# MnModel Phase 4

Project Summary and Statewide Results

Elizabeth Hobbs

June 24, 2019

# Contents

# Introduction

The Minnesota Department of Transportation (MnDOT) budgets over one million dollars annually for the identification and evaluation of historic and archaeological resources (historic properties) that are threatened by transportation related undertakings. Since 1997, MnDOT has been using its statewide archaeological predictive model, MnModel, to determine the most probable locations for archaeological resources and guide archaeological surveys conducted to meet the requirements set forth in Section 106 of the National Historic Preservation Act (NHPA) of 1966, as amended.

The MnModel project began in 1995 and completed three phases by 1998, culminating in the release of Phase 3 models (Hobbs 2003; Hudak et al. 2002). These models were used by MnDOT cultural resource professionals for more than twenty years. Although Phase 3 models performed well, the goal was always to improve the archaeological and environmental data, refine the modeling procedures, and improve model performance. Phase 4 of MnModel was undertaken to improve the accuracy of the predictive models used by MnDOT. Improvements are based on the acquisition and development of better GIS data since the completion of Phase 3 models in 1998. The majority of the effort in Phase 4 was devoted to data development. Good models cannot be built from poor data. Improved statistical procedures that have become available only in the last ten years also contributed to the success of the Phase 4 models.

## Archaeological Predictive Models

Archaeologists have been using Geographic Information Systems (GIS) since the 1980s to create models predicting archaeological site locations. Models have been developed for both large and small regions, primarily for Cultural Resource Management (CRM) purposes. In the US, the use of predictive models to avoid potential archaeological sites and to determine where to survey is considered a 'good faith effort' in fulfilling the requirements of Section 106 of the National Historic Preservation Act of 1966.

Volumes have been published on the subject, including methodologies (Ejstrud 2003; Judge and Sebastian 1988; Kvamme 1988; Parker 1985; Warren 1990; Westcott and Brandon 2000; Van Leusen and Kamermans 2005; Verhagen 2007), and critiques (Mehrer and Wescott 2006). Modeling methods can generally be categorized as either deductive or inductive. Deductive models, sometimes called 'expert systems' models, generally begin with theory to deduce where sites might be located. Inductive models begin with data, looking for patterns from known site observations that may predict where undiscovered sites might be found. These models use statistical methods to predict potential site locations based on observed correlations between known sites and a suite of environmental variables. MnModel is an inductive model, as are other models that cover very large areas.

Several authors have reviewed the state of archaeological predictive modeling both in the US (Kohler 1988; Kvamme 2006, 2011; Thoms 1988) and elsewhere (Kamermans 2011; van Leusen 1996). Aside from MnModel, statewide models have been produced for Washington (Kauhi and Markert 2009) and Pennsylvania (Harris et al. 2015). A model for North Carolina was begun (Madry et al. 2003) but not completed. Numerous models have been produced for smaller regions (Fry et al. 2004; Kvamme 1992; Warren and Asch 1996; Westcott and Kuiper 2000). Many of these are reported only in the grey literature (e.g. Cassell et al. 1997).

Critiques of inductive predictive modeling are several and have been identified by Kamermans (2011) as:

- Use of site data collected by non-random sampling. For models of small areas, it may be possible to conduct a truly random sample for archaeological sites, then to use the data collected to build a model of site location. For large areas, however, this is not feasible. Large area models rely on available archaeological site data collected over a number of years by many different teams of archaeologists. Surveys were conducted in locations where archaeologists expected to find sites or where a construction project required that any historic resources present be discovered. Even within project areas, surveys are seldom random, though in small project areas the entire area is sometimes surveyed. Because very large numbers of sites are required for statistical modeling, modelers cannot afford to be use only sites from random or 100 percent surveys. The implications of this non-random sample are that we cannot know for certain whether areas predicted to have low probability for sites are actually unlikely to have sites or if they simply have not been adequately surveyed. MnModel addresses this dilemma by also modeling survey distributions (see below).

- Combining sites of different time periods and functions in the same model. This is another problem that is necessitated by the inadequacies of the archaeological database. In Minnesota 45 percent of recorded sites have been assigned a temporal affiliation and only 29 percent have been assigned a function. Once these subsets of the database have been further subdivided into individual time periods and functions, there would be far too few sites in any group for modeling.

- Failing to consider how 'proxy' variables contribute to site location. By 'proxy' variables, Kamermans (2011) and Kohler and Parker (1986) are referring to variables that may have a number of different meanings to people looking to select a site for habitation or other use. Kohler and Parker (1986) use the example of the 'elevation', which may be a proxy for micro-climate, soil moisture, visibility, vegetation type, or other important site-selection considerations. In MnModel Phase 3 we certainly observed the strong effect of proxy variables when some high probability areas proved to coincide with the shores of drained lakes, about which the model had no information. In MnModel Phase 4 we have attempted to address this issue by deriving more complex variables from elevation to make some of these site selection criteria more explicit. These new variables include visibility (based on viewshed analysis), Topographic Position Index, Topographic Wetness Index, and Shelter Index. We have addressed the issue of drained lakes by developing a historic hydrographic model.

- Low spatial resolution. Early models were limited both by computing capacity and the scale of available environmental data. Kohler and Parker (1986) consider a cell size less than one acre to be very small and 1.5 ha cells to be high resolution. MnModel Phases 3 and 4 both model 30 m cells, though functional resolution is really dependent on the source data. When MnModel Phase 3 was developed, 1:24,000 scale data were just becoming available for terrain and hydrographic features. Geomorphology, soils, and vegetation data were at scales ranging from 1:100,000 to 1:500,000. In Phase 4, we have improved the resolution of all data sources.

- Inappropriate statistical tools. A wide variety of statistical tools have been used for modeling. It is beyond the scope of this report to review them all. Suffice it to say that MnModel has relied on the advice of professional statisticians to develop our statistical methods. Phase 3 models used stepwise

multiple logistic regression, which seemed to be the best solution available at the time. For Phase 4, Oehlert and Shea (2007) recommended bagging, a 'perturb and aggregate method'. We implemented the Phase 4 models using Random Forest, an improved version of bagging. Random Forest was also used to develop the model for Pennsylvania (Harris et al. 2015).

- Little model validation. Both MnModel Phases 3 and 4 have built models with one subset of the archaeological data and test with a separate subset. However, in both cases, final models have been built with the entire population of archaeological sites, leaving no population for further testing. The best test of any model would be random field survey. For large areas, such as Minnesota, the expense would be considerable.

- Relevance of environmental data. Most models, including MnModel Phase 3, have been developed using modern environmental data for deriving predictor variables. This is problematic because climate change and human agency over the past 10,000 years have greatly altered the environment nearly everywhere. For MnModel Phase 4 we have attempted to mitigate this problem by creating historic and prehistoric models of surface hydrography and a historic vegetation model.

- Need for cultural data. The absence of cultural data in MnModel, and likely in other models as well, is a limitation of the database. As mentioned above, Minnesota's archaeological database lacks information even on cultural affiliation for most sites.

Despite these limitations, models are useful tools for Cultural Resource Management. They identify the areas at highest risk for impacts to cultural resources. MnModel, by making the distribution of past surveys explicit, also identifies the areas for which we have the least information. They are best used as an advanced planning tool to encourage avoidance of high risk areas. They should not be used to justify archaeological survey only in the high probability areas. As long as users understand that some percentage of sites will always be found in the low probability areas and survey accordingly, models can be a powerful tool for protection of archaeological resources.

# Purpose

The primary objective of the MnModel project is to create accurate digital maps capable of alerting planners to the presence of potential pre-contact archaeological properties in accordance with the identification requirements set forth in Section 106. By using Geographic Information Systems, digital maps are created that delineate areas of high archaeological site potential based on statistical correlations between environmental attributes and known archaeological site locations. Linking this information with maps of high and low survey coverage directs where archaeological survey efforts should be concentrated. It also assists planners in avoiding areas that potentially contain cultural resources requiring costly mitigation or in weighing the cost of their disturbance against other project effects, such as wetland disturbance or socioeconomic impacts. This approach permits planners to conduct advance planning and base decisions on sound scientific findings.

# MnModel Phase 4 Goals

The purpose of MnModel Phase 4 was to improve the performance and reliability of the archaeological predictive models.  Better models would be characterized by one or more of the following:

- A higher proportion of sites predicted.  Archaeologists generally prioritize surveys in high site potential areas.  We want as many known sites as possible to fall within these areas.

- A smaller area of high site potential.  Archaeological survey is expensive.  If high site potential areas cover a large portion of the state, survey costs will be high.

- A smaller area of 'unknown' site potential.  We label areas as 'unknown' if both site probability and survey probability are low.  We do not know if few sites have been found there because site potential really is low or simply because very little survey has been done.

- Greater model stability.  We have more confidence in models that do not vary much if we use different parts of the same dataset to build them.

In the Model Performance section of this report, we evaluate the project's ability to meet these goals.

# Phase 4 vs. Phase 3

## Conceptual Framework

Phase 4 models are built on the same conceptual framework as the original modeling (Hudak et al. 2002). The underlying assumption of the models is that the most important factors controlling pre-contact hunter-gatherer settlement and activity location decisions were physical and biotic attributes of the landscape (Dalla Bona 1994; Kohler and Parker 1986; Kvamme 1985). Phase 4 places greater emphasis on the geomorphological paleo-landscape, historic and prehistoric surface hydrography, and historic vegetation to improve model accuracy. Phase 4 also benefits from the larger population of archaeological site and survey data collected in the past twenty years.

## Similarities and Differences

When Phase 3 of MnModel was developed, statewide GIS data were just becoming available for Minnesota. Few datasets were as good as 1:24,000 scale.  Models were developed on UNIX workstations that were less powerful than today's personal computers.  It was not possible to calculate visibility or cost distances for more than a handful of points at a time.  Moreover, the research team had no experience with predictive modeling, and no predictive model of this magnitude had ever been attempted.

Phase 4 of MnModel provided the opportunity to take advantage of more and better data, more powerful computers, a wider array of GIS and statistical functions, and all of the lessons learned in the course of Phase 3 (Table 1).  Shortcomings of Phase 3 included:

- Inaccurate site point locations. Point locations were based on UTM coordinates recorded in the MN State Historic Preservation Office database. Sources of error included misinterpretation of the coordinates on paper USGS quad sheets, mistakes made when entering these coordinates on site forms, and typographical errors when entering the coordinates into the database. Many of these errors were not discovered until after modeling was complete. A total of 74 Phase 3 sites had their locations corrected prior to inclusion in the Phase 4 database. An additional 40 Phase 3 sites were removed from the database because they lacked a prehistoric component.

- Incomplete survey data. Survey mapping for MnModel Phase 3 was limited, first, to surveys considered 'probabilistic' by the project archaeologists (Gibbon et al. 2002). In reality, only the surveys conducted for MnModel in the summer of 1996 used a probabilistic survey design. Additional 'non-probabilistic' surveys were included if they met specified criteria (Gibbon et al 2002: Table 5.6). These surveys were a small portion of those that had been conducted at the time.

- Low resolution terrain data. This is particularly a problem in the very flat Lake Agassiz Basin. Moreover, the data available at the time contained errors (Hobbs 2002a) that produced false high probabilities in the resultant models.

- Modern surface hydrography. Euro-American settlement of Minnesota has greatly altered the surface water features, particularly in the agricultural regions of the western part of the state. Many wetlands and lakes have been drained. In Phase 3, we had no source data for these features that were likely to have been important to prehistoric people. Nevertheless, the Phase 3 models sometimes picked up topographic differences between lake beds and their former shorelines and relied on these differences to identify areas of high potential.

- Scale of geomorphic and vegetation data. The best geomorphic data available for Phase 3 modeling was rasterized at a 40 acre resolution. Higher resolution (1:100,000 statewide and 1:24,000 for select river valleys) data became available near the end of the project, but too late for our use. Historic vegetation data were digitized from paper maps at 1:500,000 scale. Not only were they very generalized, they registered poorly with the hydrography and terrain.

- Euclidean distance. At the time Phase 3 was developed, it was not feasible to calculate cost distances from every cell to the closest feature of interest, such as the nearest lake. The Euclidean distances calculated may have been acceptable near target features, but with increasing distance there is nothing to keep the line from crossing steep terrain or water bodies.

- Limited processing capacity. The ArcINFO UNIX workstations used for Phase 3 could not have performed many of the raster procedures used to develop Phase 4 variables (visibility, topographic wetness index, topographic position index, cost distances) or to apply the random forest models to the regional point grids.

**Table 1: Summary of differences between MnModel Phase 3 and MnModel Phase 4**

| Measure | Phase 3 | Phase 4 |
|---|---|---|
| Number of archaeological site points | 6,828 | 9,611 |
| Number of archaeological survey points | 23,443 | 37,256 |
| Sampling Method | Value at Point | Mean or Most Frequent Value within Polygon |
| Resolution of digital elevation data | 30 m | 10 m |
| Scale of geomorphic data | 40 acres (source data 1:1,000,000) | 1:5,000 – 1:100,000 |
| Source of surface water data | Modern | Historic Model |
| Scale of historic vegetation data | 1:500,000 | ~1:24,000 (model) |
| GIS Software | Arc/Info GRID | ArcGIS Spatial Analyst |
| Statistical Software | S-Plus | R |
| Distance measure used | Euclidean | Least-Cost Path Distance |
| Statistical Model Used | Multiple Stepwise Logistic Regression | Random Forest |

## Units of Analysis

The basic unit of analysis in MnModel is a 30-meter square cell or parcel of land. This provides the models with a resolution that is comparable to mapping at a scale of 1:24,000.  Early in Phase 4 we considered using a 10-meter square cell to be consistent with the updated Digital Terrain Model (DTM).  We quickly determined that the statistical software would not be able to handle models at this resolution.  Moreover, since no data besides the DTM approached this resolution, 30-meter cells are more appropriate for the analysis.  Thus model resolution did not change between Phases 3 and 4.

## Regionalization

Because Minnesota is a large state (218,601 km$^2$ or 85,254 mi$^2$) with considerable environmental variation, it is necessary to model by smaller, relatively homogeneous regions then mosaic the regional models into a statewide model.  Phases 1-3 of MnModel experimented with different regionalization schemes. Phases 1 and 2 divided the state into Archaeological Resource Subregions (Anfinson 1990).  County boundaries rather than natural features were used for the borders of some of these regions, resulting in abrupt changes in probability values along these borders.  More natural boundaries, based on Ecological Classification System (ECS) subsections (Hanson and Hargrave 1996) were adopted in Phase 3 (Hobbs 2002a).  Minnesota's ECS is part of a hierarchical national system of classification (Cleland et al. 1997) based on climate, geomorphology, terrain, soils, and vegetation.

Phase 4 regions were also based on ECS subsections, but with some changes.  First, the boundaries of the subsections were revised based on higher resolution mapping developed for classifying Minnesota's landtype associations, the next lower level in the classification hierarchy (Figure 1).  Second, some of the modeling regions based on combined subsections were revised.  There were 24 ECS subsections in the original classification and are now 26 ECS subsections.  Subsection sizes and the number of sites they contain vary considerably (Table 2), as do site frequencies.  As we need about 100 sites for analysis and modeling, some of the smaller subsections must be combined with each other or with adjacent, larger subsections.  Site numbers are a function of several factors:

- Intensity of survey.  More developed parts of the state (Anoka Sand Plain, Big Woods, St. Paul-Baldwin Plains & Moraines) tend to have a higher density of surveys.  The exception would be the National Forests (particularly within Chippewa Plains and Border Lakes), which have active archeological survey programs.
- Terrain and hydrography: Lowest site frequencies are in regions where swamps and bogs occupy large parts of the landscape (Agassiz Lowlands, Tamarack Lowlands).  Not only are these areas less well suited for habitation, forests and wetlands make access for archeological survey difficult.
- Suitability for prehistoric habitation: The Rochester Plateau is an area of high terrain with no lakes and only the upper reaches of streams.  Reliable water supplies may have limited where people could live. The Blufflands, on the other hand, has close access to the Mississippi River, a major prehistoric transportation corridor.  It supported a large prehistoric population and, probably for that reason, has also been well surveyed.  On average, known site frequencies in Minnesota are quite low (0.044 sites per km$^2$).  Yet in places more suitable for habitation, frequencies may be as high as 0.917 sites per km$^2$ (Table 2).

**Figure 1: Revised ECS Subsections for Minnesota, with Outlines of Original ECS Subsections for Comparison**



Aspen Parklands

Agassiz Lowlands

Littlefork-Vermilion Uplands

Border Lakes

Nashwauk Uplands

Laurentian Uplands

Chippewa Plains

Red River Prairie

Toimi Uplands

North Shore Highlands

St. Louis Moraines

Tamarack Lowlands

Pine Moraines & Outwash Plains

Glacial Lake Superior Plain

Hardwood Hills

Mille Lacs Uplands

Anoka Sand Plain

St. Croix Moraine

St. Paul-Baldwin Plains & Moraines

Minnesota River Prairie

Big Woods

Coteau Moraines

Rochester Plateau

Inner Coteau

Oak Savanna

Blufflands

0   25   50        100 Miles
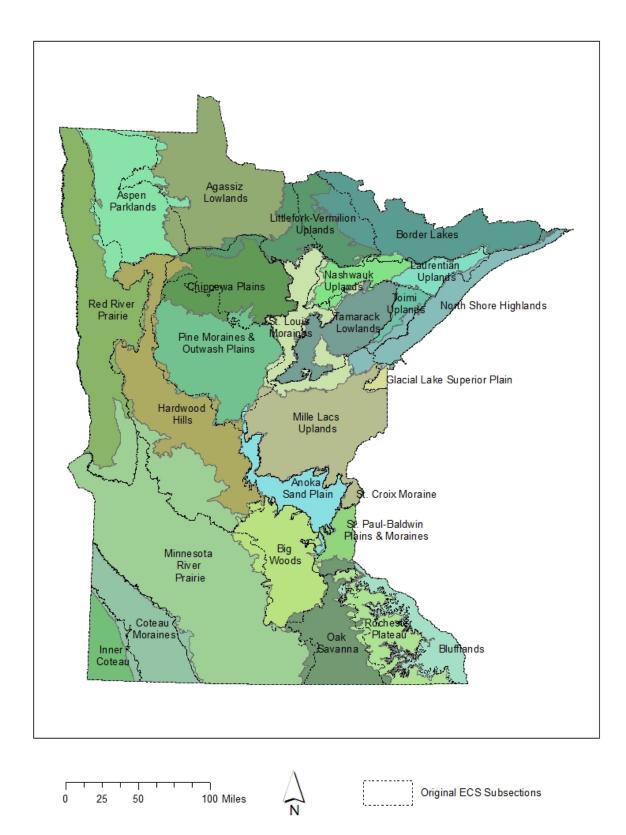
N

Original ECS Subsections

**Table 2: Assignment of ECS Subsections to Phase 4 Modeling Regions**

| Current ECS Subsection | Phase 4 Modeling | Area (km²) | Number of Phase 4 Sites | Site Frequency (sites/km²) |
|---|---|---|---|---|
| Agassiz Lowlands | AGLV | 14800 | 41 | 0.003 |
| Anoka Sand Plain | ANOK | 4857 | 513 | 0.106 |
| Aspen Parklands | ASPK | 11767 | 120 | 0.010 |
| Big Woods | BGWD | 8951 | 941 | 0.105 |
| Border Lakes | BDLK | 11219 | 1174 | 0.105 |
| Chippewa Plains | CHIP | 8913 | 846 | 0.095 |
| Coteau Moraines | COTM | 8283 | 376 | 0.045 |
| Glacial Lake Superior Plain | MLAC | 444 | 16 | 0.036 |
| Hardwood Hills | HRDH | 14152 | 532 | 0.038 |
| Inner Coteau | ICOT | 3144 | 127 | 0.040 |
| Laurentian Uplands | NSUP | 2296 | 24 | 0.010 |
| Littlefork-Vermilion Uplands | AGLV | 6686 | 50 | 0.007 |
| Mille Lacs Uplands | MLAC | 13715 | 520 | 0.038 |
| Minnesota River Prairie | MNRP | 37728 | 1450 | 0.038 |
| Nashwauk Uplands | NSUP | 3278 | 26 | 0.008 |
| North Shore Highlands | NSHH | 5997 | 230 | 0.038 |
| Oak Savanna | OSAV | 7363 | 239 | 0.032 |
| Pine Moraines & Outwash Plains | PINE | 12245 | 759 | 0.062 |
| Red River Prairie | REDR | 15988 | 330 | 0.021 |
| Rochester Plateau | PLAT | 5503 | 125 | 0.023 |
| St. Croix Moraine | MLAC | 12 | 11 | 0.917 |
| St. Louis Moraines | STTA | 6670 | 226 | 0.034 |
| St. Paul-Baldwin Plains & Moraines | STPB | 1876 | 178 | 0.095 |
| Tamarack Lowlands | STTA | 6127 | 22 | 0.004 |
| The Blufflands | BLUF | 5214 | 691 | 0.133 |
| Toimi Uplands | NSUP | 1373 | 48 | 0.035 |

# Methods

Archaeological predictive modeling requires a series of steps from data development, derivation of variables, and sampling, to statistical analysis and modeling, and finally model evaluation and classification.  MnModel Phase 4 methods are detailed in several documents, so will only be summarized here.

## Data Development

Archaeological predictive modeling is premised on the assumption that archaeological site locations bear some relationship to the natural environment, as prehistoric populations depended on the environment for their livelihood.  These models are necessarily environmentally deterministic, since we have little cultural data about the ancient populations to incorporate into our models.  Thus both archaeological and environmental data are necessary for modeling.

The first priority for Phase 4 was to more accurately represent the locations of known prehistoric archaeological sites and archaeological surveys.  Digitizing these as polygons allowed us to better understand the range of environments in which they are found.  Moreover, both populations had increased dramatically since the Phase 3 models were developed.

The environmental data used in Phase 3 had two problems.  They were low resolution, and they reflected modern conditions.  A goal for Phase 4 was to improve data resolution and to better reflect conditions prior to Euro-American settlement.

A number of projects were undertaken to improve the data.  These are summarized below.  More detail can be found on the MnModel web site.

### Known Archaeological Sites

The number of archaeological sites available for modeling increased 41 percent between Phases 3 and 4, from 6,828 to 9,602.  From 2012 to 2015, MnDOT contractors digitized archaeological site polygons from records at the Minnesota State Historic Preservation Office (MnSHPO) and Minnesota Office of the State Archaeologist (MnOSA).  The US Forest Service (Chippewa and Superior National Forests) generously shared their data. The resultant archaeological database was refined at Mankato State University (MSUM) to be suitable for modeling.  This involved, first, removing sites we cannot expect to predict: sites with only historic contexts, single artifacts, and sites dependent on resources not represented in our environmental data (particularly sites, such as quarries and rock art, requiring a rock outcrop).  In addition, it was necessary to remove overlapping and duplicate polygons.  The initial database consisted of 8,836 polygons ranging in size from 7 $m^2$ to 5.2 $km^2$.  These sites were used to develop the first round of Phase 4 models.  In 2019, MSUM added 766 more site polygons to the database. Some of these had inadvertently been removed from the original database, while others were recorded between 2015 and 2018.  The augmented database was used to develop the final Phase 4 models.

## Archaeological Surveys

From 2012 to 2015, MnDOT contractors digitized archaeological survey polygons from records at MnSHPO and MnOSA.  Survey data are of varying quality.  Some areas mapped were thoroughly surveyed using modern methods.  Other areas were simply project areas within which some, but not all, areas were surveyed.  To indicate the level of confidence that the surveyors actually surveyed the entire mapped polygon, a confidence value was assigned.  Survey density varies throughout the state (Figure 2), with the highest concentrations in the Twin Cities metropolitan area and Chippewa National Forest.

We prepared the data for modeling by removing overlapping survey polygons, giving preference to polygons with higher confidence values, and splitting very large polygons into smaller units for modeling (Hobbs et al. 2019b).  The resultant database included 36,297 surveyed polygons with a maximum size of four square kilometers.  Note that this does not indicate the number of surveys.  Surveys can cover very large areas, and a single survey can be represented by multiple polygons.  In both Phases 3 and 4, individual surveys polygons were represented in the modeling database by multiple points, though with different spacing standards.  Compared to Phase 3, the number of surveyed points increased 59%, from 23,443 to 37,249.
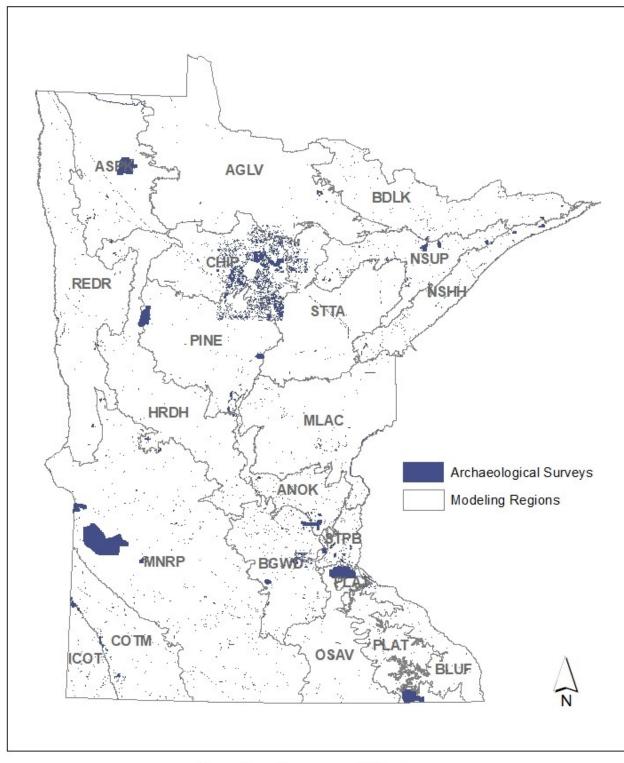
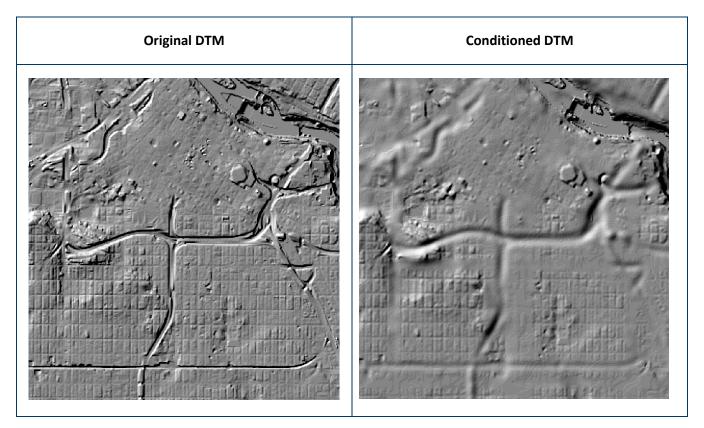**Figure 2: Archaeological Surveys Used for Phase 4 Modeling**

## Terrain Model

Phase 3 modeling utilized 30-meter resolution digital elevation models (DEMs) from US Geological Survey (USGS), which were available for only 81% of the state, supplemented by 1:250,000 scale elevation data for the remainder of the state. Some of the USGS DEMs suffered from distortion and elevations did not always match along edges of quad sheets (Hobbs 2002b). These issues were apparent in the resulting predictive models, as striping and as high probability areas along false 'ridges' at quad sheet boundaries.

In 2012, we created a ten-meter resolution digital terrain model (DTM) from one-meter resolution LiDAR data. We quickly discovered that modern infrastructure and surface disturbances are quite prominent in the higher resolution data (Figure 3) and expected these would be reflected in the models, just as the terrain data errors were in the Phase 3 models.
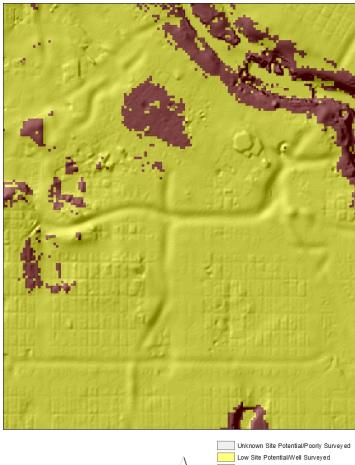
**Figure 3: 10 m DTM Before and After Conditioning**

| Original DTM | Conditioned DTM |
|---|---|
|  |  |

In 2017, we updated the terrain data and developed procedures to restore the DTM to something approximating a pre-modern surface. This process involved several procedures. We used bathymetric data from the Minnesota Department of Natural Resources (MnDNR) to replace lake plains for large lakes. This provided us with relief in places that have been flooded by reservoirs. We replaced a portion of the Mesabi Iron Range with topographic data from 1899 (Lively et al. 2002), restoring the pit mines to a more natural surface. To minimize the effects of disturbance and infrastructure, we buffered existing features (roads, ditches, gravel pits, railroads, airports) extracted from the MnDOT route centerlines, gSSURGO soils data, and the National Wetlands Inventory (NWI), then used the composite buffer to remove these areas from the DTM. We then developed a

set of custom processing tools in Python to search for the resulting "NoData" cells and replace these cells with values.  The procedure used a dynamic cut-fill process that referenced the existing terrain using multiple, iterative passes to fill in the 'No Data' areas one row of cells per pass starting along the outermost edge.  The goal of this procedure was to raise ditches and lower road crowns to a calculated local mean elevation approximating the original terrain surface.   A secondary goal was to reduce the slopes within the replacement zones to less than 15 degrees, as MnModel Phase 3 models were found to have a sensitivity to slopes of 15 degrees or greater.  Finally, the TauDEM Pit Remove tool was run to remove pits from the DTM.  Results were not as effective as we had hoped they might be (Figure 3).  This conditioned DTM was used to derive a suite of predictor variables (Hobbs, Walsh, and Hudak 2019).  Despite the modest success in smoothing modern features, the final models show no apparent artifacts from infrastructure (Figure 4).

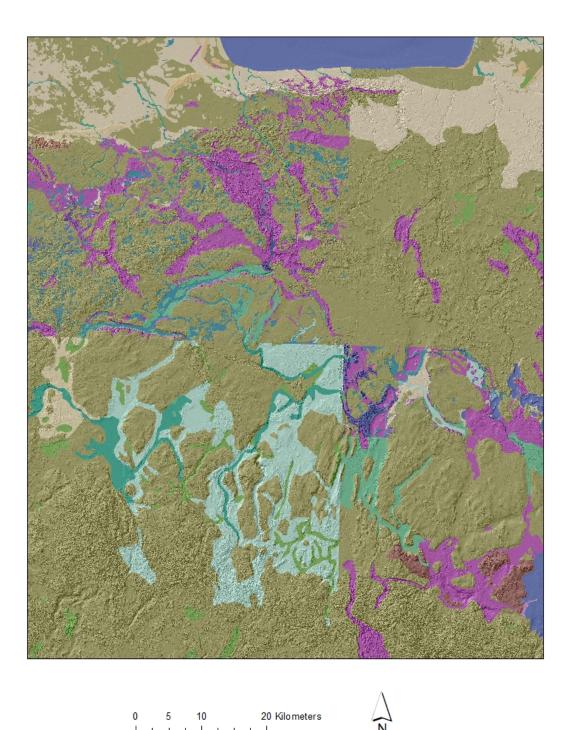**Figure 4: Phase 4 Survey Implementation Model vs. Conditioned DTM**



Legend:
- Unknown Site Potential/Poorly Surveyed
- Low Site Potential/Well Surveyed
- High Site Potential/Poorly Surveyed
- High Site Potential/Well Surveyed

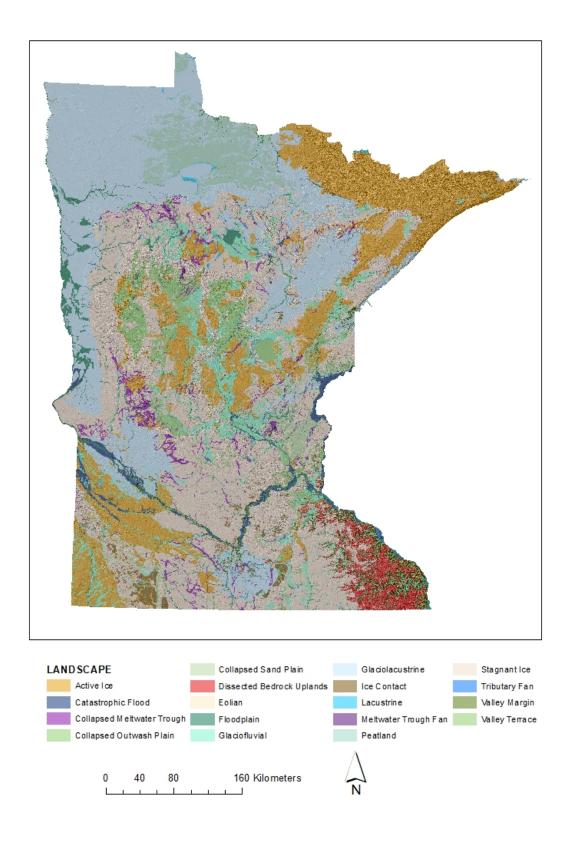0    0.275    0.55         1.1 Kilometers

N

## Landscape Model

Geomorphic variables used in Phase 3 were limited and were derived from State Soil Atlas sheets rasterized at a 40-acre resolution (Hobbs and Nawrocki 2002). MnDNR published their 1:100,000 scale statewide landforms data while Phase 3 was underway, but too late for our use. Landform Sediment Assemblage (LfSA) mapping (Hudak and Hajic 2002), begun in Phase 1 of MnModel, was also not ready in time to contribute to the Phase 3 models. For Phase 4, we extended LfSA mapping to the Mississippi River between St. Cloud and St. Paul (Hajic and Hudak 2002), the Anoka Sand Plain (Hajic, Hudak, and Walsh 2009), and the Mississippi River between St. Paul and Iowa and the Zumbro River (Hajic, Hudak and Walsh 2011).

From 2014 to 2016, we created a statewide landscape model by re-classifying and mosaicking the LfSA mapping, digital maps of various parts of the state published by Minnesota Geological Survey (MGS), and, where no other data were available, the geomorphic map published by MnDNR. Many of these data sources were surface geology maps and not specifically geomorphic in nature. That said, many of these surface geology maps did include geomorphic landforms as part of their mapping process. The upland data sources often conflicted with respect to glacial geology and phases, and it was up to the MnDOT team to make a best judgement case in favor of one data source over the other where the two data sources edge-matched and beyond. Our revision and reclassification of the original data sources did not include reshaping or redrawing any of the geospatial data (i.e., lines and polygons). On rare occasion, polygons of like tabular values on either side of an edge-matched seam were joined to help reduce both editing time and the appearances of straight-line edges.

Each of the mosaicked portions were categorized by Region, Region Name, Subregion, Subregion Name, Landscape, Landform, and Mantle. Some reinterpretations at and near boundaries with adjoining model segments were made to provide for a more logical transition at the mosaicked seams. Nevertheless, it was impossible to make all sources consistent with each other, and edge effects are apparent in some places (Figure 5). This challenging project created a statewide map of varying map scales but with a consistent hierarchical classification. In all, the model maps 89 landforms that are nested within 16 named regions, 220 named subregions, and 18 landscapes (Figure 6).

**Figure 5: Landform Discrepancies Where Different Source Data Meet, MnModel Phase 4 Landscape Model**

**Figure 6: MnModel Phase 4 Landscape Model: Landscapes**



LANDSCAPE
- Active Ice
- Catastrophic Flood
- Collapsed Meltwater Trough
- Collapsed Outwash Plain
- Collapsed Sand Plain
- Dissected Bedrock Uplands
- Eolian
- Floodplain
- Glaciofluvial
- Glaciolacustrine
- Ice Contact
- Lacustrine
- Meltwater Trough Fan
- Peatland
- Stagnant Ice
- Tributary Fan
- Valley Margin
- Valley Terrace

0  40  80  160 Kilometers
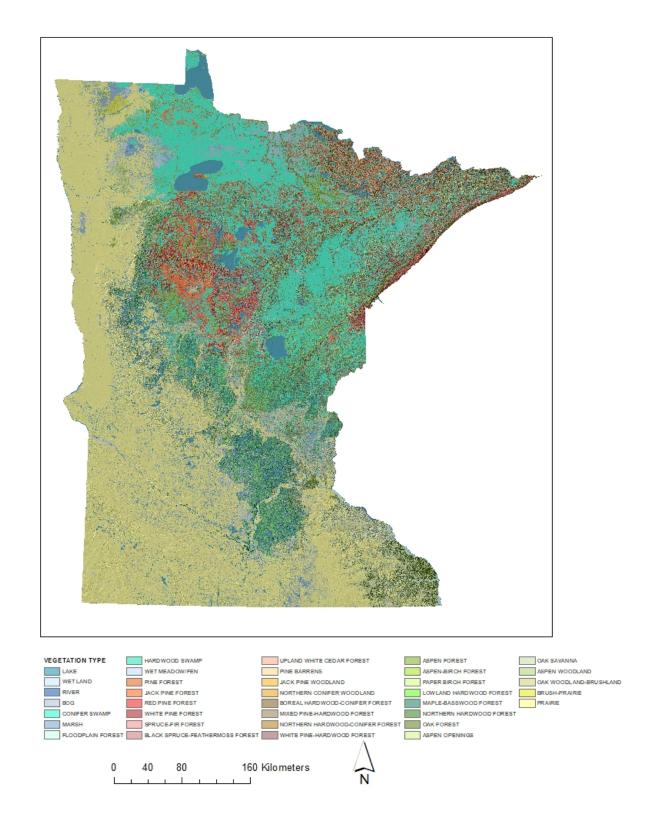
N

## Vegetation Model

In 1930, Francis J. Marschner produced a map of Minnesota vegetation complied from the Public Land Survey notes (Marschner 1974). This 1:500,000 scale map was later digitized by MnDNR and used as the source of vegetation variables for MnModel Phase 3. The digitized Marschner map has several problems as a data source for modeling. First, it is very generalized. Many features are mentioned in the surveyors' notes and illustrated on the plat maps but do not appear on the Marschner map. Second, Marschner's methods were not documented, and his vegetation classification scheme is not ideal for our purposes. Finally, and most important, the map does not register well with terrain. Most conspicuously, lakes and wetlands do not overlay their basins.

In 2013, we created a statewide mosaic of the scanned and georeferenced Public Land Survey plat maps and digitized polygons of hydrographic and vegetation features. Vegetation observations and species of bearing trees made by the surveyors had been extracted to section and quarter section corner points and published by MnDNR in 1997. More recently, John Almendinger (MnDNR) made available surveyor line notes extracted to section lines for the northern half of the state. We reclassified the vegetation point data using a classification system developed by MnDNR (Aaseng 1993). The line and bearing tree data were used as needed to help make decisions for classifying the corner points.

The classified point data were used to develop a high resolution statistical model of historic vegetation (Hobbs 2019). The modeling procedures were very much like those used to create the archaeological predictive models (Landrum and Hobbs 2019). Modeling was done for the same modeling regions with one exception: ICOT and COTM were combined to provide enough sample points of rare vegetation types (primarily wetlands) for modeling. The set of predictor variables used was smaller, since distance to water and certain other variables, such as visibility, are not relevant to vegetation. However, additional soil variables were included. The vegetation model is a classification model, so no background points were necessary. The Random Forest output predicts the most likely vegetation type. To ensure that the lakes and rivers were as accurate as possible, these were inserted from the source data for the historic hydrographic model (Public Land Survey plat map and some modern lakes and rivers). This left areas classified by the model as lakes or rivers that were not mapped as such by surveyors. These tend to occur on floodplains, in reservoir basins, and in depressions. We classified these as 'wet land.' They may have been occupied by standing or intermittent water in the past but were not mapped as such on the PLS plat maps.

At best, the MnModel Phase 4 historic vegetation model is an approximation of potential natural vegetation at the time of the Public Land Survey (Figure 7). It is limited by the surveyors' vocabulary and familiarity with Minnesota vegetation, by the environmental data used to develop the predictive variables, and by the lack of information about historic disturbances, particularly fire, that helped shape vegetation patterns. Ultimately, 36 vegetation types are included in the model. Accuracy is highest for the dominant vegetation types and becomes very low for rare vegetation types. In southwestern Minnesota, trees were so rare that all woody vegetation types (floodplain forest and oak savanna, for example), though present, were dropped from the model altogether. Likewise, shrub swamps were infrequently mentioned by surveyors, so they do not appear in the model.

**Figure 7: MnModel Phase 4 Historic Vegetation Model: Vegetation Types**



| VEGETATION TYPE | HARDWOOD SWAMP | UPLAND WHITE CEDAR FOREST | ASPEN FOREST | OAK SAVANNA |
|---|---|---|---|---|
| LAKE | WET MEADOW/FEN | PINE BARRENS | ASPEN-BIRCH FOREST | ASPEN WOODLAND |
| WETLAND | PINE FOREST | JACK PINE WOODLAND | PAPER BIRCH FOREST | OAK WOODLAND-BRUSHLAND |
| RIVER | JACK PINE FOREST | NORTHERN CONIFER WOODLAND | LOWLAND HARDWOOD FOREST | BRUSH-PRAIRIE |
| BOG | RED PINE FOREST | BOREAL HARDWOOD-CONIFER FOREST | MAPLE-BASSWOOD FOREST | PRAIRIE |
| CONIFER SWAMP | WHITE PINE FOREST | MIXED PINE-HARDWOOD FOREST | NORTHERN HARDWOOD FOREST | |
| MARSH | SPRUCE-FIR FOREST | NORTHERN HARDWOOD-CONIFER FOREST | OAK FOREST | |
| FLOODPLAIN FOREST | BLACK SPRUCE-FEATHERMOSS FOREST | WHITE PINE-HARDWOOD FOREST | ASPEN OPENINGS | |

0    40    80         160  Kilometers

N

## Hydrographic Model

The first MnModel hydrographic modeling procedures were developed in 2008 (Stark et al. 2008). This model utilized modern hydrographic data (NWI), digital soils data (pre-gSSURGO), and LfSA data to model locations of historic and prehistoric surface water. The modelers also experimented with surface water polygons digitized from GLO maps and with restorable wetlands inventory data from MnDNR, but neither of these was widely available at the time. The shortcomings of the models produced by these procedures were that there was no way to distinguish types of water bodies or to distinguish between historic and prehistoric water bodies. Moreover, the tools developed would not run with gSSURGO attribute tables.

We updated the hydrographic model in 2018 (Hobbs et al. 2019a) to take advantage of the newly digitized data from the Public Land Survey maps, updated digital soils data (gSSURGO), statewide geomorphic data from the MnModel Phase 4 Landscape Model, and wetland distributions as modeled by the MnModel Phase 4 historic vegetation model. The updated procedures produce two models, one for historic hydrography and one for prehistoric hydrography. The available data are not sufficient to make finer distinctions between time periods.

The historic hydrographic model depicts surface water features (lakes, rivers, wetlands, floodplains) at the time of the Public Land Survey. Lakes are pulled from two sources: those lakes mapped on the Public Lands Survey plat maps and additional lakes from modern sources. Lakes from modern sources were used only if they did not touch section lines, since any historic lake crossing or near a section line would have been observed and mapped by the surveyors. Rivers were also taken from both PLS plat maps and modern sources. A few large rivers on the plat maps were 'meandered' or surveyed, and these were always used. Smaller rivers were often drawn on the plat maps in a generalized way and only approximated the river's course. Modern rivers were substituted for these as long as they had not been dammed, straightened, or diverted. Modern floodplains were mapped only where wetlands were not otherwise present. It was the opinion of the project geomorphologist that modern floodplains were likely also floodplains historically.

Wetland polygons on the historic plat maps are not accurately mapped. Locations where wetlands are crossed by section lines, however, can be considered accurate. Modern wetlands do not correspond well with the wetlands observed by surveyors, either spatially or by type. We compared the efficacy of the vegetation model and modern hydrographic data for predicting a set of wetland points taken from the PLS line data in the BGWD modeling region. These points were on the margins of the region and had not been used to develop the vegetation model. The historic vegetation model was much better at predicting the wetland types identified by the surveyors at these points than were the modern hydrographic data (Table 3). For this reason, we decided to use the vegetation model wetlands for the historic hydrographic model.

**Table 3: Prediction of Historic Wetland Test Points by Vegetation Model and Modern Hydrographic Data**

| PLS Point Value | Floodplain Forest | Lake | Marsh | Meadow/Fen | River | Swamp |
|---|---|---|---|---|---|---|
| Predicted by Vegetation Model | 75% | 89% | 75% | 67% | 82% | 74% |
| Predicted by Modern Data | 20% | 54% | 32% | 1% | 11% | 14% |

The prehistoric hydrographic model depicts surface water features over a very long time period.  Because the level of uncertainty is so much greater for the prehistoric period, only four categories of features are mapped: wetlands, shores, lakes, and floodplains.  Lakes from the historic hydrographic model were augmented with lake beds identified from geomorphic and soils sources.  No attempt was made to map prehistoric river channels.  In addition to including historic floodplains, terraces, alluvial fans, and other floodplain features identified by the project geomorphologist were identified and mapped as prehistoric floodplains.  Wetland locations are mapped based on soil characteristics, as defined by the original hydrographic model (Stark et al. 2008).  A number of soil Great Groups not present in the original data were added to allow wetlands to be mapped in parts of the state not included in the original research.  Where soils data were absent, we used wetland areas from the historic hydrographic model.  Shores were mapped from soils and geomorphic data.  These were included because they were occasionally present without an adjacent wetland soil or lake bed.

The historic and prehistoric hydrographic models bracket the span of years when Minnesota was peopled, but not yet settled by Euro-Americans.  The prehistoric model illustrates the vast extent of water at the end of the glacial era, including the floodplains of glacial meltwater rivers.  The historic model indicates much more extensive wetlands than at present, at least in the agricultural and developed parts of the state, but much less than the prehistoric model (Figure 8).  One might think of the prehistoric model as a 'maximum extent' model and of the historic model as a 'minimum extent' model.  The hydrographic conditions at any particular point in time would be found somewhere in between.
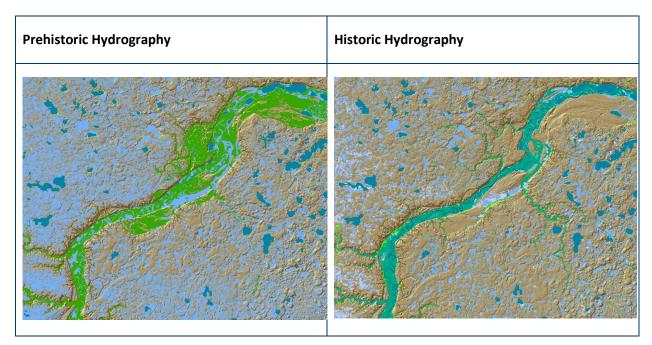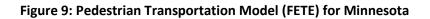
**Figure 8: Comparison of Extent of Surface Water in Prehistoric and Historic Hydrographic Models**

| Prehistoric Hydrography | Historic Hydrography |
|---|---|
|  |  |

## Pedestrian Transportation Model

Prehistoric transportation in Minnesota was by foot and by water. Pedestrian transportation routes can be predicted using least cost-path analyses supported by terrain data and resistance values based on obstacles such as vegetation and water bodies. Devin White (Sandia Ridge National Laboratory) developed a pedestrian transportation model for Minnesota for this project. This 'From Everywhere To Everywhere' (FETE) model is based on White and Barber's (2012) models of pedestrian transportation networks in Mexico. Values in the model indicate the number of least-cost paths that cross each cell when least-cost paths are calculated from every cell to every other cell. Resistance values for the model were calculated based on a 30 m resolution DTM and Marschner's (1974) presettlement vegetation map. Water bodies were assumed to be barriers.

The model shows a dense network of potential transportation routes, with red colors signifying the most-traveled routes and green the least-traveled routes (Figure 9). One of the most notable features of this model is the convergence of multiple long-distance paths in two notable locations: Fort Snelling and Traverse des Sioux. The site of Fort Snelling, at the confluence of the Mississippi and Minnesota Rivers, was known as 'Bdote' to the Dakota and featured in their creation stories. Traverse des Sioux was a major river crossing for voyageurs; Native Americans used the same ford prehistorically. This striking coincidence underscores the importance of terrain for shaping travel and trade.

**Figure 9: Pedestrian Transportation Model (FETE) for Minnesota**



Fort Snelling and Traverse des Sioux

0    40    80         160  Kilometers

N

**Soils Data**

Soils data (gSSURGO) are available for most of Minnesota from the Natural Resources Conservation Service (NRCS). Even where soils data are present, there are many gaps in coverage. These include missing variable values within water bodies, disturbed areas (e.g. gravel pits or mines), and urban areas. Some variables simply were not reported for all map units. In some cases, missing data can be extracted from map unit names or other text fields. We were able to find additional information for soil Great Groups and disturbance factors in this way.

Most soil variables are not associated directly with soil map units. This complicates the use of the gSSURGO data. For example, a number of the variables of interest are recorded by soil components. Each mapunit typically has multiple components. Since one value reported is the percentage of each component within the mapunit, we were able to develop procedures for extracting the most common or dominant value of a variable from each mapunit's components (Brown et al. 2019). Where variables are reported by soil horizons, we used the value for the surface horizon. After evaluating the available variables, we selected only those that appeared to have adequate geographic coverage (Hobbs, Walsh, and Hudak 2019).

We supplemented the gSSURGO data with drainage and productivity indices provided by Michigan State University (Schaetzl et al. 2009). Because the drainage index (DI) is primarily derived from the soil's taxonomic characteristics, it may provide different information than the gSSURGO drainage variable per se. Likewise, the productivity index (PI) indicates the relative amount of nutrients available to plants.

# Variable Derivation

Considerable effort was expended to collect data and develop models of historic and prehistoric terrain, landforms, hydrography, and vegetation to support modeling of prehistoric archaeological sites. Before modeling, however, these data and models must be transformed into variables thought to be associated with prehistoric land use, resource acquisition, or settlement patterns.

For example, the digital terrain model records elevation, and elevation may be a significant correlate with site location. But other variables can be derived from elevation, and these may further elucidate aspects of site distributions. These include slope, aspect, surface curvature, relative elevation, shelter index, topographic position index, topographic wetness index, and visibility.

The locations of lakes, rivers, and wetlands are recorded in the hydrographic models, but for site location it is more important to know the distance of a site to the nearest hydrographic feature, as no one wants to carry water very far. Thus most hydrographic variables are based on least-cost path distances to various types of water bodies.

The local type of vegetation may be important for the food, fuel, and material resources it can provide, but the diversity of vegetation in proximity to a site tells us about the full range of resources available. Likewise, proximity to transportation, whether by land or by water, is an advantage. For Phase 4, we included measures of distance to three levels of pedestrian pathways from the FETE model. Finally, soil drainage, flooding, hydric characteristics, and productivity may influence where people choose to live.

The variables developed for MnModel Phase 4 are listed in Appendix A of this report. Their derivation and measures of their performance in the models are fully documented in Hobbs, Walsh and Hudak (2019).

## Sampling

Sampling refers to the association of the environmental variables with the locations of both archaeological sites (or surveys) and other locations (background points) that are not known to be sites or surveys. This creates a database that can be analyzed to determine how site/survey and non-site/non-survey locations differ with respect to the environmental factors measured.

Sampling data represented as points is simply a matter of overlaying each point on the environmental variables and recording the single variable value at the point's location. However, if the point represents something large, such as an archaeological site or survey polygon, the data value at one point may not be representative of the entire polygon. By converting polygons to rasters, it would be possible to sample the value of each variable at each cell within each polygon, but this would introduce problems of spatial autocorrelation into the analysis.

To better represent site and survey polygons, sampling techniques were developed to capture and summarize information about the entire polygon and attach those data to a single point for each polygon (Hobbs et al. 2019b). Entire site polygons were summarized to a single point, but survey polygons larger than 4 km$^2$ were segmented to achieve a point sampling interval of 2,000 m. Categorical variables and numeric variables with a limited number of possible values are represented by the value associated with the majority of 30 m cells within the polygon. Numeric variables with many possible values are represented by the mean value of all of the cells within the polygon. The majority and mean values calculated for the polygons are attached to the polygon's centroid point for representation in the modeling database.

Procedures were developed in ArcGIS to generate a grid of background points at locations at least 1500 m distant from where sites or surveys were present. The ratio of background points to sites was between 2:1 and 3:1. The ratio of background points to surveys as between 1.5:1 and 2.5:1. Background points were not associated with polygons. Variable values at background points were simply the value at the point itself. The background points were combined with the site or survey centroids to create regional sample populations for analysis and modeling. These were exported from ArcGIS as a .csv file.
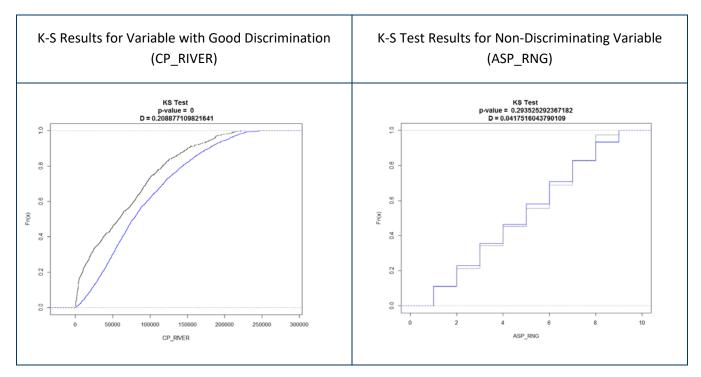
Finally, a point feature class (the 'prediction points') was created from every cell in a 30 m resolution raster of each region and used to sample all of the environmental variables. X,Y coordinates were added to this point feature class, and the table was exported in .csv format. This dataset is used in R to calculate the final model values for the region, which can then be imported back into ArcGIS as a raster. It should be noted that each modeled region includes a 10 km buffer so that measurement of path distance from locations on the region's boundary can include destination features outside the region if they are closest.

## Exploratory Data Analysis

MnModel Phase 4 statistical procedures include several steps to describe the data and to clean up the database by removing records and variables that may cause problems (Landrum et al. 2019). R provides tables of standard descriptive statistics for each variable for both site/survey points and background points and produces

histograms of the values of each variable for sites/surveys and background points. It then performs analysis providing useful information for making decisions about cleaning up the data. These include:

- Identifying and removing NULL values, either by removing the records containing NULL values or, if there are too many of these records, by removing the variables that contain the NULLs.

- Evaluating categorical variables to determine whether the values present are the same in both the training and predictive data. Any values in the predictive data that are not found in the training data must be reclassified, as R will not be able to run a model on predictive data with unexpected categorical values. At the same time, values with very low numbers of records must also be reclassified to ensure that these rare values do not cause problems.

- Evaluating differences in variable distributions between sites/surveys and background points using a Kolmogorov-Smirnov Test (K-S Test) and removing variables which fail to distinguish between the cases. This test also produces graphs of the distribution of each variable's values for sites/surveys and background points (Figure 10).

- Identifying variables that display near-zero variance and removing these from the data.

- Listing variables with skewed distributions. No action was taken regarding these variables as 'tree' methods such as Random Forest perform well with skewed data. This is fortunate, as a very large number of MnModel variables have skewed distributions. However, this could be useful if other modeling procedures are used.

- Identifying correlated pairs of variables using Spearman's Rank Correlation test.

- Calculating the Variable Inflation Factor (VIF) to help make decisions about which of the correlated variables should be removed from the dataset.

- Running a Chi-Squared Test of Independence on the categorical variables. Again, we did not act on the result of these tests, which showed that our categorical variables are nearly always related.

**Figure 10: Example of K-S Test Plots**



| K-S Results for Variable with Good Discrimination (CP_RIVER) | K-S Test Results for Non-Discriminating Variable (ASP_RNG) |
|---|---|

# Statistical Modeling

Efforts to update modeling procedures began in 2007 with an evaluation of several new statistical techniques (Oehlert and Shea 2007). We adapted these procedures in 2018, migrating from the commercial S-Plus statistical software to R, an open-source statistical programming language. Oehlert and Shea (2007) recommended using a tree classification technique called Bootstrap Aggregating or 'Bagging.' The procedure adopted in Phase 4 is 'Random Forests' (Breiman 2001), a newer procedure based on Bagging. The selection of Random Forest is consistent with Oehlert and Shea's (2007) recommendations and follows its successful use to develop Pennsylvania's archaeological predictive model (Harris, Kingsley, and Sewell 2015). The Random Forest procedure creates multiple tree models for the data set and calculates predictions based on the results of all the trees.

The statistical procedures developed include procedures for 'tuning' the models. Random Forest models are optimized in R by determining the optimum number of trees to build (ntree value) and the optimum number of variables to examine at each node (mtry). Finally, they provide measures of model performance. The optimized model is evaluated by how well it predicts a test population of site and background points not used to build the model. Finally, the optimized model is applied to the dataset representing a 30 m grid of points that have been used to sample all of the environmental variables. As this dataset contains x,y coordinates for each point, it can be exported from R in .csv format and imported into ArcGIS as a raster dataset (ArcGIS grid format).

Because statistical procedures cannot handle missing data, we developed procedures for two sets of models: one without soils variables (Model A) and one with soils variables (Model B). Model A uses all sample points and all variables except soil variables. Model B uses soil variables, but only the sample points that did not have NULL

values for any soil variables.  As a consequence, Model B always has NULL values where soils data are missing. To complete Model B, Model A values are used to replace the NULL values, creating a composite model (Model C).  Models A and C are then evaluated to select the best model for each region.

In the first round of Phase 4 modeling, models were run on a 75 percent subset of the sample population (training data) and tested on the other 25 percent.  Two site models (one Model A and one Model B) and two survey models (one Model A and one Model B) were run for each of the twenty modeling regions.  These models were exported to ArcGIS and composite models (Model C) were created.  Models A and C were then evaluated for their ability to predict the entire sample population.  Both site and survey models performed extremely well - much better than the Phase 3 models. The site models demonstrated that it was possible to select a threshold between high and low site probability that captured 95 percent of the sample population in the high probability category.  However, when evaluated with a test population of new sites, performance fell short of our goals in most regions.

A second round of modeling was used to take advantage of the additional site data and to refine modeling procedures.  These final Phase 4 models were developed using jackknife procedures that provided us with improved measures of the models' abilities to predict locations of sites not used to build the models.  Whereas the original procedures built only one model of each type using 75 percent of the data and tested using the other 25 percent, the final models built four models of each type.  Each model run was randomized to use a different 75 percent training sample and 25 percent test sample.  The model evaluation measures for the four models can then be examined to better understand how stable the models are across different subsets of the data (see Figures 16 and 17 below).  Finally, a model is built using the entire database.  This model cannot be evaluated the same way, since there is no test population, but we can assume that it is at least as strong as the four previous models.  It may be stronger, since it is based on more site information.  This is particularly important in regions where site populations are small.

Since we had no additional survey information after completing the first round of modeling, and since survey populations are already relatively large, we did not develop any additional survey models.  When more surveys are digitized, it would be worthwhile to update the Phase 4 survey models using the refined modeling procedures.

## Model Classification

The raster models produced by these procedures contain floating point values between 0 and 1.  To be useful, they must be classified into high probability and low probability areas.  In Phase 3, we divided models into high, medium, and low probability classes.  For Phase 4 we decided to simplify the classification and use only high and low probability.  This was possible because the models are much more accurate and precise than the Phase 3 models (see below).

The R procedures include determination of an 'optimum' threshold for making this distinction.  We must keep in mind, however, that this threshold is determined based on the performance of only 25 percent of the database. It may not provide the best threshold when the database as a whole is considered.  In the first round of modeling, models were classified using this cutoff and then evaluated.  These will be referred to as the 'Maximum Accuracy' models.

In Phase 3, the threshold value was selected to insure that 85 percent of all sites in the database were within the high/medium probability areas (Hobbs et al. 2002). Thus, model sensitivity was held more or less constant across the state and the area classified as high/medium probability varied. For round one of Phase 4, we evaluated thresholds to insure capture of 85 percent, 90 percent, and 95 percent of the total site sample in the high probability areas. We will refer to these models as the 'Target Sensitivity Models.' Keep in mind that the Maximum Accuracy and Target Sensitivity models are merely different classification rules applied to the same numeric models.

Maximum Accuracy cutoffs and Target Sensitivity cutoffs were applied to Site Models A and B to create classified raster models. At this point, there were four versions of Site Model A (Maximum Accuracy, 85 percent sensitivity, 90 percent sensitivity, and 95 percent sensitivity) and the same four versions of Site Model B. Only Maximum Accuracy cutoffs were applied to Survey Models A and B. We then created composite models (Model C) from models A and B. Because every Site Model B has some cells with NULL values where soils data are missing, those values must be replaced from values in Site Model A. Composites of each version of the classified Target Sensitivity models were created. Composites of Maximum Accuracy Survey Models A and B were also created.

After evaluating the first models, we determined that the 95 percent Target Sensitivity site models were the best choice for use by archaeologists. This is the only classification we applied to the final site models.

# Model Evaluation

Model performance was evaluated in both R and ArcGIS. Statistical procedures in R measure the ability of each model built from 75 percent of the sample to predict the other 25 percent. However, there are a number of important measures that cannot be generated in R. These include evaluations of the portion of land area classified as high probability and the number of sites, surveys, and background points in the total population that fall into the high probability category.

### Evaluation Measures in R

The R statistical procedures report a number of performance measures. These are the best measures to consult to determine the model's ability to predict as yet undiscovered sites. If the models are stable, we would expect the results to be relatively consistent between different models created for the same region. Variations in these measures may indicate the presence of outliers in the data. Small numbers of sites in locations that are quite different from the majority of sites in the sample may not be well predicted when they fall into the test population. This is more likely to happen when there is a small overall site population, as the outliers can then have a greater impact.

For each model, R generates graphs of variable importance (Figure 11). These plots indicate which variables are most responsible for model results. The plot on the left indicates the expected percent increase in the Mean Square Error (MSE) of the model if each variable is removed. The plot on the right plots the variables contribution to 'node purity', which relates to the loss function by which best splits are chosen. More useful variables achieve higher increases in node purity. In Figure 11, the variable *Path Distance to Nearest Historic Wetland* (CP_WETLAND) is by far the most important variable in the model, followed by *Path Distance to*

*Nearest Large Historic Lake* (CP_LLK), *Elevation* (ELEV), and *Topographic Position Index within 250 Meters* (TPI250).  The variable important plots are generated for all models and do not depend on an evaluation of the model's ability to predict the test population.  R also provides tables with the numeric values used to create these plots.

**Figure 11: Example of Variable Importance Plots**



Models based on 75 percent training samples are evaluated using a Confusion Matrix.  This is a table that displays the distribution of site/survey points and background points between the high and low probability classes of the model (Table 4).  A 'threshold' value is required for this matrix to be generated.  That is, the user must specify a model value that represents the dividing point between high and low probability.  The default value is 0.5 (possible model values range from zero to one).  In the preliminary round of models, all models were evaluated using the default value, as is the example in Table 4.  In the second round of modeling, models were evaluated using the 'maximum accuracy cutoff' value calculated by R.
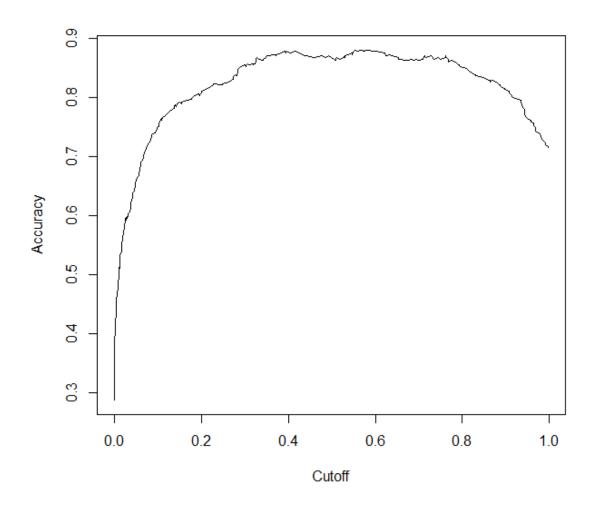
**Table 4: Example of a Confusion Matrix**

| Test Population | Actual Sites | Actual Background Points |
|---|---|---|
| Predicted To Be Sites | 84 | 25 |
| Predicted To Be Non-Sites | 27 | 290 |

R reports a number of measures of model performance based on the Confusion Matrix results. These include:
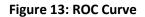
- Accuracy: This is the proportion of both sites and background points correctly predicted. In the example above, accuracy is 0.8653. R also reports a 95 percent confidence interval around the accuracy value, indicating the likely minimum and maximum accuracy estimates. The narrower this confidence interval, the more confident the accuracy estimate. In this example the accuracy reported may not be the best achievable by the model since the threshold between high and low probability was determined by the default value (0.5) rather than the 'maximum accuracy threshold' value. After constructing this matrix, R suggested a 'maximum accuracy threshold' value of 0.5803333, suggesting that using this cutoff would achieve a model accuracy of 0.8808290. The relationship between model accuracy and the threshold or cutoff value used for model classification is illustrated in Figure 12.

- No Information Rate (NIR): The 'best guess' of the chance that sites are absent given no information beyond the overall distribution of sites and non-sites. It is calculated as the portion of the sample that is background points. Ideally, the NIR should be lower than the accuracy estimate. This estimate would be more accurate if our sample was based on a random survey, but it is not. Our sampling procedures try to achieve a ratio of 2:1 to 3:1 between background points and site or survey points. A high ratio of background points would likely result in the background points swamping the model. In reality, though, the portion of the landscape occupied by sites and surveys is very low, so the NIR reported by R is over-estimated. A better estimate for us might be the fraction of the surveyed area in which no sites were found.

- Cohen's Kappa: Kappa measures how well the model performed compared to how well it would have performed by chance. Kappa should be high if there is a large difference between accuracy and the NIR. Again, this measure depends on the ratio between sites and background points within the sample to estimate how likely it would be to find a site by chance. It is therefore likely that the estimates of Kappa provided are under-estimated.

- Sensitivity: Also known as the True Positive Rate (TPR), sensitivity is simply the proportion of sites in the test population that are correctly predicted. This is the most important measure for us, as failing to predict that sites may be present can be very costly. In the example above, sensitivity is 0.7568, which is lower than the sensitivity that could be achieved with a different threshold value.
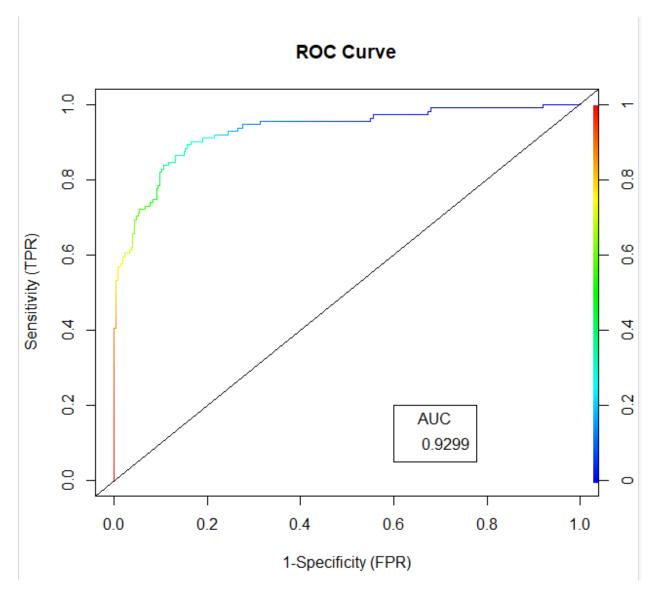
- Specificity: Also known as the True Negative Rate (TNR), specificity is the proportion of background points correctly predicted. In the example above, specificity is 0.9091. That this value is much higher than sensitivity reflects the large number of background points in the sample population compared to sites. The model is much better at predicting non-site locations than site locations, at least at the threshold selected. This illustrates the possible consequences of increasing the number of background points even further.

- Positive Predictive Value (PPV): Also known as model precision, this measures the portion of locations the rule classifies as sites that actually contain sites. In the example above, the PPV is 0.7706.

- Negative Predictive Value (NPV): This is the fraction of locations the rule classifies as non-sites that are non-sites. The NPV of the example model is 0.9025.

- Prevalence: Presumably, prevalence is the fraction of the landscape that actually contains sites. In R, this is calculated as the True Positives plus the False Negatives divided by the Total Sample. Once again, since we use a set ratio of background points to sites, our sample is not reflective of what one would find doing a random archaeological survey.

- Detection Rate: The fraction of true positives in the total sample. If our sample was based on a random archaeological survey, it would be the chance that we would find archaeological sites whenever we surveyed. Since we have a set ratio of sites to non-sites, however, it is an over-estimate.

- Detection Prevalence: The fraction of true positives plus false positives in the sample. If we were using a random sample, this would essentially be the fraction of the landscape classified as high site or survey potential.

- Balanced Accuracy: An accuracy measure that is used if the classes are not balanced (equal numbers of sites and non-sites), as is the case in our data. It is calculated as the sum of sensitivity and specificity divided by 2. In the example above, balanced accuracy is 0.8329.

**Figure 12: Relationship between Model Accuracy and the Cutoff Value Selected for Classification**



ROC Curves provide another way to evaluate model performance (Figure 13). These curves graph the relationship between model sensitivity (the True Positive Rate or TPR) and the False Positive Rate (FPR), which is the inverse of specificity. In the ROC curve, TPR and FPR are graphed at a full range of cutoff values. Ideally, we select a cutoff where the TPR is high and the FPR is low. However, as we achieve higher TPR values, FPR values invariably rise. The straight diagonal line on the graph represents where the curve would be by chance alone. We want our model to be graphed to the left of this diagonal. In a perfect model, the curve would go straight up the left-hand site of the graph, with a TPR of 1.0 and FPR of 0. A simple measure of the model's performance is the area under the curve (AUC). The total area of the box is 1.0. In a perfect model, the AUC would be 1.0. In the random model, the AUC is 0.5. We want our AUC values to be as close as possible to 1.0.

**Figure 13: ROC Curve**



## Evaluation Measures in ArcGIS

As noted above, the ratio between our site or survey points and background points is not representative of the landscape as a whole. For this reason, some of the evaluation measures reported by R are not useful for us. Moreover, we cannot use R to estimate the portion of the landscape that would be occupied by high probability in our models. Also, R evaluates only Model A and Model B. Because Model C is a composite constructed in ArcGIS, we must evaluate it in ArcGIS. Finally, the evaluation in R depends on a test population. That prevents us from using R to evaluate models run on the total sample. Yet models trained on the full sample should perform better than models trained on only 75 percent of the sample. This is particularly true in regions where site numbers are low to begin with. For these reasons, we perform additional evaluations of the models after they are exported to ArcGIS.

We first use standard sampling procedures to extract raw (floating point) model values from the model rasters to the sample points (site or survey points and background points).  We use the point attribute tables in ArcGIS to determine the cutoff values for each of the target sensitivity models and create the classified models.  We then repeat the sampling procedures to attach high and low probability classes from the classified models to the same sample points  After recording numbers of sites/surveys, background points, and raster cells in each category in a spreadsheet, we are able to calculate several performance measures.  These measures illustrate how well the models perform with respect to the entire dataset from which they were derived.  This is not the same as their ability to predict undiscovered sites.  The simple performance measures calculated are:

- Sensitivity:  percent of sites falling in high probability areas (True Positives/True Positives + False Negatives)

- Specificity: percent of background points falling in low probability areas (True Negatives/False Positives + True Negatives)

- Overall Accuracy: (True Positives + True Negatives)/Total Sample

- Balanced Accuracy: (Sensitivity + Specificity/2)

- Extent of high probability: percent of raster cells classified as high probability

- Extent of low probability: 1 – extent of high probability

- GAIN: 1 – (% of region in high probability/% of sites in high probability) (Kvamme 1988)

- Error Rate (Misclassification Rate):  (False Positives + False Negatives)/Total Sample

In addition, we can use GIS to estimate some of the performance measures that are unreliable when obtained from the sample data alone.  This is because we can estimate Prevalence, the *a priori* chance of finding a site, by overlaying site polygons with survey polygons.  The fraction of the surveyed area occupied by sites is an estimate of the *a priori* chance of finding a site, and the fraction of the surveyed area lacking sites is an estimate of the No Information Rate (NIR).  Given these estimates, we can calculate estimates for Kappa, PPV, and NPV that are more realistic than possible from the site and background point data alone.  Except for the estimates of Prevalence and NIR, all equations are based on on-line documentation of how R calculates the same measures.

- Prevalence: fraction of the landscape that actually contains sites, also known as the *a priori* probability of finding a site (area of site polygons within survey polygons/area of survey polygons)

- No Information Rate (NIR): fraction of the landscape that actually does not contain sites (1 – Prevalence)

- Observed Proportionate Agreement (Po): Same as Overall Accuracy

- Probability of Random Agreement on Site Presence (P1): Prevalence * Extent of high probability

- Probability of Random Agreement on Site Absence (P0): NIR * Extent of low probability

- Overall Random Agreement Probability (Pe): P1 + P0

- Kappa: a measure of how well the model performed compared to how well it would have performed by chance [(Po-Pe)/(1-Pe)].

- Positive Predictive Value (PPV): Fraction of locations the rule classifies as sites that actually contain sites [(Sensitivity * Prevalence)/((Sensitivity * Prevalence) + ((1 – Specificity) * (1 – Prevalence)))]

- Negative Predictive Value (NPV): Fraction of locations the rule classifies as non-sites that are non-sites [(Specificity * (1-Prevalence))/(((1 – Sensitivity)* Prevalence) + ((Specificity )* (1-Prevalence)))]

- Unexpected Discovery Rate (UDR): The fraction of times we find a site where we don't expect it. Also known as the 'oops' rate. (1-NPV)

- Positive Predictive Gain (PPG): How many times better the model is at discovering sites than a random survey would be (PPV/Prevalence)

- Negative Predictive Gain (NPG): How much less likely we are to discover a site at a location labeled a non-site using the model than if we were surveying randomly. (UDR/Prevalence)

- Detection Rate: Percentage of the total sample accurately predicted as positive (True Positives/Total Sample)

- Detection Prevalence: Percentage of the total sample predicted to be positive (True Positives + False Positives/Total Sample)

- Precision: Proportion of positive predictions that are true (True Positives/True Positives + False Positives)

## First vs. Second Round of Modeling

We initially evaluated models classified using the suggested 'maximum accuracy threshold' as well as models that met our criteria for targeted sensitivity levels. Model A versions (no soils data) are compared to their corresponding Model C version (composite models partially based on soils data) to determine which performed better. For site models, we decided to select the best model from each pair based on sensitivity, since it is most important to not miss sites if they are present. If sensitivity values were the same we looked at the model's specificity, selecting the model that minimized false positives.

In the first round of modeling, we expected that the 95 percent Target Sensitivity site models might have unacceptably low specificity values. This was not the case. Specificity values of the best 95 percent site models were never lower than 0.93. Thus it was clear that the 95 percent Target Sensitivity site models optimized both specificity and sensitivity and should be the models to release for use. Of these, we selected thirteen regional site models created without soils data (Model A) and seven created with soils data (Model C).

The same procedures were followed to evaluate the Maximum Accuracy survey models. The best survey model was deemed to be the one that maximized specificity, since identifying a location as likely to have been surveyed and not surveying it is potentially more costly than re-surveying an area that has been surveyed. Because the

Maximum Accuracy models performed so well on this measure, with specificity values ranging from 0.97 to 1.0, there was no reason to create Target Specificity survey models.

After completing these models we discovered a population of 766 sites that had not been used to build the initial models. We evaluated how well the models predicted these sites. These results, discussed in the Model Performance section below, prompted us to refine our modeling procedures to develop a second set of site models. In the second round of modeling, four preliminary models were run from four different 75 percent subsets of the database and tested in R with the other 25 percent of the data. The 'maximum accuracy threshold' was substituted for the default threshold value (0.50) prior to running model evaluation procedures in R to achieve more realistic measures of how the final models would perform. The results of the four models were compared to assess the magnitude and range of model predictions and to gauge model stability. We then ran a final model with the entire dataset. Because there was no test population available for this model, we performed all model evaluation outside of R using the entire database. We expect that, in terms of predictive power, the final model would perform at least as well and probably better than the four preliminary models. However, the measures assessing the final model do not assess its ability to predict new sites as much as its ability to describe the site population from which it was created.

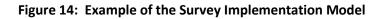## Create and Classify the Survey Implementation Model

As a final step, the best 95 percent Target Sensitivity site models were combined with the best Maximum Accuracy survey models to create the survey implementation models. Survey implementation models are classified into four categories (Figure 14):
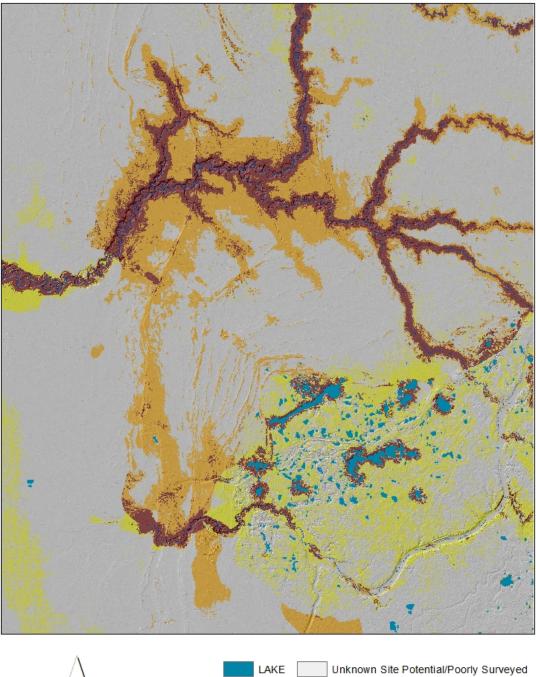
- Unknown Site Potential/Poorly Surveyed
- Low Site Potential/Well Surveyed
- High Site Potential/Poorly Surveyed
- High Site Potential/Well Surveyed

The 'Unknown Site Potential' category covers areas where site potential is 'low' but, because these areas are poorly surveyed, we have no way of knowing whether sites are absent because the area is unsuitable for sites or simply because no one has looked for sites there.

By comparison, Phase 3 models had nine classes, which could be confusing. For archaeologists familiar with the Phase 3 models, the following correspondences may be helpful:

- Phase 4 'High Site Potential/Well Surveyed' category corresponds to Phase 3 'High', 'Medium', 'Possibly High', and 'Possibly Medium' categories.
- Phase 4 'High Site Potential/Poorly Surveyed' category corresponds to Phase 3 'Suspected High' category.
- Phase 4 'Low Site Potential/Well Surveyed' category corresponds to Phase 3 'Low' and 'Possibly Low' categories
- Phase 4 'Unknown Site Potential/Poorly Surveyed' category corresponds to Phase 3 'Unknown' category.

**Figure 14: Example of the Survey Implementation Model**



Legend:
- LAKE
- RIVER
- Unknown Site Potential/Poorly Surveyed
- Low Site Potential/Well Surveyed
- High Site Potential/Poorly Surveyed
- High Site Potential/Well Surveyed

0  3.5  7  14 Kilometers

N

# Model Performance

Model performance for the entire state is summarized in this section.  Comparisons to Phase 3 models are based on the Phase 3 models' performance as evaluated by the Phase 4 data.

## Site Potential Model

### Initial Models

The initial statewide Phase 4 site potential model developed performed exceedingly well on the sample database of 8836 sites.  In all measures, it performed 10-14 percent better than the Phase 3 model (Table 5).  The Phase 4 model had particularly high specificity values – the ability to correctly predict which modeled locations were 'non-sites.'  This model specificity was instrumental in reducing the area classified as high sites potential and elevating both GAIN and overall model accuracy.

**Table 5. First Phase 4 Site Model Performance Evaluated Using Sample Data**

| Performance Measure | Phase 3 | Phase 4 |
|---|---|---|
| Sensitivity (True Positive Rate) | 0.84 | 0.95 |
| Specificity (True Negative Rate) | 0.83 | 0.97 |
| Overall Accuracy | 0.84 | 0.96 |
| % Region High Probability | 0.23 | 0.13 |
| GAIN | 0.73 | 0.86 |

Unfortunately, this model was able to predict only 79 percent of the newly discovered population of 766 sites (test population).  The Phase 3 models predicted 76 percent of these sites, so the improvement in model predictive power (sensitivity) was not impressive.  However, the improvement in specificity, the portion of the landscape classified as high probability, and model GAIN were still much better.

The initial model's ability to predict the test population varied considerably between regions.  However, when the new sites were combined with the original dataset, the percent of all sites in the high site potential area decreased by only one percentage point (Figure 15).  The size of the test population within each region varied from one to 179 sites.  Small numbers of sites clearly skewed test results.  Where only one new site was present, and it was predicted, the model's sensitivity was 1.0.  Where only seven new sites were present and two were

predicted, the sensitivity was 0.29.  This should serve as a cautionary tale for users of the model.  Though the site model may perform very well overall, it will be quite variable with respect to individual surveys.  Some portion of the total site population will always be found in the low and unknown site potential areas.
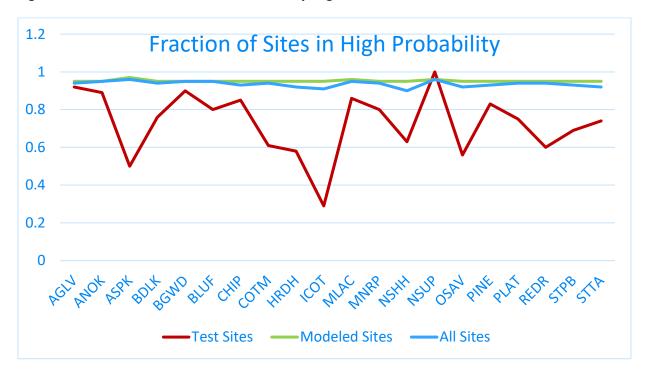
**Figure 15: Performance of Initial Site Models by Region**



## Intermediate Models

After combining the newly discovered sites with the site population used to develop the initial models, regional models were developed with updated procedures.  This section discusses results of the four intermediate models from randomly partitioned data that were evaluated in R for their ability to predict a 25 percent test sample.

Performance of the four intermediate models in each region varied.  Sensitivity, the portion of sites correctly predicted, ranged from 0.536 to 0.96 for models without soils data (Figure 16) and from 0.417 to 0.976 for models with soils data.  Mean sensitivity was 0.80 without soils data and 0.77 with soils data.  Values varied between the four models for each region, particularly within regions with smaller samples (AGLV, ASPK and ICOT in Figure 16).
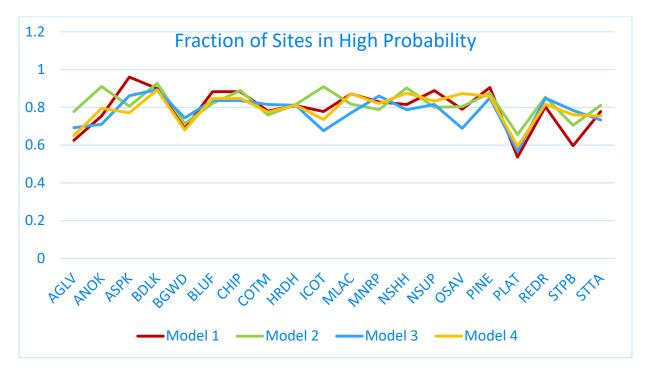
**Figure 16: Sensitivity of All Intermediate Site Models without Soils Data**



Specificity values, the fraction of background points correctly predicted to be non-sites, were uniformly high (Figure 17), ranging between 0.95 and 1 for models without soils data and between 0.87 and 1 for models with soils data. High specificity values are likely a consequence of the relatively high ratio of background points to sites. The model has more information about background point locations and is thus better able to predict them. These high specificity values have a strong effect on overall accuracy of the models, which ranges from 0.877 to 0.979 for the models without soils data and from 0.822 to 0.989 for models with soils data. On all of these measures, the low minimum values for models with soils data are from regions where soils data are absent for large areas resulting in much smaller sample sizes for modeling.

Performance of all initial models is summarized in Table 5. In the end, it shows that the models with and without soils data differ very little, on average, in how well they predict sites.
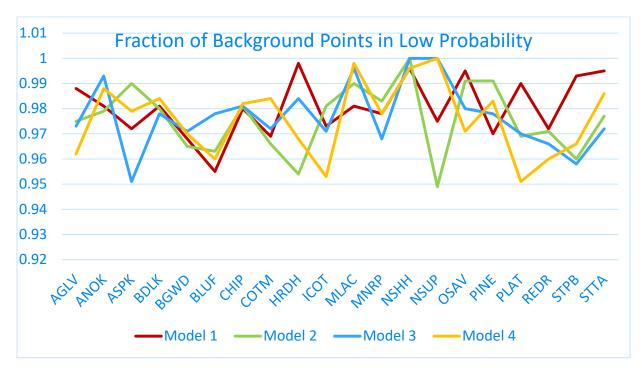
**Figure 17: Specificity of All Intermediate Site Models without Soils Data**



**Table 5: Mean Performance Measures for All Initial Site Models**

| Measure | Model A (No Soils Data) | Model B (Soils Data) |
|---|---|---|
| Sensitivity | 0.80 | 0.77 |
| Specificity | 0.98 | 0.98 |
| Accuracy | 0.94 | 0.94 |
| 95% CI Range | 0.07 | 0.07 |
| ROC AUC | 0.96 | 0.95 |
| Positive Predictive Value | 0.91 | 0.92 |
| Negative Predictive Value | 0.94 | 0.94 |
| Balanced Accuracy | 0.89 | 0.88 |

## Final Models

The final statewide site model (Figure 18) is based on a total population of 9602 sites, including both the sites modeled in the initial stage and the newly discovered sites.
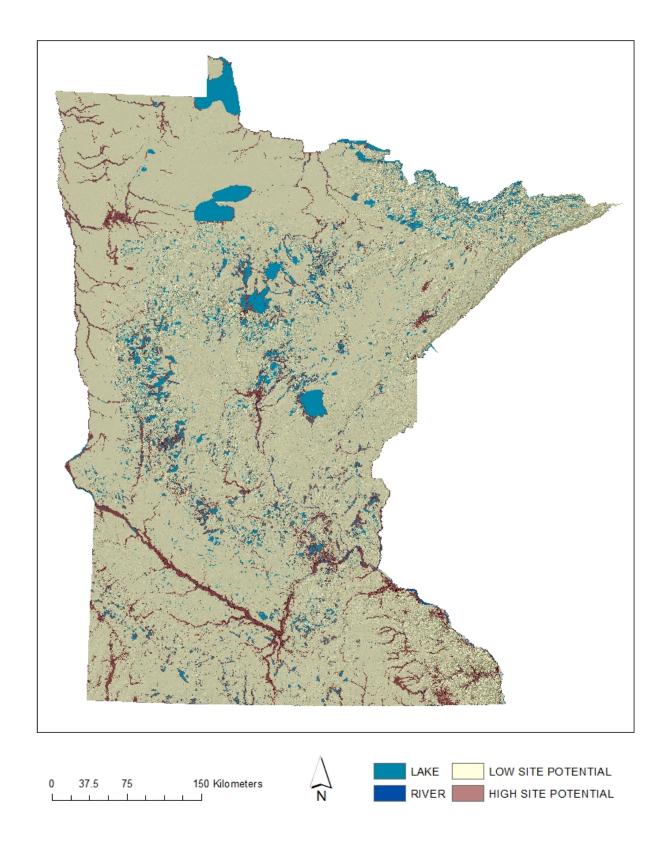
Since the final Phase 4 site models were developed from the entire dataset, their predictive abilities should be at least as good as reported in Figure 16 and are likely to be better. Because the final models are classified using a 95 percent Target Sensitivity threshold, they should be much better. However, without a test dataset we cannot provide a more accurate estimate. The following discussion is based on how well these models represent the entire dataset and how they compare to Phase 3 and Phase 4 initial models.

Table 6 summarizes model performance for the entire state. The measures in the left columns are simple calculations based on overlays of the site and background points on the classified model. Measures on the right incorporate an estimate of prevalence. Prevalence (or *a priori* probability) was calculated in ArcGIS as the fraction of surveyed polygons that contained site polygons. This varies widely between regions, from a minimum of 0.0006 to a maximum of 0.0826. The statewide average (0.0116) is consistent with estimates based on surveys conducted for MnModel Phase 3 (Gibbon et al. 2002, Table 5.3). This prevalence value is then used to estimate the remaining measures. These should be more reliable estimates than those provided by R for the intermediate models.

**Table 6: Performance Measures for Final Statewide Site Potential Model**

| Measure | Statewide Model | Measure | Statewide Model |
|---|---|---|---|
| Sensitivity | 0.95 | Prevalence | 0.0116 |
| Specificity | 0.99 | Kappa | 0.84 |
| Accuracy | 0.98 | Positive Predictive Value | 0.65 |
| Balanced Accuracy | 0.97 | Negative Predictive Value | 1.0 |
| Fraction of Area in High Probability | 0.09 | Unexpected Discovery Rate | 0.01 |
| GAIN | 0.91 | Positive Predictive Gain | 84.35 |
| Precision | 0.98 | Negative Predictive Gain | 1.19 |

**Figure 18: MnModel Phase 4 Site Potential Model**



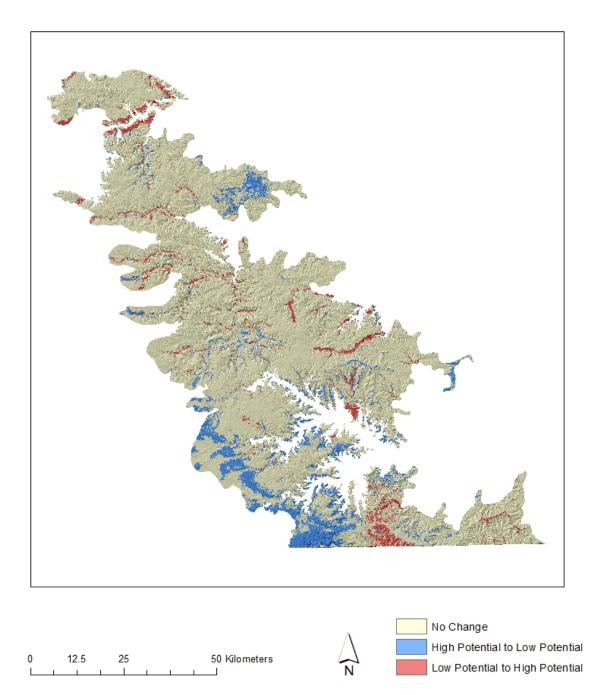| | | |
|---|---|---|
| LAKE | | LOW SITE POTENTIAL |
| RIVER | | HIGH SITE POTENTIAL |

0    37.5    75    150 Kilometers

N

As expected, the map of the final Phase 4 site potential model shows the strong relationship between site locations and sources of fresh water (Figure 18). Yet distance from water alone is insufficient for predicting site distributions. Some water bodies are surrounded by broad bands of high site potential, while others by only very narrow bands, implying the effects of terrain, landforms, vegetation, and accessibility. High site potential is associated with some terrain features, such as escarpments (Figure 19), with resources such as wild rice, with river terraces, and with proximity to transportation corridors. Unlike the Phase 3 models, there are few, if any, artifacts of high site potential associated with modern infrastructure or data errors. More detailed information about the variables used in the models and their performance can be found in Hobbs, Walsh and Hudak (2019).

**Figure 19: High Site Potential Associated with Escarpment**



The final site model is 93 percent identical to the initial statewide site model. With the 8.5 percent increase in the site population and an increase in the ratio of background points to sites, six percent of the state shifted from high probability to low probability while only one percent shifted from low probability to high probability. These changes in the model were not necessarily in proximity to a critical mass of new sites. In the PLAT region, for example, there were only four new sites, a three percent increase in the number of sites modeled. Three of these new sites were already predicted by the initial site model. Presumably, then, there was very little new site information contributing to the final model. Yet considerable changes were observed (Figure 20).

**Figure 20. Changes in Site Potential Classes between the Phase 4 Initial and Final Models for the PLAT Region**

In the PLAT region, shifts from high to low site potential generally occurred on ridges and steep valley sides. Shifts from low to high site potential were similarly on ridges and hillslopes, but in different parts of the region, including around the single new site that was not predicted by the initial model.

If increased site information cannot explain these changes, one must look at the only other variable, the number and distribution of background points. It may be that having more information about non-sites reduces the probability values of a large number of cells without actually increasing the probability of other cells. When the model is classified, these other cells may end up in the high probability category since cells that were previously between them and higher probability cells have been demoted. Certainly the greatest improvement in the final model is in the prediction of non-sites (Table 7). This raises interesting questions about designing the archaeological sample, including the optimum ratio of background points to sites, particularly in regions with low site numbers. There is clearly a trade-off in these models between sensitivity (site prediction) and specificity (non-site prediction).

**Table 7. Comparison of Initial and Final Models with Respect to Prediction of Sites and Background Points**

| Change in Probability Category | No. Sites | % Sites | No. Background Points | % Background Points |
|---|---|---|---|---|
| No Change | 9057 | 94.2 | 30717 | 94.2 |
| High Potential to Low Potential | 205 | 2.1 | 1859 | 5.7 |
| Low Potential to High Potential | 349 | 3.6 | 20 | 0.0006 |

Statewide only 3.6 percent of all sites were found in areas newly classified as high potential, and this is offset by the 2.1 percent of sites that were predicted by the initial model and not by the final model (Table 7). Nevertheless, the very low percentages of sites whose prediction is affected by changes in the model provides further evidence of model stability. Moreover, it suggests that future modeling efforts may not make very large differences in site prediction unless very large numbers of sites are added to the database.
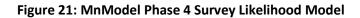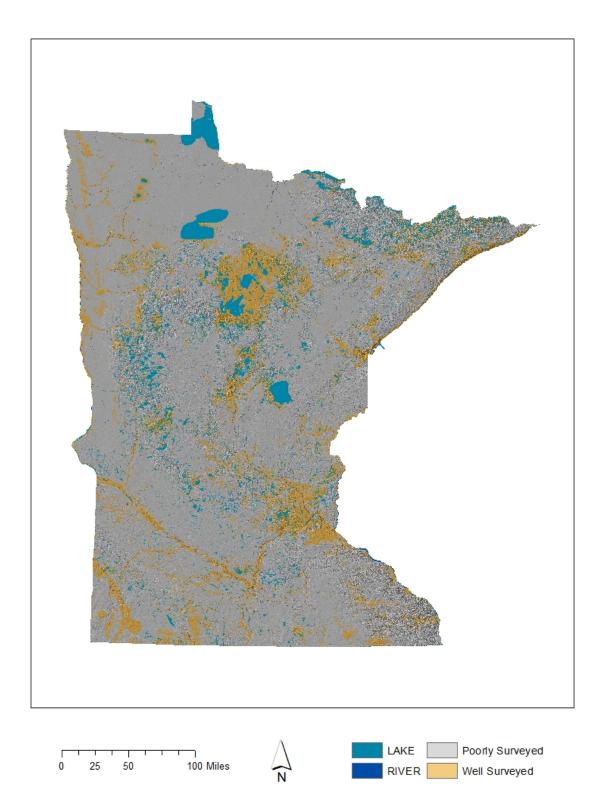
## Survey Likelihood Model

Phase 4 survey likelihood model performance increased dramatically over Phase 3 (Table 8). Statewide, 82 percent of surveys were in the high survey likelihood area. More important, 99 percent of non-surveyed points were in the low likelihood area. This very high model specificity value is responsible for the 0.93 overall accuracy of the Phase 4 survey model.

Compared to Phase 3, a much lower percentage of the land area is classified as likely to have been surveyed, further emphasizing the fact that our archaeological database is far from being a random sample. Because of our current understanding of site distributions, archaeologists prioritize surveys near water (Figure 21). To
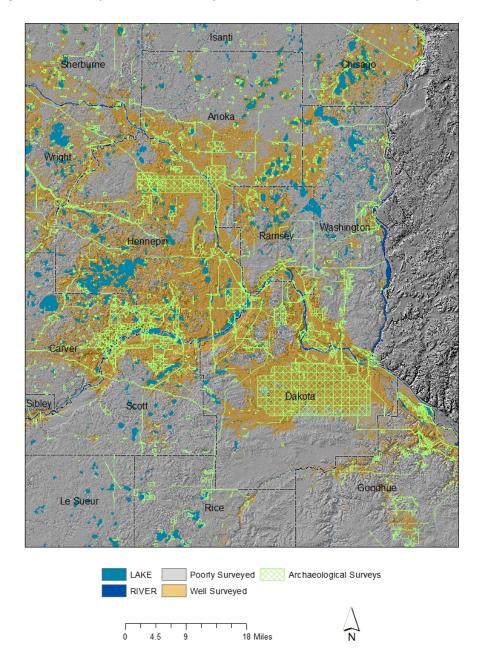
illustrate this bias, the mean path distance to a large lake is 9,754 for sites and 9,763 for surveys.  Nevertheless, surveys associated with large projects, such as highways and pipelines, should ensure that some areas away from water are surveyed.  Still, because of the nature of the projects, these surveys are likely to follow terrain that is relatively level and that may have other characteristics that make the survey locations rather predictable.

**Table 8: Phase 3 vs. Phase 4 Statewide Survey Model Performance**

| Performance Measure | Phase 3 | Phase 4 | Improvement |
|---|---|---|---|
| Sensitivity (True Positive Rate) | 0.68 | 0.82 | 0.14 |
| Specificity (True Negative Rate) | 0.65 | 0.99 | 0.34 |
| Overall Accuracy | 0.66 | 0.93 | 0.27 |
| % Region High Probability | 0.43 | 0.16 | 0.27 |
| GAIN | 0.36 | 0.80 | 0.44 |

**Figure 21: MnModel Phase 4 Survey Likelihood Model**



LAKE    Poorly Surveyed

RIVER    Well Surveyed

0    25    50        100 Miles

N

Notable exceptions to the 'survey near water' bias are found in the Chippewa National Forest and the Minneapolis-St. Paul metropolitan area. The U.S. Forest Service has surveyed fully 41 percent of their dry land area, producing a large 'well surveyed' area surrounding Leech, Cass, and Wininibigoshish Lakes in north-central Minnesota (Figure 21). In the metropolitan area, intensive surveying has been associated with development. Northern Dakota and southwestern Hennepin Counties are particularly well-surveyed (Figure 22).

**Figure 22: Metropolitan Area Survey Distributions and Phase 4 Survey Likelihood Model**



The survey likelihood model is meant to serve as a 'heads-up' to archaeologists about which types of environments have been well-surveyed and which have not. The 'well surveyed' category does not mean,

however, that any particular location has actually been surveyed (Figure 22). Archaeologists must consult the digitized survey boundaries as well as additional survey reports that have not yet been digitized.

## Survey Implementation Model

The site potential model and the survey likelihood model are combined to create the survey implementation model. (Figure 23) This composite model provides archaeologists with information about both site and survey distributions, allowing them to make informed decisions about where future surveys are needed.

When Phase 4 survey implementation models are compared to Phase 3 models, the higher specificity of the Phase 4 survey models are apparent. Both categories characterized by low survey likelihood have increased in area, while both categories with high survey likelihood have decreased (Table 9). The higher sensitivity of the Phase 4 site model is apparent when looking at the difference between the Phase 3 and 4 models where both site potential and survey likelihood are high. Maps comparing the Phase 3 and Phase 4 survey implementation models for each region are provided in the appendix to this report.

**Table 9. Phase 3 vs. Phase 4 Statewide Survey Implementation Model Performance**

| Site Potential | Survey Likelihood | Phase 3 (% Region) | Phase 4 (% Region) | Difference |
|---|---|---|---|---|
| High | High | 19.4 | 6.5 | -12.9 |
| High | Low | 3.5 | 5.9 | + 2.4 |
| Low | High | 23.9 | 9.5 | -14.4 |
| Low | Low | 49.5 | 74.4 | + 24.9 |
| Water, Mines, Steep Slopes | N/A | 3.7 | 3.7 | 0 |

**Figure 23: MnModel Phase 4 Survey Implementation Model**



LAKE    Unknown Site Potential/Poorly Surveyed    High Site Potential/Poorly Surveyed

RIVER    Low Site Potential/Well Surveyed    High Site Potential/Well Surveyed

0   25   50    100 Miles

N

# Meeting Phase 4 Goals

## Higher proportion of sites predicted

The proportion of sites correctly predicted by the models is known as sensitivity. Phase 4 models met and exceeded expectations for predicting a larger percentage of archaeological sites (Figure 24). . As a point of comparison, Phase 3 models were originally classified so that the high and medium probability classes captured 85% of the modeled site at that time. For purposes of this discussion, the high and medium probability classes have been combined and the percentage of sites predicted is based on the current modeled population. As Figure 24 shows, Phase 3 models for some regions performed better than others in predicting the current site population. However, none performed as well as the Phase 4 models, which consistently predict 95% of the modeled site population.

**Figure 24: Sensitivity of Phase 4 vs. Phase 3 Models, By Region**



In addition to predicting site locations with confidence, it is important to keep the percentage of non-sites falsely predicted to be sites as low as possible. Phase 4 models performed much better than Phase 3 models at this task (Figure 25). Only one Phase 4 model has a false positive rate higher than five percent.

**Figure 25: False Positive Rates of Phase 4 vs. Phase 3 Models, By Region**



## Reduce the area of high site potential

The very low false positive rates of the Phase 4 models can be attributed to their spatial precision. In all but one region, the percent area occupied by high site potential has been reduced (Figure 26). In Phase 3, 23 percent the state was categorized as high site potential. In Phase 4, only 13 percent of the state is high potential.

**Figure 26: Percent Area High Probability, Phase 4 vs. Phase 3 Models, By Region**



## Reduce the area of unknown site potential

'Unknown site potential' refers to areas classified as both low site potential and low likelihood of having been surveyed.  We call these areas 'unknown' because, without surveys, we cannot say for certain whether sites are not reported because they are not there or whether they are not reported simply because no one has looked there.  MnModel Phase 4 did not achieve the goal of reducing the area of unknown site potential.

To make the comparison between Phases 3 and 4 valid, the same water bodies, steep slopes, and mines were removed from the Phase 4 models that were removed from the Phase 3 models.  Excluding those removed areas, the 'low site probability/well surveyed' areas decreased from 24 percent in Phase 3 to 9 percent in Phase 4 as the survey models have become more precise.  Consequently, the 'unknown' area in the survey implementation model increased from 50 percent in Phase 3 to 74 percent in Phase 4.  Even with the greatly expanded archaeological survey database, the pattern of surveys was predictable enough in each region that precise models could be created.  These results are completely attributable to the improved archaeological surveys database and more robust statistical methods.

The key to reducing the area of unknown site potential is to reduce the area classified as 'poorly surveyed'.  Where the density of surveys is high, the extent of the area modeled as 'well surveyed' is high.  Chippewa National Forest is the best example of this (Figure 27).  The Forest Service has conducted 1,897 square kilometers of survey within their 6,462 square kilometer forest.  Since 39 percent of the forest is covered by lakes and wetlands, a full 41 percent of the dry land within the forest boundaries has been surveyed.  Consequently, the model for this area classifies 82 percent of the dry land as 'well surveyed' and only 18 percent as 'poorly surveyed'.

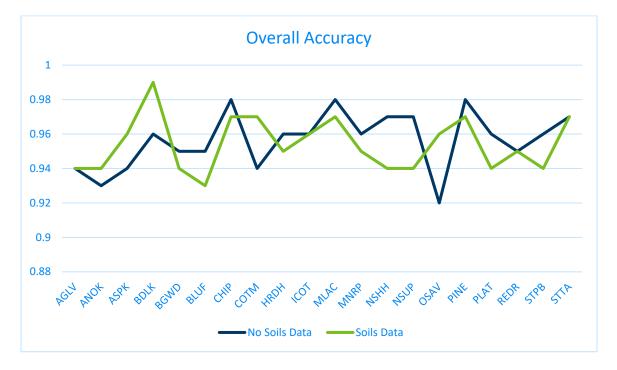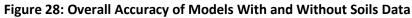**Figure 27.  Survey Model for Chippewa National Forest**



## Greater model stability

Model stability was evaluated in Phase 3 by running multiple models with different subsets of the archaeological database and evaluating their performance.  Models were considered to be 'stable' if the several models run consistently predicted similar percentages of sites.  For Phase 4, we did not run multiple models for each region with the same environmental variables.  One model was run without soil variables and one with soil variables.  Differences between these models can be attributed to either the differences in the site population (models run

with soils variables always had a smaller site population to draw from) or differences in the predictive variables available. With that in mind, we can get a rough estimate the stability of the Phase 4 models.

Figure 28 illustrates the overall accuracy (percentage of sites and background points correctly predicted) by the two models for each region. On average, there is only a difference of two percentage points. The maximum difference is four percentage points. Neither type of model consistently performs better than the other. Thus, we can consider these models to be stable.
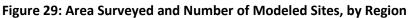
**Figure 28: Overall Accuracy of Models With and Without Soils Data**



# Potential for Improving Models

## New Archaeological Data

The greatest potential for improving the models lies in the acquisition of more archaeological data. Fewer than 200 sites were available for modeling in AGLV, ASPK, ICOT, NSUP, OSAV, PLAT, and STPB. The site models for these regions would likely be more precise and more stable with larger sample sizes. However, it likely will be difficult to find many more prehistoric sites in some of these regions. Much of AGLV is in wetlands and is not surveyable. NSUP is covered with forest, but has not been terribly under-surveyed. Most of STPB is urbanized and already has a larger portion surveyed than most regions. As Figure 29 shows, there is not always a relationship between the amount of survey and the number of site points available for modeling. At some point, there is a diminishing return on surveys, as seen in the CHIP region. Other regions, such as BDLK, BGWD, and BLUF seem to have a very high return on surveys.

**Figure 29: Area Surveyed and Number of Modeled Sites, by Region**



One way to understand this would be to estimate the *a priori*, or random, chance of finding a site. If all site locations were known, we could measure the *a priori* probability as a ratio between the number or area of sites and the area of the region. Not knowing all site locations, assessing the true *a priori* chance of finding a site would require a very large random survey. Lacking that, we will attempt an estimate from the available data. To estimate *a priori* probabilities, we intersected all prehistoric site polygons (including site types we did not model) with all survey polygons. The area of intersection is our estimate of the chance of finding a site. It is likely an overestimate of the true *a priori* probability of site occurrence in Minnesota, as surveys are more often conducted where sites are expected than where they are not expected. The survey models reported here show very clearly their non-random nature. Still, this crude estimate is informative.

Figure 30 graphs the estimated *a priori* probabilities for sites in each modeled region and allows us to compare them to site density and survey intensity. We can see that *a priori* probabilities vary from region to region. They are notably higher in the BLUF region, where site numbers are high despite a rather low portion of the region having been sampled. We would expect *a priori* probabilities and site density to track fairly well, and they do for most regions. A few anomalies stand out, however. The site density in STPB is higher than expected by the *a priori* probability estimate. This may be attributable to the higher percentage of the region that has been surveyed. Yet the data for CHIP hint that there may be a point of diminishing returns from surveys, where the number of sites identified does not increase greatly with the high survey intensity.
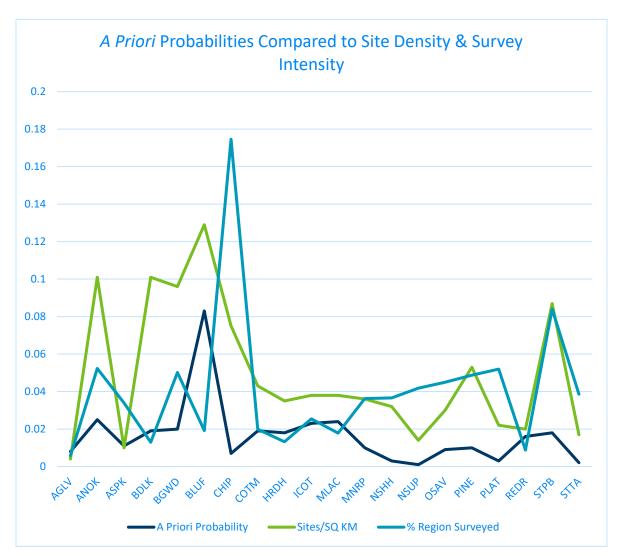
**Figure 30: Estimated a Priori Probabilities for Prehistoric Sites Compared to Site Density and Survey Intensity**



We can speculate that the regions with less separation between the *a priori* estimate and site density are those that might benefit from more survey. These include AGLV, ASPK, and REDR. Both PLAT and NSUP are relatively well-surveyed; their low site numbers are likely attributable to low *a priori* probabilities. ICOT has a reasonable density of sites, but is a very small region so that the total number of sites is low. It was combined with COTM for Phase 3 modeling, and might benefit from that in the future if more sites cannot be found.

That said, the estimated *a priori* probabilities reported here, with the exception of CHIP, are really probabilities of finding sites in environments where archaeologists have traditionally surveyed. Future surveys in the 'unknown' portions of the model might yield some surprises and provide new information that will change model patterns. Such surveys will also provide us with more realistic estimates of *a priori* probabilities.

Still, it will take very large increases in the site population to make a significant difference in the site probability models. Between Phase 3 and 4, the site population increased by 41 percent. The improvement in the models was dramatic (Figure 31), though not all of this can be attributed to the increase in site numbers. Environmental data and statistical procedures also improved dramatically. The final Phase 4 site population was about nine

percent larger than the population of sites used to develop the initial Phase 4 models.  The improvement in the models was quite small for most regions (Figure 31).

**Figure 31. Percent of Sites in High Probability (Sensitivity) of Site Models (Phases 3 and 4)**



Though increasing the number of archaeological sites may not improve site prediction, increasing the number of background points appears to improve specificity.  For the initial Phase 4 models the ratio of background points to sites was only about 2:1 in some regions.  We increased this ratio to 3:1 for the final models.  Having more background points provides additional information about the range and distribution of variable values within the region as well as information about where sites are not found.  This change clearly improved the specificity of the final models (Figure 32).  Moreover, it had the effect of reducing the total area of high probability (Figure 33) and improving overall model accuracy (Figure 34).
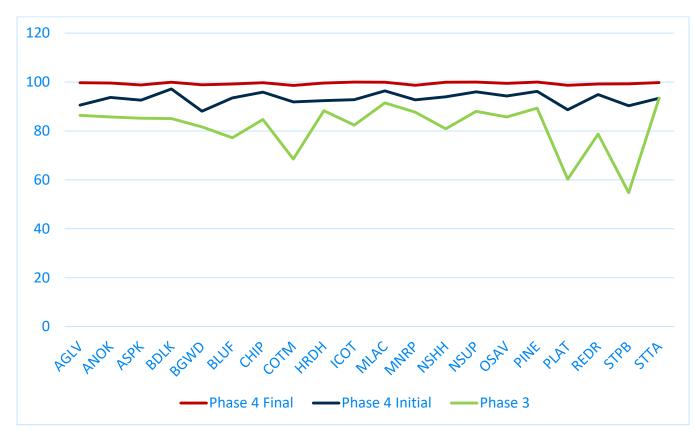
**Figure 32. Percent of Non-Sites Correctly Predicted (Specificity) in Site Models (Phases 3 and 4)**



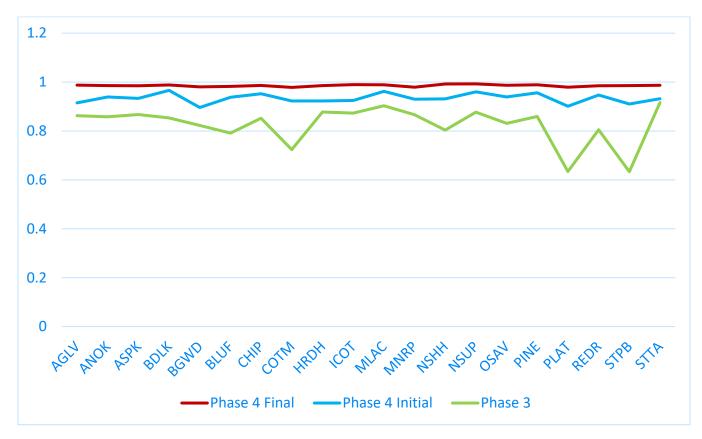**Figure 33. Percent Land Area in High Site Probability (Phases 3 and 4)**

**Figure 34. Overall Accuracy of Site Potential Models (Phases 3 and 4)**



If additional surveys cannot make large differences in the site probability models, they can make significant improvements in the survey models and, in particular, in the confidence we can place in whether areas have been adequately surveyed. As Figures 29 and 30 show, the Chippewa Plains region has been more intensively surveyed than other regions, and this is reflected in the extent of the High Survey Likelihood/Low Site Potential area in the survey implementation model (Figure 27). This will require, though, that archaeologists direct some survey effort to the areas now classified as 'unknown'. Continued survey of only the high probability portions of the map will simply reinforce, rather than improve, the current model.

## Improved Archaeological Attributes

Current models rely more on archaeological site locations than on attribute data. These models cannot distinguish between site distributions from different time periods, cultural affinities, or functions. Because the models built using statistical analysis, the types of sites that are most abundant in the database will have the greatest effect on the model. In Minnesota, these are the Woodland Tradition sites (Table 10).

**Table 10: Cultural Traditions in Minnesota**

| Tradition | Number of Sites | Dominant Regions (Number of Sites) |
|---|---|---|
| Paleoindian | 220 | MNRP (40) |
| Archaic | 625 | MNRP (108) CHIP (60) |
| Woodland | 4,298 | CHIP (520) BGWD (499) MNRP (488) |
| Plains Village | 138 | MNRP (70) |
| Mississippian | 96 | BLUF (43) |
| Oneota | 116 | MNRP (51) BLUF (32) |

Efforts to model selected segments of the site population, such as time periods, cultural affinities, or functions, would benefit from additional work to verify, correct, and augment the information now in the database by careful examination of site forms and collections. Not all sites in the current database have been associated with a tradition. Of those that have, the numbers associated with non-Woodland traditions are quite low (Table 10). If we assume that we need a sample of about 100 sites for modeling, we could model Woodland-only sites in only 13 of 20 regions. The only non-Woodland model we could build would be for Archaic sites in the MNRP region. For the same reason, it would be difficult to build statistical models based on sites in specific functional classes. Better attribute data are needed before more refined models can be developed.

# Improved Environmental Data

## Lithic Sources

Phase 4 models predict the majority of sites that are currently identified. Sites that require a rock source (rock shelters, rock art, quarries) cannot be included in the models because we have no data about the locations of rock outcrops or sources of specific types of lithic materials. Mapping of these lithic sources could improve our ability to predict locations of this class of sites.

### Historic Environmental Data

The Public Land Survey data are critical to the accuracy of the vegetation model and historic hydrographic model. These data can be further improved in several ways.

- Associate surveyors' line notes with GIS section lines in the southern half of the state. This extension of Almendinger's work will provide more detailed information for improving the vegetation and historic hydrographic models.

- Further quality control of the vegetation classification of both the existing digitized line notes and the section corners. In both cases, MnDNR used very generalized vegetation classes. Moreover, the vegetation classes are not always assigned to the correct feature. For example, when a section corner falls within a wetland, the vegetation class assigned may be that of the surrounding vegetation rather than of the wetland.

- Improve georeferencing of the scanned Public Survey plat maps. These were originally georeferenced using only township corners. Georeferencing to section corners will greatly improve their accuracy.

- Re-digitize features from Public Survey plat maps that move because of improved georeferencing.

- Digitize cultural features (trails, houses, Indian villages) from the plat maps.

### Landscape Model

The Phase 4 landscape model is based on available geomorphic data from several sources and a range of scales. More high resolution (1:24,000 scale) geomorphic mapping would improve this model and likely improve the performance of the variables derived from it.

# Advice for Users of the Model

Archaeologists using the Phase 4 survey implementation model should be aware of its limitations. Most importantly, five percent of the sites in the modeling database do not fall within the high probability areas. That implies that at least five percent of undiscovered sites will be in the low and unknown site potential areas. Moreover, there are types of sites not used for modeling that can be in any of these areas. These include historic sites, single artifacts, and sites associated with rock outcrops (quarries, rock art, rock shelters).

The model must be used in conjunction with other data. Unknown site potential areas are likely not to have been surveyed, but in some cases that may be for good reasons. Extensive wetlands or steep slopes may make areas unsurveyable. High site potential areas, on the other hand, may have experienced sufficient surface disturbance that sites are no longer present or, if present, lack integrity and context.

The model indicates potential for surface sites on surfaces that have not been disturbed. In some geomorphic situations, historic and prehistoric surfaces have been either eroded away or buried by modern sediment. The MnModel Landform Sediment Assemblage (LfSA) data and Landscape Suitability models provide information

about such surfaces as well as assessing the potential for subsurface layers to contain buried sites (Hajic et al. 2000; Hajic and Hudak 2002; Hajic et al. 2009 and 2011; Hudak and Hajic 2002).

In conclusion, the Phase 4 survey implementation model is very precise and very accurate.  It is not, however, perfect.  It is intended to be used in conjunction with other data by professional archaeologists.  Survey outside of the high site potential area is strongly encouraged – both to avoid missing the types of sites that may be found there and to provide new information for future modeling.

# Conclusions

MnModel Phase 4 models perform exceptionally well.  Site models exhibit both high sensitivity and high specificity.  Survey models are less sensitive to where surveys have been conducted, but are highly specific about which locations are not likely to have been surveyed.

The statistical procedures used are both sophisticated and robust.  Future model improvement will require further investment in the underlying data.  In particular, surveys of currently under surveyed areas are needed to confirm site absence.

# References

Aaseng, Norman E. et al.
> 1993 *Minnesota's Native Vegetation: A Key to Natural Communities. Version 1.5*. Minnesota Department of Natural Resources, Natural Heritage Program. St. Paul, MN.

Anfinson, Scott F.
> 1990 Archaeological Regions in Minnesota and the Woodland Period. In *The Woodland Tradition in the Western Great Lakes: Papers Presented to Elden Johnson*, edited by G.E. Gibbon, pp. 135-166. University of Minnesota Publications in Anthropology No. 4. Department of Anthropology. University of Minnesota, Minneapolis.

Breiman, Leo
> 2001 Random Forests.  *Machine Learning* 45(1):5-32.

Brown, Andrew, Alexander Anton, Luke Burds, and Elizabeth Hobbs
> 2019 *Tool Handbook*.  Appendix C in *MnModel Phase 4 User Guide*, by Carla Landrum et al. Minnesota Department of Transportation. St. Paul, MN.

Cassell, Mark S., Howard D. Mooers, Clark A. Dobbs, Thomas Madigan, Morgan Coville, Jeff Berry and Douglas A. Birk
> 1997 An archaeological sensitivity model of prehistoric and contact period settlement at Camp Ripley, Morrison County, Minnesota.  Institute for Minnesota Archaeology Reports of Investigation Number 397, submitted to St. Paul District, U.S. Army Corps of Engineers, St. Paul.

Cleland, D.T., P.E. Avers, W.H. McNab, M.E. Jensen, R.G. Bailey, T. King, and W.E. Russell
  1997 National Hierarchical Framework of Ecological Units.  In *Ecosystem Management Applications for Sustainable Forest and Wildlife Resources,* edited by M.S. Boyse and A. Haney, pp. 181-200. Yale University Press, New Haven, CT.

Dalla Bona, L.
  1994 Methodological Considerations.  Cultural Heritage Resource Predictive Modeling Project Vol. 4 Centre for Archaeological Resource Prediction, Lakehead University, Thunder Bay, Ontario.
  2000 *Protecting Cultural Resources through Forest Management Planning in Ontario Using Archaeological Predictive Modeling.*  Chapter 5 in *Practical Applications of GIS for Archaeologists: A Predictive Modeling Kit*, edited by Konnie L. Westcott and R. Joe Brandon, pp. 73-99.  Taylor & Francis, London.

Ejstrud, Bo
  2003 Indicative Models in Landscape Management: Testing the Methods. In *Symposium: The Archaeology of Landscapes and Geographic Information Systems: Predictive Maps, Settlement Dynamics and Space and Territory in Prehistory*, edited by Jürgen Kunow and Johannes Müller, pp. 119-134.  Forschungen zur Archäologie im Land Brandenburg 8: Archäoprognose Brandenburg I, Wünsdorf.

Fry, G.L.A., B. Skar, G. Jerpåsen, V. Bakkestuen, and L. Eristad
  2004 Locating Archaeological Sites in the Landscape: A Hierarchical Approach Based on Landscape Indicators. *Landscape and Urban Planning* 67: 97-107.

Gibbon, Guy E., Craig M. Johnson and Stacey Morris
  2002 The Archaeological Database.  Chapter 5 in *Mn/Model Final Report Phases 1-3*, edited by G. Joseph Hudak, et al. Minnesota Department of Transportation. St. Paul, MN.

Hajic, Edwin R. and Curtis M. Hudak
  2002 *Landform Sediment Assemblages in the Upper Mississippi Valley, St. Cloud to St. Paul, for Support of Cultural Resource Investigations*. Minnesota Department of Transportation. St. Paul, MN. 38 pp.

Hajic, Edwin R., Curtis M. Hudak, and Jeffrey Walsh
  2009 *Landform Sediment Assemblages in the Anoka Sand Plain for Support of Cultural Resource Investigations*.  Minnesota Department of Transportation.  St. Paul, MN.
  2011 *Landform Sediment Assemblages in the Mississippi River Valley and Selected Tributaries between the City of St. Paul and the Minnesota-Iowa Border for Support of Cultural Resource Investigations*. Minnesota Department of Transportation.  St. Paul, MN.

Hajic, Edwin R., Philip E. Paradies, and Curtis M. Hudak
  2000 *How to Construct a MnModel Landscape Suitability Model*. Minnesota Department of Transportation. St. Paul, MN.

Hanson, D.H. and B.C. Hargrave
  1996 Development of a Multilevel Ecological Classification System for the State of Minnesota. *Environmental Monitoring and Assessment* 39:75-84.

Harris, Matthew D., Robert G. Kingsley, and Andrew R. Sewell
  2015 *Archaeological Predictive Model Set: Final Report*. Commonwealth of Pennsylvania, Department of
    Transportation.  Harrisburg, PA.

Hobbs, Elizabeth
  2002a GIS Design. Chapter 4 in *Mn/Model Final Report Phases 1-3*, edited by G. Joseph Hudak, et al.
    Minnesota Department of Transportation. St. Paul, MN.
  2002b Mn/Model GIS Standards and Procedures. Appendix B in *Mn/Model Final Report Phases 1-3*, edited by
    G. Joseph Hudak et al.  Minnesota Department of Transportation. St. Paul, MN.
  2003 The Minnesota Archaeological Predictive Model. In *Symposium: The Archaeology of Landscapes and
    Geographic Information Systems: Predictive Maps, Settlement Dynamics and Space and Territory in
    Prehistory*, edited by Jürgen Kunow and Johannes Müller, pp. 141-150.  Forschungen zur Archäologie im
    Land Brandenburg 8: Archäoprognose Brandenburg I, Wünsdorf.
  2019 *Historic Vegetation Model for Minnesota: MnModel Phase 4*.  Minnesota Department of Transportation.
    St. Paul, MN.

Hobbs, Elizabeth, Andrew Brown, Alexander Anton, and Luke Burds
  2019a *Historic/Prehistoric Hydrographic Models for Minnesota: MnModel Phase 4*.  Minnesota Department
    of Transportation. St. Paul, MN.

Hobbs, Elizabeth, Andrew Brown, Alexander Anton, Jeffrey Walsh, Carson Smith, and Luke Burds
  2019b Preparing Data for Modeling.  Appendix B in *MnModel Phase 4 User Guide*, by Carla Landrum et al.
    Minnesota Department of Transportation. St. Paul, MN.

Hobbs, Elizabeth, Craig M. Johnson, and Guy E. Gibbon
  2002 Model Development and Evaluation.  Chapter 7 in *Mn/Model Final Report Phases 1-3*, edited by G.
    Joseph Hudak, et al.  Minnesota Department of Transportation. St. Paul, MN.

Hobbs, Elizabeth and Tatiana Nawrocki
  2002  Archaeological and Environmental Variables  Chapter 6 in *Mn/Model Final Report Phases 1-3*, edited by
    G. Joseph Hudak, et al.  Minnesota Department of Transportation. St. Paul, MN.

Hobbs, Elizabeth, Jeffrey Walsh and Curtis M. Hudak
  2019 *Environmental Variables: MnModel Phase 4*.  Minnesota Department of Transportation. St. Paul, MN.

Hudak, Curtis M. and Edwin R. Hajic
  2002 Landscape Suitability Models for Geologically Buried Precontact Cultural Resources.  Chapter 12 in
    *Mn/Model Final Report Phases 1-3*, edited by G. Joseph Hudak, et al. Minnesota Department of
    Transportation. St. Paul, MN.

Hudak, G. Joseph, Elizabeth Hobbs, Allyson Brooks, Carol Ann Sersland, and Crystal Phillips, Eds.
  2002 *Mn/Model Final Report Phases 1-3*.  Minnesota Department of Transportation. St. Paul, MN.

Judge, W. James and Lynne Sebastian, Eds.
  1988 *Quantifying the Present and Predicting the Past: Theory, Method, and Application of Archaeological
    Predictive Modeling*.  U.S. Department of the Interior, Denver.

Kamermans, Hans

2011 Predictive Maps in the Netherlands, Problems and Solutions. In *A Piccoli Passi: Archaeologica predittiva e preventive nell'esperienza cessenate,* edited by Sauro Gelichi and Claudio Negrelli, pp. 13-18. All'Insegna del Giglio, Florence.

Kauhi, Tonya C. and Joanne L. Markert

2015 *Washington Statewide Archaeology Predictive Model Report*. Washington Department of Archaeological and Historic Preservation, Olympia.

Kohler, Timothy A.

1988 Predictive Locational Modeling: History and Current Practice. In *Quantifying the Present and Predicting the Past: Theory, Method, and Application of Archaeological Predictive Modeling*, edited by W.J. Judge and L. Sebastian, pp. 19-59. U.S. Government Printing Office, Washington, DC.

Kohler, T.A. and S.C. Parker

1986 Predictive Models for Archaeological Resource Location. In *Advances in Archaeological Method Theory*, Vol. 9, edited by M.B. Schiffer, pp. 397-452. Academic Press, New York.

Kvamme, Kenneth L.

1985 Determining Empirical Relationships between the Natural Environment and Prehistoric Site Locations: A Hunter-Gatherer Example. In *For Concordance in Archaeological Analysis: Bridging Data Structure, Quantitative Technique, and Theory*, edited by C. Carr, pp. 208-238. Westport Press, Kansas City.

1988 Development and testing of quantitative Models. In *Quantifying the Present and Predicting the Past: Theory, Method, and Application of Archaeological Predictive Modeling*, edited by W.J. Judge and L. Sebastian, pp. 325-428. U.S. Government Printing Office, Washington, DC.

1992 A Predictive Site Location Model on the High Plains: An Example with an Independent Test. *Plains Anthropologist* 37(138): 325-428.

2006 There and Back Again: Revisiting Archaeological Locational Modeling. In *GIS and Archaeological Site Location Modeling,* edited by Mark W. Mehrer and Konnie L. Wescott, pp. 3-38. Taylor & Francis, Boca Raton.

2011 Archaeological Predictive Modeling in the USA. In *A Piccoli Passi: Archaeologica predittiva e preventive nell'esperienza cessenate,* edited by Sauro Gelichi and Claudio Negrelli, pp. 19-26. All'Insegna del Giglio, Florence.

Landrum, Carla and Elizabeth Hobbs

2019 *Vegetation Modeling User's Guide: MnModel Phase 4*. Minnesota Department of Transportation. St. Paul, MN.

Landrum, Carla, Elizabeth Hobbs, Alexander Anton, Andrew Brown, and Luke Burds

2019 *Archaeological Predictive Modeling Guide: MnModel Phase 4.* Minnesota Department of Transportation. St. Paul, MN.

Lively, R.S., G.B. Morey, and E.J. Bauer

2002 *One hundred years of mining: alterations to the physical and cultural geography of the western half of the Mesabi Iron Range, northern Minnesota*. MGS Miscellaneous Map Series, M-118, 4 pls. Scale

1:100,000.  Paper plots. Pl. 1, land-surface topography; pl. 2, drainage and cultural features; pl. 3, topographic disturbance; pl. 4, bedrock geology.  Minnesota Geological Survey. St. Paul, MN.

Madry, Scott, Steve Gould, Ben Resnick, and Matt Wilkerson
    2003 A GIS-based Archaeological Predictive Model for the North Carolina department of Transportation. In *Symposium: The Archaeology of Landscapes and Geographic Information Systems: Predictive Maps, Settlement Dynamics and Space and Territory in Prehistory*, edited by Jürgen Kunow and Johannes Müller, pp. 161-170.  Forschungen zur Archäologie im Land Brandenburg 8: Archäoprognose Brandenburg I, Wünsdorf.

Marschner, Francis J.
    1974 *The Original Vegetation of Minnesota.* Compiled from U.S. General Land Office Survey notes.  North Central Forest Experiment Station, Forest Service, U.S. Department of Agriculture.

Mehrer, Mark W. and Konnie L. Wescott, Eds.
    2006 *GIS and Archaeological Site Location Modeling*.  Taylor & Francis, Boca Raton.

Oehlert, Gary W. and Brian Shea
    2007 *Statistical Methods for MnModel Phase 4: Final Report*.  Research Services Section, Minnesota Department of Transportation.  St. Paul, MN.

Parker, Sandra
    1985 Predictive Modeling of Site Settlement Systems Using Multivariate Logists.  In *For Concordance in Archaeological Analysis: Bridging Data Structure, Quantitative Technique, and Theory*, edited by C. Carr, pp. 173-207. Westport Press, Kansas City

Schaetzl, Randall J., Frank J. Krist, Jr., Kristine Stanley, and Christina M. Hupy
    2009 The natural soil drainage index: an ordinal estimate of long-term soil wetness.  *Physical Geography* 30:383-409.

Stark, Stacey L., Patrice M. Farrell, and Susan C. Mulholland
    2008 *Methods to Incorporate Historic Surface Hydrology Layer in Mn/Model [Phase 4] Using Existing Geographic Information System Data*.  Minnesota Department of Transportation. St. Paul, MN.

Thoms, P. Martijn
    1996 GIS and Locational Modeling in Dutch Archaeology: A Review of Current Approaches.  In *New Methods, Old Problems: Geographic Information Systems in Modern Archaeological Research*, edited by Herbert D. G. Maschner, pp. 177-197. Occasional Paper No. 23, Center for Archaeological Investigations, Southern Illinois University, Carbondale.

Van Leluwen, Robert E.
    1990 Predictive modelling in archaeology: a primer.  In *Interpreting Space: GIS and Archaeology*, edited by Kathleen M.S. Allen, Stanton W. Green, and Ezra B.W. Zubrow. Taylor & Francis, Bristol.

Van Leusen, Martijn and Hans Kamermans, Eds.
  2005 *Predictive Modelling for Archaeological Heritage Management: A Research Agenda*. Nederlands Archaeologischy Rapporten 29.  Rijksdienst voor het Oudheidkundig Bodemonderzoek, Amersfoort.

Verhagen, Philip
  2007 *Case Studies in Archaeological Predictive Modelling*. Archaeological Studies Leiden University, Leiden University Press, the Netherlands.

Warren, Robert E.
  1990 Predictive modelling in archaeology: a primer.  In *Interpreting Space: GIS and Archaeology*, edited by Kathleen M.S. Allen, Stanton W. Green, and Ezra B.W. Zubrow. Taylor & Francis, Bristol.

Warren, Robert E. and David L. Asch
  1996 A Predictive Model of Archaeological Site Location in the Eastern Prairie Peninsula, Illinois.  Illinois State Museum Research & Collections Center, Springfield.

Westcott, Konnie L. and R. Joe Brandon, Eds.
  2000 *Practical Applications of GIS for Archaeologists: A Predictive Modeling Kit*.  Taylor & Francis, London.

Westcott, Konnie L. and James R. Kuiper
  2000 *Using a GIS to Model Prehistoric Site Distributions in the Upper Chesapeake Bay.*  Chapter 4 in *Practical Applications of GIS for Archaeologists: A Predictive Modeling Kit*, edited by Konnie L. Westcott and R. Joe Brandon, pp. 59-72.  Taylor & Francis, London.

White, Devin A. and Sarah B. Barber
  2012 Geospatial modeling of pedestrian transportation networks: a case study from pre-Columbian Oaxaca, Mexico. *Journal of archaeological Science* 39: 2684-2696.

# Appendix A: Predictor Variables Used in MnModel Phase 4

**Table A1: Complete List of MnModel Phase 4 Predictor Variables**

| VARIABLE | DEFINITION |
|----------|------------|
| ASP_RNG | Aspect range |
| CP_BOG | Path distance to nearest historic bog |
| CP_FLOOD | Path distance to nearest historic floodplain |
| CP_INT | Path distance to nearest intermittent stream |
| CP_LAKE | Path distance to nearest historic lake |
| CP_LLK | Path distance to nearest large historic lake |
| CP_MAJPATH | Path distance to nearest major pedestrian transportation route |
| CP_MAJRIDGE | Path distance to nearest major ridge or divide |
| CP_MARSH | Path distance to nearest historic marsh |
| CP_MEADOW | Path distance to nearest historic wet meadow or fen |
| CP_MEDPATH | Path distance to nearest medium pedestrian transportation route |

| VARIABLE | DEFINITION |
|---|---|
| CP_MINPATH | Path distance to nearest minor pedestrian transportation route |
| CP_MINRIDGE | Path distance to nearest minor ridge or divide |
| CP_PEREN | Path distance to nearest perennial stream |
| CP_PFLOOD | Path distance to nearest prehistoric floodplain |
| CP_PLAKE | Path distance to nearest prehistoric lake |
| CP_PLLK | Path distance to nearest large prehistoric lake |
| CP_PWET | Path distance to nearest prehistoric wetland |
| CP_RICE | Path distance to nearest wild rice location |
| CP_RIVER | Path distance to nearest river |
| CP_SWAMP | Path distance to nearest historic swamp |
| CP_WAT | Path distance to nearest historic surface water (of all types) |
| CP_WET | Path distance to nearest historic 'wet' land |
| CP_WETLAND | Path distance to nearest historic wetland (of any type) |

| VARIABLE | DEFINITION |
|----------|------------|
| CURV | Surface Curvature |
| DI | Drainage Index |
| DRAIN | Soil drainage |
| ELEV | Elevation |
| FFD_R | Frost-free days |
| FLDFRQD | Flood frequency |
| HYDGRPDCD | Hydric Group (dominant condition) |
| HYDPRS | Hydric soil presence |
| HZDEP | Depth of surface soil horizon |
| ISLAND | On an island |
| LFORM | Landform |
| LSCAPE | Landscape |
| MAJ_SIZE | Size of major watershed |

| VARIABLE | DEFINITION |
|----------|------------|
| MIN_SIZE | Size of minor watershed |
| ORD_STRM | Order of nearest stream |
| PATH_ORD | Order of nearest pedestrian transportation route |
| PI | Productivity Index |
| REL | Relative Elevation |
| REL90 | Relative Elevation within 90 meters |
| RGH | Surface Roughness |
| RGH90 | Surface Roughness within 90 meters |
| SHELTER | Shelter Index |
| SLOPE | Percent Slope |
| TPI1000 | Topographic Position Index within 1000 meters |
| TPI1MI | Topographic Position Index within one Mile |
| TPI250 | Topographic Position Index within 250 meters |

| VARIABLE | DEFINITION |
|---|---|
| TPI5MI | Topographic Position Index within five miles |
| TPI90 | Topographic Position Index within 90 meters |
| TWI | Topographic Wetness Index |
| VEGDIV10K | Vegetation diversity within ten km |
| VEGDIV1K | Vegetation diversity within one km |
| VEGDIV5K | Vegetation diversity within five km |
| VEGMOD | Historic vegetation type |
| VISIBLE | Visibility |
| WETSOIL | On a wetland soil |

# Appendix B: Comparisons of Phase 3 and Phase 4 Survey Implementation Models by Region

This appendix consists of maps of each Phase 4 modeling region comparing Phase 3 and Phase 4 survey implementation models. For purposes of comparison, the Phase 3 models have been symbolized according to the same criteria as the Phase 4 models. The legend for all maps is provided here:

## MnModel Phase 4 Survey Implementation Model

Unknown Site Potential/Poorly Surveyed
Low Site Potential/Well Surveyed
High Site Potential/Poorly Surveyed
High Site Potential/Well Surveyed

Historic lakes and rivers from the MnModel Phase 4 historic hydrographic model are displayed over both Phase 3 and Phase 4 models. Since modern water bodies were incorporated into the Phase 3 models, these may be apparent on the Phase 3 models where they do not overlap historic water bodies. Mines (in black) and steep slopes (in purple) were also incorporated into the Phase 3 models, but do not appear on the Phase 4 models.

The percentage of all Phase 4 sites predicted by each model is reported on the maps. Phase 3 models were designed to predict 85 percent of sites known at that time. Statistics presented here show what percentage of sites now known are predicted by those models.

Kvamme's GAIN statistic is reported on the maps. The GAIN statistic is a commonly used measure of model performance. GAIN is calculated as:

$$GAIN = 1 - (\% \text{ of region in high probability}/\% \text{ of sites in high probability})$$
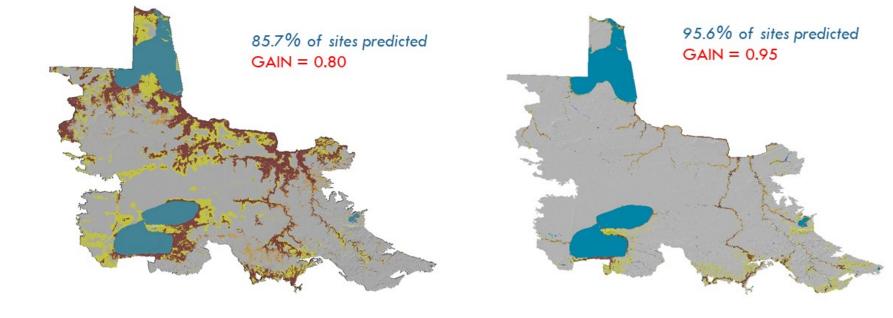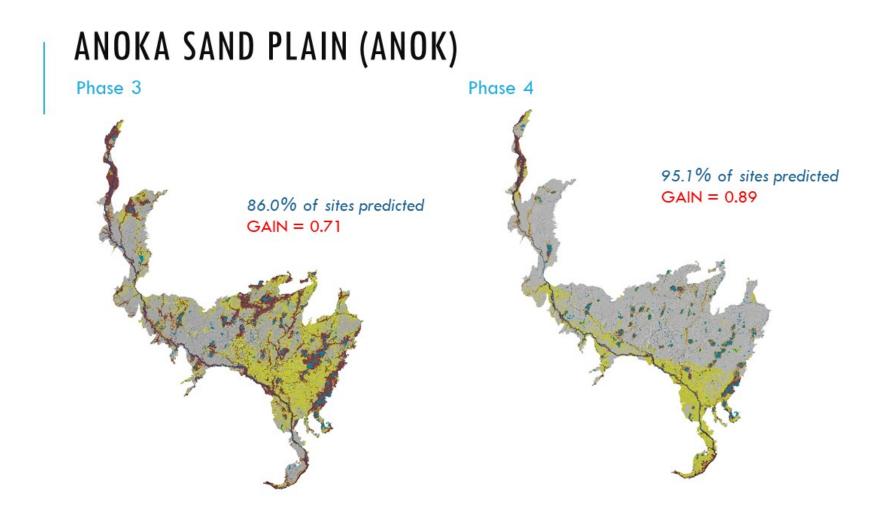
# STATEWIDE SURVEY IMPLEMENTATION MODEL

Phase 3

Phase 4

# AGASSIZ LOWLANDS/LITTLEFORK-VERMILION UPLANDS (AGLV)

Phase 3

Phase 4



85.7% of sites predicted
GAIN = 0.80

95.6% of sites predicted
GAIN = 0.95

# ANOKA SAND PLAIN (ANOK)
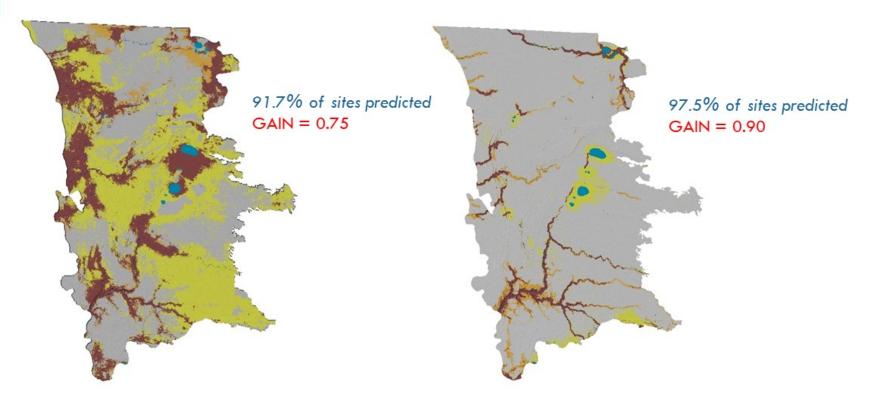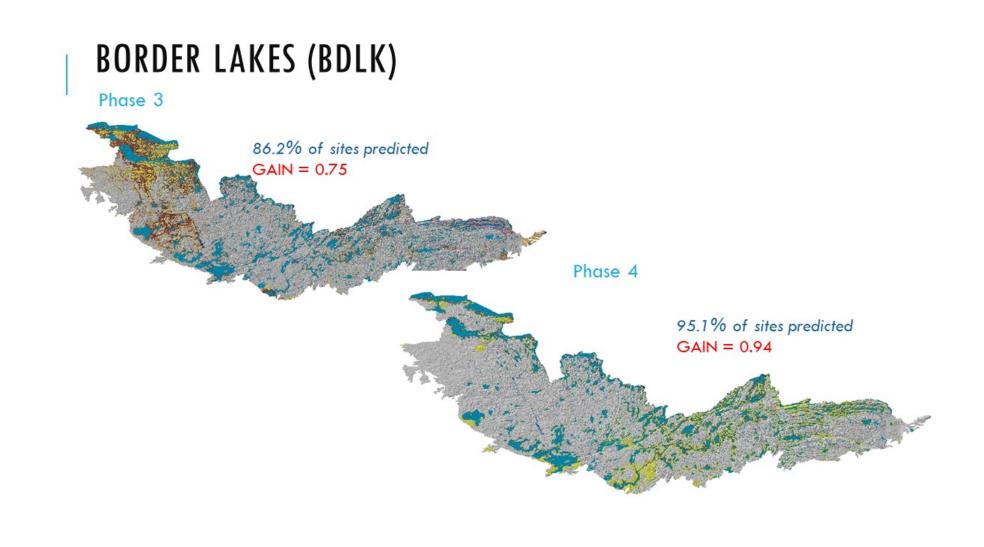
Phase 3

Phase 4



86.0% of sites predicted
GAIN = 0.71

95.1% of sites predicted
GAIN = 0.89

# ASPEN PARKLANDS (ASPK)

Phase 3

Phase 4

91.7% of sites predicted
GAIN = 0.75

97.5% of sites predicted
GAIN = 0.90

# BORDER LAKES (BDLK)

Phase 3

86.2% of sites predicted
GAIN = 0.75

Phase 4

95.1% of sites predicted
GAIN = 0.94

# BIG WOODS (BGWD)

Phase 3

Phase 4

84.1% of sites predicted
GAIN = 0.59

95.3% of sites predicted
GAIN = 0.83

# BLUFFLANDS (BLUF)

Phase 3

85.2% of sites predicted
GAIN = 0.58

Phase 4

95.1% of sites predicted
GAIN = 0.83

# CHIPPEWA PLAINS (CHIP)

## Phase 3

*86.9% of sites predicted*
GAIN = 0.69

## Phase 4

*95.0% of sites predicted*
GAIN = 0.93

# COTEAU MORAINES (COTM)

Phase 3

Phase 4

85.6% of sites predicted
GAIN = 0.55

95.2% of sites predicted
GAIN = 0.89

# HARDWOOD HILLS (HRDH)

Phase 3

Phase 4



84.0% of sites predicted
GAIN = 0.78

95.1% of sites predicted
GAIN = 0.91

# INNER COTEAU (ICOT)
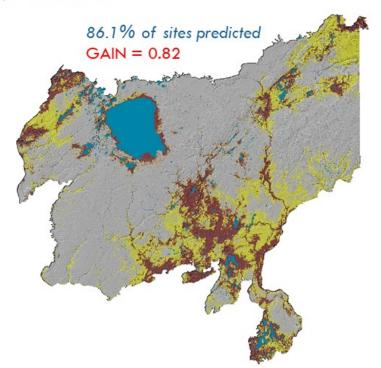
Phase 3

72.4% of sites predicted
GAIN = 0.65

Phase 4

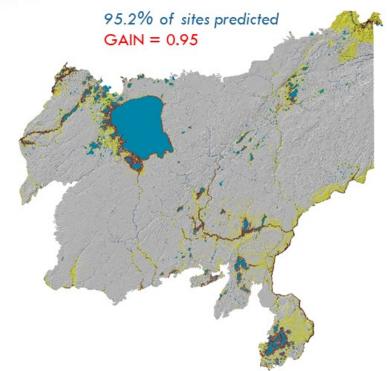95.3% of sites predicted
GAIN = 0.85

# MILLE LACS UPLANDS (MLAC)
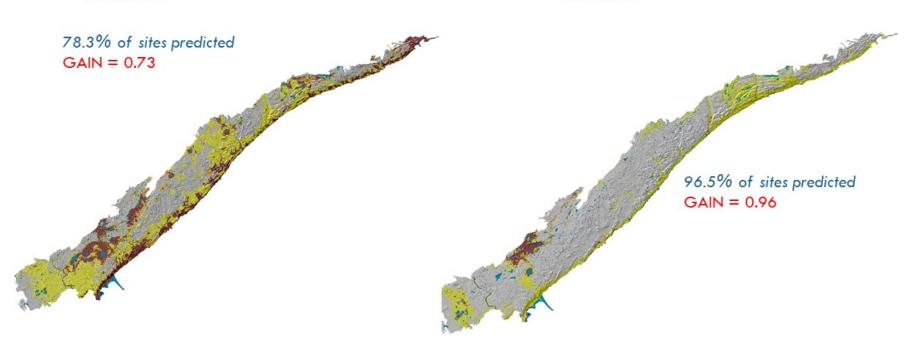## (WITH GLACIAL LAKE SUPERIOR PLAIN AND ST. CROIX MORAINE)

Phase 3

Phase 4

86.1% of sites predicted
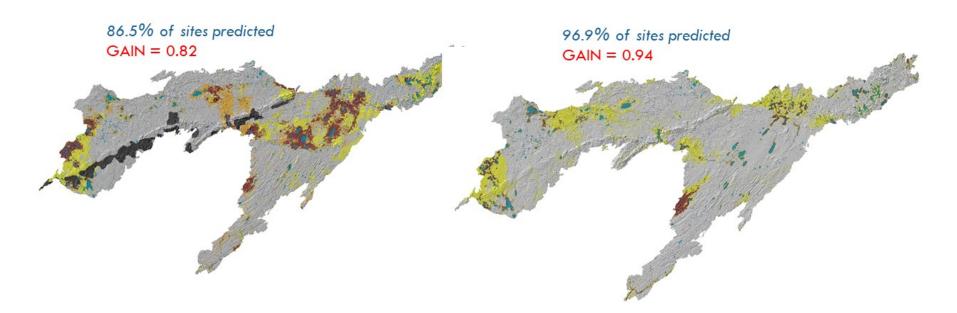GAIN = 0.82

95.2% of sites predicted
GAIN = 0.95

# MINNESOTA RIVER PRAIRIE (MNRP)

Phase 3

Phase 4



95.1% of sites predicted
GAIN = 0.88

83.0% of sites predicted
GAIN = 0.74

# NORTH SHORE HIGHLANDS (NSHH)

Phase 3

78.3% of sites predicted
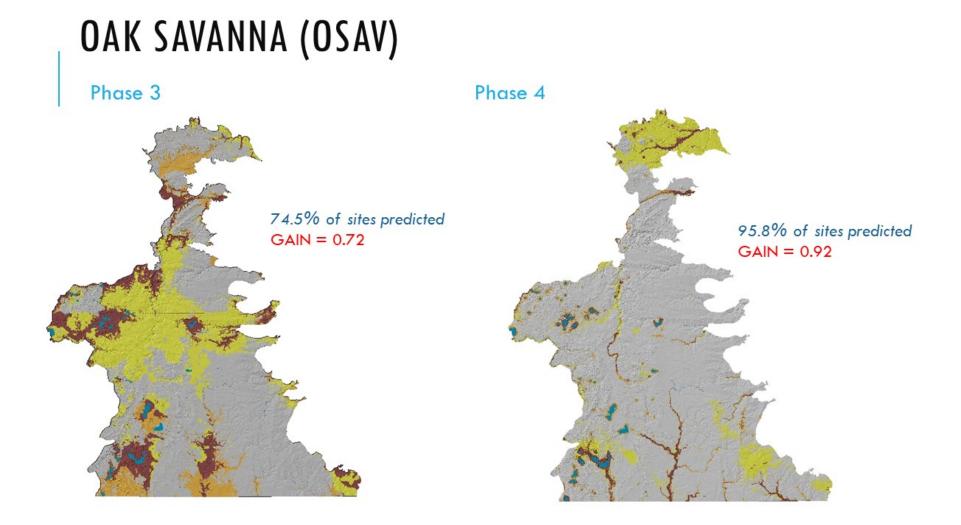GAIN = 0.73

Phase 4

96.5% of sites predicted
GAIN = 0.96

# NASHWAUK UPLANDS/LAURENTIAN UPLANDS/TOIMI UPLANDS (NSUP)

Phase 3

Phase 4

86.5% of sites predicted
GAIN = 0.82

96.9% of sites predicted
GAIN = 0.94

# OAK SAVANNA (OSAV)

Phase 3

Phase 4



74.5% of sites predicted
GAIN = 0.72

95.8% of sites predicted
GAIN = 0.92

# PINE MORAINES & OUTWASH PLAINS (PINE)

Phase 3

74.2% of sites predicted
GAIN = 0.77

Phase 4

95.1% of sites predicted
GAIN = 0.92

# ROCHESTER PLATEAU (PLAT)



Phase 3

73.4% of sites predicted
GAIN = 0.36

Phase 4

95.2% of sites predicted
GAIN = 0.82

# RED RIVER PRAIRIE (REDR)

Phase 3

86.7% of sites predicted
GAIN = 0.68

Phase 4

96.1% of sites predicted
GAIN = 0.91

# ST. PAUL-BALDWIN PLAINS & MORAINES (STPB)

Phase 3

Phase 4

92.1% of sites predicted
GAIN = 0.40

96.1% of sites predicted
GAIN = 0.78

# ST. LOUIS MORAINES/TAMARACK LOWLANDS (STTA)

## Phase 3

85.4% of sites predicted
GAIN = 0.86

## Phase 4

95.2% of sites predicted
GAIN = 0.95

# Appendix C: Acknowledgements

Phase 4 of MnModel was funded by the Federal Highway Administration, the Minnesota Department of Transportation, and Minnesota State University, Mankato.  Development of MnModel Phase 4 was a multi-year process that required the expertise and hard work of many people.

## MnModel Phase 4 Team

### Minnesota Department of Transportation

*Elizabeth Hobbs, PhD*

Elizabeth Hobbs received her Ph.D. in Geography from the University of California, Los Angeles, specializing in biogeography.  She has many years' experience in university teaching, research, and consulting.  Her professional interests include historic vegetation, human impact on vegetation, Geographic Information Systems, and archaeological predictive modeling.  She spent 19 years in the Cultural Resources Unit at MnDOT, where she has focused on the development of digital data and applications to support and streamline the Section 106 process.

Beth served as Principal Investigator for GIS for MnModel Phase 3 and as Research Director for MnModel Phase 4.  In Phase 3, Beth was responsible for conceptualizing and developing the GIS modeling procedures, implementing standards, evaluating the results, and documenting the project.  In addition to coordinating all Phase 4 research, she assembled and interpreted the historic vegetation data, conceptualized and directed the update of the hydrographic model and development of the vegetation model, and classified and evaluated the archaeological predictive models.  She is currently a Research Fellow at the AGES Lab at Minnesota State University, Mankato.

### Landform Consulting Strategies, LLC

*Curtis Hudak, PhD*

Curtis Hudak received his PhD in Geology from the University of Iowa.  His professional career has included geomorphic mapping, deep-site predictive modeling, soil and sedimentology, Geographic Information System database and model development and editing, hydrogeological studies, wetland soil delineations, geoarchaeological resource investigations, environmental site

assessments, and Lean Visual Project Management. His professional interests include geomorphology, soils, sediments, GIS-based mapping, wetland soil delineation, and natural and cultural resource assessments. Curt Hudak has been working intermittently on various archaeological predictive models since his first expert system model was developed for a natural gas pipeline corridor across the Dakotas, Minnesota and Iowa in 1979. He was Principal Investigator for Geomorphology for MnModel Phases 3 and 4.

Most recently as part of MnModel Phase 4, Curt focused on the updating of statewide digital elevation and geomorphic models; as well as the development of the prototype and pilot statistical models. He was also the project geomorphologist that managed, interpreted, and mosaicked the MnModel Phase 4's statewide geomorphology database from multiple sources including both high and low resolution surface geology and pre-existing geomorphic data.

## Whirrx LLC

*Jeffrey J. Walsh C.P., GISP*

Prior to co-founding Whirrx LLC in 2018 and serving as its CTO, Jeff Walsh had 20 years' experience in GIS and Remote Sensing. Jeff specializes in advanced terrain processing and extending GIS and Photogrammetric workflows through scripts and coding. Jeff is currently focused on building and refining unmanned sUAS systems capable of attaining design-grade vertical accuracy. Jeff has a Professional Masters in GIS & Remote Sensing from the University of Minnesota, is an ASPRS Certified Photogrammetrist, a GISP, and a Private Pilot.

Jeff was primarily responsible for the MnModel Phase 4 digital terrain model. He developed procedures for smoothing the effects of infrastructure and deriving terrain variables. He was also responsible for assembling the statewide mosaic of the Public Land Survey Plat maps and digitizing plat map features.

## Wood

*Carla Landrum, PhD*

Carla Landrum obtained her PhD from the University of Kentucky in 2013. She has over 11 years of experience in geostatistical (2D and 3D), geospatial, statistical, and time series modeling with a focus in resource

management and environmental remediation. Carla has authored and co-authored peer-reviewed publications and traveled internationally presenting technical seminars on the topics of agriculture water resource management, soil and groundwater remediation, geostatistics, and space-time data analysis. She was active in publishing industry-standard guidance as a member of the Interstate Technology and Regulatory Council - Geostatistics for Remediation Optimization Group and currently leads an international geostatistics professional practice network within Wood (formerly Amec Foster Wheeler).

Carla was the statistician for Phase 4 of MnModel. She migrated procedures developed earlier from S-Plus to R and developed additional procedures for both vegetation and archaeological predictive modeling.

### Carson Smith, MS

Carson Smith received his M.S. in Geography from Minnesota State University, Mankato. He participated in the development of the Landscape Model and was responsible for the preparation of soils data and soils variables.

## Minnesota State University, Mankato

### Andrew Brown, MS

Andy Brown is the archeological database developer and researcher in the Archeology, Geography, and Earth Sciences (AGES) Laboratory at Minnesota State University, Mankato. He earned a B.S. in Anthropology in 2010, a M.S. in Anthropology in 2016 at M.S.U., Mankato, and a graduate Geographic Information Science Certificate in 2016. He first worked on MnModel in 2013 and became fully involved in 2016. His research interests include archeological databases, data management, G.I.S. analysis, programming, and 3D modeling.



### Alexander Anton

Alec Anton is a master's student at Minnesota State University, Mankato. Before working on the MnModel project, he worked in cultural resources management from 2015-2017 and earned a B.A. in Anthropology in 2014 from the University of South Dakota. Alec contributed to preparing the archeological sites feature class, assisted with writing Python code to update the archeological surveys feature class, and ran regional predictive models using R Studio. His research interests include Midwestern and Great Plains archeology and GIS.

*Luke Burds*

Luke Burds earned his B.A. in Geography from the University of Wisconsin - Eau Claire in May 2018 and is pursuing a Master's of Science in Applied Anthropology at Minnesota State University Mankato.  He began working on the phase 4 MnModel project in August 2018 as a research assistant, just as data preparation was wrapping up and the modeling phase was beginning. His research interests include geophysical investigations for archeological prospecting, primarily using ground penetrating radar, GIS, and historical archeology.



*J.T. Salfer, MS*

J.T. Salfer worked on the preparation of the hydrographic data while a Master's student at Minnesota State University, Mankato.

*Woo Suk Jang, PhD*

Dr. Woo Suk Jang is an Associate Professor and Graduate Coordinator in the Geography Department at Minnesota State University, Mankato, where he teaches courses in Geographic Information Systems, Urban Geography, and Geospatial Technologies.  Woo supervised the student team at Minnesota State University, Mankato, and provided valuable technical and administrative assistance.

*Ron Schirmer, PhD*

Ron Schirmer received his PhD in Archaeology from the University of Minnesota.  He is a Professor in the Anthropology Department at Minnesota State University, Mankato (MSUM), where he teaches courses in anthropology, archaeology, and ethnobotany.  Ron put together the funding that allowed MSUM to participate in MnModel Phase 4 and provided valuable advice along the way.

## University of Minnesota, Twin Cities

*Gary W. Oehlert, PhD*

Dr. Gary Oehlert is a Professor in the School of Statistics at the University of Minnesota. Gary was the statistician for MnModel Phases 1-3. He developed updated statistical methods for MnModel Phase 4 that became the foundation of the procedures used to develop the final predictive models.

*Brian Shea*

Brian Shea assisted with the development of Phase 4 statistical methods.

## University of Minnesota, Duluth

*Stacey L. Stark, MS, GISP*

Stacey Stark is Director of the Geospatial Analysis Center (GAC) at the University of Minnesota Duluth (UMD). She obtains grants and contracts to support GAC while providing undergraduate GIS students with hands-on experience. Currents areas of work include mapping and analysis related to community resilience and hazard mitigation. http:\\scse.d.umn.edu\gac.

Stacey served as the project manager for development of the prehistoric hydrography layer for MnModel Phase 4. She trained and supervised students to georeference and digitize US General Land Office Survey (GLO) plat maps for the study regions to indicate water features present circa 1854-1858. Stark tested and implemented GIS methodology developed with partners and supervised the development of the GIS model for automated identification of the prehistoric hydrography.

*Patricia Farrell, PhD*

Pat Farrell received her PhD in Geography from the University of Cincinnati in 1997. She is a Professor of Geography at the University of Minnesota Duluth where she teaches physical geography courses, including Soils, Weather and Climate, and Environmental Conservation. Her research and field work has focused on signatures of human activity written in soil and sediment records.

Pat contributed to the hydrographic model by ascertaining how SSURGO soil data could be used to identify past wetlands, based on relict signatures, such as redoximorphic features, which were then filtered by hydric soil ratings and drainage classes.

*Susan C. Mulholland, PHD, RPA*

Sue Mulholland is President of the Duluth Archaeology Center (a private CRM company) and adjunct Assistant Professor of Anthropology at the University of Minnesota Duluth. Mulholland has served as Principal

Investigator/Project Director on numerous CRM projects in Minnesota, Wisconsin, and North Dakota, including Phase I surveys, Phase II site evaluations, Phase III site mitigations, site monitoring, and literature/background reviews.

Sue's experience with CRM projects provided input to the hydrographic model from the viewpoint of an end-user and of an archaeologist. She consulted on the development of the methodology for the identification of prehistoric water features and correlation to other environmental characteristics pertinent to predictive modeling of archaeological site locations.

## Sandia National Laboratories

### Devin White, PhD

Devin White (Ph.D. 2007, University of Colorado) leads the Autonomous Sensing and Perception Department at Sandia National Laboratories in New Mexico and is a Research Assistant Professor of Anthropology at the University of Tennessee, Knoxville. His research interests include machine intelligence, photogrammetry, remote sensing, computer vision, imaging science, geographic information science, computational social science, high performance computing, modeling human movement and visibility across landscapes, complex adaptive systems, and Southwestern archaeology.

For MnModel Phase 4, Devin applied a least cost pedestrian transportation model that he developed which does not require knowledge of the origins and destinations for travel, known as From Everywhere To Everwhere (FETE), to the entire state of Minnesota. The output was a map of the most likely locations for high traffic corridors, which could be used to inform the larger predictive model with respect to site location.

## Summit EnviroSolutions/DDMS

Molly O'Brien, Heidi Gaede, and Laurie Ollila (Summit Envirosolutions/DDMS) developed procedures for scanning and digitizing archaeological sites and surveys and directed the first major effort to digitize the archaeological data.

## ProWest & Associates, Inc.

### Annette Theroux

Annette Theroux has a GIS technical background and a 25-year track record in government GIS. She has partnered with MnDOT to lead numerous projects with an emphasis on aligning projects to business goals, managing stakeholders and managing change.  As President and CEO of Pro-West & Associates, a Minnesota-based GIS consulting firm, Annette served as Project Manager on the second project to digitize archaeological site and survey data.

### Brandon Crissinger

Brandon Crissinger is Vice-President and Chief Operations Officer at Pro-West & Associates.  Brandon has worked extensively with MnDOT and many of Pro-West's government clients to create robust GIS data and systems that solve complex organizational challenges.  As the company's Data Development Manager, he served as Project Supervisor on the second project to digitize archaeological site and survey data.  Brian was

responsible for refining the process, ensuring project requirements were met, and managing and training team members.