

# Deep Spatial Learning for Forensic Geolocation with Microbiome Data

**Neal S. Grantham**  
**North Carolina State University**

**JSM Chicago**  
**August 3<sup>rd</sup>, 2016**

Joint work with

 Brian Reich (NCSU Statistics)

 Eric Laber (NCSU Statistics)

In collaboration with

 Rob Dunn (NCSU Biology)

# What is a microbiome?

Community of microbial organisms occupying an ecological niche.

Next-generation sequencing technologies make possible efficient identification of these microbes at affordable cost<sup>1</sup>.

Huge interest in understanding microbiomes as they relate to human health, diet, agriculture, environment, forensics, etc.

---

<sup>1</sup> Metzker. (2010) [Sequencing technologies—the next generation](#). Nature Reviews Genetics.

# The White House Launches the National Microbiome Initiative<sup>2</sup>

Half a billion dollars pledged,  
with three major goals:

- ☞ collaboration,
- ☞ developing better tools for  
studying microbiomes, and
- ☞ recruitment.

---

<sup>2</sup> [The Atlantic article](#) by Ed Yong, photo by Jim Young / Reuters



# Statistical challenges

Microbiome data are...

- ☞ high-dimensional,
- ☞ sparse,
- ☞ over-dispersed, and
- ☞ possess complex dependence structure.

Many exciting opportunities for research with microbiologists.

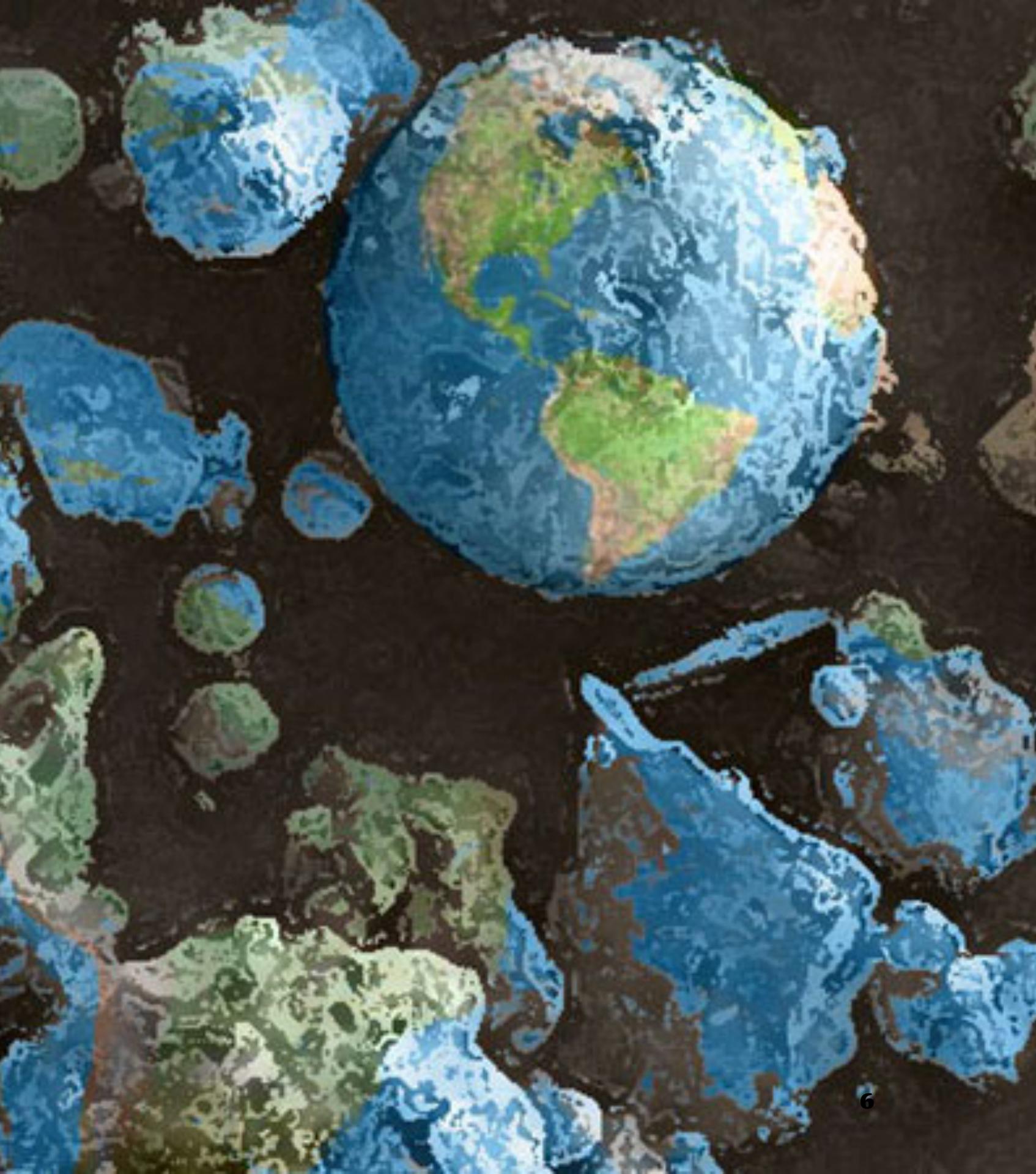
# Motivating dataset

[Wild Life of Our Homes](#), a public science project by Rob Dunn Lab.

- ☞ Dust samples collected from outer door frames of  $n \approx 1,000$  homes across the continental U.S.
- ☞ DNA sequencing revealed  $p \approx 50,000$  fungi species<sup>3</sup>.

---

<sup>3</sup> Well, operational taxonomic units (OTUs) to be precise.



Research question:

Does the microbiome composition of an ambient dust sample inform its geographic origin?

# Our approach:

Build a model to estimate the unknown origin  $\mathbf{s}$  of a dust sample conditional on its known microbiome composition  $\mathbf{x}$ .

# Initial model and direction

Naive Bayes discriminant analysis estimates origin of dust samples with a median prediction error of 230 kilometers<sup>4</sup>.

Here, we develop a new model based on

- ① spatial point pattern theory, and
- ② deep learning.

---

<sup>4</sup> Grantham et al. (2015) [Fungi identify the geographic origin of dust samples](#). PLOS One.

# 1. Spatial point pattern theory

Assume a spatial point pattern<sup>5</sup> intensity surface  $\lambda(s \mid x)$  has non-homogeneous Poisson process likelihood

$$\mathcal{L} [\lambda(s \mid x); \{(s_i, x_i)\}_{i=1}^n] = \exp \left[ - \int_{\mathcal{D}} \lambda(s \mid x) ds \right] \prod_{i=1}^n \lambda(s_i \mid x_i).$$

May select a parametric model for  $\lambda(s \mid x)$ .

---

<sup>5</sup> Gelfand et al. (2010) [Handbook of Spatial Statistics](#). CRC Press.

# 1. Spatial point pattern theory

Liang et al.<sup>6</sup> propose a log Gaussian process (GP),

$$\lambda(s \mid x) = \pi(s) \exp[\beta' x + w(s)]$$

with  $w(s) \sim GP$ ,  $\pi(s)$  a population offset, and  $\beta$  unknown.

However, no closed-form solution to integral in  $\mathcal{L}$ .

---

<sup>6</sup> Liang et al. (2008) [Analysis of Minnesota colon and rectum cancer point patterns with spatial and nonspatial covariate information](#).  
The Annals of Applied Statistics.

# 1. Spatial point pattern theory

Numerical approximation possible with a Monte Carlo algorithm using a knot-based predictive process<sup>6</sup>.

Effective, but difficult to implement in practice:

- ☞ Requires careful knot construction.
- ☞ Performance suffers for large  $n$ , larger  $p$ .

---

<sup>6</sup> Liang et al. (2008) [Analysis of Minnesota colon and rectum cancer point patterns with spatial and nonspatial covariate information](#). The Annals of Applied Statistics.

## 2. Deep learning

Well-suited for high-dimensional data with complex structure<sup>7</sup>.

Rich, GPU-enabled software libraries available:

- ☞ Theano (Python)
- ☞ Torch7 (Lua)
- ☞ Tensorflow (C++)

---

<sup>7</sup> LeCun et al. (2015) [Deep learning](#). Nature.

## 2. Deep learning

Let  $\mathcal{P} = \{P_k\}_{k=1}^K$  denote a partition of spatial domain  $\mathcal{D}$ .

Represent every  $s$  by the region to which it belongs.

Two major benefits:

- ☞ Avoids costly approximation of integral in  $\mathcal{L}$ .
- ☞ Reframes estimation as a supervised classification problem.

## 2. Deep learning

For regions  $k = 1, \dots, K$ ,

$$Pr(s \in P_k \mid x) = \frac{\exp[f_k(x)]}{\sum_{l=1}^K \exp[f_l(x)]} \text{ where } f_k : \mathcal{X} \rightarrow \mathbb{R}.$$

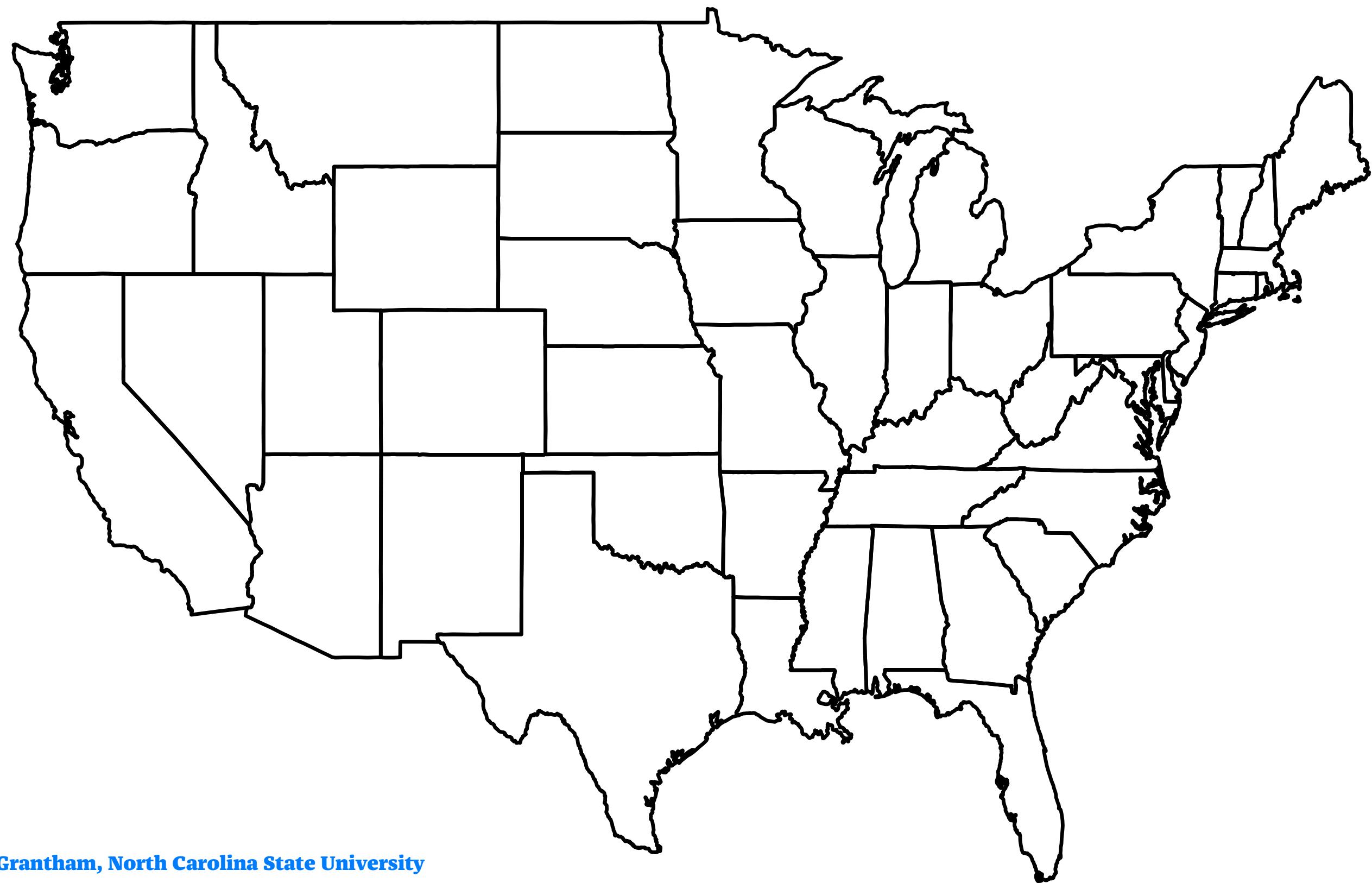
Estimate  $f_1(\cdot), \dots, f_K(\cdot)$  by training a deep neural network (DNN) on  $\{(s_i, x_i)\}_{i=1}^n$  with categorical cross-entropy cost function

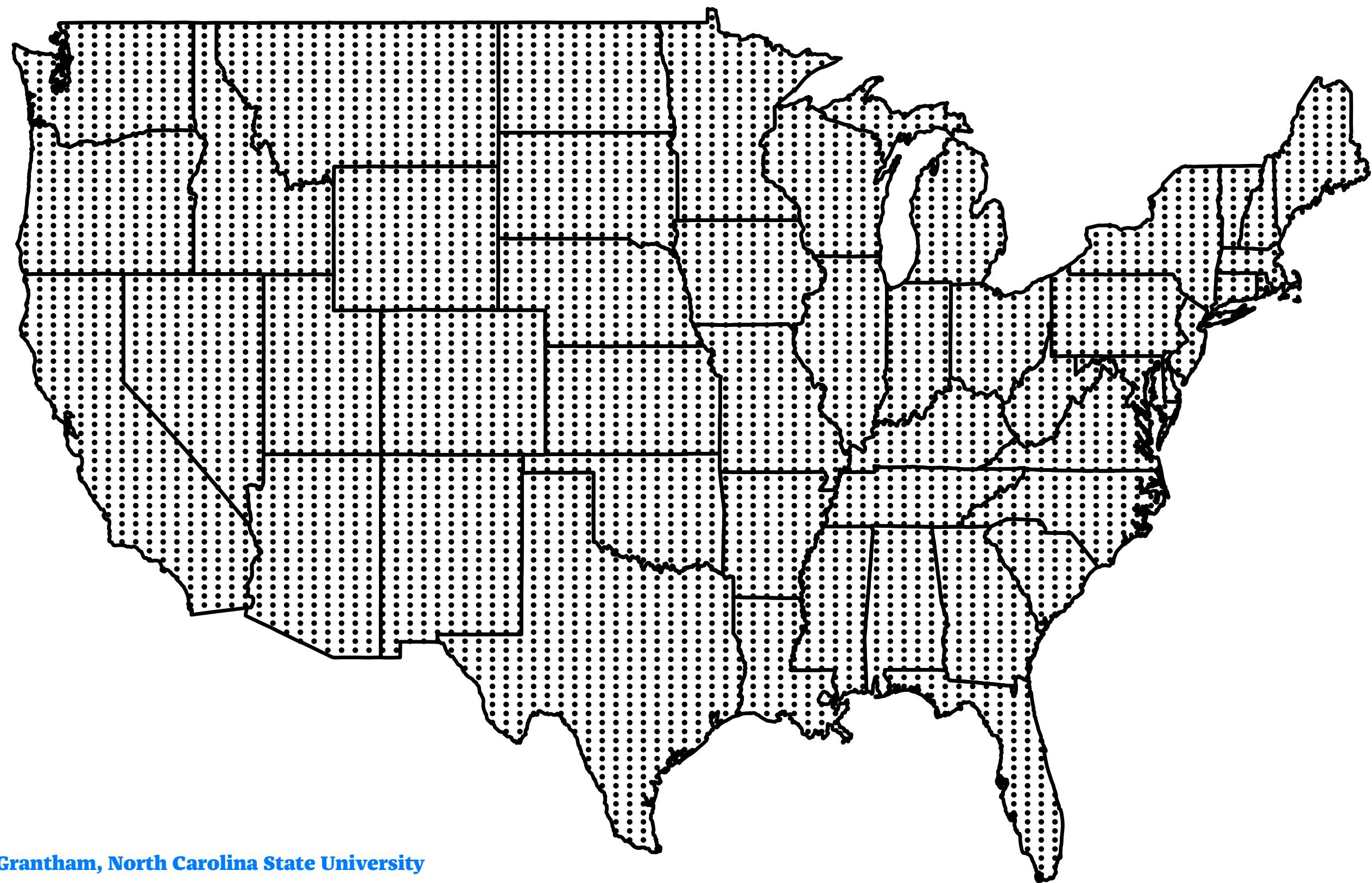
$$C = - \sum_{i=1}^n \sum_{k=1}^K \log[Pr(s_i \in P_k \mid x_i)] I(s_i \in P_k).$$

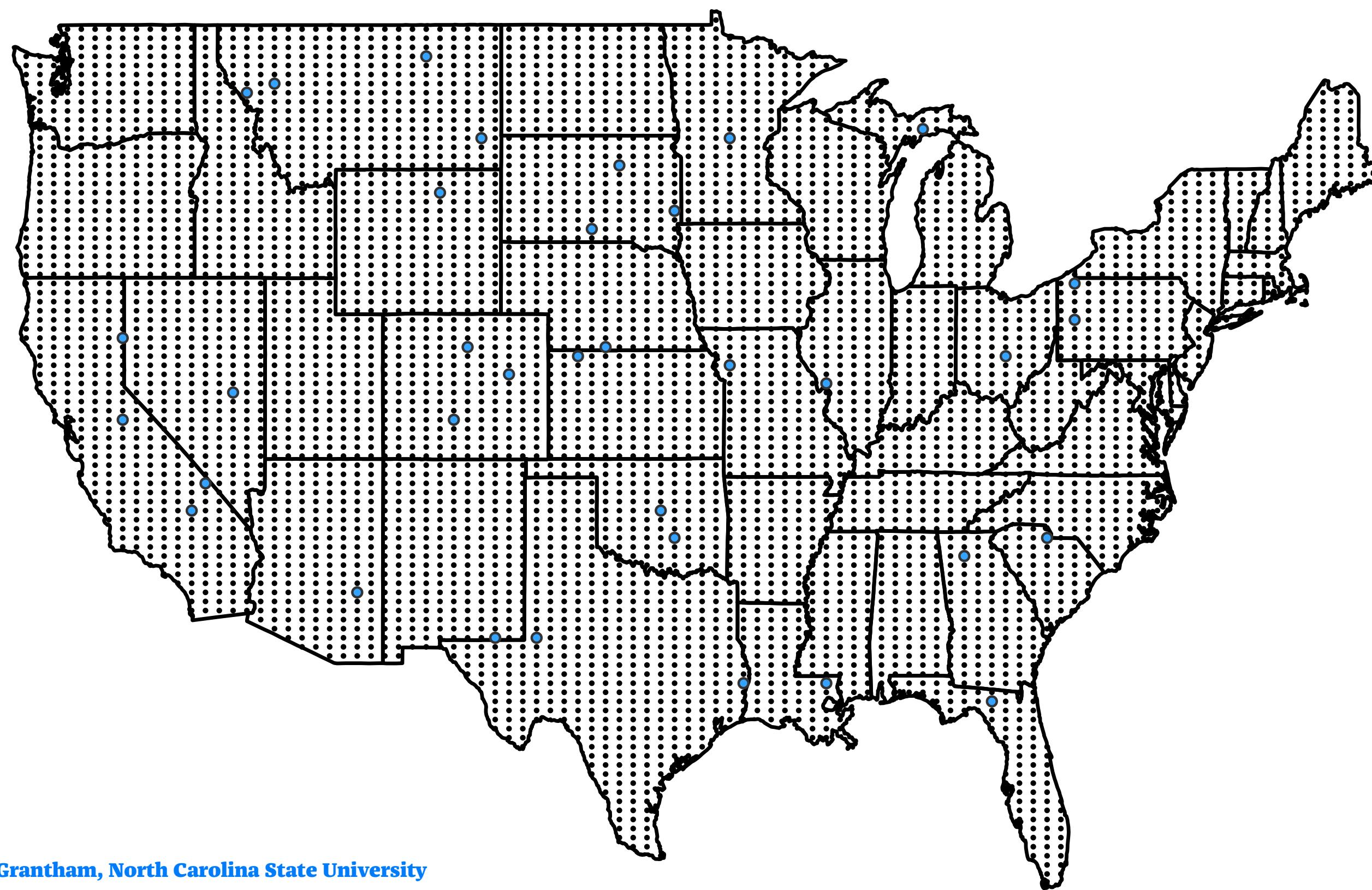
# Our "Deep Space" algorithm

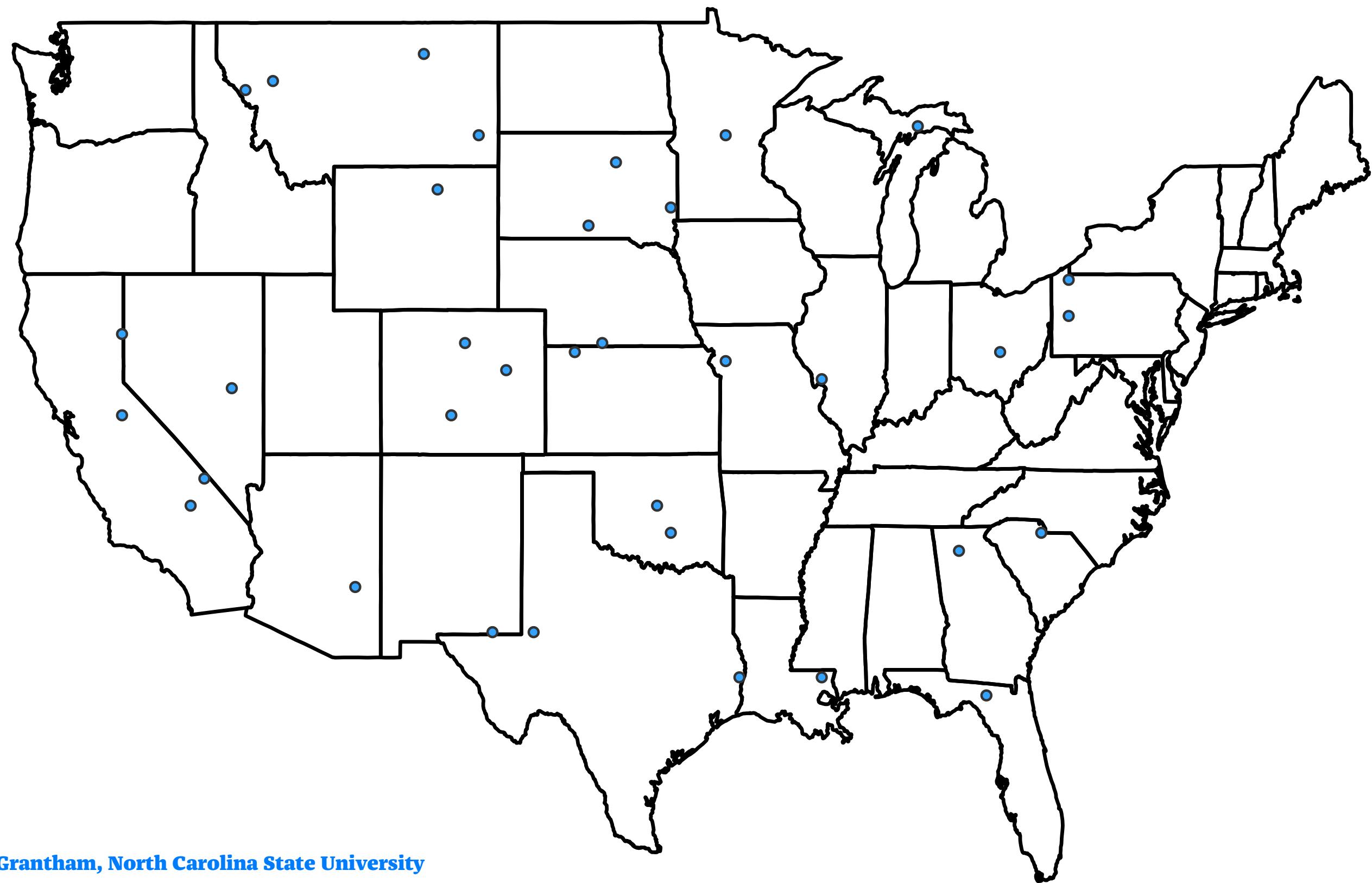
- ① Generate random Voronoi partition  $\mathcal{P}$  over  $\mathcal{D}$ .
- ② Train DNN on available data from these regions.
- ③ Repeat steps 1 & 2  $N$  times to develop a diverse collection,  $\mathcal{M}$ , of trained DNNs.
- ④ Predict most likely origin  $\hat{s}$  by averaging over DNNs in  $\mathcal{M}$ .

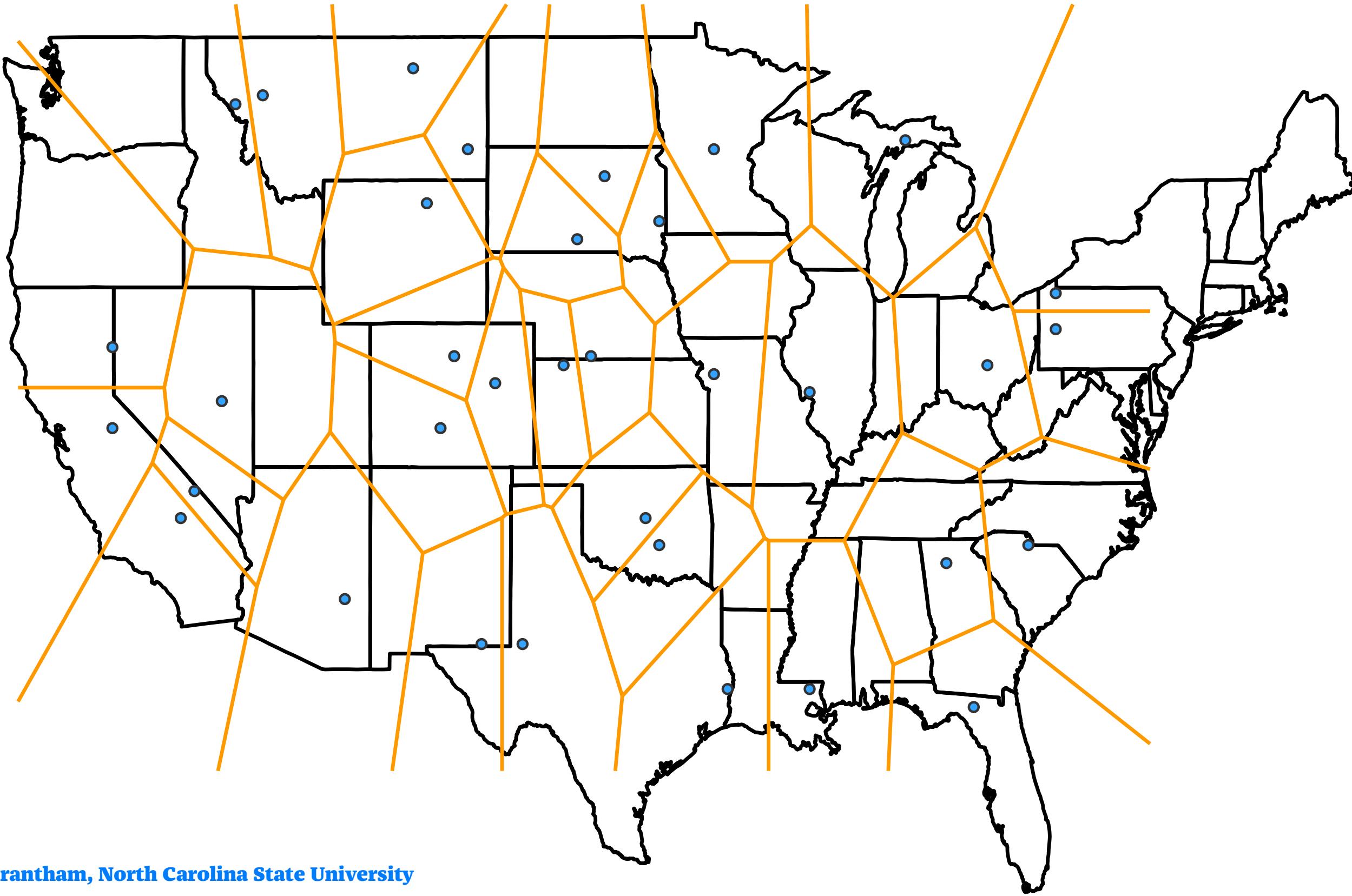
1. Generate random Voronoi partition  $\mathcal{P}$  over  $\mathcal{D}$ .

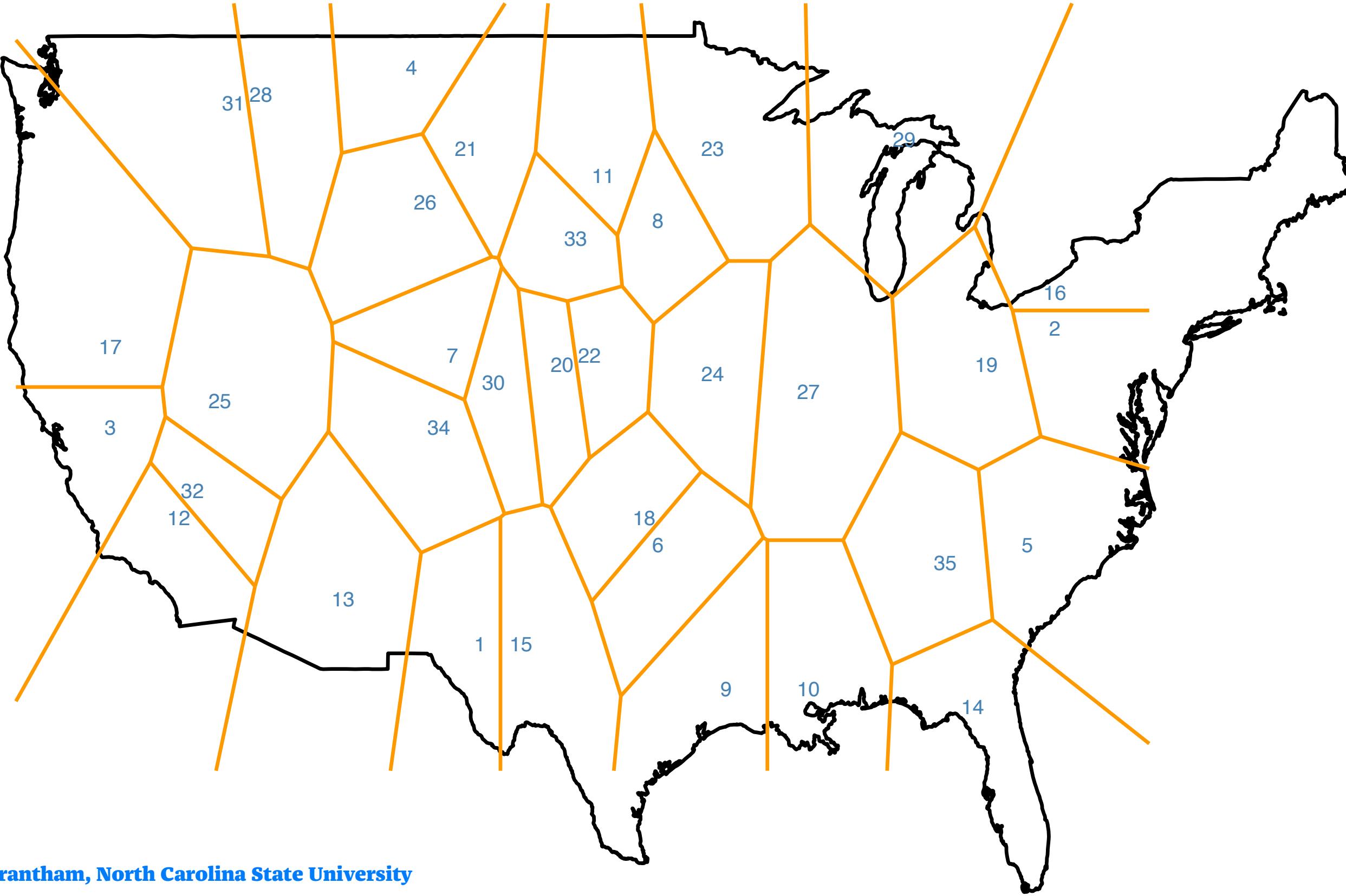




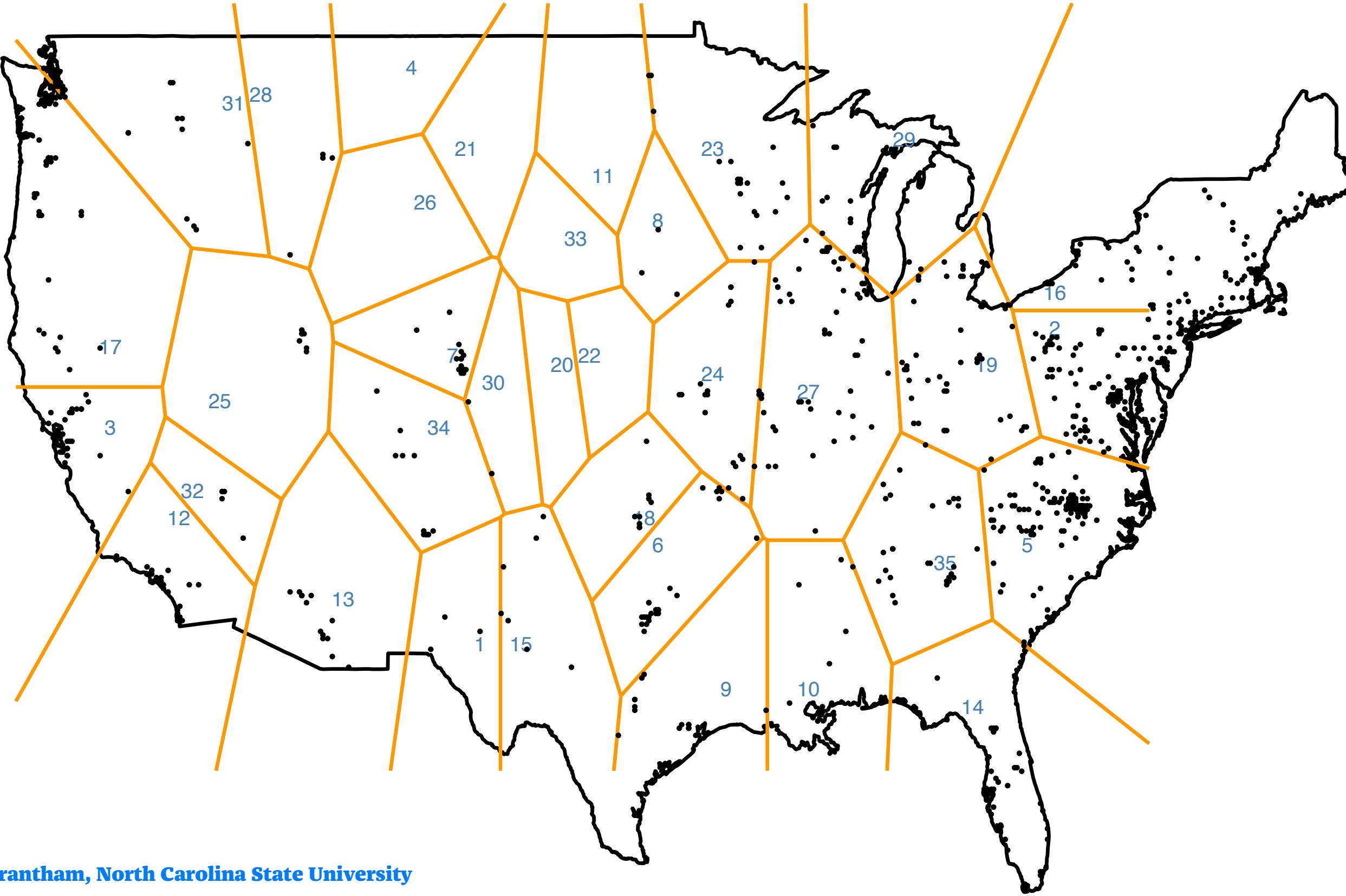


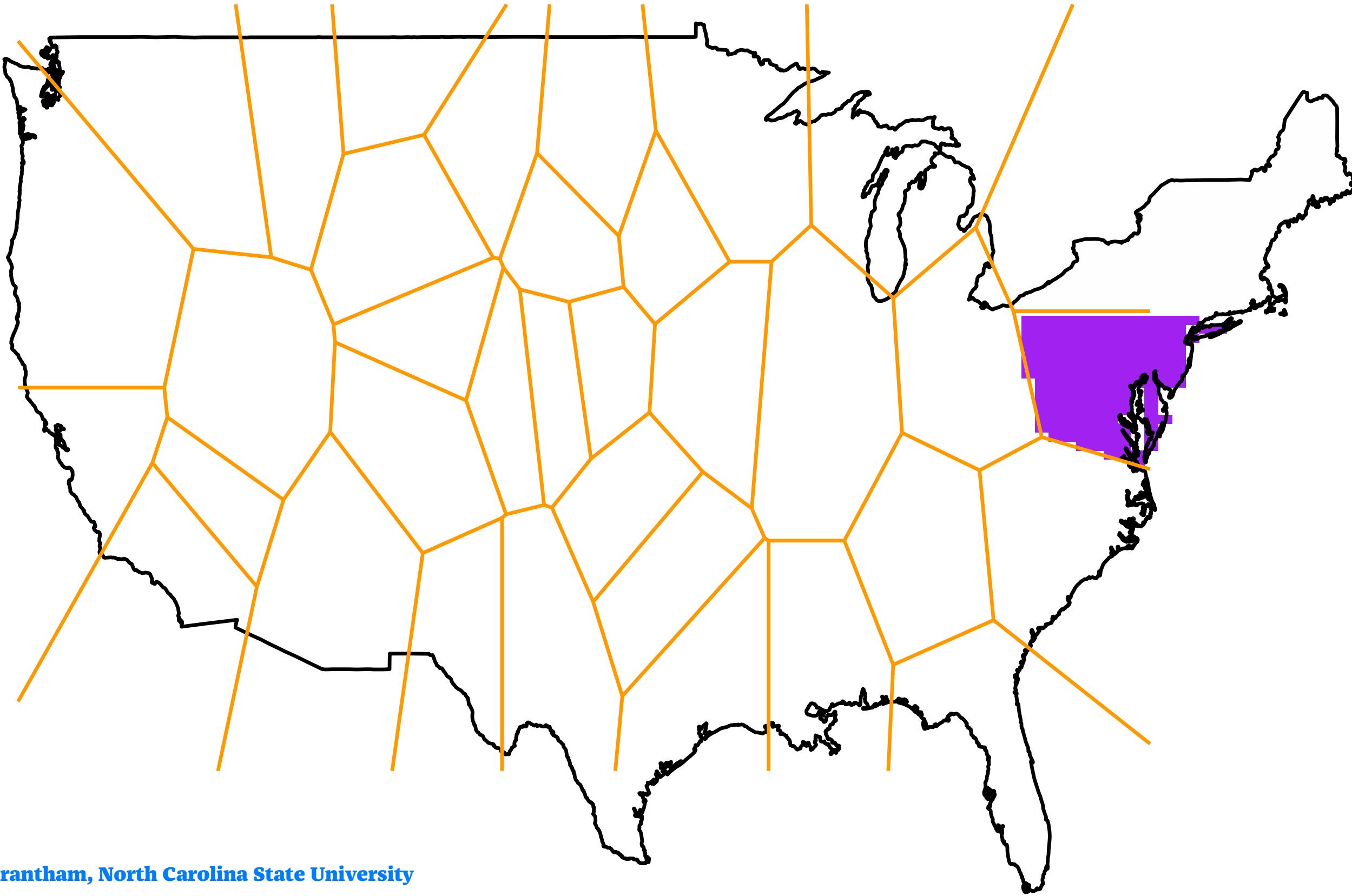




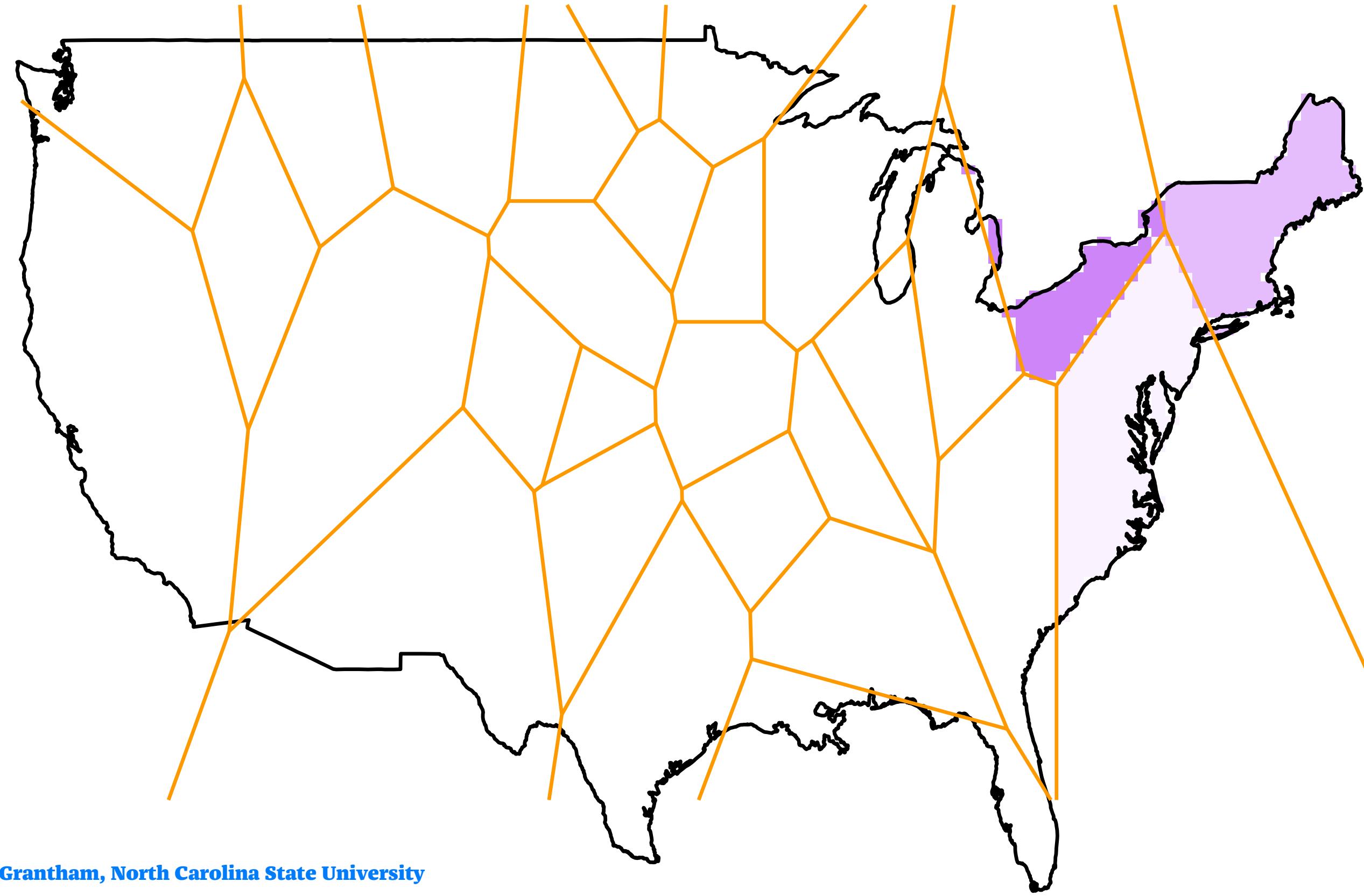


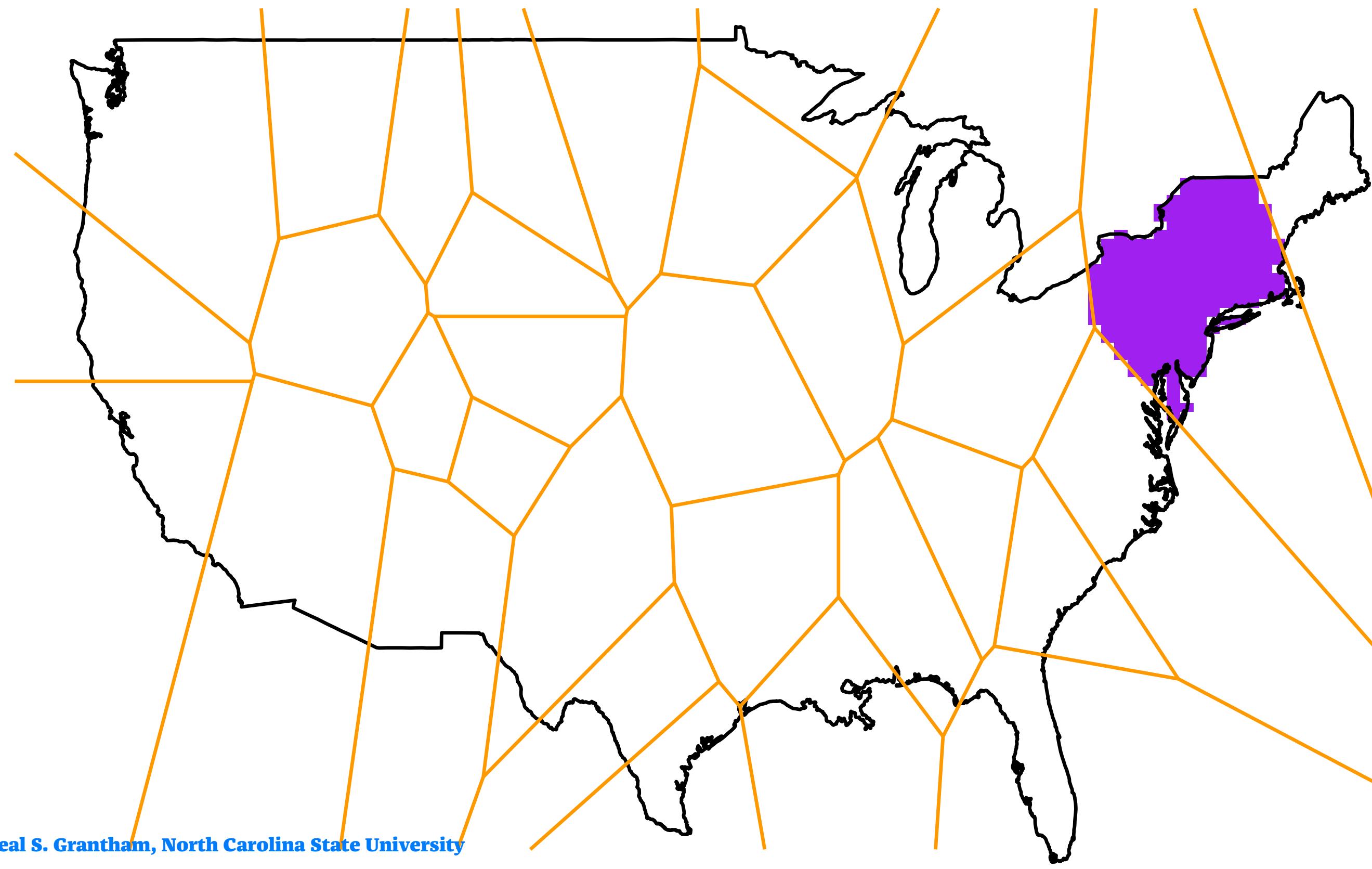
2. Train DNN on available data  
from these regions.

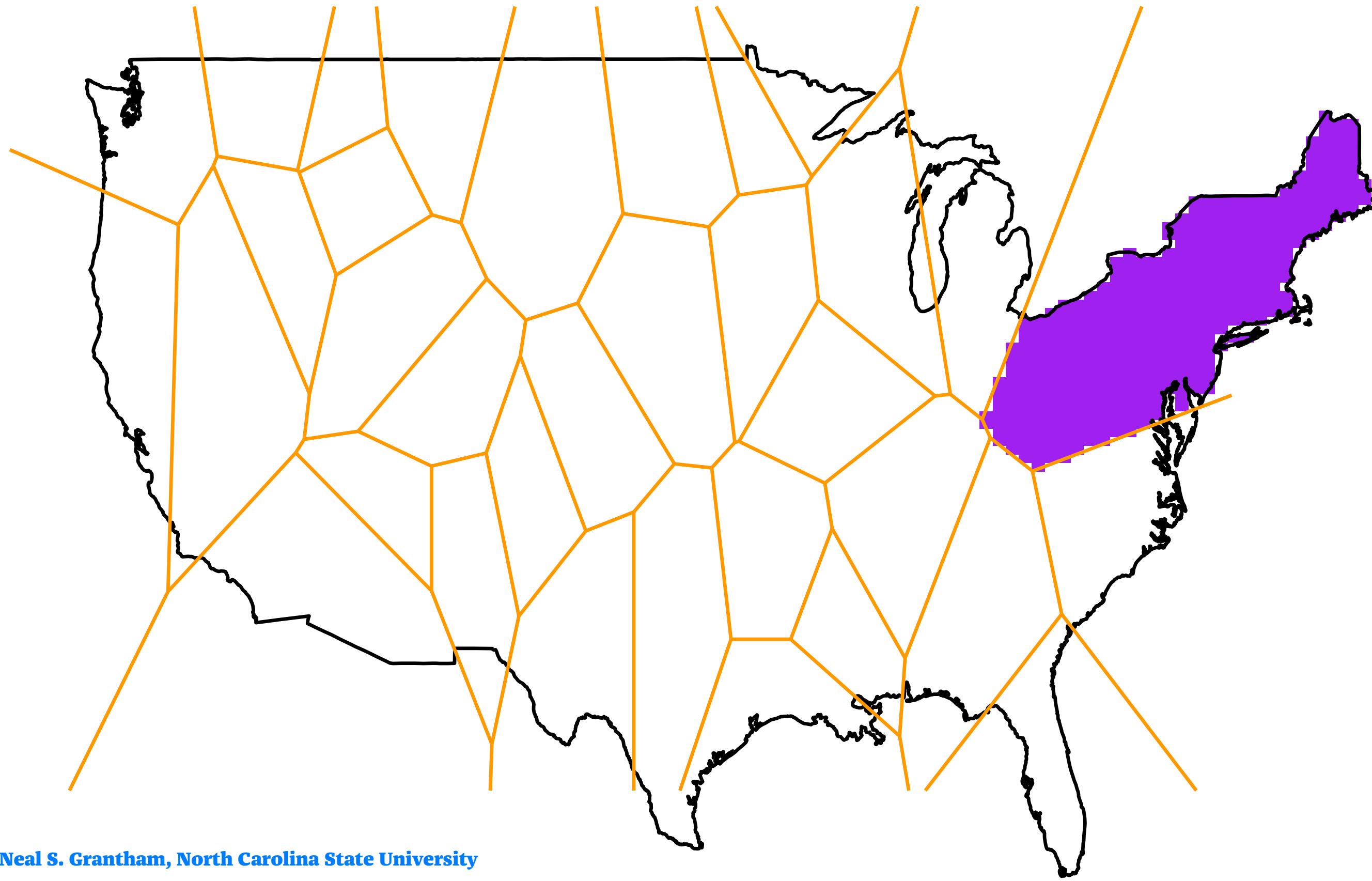


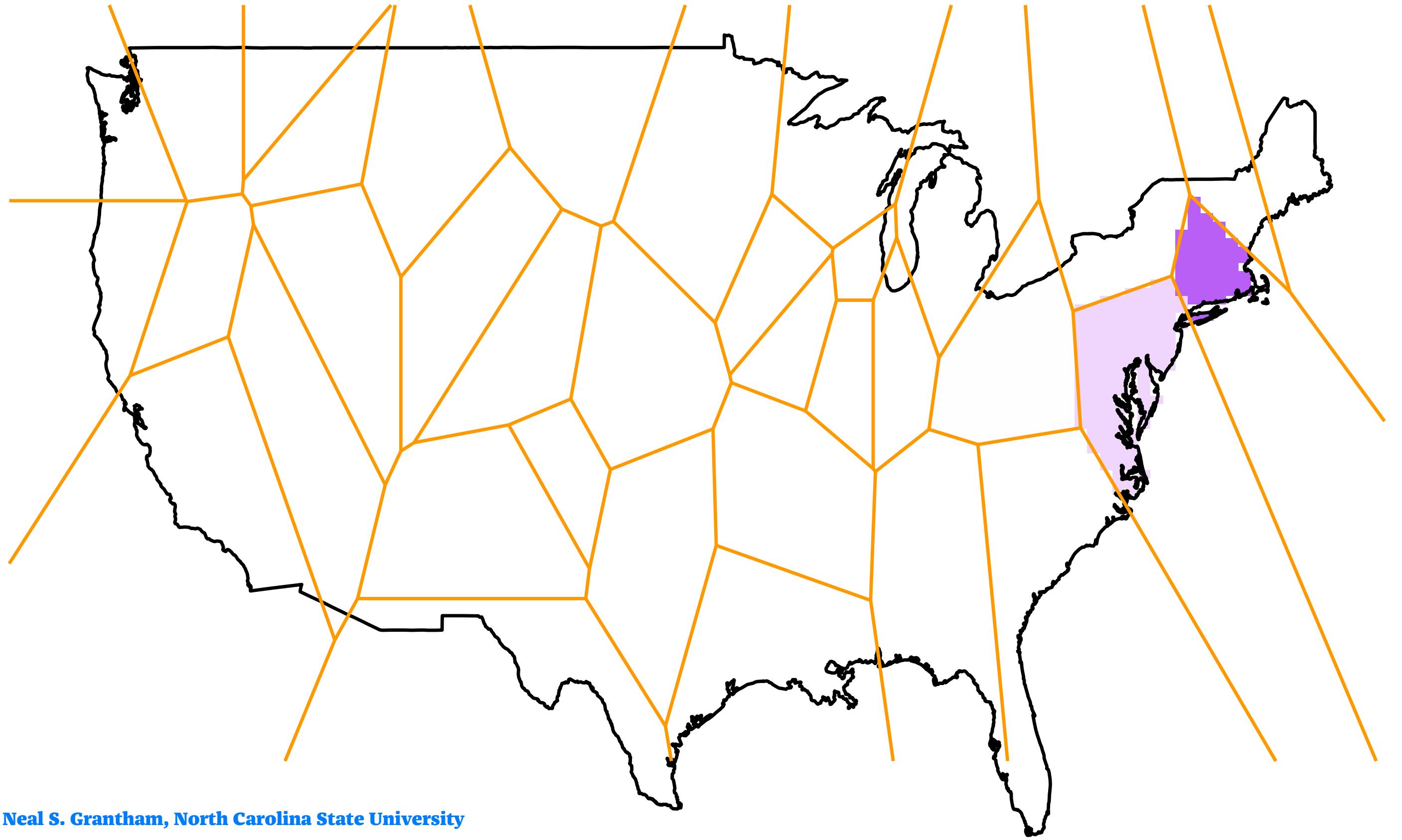


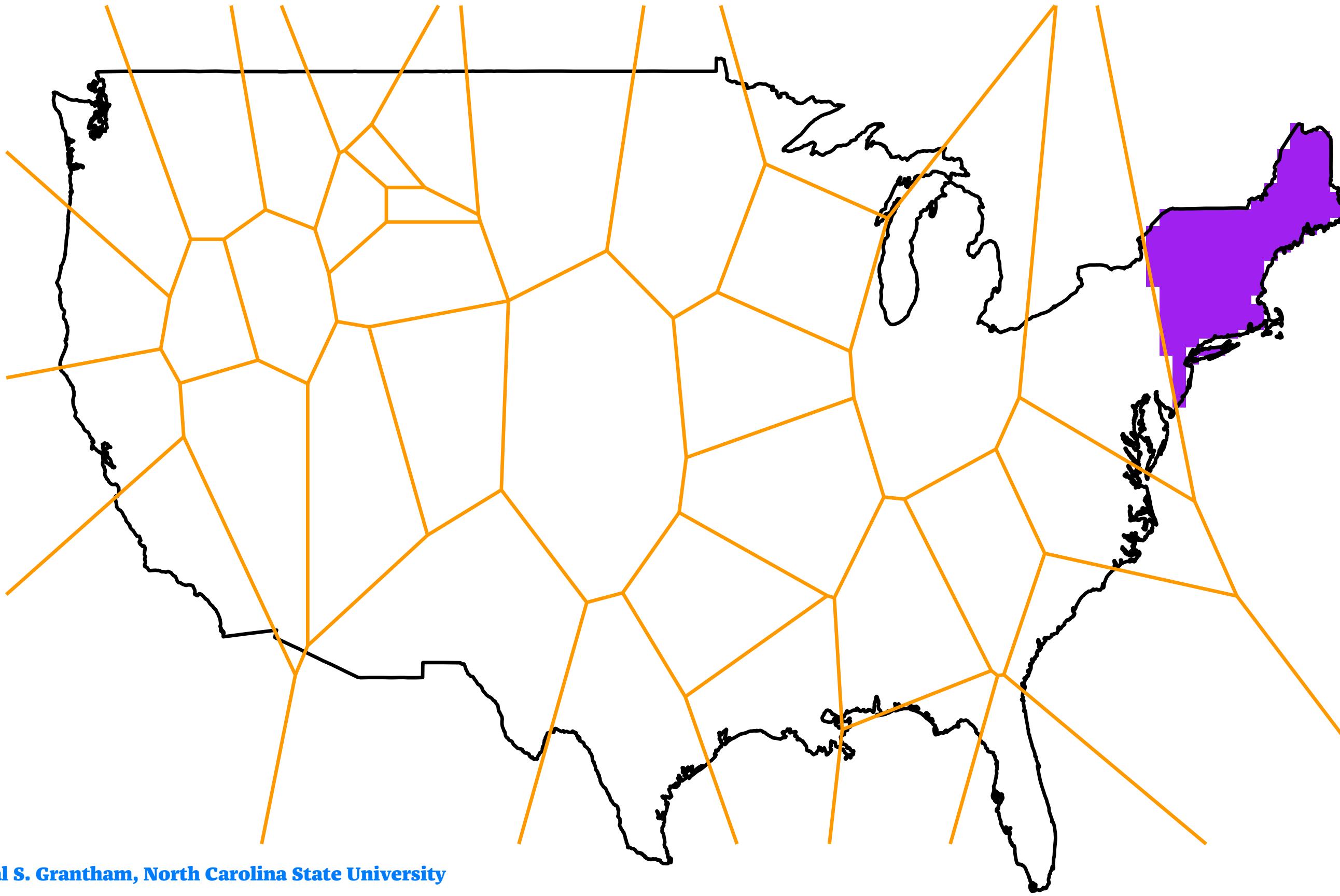
3. Repeat steps 1 & 2  $N$  times to develop a diverse collection,  $\mathcal{M}$ , of trained DNNs.

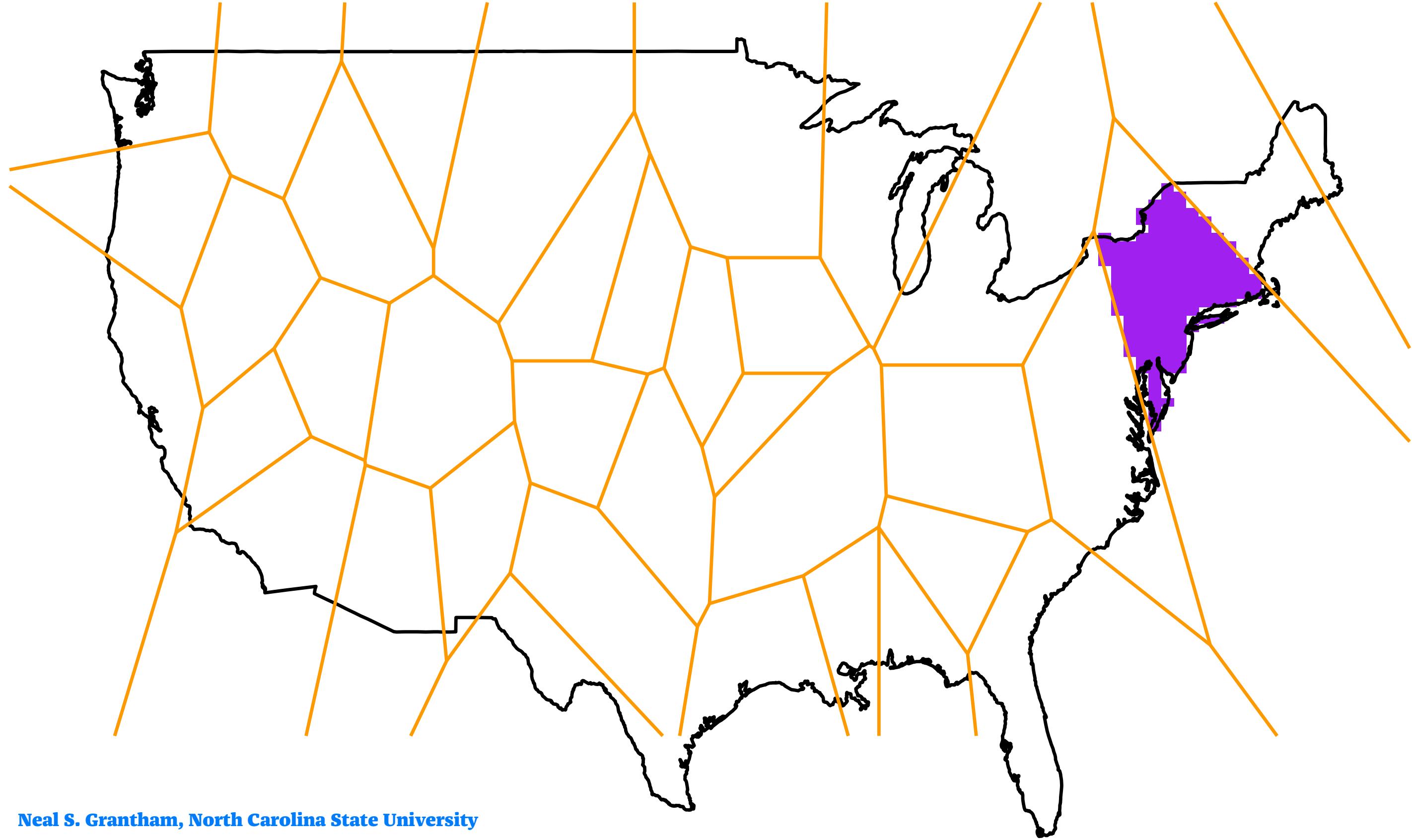


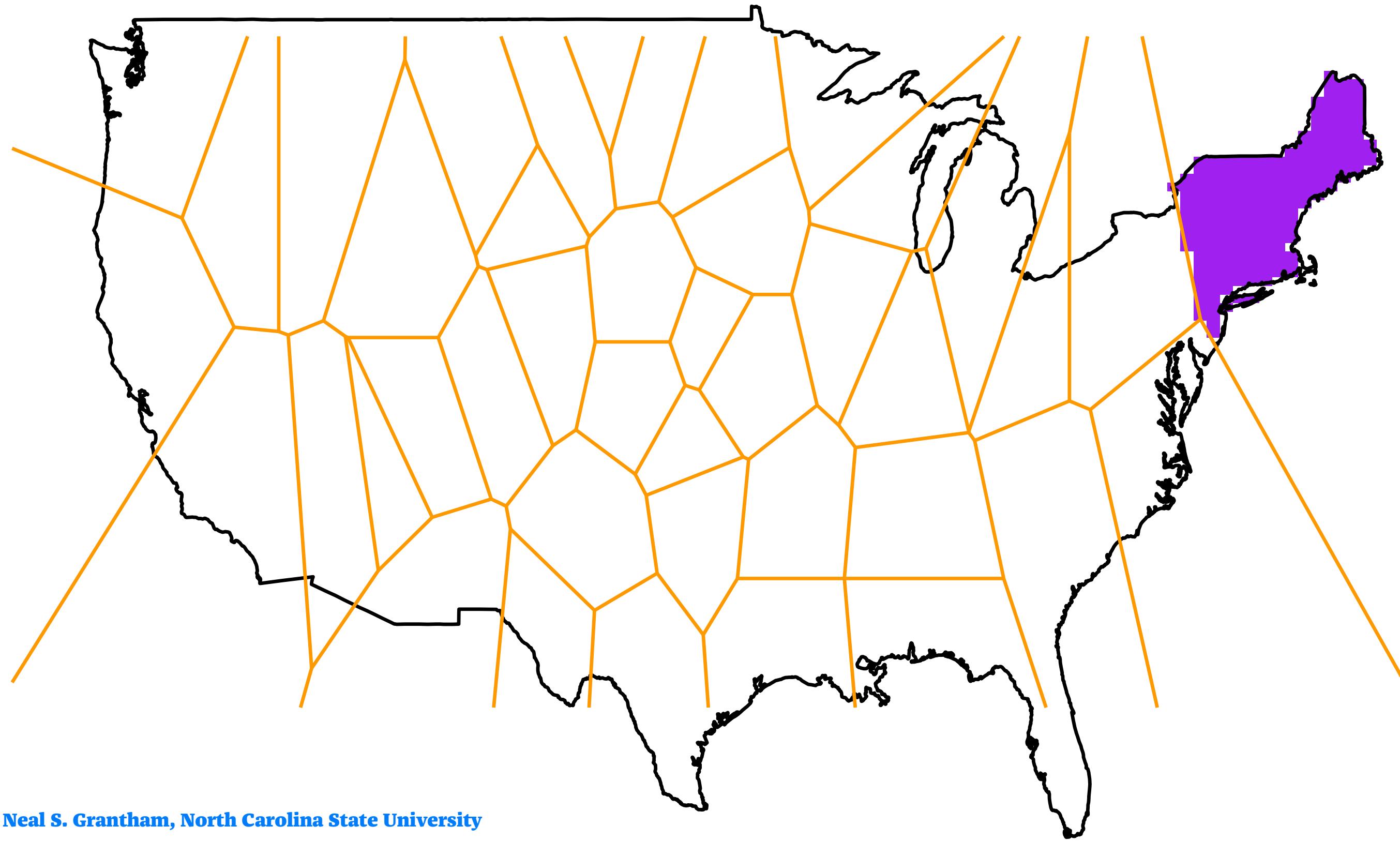


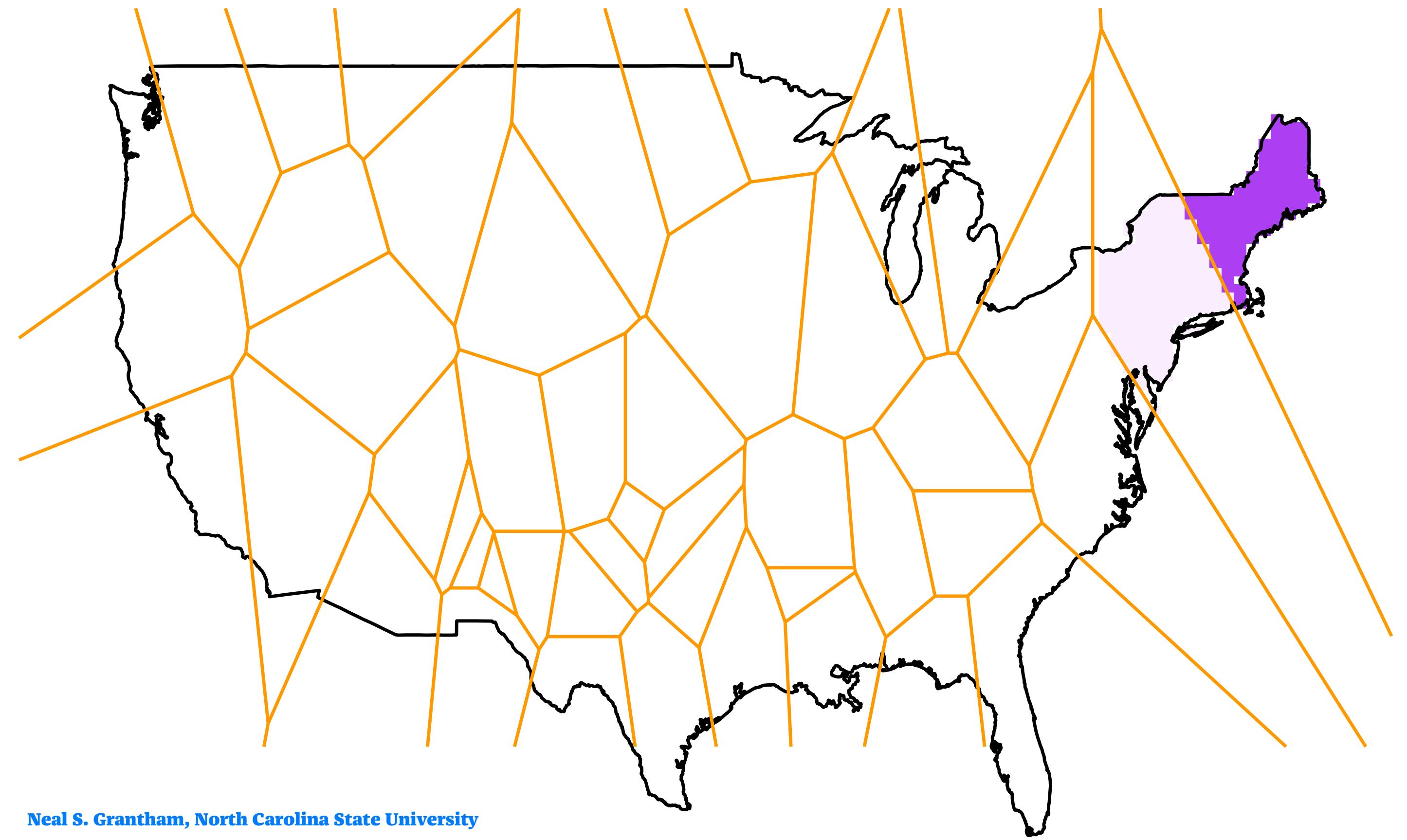


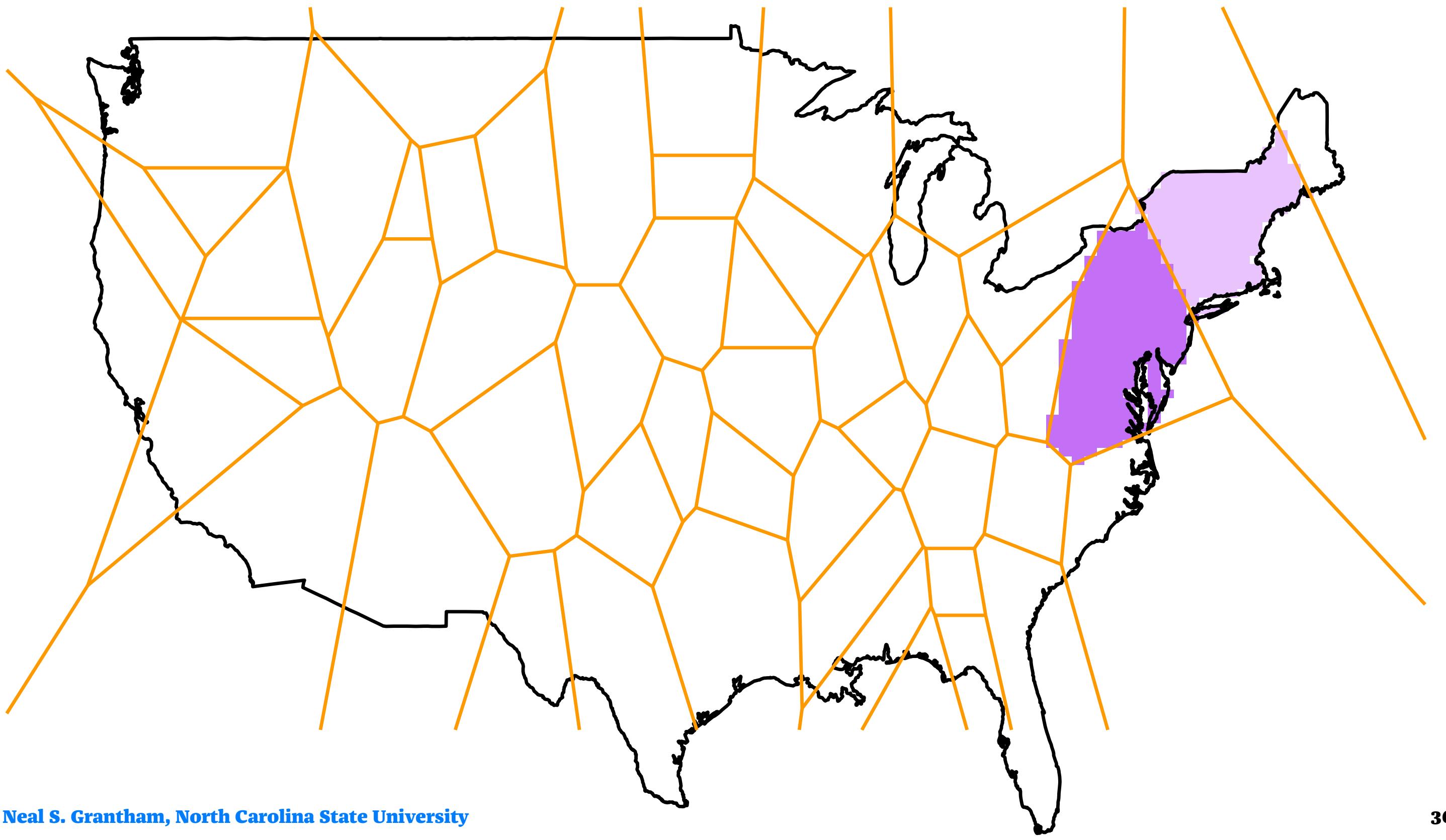












4. Predict most likely origin  $\hat{s}$  by averaging over DNNs in  $\mathcal{M}$ .

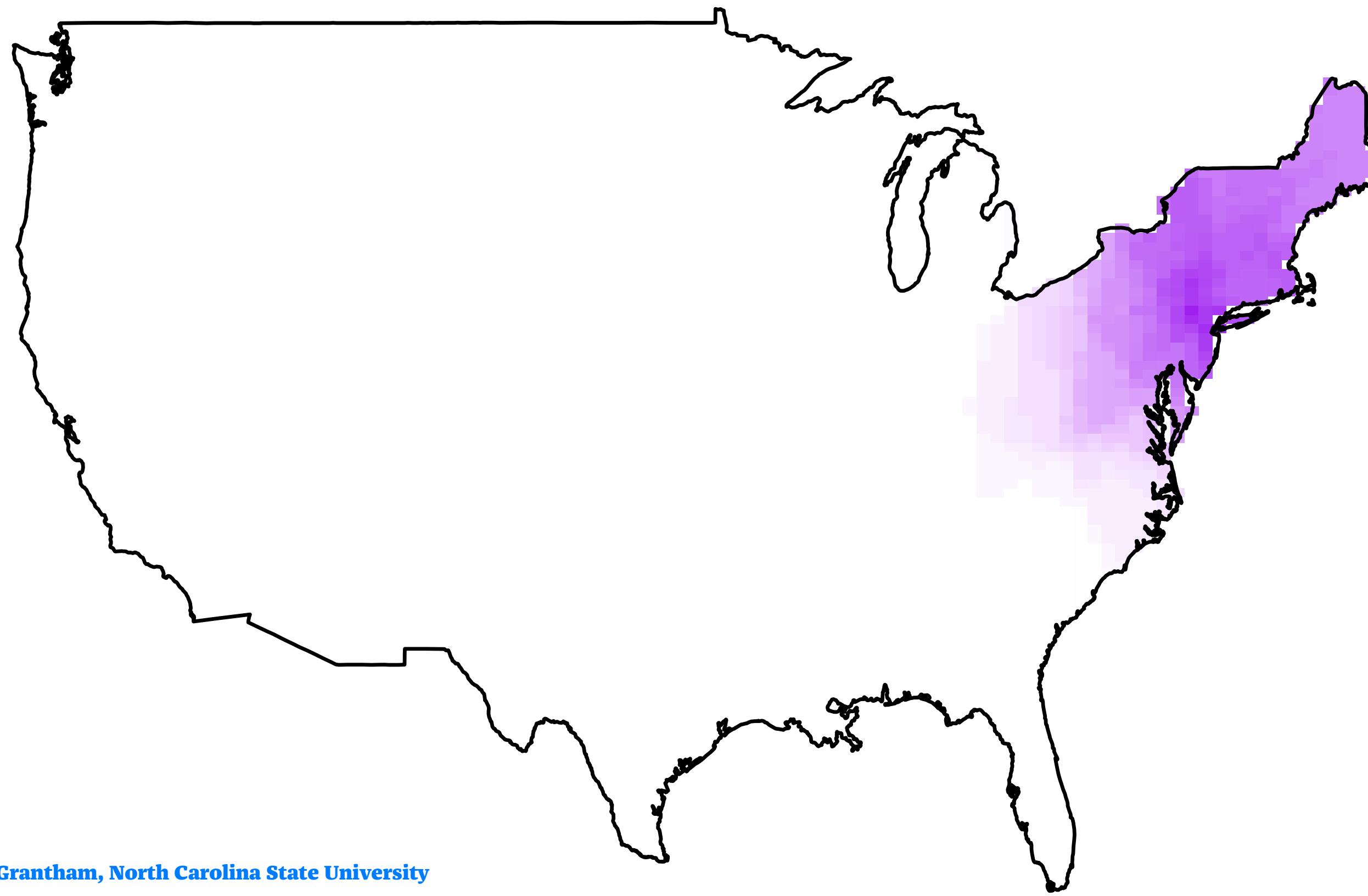
# Geolocation

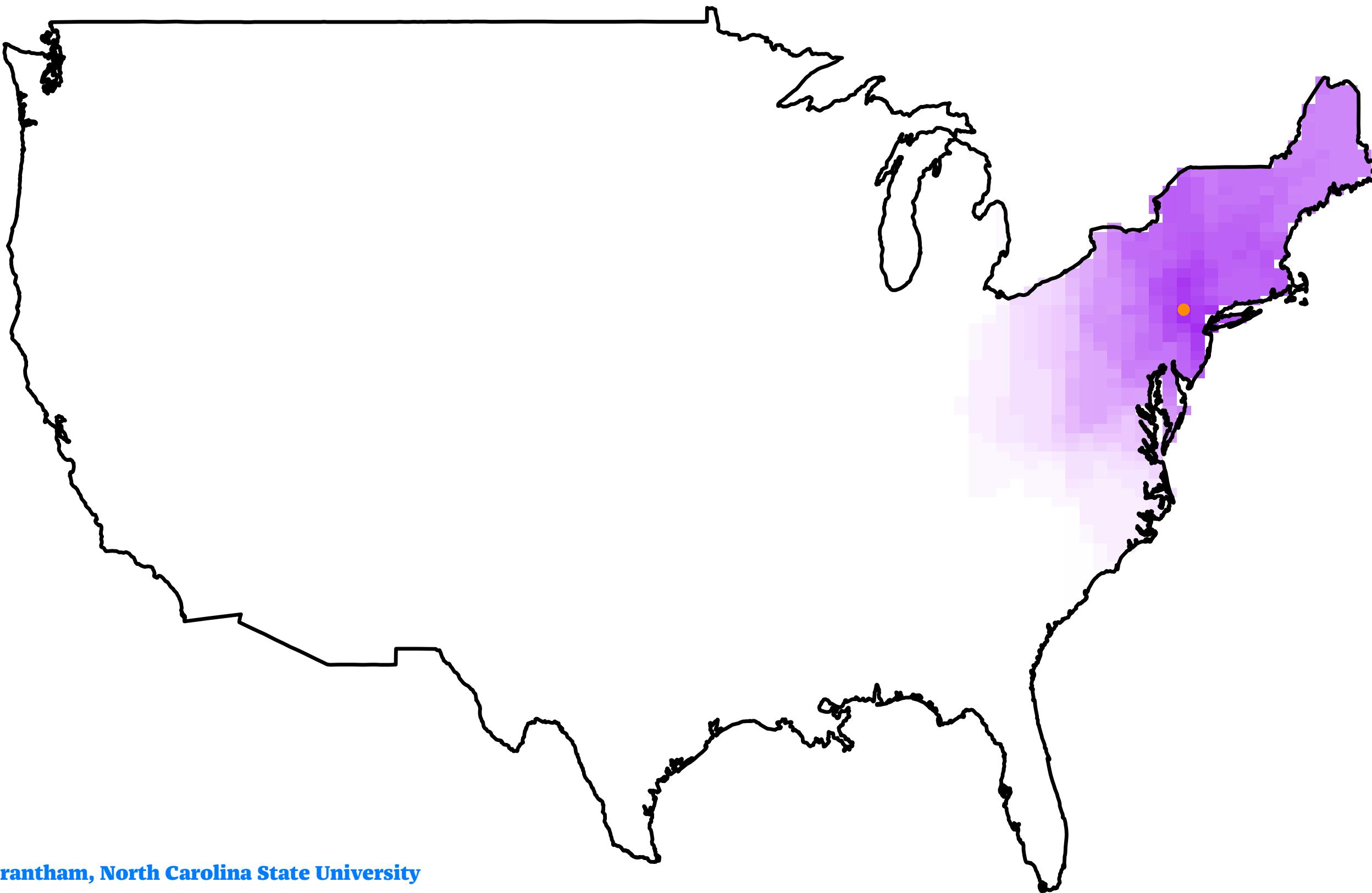
A sample with microbiome  $\mathbf{x}$  is most likely to have originated from

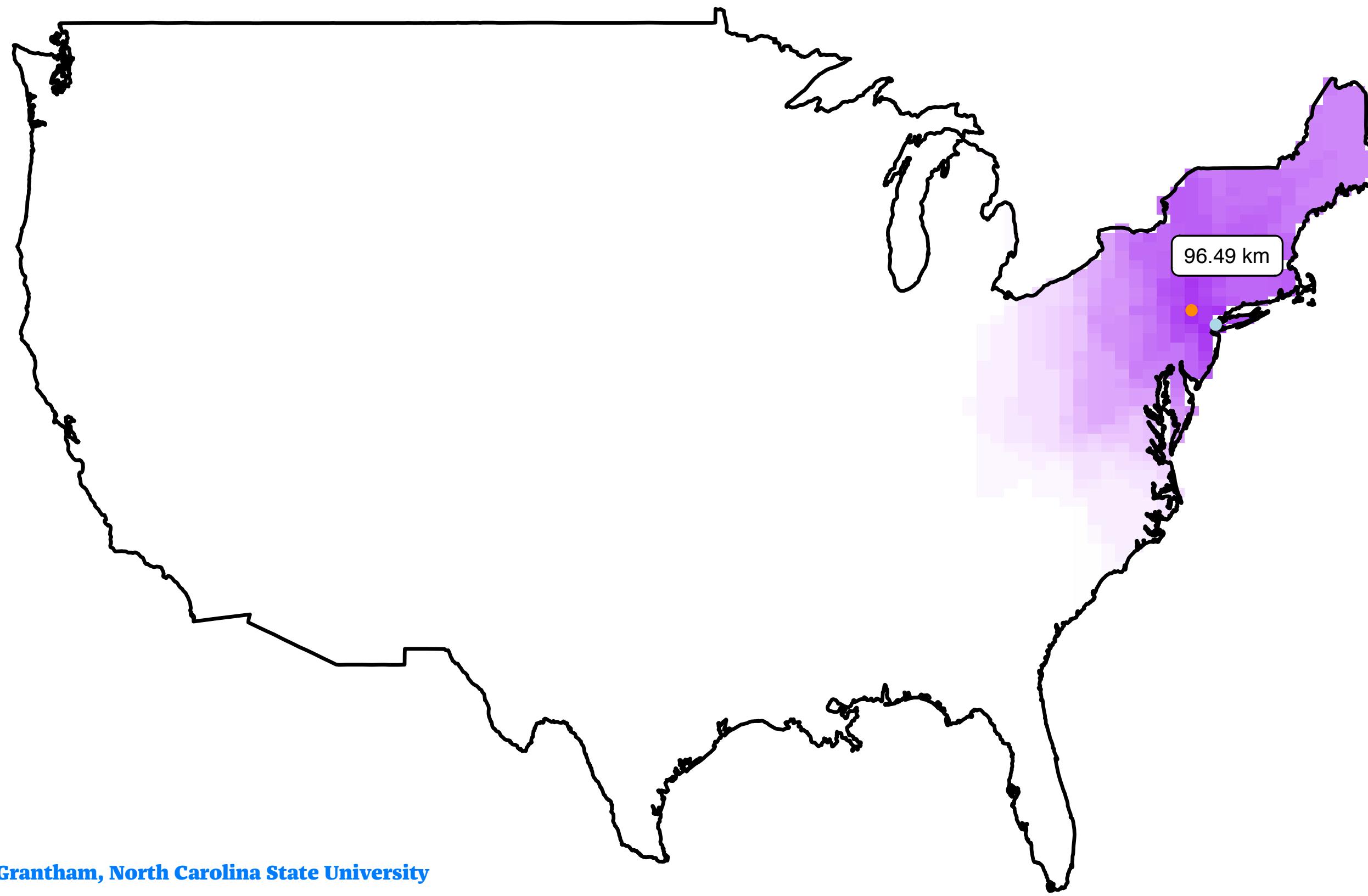
$$\hat{s} = \arg \max_{s \in \mathcal{D}} g(s \mid \mathbf{x}, \mathcal{M})$$

where  $g(\cdot)$  is the geolocation function given by

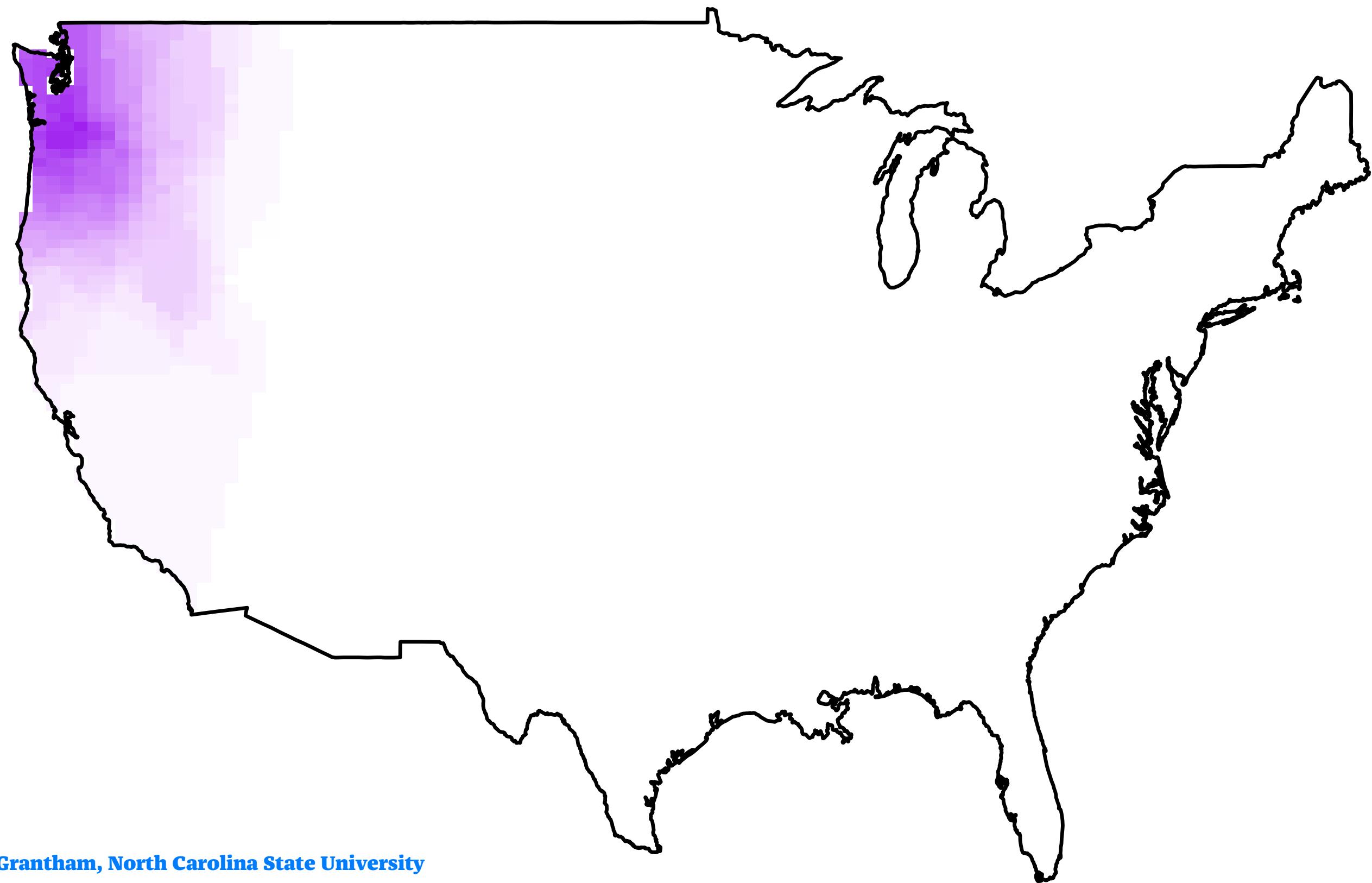
$$g(s \mid \mathbf{x}, \mathcal{M}) = \frac{1}{N} \sum_{j=1}^N \sum_{k=1}^{K_j} \frac{1}{|P_{jk}|} Pr(s \in P_{jk} \mid \mathbf{x}) I(s \in P_{jk}).$$

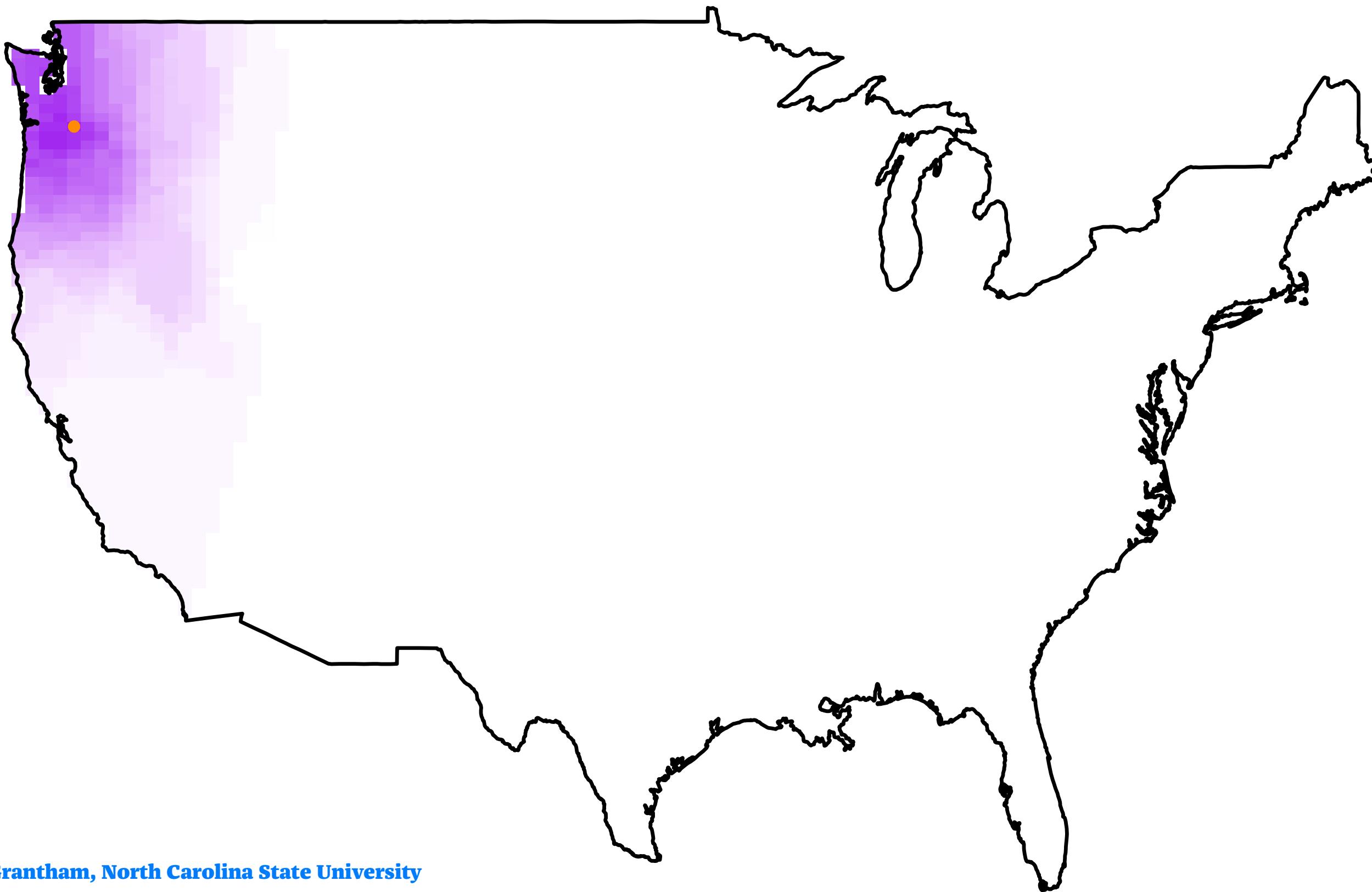


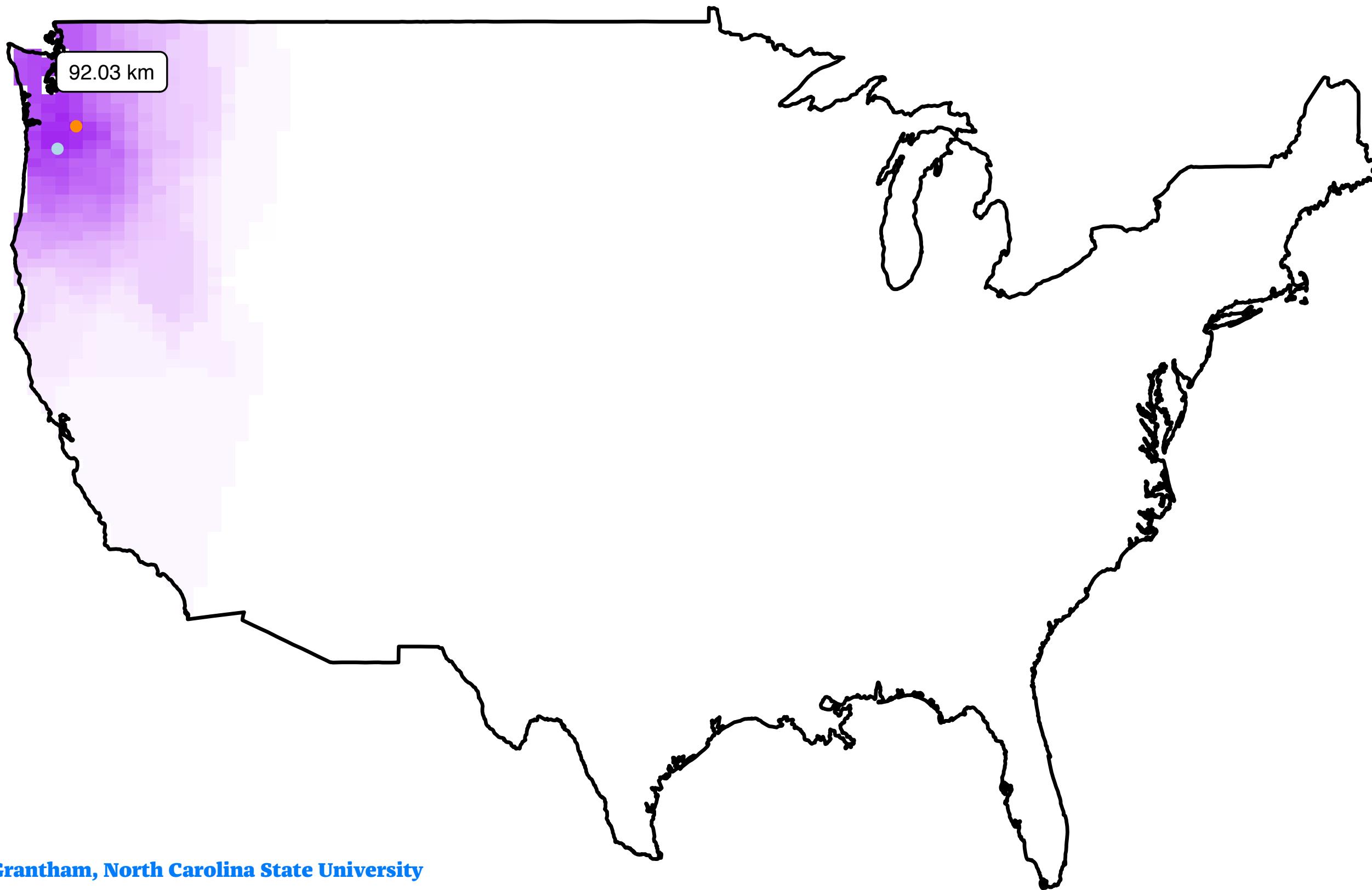


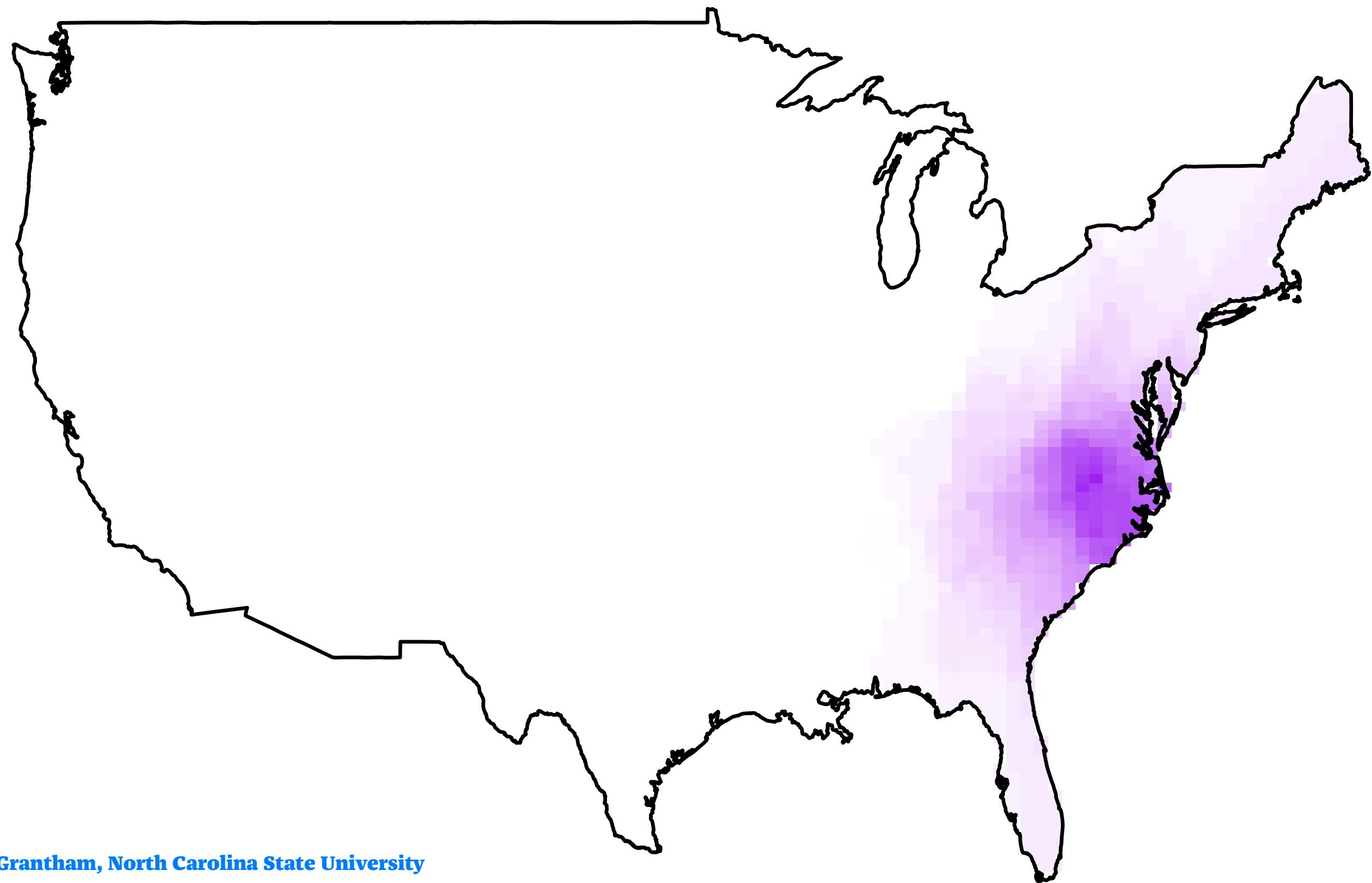


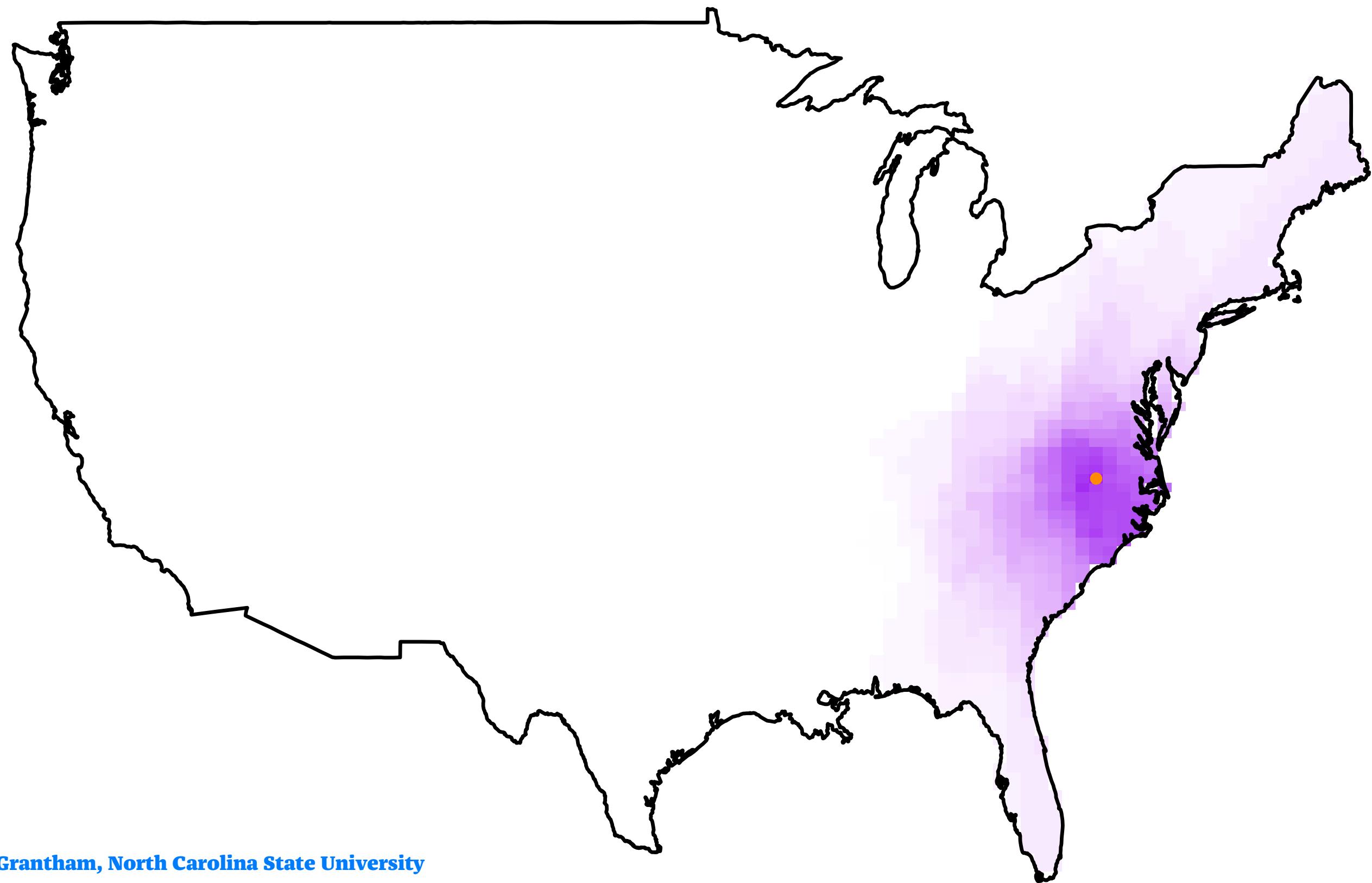
# More test sample predictions...

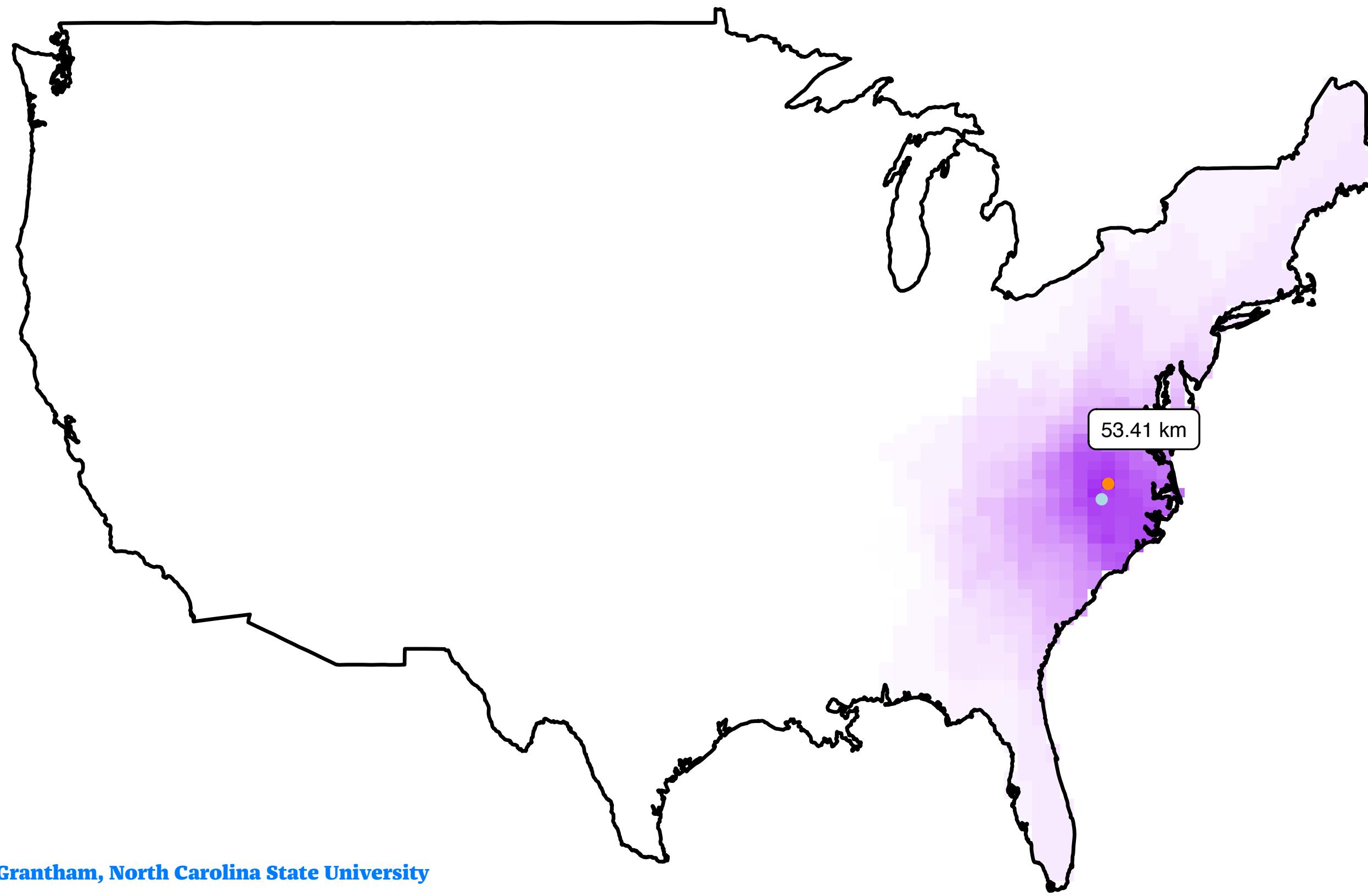


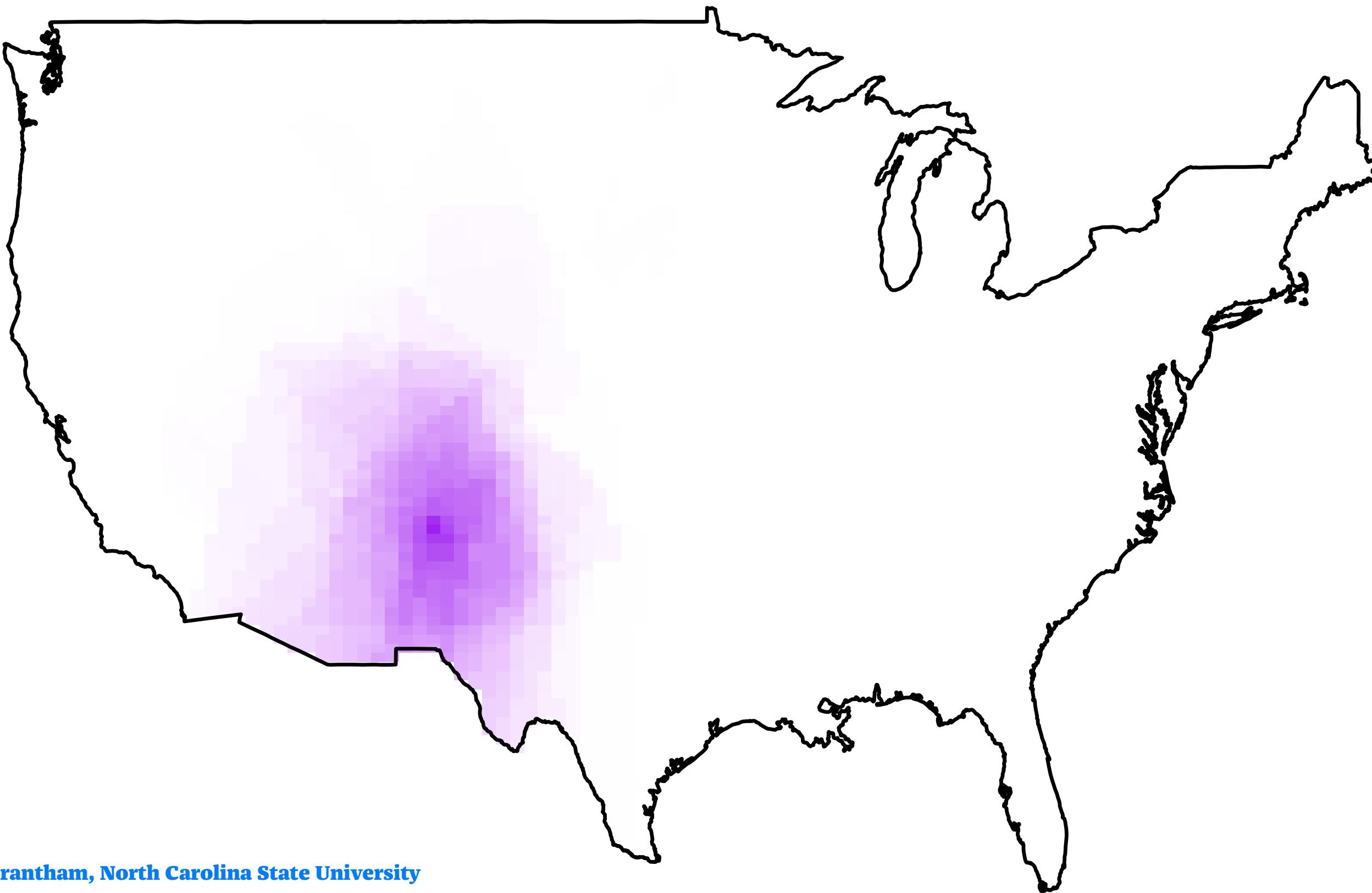


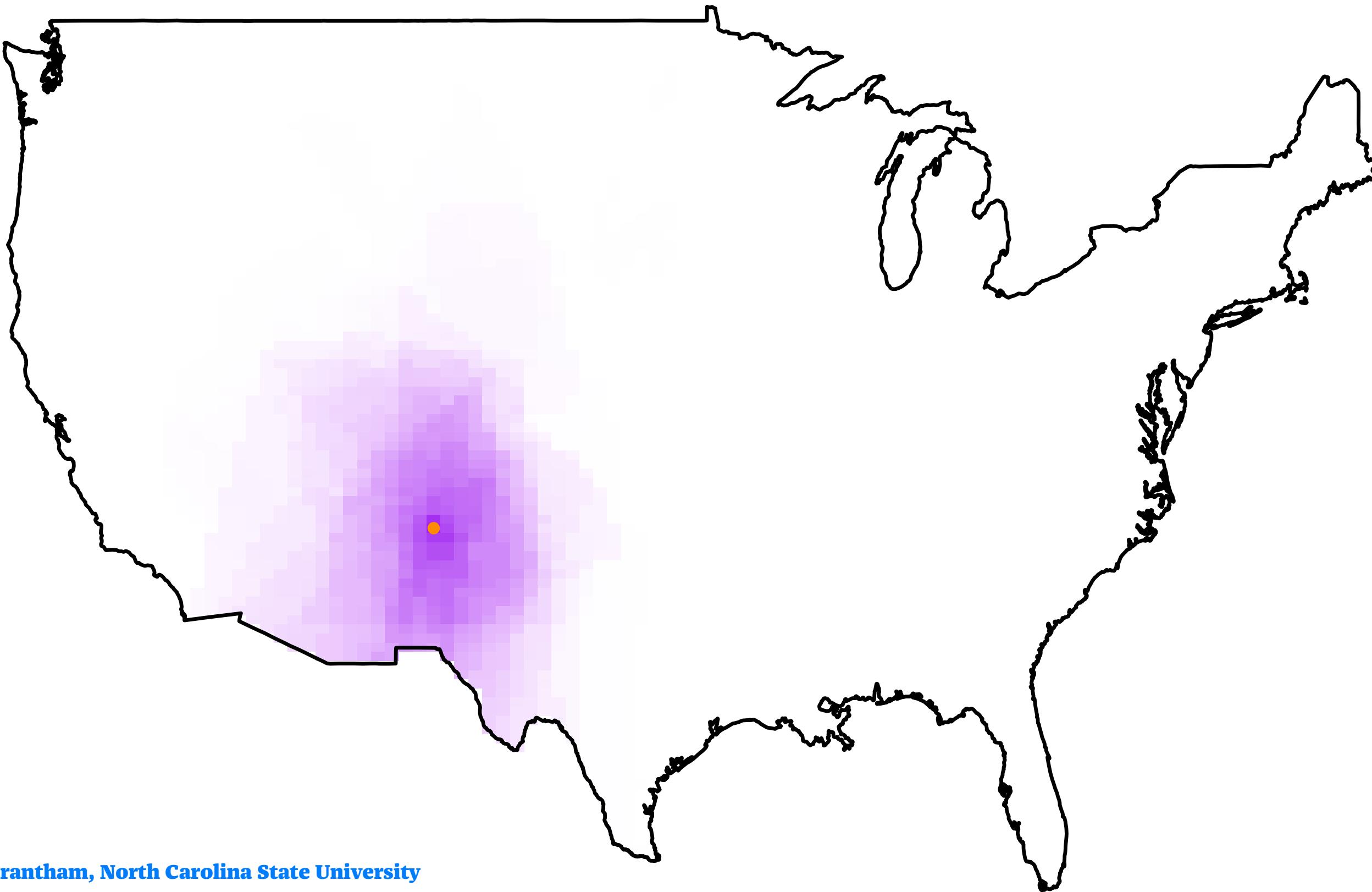


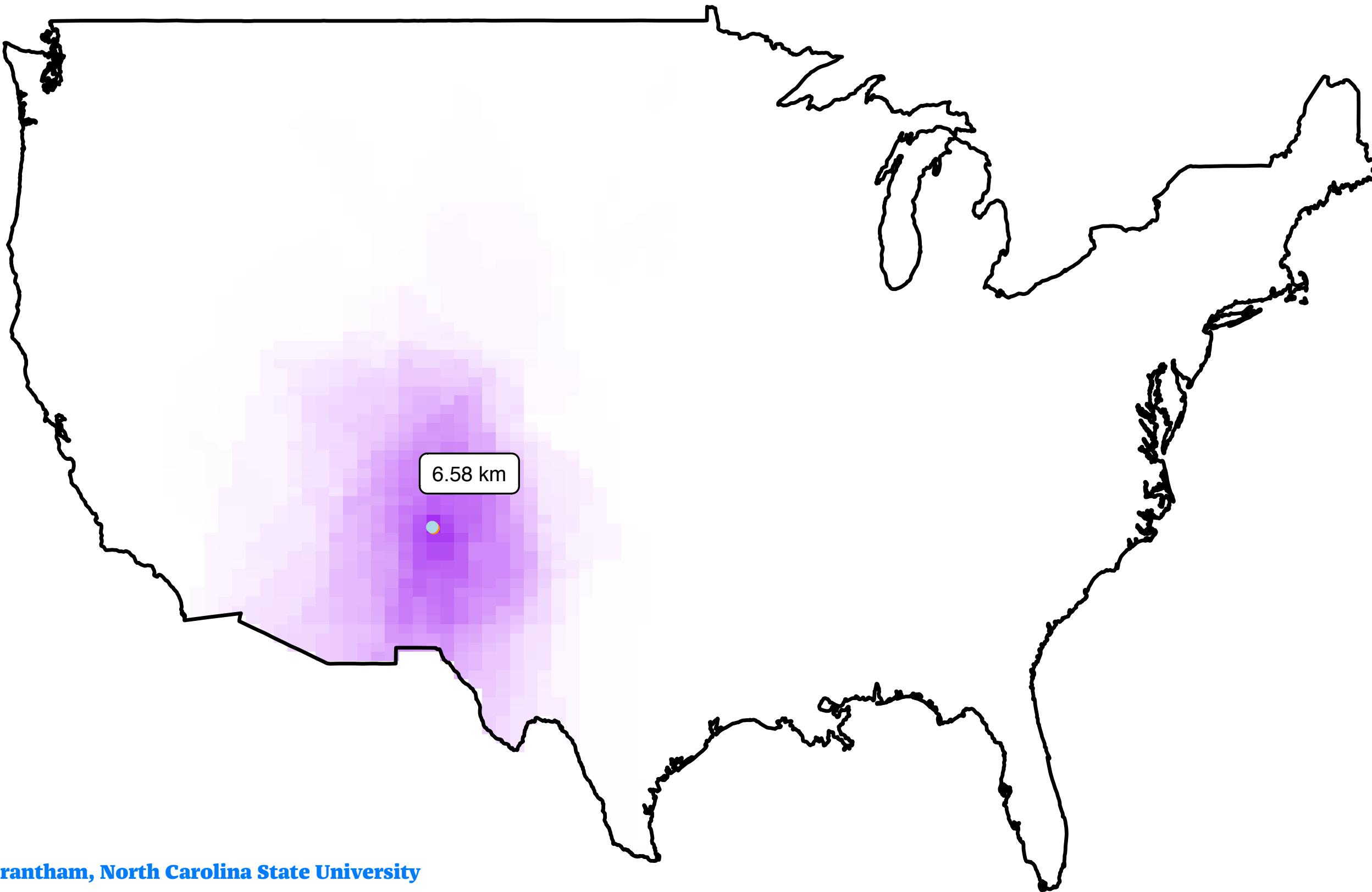


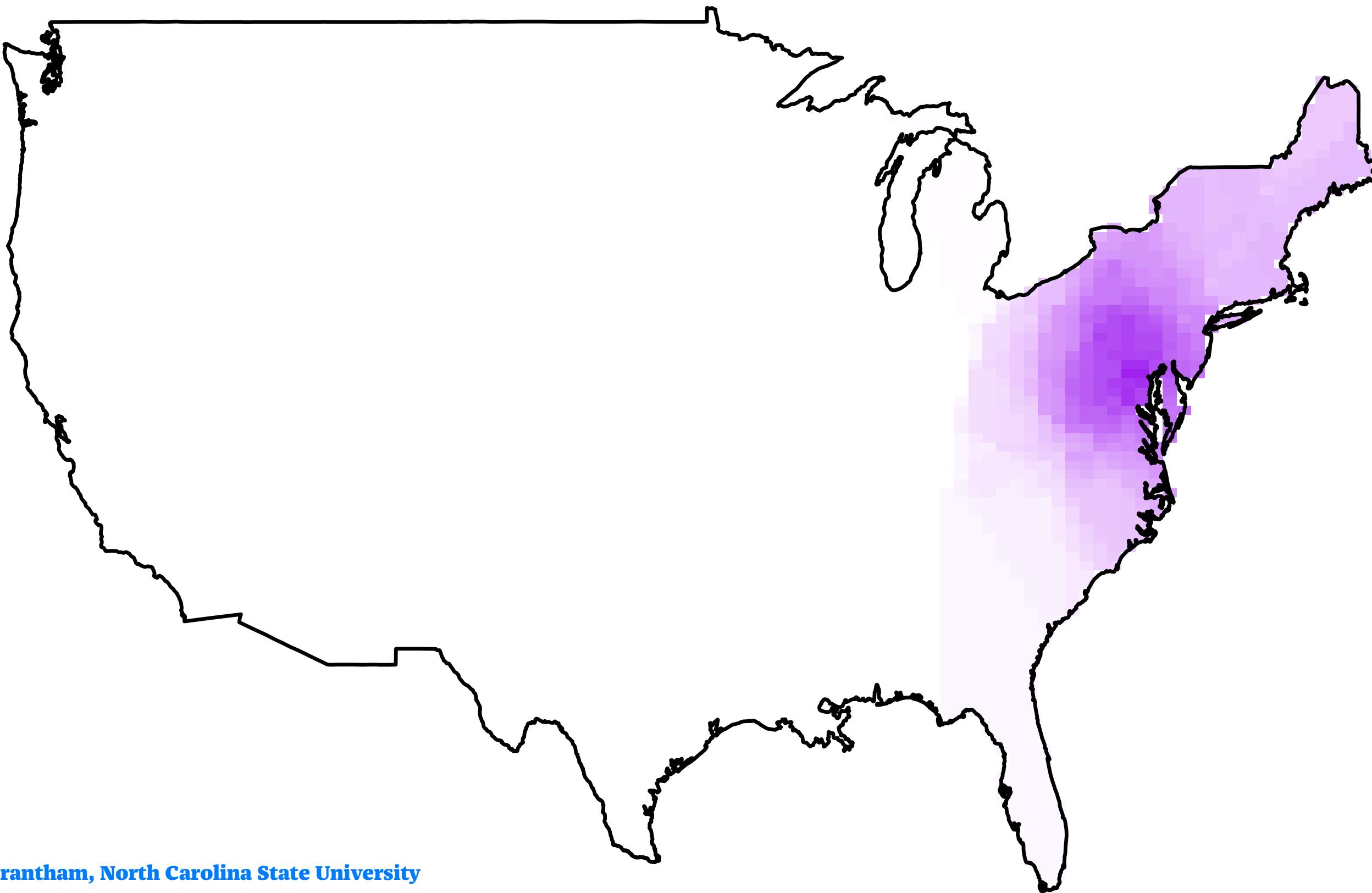


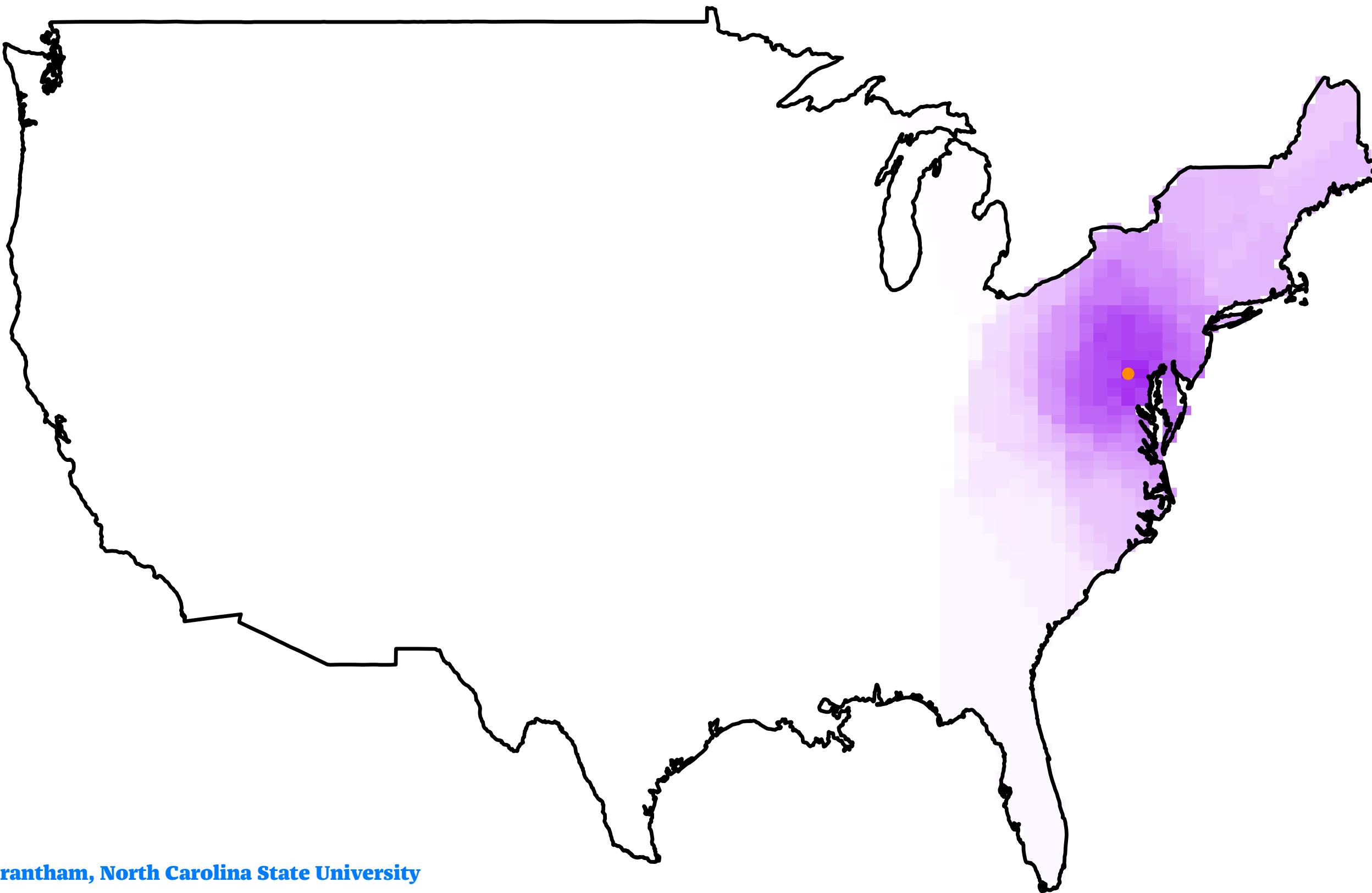


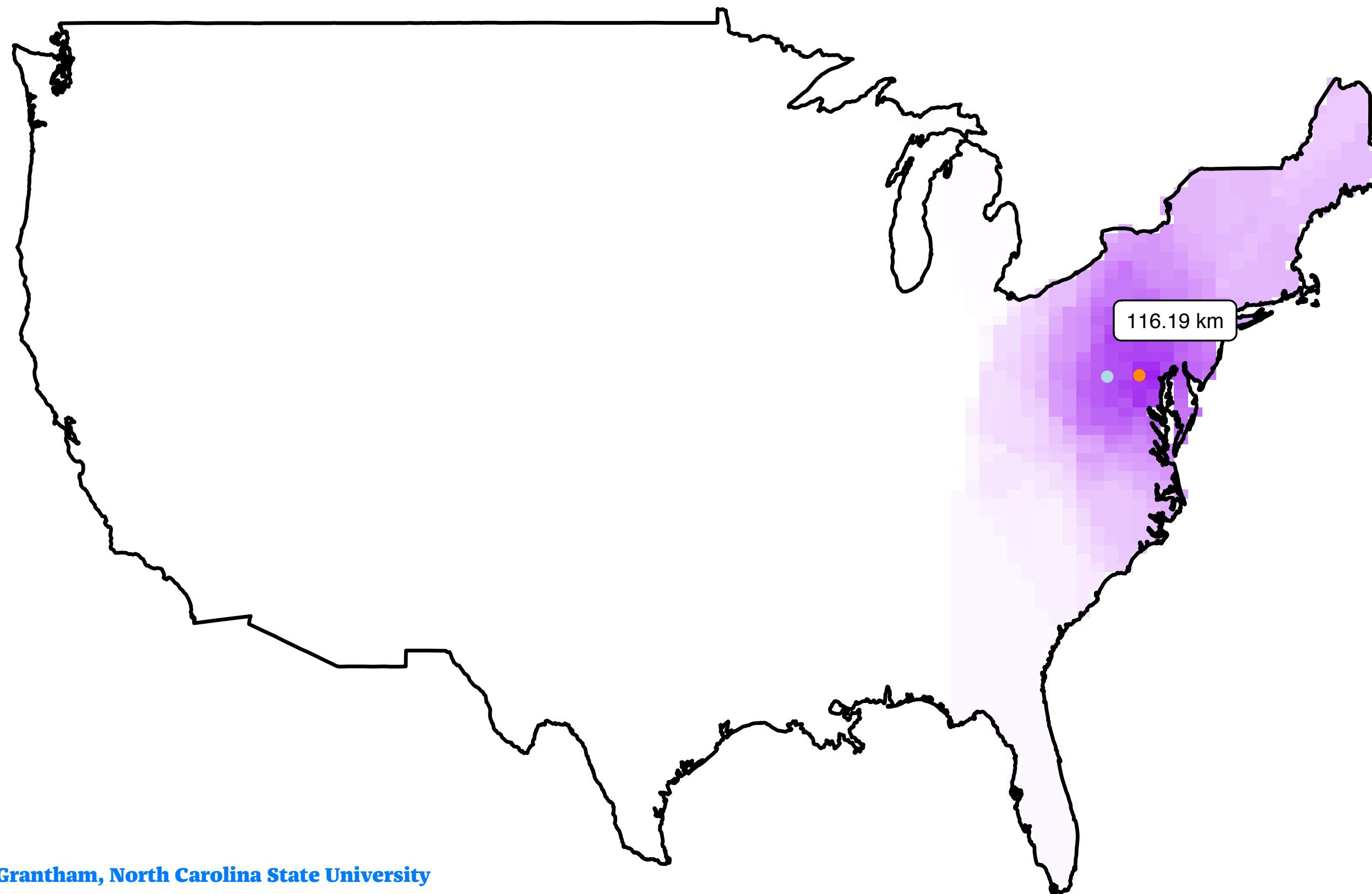




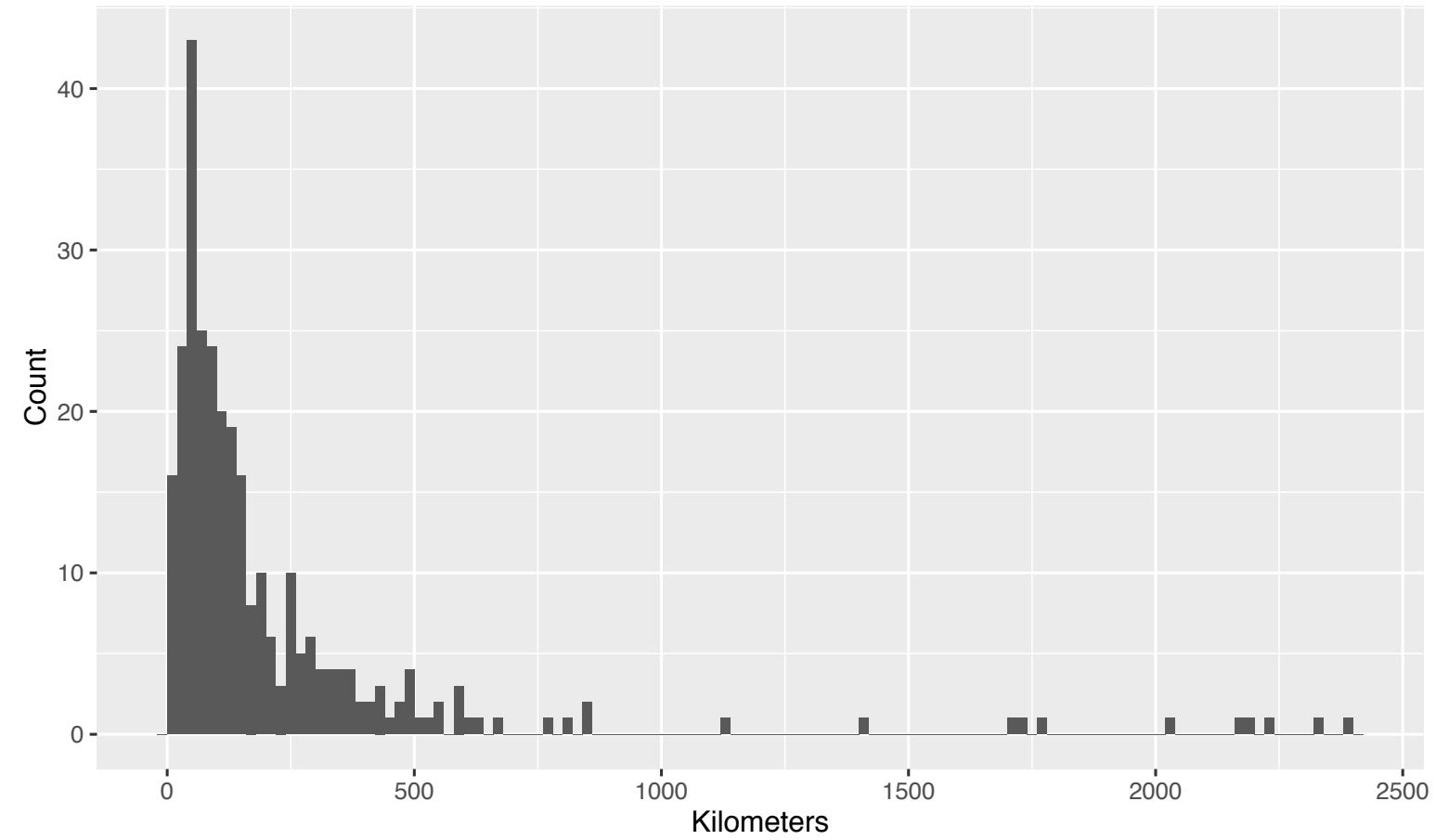








# Prediction errors



Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
1.751	53.990	114.700	230.000	245.500	2394.000

# Further work

Develop hypothesis testing framework to test if a sample originates from a particular county, state, etc.

Which fungi are most endemic to different biogeographies?  
Inference in deep learning is an active area of research.

New applications? The algorithm is not restricted to these data.

Python module `deepspace` to be made available at [github.com/  
nsgrantham/deepspace](https://github.com/nsgrantham/deepspace) upon publication.

# Thank You! Questions?

Slides available at [nsgrantham.github.io/documents/jsm-2016.pdf](https://nsgrantham.github.io/documents/jsm-2016.pdf)

Contact me:

- 👉 Website: [nsgrantham.github.io](https://nsgrantham.github.io)
- 👉 Github: [github.com/nsgrantham](https://github.com/nsgrantham)
- 👉 Twitter: [twitter.com/nsgrantham](https://twitter.com/nsgrantham)