# Using concepts in natural language processing to determine the composition of artifact residues

Tung Tho Nguyen[1, *] and Korey J. Brownstein[2, *]

[1]The University of Chicago, Department of Mathematics, Chicago, Illinois, 60637, USA
[2]The University of Chicago, Department of Molecular Genetics and Cell Biology, Chicago, Illinois, 60637, USA

*Correspondence:
Tung Tho Nguyen
tungnt@uchicago.edu

Korey J. Brownstein
kbrownstein@uchicago.edu

**Abstract**

Determining the composition of artifact residues is a central problem in ancient residue

metabolomics. Traditionally, this is done by comparing mass spectra features in common with an

experimental pipe and a sample pipe (classical method). While this method is simple and

straightforward, its prediction capabilities might be inaccurate. Here, we introduce a novel

approach based on ideas from the field of natural language processing to solve this problem. We

tested our strategy on a set of modern clay pipes. To limit biases, we were not provided

information on which plant species had been smoked in which clay pipe. The results indicate that

our algorithms performed 12.5% better than the previously published classical method.

**Introduction**

Metabolomics is the systematic quantitative and qualitative study of small molecules (or mass spectral features) in biological systems. Brownstein *et al.* (2020) expanded upon this field with their ancient residue metabolomics-based method. Albeit the mass spectral features in this study were not derived from biological systems, they were residues left behind from biological processes; i.e., originating from plants including several *Nicotiana* species that were smoked by indigenous peoples. Before the Brownstein *et al.* (2020) study, ancient residue analysis relied on the biomarker approached. However, the biomarker approach failed to distinguish between related species, leaving open questions about the relationship between plants and people. For ancient residue metabolomics, all compounds are of interest improving the resolution of determining which plants species had been smoked in a particular pipe (Brownstein *et al.*, 2020).

In short, data from hyphenated chromatography instruments (such as gas chromatography- and liquid chromatography-mass spectrometer) are processed and aligned in MZmine 2 (Pluskal *et al.*, 2010), Progenesis QI (Waters Corporation, Milford, MA, USA), or another "omics" software. Afterwards, these data are exported from the software and then processed manually as described in the following passage from Brownstein *et al.* (2020): "[T]he dataset was exported into Microsoft Excel and mass spectral features shared with the blank [samples] were removed from the analysis. The [three solvent] extracts from each sample were combined into a single compound list and compounds with no abundance values were removed. The [ancient samples] were then compared to the [experimental samples using a Venn diagram]." This final component of processing and analyzing the datasets determines which plant species may have been used in an ancient artifact. Because this step requires a manual process, it can introduce errors and is

53  time consuming. Various metabolomics platforms exist including MetaboAnalyst 4.0 (Pang *et*

54  *al.*, 2021) or XCMS Online (Tautenhahn *et al.*, 2012); however, these platforms are limited in

55  their ability to process datasets from ancient residue studies. Therefore, we introduce a novel,

56  automated method for determining the composition of organic residues in modern smoking

57  pipes. Our approach is inspired by techniques and ideas from the field of natural language

58  processing (NLP).

59

60  **Materials and Methods**

61  A novel functionality of our approach is to introduce an entirely new method to compare the

62  mass spectra feature similarities between experimental and ancient artifacts (i.e., sample groups).

63  Our approach is inspired by advances in NLP (Cohen *et al.*, 2004; Cong *et al.*, 2017; Goodfellow

64  *et al.*, 2016; Young *et al.*, 2018). Here we use the following analogy:

65

66                              Words ←→ Mass Spectra Features

67                              Documents ←→ Sample Groups

68

69  The standard technique in NLP is to first transform the original data into the term frequency-

70  inverse document frequency (TF-IDF) matrix (Goodfellow *et al.*, 2016). This transformation

71  helps to resolve the fact that some substances appear more often than others. More precisely, the

72  importance of a term is not solely determined by its frequency in a text (TF) but also how rare

73  this term is in other texts in the corpus (IDF). Let us recall these terminologies mathematically.

74  Term frequency refers to the frequency of a word in a particular document:

75

$$tf(w, d) = \frac{\text{count of w in d}}{\text{number of words in d}}$$

The inverse of the document frequency which measures the informativeness/prevalence of term t

$$idf(w) = \log\left(\frac{N}{(df(w) + 1)}\right),$$

where N is the number of documents and df(w) is the number of documents containing w. IDF score depends on the occurrence of terms and not on their numerical frequencies. Once the TF-IDF is computed, we can then use cosine similarity to compare two different groups. Recall that for two vectors $v, w$ their cosine similarity is defined to be cosine of the angle $\theta$ between them, namely

$$\text{similarity} = \cos(\theta) = \frac{\langle v, w \rangle}{||v||||w||}$$

Here, $\langle v, w \rangle$ is the inner product of $v, w$ and $||v||, ||w||$ is the Euclidean norm of $v, w$. We note similarity score ranges from -1 meaning exactly opposite to 1 meaning exactly the same, with 0 indicating orthogonality, while in-between values indicate intermediate similarity or dissimilarity.

*Artemisia ludoviciana* Nutt. (*Alu*) leaves, *Arctostaphylos uva-ursi* (L.) Spreng. (*Auv*) leaves, *Cornus sericea* L. (*Cse*) bark, *Gaultheria shallon* Pursh (*Gsh*) leaves, *Lobelia inflata* L. (*Lin*) leaves, *Nicotiana attenuata* Torr. ex S. Watson (*Nat*) leaves, *Nicotiana glauca* Graham (*Ngl*)

98    leaves, *Nicotiana obtusifolia* M. Martens & Galeotti (*Nob*) leaves, *Nicotiana quadrivalvis* Pursh

99    (*Nqu*) leaves, *Nicotiana rustica* L. (*Nru*) leaves, *Nicotiana tabacum* L. (*Nta*) leaves, *Rhus glabra*

100   L. (*Rgl*) autumn leaves, *Salvia sonomensis* Greene (*Sso*) leaves, *Taxus brevifolia* Nutt. (*Tbr*)

101   needles, and *Verbascum thapsus* L. (*Vth*) leaves were collected, freeze-dried for 3 days, and

102   crushed for experimental smoking. American Spirit (AmSp) tobacco (Santa Fe Natural Tobacco

103   Company, Oxford, NC, USA) was purchased from a local grocery store in Pullman, Washington,

104   USA.

105

106   The plant materials and AmSp were smoked following the experimental conditions detailed in

107   Brownstein *et al.* (2020). To limit biases, the authors did not know which plant species had been

108   smoked in which clay pipe. After the samples were analyzed by liquid chromatography-mass

109   spectrometry and processed in MZmine 2 (Pluskal *et al.*, 2010) following the parameters

110   described in Brownstein *et al.* (2020), the data were exported into .csv files. Python libraries,

111   such as Sklearn and Pandas, were then used to apply the TF-IDF computation scores to these

112   datasets.

113

114   **Results and Discussion**

115   We used Python to write the scripts because of the availability of several useful data analysis,

116   machine learning, and deep learning libraries. All the scripts and datasets are freely available on

117   GitHub: https://github.com/tungprime/NLP_and_composition_of_artifact_residues. Our script

118   automates the classical method described in Brownstein *et al.* (2020), as well as utilizes recent

119   advances in machine and deep learning to better predict which plant species had been smoked in

120   a particular artifact (new method). As shown in Table 1, the new method predicts that CP1 was

121    most likely smoked with *Nta* (0.0370). Table 2 summarizes the model predictions of the classical

122    and new methods, and the key provides the expected results. In fact, CP1 was smoked with *Nta*.

123    While the classical method only predicted four out of eight (50.0%) of the samples correctly, the

124    new method performed slightly better, i.e., it classified five out of eight (62.5%) of the samples

125    correctly (Table 2).

126

127    **Table 1.** Similarity scores of clay pipe 1 (CP1) smoked with an unknown plant sample. Sixteen
128    (16) different experimental pipes smoked with only one of the plant species or AmSp were
129    individually compared to CP1. The top five scores were only included in the table. *Nta* is the
130    most likely candidate smoked in CP1. *Nat*, *Nicotiana attenuata*; *Ngl*, *Nicotiana glauca*; *Nob*,
131    *Nicotiana obtusifolia*; *Nta*, *Nicotiana tabacum*; and AmSp, American Spirit.

|  | *Nta* | *Ngl* | *Nat* | AmSp | *Nob* |
|---|---|---|---|---|---|
| Similarity scores | 0.0370 | 0.0287 | 0.0232 | 0.0185 | 0.0170 |

132

133

134    Contamination is a significant concern for ancient residue metabolomics (Damitio *et al.*, 2021;

135    Zimmermann *et al.*, 2021). For instance, residues from commercial tobacco smoke may

136    contaminate the surface of artifacts at excavation sites or on display at a museum. Thus, we

137    included AmSp in our study as a contaminate control. With the contaminate control (0.0185), we

138    were still able to accurately determine the composition of CP1 (Table 1) and the other clay pipes.

139    Utilizing contaminate controls will improve confidence in whether or not a particular artifact had

140    been smoked with an endemic tobacco. Furthermore, our new method will enable researchers to

141    confidently determine if the caffeine present in/on an artifact resulted from ancient cacao or holly

142    brewing practices instead of modern contaminates from caffeinated beverages such as coffee

143    (King *et al.*, 2017; Washburn *et al.*, 2014).

144

145    **Table 2.** Predicted plant species in each clay pipe (CP). Sixteen (16) different experimental pipes
146    smoked with only one of the plant species or AmSp were compared individually to each CP. The

147 expected results are the plant species listed under key. *Alu*, *Artemisia ludoviciana*; *Auv*,
148 *Arctostaphylos uva-ursi*; *Cse*, *Cornus sericea*; *Lin*, *Lobelia inflata*; *Nat*, *Nicotiana attenuata*;
149 *Ngl*, *Nicotiana glauca*; *Nob*, *Nicotiana obtusifolia*; *Nqu*, *Nicotiana quadrivalvis*; and *Nta*,
150 *Nicotiana tabacum*.

| Clay pipe (CP) | Classical method (number of mass spectra features shared with the experimental pipe) | Key | New method (similarity score) |
|---|---|---|---|
| CP1 | *Auv* (3) | *Nta* | *Nta* (0.0370) |
| CP2 | *Nat* (6) and *Nqu* (6) | *Nqu* | *Nqu* (0.1106) |
| CP3 | *Nob* (10) | *Nob* | *Nob* (0.1145) |
| CP4 | *Nat* (4) | *Alu* | *Nta* (0.0934) |
| CP5 | *Ngl* (6) | *Ngl* | *Ngl* (0.0884) |
| CP6 | *Lin* (10) | *Lin* | *Ngl* (0.0844) |
| CP7 | *Auv* (13) | *Auv* | *Auv* (0.0735) |
| CP8 | *Auv* (9) and *Cse* (4) | *Auv* and *Nta* mixture | *Cse* (0.0813) and *Auv* (0.0617) |

151

152

153 It was also revealed that neither the new nor classical methods could accurately predict that CP8

154 had a mixture of *Auv* and *Nta* (Table 2). Both methods partially predicted the composition of

155 CP8. Though the classical method performed slightly better because it had more hits for *Auv* than

156 *Cse* (Table 2). Nonetheless, the experimental pipes compared to CP8 had only been smoked with

157 *one* plant species. It is possible that training the new method with experimental pipes smoked

158 with complex mixtures may improve the likelihood of predicting if a pipe had been smoked with

159 more than one plant species.

160

161 **Conclusion**

162 Machine and deep learning have been vital tools for solving problems in biology where

163 traditional methods seem inadequate or are time-consuming. Combining conventional and

164 machine learning-based methods to process and analyze data helps researchers gain a more in-

165 depth analysis of their data. The new method was used to predict which plant species had been

166 smoked in modern clay pipes, and we believe this method can be applied to ancient smoking

167    pipes, brewing vessels, and other artifacts.

168

180

**References**

Brownstein, K. J., Tushingham, S., Damitio, W. J., Nguyen, T., and Gang, D. R. (2020). An ancient residue metabolomics-based method to distinguish use of closely related plant species in ancient pipes. *Frontiers in Molecular Biosciences* 7, 133.

Cohen, K. B., and Hunter, L. (2004). Natural language processing and systems biology (pp. 147–175). In: Dubitzky, W., and Pereira, F. (eds.), *Artificial intelligence and systems biology*. Springer, Dordrecht, Netherlands.

Cong, Y., Chan, Y. B., Phillips, C. A., Langston, M. A., and Ragan, M. A. (2017). Robust inference of genetic exchange communities from microbial genomes using TF-IDF. *Frontiers in Microbiology* 8, 21.

Damitio, W. J., Tushingham, S., Brownstein, K. J., Matson, R. G., and Gang, D. R. (2021). The evolution of smoking and intoxicant plant use in ancient Northwestern North America. *American Antiquity*, 1–19.

Goodfellow, I., Bengio, Y., Courville, A., and Bengio, Y. (2016). *Deep Learning*. MIT Press, Cambridge, Massachusetts, USA.

King, A., Powis, T. G., Cheong, K. F., and Gaikwad, N. W. (2017). Cautionary tales on the identification of caffeinated beverages in North America. *Journal of Archaeological Science* 85, 30–40.

Pang, Z., Chong, J., Zhou, G., de Lima Morais, D. A., Chang, L., Barrette, M., Gauthier, C., Jacques, P. É., Li, S., and Xia, J. (2021). MetaboAnalyst 5.0: narrowing the gap between raw spectra and functional insights. *Nucleic Acids Research* 49, W388–W396

Pluskal, T., Castillo, S., Villar-Briones, A., and Orešič, M. (2010). MZmine 2: modular framework for processing, visualizing, and analyzing mass spectrometry-based molecular profile data. *BMC Bioinformatics* 11, 395.

Tautenhahn, R., Patti, G. J., Rinehart, D., and Siuzdak, G. (2012). XCMS Online: a web-based platform to process untargeted metabolomic data. *Analytical Chemistry* 84, 5035–5039.

Washburn, D. K., Washburn, W. N., Shipkova, P. A., and Pelleymounter, M. A. (2014). Chemical analysis of cacao residues in archaeological ceramics from North America: considerations of contamination, sample size and systematic controls. *Journal of Archaeological Science* 50, 191–207.

Young, T., Hazarika, D., Poria, S., and Cambria, E. (2018). Recent trends in deep learning based natural language processing. *IEEE Computational Intelligence Magazine* 13, 55–75.

225    Zimmermann, M., Brownstein, K. J., Díaz, L. P., Aragón, I. A., Hutson, S., Kidder, B.,
226    Tushingham, S., and Gang, D. R. (2021). Metabolomics-based analysis of miniature flask
227    contents identifies tobacco mixture use among the ancient Maya. *Scientific Reports* 11, 1–11.