so that the axes are now the principal components.

PCA is a method that brings together a measure of how each variable is associated with one another and at the same time shows the directions in which our data are dispersed.

In this case, R was used to analyze and visualize the multi-dimensional SERS data sets by using the built in PCA function: "prcomp" and open source library "ggbiplot" (with extended functionality for labeling groups, drawing a correlation circle, and adding normal probability ellipsoids). Dimensions in the analyses were the intensities at five different wavenumbers determined by decision tree analysis to be necessary

and sufficient for differentiating among time points. PCA transforms the original variables into a set of linear combinations (principal components: PC), which allows the retention of the data variability, while examining them independently in a weighted fashion of decreasing order of variance. Variability in the vectors gathered by the PCs was calculated and the largest two PCs were plotted. Data from each time point was considered a separate group and PCA was done to maximize the between-group variation.

*2.2.2 Decision tree*

A decision tree is a decision support tool that uses a tree like model and it belongs to the family of supervised learning algorithms. It can be used to solve regression and classification problems, which is important when we choose which peak(s) are vital in spectral differentiation.

For a binary classification problem, entropy measure is calculated for information gain, which calculates the expected reduction in entropy due to sorting on the attribute.

$$\text{H(X)} = \mathbb{E}_X[I(x)] = -\sum_{x \in \mathbb{X}} p(x) log p(x)$$

In the formula,

H(X) is the entropy to measure the amount of uncertainty in the data set;

p(x) is the proportion of the number of elements in class x to the number of elements in X.

$$\text{IG(A, X)} = H(X) - \sum_{t \in T} p(t)H(t)$$

In this formula,

H(X) is the entropy of set X;

H(t) is the entropy of subset t.

Decision tree algorithm construct decision tree based on features that have highest information gain. In our case, decision trees were produced using R with the "rpart" package (generating classification and regression trees). The classification method in the "rpart" function was used to produce the tree information and the "prp" function in "rpart.plot" library was used to visualize the decision trees.

### 2.2.3  Experimental settings

*Circular dichroism spectroscopy (CD)*

LMW Aβ was prepared at a concentration of 20 μM in sodium phosphate, pH 7.4, and incubated without agitation at 37 ℃ in a 1 mm path-length quartz cuvette (Hellma, Forest Hills, NY, USA). CD spectra then were acquired periodically with a J-810 spectropolarimeter (JASCO, Tokyo, Japan). Spectra were recorded from 195-260 nm

at 0.2 nm resolution with a scan rate of 100 nm/min. Ten scans were acquired and averaged for each sample.

*Transmission electron microscopy (TEM)*

Five μL of Aβ42 (20 μM) were removed at the time of each CD measurement and then spotted onto carbon-coated Formvar grids (Electron Microscopy Sciences, Hatfield, PA, USA). After 2 min, the droplet was displaced with an equal volume of 1% (w/v) filtered (0.2 μM) uranyl acetate in water (Electron Microscopy Sciences). This solution was wicked off and then the grid was air-dried. All grids were coded at the end of the time course so that the operator of the electron microscope did not know what samples were being imaged. Electron microscopy was done using a JEOL 1200 EX transmission electron microscope with an accelerating voltage of 80 kV, which is typical for protein examination. Digital images were analyzed with ImageJ 1.50d, using the "measure tool" to calculate dimensions, and unblinded after the analysis was complete.

*Raman Measurement*

Immediately following solubilization, 20 μL aliquots of Aβ40 or Aβ42 were applied to a graphene-coated, pyramidal gold hybrid platform and dried *in vacuo*. Spectra were acquired using a Renishaw inVia microscope under ambient conditions. The excitation wavelength was 785nm and the He-Ne laser power was 0.5 mW. The 785

nm laser was chosen due to the relatively lower photon energy of excitation, which avoids thermal degradation of biomaterials. The grating used was 1800 lines/mm, and the objective lens used was 50×. We scanned the entire region on the platform occupied by the samples ($\approx$24 μm $\times \approx$30 μm) using Raman mapping with a step size of 3 μm (i.e., independent areas of 9 μm$^2$ each). Eighty spectra were acquired for each sample. This process controls for acquisition of spectra unrepresentative of the average spatial orientation or assembly state of the peptide, two factors that can affect peak location and intensity, and which become problematic if spectra are acquired from only one or a few areas of the platform. For Raman measurements done in parallel with CD, 10 μL aliquots were taken from the CD cuvette at the time of each CD measurement, applied to the platform, and then spectra were acquired, essentially as described above. Raman data were analyzed using Renishaw WiRE 4.2 software, which automatically subtracts the baseline signal and removes noise. Peak intensities in each spectrum were normalized to the graphene G peak to enable spectral comparisons among samples.

*t-test*

Paired *t*-tests, as implemented in R as "t.test" and using the key parameter "p.value," were done to assess the statistical significance of differences between the centroids of clusters in the PC1 dimension, in which the largest variances were observed. The

analyses were performed on centroids in the PC2 dimension if no significant differences were observed in PC1. Significance was defined as $p < 0.05$.

## 2.3   SERS analysis of unassembled Aβ40 and Aβ42

Our initial experiments sought to establish the spectral characteristics of low molecular weight (LMW; see Methods) Aβ40 and Aβ42. To do so, these peptides were prepared at concentrations of 20 μM in 10 mM sodium phosphate, pH, 7.4, and applied to a unique hybrid SERS platform. This platform consists of a hexagonal array of gold pyramids overlaid with a single molecular layer of graphene[10]. SERS spectra were acquired in the wavenumber range of 550-1800 cm$^{-1}$. Graphene produces two characteristic peaks, the D- and G-peaks, at ≈1350 and ≈1580 cm$^{-1}$, respectively[12]. The graphene G-peak height depends directly on the local electromagnetic field (EM) intensity (within the area of illumination of a tightly focused laser beam ~1 μm in diameter), which can vary substantially among Raman active locations (hot spots) on the platform. Normalization of peak heights at a particular hot spot to the graphene G-peak height thus provides the means to accurately determine analyte signal intensities.

The graphene D-peak arises from the disordered structure of graphene, including broken carbon-carbon bonds and folds formed from the nearly planar graphene overlaid on the pyramidal platform surface, both of which can be byproducts of the fabrication process. The presence of disorder in the sp$^{2}$-hybridized carbon system leads to the

appearance of the D-peak peak[12,13]. The D-peak intensity depends on the polarization

direction of the laser beam relative to that of the graphene fold direction[13], thus its

provenance differs from that of the G-peak. Neither the D-peak nor the G-peak arise

from the protein analyte. However, these peaks do occur in the higher wavenumber

portion of the amide II band (1510-1580 cm$^{-1}$) and the lower wavenumber portion of

the amide I band (1600-1700 cm$^{-1}$) and thus can obscure some protein vibrations related

to secondary structure. As will be shown below, our method of data analysis does not

depend upon these obscured signals.



**Fig 2.2.** SERS analysis of Aβ40. SERS spectra of Aβ40 (red) and Aβ42 (turquoise) are shown.
Wavenumbers are listed above each peak. Graphene D- and G-peaks, at 1350 and 1580 cm-1,
respectively, are signified by letters. Peaks signified by asterisks are likely due to cosmic rays,

as the peak height-to-width ratios are extremely large. The data presented are representative of two independent experiments.

Spectra for Aβ40 and Aβ42, which have been baseline subtracted and normalized to the G-peak, are shown in Fig. 2.2. Predominant peaks in the Aβ40 spectrum occurred at 935, 1000, and 1124 cm$^{-1}$. Less intense peaks at 559, 575, 823, 850, 982, and 1450 cm$^{-1}$ were observed reproducibly in the Aβ40 spectra. The Aβ42 spectrum had clearly observable, but smaller, peaks at 935, 1000, and 1124 cm$^{-1}$. The 823 cm$^{-1}$ peak observed in the Aβ40 spectrum was not seen in the Aβ42 spectrum and the 935 cm$^{-1}$ peak shoulders at 850 and 982 cm$^{-1}$ were substantially smaller. However, the peak at 1450 cm$^{-1}$ was more pronounced.

**Table 2.1.** Aβ40 and Aβ42 Raman peak positions, assignments, and intensities. Peaks have been assigned to specific bond resonances based on published data.

| Wavenumber | 559 | 575 | 823 | 850 | 935 | 982 | 1087 | 1124 | 1450 |
|---|---|---|---|---|---|---|---|---|---|
| Peak assignment | Aliphatic | C-C stretching | Out-of-plane ring breathing or Tyr | Amino acid single bond | n(C-C) of protein backbone or Gly | C-C stretching β-sheet or Phe | Lys or Asn | Val or Ile | CH$_2$ bending or Phe |
| Aβ40 (average) | 0.218 | 0.163 | 0.021 | 0.156 | 0.445 | 0.176 | 0.032 | 0.142 | 0.003 |
| Aβ40 (SD) | 0.0587 | 0.0734 | 0.0529 | 0.0619 | 0.0722 | 0.0338 | 0.0155 | 0.0419 | 0.0311 |
| Aβ42 (average) | 0.129 | 0.077 | 0.000 | 0.086 | 0.289 | 0.086 | 0.017 | 0.051 | 0.027 |
| Aβ42 (SD) | 0.0319 | 0.0392 | 0.0000 | 0.0291 | 0.0613 | 0.0233 | 0.0329 | 0.0519 | 0.0411 |

The observation in the Aβ40 and Aβ42 spectra of peaks at identical wavenumbers is expected because the primary structures of the two peptides also are identical, with the exception of the Ile-Ala dipeptide at the C-terminus of Aβ42. However, the conformational states of the peptides during oligomerization and fibril formation have been shown to differ[14]. In addition, different conformers may be oriented differently with respect to the graphene surface[15]. These factors likely explain the fact that the peak intensity profiles of Aβ40 and Aβ42 are distinct. To more fully understand the significance of the distinct patterns of spectral intensities, we performed unbiased multivariate analysis using principal component analysis (PCA), reasoning that PCA might enable the differentiation of Aβ isoforms and assembly states[16,17]. We parameterized the analysis using nine normalized peak intensities (Table 2.1).

We then performed PCA with each vector having the same variance and found that principal components (PC) 1 and 2 accounted for 57.8% and 15.0%, respectively, of the variance in the data. The cumulative percentage of 72.8% means that the first two principal components account for the majority of the variance in the system. We note that other components account for no more than 5% each of the total variance, so their inclusion in our analyses would not alter our conclusions. PCA analysis can produce statistically significant results when $n>5p$, where $n$ is number of nodes and $p$ is number of vectors used [1]. In our case, using 80 spectra and nine vectors, $n>11p$, so we are

confident that the first two components accurately account for ~75% of the total system variance.

A graph of the results using PC1 and PC2 as axes revealed that the data from Aβ40 and Aβ42 clustered in two distinct regions (Fig. 2.3). The Aβ40 cluster was substantially smaller than the Aβ42 cluster, which suggests that its component conformers were more homogenous structurally than the conformers in the Aβ42 cluster. The Aβ40 cluster displays similar variance in PC1 dimension compared to the Aβ42 one. However, its variance in the PC2 dimension was approximately twice that of Aβ42. The equations specifying PC1 ($\vec{C}_1$) and PC2 ($\vec{C}_2$) provide an explanation for the cluster locations and shapes.



**Fig 2.3.** PCA analysis. Plot of principal components 1 and 2 from analysis of unassembled Aβ40 (salmon) and Aβ42 (turquoise). Ellipses surrounding clusters enclose 67% of the data, indicating the majority of the data points are in the ellipse. The brown arrows are the projections of the vectors in PC space.

$$\vec{PC_1} = -0.27\vec{V_1} - 0.31\vec{V_2} - 0.35\vec{V_3} - 0.38\vec{V_4} + 0.42\vec{V_5} - 0.33\vec{V_6} + 0.24\vec{V_7} +$$

$$0.26\vec{V_8} - 0.39\vec{V_9}$$

$$\vec{PC_2} = 0.15\vec{V_1} + 0.21\vec{V_2} + 0.08\vec{V_3} + 0.11\vec{V_4} - 0.08\vec{V_5} + 0.23\vec{V_6} + 0.65\vec{V_7} +$$

$$0.64\vec{V_8} + 0.10\vec{V_9}$$

Vectors ($\vec{V}$) 1-9 represent the peak intensities of SERS peaks at 559, 575, 823, 850, 935, 982, 1000, 1124 and 1450 cm$^{-1}$, respectively. A key difference between the two principal component vector equations is the absolute value of the coefficients of vectors $\vec{V_7}$ and $\vec{V_8}$ (the peak heights at 1000 and 1124 cm$^{-1}$, respectively), which are substantially larger in the case of Aβ42 compared to Aβ40. This observation, which is not immediately apparent from analysis of peak intensities alone (Table 1), explains why the variance in PC2 space is twice as large for Aβ42 than it is for Aβ40. Resonances at 1000 cm$^{-1}$ and 1124 cm$^{-1}$ are produced by Lys and Asn, and by Val and Ile, respectively (Fig. 2.4). The presence of the additional Ile at the C-terminus of Aβ42 likely is an explanation for at least a portion of the increased magnitude of $\vec{V_8}$. Conformational effects due to the distinct conformational dynamics of Aβ42 may also contribute to the differences in vector magnitudes.

**Fig 2.4.** Examples of vibrational modes in Aβ42. PDB 5KK3 (amino acids 11-42) was used to illustrate possible locations within Aβ42 that could lead to the vibrational modes. Amino acids proposed to contribute to the Raman signal are indicated as follows: Phe: salmon (982 cm$^{-1}$, 1450 cm$^{-1}$), Val: cornflower blue (850 cm$^{-1}$, 1124 cm$^{-1}$), Ile: light sea green (850 cm$^{-1}$, 1124 cm$^{-1}$), Gly: sky blue (935 cm$^{-1}$), Lys and Asn: medium purple (1000 cm$^{-1}$). Out-of-plane ring breathing from Phe could contribute to the 823 cm$^{-1}$ peak. C-C stretching in amino acids could give rise to the peak at 850 cm$^{-1}$. C-C stretching of the protein backbone (935 cm$^{-1}$) and CH$_2$ bending in amino acids (1450 cm$^{-1}$) peak. Black arrows in each inset indicate locations within the peptide that could give rise to these modes.

## 2.4   SERS analysis of Aβ assembly

We next sought to establish whether SERS could distinguish different stages of Aβ assembly. We characterized assembly stages by performing SERS in parallel with

circular dichroism (CD) spectroscopy and transmissions electron microscopy (TEM). We thus obtained CD, TEM, and SERS data from the same sample aliquots. The data shown are representative of four independent experiments. CD complements SERS by providing information in spectral regions obscured by the graphene D- and G-peaks. It also allows real time monitoring of secondary structure changes in assembly reactions *in hydro*. CD spectra were acquired immediately after initiation of assembly reactions of 20 μM Aβ42 in 10 mM sodium phosphate, pH 7.4, at 37 ℃. The spectra were consistent with statistical coil (SC) structure, as indicated by a minimum in molar ellipticity $[\Theta]$ at 198 nm and a gradual increase in $[\Theta]$ as wavelength increased toward 260 nm (Fig. 2.5). A concerted time-dependent increase in $[\Theta]_{198}$ and decrease in $[\Theta]_{218}$ were consistent with β-sheet formation. Maximal β-sheet content was observed at 120 h.



Fig. 3

**Fig 2.5.** Circular dichroism spectroscopy. Aβ42 was incubated at a concentration of 20 μM in 10 mM sodium phosphate, pH 7.4, at 37 ℃. CD spectra then were acquired periodically. Overlapping spectra have not been presented so as to make visualization of the time-dependence of spectral changes easier. The data shown are representative of four independent experiments. Spectral colors represent different time points, which are specified in hours in the box to the right.

These data were consistent with the SC→β-sheet secondary structure transitions that occur during Aβ fibril formation[19]. Negative stain EM done in parallel with the CD studies confirmed that fibril assembly was occurring (Fig. 2.6). The starting peptide solution (0 h), which displayed statistical coil secondary structure, contained only small globular structures of ≈8 nm diameter. Short fibrils of width ≈10 nm were observed at 6 h, during the coil→β-sheet transition period. When maximal β-sheet was observed, long fibrils were present. The morphologies of these structures did not change substantially after 24 h.



Fig. 4

| 0 h | 6 h | 120 h |

**Fig 2.6.** Transmission electron microscopy. Aβ42 (20 μM in 10 mM sodium phosphate, pH 7.4) was incubated at 37 ℃ for a total of 192 h. Panels shown are representative of the sample morphologies observed at the indicated time point.   Scale bars are 100 nm.

**Fig 2.7.** SERS analysis of Aβ42. Spectra were acquired periodically from the same samples used for EM. Spectra from different time points are shown in a fence plot. Axes are wavenumber (cm⁻¹), time (h), and intensity (AU). We show spectra up to and including 24 h, after which all spectra were identical, within experimental error. Numbers above peaks specify their wavenumbers. Graphene peaks are denoted by letters. The data shown are representative of four independent experiments.

SERS spectra were acquired periodically from 0-168 h. Eighty spectra were taken at each time point and the intensities of the peaks were normalized to the graphene G peak (1584 cm⁻¹) and then averaged. Fig. 2.7 shows a plot of the averaged spectra from nine different time points. Spectra acquired after 24 h were identical to those at 24 h. The graphene G peaks for each spectrum have the same intensities because of normalization. Time-dependent spectral changes were observed for peaks produced by Aβ. At 0 h, the

spectrum (red) was dominated by the graphene D and G peaks, but some intensity at 935 cm$^{-1}$ was observed. At 1 h (orange spectrum), the intensity of the 935 cm$^{-1}$ peak had increased substantially and peaks now also could be seen at 850, 1000, 1087, 1124, ≈1175, and 1460 cm$^{-1}$. The intensity of the 935 cm$^{-1}$ peak relative to the other peaks was lower at 2 h (dark green spectrum). Relative peak intensities differed at a number of time points, which suggested that unique populations of assemblies were being detected.



**Fig 2.8.** Decision tree. Spectral data to be included in PCA analyses of Aβ assembly were determined using decision trees. All of the spectral intensity data at each time point was considered. Recursive cost function analysis determined which nodes would be useful in spectral differentiation. Nodes are presented with the least costly at the top and lower ranked nodes below.

To determine whether the spectral differences observed could distinguish different assembly states, we performed PCA analyses to visualize the data. We employed decision trees to identify peaks that would be most useful in differentiating among time points[20]. Decision trees are particularly suitable for this purpose as they provide simple measures of attribute importance that can be used to rank the importance of the peaks. The C4.5 algorithm was used to prioritize key peaks after examination of all non-graphene peaks with signal-to-noise ratios greater than 10, which comprised peaks at wavenumbers 559, 575, 639, 650, 671, 823, 850, 935, 982, 1000, 1087, 1124, 1190, 1474 and 1612). Nodes in the decision trees (Fig 2.8) were generated by choosing the attribute of the data that most effectively split the tree with the highest normalized information gain. These nodes comprised peak intensities at 823, 850, 935, 1000 and 1124 cm$^{-1}$.



**Fig 2.9.** PCA analysis. Plot of principal components 1 and 2 from analysis of spectra acquired during Aβ42 assembly.

We performed PCA analysis using intensities observed at the five different wavenumbers. PC1 and PC2 contributed 77.7% and 8.5%, respectively, to the total variance of the data. Cluster analysis in the PC1:PC2 plane (Fig. 2.9) revealed a striking time-dependence, and thus assembly-stage dependence, of the locations of the data clusters. Each cluster comprises 67% of the data from a particular time point. Centroids for each cluster in the PC1 dimension was determined by averaging *all* data points. For ease of examination, the clusters have been delimited by color-coded elliptical boundaries, each color representing a specific time. To determine whether the differences in centroid positions were significant, paired Student's t-tests were performed among all pairs of centroids at all times (Table 2.2). With the exception of differences between 4 and 6 h, 6 and 8 h, and 24 and 48 h, all differences were highly significant ($10^{-39} < p < 10^{-3}$). However, if we consider differences in the PC2 dimension, the p-values of 4h vs. 6 h and 6 vs. 8 were highly significant (p=0.002 and p=0.037, respectively). The centroids at 24 and 48 h were almost identical (p=0.99), as they were in the PC1 dimension (p=0.94). This suggests that the assembly process was complete by 24 h.

**Table 2.2.** Significance of differences in centroid positions. Paired Student's *t*-tests were performed to determine if positions of cluster centroids in the PC1 dimension were significantly different (p<0.05). The table lists the *p*-values obtained. Insignificant differences are bolded. "**X**" indicates positions of self-comparison.

| Time | 0 | 1 | 2 | 3 | 4 | 6 | 8 | 12 | 24 | 48 |
|------|---|---|---|---|---|---|---|----|----|----|
| 0 | X | | | | | | | | | |
| 1 | $5.7 \times 10^{-10}$ | X | | | | | | | | |
| 2 | $5.6 \times 10^{-17}$ | $5.1 \times 10^{-14}$ | X | | | | | | | |
| 3 | $6.0 \times 10^{-22}$ | $7.1 \times 10^{-27}$ | $4.9 \times 10^{-12}$ | X | | | | | | |
| 4 | $1.5 \times 10^{-28}$ | $6.7 \times 10^{-34}$ | $1.1 \times 10^{-25}$ | $1.5 \times 10^{-16}$ | X | | | | | |
| 6 | $9.0 \times 10^{-29}$ | $5.2 \times 10^{-32}$ | $1.9 \times 10^{-24}$ | $1.4 \times 10^{-16}$ | **0.17** | X | | | | |
| 8 | $4.6 \times 10^{-30}$ | $3.5 \times 10^{-35}$ | $6.7 \times 10^{-28}$ | $6.1 \times 10^{-20}$ | 0.02 | **0.49** | X | | | |
| 12 | $7.3 \times 10^{-35}$ | $1.3 \times 10^{-47}$ | $2.4 \times 10^{-42}$ | $3.6 \times 10^{-35}$ | $1.2 \times 10^{-9}$ | $2.2 \times 10^{-5}$ | $1.6 \times 10^{-4}$ | X | | |
| 24 | $1.2 \times 10^{-39}$ | $3.3 \times 10^{-63}$ | $2.9 \times 10^{-62}$ | $9.0 \times 10^{-57}$ | $5.6 \times 10^{-23}$ | $2.0 \times 10^{-15}$ | $2.7 \times 10^{-15}$ | $2.2 \times 10^{-10}$ | X | |
| 48 | $1.2 \times 10^{-39}$ | $2.9 \times 10^{-63}$ | $2.4 \times 10^{-62}$ | $7.3 \times 10^{-57}$ | $4.8 \times 10^{-23}$ | $1.8 \times 10^{-15}$ | $2.4 \times 10^{-15}$ | $1.8 \times 10^{-10}$ | **0.94** | X |

The fact that nine clusters were located at different positions in the PC1:PC2 plane shows that at least nine different assembly states were differentiated by SERS/PCA. It is possible that more unique states exist. These could be determined by sampling the assembly process at additional times prior to 24 h. We interpret the clustering as indicative of populations of assemblies that differ between themselves in both the types and relative amounts of different conformers present. Each centroid thus represents a population-average conformer. Overlaps among clusters indicate some population-average conformational similarity. The time-dependence of ellipse centroid

position displayed an amplitude in the PC1 dimension that was ~3-fold larger than that in the PC2 dimension.

Analysis of the component vectors comprising PC1 and PC2 shows that, from 0 h to 4 h, the positions of ellipse centroids within PC2 space are determined primarily by decreases in $\vec{V_4}$, corresponding to peak intensity at 1000 cm$^{-1}$. This suggests that Lys or Asn residues, which resonate at 1000 cm$^{-1}$, are oriented with the graphene surface in a manner that is sub-optimal with respect to Raman signal production. This orientation difference likely reflects conformational changes during peptide assembly. Fig2.4 highlights regions of the Aβ structure wherein vibrational mode intensity differences are noted. Between 4-8 h, ellipse position is determined primarily by $\vec{V_3}$, indicating an increase in the intensity of the 935 cm$^{-1}$ Raman peak, which is produced by carbon-carbon bond resonances[21]. As discussed above, this peak intensity change likely reflects changes in the interaction of the peptide with the graphene due to peptide assembly. After 8 h, the predominant contributor to ellipse position again is $\vec{V_4}$, which shows that the peak intensity of the 1000 cm$^{-1}$ Raman peak decreases. After 24 h, ellipse position and shape do not change substantially, suggesting that the structures of the assemblies producing the Raman spectra are end-state products. This supposition was consistent with results of the EM analysis, which showed no substantial morphological changes after 24 h. A general feature of the time-dependence of ellipse position is that it increases monotonically in the PC1 dimension. This shift is primarily due to the change

in $\vec{V_1}$, $\vec{V_2}$, and $\vec{V_5}$ during peptide assembly. Brown arrows in Fig. 2.9 are the projections of vectors from the initial five-dimensional space into the PC1:PC2 plane. The ~830 cm$^{-1}$ (sideband at ~850 cm$^{-1}$) and 1124 cm$^{-1}$ vibrations are characteristic of Tyr, and of Val and Ile, respectively. The most likely explanation for the increased peak intensities observed at these wavenumbers is the orientation of peptide segments containing these amino acids relative to the nanopyramids. The two most important orientational factors are the proximity of a peptide segment to a hot spot and the conformation of the peptide at that location. Both factors determine peak intensities because lower EM enhancements occur outside the hot spots and the tertiary and quaternary structures of peptide monomers and higher-order assemblies affect the proximity of the resonant bonds to the surface of the hot spot[10]. This is critically important because of the distance dependence of the SERS signal[10]. Through an analysis of the peaks not traditionally thought to report on secondary structure *per se*, we were indeed able to distinguish changes in the structures formed by Aβ42 during aggregation. The complexity of the protein spectra contains a vast array of information, with individual amino acids contributing both within and outside of the amide I, II or III regions[22].

## 2.5 Correlation of CD, TEM, and PCA analyses

When we compare the PCA data (Fig. 2.9) with those obtained by CD (Fig. 2.5), we note that from 0-4 h, the contributions to the CD spectra of their SC component increases monotonically from $[\Theta]_{198}$ = -50 to -38 deg cm$^2$ dmol$^{-1}$, consistent with a

peptide folding process. During this time period, no negative inflection between 215-220 cm$^{-1}$ (β-sheet wavenumbers) is observed and TEM images reveal that fibril formation is initiated. A monotonic decrease in PC2 corresponds to these events. During the first 3 h, the sizes of the ellipses decrease as well, which is consistent with a folding process that decreases the conformational space of the peptide. Taken together, the data suggest that the decrease in PC2 during this time period is indicative of initial Aβ self-association leading to small oligomers and fibril nuclei. As fibril growth occurs, the heterogeneity of assemblies increases, which explains why the 4 h cluster is larger in area than those at 1-3 h. This growth period is reflected in a modest increase in the rate of change in $[\Theta]_{198}$ (~4 deg cm$^2$ dmol$^{-1}$ h$^{-1}$ compared to an initial rate of ~3 deg cm$^2$ dmol$^{-1}$ h$^{-1}$). In addition, between 6-8 h, a negative inflection at $[\Theta]_{216}$ appears, which monotonically decreases over time, consistent with the increased β-sheet secondary structure produced by fibril formation. Ellipse position in the PC2 dimension rises concurrently (Fig. 2.9). From 8-24 h, progressive increases in β-sheet (CD) and fibril content (TEM) correlated with decreases in cluster position in the PC2 dimension. Increasing $\overrightarrow{V_4}$ magnitudes were the prime contributor to the monotonic decrease in ellipse position in the PC2 dimension. We note that the centers of the ellipses of 24 h and 48 h were in essentially the same position in PC1:PC2 space, but ellipse area decreased markedly during this time period, which suggests increased structural homogeneity of the peptide assemblies. This effect may be related to fibril aggregation, which commonly is observed following fibril growth phases[23].

## 2.6 Conclusion

We show that Raman spectra obtained using a graphene-gold hybrid plasmonic platform, in combination with PCA analysis, enables facile distinction between Aβ40 and Aβ42, the peptide isoforms associated with classical vascular AD (Aβ40) and parenchymal (Aβ42) plaques, respectively, in AD. We further show that the approach is capable of revealing assembly-dependent changes in peptide conformation and self-association. Correlation of these spectral changes with CD and TEM data allow regions in PCA space to be linked to specific populations of Aβ assemblies. What may be particularly important is the observation of a minimum of nine differentiable clusters within PCA space, which reflect at least nine differentiable assembly states in the fibril formation pathway. Because spectral changes can be linked to changes in resonances of specific amino acids within the Aβ peptide, future sited-directed amino acid substitution studies of these individual states may provide new insights into the roles of different amino acids in stabilizing or destabilizing these states. Thus, coupled with the label-free, single molecule sensitivity of SERS, the SERS/PCA approach should prove useful for determining structure activity relationships, suggesting target sites for drug development, and for testing the effects of such drugs on the assembly process. The approach also could be of value in other systems in which assembly-dependent changes in protein structure correlate with the formation of toxic peptide assemblies.

## 2.7 References

1.    Roychaudhuri, R., Yang, M., Hoshi, M. M., and Teplow, D. B. (2009) Amyloid β-protein assembly and Alzheimer disease. J Biol Chem 284, 4749-4753

2.    Kneipp, K., Wang, Y., Kneipp, H., Perelman, L. T., Itzkan, I., Dasari, R., and Feld, M. S. (1997) Single molecule detection using surface-enhanced Raman scattering (SERS). Physical Review Letters 78, 1667-1670

3.    Beier, H. T., Cowan, C. B., Chou, I.-H., Pallikal, J., Henry, J. E., Benford, M. E., Jackson, J. B., Good, T. A., and Coté, G. L. (2007) Application of surface-enhanced Raman spectroscopy for detection of beta amyloid Using nanoshells. Plasmonics 2, 55-64

4.    Benford, M. E., Chou, I.-H., Beier, H. T., Wang, M., Kameoka, J., Good, T. A., and Cot, G. L. (2008) In vitro detection of β amyloid exploiting surface enhanced Raman scattering (SERS) using a nanofluidic biosensor. Proc. SPIE 6869, 68690W–68690W–68695

5.    Bhowmik, D., Mote, K. R., MacLaughlin, C. M., Biswas, N., Chandra, B., Basu, J. K., Walker, G. C., Madhu, P. K., and Maiti, S. (2015) Cell-membrane-mimicking lipid-coated nanoparticles confer Raman enhancement to membrane proteins and reveal membrane-attached amyloid-β conformation. Acs Nano 9, 9070-9077

6.    Choi, I., Huh, Y. S., and Erickson, D. (2012) Ultra-sensitive, label-free probing of the conformational characteristics of amyloid β aggregates with a SERS active nanofluidic device. Microfluidics and Nanofluidics 12, 663-669

7.    Voiciuk, V., Valincius, G., Budvytyte, R., Matijoska, A., Matulaitiene, I., and Niaura, G. (2012) Surface-enhanced Raman spectroscopy for detection of toxic amyloid beta oligomers adsorbed on self-assembled monolayers. Spectrochim Acta A Mol Biomol Spectrosc 95, 526-532

8.    Nabers, A., Ollesch, J., Schartner, J., Kotting, C., Genius, J., Hafermann, H., Klafki, H., Gerwert, K., and Wiltfang, J. (2016) Amyloid-β-secondary structure distribution in cerebrospinal fluid and blood measured by an immuno-infrared-sensor: A biomarker candidate for Alzheimer's disease. Anal Chem 88, 2755-2762

9.    Larkin, P. (2011) Infrared and Raman spectroscopy: Principles and spectral interpretation, Elsevier, Waltham, San Diego, Oxford, Amsterdam

10.  Wang, P., Liang, O., Zhang, W., Schroeder, T., and Xie, Y. H. (2013) Ultra-sensitive graphene-plasmonic hybrid platform for label-free detection. Adv Mater 25, 4918-4924

11.  Yu, X., Hayden, E. Y., Wang, P., Xia, M., Liang, O., Bai, Y., Teplow, D. B., and Xie, Y.-H. (2017) Quantification characteristics of a graphene-gold hybrid plasmonic SERS platform and its use in studies of Alzheimer amyloid β-protein. submitted

12. Ferrari, A. C., Meyer, J. C., Scardaci, V., Casiraghi, C., Lazzeri, M., Mauri, F., Piscanec, S., Jiang, D., Novoselov, K. S., Roth, S., and Geim, A. K. (2006) Raman spectrum of graphene and graphene layers. Physical Review Letters 97

13. Wang, P., Zhang, W., Liang, O., Pantoja, M., Katzer, J., Schroeder, T., and Xie, Y. H. (2012) Giant optical response from graphene--plasmonic system. ACS Nano 6, 6244-6249

14. Bitan, G., Kirkitadze, M. D., Lomakin, A., Vollers, S. S., Benedek, G. B., and Teplow, D. B. (2003) Amyloid β-protein (Aβ) assembly: Aβ40 and Aβ42 oligomerize through distinct pathways. Proc Natl Acad Sci U S A 100, 330-335

15. Medek, A., Hajduk, P. J., Mack, J., and Fesik, S. W. (2000) The use of differential chemical shifts for determining the binding site location and orientation of protein-bound ligands. J Am Chem Soc 122, 1241-1242

16. Abdi, H., and Williams, L. J. (2010) Principal component analysis. Wiley interdisciplinary reviews: computational statistics 2, 433-459 %@ 1939-0068

17. You, Z. H., Lei, Y. K., Zhu, L., Xia, J., and Wang, B. (2013) Prediction of protein-protein interactions from amino acid sequences with ensemble extreme learning machines and principal component analysis. BMC Bioinformatics 14 Suppl 8, S10

18. Reynaldo, J., and Santos, A. (1999) Cronbach's alpha: A tool for assessing the reliability of scales. J Extension 37, 1-5

19. Hayden, E. Y., Yamin, G., Beroukhim, S., Chen, B., Kibalchenko, M., Jiang, L., Ho, L., Wang, J., Pasinetti, G. M., and Teplow, D. B. (2015) Inhibiting amyloid β-protein assembly: Size-activity relationships among grape seed-derived polyphenols. J Neurochem 135, 416-430

20. Quinlan, J. R. (1993) C4.5 : programs for machine learning, Morgan Kaufmann Publishers, San Mateo, Calif.

21. Movasaghi, Z., Rehman, S., and Rehman, I. U. (2007) Raman spectroscopy of biological tissues. Appl Spectr Rev 42, 493-541

22. Movasaghi, Z., Rehman, S., and Rehman, I. U. (2007) Raman spectroscopy of biological tissues. Appl Spectrosc Rev 42, 493-541

23. Qiang, W., Yau, W. M., and Tycko, R. (2011) Structural evolution of Iowa mutant β-amyloid fibrils from polymorphic to homogeneous states under repeated seeded growth. J Am Chem Soc 133, 4018-4029

24. Walsh, D. M., Lomakin, A., Benedek, G. B., Condron, M. M., and Teplow, D. B. (1997) Amyloid β-protein fibrillogenesis. Detection of a protofibrillar intermediate. J Biol Chem 272, 22364-22372

# Chapter 3

# Quantification Capability

## 3.1 Introduction

### *3.1.1 Quantification in bio-sensing*

The development of sensitive techniques for the detection and quantitative analysis of bio-molecules is important for trace element detection, environmental monitoring, and early stage diagnosis and treatment of diseases[1-4].

According to US Food and Drug Administration (FDA), bioanalytical method development involves optimizing the procedures and conditions involved with extracting and detecting the analytes and 10 bioanalytical parameters are listed for the optimization:

Reference standards, critical reagents, calibration curve, quality control samples (QCs), selectivity and specificity, sensitivity, accuracy, precision, recovery and stability of the analyte in the matrix.

Accuracy, precision and recovery are vital for the procedure and are highly related to sample quantification. Evaluating the accuracy and precision across the quantitation range is essential to determine whether the method is ready for validation. Having good quantification capability also involves analyzing replicate QCs at multiple concentrations across the assay range. Quantification becomes a good screening criterion for bioanalytical methods.

A wide range of methods provide quantification capability. Detection methods include high-performance liquid chromatography (HPLC)[5], liquid chromatography mass spectrometry (LCMS)[6], Ring Resonator biosensor and enzyme-linked immunosorbent assays (ELISA)[7]. A limit of detection (LOD) of 0.1 ng/mL has been achieved with these platforms. Table 3.1 shows limit of detection as well as the limitation for these popular quantification technologies.

**Table 3.1**. Summary of current quantification methods

| Method | Limit of detection | Limitations | Ref |
|---|---|---|---|
| **HPLC** | 0.1ng/uL for vitamin B12 | Less separation efficiency. | Luo 2016 |

| | | | |
|---|---|---|---|
| **LCMS** | 5nmol/L for polyglutamates | Lack of traceability leads to imprecision | E Den Boer 2013 |
| **ELISA** | 5ng/mL for urine sample | Limited by antigen in the sample | Dheda 2010 |
| **Ring Resonators** | 112nm/RIU for Aflatoxin | Output uncertainty | Guider 2015 |

## 3.1.2 Quantification using SERS

Surface enhanced Raman scattering (SERS) is a method that has gained increasing notice because of its ability to achieve single molecule detection with high molecular specificity[8-11] without the use of biological labels. Recent advances in nanotechnology have led to many SERS-based analytical applications. For example, self-assembled monolayer (SAM)-coated colloidal gold platforms are able to detect Rhodamine 6G in the concentration range of 0.1-5 μM[7]. Metallic glassy nanowire arrays (MGNWAs) have a dynamic range of 1-10 nM for Rhodamine B[12]. Gold nanoparticles allow detection of glucose in the concentration range 0.5-32 mM[13]. A summary for some current SERS quantification platforms is listed in Table 3.2.

**Table 3.2.** Summary for current SERS quantification platform.

| Platform type | Analyte | Dynamic range | Ref |
|---|---|---|---|

| | | | |
|---|---|---|---|
| Ag nanoparticle | Nicotinamide | 0.1-1mM | Castro JL. El al (2013) |
| Au nanoparticle | Glucose | 0.5-32 mM | Quyen el al (2013) |
| Graphene oxide and AgNP | R6G | 1nM-10uM | Ren el al (2014) |
| Sandwich assay with AuNP | MUC4 (protein) | 0.01-10ug/mL | Wang et al (2011) |
| SERS tagged Ag colloidal | Mouse IgG | 0.1-3ng/mL | Dou et al (1998) |

The listed quantification methods are mainly based on the linear relationship between SERS peak intensity with analyte concentration. However, SERS is a complicated method and peak intensity may change due to a wide range of parameters (include but not limited to concentration, incident light intensity, hotness of local hotspot .etc.).This means that signal intensities measured at different points on a surface vary even if analyte concentration is constant. SERS signal intensity *per se* does not have a one-to-one correlation with analyte concentration.

Such correlation could be established only if the EM field intensity at an individual plasmonic hot spot could be determined independent of analyte concentration. Prior studies have employed the assembly of marker molecules as internal standards. However, the non-planar topology typical of plasmonic surfaces, which can create local inhomogeneities in marker concentration significantly confound the situation. In

addition, the relatively large size of typical SERS internal standards (~50 nm)[7,15] relative to hot spot size can interfere with adsorption of the actual analyte at the hot spot, precluding Raman signal production by the analyte.

Inhomogeneity of the platform together with the large internal standard size jeopardize the simple proportionality between intensity versus concentration and thus making the analysis argumentative if no reference is presented for the quantification platform.

If a built-in marker of local electromagnetic (EM) field intensity, which is known to vary substantially from location to location on plasmonic surfaces currently in use[14] is created and combined to the platform, signal intensity can solely depend on the analyte concentration as the local EM field intensity has been ruled out from the matrix. We report here for the first time, the quantification ability of an ultra-sensitive graphene-plasmonic hybrid platform that largely eliminates these sources of variability. The platform incorporates a single atomic layer of graphene overlaying a gold surface consisting an array of pyramids, each of which is ~200 nm in width and height. The periodicity and size of individual pyramids are chosen for the optimization of surface plasmon resonance by laser excitation at 785 nm. The novelty of this work is that the single layer graphene serves as an internal standard and provides Raman signal enhancement, which allows for accurate quantitation, not previously feasible [2]. SERS enhancement results primarily through an EM mechanism (up to $10^{14}$)-fold. A chemical

mechanism may also contribute to the enhancement, but to a much lesser degree (∼10–100-fold)[3] . We use the amyloid β-protein (Aβ), a well-studied pathologic agent of Alzheimer's disease[17-19], as an example biologically relevant analyte to assess the potential of the hybrid platform for quantification and subsequent study of Aβ assembly dynamics. We observe two complementary quantification modes. The first (high analyte concentration regime) relates analyte concentration to the SERS peak intensity at individual SERS hotspots. The second (low concentration regime) relates analyte concentration to the probability of observing any Raman signal at any hotspot. In combination, these two modes enable analyte detection in a concentration range spanning seven orders of magnitude ($10^{-18}$-$10^{-15}$M, $10^{-13}$-$10^{-11}$M).

Chapter 3.2 present the rational of choosing dried Amyloid beta 42 as sample analyte and experimental details. Chapter 3.3 demonstrates two different types of quantification mechanism using SERS hybrid platform. Chapter 3.4 details a novel method to statistically analyze the number of monomer in each hotspot for extremely detailed analysis.

## 3.2   Experimental procedure

### 3.2.1 Hybrid platform as a quantification media

The graphene-gold hybrid platform fabrication is based on sphere lithography, as previously reported[20]. The periodic gold nano-pyramid structure with tunable size and

sharpness is fabricated by a wafer-scale bottom-up templating technology. Spin-coated on (001) silicon wafers, close-packed monolayer polystyrene balls with a diameter of 200 nm serve as templates. Monolayer graphene is grown by chemical vapor deposition (CVD) and the solution transferred onto the gold tipped surface using polymethyl methacrylate (PMMA) backing followed by PMMA removal subsequent to the transfer. Platforms can be fabricated with user-determined areas. We typically use platforms of ~1 cm$^2$. The pyramids form a quasi-periodic array of hexagonal arrangement that is uniformly distributed across the entire sample surface of 1 cm $\times$ 1 cm area. Because of the way the pattern is generated (self-assembly of polystyrene balls), variations in the spacing between pyramids, and the sizes of the pyramids themselves, can vary. This variance has been estimated to be ±30 nm.

To grow the graphene, a 25 μm thick copper foil from Alfa Aesar (catalog #13382) is cut into a 5 cm squares. The copper foil is loaded onto the center of a quartz CVD chamber, the furnace is heated to ≈1025°C under the flow of hydrogen gas (~1000 sccm). After 30 min annealing, the CVD growth was carried out with 20 Torr total pressure with methane gas (~20 sccm) and hydrogen gas (~1000 sccm) for 15 minutes. The chamber then was cooled to room temperature.

Scanning Electron Microscopic (SEM) analysis was performed using FEI Nova Nano SEM 230 instrument, an accelerating voltage of 10 kV, and a beam current of 0.14 A. After the production of the gold pyramids, as described above, the substrate was

mounted onto an SEM stub using double sided adhesive tape. Imaging was performed at magnifications ranging from 30,000 to 200,000.



**Fig 3.1.** SEM image of the gold nano pyramids. The magnifications are 200,000 (left panel) and 30,000 (right panel).

We used scanning electron microscopy (SEM) to examine the surface morphology of the hybrid platfomrs after facrication to confirm the presence of the pyramid structures. We observed that the pyramids form a quasi periodic array of hexagonal arrangement uniformly distributed over the surface (Fig 3.1). Some variation ($\pm$ 30nm) is observed in the spacing between pyramids and the sizes of the pyramids, as expected based on the fabrication process.

3.2.2 Amyloid $\beta$ as a quantification analyte

To examine the quantification characteristics of the hybrid platform, we studied the 42-amino acid form of A$\beta$, A$\beta$42, which is thought to be a seminal pathologic agent in

AD and is an important disease biomarker[24].

Aβ was synthesized in the UCLA Biopolymer facility and then purified and characterized, as described previously. Briefly, peptide synthesis was performed on an automated peptide synthesizer (model 433A, Applied Biosystems, Foster City, CA, USA) using 9-fluorenylmethoxycarbonyl-based methods on preloaded Wang resins[21]. Aβ was purified to >97%, using reverse-phase high-performance liquid chromatography (HPLC). Quantitative amino acid analysis and mass spectrometry yielded the expected composition and molecular weight. Purified peptides were stored as lyophilizates at –20 ℃.

Aβ was prepared by dissolution in 10% (v/v) 60mM NaOH, 45% (v/v) Milli-Q water, and 45% (v/v) 22.2 mM sodium phosphate, pH 7.4, to yield a nominal Aβ concentration of 1 mg/mL in 10 mM sodium phosphate, pH 7.4. The Aβ solution then was sonicated for 1 min in a bath sonicator (Branson Model 1510, Danbury, CT, USA) and filtered through a prewashed 30,000 molecular weight cut-off Microcon centrifugal filter device (Millipore, Billerica, MA, USA) for 15 min at 16,000 $\times$ g. The concentration of Aβ in the eluate was determined using UV absorbance ($\varepsilon_{280} = 1280$ cm$^{-1}$ M$^{-1}$). The peptide was diluted with 10 mM sodium phosphate, pH 7.4, to a final concentration of 20 μM before use. Serial dilutions then were done in 10 mM sodium phosphate, pH 7.4. All measurements were performed at 22 ℃. This protocol reproducibly yields aggregate-free Aβ monomer in rapid equilibrium with low order oligomers, which is

termed low molecular weight Aβ[22].

For all test done, a 20 µl volume of Aβ42 was pipetted onto the center of the platform and then immediately dried *in vacuo*. Raman spectra were measured using a Renishaw inVia microscope under ambient conditions. Excitation was accomplished using a GaAlAs diode laser of wavelength 785 nm. A laser power of 0.5 mW, a grating of 1800 lines/mm, and an objective lens of 50× were used. A step size of 200 nm was used for Raman mapping. Raman data were analyzed using Renishaw WiRE 4.2 software. Strong hotspots appear in between pyramids and at their apices.

We initially applied a 20 µM solution of freshly prepared, unaggregated, low molecular weight Aβ[25] to our platforms. Aβ42 is known to aggregate into oligomers and fibrils over time so to ensure that our starting samples did not aggregate during preparation for SERS we prepared the samples rapidly (<10 min) at low temperature (4 ℃). We acquired spectra immediately after preparation and periodically thereafter. Examination of these spectra showed that no observable aggregation occurred[28] within 10 min[27]. (spectral figure: Fig 3.2, principal component analysis: Fig 3.3).

**Fig 3.2.** SERS spectra of Aβ42. spectra were acquired periodically during a 2 h incubation of a 20 μM solution of Aβ42 in 10 mM sodium phosphate buffer, pH 7.4, at 37 ℃.



**Fig 3.3.** Principal component analysis revealed four distinct clusters (the 0 and 10 min data clustered together suggesting that no significant spectral changes occurred between these times).

In addition, when experiments at Aβ concentrations in the sub-micromolar regime are done, rates of simple collision-induced aggregation or nucleation-dependent polymerization are so low that no substantial aggregation occurs. Fig 3.4 shows a typical SERS spectrum of Aβ42. We also acquired Raman spectra for the sodium phosphate buffer without Aβ42, and did not observe any Raman peaks. This indicated that all the peaks we observed are from graphene or Aβ42. Several characteristic Raman peaks were observed, including those due to Tyr (823 and 850 cm$^{-1}$), carbon-carbon (C-C) stretching (935 cm$^{-1}$), Phe (982 and 1450 cm$^{-1}$), Lys or Asn (1087 cm$^{-1}$), Val or Ile (1124 cm$^{-1}$), and graphene D (1350 cm$^{-1}$), and G peaks (1580 cm$^{-1}$)[29,30]. We note that an amide I peak (1650 cm$^{-1}$) appears near the graphene G peak, but no overlap is observed.



**Fig 3.4.** Aβ42 was prepared at a concentration of 20 μM, pH 7.4, and applied to the platform. Abscissa indicates wavenumber (cm$^{-1}$). Peaks were assigned based on published result.

Wavenumber assignments are: 559, aliphatic; 575, C-C bond stretching mode; 823, out-of-plane ring breathing vibration or double Tyr (Tyr2); 850, single bond stretching for Tyr and Val; 935, number of carbon -carbon bonds of protein backbone or Gly; 982, C-C stretching in β-sheets or part of Phe; 1000, Lys or Asn; 1124, Val or Ile; 1360, graphene D peak;1450, CH2 bending or Phe; 1580, graphene G peak.

## 3.2.3 Graphene as a quantification gauge

To establish the quantitative ability of the hybrid platform, Aβ42 at concentrations ranging from $10^{-21}$–$10^{-9}$ M was applied to the substrate and spectra were acquired (Fig 3.5). The spectra were normalized to the graphene G peak (1580 cm$^{-1}$) to account for any variation in the local electromagnetic field intensity among the various hot spots, allowing us to correlate the Raman peak intensity with Aβ42 concentration.

The graphene G peak intensity is correlated to both the graphene configuration and EM field intensity. Prior to graphene transfer to the pyramid substrate, two confirmations were carried out to ensure that it existed exclusively as a single atomic layer[31]. Changes in the graphene G peak intensity thus should arise solely from changes in the EM field and thus can be used to normalize protein peak intensities obtained across the substrate surface.

**Fig 3.5.** Spectra of Aβ42 at concentrations ranging from $10^{-13}$–$10^{-9}$ M. Spectra from concentrations of $10^{-15}$ and $10^{-14}$ were obtained but they are not shown because they are essentially flat in this representation.

The graphene G peaks among spectrum superimpose as a result of the normalization. The 1360 cm$^{-1}$ peak is the graphene D peak, which results from the breathing modes of $sp^2$ atoms in the carbon ring structure[32]. The D-peak is related to defects in the graphene, especially graphene folds formed when the nearly planar graphene is overlaid on the pyramids of the platform. As such, the D-peak is a function of surface topology and not suitable for use as a normalization signal.

## 3.3 Quantification mechanism

We present here two types of quantification mechanism as shown in Fig 3.6. In analyte concentration regimes in which essentially all hot spots contain at least one analyte molecules, Raman signal intensity depends on analyte number concentration.

**Fig 3.6.** (A-C) Dependence of Raman signal intensity on analyte number concentration at hot spots. (D-E) Hot spot occupancy versus analyte concentration. Analyte molecules within a hot spot are shown in red. Analytes outside of hot spots are blue.

(A). Signal intensity thus increases with analyte concentration (B) until hot spots are saturated with analytes (C), at which time accurate determination of concentration is no longer possible because not all analytes are associated with hot spots. (D-E) Hot spot occupancy versus analyte concentration. As illustrated in (D), at lower concentrations, hot spot occupancy is <100% and peak intensities begin to correlate with the probability of an individual protein molecule being collocated with a hot spot, as opposed to the number of molecules at each hot spot (as in the high concentration regime). In the low concentration regime, quantification is accomplished by determination of occupancy frequency *per se*. Panel (E) illustrates a

concentration regime in which most or all hot spots contain at least one analyte molecule. Hot spot occupancy thus is ≅100% and signal intensity depends on the number of analyte molecules. Analyte molecules within a hot spot are shown in red. Analytes outside of hot spots are blue.

### 3.3.1 Intensity based quantification (high concentration)

We observed no qualitative differences among the spectra obtained at different Aβ42 concentrations. Instead, as expected, a direct relationship between peak intensity and concentration was seen. At a single hot spot, at which multiple analytes can bind, Raman signal intensity is the sum of the individual intensities of all the Raman active analytes present. Increases in signal intensity with analyte concentration thus are observed until the limited volume of the hot spot is fully occupied by analyte molecules, after which increases in analyte concentration do not lead to increased peak intensity. This is seen clearly in Fig 3.7, in which the concentration-dependence of peak intensity at 935 cm$^{-1}$ is shown.

**Fig 3.7.** Normalized peak intensity (AU) of the 935 cm$^{-1}$ peak. All points are the averages of three replicates. Red bars signify standard deviations. If error bars are not visible, this indicates that the size of the standard deviation is less than the size of the data point.



**Fig 3.8.** Log-log plot of the data from panel B. The blue line was produced by linear regression analysis (R=0.97).

The data produce a sigmoidal curve within which a quasi-linear region is seen extending from $\sim 10^{-13}$–$10^{-11}$ M. The linearity within this region is more apparent from examination of a log-log plot (Fig 3.8), which we utilize because of the very broad

59

concentration regime studied. Linear regression analysis of these data yields a straight line with a correlation coefficient of 0.97.

From the peaks emanating from protein, the 935 cm$^{-1}$ peak has the highest relative peak intensity and the lowest signal/noise ratio. As such, experimental noise has less impact on its intensity, and thus the linear relationship between intensity of the 935 cm$^{-1}$ peak and concentration provides a more accurate quantitation than the use of more isolated vibrational modes at 1087 cm$^{-1}$ and 1124 cm$^{-1}$ (Fig 3.9).

**Fig 3.9.** Concentration-dependence of Raman signal intensities at 1087 cm$^{-1}$ and 1124 cm$^{-1}$. Log-log plot of normalized peak intensity (AU) of the 1087 cm$^{-1}$ peak (left panel) and 1124 cm$^{-1}$ (right panel). Data points are the average of three independent experiments, for which each includes >200 individual scans. Black bars signify standard deviations. The red line was produced by linear regression analysis for the points between 10$^{-13}$ and 10$^{-11}$ M ($R^2$=0.85 and =0.76 for 1087 cm$^{-1}$ and 1124 cm$^{-1}$ respectively).

We observed increasing protein concentration towards the perimeter of the applied droplet, induced by the surface tension of the liquid during drop casting, likely due to the "coffee ring" effect. We find that the concentration change is not high enough to influence the linear relationship between protein concentration and peak intensity (c.f. error bars on Fig 3.7) One explanation for the linearity of increasing SERS intensity with increasing concentration is that many protein molecules can fit within a single hot spot before it is filled, which would be difficult to observe if the hot spot size is closer to the size of a single analyte molecule.

## 3.3.2 Frequency based quantification (low concentration)

As analyte concentration decreases, not all hot spots will have adsorbed analytes and a direct relation-ship between peak intensities and analyte concentration does not exist.

For this reason, instead of quantifying signal intensities at individual hot spots, we implement a quantification method that considers instead the frequency of hot spots



from which Raman spectra signals of the analytes are detectable. To determine this frequency, we scanned relatively large areas (~50 μm × 50 μm) of the hybrid platform using Raman mapping at a step size of 1 μm (i.e., a 1 μm2 area of pixels for each measurement).

**Fig 3.10.** Intensity mapping of the 935 cm$^{-1}$ peak at concentrations of $10^{-13}$, $10^{-15}$, and $10^{-17}$ M. The step size of the mappings was 1 μm and 2600 spectra were acquired at each concentration.

This scanning encompasses the entire area of the original droplet, including the outer perimeter and inside of the dried ring, so that we get a representative sampling of the protein concentration. We performed this scanning on platforms on which we applied Aβ42 in concentrations ranging from 10-18–10-10 M. Fig 3.10 shows heat maps of the intensity data collected at concentrations of $10^{-13}$, $10^{-15}$, and $10^{-17}$ M. Inspection reveals

a substantial concentration-dependent decrease in frequency. A plot of the frequency

distribution (Fig 3.11) shows that no signals were observed at Aβ42 concentrations of

10-20 or 10-19 M. A direct relationship between frequency and concentration was

observed was observed between $10^{-18}$–$10^{-15}$ M (Fig 3.11, solid line).



**Fig 3.11.** A log-log plot of concentration (M) versus hot spot signal frequency (%) determined in the concentration range of $10^{-20}$–$10^{-10}$ M. For ease of visualization, points at $10^{-20}$ M and $10^{-19}$ M, which had zero intensity, are plotted with frequencies of 0.0001%. Solid line shows result of linear regression analysis in the concentration regime $10^{-18}$–$10^{-15}$ M, inclusive (correlation coefficient R = 0.97).

Above $10^{-15}$ M, a concentration regime is encountered in which the majority of hot

spots have at least one Aβ peptide and increasing Aβ concentration results in an

increase in the number of peptides per pyramid but not in a substantial increase in the

percentage of pyramids with at least one peptide (Table 3.3). Table 3.3 shows the

frequency of observable spectrum in each of the reported concentrations, from two

in-dependent experiments, which is the number of detectable signals divided by the total number of scans across the area examined. If analyte concentration is within the transition region between partial and full hot spot occupancy, simple dilution will allow accurate quantification based solely on occupancy frequency. We note that Pérez-Ruiz et al., in studies determining tau concentrations [33], also have successfully employed a frequency approach (cf. Fig. 3.10 and 3.11 of this manuscript with Fig. 5 of Pérez-Ruiz et al.). This approach enabled attomolar limits of detection depending on whether samples were prepared in buffer (24 aM) or blood plasma (55 aM). Coupled with analogue measurements at higher concentrations, a dynamic concentration range of six orders of magnitude could be obtained. These capabilities compare favorably with our own—a dynamic range of 7 orders of magnitude and a limit of detection of 1 aM.

**Table 3.3:** Hot spot signal frequency versus Aβ42 concentration.

| Concentration (log M) | -20 | -19 | -18 | -17 | -16 | -15 | -14 | -13 | -12 | -11 | -10 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Frequency (Exp. #1) | 0.000 | 0.000 | 0.023 | 0.470 | 4.12 | 38.4 | 72.2 | 84.9 | 95.2 | 100 | 100 |
| Frequency (Exp. #2) | 0.000 | 0.000 | 0.011 | 0.102 | 3.85 | 28.7 | 59.2 | 87.2 | 96.0 | 100 | 100 |

The frequencies (the number of hotspots with a detectable signal, divided by the total

number of scans) of observable signals from hot spots were determined over a concentration range of $10^{-18}$–$10^{-10}$ M in two independent experiments. Frequency (%) = $100 \times$(number of hot spots with detectable signal)/(total number of scans).

It should be pointed out that the true detection limit is when the concentration of the analyte becomes so low that the probability of a single molecule existing within a 20 L droplet is <1. This probability is 0.37 at 10-19 M. This explains why, in practice, no A   signal was observed at lower concentrations. In theory, there is a finite probability that a single molecule will exist within an applied volume at any concentration and that a spectroscopist can spend sufficient time to find its location and signal on the platform.

## 3.4 Hotspot intensities

The signal intensities observed at 935 cm$^{-1}$ in samples analyzed at concentrations of $10^{-17}$, $10^{-15}$, and $10^{-13}$ M were incorporated into a data table (see text for rationale), each element of which represented the intensity from a single hot spot. The total numbers of hot spots at which signals were observed, $n_i$, were 14 ($10^{-17}$ M), 100 ($10^{-15}$ M), and 200 ($10^{-13}$ M). To produce a histogram of intensities, individual intensities were binned using a bin size of 100 μAU. Frequencies were calculated according to the formula $f_i = n_i/n_t \times 100$; in which $f_i$ is percent frequency of occurrence of intensity $i$, $n_i$ is number of observations of intensity $i$, and $n_t$ is total number of intensity observations. The weighted average intensity for the histogram envelope observed in the $10^{-17}$ M sample was calculated according to the formula $I_{avg} = \sum_{a=1*10^{-4}AU}^{\infty} i_a * n_i/n_t$ . Plots

were done using Origin v8.4.

We next sought to determine, in the frequency regime of concentration, how intensity was related to number of analyte molecules per hot spot. To do so, we created histograms of normalized spectral intensities at the same concentrations in Fig. 3.12 ($10^{-17}$, $10^{-15}$, and $10^{-13}$ M). We used the intensity of the 935 $cm^{-1}$ peak (C-C bonds) for this purpose, as this peak had been used to quantify Aβ concentration. Our expectation was that the lowest observed intensity should be produced by a single analyte molecule and that subsequent signal intensities should be integer multiples of that lowest intensity. In Fig. 5, we observed a single node with an average intensity of 1.3 x $10^{-4}$ AU. At a concentration of $10^{-15}$ M, this node also was observed, in addition to prominent nodes at intensities that were double or triple that intensity. This shifting of the overall frequency distribution to higher intensities was seen at $10^{-13}$ M as well, a concentration that produced nodes (blue arrows) at intensities that were 23-30-fold larger than the lowest intensity node, consistent with the conclusion that this distribution reflected hot spots containing 23-30 analyte molecules. When we compared the average signal intensities for each node envelope with those predicted based on multiples of 1.3 x $10^{-4}$ AU per monomer, we observed remarkable agreement (mean and standard deviation of the differences was $0.04 \pm 0.09 \times 10^{-4}$ AU; Table 3.4). These data support the conclusion that we are, at minimum, able to differentiate signals produced by 1-30 peptides per hot spot.

Raman intensities of proteins depend not only on analyte quantity per hot spot, but on the structure of the protein, its orientation relative to the hot spot surface, and the electromagnetic field intensity. Our normalization procedure controls for the latter factor. The variation in the former two factors is reflected in the widths of the overall intensity envelopes observed in the histograms. These increase with concentration, but even at a concentration of $10^{-13}$ M, we see that the envelopes reflect a discrete range of analyte numbers, as opposed to including intensities from the continuum of possible analyte numbers per hot spot. This likely reflects the fact that the application and binding of protein to the matrix of pyramids is consistent with simple laws of mass action.

**Fig 3.12.** Hot spot intensities. The graphene normalized signal intensities of the 935 cm-1 Raman peak acquired at Aβ concentrations of 10-17, 10-15, and 10-13 M are presented in histograms. Axes are frequency (ordinate) and normalized intensity (abscissa). Numbers at blue arrows signify the number of monomers producing the observed intensities.

**Table 3.4.** Agreement of predicted and observed weighted average hot spot intensities. Leftmost column displays number of Aβ monomers per hot spot. "Predicted Intensity" is number of peptides times intensity per peptide (1.3). "Observed Intensity" is experimentally observed intensity.

| Peptides per hot spot | Predicted Intensity[a] | Observed Intensity[a] |
|:---:|:---:|:---:|
| 1 | 1.3 | 1.3 |
| 2 | 2.6 | 2.5 |
| 3 | 3.9 | 3.7 |
| 4 | 5.2 | – |
| 5 | 6.5 | 6.3 |
| 6 | 7.8 | 7.7 |
| 7 | 9.1 | 9.0 |
| 8 | 10.4 | 10.5 |
| 9 | 11.7 | 11.7 |
| 10 | 13.0 | 12.8 |
| 11 | 14.3 | 14.2 |
| 12 | 15.6 | 15.6 |
| 13 | 16.9 | 16.8 |
| 14 | 18.2 | – |
| 15 | 19.5 | – |
| 16 | 20.8 | – |
| 17 | 22.1 | – |
| 18 | 23.4 | – |

| | | |
|---|---|---|
| 19 | 24.7 | – |
| 20 | 26.0 | – |
| 21 | 27.3 | – |
| 22 | 28.6 | – |
| 23 | 29.9 | 30.0 |
| 24 | 31.2 | 31.1 |
| 25 | 32.5 | 32.5 |
| 26 | 33.8 | 33.7 |
| 27 | 35.1 | 34.9 |
| 28 | 36.4 | 36.4 |
| 29 | 37.7 | 37.8 |
| 30 | 39.0 | 38.9 |

## 3.5 Conclusion

This work demonstrates the quantification ability of the graphene-gold hybrid SERS platform using Raman mapping. The platform exhibits a linear relation between peak intensity and concentration at single hot spots (high analyte concentration), as well as a linear relationship between detection frequency and analyte concentration when scanning multiple hotspots (low analyte concentration). The platform is capable of single-molecule detection. The useful dynamic range of the hybrid platform of seven orders of magnitude (3 orders of magnitude for higher concentration and 4 orders of

70

magnitude of lower concentration) offers the possibility that the platform could be useful in a broad range of applications such as early stage diagnosis of Alzheimer's disease.

## 3.6 References

(1)  Coffman VC, Wu J-Q Trends Biochem Sci 2012, 37, 499.

(2)  Hayashi Y, Matsuda R, Maitani T, Imai K, Nishimura W, Ito K, Maeda M Anal Chem 2004, 76, 1295.

(3)  Navratil M, Mabbott GA, Arriaga EA Anal Chem 2006, 78, 4005.

(4)  Pitt JJ Clin Biochem Rev 2009, 30, 19.

(5)  Lemma T, Saliniemi A, Hynninen V, Hytönen VP, Toppari JJ Vib Spectrosc 2016, 83, 36.

(6)  Stephen KE, Homrighausen D, DePalma G, Nakatsu CH, Irudayaraj J Analyst 2012, 137, 4280.

(7)  Lorén A, Engelbrektsson J, Eliasson C, Josefson M, Abrahamsson J, Johansson M, Abrahamsson K Anal Chem 2004, 76, 7391.

(8)  Antonio KA, Schultz ZD Analytical chemistry 2013, 86, 30.

(9)  Hudson SD, Chumanov G Anal Bioanalytic Chem 2009, 394, 679.

(10) Kneipp K, Kneipp H, Kneipp J Accounts of chemical research 2006, 39, 443.

(11) Xie W, Schlücker S Physical Chemistry Chemical Physics 2013, 15, 5329.

(12) Liu X, Shao Y, Tang Y, Yao K-F Scientific reports 2014, 4.

(13) Wu Z-S, Zhou G-Z, Jiang J-H, Shen G-L, Yu R-Q Talanta 2006, 70, 533.

(14) Colthup NB, Daly LH, Wiberley SE Introduction to infrared and Raman spectroscopy; 3rd ed.; Academic Press: Boston, 1990.

(15) Zhang D, Xie Y, Deb SK, Davison VJ, Ben-Amotz D Analytical chemistry 2005, 77, 3563.

(16) Wang P, Liang O, Zhang W, Schroeder T, Xie YH Advanced Materials 2013, 25, 4918.

(17) Beier HT, Cowan CB, Chou I-H, Pallikal J, Henry JE, Benford ME, Jackson JB, Good TA, Coté GL Plasmonics 2007, 2, 55.

(18) Bhowmik D, Mote KR, MacLaughlin CM, Biswas N, Chandra B, Basu JK, Walker GC, Madhu PK, Maiti S ACS Nano 2015, 9, 9070.

(19) Chou IH, Benford M, Beier HT, Cote GL, Wang M, Jing N, Kameoka J, Good TA Nano Lett 2008, 8, 1729.

(20) Wang P, Xia M, Liang O, Sun K, Cipriano AF, Schroeder T, Liu H, Xie Y-H Anal Chem 2015, 87, 10255.

(21) Walsh DM, Lomakin A, Benedek GB, Condron MM, Teplow DB J Biol Chem 1997, 272, 22364.

(22) Teplow DB Methods in Enzymology 2006, 413, 20.

(23) Younkin SG Annals of Neurology 1995, 37, 287.

(24)Andreasen N, Minthon L, Davidsson P, Vanmechelen E, Vanderstichele H, Winblad B, Blennow K Archives of Neurology 2001, 58, 373.

(25) Hayden EY, Yamin G, Beroukhim S, Chen B, Kibalchenko M, Jiang L, Ho L, Wang J, Pasinetti GM, Teplow DB J Neurochem 2015, 135, 416.

(26) Roychaudhuri R, Yang M, Hoshi MM, Teplow DB J Biol Chem 2009, 284, 4749.

(27) Yu X, Hayden EY, Xia M, Liang O, Cheah L, Teplow DB, Xie YH Protein Sci 2018, 27, 1427.

(28) Hellstrand E, Boland B, Walsh DM, Linse S ACS Chem Neurosci 2010, 1, 13.

(29) Movasaghi Z, Rehman S, Rehman IU Applied Spectroscopy Reviews 2007, 42, 493.

(30) Tu AT, Tu A Raman spectroscopy in biology: principles and applications; John Wiley & Sons: New York, 1982.

(31) Wang P, Zhang W, Liang O, Pantoja M, Katzer J, Schroeder T, Xie YH ACS Nano 2012, 6, 6244.

(32) Ferrari A, Meyer J, Scardaci V, Casiraghi C, Lazzeri M, Mauri F, Piscanec S, Jiang D, Novoselov K, Roth S Physical review letters 2006, 97, 187401.

(33) Pérez-Ruiz E, Decrop D, Ven K, Tripodi L, Leirs K, Rosseels J, van de Wouwer M, Geukens N, De Vos A, Vanmechelen E, Winderickx J, Lammertyn J, Spasic D Analytica Chimica Acta 2018, 1015, 74.

# Chapter 4

# Orientation Dependence

## 4.1 Introduction

The detection and identification of biomarkers is vital for early diagnosis of disease[1-3]. Yet, for many diseases, there exists few useful biomarkers that can be identified, which requires a high sensitive and high specific detecting method. Surface-enhanced Raman Scattering (SERS) is a powerful analytical technique that is routinely used in identifying single molecules with high specificity[4,5]. With its extremely large scattering cross sections ($10^{-17}$-$10^{-16}$ cm$^2$/molecule), SERS is often adopted in biological research[4]. The well-studied Raman amide peaks have been leveraged to infer protein secondary structure[6], which is directly related to protein functionality. SERS detection is free of labels[9,10], a feature that is important for studying proteins. Labels inevitably change the molecular structure[7-8], which is vital to protein functionality. As a result, labels change the molecular structure and prevent the study of conformational state of proteins.

The ability of SERS to distinguish minute differences between protein types has led to the fundamental question of whether the differences in SERS spectral features are caused by difference in molecular type or an array of other parameters including orientation, protein-protein interaction, etc. Intuitively, one could expect that even for a pure sample of a specific type of biological molecules whose structure tends to be much more complex than small inorganic molecules, SERS spectra would vary by a large extent due to factors such as their random orientation on SERS surface. In this sense, it is important to understand whether such variation in spectral features will overshadow the differences between molecular species leading to degraded specificity. This study aims to answer this question focusing on one of the several variables, namely orientation dependence.

The study of orientation of individual amino acid is made possible by the single-molecule sensitivity of our SERS platform[11,12]. By exploiting the extremely large effective cross sections of SERS, it is possible for us to obtain SERS spectra with good signal-to-noise ratio from individual proteins/peptides. Objective verification of single-molecule detection is ensured by bi-analyte approach[13,14] commonly employed in biological research (see Method section for details).

Here we report the orientation dependence of bio-molecules at single molecule sensitivity by comparing the standard deviations of relative peak intensity from individual molecules randomly placed over SERS substrate surface. Au nano pyramid

arrays with well-controlled dimensions were fabricated and used as SERS substrates rendering single molecule sensitivity. The capability of SERS to probe the orientation difference of individual molecules is demonstrated first followed by the observation of the variation of peak intensity from analytes consisting of multiple bio-molecules to show the orientation dependence. It is shown that despite of the variation of protein types, there is a quantitative trend of decreasing standard deviation in peak intensity with increasing molecule size. It is interesting to note that the variation of the intensity of the peaks derived from Amide III remain low for all proteins presumably because the secondary structure of protein is determined by all the bonding which orientation counteract against one another, indicating that secondary structure of protein can be characterized by SERS with high specificity. Simulation results corroborate the experimental observation in that the features of SERS spectrum being completely independent (to within experimental uncertainty) of molecule orientation for large molecules. This last point is the most important outcome of the current study and is of critical importance to explaining the exceedingly high specificity of SERS of proteins reported by researchers worldwide.

To our knowledge, this is the first time orientation dependence of Raman spectroscopy has been systematically studied at single molecule level. Considering the size of molecules exist in human body, the orientation dependence study sets a solid foundation for not only *in vitro* differentiation and secondary structure study

capability using SERS, but potential *in vivo* detection including disease diagnostic as well.

In this chapter, a detailed analysis on the influence of orientation on protein differentiation is presented. In Chap 4.2, we present the method of bi-analyte analysis and standard deviation analysis. In Chap 4.3, we use physics method and simulation analysis to explain the possibility of doing differentiation with orientation difference. Chap 4.4 use experimental data to validate our theory and Chap 4.5 further expand the theory into showing the stability of secondary structure analysis using SERS.

## 4.2 Experiment procedure

### 4.2.1 Bi-analyte analysis

To test the orientation dependence of SERS, we first need to make sure that the analyte we are testing is single molecule concentration. One of the most well accepted methods to validate single molecule is bi-analyte analysis. Bi-analyte analysis pin down unambiguous proof single-molecule sensitivity by using two analyte molecules at low concentration.

The low concentration of each of the 2 analytes suggests that, statistically speaking, there cannot be much more than 1 molecule per colloid. If we are able to observe three different types of spectrum (for each analyte and for the two combined), we know we have not reached single molecule concentration. As we further lower the

concentration, only two types of spectrum appear. Statistically speaking, it is unlikely that each of the colloid contains more than one molecule and none of the collides has both types of the molecule.

In our case, solutions containing two types of protein of similar size and concentration were mixed, deposited on the gold nano-pyramid substrate and let dry. Raman measurement was performed. The bi-analyte analysis is a contrast based spectroscopic technique that monitors the spectra of two types of protein at the same time. As the solution being progressively diluted, the frequency of both proteins exist at the same hotspot decreases monotonically until only one of the two protein spectra being present at each hotspot. In our experiments, this occurred at a concentration of $10^{-9}$ M. The bi-analyte method is by far the only reliable method of making sure that a spectrum is truly from one single molecule.

As a result, we can use a mixture of the two analytes to circumvent many problems associate with the uncertainty of single molecule detection and further understand the orientation dependence of each molecule.

### 4.2.2 Standard Deviation analyisis

Standard deviation is a number used to tell how each sample in the measurement spread out from the mean of the group. The smaller the number is, the more close to the average the samples are. In our case of SERS analysis, if for one molecule, the standard

deviation of all the measurements is low, then we know the spectrum of the sample is relative stable and easy to be distinguished (from other types of molecule).

To calculate the standard deviation, we use the standard formula:

$$\sigma = \sqrt{\frac{1}{N}\sum_{i=1}^{N}(x_i - \mu)^2}$$

In this formula, N means the total sample number; μ stands for the average of the entire sample tested.

We try to understand the change of the standard deviation of SERS spectrum and especially use peak intensity as the main character tested. So we try to calculate the standard deviation of each characteristic peak and evaluate whether the σ is small enough for us to know that the protein sample is distinguishable from the other proteins.

Simulation of standard deviation was done to understand the theoretical limitation of the physics model using the computer language R with the function "runif" for generating random numbers and function "sd" for calculating the standard deviation. Each of the random numbers was generated between 0 and 1 representing the cosine value of the angle between EM field polarizations and bonding direction. 1000 random numbers were created for each group to calculate the standard deviation. Potential error of this simulation may arise from the amino acid orientation preference. E-field vector of plasmon resonance being along the surface norm whereas that of the amino acids

might have a preference of being perpendicular to the surface norm due to protein-surface interaction. Multiple groups of random number were created to simulate the increase number of same bonding. Parameter "geom intensity" was added in "qplot" function of "ggplot2" package to plot the data distribution for each group.

*4.2.3 Experimental settings*

Proteins: Neurotensin (8-13) (6 amino acid) and Substance P (7-11) (5 amino acid) were purchased from GenScript (Piscataway, NJ). Human MANF/ARMET Protein (182 amino acid) and Human CD137/4-1BB Protein (225 amino acid) were purchased from Sino Biological (Wayne, PA). All proteins were of the highest purity available.

Substrate Fabrication: The gold nano-pyramid platform is based on polystyrene (PS) sphere lithography. Readers interested in learning about the details of the substrate fabrication are referred to published literature [5]. PS sphere was used to self-assembly and form monolayer on $SiO_2$/Si wafer surface to create hexagonal patterns. The sample was then dry etched by $O_2$ plasma (200W, 50s) to reduce PS sphere size. The PS spheres are used as the mask for plasma etching to remove the exposed SiO2 film. The substrate is then etched in KOH solution (60%) for 2min to form pyramidal structures on Si surface with patterned SiO2 as etch mask. 200nm Au film was deposited on the mode and epoxy was used to glue the Au film on another wafer. Au nano-pyramids with a 200nm size can be peeled off as the gold nano-pyramid platform.

Raman measurement: 20μL of protein solution was pipetted onto the center of the platform and dried in a fume hood. Raman spectra and mapping of molecules were carried out using a Thermo Fisher DXR 2xi Raman Imaging Microscope under ambient condition. The excitation wavelength is 785nm from a GaAlAs diode laser. The power of the laser was kept at 0.5mW and the spectroscope was accomplished with a 300 lines/mm gating. An objective lens of 50x (Long Working Distance) were used. A step size of 300nm was used for mapping. Raman data were analyzed using Renishaw WiRE 4.2 software.

## 4.3 Theoretical analysis for protein orientation dependence

### *4.3.1 Physical explainations*

**Fig 4.1.** Schematic diagram of the polarization configuration in SERS. Red arrow indicate the electric-field direction, which is perpendicular to the substrate surface. Blue bonds are the same type for both figures. (a) Light incidence where a single bond has $\theta_0$ solid angle with E-field direction. (b) Light incidence where multiple bonds interact with E-field and form multiple angles.

We show the polarization configuration in Raman spectroscopy. The Raman peak intensity is proportional to the incident light intensity as well as the cosine of the angle between the polarization direction of the EM field of local plasmon resonance and the orientation of the Raman active bond. As normalization was done to rule out the impact of laser intensity variation and plasmon resonance intensity, the relation shown in Fig 4.1 can be obtained for the Raman peak intensity of individual Raman active bonds. When multiple bonds of the same type exist within one hotspot, multiple angles contribute additively to the observed peak intensity. As the number of bonds increases, the values of cosine function of different angles counteracts each other, leading to a convergence in standard deviation.

## 4.3.2   Simulation results



**Fig 4.2.** Standard deviation (SD) value of the peak intensity for Tyr ($850cm^{-1}$ and $1360\ cm^{-1}$), C-N ($1140\ cm^{-1}$), $CH_3$ ($1386\ cm^{-1}$) and Phe ($1000\ cm^{-1}$) peaks for the four types of protein. The SD values of different peak intensity for the two peptides (Substance 5 amino acid, Neurotensin 6 amino acid) and two proteins (MANF 182 amino acid, CD137 255 amino acid). The SD value decrease from substance to MANF, and almost remain steady between MANF and CD137.

We performed simulation (see materials and methods) to validate the physics model (Fig 4.2.). Amino acid number from 1 to 13 were simulated using 13 groups of random cosine number (each group consists of 1000 random number from 0-1) to simulate the random orientation of Raman active bonds on SERS substrate relative to the incident light direction. The simulation was repeated five times to ensure reproducibility. We then calculated and plotted the standard deviation to compare with experimental results. Relative standard deviation so obtained ranges from 0.32 to 0.07 for random angles

with a sharp decrease as amino acid number increases from 1 to 7 and the decrease becomes slower and eventually stops at 12 amino acids. From the three distributions at amino acid number 1, 6 and 13, we attribute the decrease in standard deviation to the trend of uniformity in summation of the values of cosine function with random angles.

4.4   Orientation dependence of protein differentiation

Our initial experiment sought to establish the spectral characteristics of multiple types of protein at single molecule concentration. The comparison between orientation-dependent protein Raman signatures is valid only after single molecule detection is established. Bi-analyte analysis, being the only reliable method to ensure that a spectrum being indeed from one single molecule, is utilized for both peptides (Neurotensin (8-13) and Substance P (7-11)) and proteins (Human MANF/ARMET Protein and Human CD137/4-1BB Protein). At around $10^{-9}$ M concentration, we observed spectral responses from over 85% of SERS hotspots to be of only one type of bio-molecule, firmly proving the single molecule detection capability. 100 Raman spectra were collected from individual molecules of each sample. All spectra were normalized by the highest peak to rule out the impact of incident light intensity. Figure 1 and 2 show the average Normalized Raman intensities of peptides and proteins from bi-analyte analyses, respectively. Common predominant peaks in proteins occurred at 850 (Tyr), 1000 (Phe), 1140 (C-N bonding), 1360 (Tyr) and 1386 ($CH_3$) cm$^{-1}$. The orientation difference does not affect the position of Raman peaks as evidenced by our

experimental results shown in fig.1 and 2. The only factor leading to the normalized

spectral difference between individual molecules of the same protein is orientation.



**Fig 4.3.** SERS of substance P (green) , neurotensin (red), and a mixture of the two peptides (blue) on hybrid platform. Feature wavenumbers are listed above the peaks for individual amino acids. Bi-analyte analyses are done to validate that individual spectra are indeed from single molecules.

**Fig 4.4.** SERS of CD137 (red), MANF (blue), and a mixture of the two proteins (green) on a hybrid platform. Feature wavenumbers are listed above the peaks for individual amino acids. Bi-analyte analyses are done to validate that individual spectra are indeed from single molecules.

Standard deviation was used to show the intensity variation of individual protein molecules. We define the change in intensity using the absolute change in intensity divided by the normalized average peak intensity to make all peaks comparable. For each type of protein we tested, standard deviation of peak intensity was calculated for multiple wavenumbers, including Tyr ($850 cm^{-1}$ and $1360\ cm^{-1}$), Phe ($1000\ cm^{-1}$), C-N bonding ($1140\ cm^{-1}$) and CH3 ($1386\ cm^{-1}$) were calculated (Fig 4). The X- axis of the figure represents the atomic weight of the substances ranging from Substance P (5 amino acid (aa)), Neurotensin (6aa) to Human MANF (182aa) and Human CD137 (255aa). In Substance P peptide, Tyrosine only appears once in the peptide structure

and the standard deviation for the Tyrosine characteristic peak is 0.25. The two types of protein (Human MANF and Human CD137) have multiple Tyrosine in their protein structure, showing a low standard deviation of 0.06-0.07. Similar situation happens to the amino acid Phenylalanine. The standard deviation for Phe in Neurotensin peptide (2 Phe in its structure) is 0.18 while the standard deviation is 0.06-0.07 for the two proteins. For the two peptides (Substance P and Neurotensin), the reason for their difference in amino acid characteristic peak is that Substance P contains only 1 Tyr while Neurotensin contains 2 Phe, leading to the two amino acid orientation counteract against one another. MANF and CD137 shows almost same deviation difference in peak corresponding to the 2 types of amino acid (Tyr and Phe). The two types of bonds ($CH_3$ and C-N) appears in every amino acid, and this leads to smaller standard deviation than the two amino acids, which appears less frequently in the proteins. A clear trend of convergence in standard deviation with increasing protein size is observed. It can be concluded that orientation has less impact on spectral signature as the size of the molecules becomes larger.

**Fig 4.5.** Standard deviation (SD) value of the peak intensity for Tyr ($850 cm^{-1}$ and $1360 cm^{-1}$), C-N ($1140 cm^{-1}$), $CH_3$ ($1386 cm^{-1}$) and Phe ($1000 cm^{-1}$) peaks for the four types of protein. The SD values of different peak intensity for the two peptides (Substance 5 amino acid, Neurotensin 6 amino acid) and two proteins (MANF 182 amino acid, CD137 255 amino acid). The SD value decrease from substance to MANF, and almost remain steady between MANF and CD137.
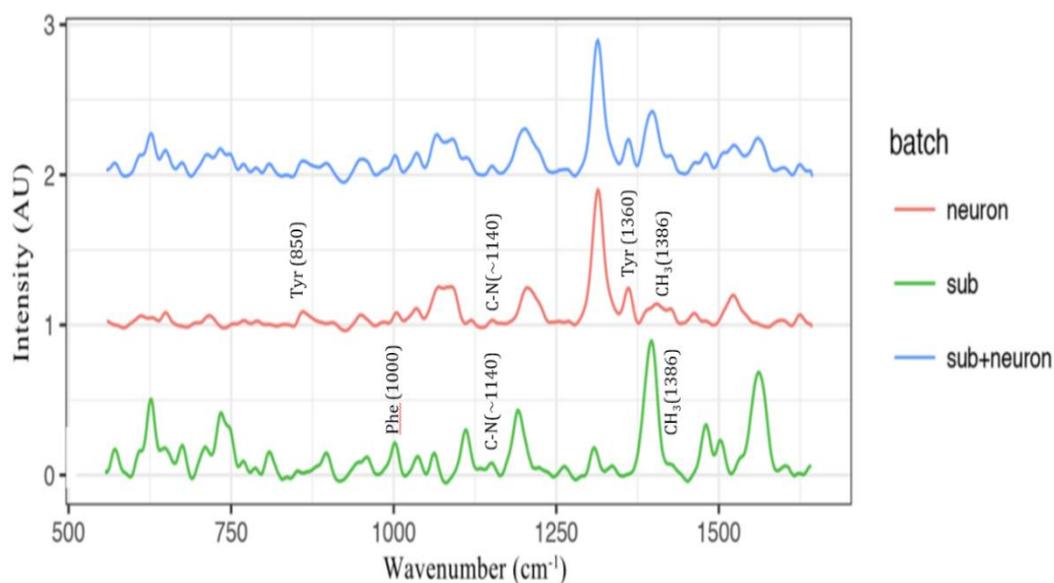
## 4.5  Orientation dependence of secondary structure

Amide III peaks containing secondary structure information of proteins such as protein-protein interaction and conformational changes are of vital importance for the study of protein functionality. By zooming in on the spectral region of 1180-1320 $cm^{-1}$ of the averaged Raman spectroscopy of all four proteins at single molecule concentration, Amide III peak is clearly visible as shown in Fig.6a. Peak corresponding to β-sheet ($1230 cm^{-1}$), random ($1250 cm^{-1}$) and α-helix ($1320 cm^{-1}$) were all marked in the figure.

**Fig 4.6.** SERS spectrum of Amide III (1180-1320cm$^{-1}$) for substance, neurotensin, MANF and CD137. Key peak related to secondary structure of protein: α-helix, β-sheet and random structure are shown in the spectra.

When evaluating the orientation impact on secondary structure detection of protein using SERS, the standard deviation was calculated for the three important Amide III peaks as is shown in Fig. 6b. The highest value (0.062) of standard deviation appears for the substance β-sheet (1230 cm$^{-1}$) peak with the comparable value of the α-helix (1320 cm$^{-1}$) peak and (0.061) of Neurotensin. All other peaks are of the standard deviation value lower than 0.06, which is below experimental uncertainty, and the value does not show amino acid number (molecular weight) dependence. Compared with the simulation results in Fig. 5, all the standard deviation values of Amide III are below the convergence value calculated by simulation. This shows that the Amide III

peaks are stable for all protein tested. This observation indicates that SERS is a

promising method in understanding protein secondary structure.



**Fig 4.7.** the SD for the three peaks respectively for the two peptides and two proteins.

## 4.6 Conclusion

We have shown through experimental study and simulation that for biomedically

relevant protein, the distribution of the various angles of amino acids over SERS

surfaces is not expected to degrade the specificity of protein identification. The four

types of protein all show a high degree of uniformity in protein signature.

Hundreds of SERS spectra from randomly oriented single molecules are

systematically

studied with emphasis on the impact of molecular orientation on SERS spectral

features. The single molecule nature of each of the SERS spectrum is ensured by

utilizing bi-analyte analysis and made possible by the single-molecule sensitivity of

the SERS substrate employed. The molecules are dispersed randomly over the SERS substrate surface with their orientation being distributed stochastically. The standard deviation in peak intensity of common peaks of amino acids is shown to be dependent monotonically on the size (and thus the number of amino acids) of individual protein molecules with orders of magnitude smaller standard deviation from proteins containing more than 13 amino acids. Numerical simulation has shown that orientation has little impact on SERS peak intensity variation of protein with more than 13 amino acids and have no impact on secondary structure detection regardless of protein size, consistent with the experimental observations. This study offers a plausible explanation of the reason that the random orientation of biomedically relevant proteins does not hinder the repeatedly demonstrated high level of specificity of SERS in identification of biological entities including proteins, exosomes and cells. This work provides a solid foundation for SERS to become an important tool for bio-medical applications.

## 4.7 References

1    Sebastiaan Engelborghs, Karen De Vreese, Tom Van de Casteele, Hugo Vanderstichele, Bart Van Everbroeck, Patrick Cras, Jean-Jacques Martin, Eugeen Vanmechelen, and Peter Paul %J Neurobiology of aging De Deyn,    29 (8), 1143 (2008).

2    Guy M McKhann, David S Knopman, Howard Chertkow, Bradley T Hyman, Clifford R Jack Jr, Claudia H Kawas, William E Klunk, Walter J Koroshetz, Jennifer J Manly, Richard %J Alzheimer's Mayeux, and dementia,    7 (3), 263 (2011).

3    P Vemuri, HJ Wiste, SD Weigand, LM Shaw, JQ Trojanowski, MW Weiner, DS Knopman, RC Petersen, and CR %J Neurology Jack,    73 (4), 287 (2009).

4    Katrin Kneipp, Yang Wang, Harald Kneipp, Lev T Perelman, Irving Itzkan, Ramachandra R Dasari, and Michael S %J Physical review letters Feld,    78 (9), 1667 (1997).

5    Shuming Nie and Steven R %J science Emory,    275 (5303), 1102 (1997).

6    Sebastian %J ChemPhysChem Schlücker,    10 (9-10), 1344 (2009).

7    JL Lippert, D Tyminski, and PJ %J Journal of the American Chemical Society Desmeules, 98 (22), 7075 (1976).

8    Nakul C Maiti, Mihaela M Apetri, Michael G Zagorski, Paul R Carey, and Vernon E %J Journal of the American Chemical Society Anderson,    126 (8), 2399 (2004).

9    Lucas A Lane, Ximei Qian, and Shuming %J Chemical reviews Nie,    115 (19), 10489 (2015).

10   Wei Xie, Bernd Walkenfort, and Sebastian %J Journal of the American Chemical Society Schlücker,   135 (5), 1657 (2012).

11   Pu Wang, UCLA, 2015.

12   Xinke Yu, Eric Y Hayden, Ming Xia, Owen Liang, Lisa Cheah, David B Teplow, and Ya-Hong %J Protein Science Xie,    (2018).

13   Jon A Dieringer, Robert B Lettan, Karl A Scheidt, and Richard P %J Journal of the American Chemical Society Van Duyne,    129 (51), 16249 (2007).

14   Eric C Le Ru, Matthias Meyer, and Pablo G %J The journal of physical chemistry B Etchegoin,    110 (4), 1944 (2006).

# Chapter 5

# Cerebrospinal fluid (CSF) diagnostic using SERS

## 5.1   Introduction

### 5.1.1   CSF and Alzheimer's disease

Alzheimer's disease (AD) has affected several million people all over the world. The disease affects older adults and is the most common cause of dementia[1-4]. No cure or disease-modifying therapy exists[5] currently and the disease inevitably progresses in all patients[6-8]. Besides, there is no currently a well-established single biomarker test to diagnose AD. Current diagnoses rely on medical history, cognitive testing, and a variety of biomarkers including brain imaging, proteins in cerebrospinal fluid (CSF),

proteins in blood, and genetic profiling[7,9,10]. As a result, diagnosing AD patients is a lengthy and costly process, which impediment patient care and increase healthcare cost. The development of a biomarker that allows detection of AD during the pre-symptomatic phase is critical to the discovery and development of effective AD diagnoses and treatments[11-13].

Cerebrospinal fluid is a liquid that surround brain and spinal cord and is frequently used for neurologic disorder disease diagnostic. As the fluid travels, it picks up supplies from the blood and gets rid of wastes from brain cells. Diseased CSFs can carry different contents compared to normal ones and thus can be used as a biomarker for multiple diseases.

Body fluid such as CSF, plasma and urine are considered as potential source for biomarkers for AD diagnostic. CSF biomarkers are of the largest interest due to their low cost, minimal invasiveness, and high diagnostic accuracy. Compared to blood, CSF is in direct contact with the extracellular space of the brain due to the existence of blood brain barrier and can reflect the biochemical changes in side the brain. Some biomarkers such as Aβ42 and Tau protein exists in CSF[12,14-16], and other unknown markers or a combination of multiple proteins might work together as a disease prediction signal. Some of the state to art biomarker progress are listed in the table 5.1. The limited number of published studies on using multi-model approaches yield only limited success. These findings highlight the importance of the highly heterogeneous

nature of AD pertaining to biomarker discovery. This leads to the proposed focus on

multi-modal biomarker discovery approach.

**Table 5.1.** Biomarkers used for disease diagnostic in CSF.

| Bio-marker | Advantages | Limitations | Ref |
|---|---|---|---|
| Aβ42, t-tau, p-tau/ <br><br> Aβ42, t-tau/ <br><br> Aβ42 | 1. Can correlate AD directly, <br><br> 2. Highly sensitive and specific, <br><br> 3. Can detect AD progression. | 1. Invasive, sample has to be collected by lumbar puncture, <br><br> 2. Accuracy of diagnostic is not ideal. | |

## 5.1.2 Machine learning and diagnostic

Machine learning, together with artificial intelligence, has been providing inexpensive

and available means to improve the healthcare condition worldwide. In recent years,

advanced computational methods have been employed to meet the needs of sensitive

and fast diagnostic[17,18].

Some mature technologies have been applied in disciplines such as oncology,

pathology and rare diseases. Stanford University researchers have trained an

algorithm to diagnose skin cancer using deep learning (CNN), by training the

algorithm with 130,000 images of skin lesions representing over 2,000 diseases. New

platforms have appeared to diagnose rare disease based on facial dysmorphic features. It is currently available only to trained clinicians to prevent false positive, instead of having own diagnostic capability, but is still a huge step forward on the diagnostic frontier[19-21].

Here in this work, we try to apply advanced machine learning tools in the diagnostic of Alzheimer's disease. Detailed methods used in the study will be detailed described in Experimental procedures section.

## 5.1.3 Current diagnostic methods for AD

Mini-mental state examination (MMSE) test is used extensively in clinical and research settings to measure cognitive impairment. Any score greater than or equal to 24 points (out of 30) indicates a normal cognition. Clinical dementia rating (CDR) is a numeric scale used to quantify the severity of symptoms of dementia. The composite rating of the score is shown in Table 5.2.

**Table 5.2.** Interpretation of CDR score.

| Composite Rating | Symptoms |
|---|---|
| 0 | None |
| 0.5 | Very mild |
| 1 | Mild |

| | |
|---|---|
| **2** | Moderate |
| **3** | Severe |

The CDR test has shown a very high reliability, and appears to be a reliable and valid tool for assessing and staging dementia. However, the test takes a long time and it is not possible to capture changes over time.

## 5.2 Experimental procedures

### 5.2.1 Patient profile

Thirty CSF samples were obtained from University of California, Irvine, Institute for Memory Impairment and Neurological Disorders, Alzheimer's Disease Research Center (UCI MIND, ADRC), and John Ringman, MD, University of Southern California. Several characteristics of the patients CSF have already been measured, using standard procedures among ADRC's. These include the levels of Aβ42, total-tau, and phospho-tau, as well as the Mini-mental state exam (MMSE) and the clinical dementia Rating (CDR). A statistical summary of the patient data is shown in Table 5.3. This is the patient sample we are able to acquire for now.

**Table 5.3.** Patient information summary

| | Healthy | Dementia | FAD+ | FAD - |
|---|---|---|---|---|
| # of cases | 10 | 9 | 5 | 4 |

| | | | | |
|---|---|---|---|---|
| Male/ Female | 3/7 | 4/5 | 3/2 | 3/1 |
| Age (years) | 76.6 (+/- 5.5) | 79 (+/- 4.9) | 36 (+/- 12.9) | 34 (+/- 14.8) |
| Adjusted age | NA | NA | -10 (+/- 10.6) | NA |
| CSF Aβ42 (pg/mL) | 645.6 (+/- 353.0) | 375.9 (+/- 305.8) | 186.2 (+/- 60.4) | 418.8 (+/- 174.9) |
| CSF Total-Tau (pg/mL) | 364.9 (+/- 265.3) | 570.6 (+/- 529.4) | 516.9 (+/- 363.3) | 312.1 (+/- 266.8) |
| CSF phospho-Tau (pg/mL) | 83.8 (+/- 43.7) | 87.2 (+/- 43.6) | 99.2 (+/- 50.8) | 73.7 (+/- 39.8) |
| MMSE (0-30) | 29.9 (+/- 0.3) | **19.6** (+/- 3.6) | **25** (+/- 7.9) | 28.8 (+/- 0.5) |
| CDR-Sum of boxes (0-18) | 0.1 (+/- 0.2) | 9.1 (+/- 2.0) | 1.6 (+/- 3.0) | 0.25 (+/- 0.5) |
| CDR-global (0-3) | 0.1 (+/- 0.2) | 1.44 (+/- 0.5) | 0.2 (+/-0.45) | 0.13 (+/- 0.25) |

## 5.2.2 Desalting process and ziptip

The crystallization of salt in CSF makes it relatively difficult to get uniform SERS spectra. As is shown in Fig 5.1, majority of the surface is covered with crystallization and only ~40% of the places can obtain SERS signal (red spectra in Fig 5.2). Desalting can be a relatively easy and stable method to increase the yield of SERS results.

**Fig 5.1.** Optical microscope image of CSF under 50 x magnifications.



**Fig 5.2.** Spectral mapping of the grid area in Fig 5.1.

The whole process is completed in the following steps: Aspirate 10 μL wetting solution into tip and dispense to waste. Repeat. Apsirate equilibration solution into tip and dispense to waste. Repeat. Bind peptides to ZipTip pipette tip by aspirating and dispensing 3-7 cycles (simple mixtures), up to 10 cycles (complex). Aspirate washing solution and dispense to waste. Repeat wash once. A 5% methanol in 0.1% TFA/water wash can improve desalting efficiency. Dispense 1-4 μL of elution solution into clean 0.5 mL Eppendorf microcentrifuge tube using a standard pipette tip (Note: if μ-C-18,

dispense 0.5-2 µL of elution solution). Aspirate and dispense eluant through ZipTip at least 3 times without introducing air. Sample recovery can be improved by increasing elution volume to 10 µL (but at expense of concentration).

To further make sure that the desalting process does not change the spectral features of SERS spectrum, a comparison between the original spectrum (CSF diluted 100 times) and the ziptip result is done. Hierarchical clustering algorithm (HCA, see next section for more detail) is done to do the validation. 32 different CSF sample is tested using SERS and an averaged spectrum is calculated for each of them. We further did ziptip desalting for one of them (unlabeled, double blind) and did SERS with same condition and also calculated the average spectrum (sample 34). A comparison group is made for the same unknown CSF sample without the dilution (sample 33) HCA analysis is done for the 34 spectrum. We take the intensity at each wavenumber as one dimension and did HCA with the high dimension data. The result is shown in Fig 5.3, indicating that the ziptip sample is the same as sample 18, and we can further know that the spectral change after ziptip is even smaller than dilution. As a result, we are able to safely use ziptip for fast and accurate CSF analysis.

**Fig 5.3.** HCA analysis of the ziptip result (sample 32 and 34) with 30 different CSF samples.

## 5.2.3. Machine learning analysis

Hierarchical clustering algorithm (HCA)

HCA is a classification method which builds a hierarchy of clusters by merging or splitting clusters in a greedy manner. In order to decide which clusters should be combined (for agglomerative method) or how the cluster should be split (for divisive

method), a distance matrix is formed to define the similarities between nodes and clusters.

The least dissimilar pair is defined by the distance matrix:

$$d[(r),(s)] = min\ d[(i),(j)]$$

where the minimum is over all pairs of clusters in the current clustering.

Clusters with the highest similarities are merged and form a single cluster and the distance matrix can be further updated. The proximity between the new cluster, denoted (r,s) and old cluster (k) is then defined as:

$$d[(k),(r,s)] = min\ d[(k),(r)], d[(k),(s)]$$

The analysis in this work is done using R. The function "hclust" is used to do the HCA analysis and the distance matrix is calculated by the function "dist". Linkage method is changed by adding "method" in the "hclust" function. The result is plotted in a dendrogram format.

*5.2.4 CNN analysis*

Considering the collected SERS spectrum has only one dimension which covers the entire spectrum of interest, we employed a one-dimensional CNN to process and classify the SERS spectral data. The convolutional layers of our model use ReLU nonlinear activation function, the convolutional layers are connected by max-pooling

layers which down-sample the feature maps, the output of last max-pooling layer is fed

to two consecutive fully-connected layer to give the final classification result. During

the training process, we used the scalar sum of weighted losses to train the CNN model.

To increase the training set, we performed the following data augmentation methods on

the training data: (i) Random shifting of each spectrum by a few (1~2) wavenumbers.

(ii) Introduction of a random noise onto each spectrum. (iii) We also randomly

produced linear combinations of spectra collected from the same mapping procedure.

The Adagrad algorithm was used to train the model, and early stopping was applied to

prevent overfitting.

## 5.3    Reproducibility analysis

Being the first time to use the hybrid SERS platform in human fluid, we need to first

validate that we have the ability to differentiate each individual. To prove that the

platform can work as a clinically viable assay of cerebral spinal fluid (CSF) for

diagnosis of neuro-degenerative diseases, 3 replicates of CSF sample from 5

individuals are prepared for classification. The study was designed to be double-blind

in which the identities of the samples are kept anonymous to the personnel conducting

the measurements and data analysis.

Each of the fifteen tubes of CSF samples was diluted by a factor of 100 and applied to

the hybrid SERS platform. SERS spectra were acquired in the wavenumber range of

550-1650 cm$^{-1}$. SERS mappings with a step size of 3μm (i.e., independent areas of 9 μm$^2$ each) are done with over 80 spectra being obtained for each sample. The 80 spectra are then averaged to one representative spectrum. The averaging process by nature allows the biological variability of the patient CSFs be represented in a linearly proportional fashion, though the actual spectral spread of each sample has been lost.

Direct observations of figure 1 show clearly that multiple samples share similar spectral features. The following samples share similar features in terms of the most intensive peaks: sample 1 and 4 (1405 cm$^{-1}$, 1439 cm$^{-1}$), sample 2, 6, and 12 (~750 cm$^{-1}$, 913cm$^{-1}$, 1239 cm$^{-1}$, 1578 cm$^{-1}$), sample 5 and 7 (~1100cm$^{-1}$, 1295cm$^{-1}$, 1420 cm$^{-1}$), and sample 8 and 9 (816 cm$^{-1}$ and 1297 cm$^{-1}$).

It is difficult to quantitatively determine the degree of similarity and dissimilarity between the 15 averaged spectra using naked eyes. To analyze the grouping of the samples, we use hierarchical clustering algorithm (HCA) to more scientifically group the replicates. Integration normalization is done when comparing samples. To do the HCA, we use intensity of each wavenumber of the averaged spectra as a dimension and try to group the samples using single linkage of aggregation hierarchical clustering.

**Fig 5.4.** SERS spectra of the 15 unknown CSF samples.

Figure 5.4 shows that we are able to find five clear groups of samples by the peak intensity of each wavenumber. The degree of similarities between the various samples is indicated by the proximity to zero of the lines connecting them. The closer the connecting line to 0, the more similarities the two samples share.

**Fig 5.5.** HCA analysis for the 15 unknown CSF samples with 100% accurate clustering result.

Comparison between the sample grouping of the double-blind study shown in Fig 5.5 and the clinical diagnosis of the patients from whom the CSF samples were obtained shows that the accuracy of the grouping performed under double-blind condition is **100%**. This outcome shows convincingly that SIM has the promise of becoming a clinically viable assay for diagnosis of neuro-degenerative diseases, for which the only means to date is psychoanalysis (with debatable accuracy).

The results of this double-blind study represents a giant step forward in establishing the power of the platform in analyzing patient samples despite of biological variability.

To further examine the uniformity of each sample, we first performed training and evaluation in a leave-one-out approach. To be specific, each time, one spectrum was

picked out and used for testing while the rest data was used for training. The procedure was repeated until every spectrum was left out once, then the average accuracy across all the data was computed. The result shows an accuracy of 97.7% for spectrum of normal sample and 93.3% for spectrum of abnormal sample, which suggests there is a good uniformity over the SERS spectral data of normal and abnormal sample, and those two types of sample are differentiable.

## 5.4   CSF as a diagnostic media

We also performed leave-one-group-out evaluation on the dataset because leave-one-out evaluation cannot determine whether our model is capable of exploring the relationship between the SERS testing results on a subject's CSF sample and the subject's diagnosed syndrome (normal or abnormal), because the spectral data collected from the same CSF sample is likely to be dependent on that individual sample, thus forming a group of dependent data, so for a unseen sample that is to be classified, the correlation information of its spectral data should not be given to the model during training. Therefore, for each round of evaluation, we need to make sure that all the spectral data in the test set comes from groups that are not represented at all in the corresponding train set. Since we have a total of 17 CSF samples, leave-one-group-out evaluation is a suitable approach for us to know whether our model generalizes well on the unseen samples. In detail, during each evaluation, the entire SERS spectral data collected from one CSF sample was used for testing while the rest data was used for

training, this kind of evaluation was repeated until every group of spectral data was left out once.

We need to point out here that the data set is much smaller than usual machine learning training set and this experiment is for preliminary test only. Next step of the experiment is described in future work part. The data set will be enlarged and more accurate diagnosis result will be used (post-mortem).

The final classification result also takes the group dependency into consideration. In each round of evaluation, after the trained CNN made predictions on every testing spectrum, all the predictions were then combined through a majority vote to produce the final prediction, i.e., the class with a higher percentage of predictions was considered to be the predicted class of the testing sample. We used the percentage of predictions leading to the predicted class as a score to represent the likelihood that the testing sample belongs to the predicted class.

We tested on all 17 CSF samples and an overall 94% diagnostic rate has been achieved. Among which, normal sample has an accuracy of 8/8 (100%) and an average score of 89.2, diseased sample has an accuracy of 8/9 (88.9%) and an average score of 72.0, as is shown in Table 5.4.

**Table 5.4.** Test score for 17 non-FAD samples.

| Sample | Label | Score |
|--------|---------|-------|
| A | NORMAL | 91.85 |
| F | NORMAL | 85.94 |
| G | NORMAL | 84.09 |
| H | NORMAL | 92.86 |
| I | NORMAL | 80.22 |
| M | NORMAL | 90.91 |
| X | NORMAL | 87.5 |
| AA | NORMAL | 100 |
| B | DEMENTIA | 54.55 |
| N | DEMENTIA | 88.64 |
| O | DEMENTIA | 86.11 |
| R | DEMENTIA | 31.25 |
| S | DEMENTIA | 100 |
| U | DEMENTIA | 89.80 |
| W | DEMENTIA | 57.81 |
| Y | DEMENTIA | 65.08 |
| AB | DEMENTIA | 75 |

Within this table, we are able to see that sample R is the only one with a prediction score less than 50, indicating we are not able to tell R is a diseased sample. To further understand this sample, we refer to her cognitive test score. Her MMSE score is 25, which should be diagnosed as normal. However, her CDR test score is far away from normal (9 for CDRSUM and 1 for CDRGLOB). These scores present mixed

information for diagnostic and inevitably influence our test result.

To further understand more complicated situation, such as FAD related patients, we do the same test based on our training of normal and diseased patients. FAD negative (normal sample) has an accuracy of 4/4 (100%) and an average score of 91.25, FAD positive sample has an accuracy of 4/5 (80%) and an average score of 80.0. The result is shown in Table 5.5.

**Table 5.5.** Test score for FAD related samples

| Sample | Label | Score |
|--------|-------|-------|
| **D** | FAD(+), -19 | 49 |
| **E** | FAD(+), -5 | 100 |
| **K** | FAD(-), -11 | 85 |
| **L** | FAD(-), -17 | 100 |
| **P** | FAD(-), 0 | 85 |
| **Q** | FAD(+), 4 | 53 |
| **Z** | FAD(+), -8 | 100 |
| **AC** | FAD(+), -22 | 98 |
| **AD** | FAD(-), -18 | 95 |

Sample D has a test score of 49, indicating we are not able to accurately define whether she is diseased or not. When referring to her medical scores, we find that she has a MMSE score of 28 and CDRSUM of 0.5, indicating she has mild symptom of dementia. More analysis needs to be done in order to see the relationship between our

test score and the symptom or severity of the patients.

Again, we need to point out here that the sample number of individuals with/without FAD is smaller than the usual data set used for machine learning and will be expanded once we get larger sample size. This result is a preliminary data to show the feasibility of the platform.

## 5.5 Correlation analysis

To further understand the feasibility of our diagnostic result, we have done correlation analysis for the diagnostic index with all Alzheimer's disease related medical parameters. Several parameters are included in the analysis: sex, age, several biomarker level (Aβ42, t-Tau, p-Tau) and several cognitive test score (MMSE, CDRSUM and CDRGLOB).

To better analyze the correlation, we first need to deal with the missing data. In the clinical information provided, t-Tau information of sample N, R,R S, T, W and X are missing, and p-Tau information of sample R is missing. According to literature, the t-Tau and p-Tau level in CSF is highly correlated and from the data provided, we see a highly linear relationship between the concentration of the two types of Tau protein with a correlation factor of r=0.8883. The prediction model can be easily shown by regression:

$$P = 0.1134T + 35.28$$

Among which, P is the concentration for p-Tau protein and T is the concentration for t-Tau protein. We are able to fill in the missing t-Tau level with this correlation.

For the missing data in sample R, a detailed correlation analysis is done, showing that the highest correlated parameter with the Tau protein level is MMSE, and the factor is only r=0.47. As a result, we will eliminate sample R when doing correlation analysis. The correlation coefficients between all the biomedical parameters are shown in Fig 5.6.

**Fig 5.6.** Correlation analysis between prediction score and bio-medical matrix.

We can see from the figure that our prediction index is highly correlated to all the 3 cognitive test scores. The correlation coefficients are r = 0.79 with MMSE, r = -0.92 with CDRSUM and r = -0.88 with CDRGLOB. Considering CDR scores are taken as an accurate method of diagnosing AD, our prediction index is accurate in the diagnostic. We have already discussed that one of the limitations of CDR score is that it takes a long time to collect and our method can be a good compensation for the

cognitive test.

Besides, we can observe from the figure that the correlation between biomarkers (Aβ protein, Tau protein) and cognitive test score (MMSE and CDR) is relatively low, the highest correlation coefficient is r = 0.47 (between t-Tau and MMSE) further show that single biomarker is not accurate in AD diagnostic.

## 5.5 Conclusion

We have presented here a novel method to diagnose Alzheimer's disease with high accuracy. By using a combination of SERS platform and machine learning analysis, preliminary result of 100% reproducibility with double blind experiment and 92% accuracy in disease diagnostic is acquired. Another analysis with the exact same procedure will be done once more samples are get, which can make the result more accurate. The correlation analysis further proves that our diagnostic system is more accurate than single bio-marker analysis and can be further applied to clinical usage.

# 5.6 References

1.    Katzman, R., The prevalence and malignancy of Alzheimer disease: a major killer. Archives of neurology 1976, 33 (4), 217-218.

2.    Farrer, L. A.;   Cupples, L. A.;   Haines, J. L.;   Hyman, B.;   Kukull, W. A.;   Mayeux, R.;   Myers, R. H.;   Pericak-Vance, M. A.;   Risch, N.; Van Duijn, C. M., Effects of age, sex, and ethnicity on the association between apolipoprotein E genotype and Alzheimer disease: a meta-analysis. Jama 1997, 278 (16), 1349-1356.

3.    Selkoe, D. J., Alzheimer's disease: genes, proteins, and therapy. Physiological reviews 2001, 81 (2), 741-766.

4.    Selkoe, D. J., Alzheimer's disease is a synaptic failure. Science 2002, 298 (5594), 789-791.

5.    Alzheimer's, A., 2015 Alzheimer's disease facts and figures. Alzheimer's & dementia: the journal of the Alzheimer's Association 2015, 11 (3), 332.

6.    Schelterns, P.; Feldman, H., Treatment of Alzheimer's disease; current status and new perspectives. The Lancet Neurology 2003, 2 (9), 539-547.

7.    Desai, A. K.; Grossberg, G. T., Diagnosis and treatment of Alzheimer's disease. Neurology 2005, 64 (12 suppl 3), S34-S39.

8.    Lopez, O. L., The growing burden of Alzheimer's disease. The American journal of managed care 2011, 17, S339-45.

9.    Hardy, J. A.; Higgins, G. A., Alzheimer's disease: the amyloid cascade hypothesis. Science 1992, 256 (5054), 184-186.

10. McKhann, G. M.;   Knopman, D. S.;   Chertkow, H.;   Hyman, B. T.;   Jack Jr, C. R.;   Kawas, C. H.;   Klunk, W. E.;   Koroshetz, W. J.;   Manly, J. J.; Mayeux, R., The diagnosis of dementia due to Alzheimer's disease: Recommendations from the National Institute on Aging-Alzheimer's Association workgroups on diagnostic guidelines for Alzheimer's disease. Alzheimer's & dementia 2011, 7 (3), 263-269.

11. Bateman, R. J.;   Xiong, C.;   Benzinger, T. L.;   Fagan, A. M.;   Goate, A.;   Fox, N. C.;   Marcus, D. S.;   Cairns, N. J.;   Xie, X.; Blazey, T. M., Clinical and biomarker changes in dominantly inherited Alzheimer's disease. New England Journal of Medicine 2012, 367 (9), 795-804.

12. Olsson, B.;   Lautner, R.;   Andreasson, U.;   Öhrfelt, A.;   Portelius, E.;   Bjerke, M.;   Hölttä, M.;   Rosén, C.;   Olsson, C.; Strobel, G., CSF and blood biomarkers for the diagnosis of Alzheimer's disease: a systematic review and meta-analysis. The Lancet Neurology 2016, 15 (7), 673-684.

13. Yu, X.; Hayden, E. Y.; Xia, M.; Liang, O.; Cheah, L.; Teplow, D. B.; Xie, Y. H., Surface enhanced Raman spectroscopy distinguishes amyloid B-protein isoforms and conformational states. Protein Science 2018, 27 (8), 1427-1438.

14. Blennow, K.; Hampel, H., CSF markers for incipient Alzheimer's disease. The Lancet Neurology 2003, 2 (10), 605-613.

15. Shaw, L. M.; Vanderstichele, H.; Knapik-Czajka, M.; Clark, C. M.; Aisen, P. S.; Petersen, R. C.; Blennow, K.; Soares, H.; Simon, A.; Lewczuk, P., Cerebrospinal fluid biomarker signature in Alzheimer's disease neuroimaging initiative subjects. Annals of neurology 2009, 65 (4), 403-413.

16. Blennow, K.; Hampel, H.; Weiner, M.; Zetterberg, H., Cerebrospinal fluid and plasma biomarkers in Alzheimer disease. Nature Reviews Neurology 2010, 6 (3), 131.

17. Kneipp, K.; Wang, Y.; Kneipp, H.; Perelman, L. T.; Itzkan, I.; Dasari, R. R.; Feld, M. S., Single molecule detection using surface-enhanced Raman scattering (SERS). Physical review letters 1997, 78 (9), 1667.

18. Stiles, P. L.; Dieringer, J. A.; Shah, N. C.; Van Duyne, R. P., Surface-enhanced Raman spectroscopy. Annu. Rev. Anal. Chem. 2008, 1, 601-626.

19. Hosseini-Asl, E.; Keynton, R.; El-Baz, A. In Alzheimer's disease diagnostics by adaptation of 3D convolutional network, 2016 IEEE International Conference on Image Processing (ICIP), IEEE: 2016; pp 126-130.

20. Sarraf, S.; Tofighi, G. In Deep learning-based pipeline to recognize Alzheimer's disease using fMRI data, 2016 Future Technologies Conference (FTC), IEEE: 2016; pp 816-820.

21. Farooq, A.; Anwar, S.; Awais, M.; Rehman, S. In A deep CNN based multi-class classification of Alzheimer's disease using MRI, 2017 IEEE International Conference on Imaging systems and techniques (IST), IEEE: 2017; pp 1-6.

# Chapter 6
# Summary and Future Study

## 6.1 Summary

This thesis discusses the outstanding features of a novel graphene-Au pyramid hybrid platform and its applications in bio-medical disciplines. This SERS substrate has overcome the limitation of traditional SERS substrates and is highly bio-compatible, which opens up the possibility for bio-sensing using SERS. After combining with advanced data analysis methods, disease diagnostic has been realized with high accuracy.

In this thesis, we have validated the quantification capability as well as the specificity of SERS substrate on protein level and have further applied these benefits to differentiate Aβ peptides at multiple time stages. Alzheimer's disease patients are further distinguished from normal individuals using machine learning algorithms by testing the "fingerprints" of their CSF.

The following is a summary of the resulting work presented herein:

**In Chapter 2 and Chapter 3**, the features of hybrid platform are validated. Outstanding quantification capability as well as specificity regardless of protein orientation has been proven. Two types of mechanism has been shown at different concentration regime. At high concentration, protein concentration is proportional to normalized peak intensity due to the built in marker of graphene; at lower concentration, protein concentration is proportional to the detection frequency of the analyte considering the coverage of SERS hotspots. $10^{-18}$ M detection limit is achieved using Aβ42 and 7 orders of magnitude dynamic range is shown using the same analyte. To prove the specificity of protein using SERS, orientation dependence is tested using the hybrid platform and extremely low standard deviation (<0.3) is shown both from experimental and simulation results. The low standard deviation shows that the SERS signal is stable for proteins regardless of their orientation.

**In Chapter 4**, application of SERS hybrid platform is demonstrated using Aβ peptides. By applying principal component analysis (PCA), the two types of Aβ is differentiated using a 2 dimensional plot and the differences are further elaborated using the peak assignment. Decision tree is applied to tell the main differences between the spectra of each peptide and better show the differences from biological point of view. Monitoring the aggregation process is a major topic in Aβ studies and we benchmark the SERS spectrum change with other techniques (CD and TEM) to show that SERS is a faster and equally accurate method for protein analysis.

**In Chapter 5**, SERS application of the hybrid platform is extended to disease diagnostic. Cerebrospinal fluids of different individuals are collected and tested using the platform and analyzed using the combined platform of SERS and advanced data analyisis. Reproducibility of the system is proved using double blind experiment of 3 replicates of 5 different samples at 100% accuracy. Spectra features of 26 individuals were further tested and the diagnostic accuracy for Alzheimer's disease has reached over 90%. The correlation between

## 6.2 Direction of future studies

The potential future work on this study could focus on two aspects: (1) Apply the diagnostic system (hybrid platform together with data analysis methods) into other diseases to increase diagnostic capability; (2) Collect more patient data and expand the application into clinical usage.

### *6.2.1 Increase diagnostic capability of SERS*

Due to the outstanding capabilities of SERS hybrid platform, it is reasonable and natural next step to apply the system into the diagnostic of other diseases. Considering the ultra-high sensitivity and the molecular specificity, diseases with target biomarker can be easily distinguished.

Exosomes are extracellular vesicles that are produced in the endosomal compartment of most eukaryotic cells, and it is well studied that cell at different stages carries

exosomes with slightly different features. This is a perfect opportunity to apply SERS based diagnostic.

Some preliminary results have been acquired using exosome. We compared the spectra of exosomes from different sources: the exosome from human serum (Fig 6.1 A) and 2 different types of conditioned tissue-culture medium of a human lung cancer cell line HCC827 (Figure 6.1 B) and H1975 (Figure 6.1 C).
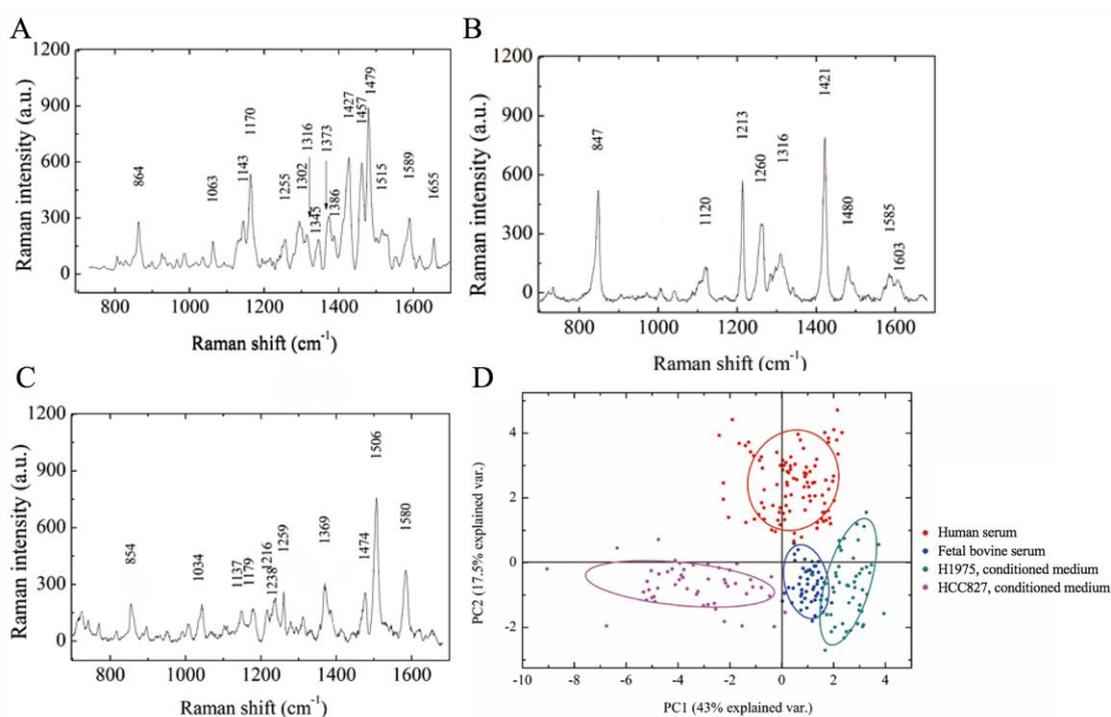


Fig 6.1. SERS analysis and PCA of exosomes from different sources. (A−C) Raman spectra of exosomes from human serum (A) conditioned medium of the lung-cancer cell line HCC827 (B), and conditioned medium of the lung-cancer cell line H1975 (C). (D) PCA of exosomes from the different sources shown in panels A−C and the spectrum shown in Figure 3C demonstrating that they are distinguishable.

Each sample showed uniquely identifiable spectral characteristics distinguished

primarily by the relative peak intensity.

We then used PCA to analyze spectral differences and similarities in ~50 Raman spectra from each sample (Figure 6.1 D). The results showed that the exosomes from all four different sources clustered into distinguishable groups with 84%. Interestingly, the largest degree of overlap was not between the two sera or the two cell lines, but between FBS and the H1975 cell line (Figure 6D). These findings suggest that analysis of exosomes from the serum of two different species, cell culture media versus serum, and cell culture media from two cancer cell lines of the same human organ, lung, can be distinguished using our platform.

These experiments, together with our previous anlaysis suggest that our platform have the potential for multiple diseases diagnostic.

## 6.2.2 Expand into clinical application

Though the capability for disease diagnostic (such as Alzheimer's disease) has been addressed using the hybrid platform, the application of the system in clinical diagnostic is still unknown. Clinical application remains challenging due to the mass pre-clinical trial data required. The clinical trials may compare our new technology to a standard diagnostic method that is already available and the whole process is defined by the Food and Drug Administration (FDA). For diagnostic regime, the clinical trial refers to the practice of looking for better ways to identify a particular

disorder or condition and a diagnostic tool with high accuracy and high reproducibility is required.

To address the key issues, our future research can focus on two parts: (1) recruit larger number of patients and increase the robustness of our diagnostic system; (2) improve the performance of our data analysis system to increase the diagnostic accuracy.

To make our diagnostic result robust and to make the diagnostic statistically reasonable, we have started our collaboration with University College of Faisalabad in Pakistan under Pakistan-U.S. Science & Technology Cooperation Program to acquire more patient sample for study purpose. As patient number grows, we are able to train the data analysis system with more information and thus making the prediction model more accurate.

Another direction for improving the system is to improve the data analysis algorithms. A basic diagnostic model have been built and proved effective, however, a lot remains to be done. Data pre-processing has remained to be problematic (noise filter, background removal, etc.) and the functionality of the current model is limited to yes/no differentiation.

A natural next step for this research is to improve the capability for the diagnostic model, including adding preprocessing functions and make more detailed diagnostic predictions, such as the stage of the disease or the subtype of the disease. These can be done with more training data and with improved analysis algorithms.