

Daily model of stream temperature for regional predictions

Daniel J. Hocking, Ben Letcher, and Kyle O'Neil

*Daniel J. Hocking (dhocking@usgs.gov), US Geological Survey, Conte Anadromous Fish Research Center, Turners Falls, MA, USA

Abstract

Set up the problem. Explain how you solve it. Tell what you find. Explain why it's the best thing ever.

Introduction

Options: Water Research, **Water Resources Research**, Freshwater Biology, Journal of Hydrology, Ecohydrology, Journal of Environmental Quality, Hydrobiologia, JAWRA

Temperature is a critical factor in regulating the physical, chemical, and biological properties of streams. Warming stream temperatures decrease dissolved oxygen, decrease water density, and alter the circulation and stratification patterns of the stream (refs). Biogeochemical processes such as nitrogen and carbon cycling are also temperature dependent and affect primary production, decomposition, and eutrophication (refs). Both physical properties and biogeochemical processes influence the suitability for organisms living in and using the stream habitat beyond just primary producers. Additionally, temperature can have direct effects on the biota, especially ectotherms such as invertebrates, amphibians, and fish [Xu *et al.*, 2010b, 2010a; Al-Chokhachy *et al.*, 2013; e.g., Kanno *et al.*, 2013]. Given commercial and recreational interests, there is a large body of literature describing the effects of temperature on fish, particularly the negative effects of warming temperatures on cool-water fishes such as salmonids. Finally, stream temperature can even affect electricity, drinking water, and recreation (see van Vliet *et al.* 2011). Therefore, understanding and predicting stream temperatures are important for a multitude of stakeholders.

Stream temperature models can be used for explanatory purposes (understanding factors and mechanisms affecting temperature) and for prediction. Predictions can be spatial and temporal including forecasting and hindcasting. Predictions across space are especially valuable because there is often a need for information at locations with little or no observed temperature data. For example, many states have regulations related to the management of streams classified as cold, cool, and warm waters (refs), but because of the tremendous number of headwater streams it is impossible to classify most streams based on observed data. Therefore, modeled stream temperature is needed to classify most streams for regulatory purposes. Forecasting can provide immediate information such as the expected temperature the next hour, day, or week as well as long-term information about expected temperatures months, years, and decades in the future. Hindcasting can be used to examine temperature

variability and trends over time and for model validation. Both forecasting and hindcasting are useful for understanding climate change effects on stream temperature regimes.

Given the importance of temperature in aquatic systems, it is not surprising that there are a variety of models and approaches to understanding and predicting stream temperature. Stream temperature models are generally divided into three categories: deterministic (also called process-based or mechanistic), stochastic, and statistical [Caissie, 2006; Benyahya *et al.*, 2007; Chang and Psaris, 2013]. Deterministic models are based on heat transfer and are often modeled using energy budgets [Caissie, 2006; Benyahya *et al.*, 2007]. The models require large amounts of detailed information on the physical properties of the stream and adjacent landscape as well as hydrology and meteorology. These models are useful for detailed site assessments and scenario testing. However, the data requirements preclude the models from being applied over large spatial extents.

Stochastic models attempt to combine pattern (seasonal and spatial trends) with the random deviations to describe and predict environmental data [Kiraly and Janosi, 2002; Sura *et al.*, 2006; Chang and Psaris, 2013]. Stochastic models of stream temperature generally rely on relationships between air and water temperature then with random noise and an autoregressive correlation, often decomposed by seasonal and annual components. These models are mostly commonly used to model daily temperature fluctuations because of their ability to address autocorrelation and approximate the near-random variability in environmental data [Caissie *et al.*, 2001; Kiraly and Janosi, 2002; Ahmadi-Nedushan *et al.*, 2007]. A limitation is that the physical processes driving temperature fluctuations are not elucidated with these models. They are generally used to describe characteristics and patterns in a system and to forecast these patterns in the future [Kiraly and Janosi, 2002]. Additionally, stochastic models rely on continuous, often long, time series from a single or a few locations. Inference cannot be made to other locations without assuming that the patterns and random deviations are similar at those locations.

As with stochastic models, statistical models generally rely on correlative relationships between air and water temperatures, but also typically include a variety of other predictor variables such as basin, landscape, and land-use characteristics. Statistical models are often linear with normally distributed error and therefore used at weekly or monthly time steps to avoid problems with temporal autocorrelation at shorter time steps (e.g. daily, hourly, sub-hourly). Parametric, nonlinear regression models have been developed to provide more information regarding mechanisms than traditional statistical models without the detail of physical deterministic models [Mohseni *et al.*, 1998]. Researchers have also developed geospatial regression models that account for spatial autocorrelation within dendritic stream networks [Isaak *et al.*, 2010; Peterson *et al.*, 2010, 2013]. However, due to the complexity of the covariance structure of network geostatistical models, they are best used for modeling single temperature values across space (e.g. summer maximum, July mean, etc.) rather than daily temperatures [Peterson *et al.*, 2007, 2010; Ver Hoef and Peterson, 2010]. Additionally, statistical machine learning techniques such as artificial neural networks have been used to model stream temperatures when unclear interactions, nonlinearities, and spatial relationships are of particular concern [Sivri *et al.*, 2007, 2009; DeWeber and Wagner, 2014].

In contrast with deterministic approaches, statistical models require less detailed site-level

data and therefore can be applied over greater spatial extents than process-based models. They also can describe the relationships between additional covariates and stream temperature, which is a limitation of stochastic models. These relationships can be used to understand and predict anthropogenic effects on stream temperature such as timber harvest, impervious development, and water control and release [Webb *et al.*, 2008]. Quantifying the relationship between anthropogenic effects, landscape characteristics, meteorological patterns, and stream temperature allows for prediction to new sites and times using statistical models. This is advantageous for forecasting and hindcasting to predict and understand climate change effects on stream temperatures. This is critical because not all streams respond identically to air temperature changes and the idiosyncratic responses may be predicted based interactions of known factors such as flow, precipitation, forest cover, basin topology, impervious surfaces, soil characteristics, geology, and impoundments [Webb *et al.*, 2008].

We describe a novel statistical model of daily stream temperature that incorporates features of stochastic models and apply it to a large geographic area. This model handles time series data of widely varying duration from many sites using a hierarchical mixed model approach to account for autocorrelation at specific locations within watersheds. It incorporates catchment, landscape, land-use, and meteorological covariates for explanatory and predictive purposes. It includes an autoregressive function to account for temporal autocorrelation in the time series, a challenge with other statistical models at fine temporal resolution. Additionally, our hierarchical Bayesian approach readily allows for complete accounting of uncertainty. We use the model to predict daily stream temperature across the northeastern United States over a 34-year time record.

Methods

Study area

Map of data locations: size = amount of data, color/shape = training-validation
- Kyle, Ana, or Matt make? See deWeber 2014 for example

Water temperature data

We gathered stream temperature data from state and federal agencies, individual academic researchers, and non-governmental organizations (NGOs). The data were collected using automated temperature loggers. The temporal frequency of recording ranged from every 5 minutes to once per hour. This data is consolidated in a PostgreSQL database linked to a web service at <http://www.ecosheds.org>. Data collectors can upload data at this website and choose whether to make the data publicly available or not. The raw data is stored in the database and users can flag problem values and time series. For our analysis, we performed some automated and visual QAQC on the sub-daily values, summarized to mean daily temperatures and performed additional QAQC on the daily values. The QAQC was intended to flag and remove values associated with logger malfunctions, out-of-water events (including first and last days when loggers were recording but not yet in streams), and days

with incomplete data which would alter the daily mean. We developed an R (ref) package for analyzing stream temperature data from our database, including the QAQC functions which can be found at <https://github.com/Conte-Ecology/conteStreamTemperature>. The R scripts using these functions for our analysis are available at https://github.com/Conte-Ecology/conteStreamTemperature_northeast.

Stream reach (stream section between any two confluences) was our finest spatial resolution for the analysis. In the rare case where we had multiple logger locations within the same reach recording at the same time, we used the mean value from the loggers for a given day. In the future, with sufficient within reach data, it would be possible to use our modeling framework to also estimate variability within reach.

Stream network delineation

Meteorological (, Climatic,) and landscape data - separate landscape if use climate data for future projections

Table of Variables - include part of the model they're in (fixed, site, huc, year)

Statistical model

Statistical models of stream temperature often rely on the close relationship between air temperature and water temperature. However, this relationship breaks down during the winter in temperature zones, particularly as streams freeze, thereby changing their thermal and properties. Many researchers and managers are interested in the non-winter effects of temperature. The winter period, when phase change and ice cover alter the air-water relationship, differs in both time (annually) and space. We developed an index of air-water synchrony ($Index_{sync}$) so we can model the portion of the year that it not affected by freezing properties. The index is the difference between air and observed water temperatures divided by the water temperature plus 0.000001 to avoid division by zero.

We calculate the $Index_{sync}$ for each day of the year at each site for each year with observed data. We then calculate the 99.9% confidence interval of $Index_{sync}$ for days between the 125 and 275 days of the year (05 May and 02 October). Then moving from the middle of the year (day 180) to the beginning of the year, we searched for the first time when 10 consecutive days were not within the 99.9% CI. This was selected as the spring breakpoint. Similarly moving from the middle to the end of the year, the first event with fewer than 16 consecutive days within the 99.9% CI was assigned as the autumn breakpoint. Independent breakpoints were estimated for each site-year combination. For site-years with insufficient data to generate continuous trends and confidence intervals, we used the mean break points across years for that site. If there was not sufficient local site information, we used the mean breakpoints from the smallest hydrologic unit the site is nested in (i.e. check for mean from HUC12, then HUC10, HUC8, etc.). More details regarding the identification of the synchronized period can be found in Letcher et al. (*in review*). The portion of the year between the spring and autumn breakpoints was used for modeling the non-winter, approximately ice-free stream temperatures.

We used a generalized linear mixed model to account for correlation in space (stream reach nested within HUC8). This allowed us to incorporate short time series as well as long time series from different reaches and disjunct time series from the same reaches without risk of pseudoreplication (ref: Hurlbert). By limited stream drainage area to $<400 \text{ km}^2$ and only modeling the synchronized period of the year, we were able to use a linear model, avoiding the non-linearities that occur at very high temperatures due to evaporative cooling and near 0 C due to phase change (ref: mohseni).

We assumed stream temperature measurements were normally distributed following,

$$t_{h,r,y,d} \sim \mathcal{N}(\mu_{h,r,y,d}, \sigma)$$

where $t_{h,r,y,d}$ is the observed stream water temperature at the reach (r) within the sub-basin identified by the 8-digit Hydrologic Unit Code (HUC8; h) for each day (d) in each year (y). We describe the normal distribution based on the mean ($\mu_{h,r,y,d}$) and standard deviation (σ) and assign a vague prior of $\sigma = 100$. The mean temperature is modelled to follow a linear trend

$$\omega_{h,r,y,d} = X_0 B_0 + X_{h,r} B_{h,r} + X_h B_h + X_y B_y$$

but the expected mean temperature ($\mu_{h,r,y,d}$) is also adjusted based on the residual error from the previous day

$$\mu_{h,r,y,d} = \begin{cases} \omega_{h,r,y,d} + \delta_r(t_{h,r,y,d-1} - \omega_{h,r,y,d-1}) & \text{for } t_{h,r,y,d-1} \text{ is real} \\ \omega_{h,r,y,d} & \text{for } t_{h,r,y,d-1} \text{ is not real} \end{cases}$$

where δ_r is an autoregressive [AR(1)] coefficient that varies randomly by reach and $\omega_{h,r,y,d}$ is the expected temperature before accounting for temporal autocorrelation in the error structure.

X_0 is the $n \times K_0$ matrix of predictor values. B_0 is the vector of K_0 coefficients, where K_0 is the number of fixed effects parameters including the overall intercept. We used **latitude, longitude, upstream drainage area, percent forest cover, elevation, surficial coarseness classification, percent wetland area, upstream impounded area, and an interaction of drainage area and air temperature**. We assumed the following distributions and vague priors for the fixed effects coefficients

$$B_0 \sim \mathcal{N}(0, \sigma_{k_0}), \text{ for } k_0 = 1, \dots, K_0,$$

$$B_0 = \beta_0^1, \dots, \beta_0^{K_0} \sim \mathcal{N}(0, 100)$$

$$\sigma_{k_0} = 100$$

??The effects of air temperature on the day of observation (d) and mean air temperature over the previous 7 days varied randomly with reach nested within HUC8, as did precipitation, the previous 30-day precipitation mean, and the interactions of air temperature and precipitation (all 4 combinations).??

$B_{h,r}$ is the $R \times K_R$ matrix of regression coefficients where R is the number of unique reaches and K_R is the number of regression coefficients that vary randomly by reach within HUC8. We assumed prior distributions of

$$B_{h,r} \sim \mathcal{N}(0, \sigma_{k_r}), \text{ for } k_r = 1, \dots, K_R,$$

$$\sigma_{r_0} = 100$$

X_h is the matrix of parameters that vary by HUC8. We allowed for correlation among the effects of these HUC8 coefficients as described by Gelman and Hill [2007].

B_h is the $H \times K_H$ matrix of coefficients where H is the number of HUC8 groups and K_H is the number of paramaters that vary by HUC8 including a constant term. In our model, $K_H = K_R$ and we assumed priors distributions of

$$B_h \sim \mathcal{N}(M_h, \Sigma_{B_h}), \text{ for } h = 1, \dots, H$$

where M_h is a vector of length K_H and Σ_{B_h} is the $K_H \times K_H$ covariance matrix.

$$M_h \sim MVN(\mu_{1:K_h}^h, \sigma_{1:K_h}^h)$$

$$\mu_1^h = 0; \mu_{2:K_h}^h \sim \mathcal{N}(0, 100)$$

$$\Sigma_{B_h} \sim \text{Inv-Wishart}(\text{diag}(K_h), K_h + 1)$$

Similarly, we allowed the some effects of some parameters (X_y) to vary randomly by year with potential correlation among the coefficients. The intercept, day of the year (day), day^2 , and day^3 all varied randomly with year such that $K_y = 4$. We assumed prior distributions of

$$B_y \sim \mathcal{N}(M_y, \Sigma_{B_y}), \text{ for } y = 1, \dots, Y$$

where M_y is a vector of length K_Y and Σ_{B_y} represents the $K_Y \times K_Y$ covariance matrix.

$$M_y \sim MVN(\mu_{1:K_y}^y, \sigma_{1:K_y}^y)$$

$$\mu_1^y = 0; \mu_{2:K_y}^y \sim \mathcal{N}(0, 100)$$

$$\Sigma_{B_y} \sim \text{Inv-Wishart}(\text{diag}(K_y), K_y + 1)$$

To estimate all the parameters and their uncertainties, we used a Bayesian analysis with a Gibbs sampler implemented in JAGS (ref) through R (ref) using the rjags package (ref). This approach was beneficial for hierarchical model flexibility and tractability for large datasets. We used vague priors for all parameters so all inferences would be based on the data.

Model validation

To validate our model, we held out 10% of subbasins (HUC8s) at random. We also held out 10% of remaining stream reaches with observed temperature data at random. Additionally, we excluded all 2010 data because it was an especially warm summer across the northeastern U.S. Therefore, we will be able to evaluate how well our model predicts across space and time. This included reaches with no data located within subbasins with and without data and how well the model predicts in warm years without data, which will be important if using this model with future climate predictions. The most challenging validation scenario was at reaches within HUC8s without any data in a year without any data. In total, **XX%** of observations and **XX%** of reaches were held out for validation.

Derived metrics

Climate change projections (future paper?)

Results

Explain what you found. Avoid blind *P-values* (or avoid *P-values* altogether)

Discussion

what we found

model separates uncertainty in estimates and predictions from variability across space and time. The random site, HUC, and year effects explicitly address spatial and temporal variability, allowing for more proper accounting of uncertainty.

lots of sensors because relatively cheap and easy to collect, but varying lengths of time at different sites. Our model incorporates sites with any length of time (a few days to decades). Sites with little data contribute less to the model but do provide some local and spatial information. The more data a location has the more informative so there is less shrinkage to the mean values. Sites with no data can be predicted based on covariate values and HUC-level random effects but do not get site-specific coefficient effects.

Disagreement (conflicting evidence? confused terminology) regarding the drivers of stream temperature

Acknowledgements

Thanks to Ethan White, Karthik Ram, Carl Boettiger, Ben Morris, and [Software Carpentry](#) for getting me started with the skills needed to [ditch MS Word](#) and produce more reproducible research.

Tables

Table 1: Example Markdown table

Name	col2	col3	col4	col5	Comments
Brook Trout	1	big	few	2.2	Ecology & life history data associated with trout
<i>Desmognathus fuscus</i>	100	small	many	0.3	Widespread salamander species

Name	Phone

example created with pander

parameter	mean	sd	sig
Intercept	17.7035	0.2486	*
AirT	2.1721	0.1472	*
7-day AirT	1.5792	0.1362	*
Development	0.1709	0.0559	*
Agriculture	-0.0583	0.0665	
Impoundment Area	0.3678	0.0660	*
AirT x Impoundment	-0.0288	0.0229	
AirT x Forest	-0.0176	0.0265	
AirT x Prcp2 x DA	-0.0036	0.0016	*
AirT x prcp30 x DA	-0.0020	0.1666	
Day	0.0506	0.1070	
Day ²	-0.5141	0.0887	*
Day ³	-0.0834	0.0778	
AR1	0.7696	0.0073	*

example created with stargazer

% Table created by stargazer v.5.1 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu % Date and time: Fri, Apr 10, 2015 - 2:29:23 PM

Table 4:

parameter	mean	sd	sig
Intercept	17.7035	0.2486	*
AirT	2.1721	0.1472	*
7-day AirT	1.5792	0.1362	*
Development	0.1709	0.0559	*
Agriculture	-0.0583	0.0665	
Impoundment Area	0.3678	0.0660	*
AirT x Impoundment	-0.0288	0.0229	
AirT x Forest	-0.0176	0.0265	
AirT x Prcp2 x DA	-0.0036	0.0016	*
AirT x prcp30 x DA	-0.0020	0.1666	
Day	0.0506	0.1070	
Day^2	-0.5141	0.0887	*
Day^3	-0.0834	0.0778	
AR1	0.7696	0.0073	*

Figures

Figure 1. Example of adding a figure.

Literature Cited

Ahmadi-Nedushan, B., A. St-Hilaire, T. B. M. J. Ouarda, L. Bilodeau, E. Robichaud, N. Thiemonge, and B. Bobee (2007), Predicting river water temperatures using stochastic models : case study of the Moisie River (Quebec , Canada), *Hydrological Processes*, 34, 21–34, doi:[10.1002/hyp](https://doi.org/10.1002/hyp).

Al-Chokhachy, R., J. Alder, S. Hostetler, R. Gresswell, and B. Shepard (2013), Thermal controls of Yellowstone cutthroat trout and invasive fishes under climate change, *Global change biology*, 19(10), 3069–81, doi:[10.1111/gcb.12262](https://doi.org/10.1111/gcb.12262).

Benyahya, L., D. Caissie, A. St-Hilaire, T. B. M. J. Ouarda, and B. Bobee (2007), A review of statistical water temperature models, *Canadian Water Resources Journal*, 32(3), 179–192.

Caissie, D. (2006), The thermal regime of rivers: a review, *Freshwater Biology*, 51(8), 1389–1406, doi:[10.1111/j.1365-2427.2006.01597.x](https://doi.org/10.1111/j.1365-2427.2006.01597.x).

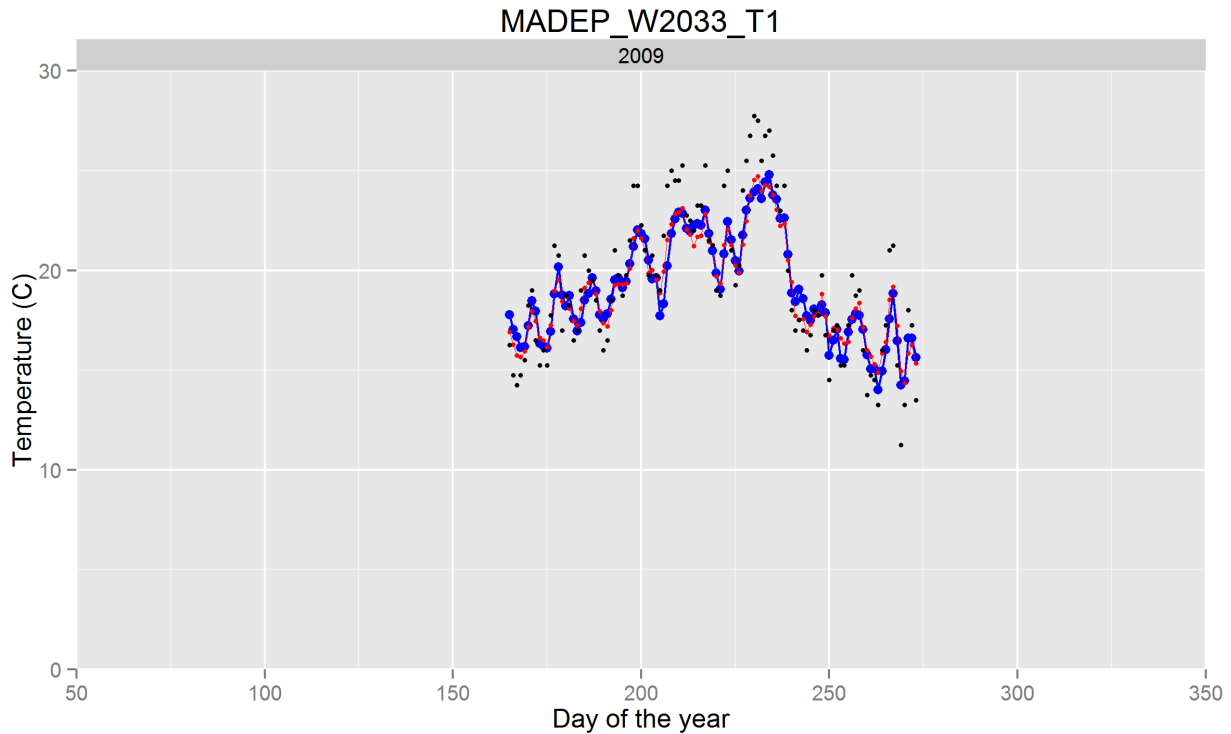


Figure 1: Figure1

Caissie, D., N. El-jabi, and M. G. Satish (2001), Modelling of maximum daily water temperatures in a small stream, *Journal of Hydrology*, 251(2001), 14–28.

Chang, H., and M. Psaris (2013), Local landscape predictors of maximum stream temperature and thermal sensitivity in the Columbia River Basin, USA., *The Science of the total environment*, 461-462, 587–600, doi:[10.1016/j.scitotenv.2013.05.033](https://doi.org/10.1016/j.scitotenv.2013.05.033).

DeWeber, J. T., and T. Wagner (2014), Predicting Brook Trout Occurrence in Stream Reaches throughout their Native Range in the Eastern United States, *Transactions of the American Fisheries Society*, 144(1), 11–24, doi:[10.1080/00028487.2014.963256](https://doi.org/10.1080/00028487.2014.963256).

Gelman, A., and J. Hill (2007), *Data analysis using regression and multilevel/hierarchical models*, Cambridge University Press, New York.

Isaak, D. J., C. H. Luce, B. E. Rieman, D. E. Nagel, E. E. Peterson, D. L. Horan, S. Parkes, and G. L. Chandler (2010), Effects of climate change and wildfire on stream temperatures and salmonid thermal habitat in a mountain river network., *Ecological applications : a publication of the Ecological Society of America*, 20(5), 1350–1371, doi:[papers2://publication/uuid/8973E71F-5D23-47C7-A085-2AB46FFD8BF0](https://doi.org/papers2://publication/uuid/8973E71F-5D23-47C7-A085-2AB46FFD8BF0).

Kanno, Y., J. Vokoun, and B. Letcher (2013), Paired stream-air temperature measurements reveal fine-scale thermal heterogeneity within headwater Brook Trout stream networks, *River Research and Applications*, 30(6), 745–755, doi:[10.1002/rra](https://doi.org/10.1002/rra).

Kiraly, A., and I. Janosi (2002), Stochastic modeling of daily temperature fluctuations,

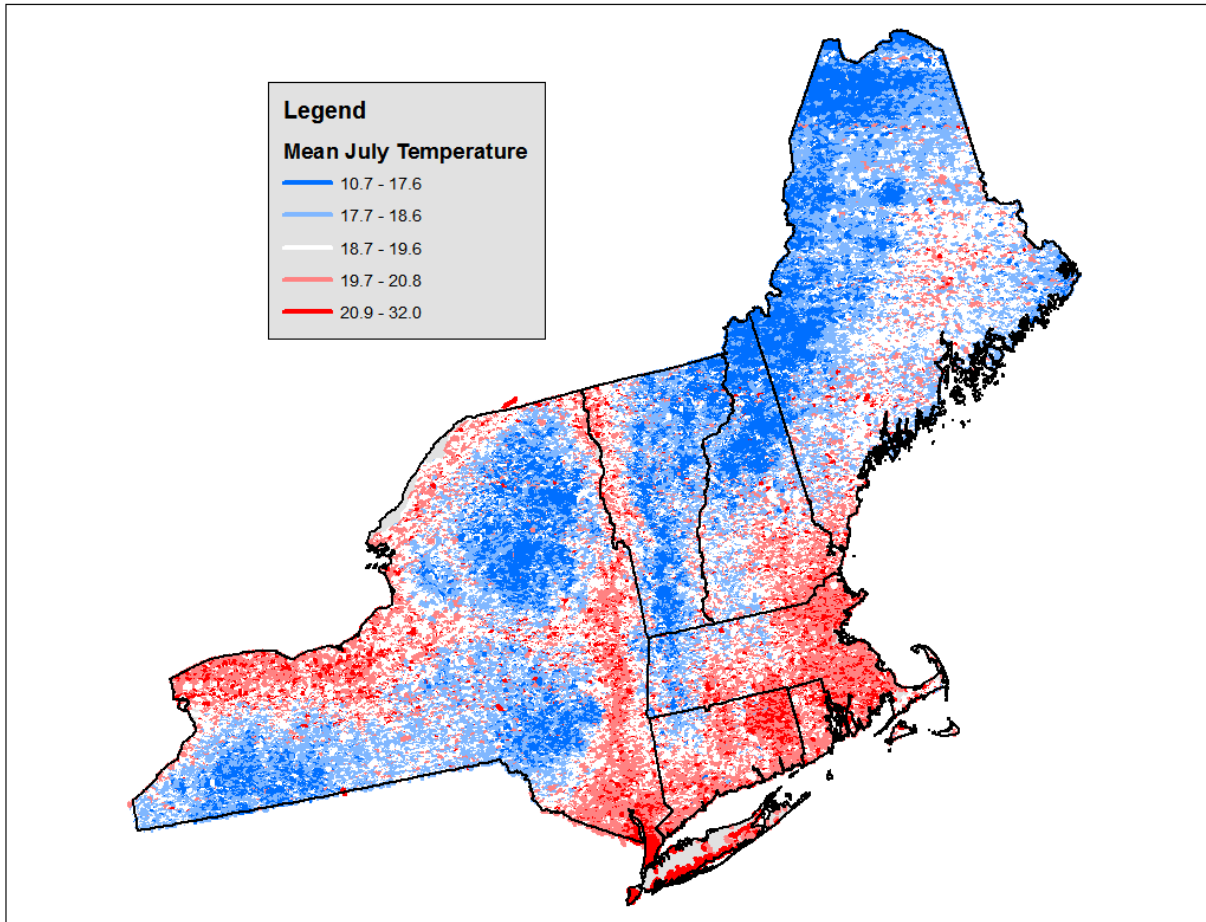


Figure 2: Figure2

Physical Review E, 65(5), 1–6, doi:[10.1103/PhysRevE.65.051102](https://doi.org/10.1103/PhysRevE.65.051102).

Mohseni, O., H. G. Stefan, and T. R. Erickson (1998), A nonlinear regression model for weekay stream temperatures, *Water Resources Research*, 34(10), 2685–2692.

Peterson, E. E., D. M. Theobald, and J. M. Ver Hoef (2007), Geostatistical modelling on stream networks: developing valid covariance matrices based on hydrologic distance and stream flow, *Freshwater Biology*, 52(2), 267–279, doi:[10.1111/j.1365-2427.2006.01686.x](https://doi.org/10.1111/j.1365-2427.2006.01686.x).

Peterson, E. E., J. M. V. Hoef, and M. Jay (2010), A mixed-model moving-average approach to geostatistical modeling in stream networks, *Ecology*, 91(3), 644–651.

Peterson, E. E. et al. (2013), Modelling dendritic ecological networks in space: an integrated network perspective., *Ecology letters*, 16(5), 707–19, doi:[10.1111/ele.12084](https://doi.org/10.1111/ele.12084).

Sivri, N., N. Kilic, and O. N. Ucan (2007), Estimation of stream temperature in Firtina Creek (Rize-Turkiye) using artificial neural network model, *Journal of Environmental Biology*, 28(1), 67–72.

Sivri, N., H. K. Ozcan, O. N. Ucan, and O. Akincilar (2009), Estimation of Stream Temperature in Degirmendere River (Trabzon- Turkey) Using Artificial Neural Network Model, *Turkish Journal of Fisheries and Aquatic Sciences*, 9, 145–150, doi:[10.4194/trjfas.2009.0204](https://doi.org/10.4194/trjfas.2009.0204).

Sura, P., M. Newman, and M. A. Alexander (2006), Daily to Decadal Sea Surface Temperature Variability Driven by State-Dependent Stochastic Heat Fluxes, *Journal of Physical Oceanography*, 36, 1940–1958.

Ver Hoef, J. M., and E. E. Peterson (2010), A Moving Average Approach for Spatial Statistical Models of Stream Networks, *Journal of the American Statistical Association*, 105(489), 6–18, doi:[10.1198/jasa.2009.ap08248](https://doi.org/10.1198/jasa.2009.ap08248).

Webb, B., D. Hannah, R. D. Moore, L. E. Brown, and F. Nobilis (2008), Recent advances in stream and river temperature research, *Hydrological Processes*, 918, 902–918, doi:[10.1002/hyp](https://doi.org/10.1002/hyp).

Xu, C., B. H. Letcher, and K. H. Nislow (2010a), Context-specific influence of water temperature on brook trout growth rates in the field, *Freshwater Biology*, 55(11), 2253–2264, doi:[10.1111/j.1365-2427.2010.02430.x](https://doi.org/10.1111/j.1365-2427.2010.02430.x).

Xu, C. L., B. H. Letcher, and K. H. Nislow (2010b), Size-dependent survival of brook trout *Salvelinus fontinalis* in summer: effects of water temperature and stream flow, *Journal of Fish Biology*, 76(10), 2342–2369, doi:[10.1111/j.1095-8649.2010.02619.x](https://doi.org/10.1111/j.1095-8649.2010.02619.x).