



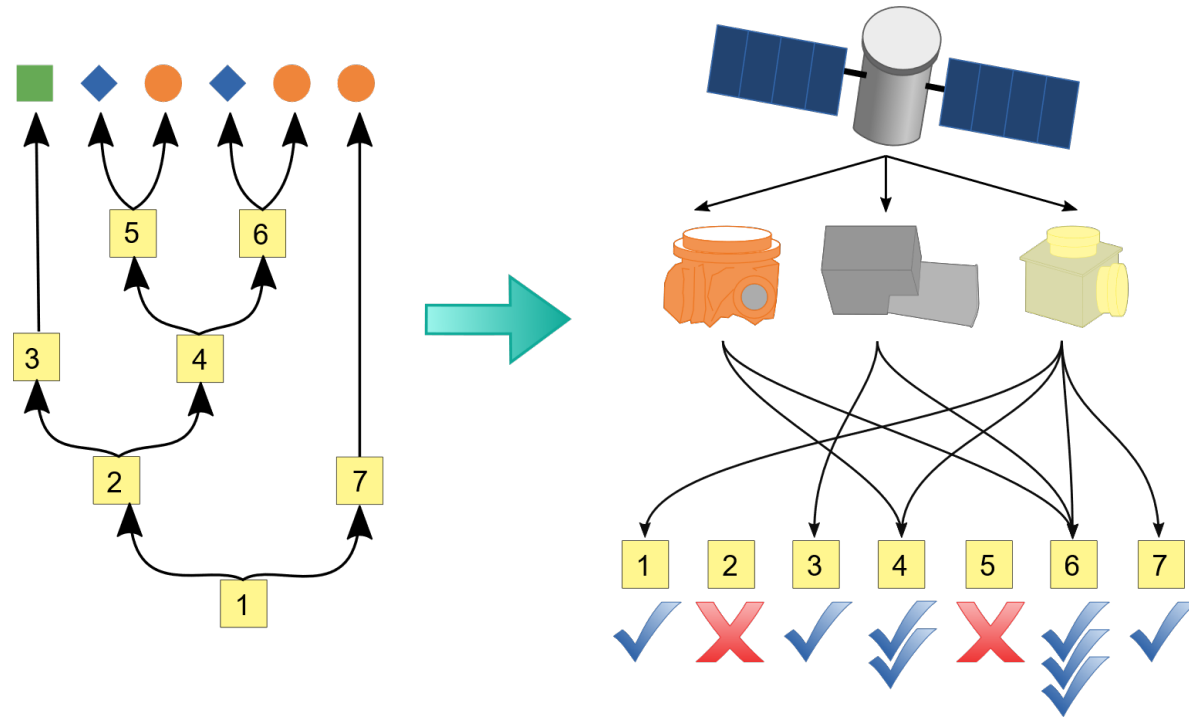
Algorithmic Classification of Raman Spectra Biosignatures: *Improving Life-Detection Confidence*

AbSciCon 2022 Conference
Atlanta, GA
May 18, 2022

Tao Sheng¹, Aarya Mishra², Jesse B. Murray³, Sunanda Sharma⁴, Diana Gentry⁵

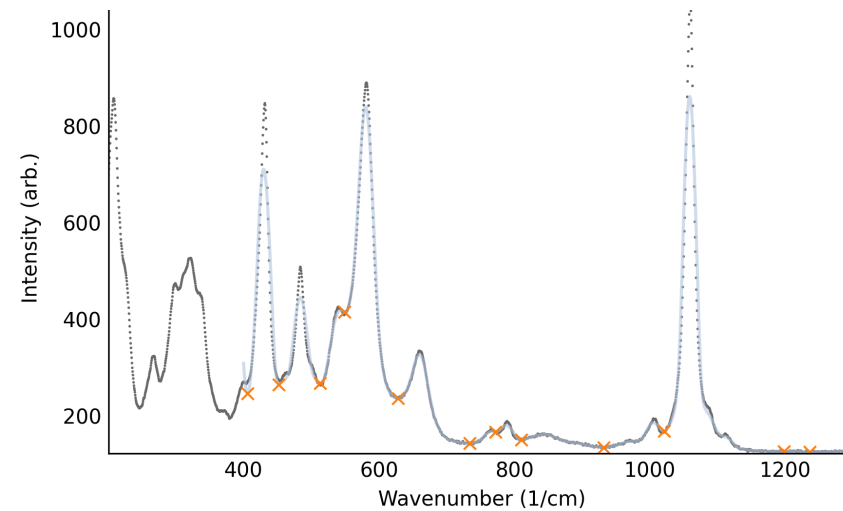
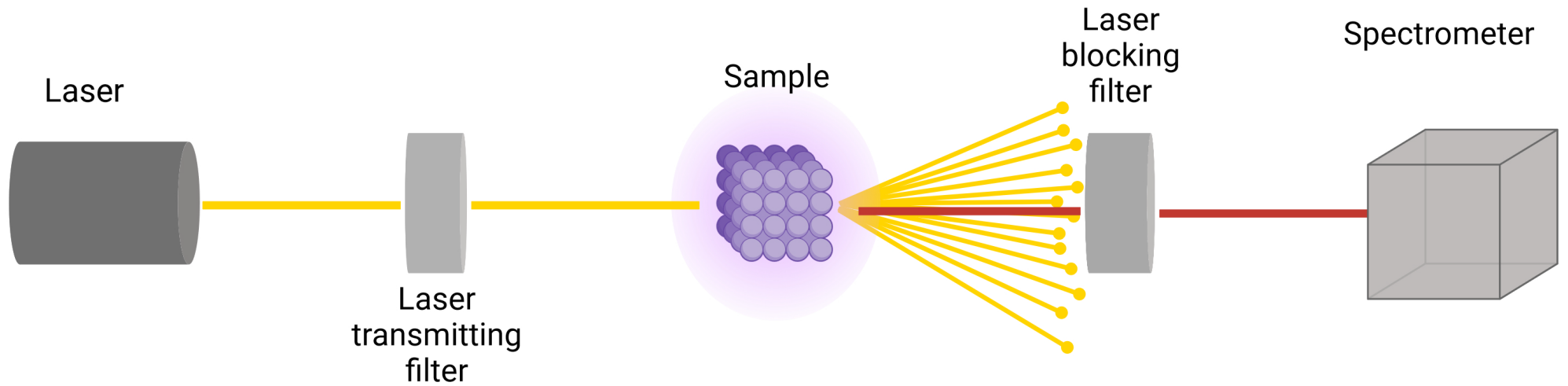
NASA Space Life Sciences Training Program¹, NASA Volunteer Internship Program¹, University of Pittsburgh¹, University of San Francisco², NASA Volunteer Internship Program², University of Oxford³, Massachusetts Institute of Technology⁴, NASA Ames Research Center⁵

Background

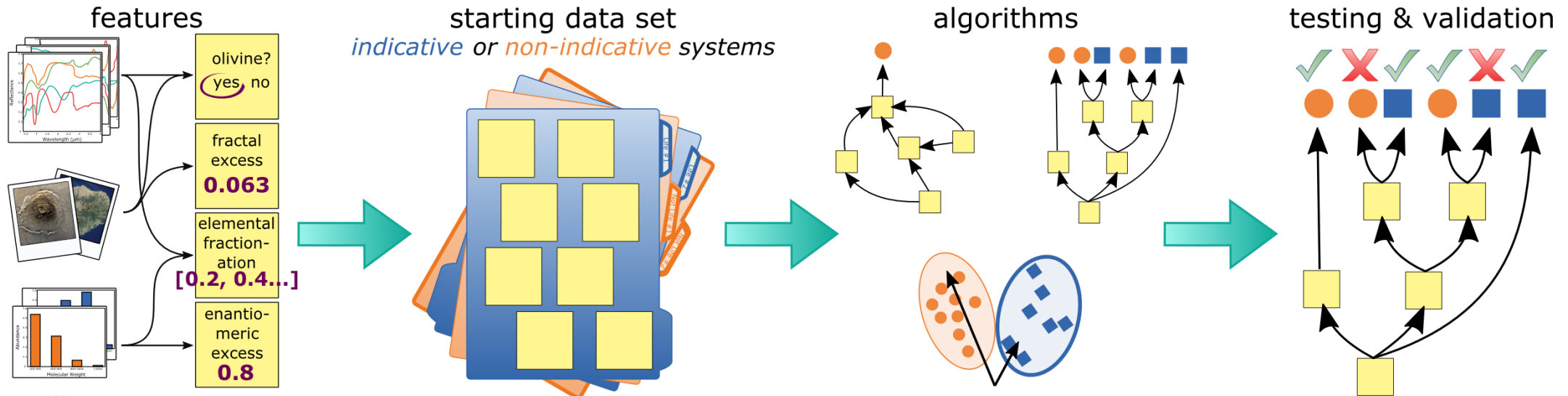


Purpose: Algorithmically combine multiple, different instrument readings of the same potential biosignature to determine whether it is indicative or not indicative of life.

Raman Spectroscopy Basics



Overall Project Workflow



- Elemental abundance
- Isotopic fractionation
- VNIR reflectance spectra

- Collecting data on various materials

- Feature engineering
- Supervised learning

- Calculate accuracy, confidence, etc.



1) Data Selection

- **Previous data types:**

- Elemental abundances, isotopic fractionation, visible and near infrared spectroscopy (VNIR)

- **Conditions**

- Data collection feasibility in the field
- Somewhat agnostic
 - Data type should not presuppose any specific biochemistry or molecular framework
 - Ex: Rover looking specifically for chlorophyll would not be agnostic
 - Ex: Rover looking for molecules that stores energy that ends up including ATP could be agnostic
 - Would be “putting the answer in the question”
- Note on elemental abundance: the dataset is not perfectly agnostic because information on C, H, O, etc. allows bias towards C and water-based life
 - A more agnostic approach was taken to look at elemental distributions; this method no longer depends on C-based life and high amounts of water



2) Data Collection

- **Curated indicative vs. non-indicative dataset**

- Distinctions unambiguous
- No edge cases (prions, etc.)
- Does not require an *a priori* definition of life

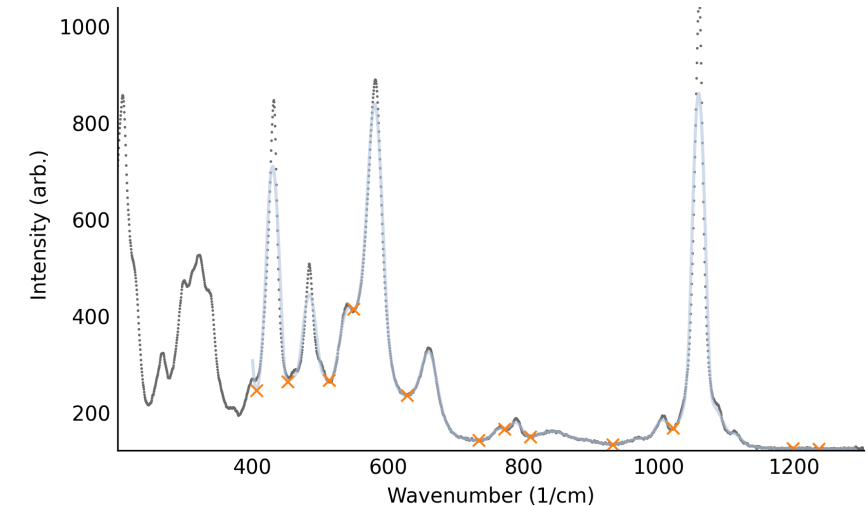
- **Raman data collection**

- Assembled from individual research articles and public domain databases
- “Standardized” to similar excitation lines (532nm or 514.5nm)
- Pre-processing: Savitzsky-Golay smoothing + interpolation, Polyfit baseline correction, Standard normal variate normalization

Class (subclass)	Class Definition	Examples
Indicative (alive)	Indicative sample of only alive materials	Vegetation, microorganisms, biofilm
Indicative (non-alive)	Indicative sample of only not-alive samples	Coral skeletons, coal, calcite
Indicative (mixed)	Indicative samples mixed with non-indicative	Non-sterilized silt, clay, or seawater
Non-indicative	Not alive and not indicative of life	Pure basalt, sand, or carbonatite

3) Feature Engineering & Algorithm Implementation

- **Feature engineering**
 - Inaccurate feature engineering = inaccurate classifier
- **How can we describe spectral / Raman data agnostically?**
 - Preliminary features: highest/lowest intensity, range, mean intensity, min/max 1st/2nd derivatives, # of peaks, # of troughs, broadest peak, mean peak width
- **Automated feature engineering**
 - Image classification promising technique
- **Supervised learning**
 - K-nearest Neighbors (KNN), Logistic Regression (LR) with L2 regularization, Support Vector Machines (SVM) with principal component, Gaussian Naïve Bayes (GNB), Random Forest (RF), Mean Votes



Example Raman Spectra of Hydroxyapophyllite, troughs detected as orange X marks

4) Testing Metrics & Validation

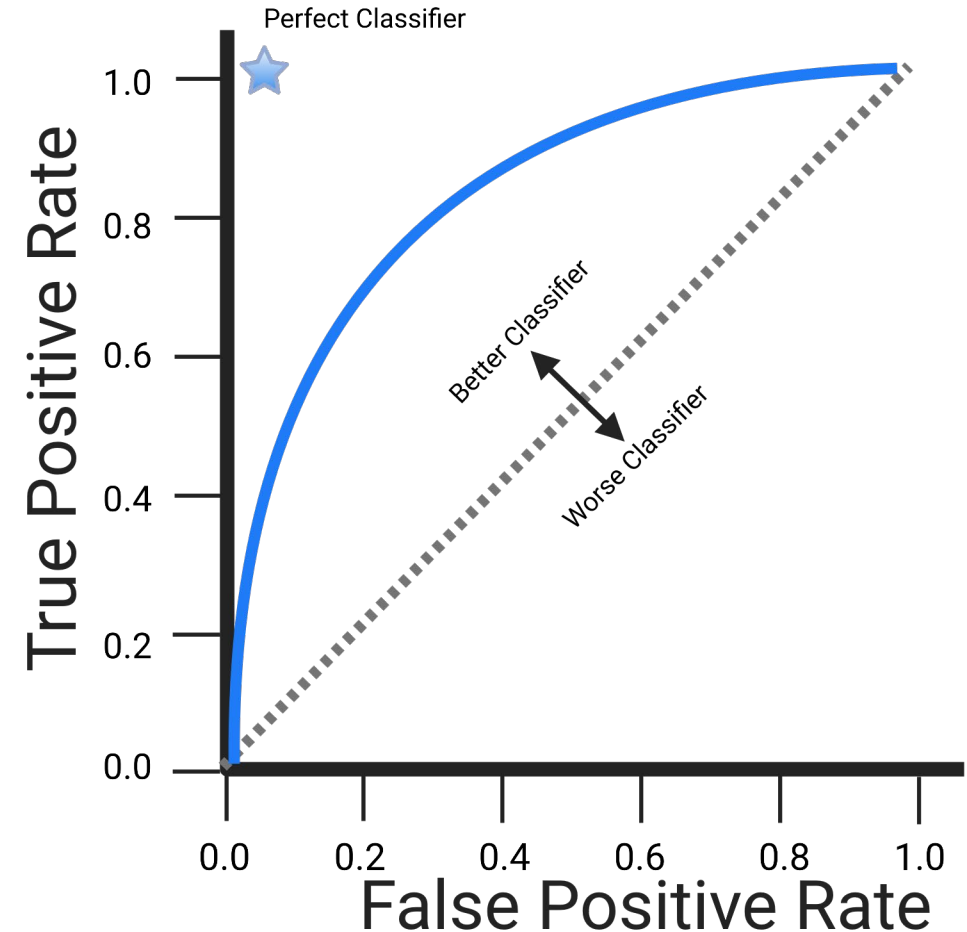
- **Accuracy metrics**

- AUC: area under receiver operating curve (ROC)
 - Probability of correct classification
- Done with 2,000 50% train-test-splits

- **False negative and false positive rates (FNR, FPR)**

- Important in mission contexts where algorithm is deciding which samples to prioritize to evaluate further – erring towards more inclusion could be better
 - Indicated by slightly higher FPR than FNR

Sample Receiver Operating Curve (ROC)



Made with Biorender

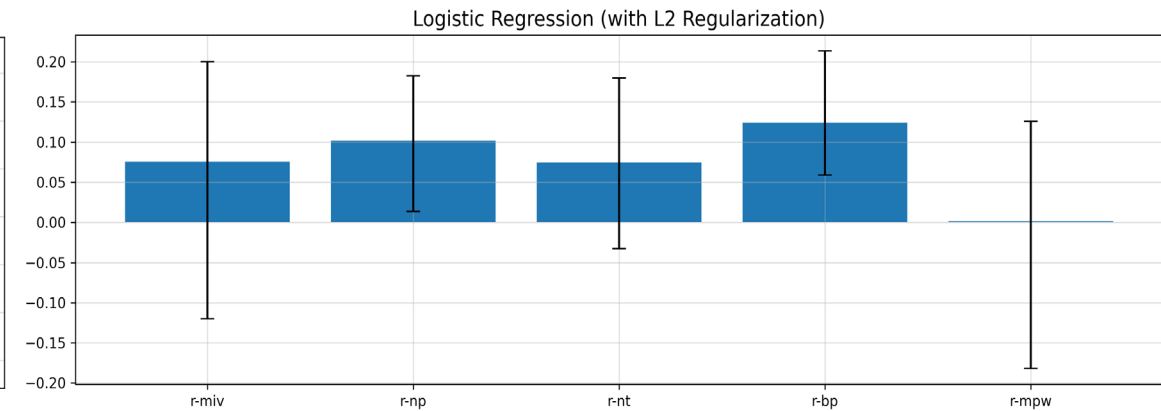
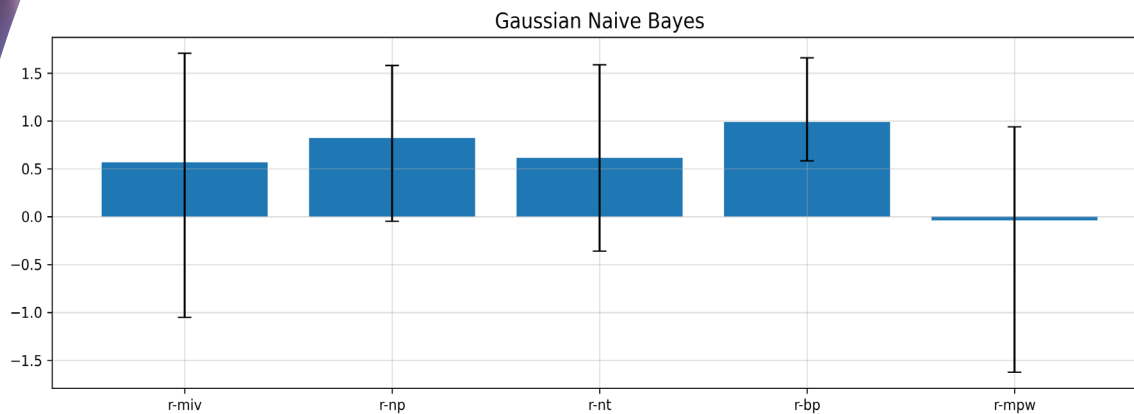


Raman Only Performance

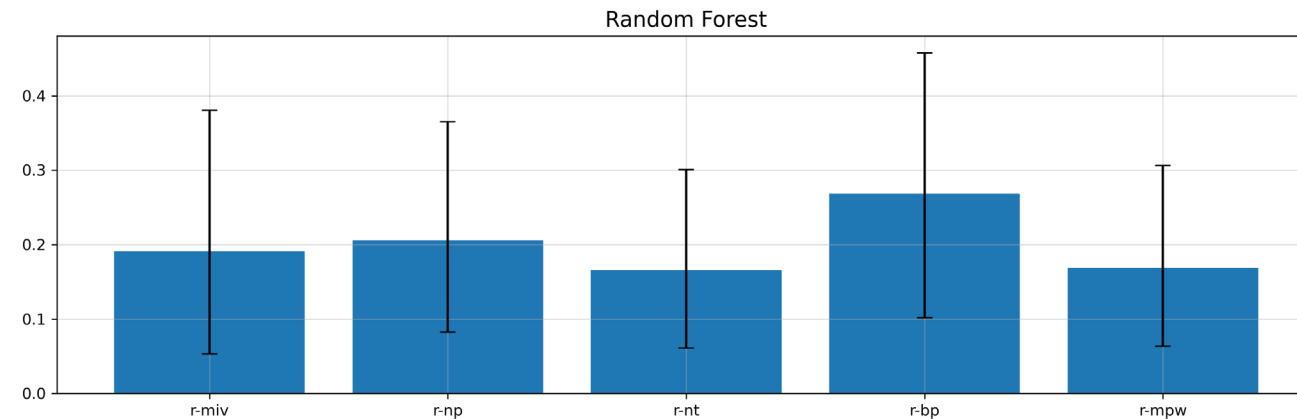
<i>Algorithms</i>	<i>Raman Only</i>	<i>Raman Only, <u>no mixed samples</u></i>
K-nearest neighbors (KNN)	0.772	0.789
Logistic regression (LR)	0.556	0.684
Support vector machine (SVM)	0.522	0.635
Gaussian naïve Bayes (GNB)	0.766	0.795
Random Forest (RF)	0.770	0.814
Mean voting performance (avg.)	0.808	0.860

Feature Importance: Raman Data

Coefficients Feature Importance

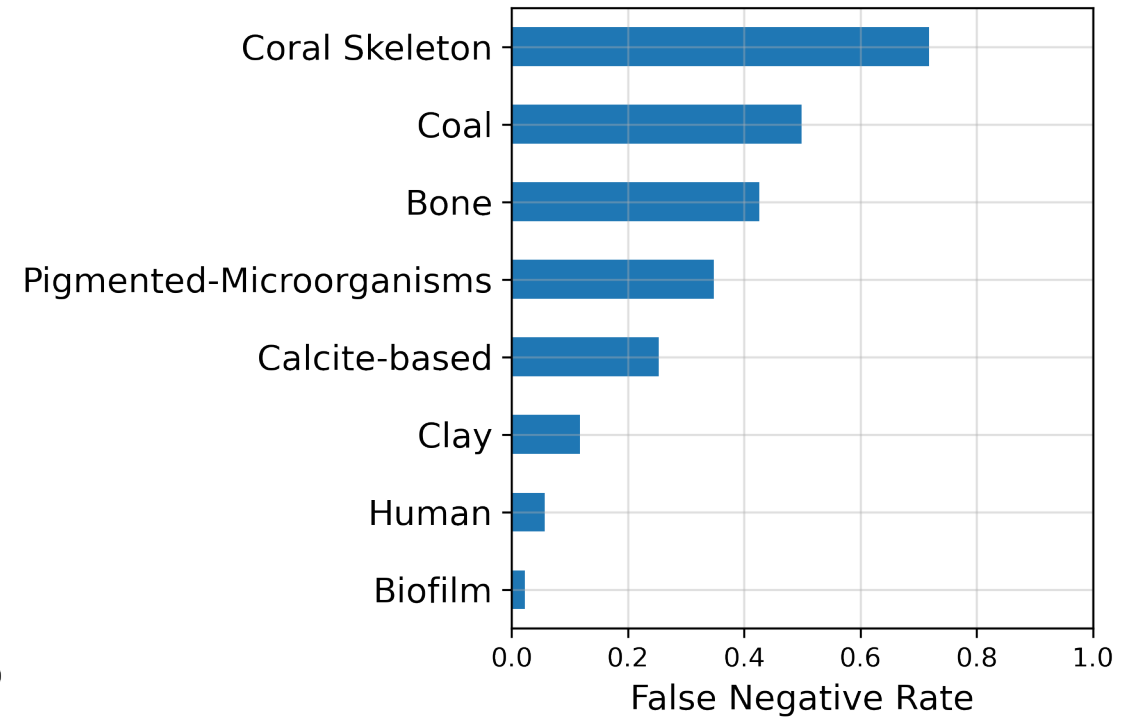
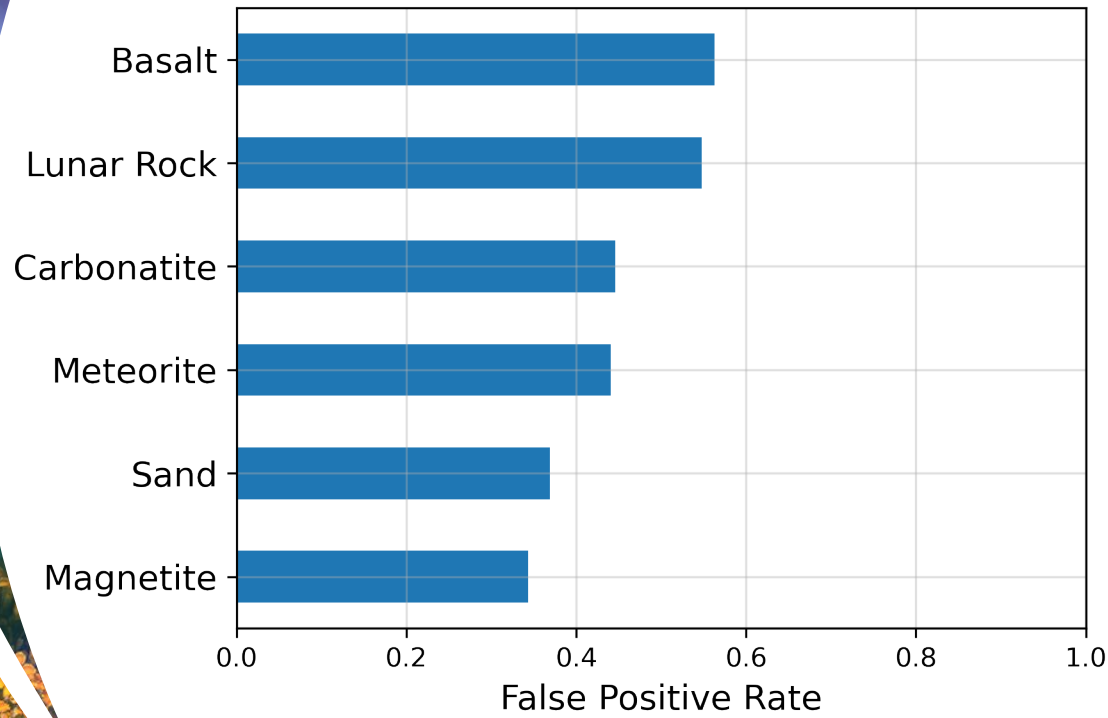


Decision Trees Feature Importance



Raman Common Misclassifications

- **Systems are disproportionally misclassified if they are indicative “mixed” or indicative not-alive**
 - Algorithm can more effectively classify between living vs. not living, but not so much when non-alive biosignatures are involved



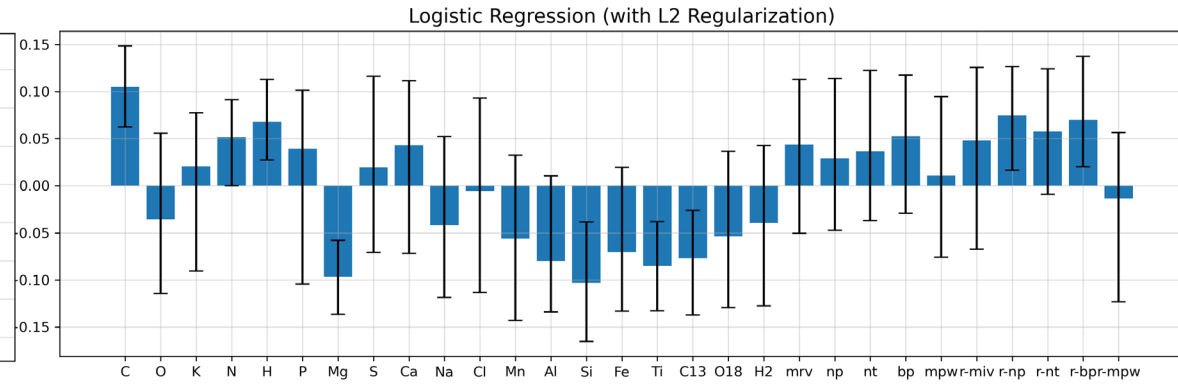
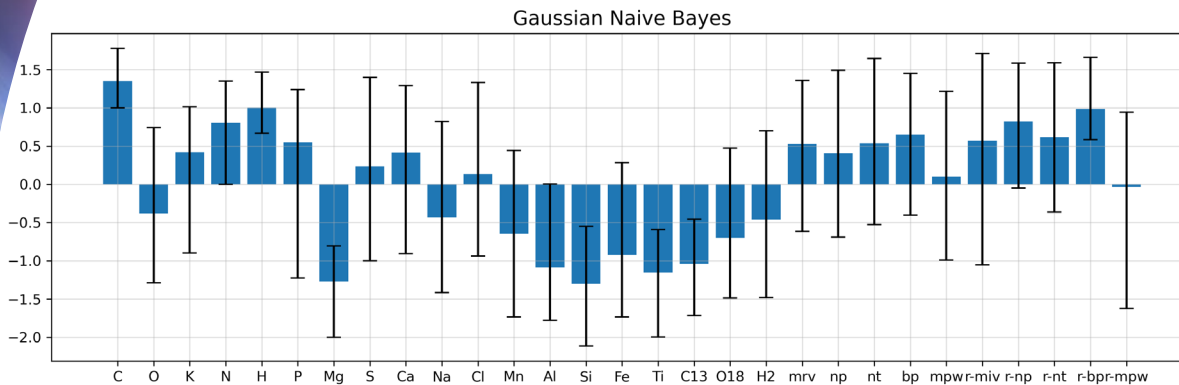


Combination Classification Performance

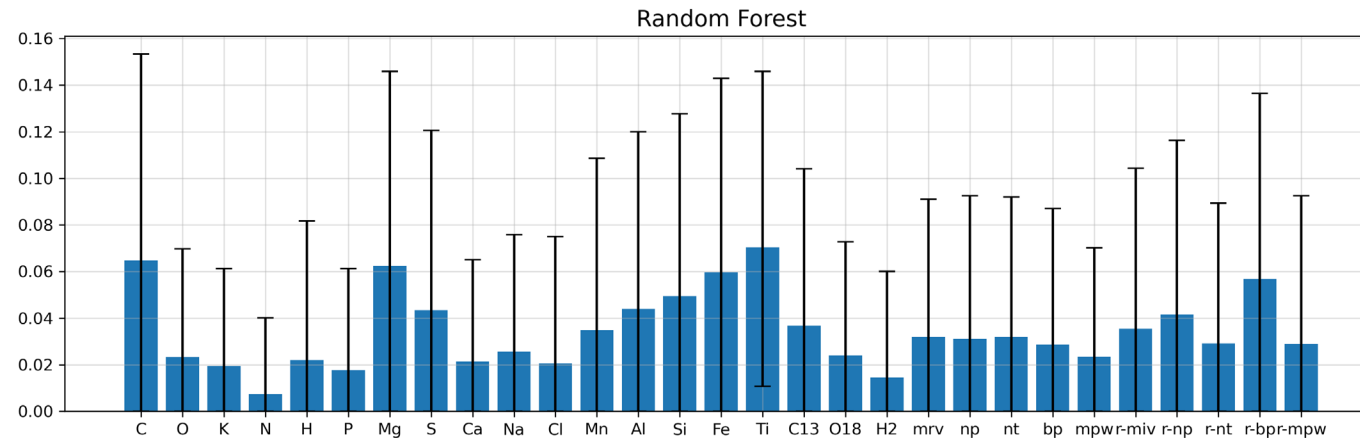
<i>Algorithms</i>	<i>Raman Only</i>	<i>Combination (Element, Isotope, Reflectance, Raman)</i>	<i>Combination (Element, Isotope, Reflectance, Raman) <u>without mixed samples</u></i>
K-nearest neighbors (KNN)	0.772	0.759	0.798
Logistic regression (LR)	0.556	0.743	0.783
Support vector machine (SVM)	0.522	0.652	0.717
Gaussian naïve Bayes (GNB)	0.766	0.693	0.714
Random Forest (RF)	0.770	0.856	0.891
Mean voting performance (avg.)	0.808	0.826	0.855

Feature Importance: Combined Data

Coefficients Feature Importance

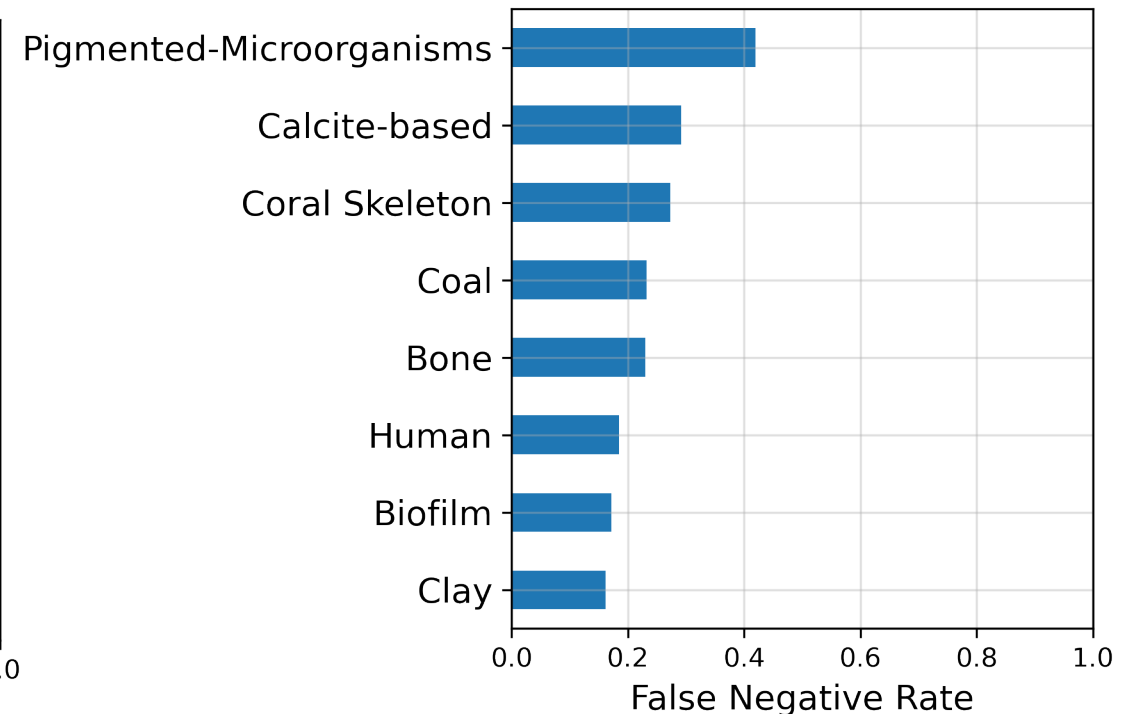
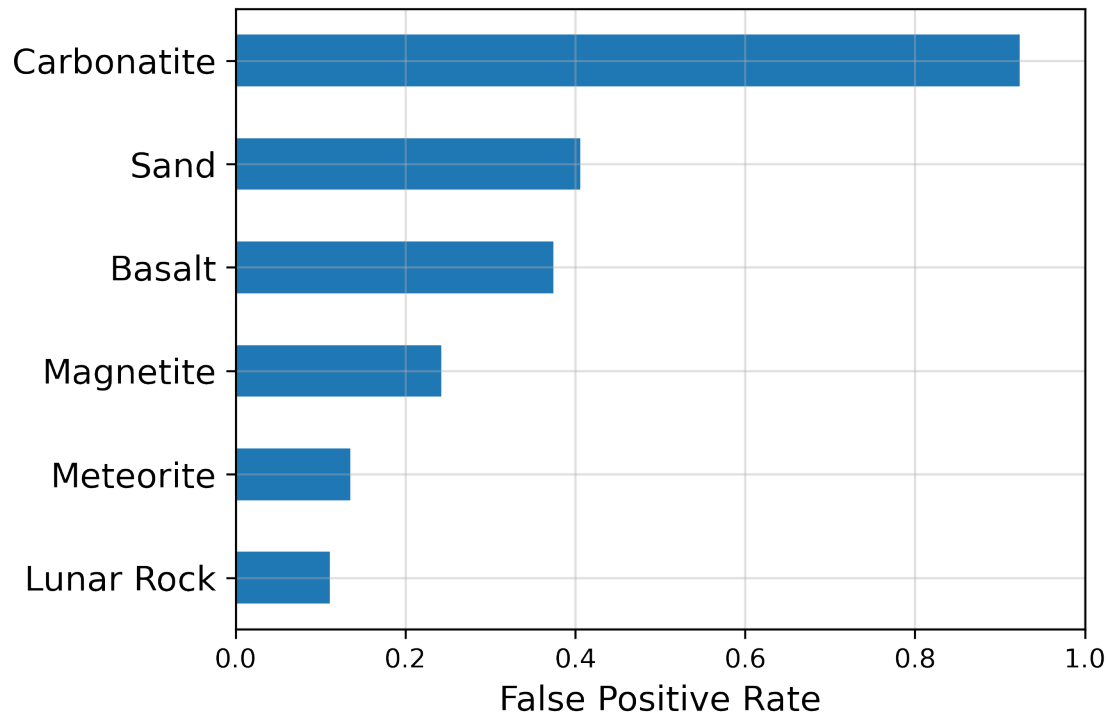


Decision Trees Feature Importance



Combined Common Misclassifications

- **Systems are disproportionally misclassified if they are indicative “mixed” or indicative not-alive**
 - Algorithm can more effectively classify between living vs. not living, but not so much when non-alive biosignatures are involved





Conclusions & Future Directions

- **Conclusions**

- Manually engineered Raman features not very discriminatory, suggests need for other feature methods
 - Combined features rely heavily on C and H based life and ^{13}C delta ratio
- Raman alone and combined performance improves without indicative-mixed samples
 - Mixed samples are more challenging to classify indicative alive or non indicative which suggests these supervised algorithms are better at life-detection but not so much with not-alive biosignatures

- **Future steps**

- Image classification – automating feature extraction and improving VNIR reflectance and Raman classification
- Further data collection of more “indicative / indicative mixed” raw Raman sample data
- Further agnosticization with elemental abundance data to elemental distributions
- Investigation into commonly misclassified systems



Acknowledgements

- Gabriela Peña Carmona and Joseph Stone and the BeING lab for early feedback
- SLSTP 2021 Management, Staffers, and fellow Research Associates
- VIP coordinators Porsche Parker and Haley Feck
- NASA Ames Project EXcellence (APEX) incubator program
- Amanda Stockton, Tori Hoehler, Tony Ricco, Michael Bicay, Dave Korsmeyer, Jaya Bajpayee, Jay Bookbinder, and Carol Carroll for valuable early feedback
- Jessica Koehne, David Mauro, Annmarie Schramm, P. Michael Furlong, and Thomas Stucky for helping shape the initial effort
- Aivaras Vilutas and Paxton Tomko for helping with the data collection and standardization

Questions?

TAS242@pitt.edu | tao.sheng@nasa.gov





References I

- Marc Neveu, Lindsay E. Hays, Mary A. Voytek, Michael H. New, and Mitchell D. Schulte. The ladder of life detection. *Astrobio*, 18(10):1–28, June 2018.
- Lafuente B, Downs R T, Yang H, Stone N (2015) The power of databases: the RRUFF project. In: Highlights in Mineralogical Crystallography, T Armbruster and R M Danisi, eds. Berlin, Germany, W. De Gruyter, pp 1-30
- El Mendili, Y., Vaitkus, A., Merkys, A., Gražulis, S., Chateigner, D., Mathevet, F., Gascoin, S., Petit, S., Bardeau, J.-F., Zanatta, M., Secchi, M., Mariotto, G., Kumar, A., Cassetta, M., Lutterotti, L., Borovin, E., Orberger, B., Simon, P., Hehlen, B., & Le Guen, M. (2019). Raman Open Database: first interconnected Raman–X-ray diffraction open-access resource for material identification. *Journal of Applied Crystallography*, 52(3), 618-625. doi: [10.1107/s1600576719004229](https://doi.org/10.1107/s1600576719004229)
- *Handbook of Raman Spectra for geology*, Laboratoire de géologie de Lyon ENS de Lyon France, 2000, www.geologie-lyon.fr/Raman/, Date Accessed: June 11, 2021
- Nicolae BUZGAR, Andrei Ionut APOPEI, Andrei BUZATU (2009) - Romanian Database of Raman Spectroscopy (<http://rdrs.uaic.ro>).



References II

- S. B. Kotsiantis, I. D. Zaharakis, and P. E. Pintelas. Machine learning: A review of classification and combining techniques. *Artificial Intelligence Review*, 26(3):159–190, November 2006.
- PDS Geosciences Node. CRISM spectral library. <https://speclib.rsl.wustl.edu>, accessed 2021-04-06.
- Raymond F. Kokaly, Roger N. Clark, Gregg A. Swayze, K. Eric Livo, Todd M. Hoefen, Neil C. Pearson, Richard A. Wise, William M. Benzel, Heather A. Lowers, Rhonda L. Driscoll, , and A. J. Klein. USGS Spectral Library Version 7. Number 1035 in Data Series. U.S. Geological Survey, 2017.
- S. Hegde, I. G. Paulino-Lima, R. Kent, L. Kaltenegger, and L. Rothschild. Surface biosignatures of exo-Earths: Remote detection of extraterrestrial life. *Proceedings of the National Academy of Sciences of the United States of America*, 112:3886–3891, March 2015.
- S. K. Meerdink, S. J. Hook, D. A. Roberts, and E. A. Abbott. The ECOSTRESS spectral library version 1.0. *Remote Sensing of Environment*, 230:1–8, September 2019.
- PDS Geosciences Node. RELAB spectral database, 2014.
<http://www.planetary.brown.edu/relabdocs/relab.htm>, accessed 2021-04-06
- C. Cooksey. Reference data set of human skin reflectance. National Institute of Standards and Technology, 2017.