

# Predicting river water temperatures using stochastic models: case study of the Moisie River (Québec, Canada)

Behrouz Ahmadi-Nedushan,<sup>1\*</sup> André St-Hilaire,<sup>1</sup> Taha B. M. J. Ouarda,<sup>1</sup> Laurent Bilodeau,<sup>2</sup>  
Élaine Robichaud,<sup>2</sup> Nathalie Thiémonge<sup>2</sup> and Bernard Bobée<sup>1</sup>

<sup>1</sup> Chair in Statistical Hydrology, INRS-ETE, Université du Québec, 490 de la Couronne, Québec G1K 9A9, Canada

<sup>2</sup> Hydro-Québec, 855 Ste-Catherine Street East, Montreal, Québec H2L 1A4, Canada

## Abstract:

Successful applications of stochastic models for simulating and predicting daily stream temperature have been reported in the literature. These stochastic models have been generally tested on small rivers and have used only air temperature as an exogenous variable. This study investigates the stochastic modelling of daily mean stream water temperatures on the Moisie River, a relatively large unregulated river located in Québec, Canada. The objective of the study is to compare different stochastic approaches previously used on small streams to relate mean daily water temperatures to air temperatures and streamflow indices. Various stochastic approaches are used to model the water temperature residuals, representing short-term variations, which were obtained by subtracting the seasonal components from water temperature time-series. The first three models, a multiple regression, a second-order autoregressive model, and a Box and Jenkins model, used only lagged air temperature residuals as exogenous variables. The root-mean-square error (RMSE) for these models varied between 0.53 and 1.70 °C and the second-order autoregressive model provided the best results.

A statistical methodology using best subsets regression is proposed to model the combined effect of discharge and air temperature on stream temperatures. Various streamflow indices were considered as additional independent variables, and models with different number of variables were tested. The results indicated that the best model included relative change in flow as the most important streamflow index. The RMSE for this model was of the order of 0.51 °C, which shows a small improvement over the first three models that did not include streamflow indices. The ridge regression was applied to this model to alleviate the potential statistical inadequacies associated with multicollinearity. The amplitude and sign of the ridge regression coefficients seem to be more in agreement with prior expectations (e.g. positive correlation between water temperature residuals of different lags) and make more physical sense. Copyright © 2006 John Wiley & Sons, Ltd.

KEY WORDS stream temperature; statistical analysis; stochastic modelling; ridge regression; Moisie River

Received 11 July 2005; Accepted 3 February 2006

## INTRODUCTION

Stream temperature is an important abiotic variable in aquatic habitat studies and has a great influence on many facets of the ecosystem (Petts, 2000). It affects many chemical and biological processes, such as dissolved oxygen concentration and growth of aquatic organisms (Edwards *et al.*, 1979; Bovee, 1982). Water temperature can be one of the factors limiting the potential fish habitat in a stream (Bovee, 1982), and changes in the stream thermal regime can significantly impact fish distribution, growth, mortality, production, habitat use and community dynamics (e.g. Edwards *et al.*, 1979; Elliott *et al.*, 1995). For instance, Hodgson and Quinn (2002) have demonstrated that the triggering of the spawning period for sockeye salmon (*Onchorhynchus nerka*) on the northwestern coast of the USA is strongly modulated by water temperature and that, when a threshold of 19 °C is

reached, spawning is interrupted as individual fish seek for thermal refuge.

The thermal regime of a watercourse is governed by the interaction of natural environment processes and human activities. Examples of the latter include thermal pollution (e.g. plant effluents), regulation (Webb and Walling, 1993) and deforestation, the latter of which has been linked to temperature rises in certain cases (Brown and Krygier, 1970; Johnson and Jones, 2000). Given the potential impact of temperature on a lotic ecosystem, it is essential to provide efficient water temperature predictive tools to water resource managers. As stream temperature standards are being developed, it is of the utmost importance that managers be provided with modelling tools that can be implemented relatively easily. An increasing number of rivers are being monitored, and these will soon have sufficiently long water temperature and flow time-series to enable development of stochastic models such as the ones presented in this study. Existing stream water temperature models are usually grouped in two broad categories: stochastic/statistical models and deterministic models (St-Hilaire *et al.*, 2000). Deterministic

\*Correspondence to: Behrouz Ahmadi-Nedushan, Department of Statistical Hydrology, INRS-EET, Université du Québec, 490 de la Couronne, Québec G1K 9A9, Canada. E-mail: Behrouz\_Nedushan@ete.inrs.ca

models are based on a mathematical representation of the underlying physics and energy budget (Morin and Couillard, 1990; St-Hilaire *et al.*, 2000). Examples of deterministic models include the US Fish and Wildlife SNTMP model (Bartholow, 1989) or its derivative, the SSTEMP model (Bartholow, 1999), the CEQUEAU hydrological and water temperature model (Morin and Couillard, 1990; St-Hilaire *et al.*, 2003), as well as a number of simpler, less generic models (e.g. Sinokrot and Stefan, 1993; Gu *et al.*, 1998; Gu and Li, 2002). Deterministic models use detailed time-series of relevant meteorological factors, such as solar radiation and wind velocity, as inputs; these are used in heat budget equations to calculate thermal exchange between the atmosphere and stream. Deterministic models can be most useful when the user wants to simulate modifications to certain terms of the heat budget, and are very useful for analysing profoundly modified watercourses and for comparing different impact scenarios (St-Hilaire *et al.*, 2000). A potential drawback of the deterministic models is in the relative complexity of development and application, as a greater number of inputs are often required, including stream geometry, hydrology and meteorology. These inputs are not always available (Caissie *et al.*, 1998). Therefore, deterministic heat budget models are not always ideal for some applications.

As an alternative, statistical or stochastic methods are also proposed for stream water temperature modelling. Simple regression-based models have been successfully used to model mean water temperature as a function of one (usually air temperature) or more independent variables (Webb and Nobilis, 1997; Mitchell, 1999). The potential advantage of a statistical approach is that it often requires less input data (e.g. only air temperature time-series in many cases). A number of statistical approaches have been tested in the past. Stefan and Preud'homme (1993) examined linear relationships between stream temperatures and air temperatures for 11 streams in the central USA. Water temperatures were shown to respond to air temperatures with time lags ranging from a few hours for small streams to 7 days for large rivers up to 5 m in depth. More recently, Pilgrim *et al.* (1998) used linear regression to relate stream water temperature to air temperature for 39 Minnesota streams. Equations were derived for daily, monthly and annual mean temperatures. Mohseni *et al.* (1998) developed a nonlinear regression model to predict average weekly stream temperatures at different locations in the USA. The nonlinear function was developed separately for the warming season and the cooling season to take heat storage effects (hysteresis) into account. They reported that, at high air temperatures ( $>25^{\circ}\text{C}$ ), the water–air temperature relationship derived from weekly mean values departs from linearity. This was attributed to increases in the moisture-holding capacity of the atmosphere, which promotes greater evaporation from the water surface and, in turn, increases evaporative cooling of the watercourse, together with enhanced back radiation as water temperatures rise (Mohseni *et al.*, 1998, 1999).

Although most statistical models only use air temperature as an exogenous variable (e.g. Stefan and Preud'homme, 1993; Mohseni *et al.*, 1998), the influence of discharge on water temperature has also been recognized. Hockey *et al.* (1982) observed that, in the Hurunui River (New Zealand), water temperature increased by about  $0.1^{\circ}\text{C}$  for each  $1\text{ m}^3\text{ s}^{-1}$  decrease in flow abstracted for irrigation. Webb *et al.* (2003) investigated the nature of the water–air temperature relationship, and its moderation by discharge for different catchments in the southwest of England. They concluded that the relationship between water and air temperatures was stronger for flows below median levels, and streamflow had a greater impact in accounting for water temperature variations in larger catchments.

The simple regression models generally perform better at weekly and monthly scales than at a daily scale (Pilgrim *et al.*, 1998; Erickson and Stefan, 2000), and errors associated with these models for daily data can be high. For example, Erickson and Stefan (2000) reported results of linear regression analysis on 39 streams each in Minnesota and Oklahoma during open water periods. At a daily time scale, the goodness of fit in terms of  $R^2$  and root-mean-square error (RMSE) were 0.71 and  $3.23^{\circ}\text{C}$  respectively for Minnesota streams and 0.77 and  $3.05^{\circ}\text{C}$  respectively for Oklahoma streams. The stochastic models are often preferred for daily time steps (Caissie *et al.*, 2001).

In this approach, the seasonal component of the signal is first removed and then time-series models (e.g. Box–Jenkins, autoregressive moving average, etc.) are fitted to water temperature residuals that represent short-term stream temperature variations. In many cases, the seasonal variation of stream water temperatures can be modelled by a Fourier series or even a simple sinusoidal function (Caissie *et al.*, 1998). The stochastic modelling of the residuals takes into account the autocorrelation structure of stream water temperature and can also account for the correlation with external variables (e.g. air temperature). Early applications of this approach include Kothandaraman (1971), who used air temperature residuals of the previous 2 days to model water temperature residuals. Cluis (1972) used a second-order Markov chain to model the short-term water temperature variations. This model also included air temperature residuals as an independent variable. Caissie *et al.* (1998) compared three different stochastic approaches to model mean daily water temperatures in a relatively small stream ( $50\text{ km}^2$  drainage area) using air temperature as the independent variable. The models of Cluis (1972) and Kothandaraman (1971) were compared with a Box–Jenkins model on the residuals. Caissie *et al.* (2001) compared a stochastic approach and regression approach (Mohseni *et al.*, 1998) to predict maximum water temperatures in the same river. The stochastic approach provided the better results at a daily time-scale and the results of the regression model showed good agreement only on a weekly basis.

Earlier work on stochastic models applied on small rivers to predict mean daily river water temperatures include only air temperatures as exogenous variables (Kothandarman, 1971; Cluis, 1972; Caissie *et al.*, 1998). Therefore, there is a need to examine the performance of these methods on larger rivers and to include streamflow as an independent variable. It must be noted that most studies concerned with the inclusion of streamflow have used either regression methods (e.g. Neumann *et al.*, 2003) or deterministic approaches (e.g. Gu and Li, 2002), and have not been based on stochastic approaches. Therefore, the objectives of the present study are: (1) To verify the efficiency of several stochastic models to predict daily mean water temperature using only air temperature as an exogenous variable in a relatively large watercourse (as opposed to a stream), using data from the Moisie River in Québec (Canada). (2) To explore the possibility of including variants of flow as additional independent variables in stochastic models. Air temperatures and flow-derived variables (i.e. streamflow indices) are used as independent variables to estimate daily stream temperature. (3) To alleviate the potential statistical inadequacies associated with collinearity by using the ridge regression.

#### DATA AND STUDY SITE

Time-series of water temperature and river flow have been measured at different locations along the Moisie

River, located on Québec's North Shore, Canada (Figure 1). The Moisie River is a relatively large river with a drainage area of 19 871 km<sup>2</sup> and annual mean discharge of 466 m<sup>3</sup> s<sup>-1</sup>. Its source is Lake Menistouc and it runs for 363 km prior to discharging into the St Lawrence Estuary. The Moisie River is one of the most important spawning grounds for Atlantic salmon (*Salmo salar*) in eastern Canada and is praised by the angling community as one of the most important salmon rivers in Québec because of the high average weight of individual adult salmon. Hourly water temperatures used in this study were gathered by Hydro-Québec, the main provincial hydroelectric authority, during the period of 1989–1998 at station TMOI0008 located at latitude 52°12'41"N and longitude 66°48'36"W at an elevation of 380 m (Figure 1). The drainage area for the portion of the catchment located above the temperature monitoring station is approximately 2500 km<sup>2</sup>. There were numerous missing values during the first part of the observation period (1989–1992); therefore, the analyses focused on the latter portion of the time-series. The daily mean water temperatures for 6 years (1993–1998) were extracted from the continuous records as averages of 24 hourly observations. Daily air temperatures were obtained for station 7047910 at Sept-Iles, operated jointly by Environment Canada and Transport Canada, and located at longitude 50°13'N and latitude 66°16'W at an elevation 50 m approximately 200 km south of the water temperature station TMOI0008. This station was the closest station

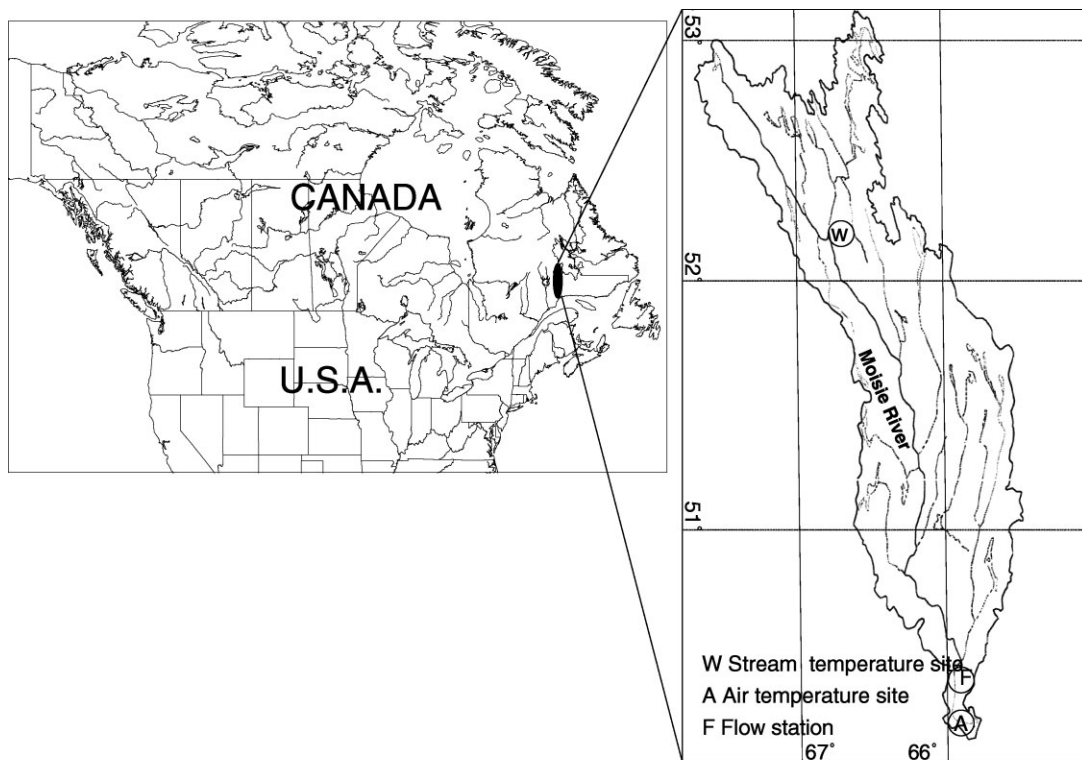


Figure 1. Map of Moisie showing the location of the water temperature site, the meteorological station and flow measuring station

to the water temperature station TMOI0008 for which the complete concomitant time-series were available. Streamflows for the same period were only available at one location, namely Québec hydrometric station 072 301 (Québec, Department of the Environment), which is located in the downstream portion of the drainage area at latitude  $50^{\circ}21'01''\text{N}$  and longitude  $66^{\circ}11'25''\text{W}$  and gauges almost 96% of the drainage basin (i.e. about  $19\,000\text{ km}^2$ ). Although this station is located in the downstream portion of the Moisie River, these flow data were used in the analysis to demonstrate the feasibility of including streamflow as an independent variable in larger rivers. A greater proximity between temperature and flow stations would likely improve the relationship, as flow may explain a greater proportion of the variance. This will have to be verified in alternate sites when the data become available.

## METHODS

The stochastic modelling of stream water temperatures  $T_w$  consists of separating the water temperatures into two different components: the long-term periodic or annual component  $T_A$  and the short-term component or residuals  $R_w$ .

$$T_w(t) = T_A(t) + R_w(t) \quad (1)$$

Previous studies (e.g. Caissie *et al.*, 1998) have compared the use of a single invariant sine function with more complex models (e.g. Fourier series) for modelling of the periodic component. It was found that both models performed equally well, since a departure from the inter-annual (seasonal) signal for a given year is often also found in the air temperature residuals, and thus taken into account in the model by considering air temperature residuals as exogenous variables. The periodic seasonal component is computed with the simple approach suggested by Cluis (1972), by fitting a sinusoidal function to the interannual daily mean time-series:

$$T_A(t) = a + b \sin \left[ \frac{2\pi}{365}(t + t_0) \right] \quad (2)$$

where  $T_A(t)$  is the seasonal component of a temperature time-series;  $a$ ,  $b$  and  $t_0$  are fitted coefficients. The first stochastic model (SM1) used is based on the work of Kothandaraman (1971), in which concomitant and lagged air temperature residuals are used as independent variables:

$$R_w(t) = \beta_1 R_A(t) + \beta_2 R_A(t-1) + \beta_3 R_A(t-2) \quad (3)$$

where  $\beta_1$ ,  $\beta_2$  and  $\beta_3$  are regression coefficients.  $R_A(t)$ ,  $R_A(t-1)$  and  $R_A(t-2)$  are air temperature residuals at times  $t$ ,  $(t-1)$  and  $(t-2)$ . The second stochastic model (SM2) is that of Cluis (1972), who used a second-order Markov chain to model the water temperature residuals. This model also uses the air temperature residuals:

$$R_w(t) = b_1 R_w(t-1) + b_2 R_w(t-2) + b_3 R_A(t) \quad (4)$$

where  $R_w(t)$ ,  $R_w(t-1)$  and  $R_w(t-2)$  are the residuals of water temperature at times  $t$ ,  $(t-1)$  and  $(t-2)$  respectively,  $R_A(t)$  is the residual of air temperature at time  $t$ ,  $b_1$  and  $b_2$  are autoregressive coefficients, and  $b_3$  is a regression coefficient that reflects the heat exchange between water and air temperature. Caissie *et al.* (1998) used a Box–Jenkins model (SM3) that uses a transfer function linking air and water residuals:

$$R_w(t) = \frac{\zeta_0}{1 - \alpha_1 B} R_A(t) + \frac{a(t)}{1 - \phi_1 B} \quad (5)$$

where  $B$  is the backward shift operator defined by  $BR_w(t) = R_w(t-1)$ ,  $\zeta_0$ ,  $\alpha_1$ , and  $\phi_1$  are estimated parameters and  $a(t)$  is a white noise series of mean equal to zero and variance  $\sigma_a^2$ . Models SM1, SM2 and SM3, which were compared by Caissie *et al.* (1998) on a much smaller ( $50\text{ km}^2$ ) basin, are implemented to verify their suitability for a larger river.

### The effect of river discharge

To investigate the potential effect of streamflow, different flow-derived variables (i.e. streamflow indices) were also considered. The results can be used to demonstrate that the inclusion of flow variables as independent variables is statistically significant. In addition to three air temperature residuals (concomitant and lags 1 and 2 days) and water temperature residuals (lags 1 and 2 days), a variety of flow-derived variables (a total of 31 variables) were also considered as input variables in the analysis. Most studies concerned with the inclusion of streamflow consider only flow, i.e.  $Q(t)$ , as an independent variable to model  $T_w(t)$ . Considering lagged flows in the analysis may improve the results if there is a delayed response of water temperatures to flow due to travel time and thermal inertia of water. We also considered additional variables to see whether the model performance could be improved. For instance, change in flow explores the possibility that stream temperature is better correlated with departure from previous flow conditions than an absolute value of discharge in cubic metres per second.

Flow-derived variables that were investigated were limited to a period of 7 days before time  $t$ , in order to ensure that only short-term variations were included and that seasonal changes (e.g. main seasonal hydrological features such as the spring flood) were excluded in modelling water temperature residuals. Flow indices included: (1) flow for the past 7 days, (2) minimum and maximum flow of 3-, 5- and 7-day periods, (3) cumulative flow for the previous 2, 3, 5 and 7 days, (4) change in flow (i.e. difference between daily flows) for past 7 days, and (5) relative change in flow for the past 7 days. The relative flow change variables are defined as

$$r_{QC}(i) = \frac{Q(t) - Q(t-i)}{Q(t)} \quad i = 1, \dots, 7 \quad (6)$$

The most adequate regression model should explain most of the variance while being the most parsimonious. It should also be a physically reasonable model

that describes the known physical relationship between water temperatures and a set of input parameters. Different approaches are available for selecting the best set of explanatory variables and building the most adequate multiple regression model. The first method is the stepwise regression; the second is the best subsets regression, which uses an overall measure of goodness of fit to compare all possible subset models (Montgomery and Peck, 1992; Ryan, 1997). The different stepwise selection algorithms have been criticized on various grounds, the most common criticism being that none of the procedures guarantees that the best subsets regression model of any size will be identified, and that the order in which independent variables enter or leave the model does not necessarily imply an order of importance for the variables. An alternative to stepwise algorithms is to consider all possible subsets of different sizes. The amount of computation required to perform all-possible-subset regression increases as the number of variables and possible models increases. Different criteria are available to compare the different subsets and to select the most adequate parsimonious model: the adjusted  $R^2$ , the Akaike information criterion (AIC; Akaike, 1974) and the Schwartz Bayesian criterion (SBC; Schwartz, 1978). All these criteria take the number of independent variables (and hence the number of parameters) into account. The adjusted coefficient of determination  $R_a^2$  is defined as

$$R_a^2 = 1 - \left( \frac{n-1}{n-p} \right) \frac{SSE}{SST} \quad (7)$$

where  $n$  and  $p$  are respectively the number of observations and the number of regression parameters, SSE is sum of squared errors and SST is the total sum of squares. The Akaike criterion is defined as

$$AIC = n \ln \left( \frac{SSE}{n} \right) + 2p \quad (8)$$

and the Schwartz Bayesian criterion is

$$SBC = n \ln \left( \frac{SSE}{n} \right) + p \ln(n) \quad (9)$$

AIC and SBC are performance metrics that balance statistical fit and model parsimony. The best model is the model that corresponds to the minimal value of these two indices.

*Multicollinearity*

The use and interpretation of a multiple regression model often depends explicitly on the estimates of the individual regression coefficients. When the explanatory variables are highly correlated, the problem of multicollinearity is present, and the computed estimates of the regression coefficients are unstable and have large standard errors (Montgomery and Peck, 1992; Afifi and Clark, 1996). The regression coefficients fluctuate when used across samples, and even a slight change in the data can result in different regression coefficients. Unstable regression coefficients prevent the use of the regression

equation for different samples (Montgomery and Peck, 1992; Afifi and Clark, 1996). Multicollinearity can lead to results in which the magnitude of some of the regression coefficients may be grossly inflated and even of the wrong sign. Variance inflation factor (VIF) analysis is a multicollinearity diagnostic that can reveal complex relationships between independent variables (Montgomery and Peck, 1992). The variance inflation factor for each independent variable is defined as

$$VIF_j = \frac{1}{1 - R_j^2} \quad (10)$$

where  $R_j^2$  is the square of the multiple coefficient of determination from the regression of variable  $j$  on all other explanatory variables. It is suggested that  $VIF > 10$  is an indication that multicollinearity may be causing problems in estimation (Chatterjee and Price, 1991).

*Ridge regression*

Ridge regression, proposed originally by Hoerl and Kennard (1970a) is designed to produce better regression estimates for correlated variables. One important feature of ordinary least squares (OLS) regression is the requirement that the estimated regression coefficients be unbiased estimates of true coefficients. In practice, an estimator that has a small bias but is substantially more precise than an unbiased estimator may be the preferred estimator if it is more stable (Hoerl and Kennard, 1970a). Ridge regression modifies OLS to allow biased estimates of regression coefficients, and provides estimates that are more robust than least OLS estimates for small perturbations in data (Hoerl and Kennard, 1970a; Chatterjee and Price, 1991). The ridge estimates are more stable in the sense that they are not affected by slight variations in the estimated data, and also forecasts of the response variable corresponding to values of the explanatory variable not included in the estimation set tend to be more accurate (Chatterjee and Price, 1991). More often than not, ridge regression is performed on standardized or so-called correlation form. The coefficients are estimated from

$$(R_{xx} + kI)\mathbf{b}^* = R_{yx} \quad (11)$$

where  $R_{xx}$  is correlation matrix of  $x$ ,  $R_{yx}$  are the correlations of  $y$  with the  $x$ ,  $\mathbf{b}^*$  is the vector of standardized ridge regression coefficients,  $k$  is called the ridge constant, and  $I$  is the identity matrix. It has been shown that there is a positive value of  $k$  for which the ridge estimates will be stable with respect to small changes in the estimation data (Hoerl and Kennard, 1970b). The strategy for choosing  $k$  is to try several successive values of  $k$  and select the values for which the regression coefficients become stable and variation inflation factors become small. Guidelines for selection of the ridge constant have been outlined by Hoerl and Kennard (1970b):

1. At a certain value of ridge constant  $k$ , the system will stabilize and have general character of an orthogonal system.

2. Coefficients will have reasonable absolute values.
3. Coefficients with improper signs at  $k = 0$  will have changed to have proper signs.
4. The residual sum of squares will not have been inflated to an unreasonable value.

Once the ridge constant is chosen and standardized coefficients are obtained, the regression coefficients are estimated using

$$b = \mathbf{b}^* \frac{S_y}{S_x} \quad (12)$$

where  $S_y$  and  $S_x$  are the sample standard deviations of  $y$  and  $x$  respectively.

#### Model evaluation and validation

The fit between simulated or predicted  $P$  and observed  $O$  water temperature was evaluated using different statistical measures: model efficiency, often called the Nash–Sutcliffe coefficient (NSC; Equation (13)), the relative mean bias (RB; Equation (14)) and RMSE (Equation (15)). These criteria were calculated for each model, for both the calibration and validation periods. The Nash–Sutcliffe coefficient of efficiency has been widely used to evaluate the performance of hydrological models (e.g. Nash and Sutcliffe, 1970; Mohseni *et al.*, 1998). Nash and Sutcliffe (1970) defined NSC, which ranges from minus infinity (poor model) to unity (perfect model), as

$$\text{NSC} = 1 - \frac{\sum_{i=1}^n (P_i - O_i)^2}{\sum_{i=1}^n (\bar{O} - O_i)^2} \quad (13)$$

where  $n$  is the number of data points. The relative mean bias RB is a dimensionless measure that is computed simply as the sum of the difference between the observed and predicted values divided by the standard deviation of the measured values; it is an expression of the bias of the model (Janssen and Heuberger, 1995):

$$\text{RB} = \frac{\sum_{i=1}^n (P_i - O_i)}{n\sigma(O)} \quad (14)$$

where  $n$  is number of observations and  $\sigma(O)$  is the standard deviation of the observed values. The RMSE is also calculated to provide a joint assessment of bias and the variance:

$$\text{RMSE} = \left[ \frac{\sum_{i=1}^n (P_i - O_i)^2}{n - k} \right]^{0.5} \quad (15)$$

where  $k$  is the number of parameters used in the model.

Time-series were divided into two sub-series, the first part for calibration and the second for validation. The first

four years' data (1993–1996) was used for calibrating the model. The 1997–1998 time-series data were then used to test and validate the calibrated model.

## RESULTS

Water temperatures in the study catchment are typical of the temperature climate in northeastern Canada. Summary statistics derived from a 6-year open water period reveal that the mean annual water temperature for the open water period varies from 9.0 to 10.5 °C (Table I). Minimum and maximum mean annual air temperatures for this period are 8.5 °C and 10.5 °C respectively. The average annual water discharge varies between 502 and 686 m<sup>3</sup> s<sup>-1</sup>. The temperature regime of Moisie is characterized by rapidly increasing temperatures in May and June, with peaks in July–August and a subsequent decline in autumn. Hence, departure and end dates of the annual cycle are strongly modulated by river icing. The river-ice season (freeze-over to break-up) starts in mid November and lasts until early May. During this period the stream temperature is around 0 °C and the daily mean water temperatures cannot respond to air temperatures to the same degree as during the open water period.

This study focuses on the open water period and covers the period from 1 May to 11 November. Six years of data (1993 to 1998) were used in the analysis. Regression analysis of water and air temperatures time-series against time did not reveal any significant trends. Seasonal components for both water and air temperature were calculated by fitting a sinusoidal function to the interannual daily mean time-series. The sine function (Equation (2)) was first fitted to the water temperature time-series of the calibration period (1993–1996) by minimizing the sum of squares error (Figure 2). Values for  $a$ ,  $b$ , and  $t_0$  were thus estimated:

$$T_A(t) = 0.23 + 16.48 \sin \left[ \frac{2\pi}{365}(t - 119.79) \right] \quad (16)$$

The estimated  $t_0$ , which corresponds to the mean time of departure from a null temperature, was 120 days (i.e. May 1). The amplitude of the sine function, the maximum of the annual water temperature cycle, is equal

Table I. Summary statistics of water temperature, air temperatures and streamflow for the open water period

	Mean temperature (°C)		Mean discharge (m <sup>3</sup> s <sup>-1</sup> )
	Water	Air	
1993	9.65	8.51	551.69
1994	9.60	9.68	685.94
1995	10.54	10.15	501.53
1996	9.42	10.04	622.31
1997	9.02	9.15	595.28
1998	10.35	10.47	533.28
1993–1996	9.68	9.81	564.28
1997–1998	9.80	9.60	590.37

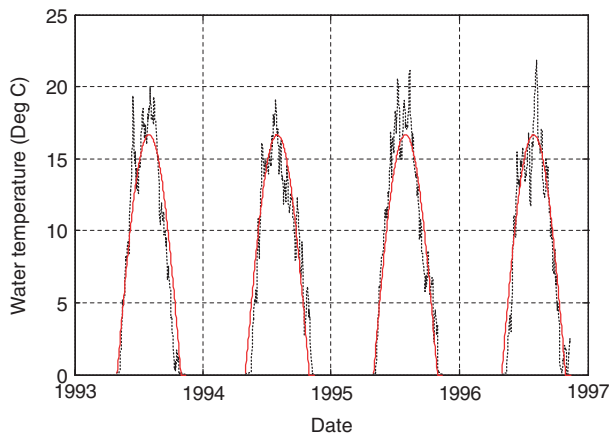


Figure 2. Daily mean water temperatures (1993–1996) and annual cycle estimated by a sinusoidal function. The dashed line represents daily mean water temperature and the solid line represents periodic sinusoidal function. This figure is available in colour online at <http://www.interscience.wiley.com/hyp>

Table II. Performance measures (RMSE, Nash coefficient) of annual component ( $T_A$ ) for every year

	$S_w$ (°C) <sup>a</sup>	$T_A$	
		RMSE (°C)	NSC
1993	6.70	2.21	0.892
1994	5.66	2.09	0.863
1995	6.14	1.87	0.907
1996	6.41	1.94	0.909
1997	5.76	1.66	0.911
1998	6.03	1.94	0.896
1993–1998	6.24	2.08	0.890

<sup>a</sup> Standard deviation of water temperature.

to 16.48 °C (Figure 2). The maximum water temperature occurs on 29 July, or day 210 (Figure 2). Water temperature exhibits a seasonal pattern of variation and the sinusoidal function fits the measured data well. The comparison on an annual basis shows that, for Moisie River, the Nash coefficient for seasonal component ( $T_A$ , Equation (16)) varies between 0.86 (in 1994) and 0.91 (in 1997; Table II). The highest variability and highest departure from long-term trend were observed in 1993. The RMSE and Nash coefficient were 2.21 °C and 0.89 respectively for this year. Equation (2) was also used to estimate the seasonal component of air temperatures for the calibration period. It should be noted that, in the case of air temperature, the sine curve was fitted for a 365-day

period and was not limited to zero and positive values. The following sine function was estimated for air temperature:

$$T_A(t) = 0.44 + 15.67 \sin \left[ \frac{2\pi}{365}(t - 116.23) \right] \quad (17)$$

The interannual mean of maximum air temperature was estimated to occur on 25 July. Comparison of Equations (16) and (17) indicates that the fitted sinusoidal function for water temperature lags the fitted sinusoidal function for the air temperature by 3.5 days (i.e.  $t_0$ ). The periodic function accounted for 86% of variance in the air temperature data.

The residuals representing short-term variations are then calculated by subtracting the seasonal components (i.e. Equations (16) and (17)) from the observed air and water temperature time-series. SM1 (Equation (3)), SM2 (Equation (4)) and SM3 (Equation (5)) were implemented to model the residuals. The estimation of the parameters for the calibration phase yielded the following equation for SM1:

$$R_w(t) = 0.201R_A(t) + 0.125R_A(t - 1) + 0.158R_A(t - 2) \quad (18)$$

Table III provides values of RMSE, NSC and RB. The NSC value for SM1 was 0.301, and the RMSE was 1.733 °C for the calibration period. Relative mean bias was 0.027. Statistical criteria values were similar for the validation period, with the exception of an increase in bias to 0.146 (Table III).

The above model does not take into account the significant autocorrelation of water temperature residuals. Hence, as suggested by Cluis (1972), another stochastic model (SM2, Equation (4)) that includes lagged water temperature residuals was used. SM2 was also adjusted to the data from the same calibration period. The estimation of parameters yielded

$$R_w(t) = 1.295R_w(t - 1) - 0.394R_w(t - 2) + 0.608R_A(t) \quad (19)$$

The NSC coefficient was 0.936, the RMSE was 0.527 °C and RB was 0.002 for the calibration period (Table III). The results show a significant reduction of RMSE and RB compared with SM1. This indicates the importance of including autoregressive components in the model.

The last model considered (SM3, Equation (5)) was a special class of Box–Jenkins models. An iterative method that minimizes the prediction error was used to estimate

Table III. Performance measures (NSC, RMSE and RB) of SM1, SM2 and SM3 for calibration period (1993–1996) and validation period (1997–1998)

Model	NSC		RMSE (°C)		RB	
	Calibration	Validation	Calibration	Validation	Calibration	Validation
SM1	0.301	0.204	1.733	1.715	0.027	0.146
SM2	0.936	0.920	0.527	0.543	0.002	0.015
SM3	0.925	0.906	0.568	0.588	0.0004	0.005

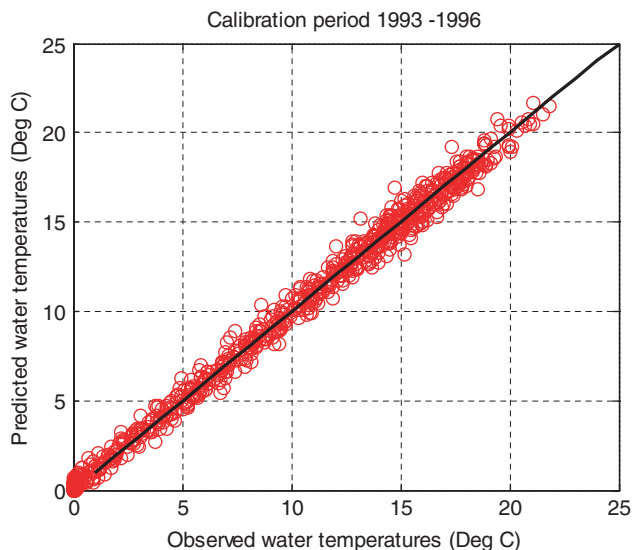


Figure 3. Scatter plot of observed and estimated water temperatures obtained by SM2 for the calibration period (1993–1996). Solid line represents 1 : 1 line. This figure is available in colour online at <http://www.interscience.wiley.com/hyp>

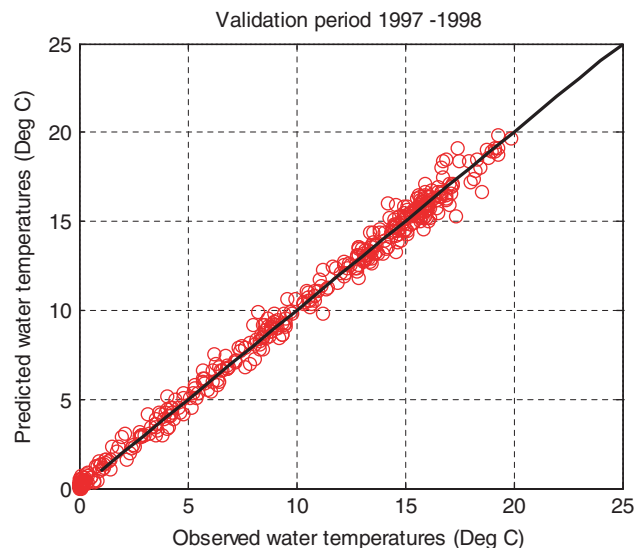


Figure 4. Scatter plot of observed and estimated water temperatures obtained by SM2 for the validation period (1997–1998). Solid line represents 1 : 1 line. This figure is available in colour online at <http://www.interscience.wiley.com/hyp>

the parameters (Matlab's System Identification Toolbox; Ljung, 1997):

$$R_w(t) = \frac{0.114}{1 - 0.702B} R_A(t) + \frac{a(t)}{1 - 0.957B} \quad (20)$$

The NSC for SM3 for the calibration period was 0.925 and the RMSE was 0.568 °C (Table III). Small relative mean bias (RB) values were observed for SM2 and SM3, and these two models slightly overestimate the water temperature residuals (Table III). NSC values were higher for the second-order autoregressive model (SM2) and Box–Jenkins model (SM3), and these two models outperformed the multiple regression (SM1). The second-order autoregressive model (SM2) is preferred over Box–Jenkins (SM3) because it has the lowest RMSE (0.527) and the highest NSC coefficient (0.936). In order to calculate the predicted water temperatures, Equation (1) was used and the seasonal component ( $T_A$ , Equation (16)) was added to predicted water temperature residuals  $R_w$  for the selected best model SM2 (Equation (19)). Figure 3 shows the scatter plot of predicted values by the model SM2 versus observed water temperatures for the calibration period. SM2 provides good estimates of measured values for the whole range of water temperatures. A good fit is also observed for the validation period (Figure 4). Time-series of observed and predicted water temperatures and model residuals obtained by SM2 for each year are presented in Figure 5. A good agreement between the measured and predicted values is observed and the difference between measured and predicted values (model residuals) of SM2 is very small (Figure 5). The maximum and minimum errors, which correspond to maximum underestimation and overestimation error for the calibration period (1993–1996), were 1.95 °C and 2.20 °C respectively. The maximum and minimum errors for the validation period (1997–1998) were 2.03 °C and 1.81 °C respectively.

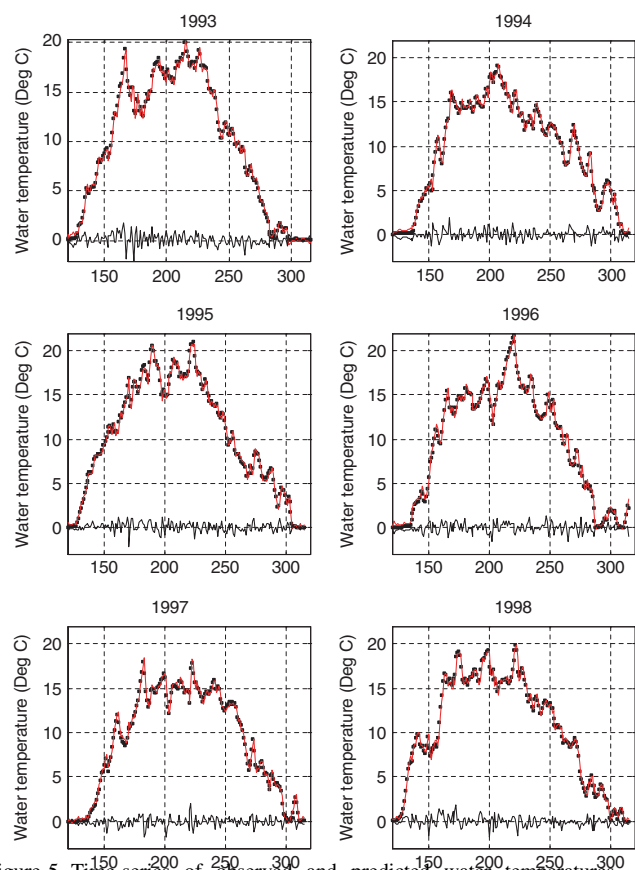


Figure 5. Time-series of observed and predicted water temperatures obtained by SM2 for every year, the squares represent the observations, solid lines represent predicted values of SM2 and model residuals. This figure is available in colour online at <http://www.interscience.wiley.com/hyp>

#### Regression methods including river discharge

A total of 36 potential explanatory variables, including autoregressive terms, lagged air temperature residuals



and 31 flow-derived variables, were considered in the analysis. Candidate models were derived using best-subsets regression with water temperature residuals as the response variable and 36 potential explanatory variables. Models including all possible sets of 1, 2, 3, . . . ,  $n$  predictors were considered to find the best-fitting models for a specified number of variables. Three different criteria, namely the adjusted  $R^2$ , the AIC and the SBC, are used to compare different models and select the most adequate parsimonious model. Models including all possible sets of four to nine variables were estimated. These models are denoted SM4, SM5, SM6, SM7, SM8 and SM9. Results for the calibration period are presented in Table IV. These results indicate that all models had similar  $R_a^2$  (between 0.939 and 0.941) and similar RMSE values between 0.502 and 0.511 °C. The lowest AIC (−1066.3) was found for the subset of seven variables (SM7) and the lowest SBC (−1030.6) value was found for the subset of five variables (SM5). The model-building principle of parsimony states that the smallest possible number of parameters should be used so as to give an adequate representation of the given data (Chatfield, 2000). The SBC will generally result in the most parsimonious model, as it utilizes a larger penalty function than AIC, suggesting a model with fewer parameters (Miller, 1990).

The performance of these models was also compared for the validation period (1997–1998) using RMSE, relative bias and the Nash coefficient described in the previous section (Table V). All models perform very well, with RMSE values of ~0.53 °C, and NSC values above 0.92. Small relative mean bias values are observed for

different models, and all these models slightly over-estimate the water temperature residuals. Although the models performed comparably for the validation period and the differences are minimal, model SM5, which has a minimum SBC for the calibration period, also has the lowest RMSE (0.521 °C) and relative bias (0.016) for the validation period. Therefore, model SM5 is selected as the most adequate model, and is defined as follows:

$$R_w(t) = 1.29R_w(t - 1) - 0.389R_w(t - 2) + 0.082R_A(t) - 0.026R_A(t - 2) - 1.534r_{QC}(1) \quad (21)$$

Results of the analysis show that relative flow change is the most important flow variable and that there is a negative relationship between water temperature and relative flow change. Improvements in performance can be obtained by including streamflow variables in the model, but the amount of explained variance associated with the addition of flow variables is modest. SM5 has an RMSE of 0.507 °C, which shows an improvement compared with SM2 (RMSE = 0.527 °C), which did not include any streamflow indices.

Once the water temperature residuals were estimated, the seasonal component ( $T_A$ , Equation (16)) was added to predicted water temperature residuals of model SM5 (Equation (21)) to calculate predicted water temperatures. Figure 6 shows the scatter plot of predicted values by the model SM5 versus observed water temperatures for the calibration period. The measured/predicted pairs are lying close to the 1 : 1 line and good estimates of measured values over the range of water temperatures are

Table IV. Performance measures (RMSE, adjusted  $R^2$ , AIC and SBC) of best subsets regression models for the calibration period

Model	Variables	RMSE (°C)	$R_a^2$	AIC	SBC
SM4, best subset of four variable	$R_w(t - 1), R_w(t - 2), R_A(t), r_{QC}(1)$	0.511	0.939	−1044.8	−1026.1
SM5, best subset of five variables	$R_w(t - 1), R_w(t - 2), R_A(t), R_A(t - 2), r_{QC}(1)$	0.507	0.940	−1053.9	−1030.6
SM6, best subset of six variables	$R_w(t - 1), R_w(t - 2), R_A(t), r_{QC}(1), r_{QC}(3), r_{QC}(6)$	0.506	0.941	−1057.6	−1029.7
SM7, best subset of seven variables	$R_w(t - 1), R_w(t - 2), R_A(t), R_A(t - 2), r_{QC}(1), r_{QC}(3), r_{QC}(6)$	0.504	0.941	−1066.3	−1029.9
SM8, best subset of eight variables	$R_w(t - 1), R_w(t - 2), R_A(t), R_A(t - 1), Res_A(t - 2), r_{QC}(1), r_{QC}(3), r_{QC}(6)$	0.502	0.941	−1066.3	−1029.0
SM9, best subset of nine variables	$R_w(t - 1), R_w(t - 2), R_A(t), Res_A(t - 2), r_{QC}(1), r_{QC}(3), r_{QC}(6), Q(t - 3), Q(t - 5)$	0.502	0.941	−1064.8	−1022.9

Table V. Performance measures (RMSE, RB and NSC) of best subsets regressions for calibration period (1993–1996) and validation period (1997–1998)

Models	No. of variables	Calibration			Validation		
		RMSE (°C)	RB	NSC	RMSE (°C)	RB	NSC
SM4	4	0.511	0.0022	0.940	0.530	0.0184	0.924
SM5	5	0.507	0.0027	0.941	0.521	0.0155	0.926
SM6	6	0.506	0.0066	0.941	0.533	0.0236	0.923
SM7	7	0.503	0.006	0.942	0.524	0.0195	0.926
SM8	8	0.502	0.006	0.942	0.522	0.0189	0.926
SM9	9	0.502	0.0055	0.942	0.523	0.0189	0.926

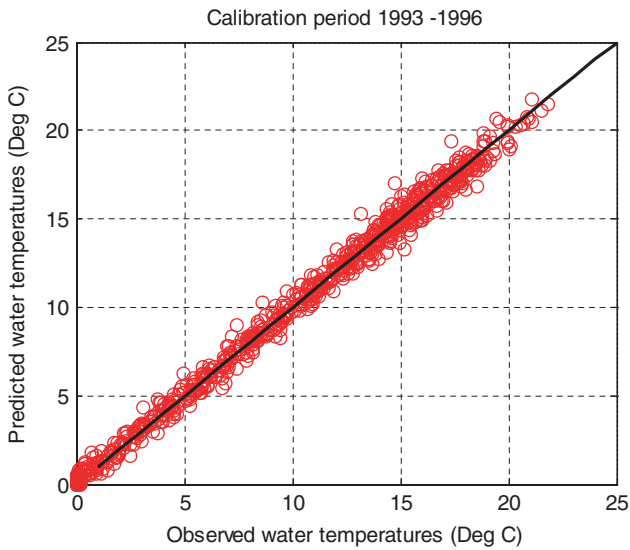


Figure 6. Scatter plot of observed and estimated water temperatures obtained by SM5 for the calibration period (1993–1996). Solid line represents 1 : 1 line. This figure is available in colour online at <http://www.interscience.wiley.com/hyp>

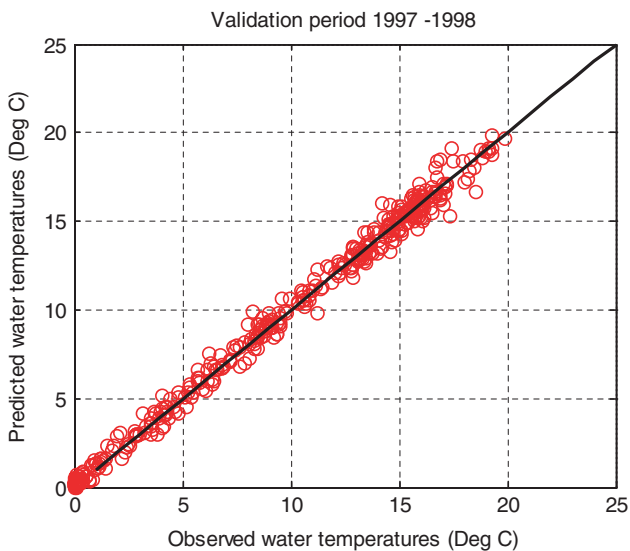


Figure 7. Scatter plot of observed and estimated water temperatures obtained by SM5 for the validation period (1997–1998). Solid line represents 1 : 1 line. This figure is available in colour online at <http://www.interscience.wiley.com/hyp>

observed. A good fit is also observed for the validation period (Figure 7). Time-series of observed and measured water temperatures and model residuals for each year are presented in Figure 8. There is a good agreement between predicted and measured values, and the differences between measured and predicted values are very small (Figure 8).

The ridge regression was applied to the best model developed by subset regression (SM5, Equation (21)) and regression coefficients  $br_1, br_2, \dots, br_5$  were calculated:

$$R_w(t) = br_1 R_w(t - 1) + br_2 R_w(t - 2) + br_3 R_A(t) + br_4 R_A(t - 2) + br_5 (r_{QC1}) \quad (22)$$

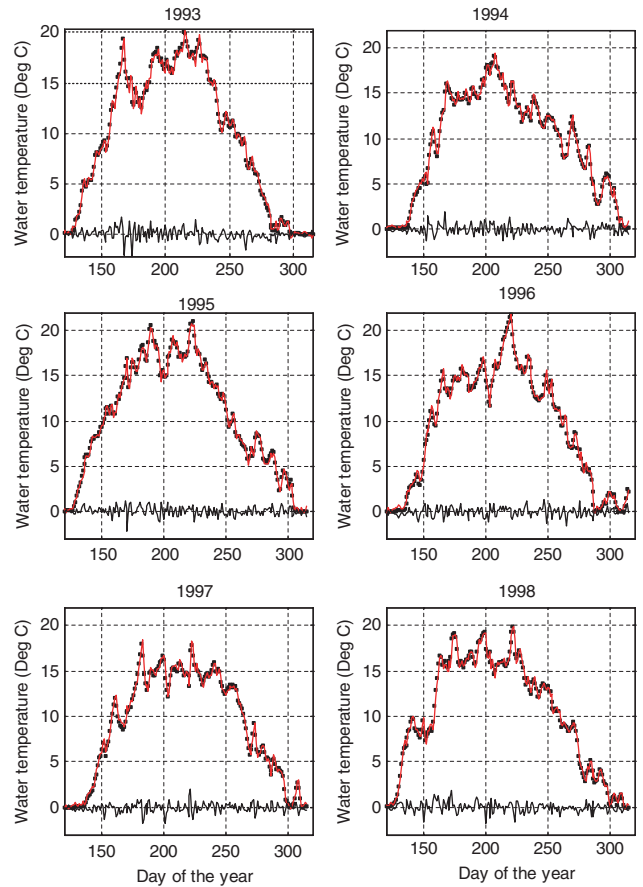


Figure 8. Time-series of observed and predicted water temperatures obtained by SM5 for every year; the squares represent the observations, solid lines represent predicted values of SM5 and model residuals. This figure is available in colour online at <http://www.interscience.wiley.com/hyp>

Figure 9 presents the ridge trace for regression coefficients of SM5. The ridge trace is a simultaneous graph of the  $p$  regression coefficients plotted as a function  $k$ , the ridge constant. The ridge constant was chosen with regard to ridge trace and following Hoerl and Kennard (1970b) guidelines and also by calculating the ridge constant for which the VIFs are close to one. The regression coefficients appear to have stabilized for values greater than  $k = 0.12$ . A more rigorous calculation was followed and a ridge constant was chosen so that the individual VIFs for independent variables in the model become close to unity. Table VI presents the VIFs of regression coefficients calculated for different ridge constants. As can be seen from Table VI, the VIFs are around unity for  $k = 0.13$ . This value was selected as the optimum ridge constant and standardized ridge regression coefficients were calculated. Equation (12) was then used and the corresponding ridge regression coefficients were then calculated:

$$R_w(t) = 0.628R_w(t - 1) + 0.186R_w(t - 2) + 0.12R_A(t) - 0.002R_A(t - 2) - 2.278(r_{QC1}) \quad (23)$$

Comparisons of ordinary least-squares (OLS) regression coefficients (Equation (21)) and ridge regression

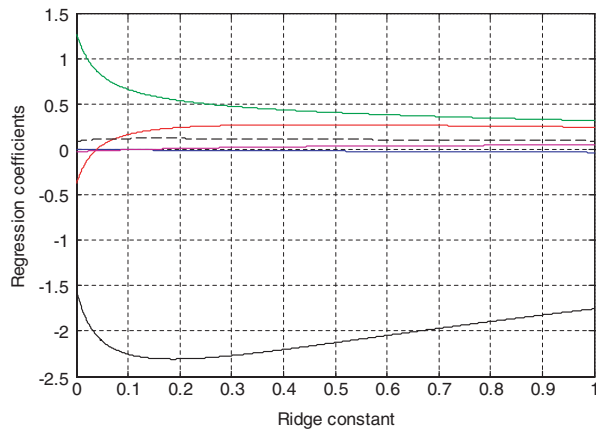


Figure 9. Ridge trace for  $0 < k < 1$ . This figure is available in colour online at <http://www.interscience.wiley.com/hyp>

Table VI. VIFs for regression coefficients of model SM5 for different ridge constants

Ridge trace $k$	VIF				
	$br_1$	$br_2$	$br_3$	$br_4$	$br_5$
0	12.78	11.78	1.41	1.46	1.17
0.01	8.42	7.80	1.32	1.41	1.13
0.02	6.00	5.59	1.26	1.36	1.10
0.03	4.52	4.24	1.21	1.32	1.07
0.04	3.55	3.34	1.17	1.28	1.04
0.05	2.87	2.73	1.13	1.24	1.01
0.06	2.39	2.28	1.10	1.20	0.99
0.07	2.02	1.94	1.07	1.17	0.96
0.08	1.75	1.69	1.04	1.13	0.94
0.09	1.53	1.48	1.01	1.10	0.92
0.1	1.35	1.32	0.98	1.07	0.90
0.11	1.21	1.19	0.96	1.04	0.88
0.12	1.09	1.08	0.93	1.01	0.86
0.13	0.99	0.99	0.91	0.98	0.84
0.14	0.91	0.91	0.89	0.96	0.82
0.15	0.84	0.84	0.87	0.93	0.81

Table VII. Performance measures for ridge and OLS regression models for the calibration period (1993–1996) and validation period (1997–1998)

Model	Calibration			Validation		
	RMSE	RB	NSC	RMSE	RB	NSC
Ridge	0.644	0.036	0.903	0.677	0.030	0.876
OLS	0.507	0.0027	0.941	0.521	0.0155	0.926

coefficients (Equation (23)) indicate that the coefficient of  $R_w(t - 2)$  is positive in ridge regression (0.186), whereas it has a negative sign in OLS regression ( $-0.389$ ). The correlation between  $R_w(t)$  and  $R_w(t - 2)$  is positive (0.86). Therefore, a positive regression coefficient is expected for this term. The ridge regression coefficient is more in agreement with prior expectation and better represents the positive relationship between

$R_w(t - 2)$  and  $R_w(t)$ . The negative sign in OLS is due to multicollinearity and high correlation ( $R = 0.95$ ) between  $R_w(t - 1)$  and  $R_w(t - 2)$ . A comparison of performance measures for ridge and OLS regression is presented in Table VII. As expected RMSE and relative bias are higher for ridge regression; however, the increase in RMSE is minimal ( $<0.15$  °C) and the relative biases for ridge regression are small for both the calibration (0.033) and the validation period (0.03). If the independent variables are orthogonal, then the regression coefficients can be used as sensitivity indices to indicate the sensitivity of the output to model inputs. The ridge model is more appropriate for this purpose. The coefficient of relative flow change in Equation (23) is  $-2.28$ . This indicates that a 100% increase of flow between day  $i$  and day  $(i - 1)$ , i.e.  $r_{QC}(1) = 0.5$  from Equation (6), would result in a decrease of water temperature by  $1.14$  °C.

### DISCUSSION AND CONCLUSION

Stochastic approaches were used to relate water temperatures to air temperatures and streamflow in the Moisie River. The analysis consisted of separating the water temperatures into two different components: a long-term seasonal component and a short-term component (residuals). The seasonal components of stream water temperatures and air temperatures were established by fitting a periodic sinusoidal function to these time-series. Residuals, or departures, from the long-term trend (seasonal variations) were calculated for both water and air temperature time-series by subtracting the seasonal components from the corresponding time-series.

Different stochastic models were used to model the water temperature residuals. The first three stochastic models used either lagged air temperature residuals or a combination of lagged air and water temperature residuals as input variables. Three different models were used: a multiple regression, a second-order autoregressive model, and a classic Box–Jenkins model. The RMSE for these models varied between  $0.53$  and  $1.70$  °C, and the second-order autoregressive model provided the best results. The results indicate that two of these models perform very well in prediction of daily stream temperatures of the Moisie River. These methods have been successfully applied on much smaller rivers, and there have been very few implementations of such models on rivers as large as the Moisie. Caissie *et al.* (1998) obtained an RMSE of  $\sim 1.28$  °C for these models when applied to a much smaller ( $50 \text{ km}^2$ ) basin. The present study indicates that stochastic models originally developed on smaller systems are adequate for larger rivers as well. The multiple regression model (SM1) had the worst performance among the models (RMSE =  $1.73$  °C). Our first impression was that this relatively poor performance may have been caused by the fact that the air temperature data used in our study were not measured in the vicinity of the water temperature station, but originated from a station located approximately 200 km further south.

However, the correlation between water and air temperature time-series was 0.89, which seems high and in range of corresponding reported values in the literature (Stefan and Preud'homme, 1993). Other studies have shown that the strength of the air–stream temperature correlation was not affected significantly for distances of the order of 200 km (Pilgrim *et al.*, 1998; Erickson and Stefan, 2000). Of course, these correlations will depend on a number of factors, including topography and land use. To investigate this further, we also performed a simple linear regression between air and water temperatures. This resulted in a Nash coefficient of 0.80 and an RMSE of 2.77 °C, which are in agreement with typical values in the literature (Erickson and Stefan, 2000; Morrill *et al.*, 2005). Hence, it appears that, in spite of a relatively important distance between air and water temperature measurements, the statistical relationship between both variables is sufficiently strong to include air temperature data in the model.

The combined influence of air temperature and streamflow on daily data was also investigated by developing stochastic models that included flow-derived variables (i.e. streamflow indices). Three different criteria ( $R_a^2$ , AIC, SBC) were used to compare models with different numbers of independent variables and to select the best parsimonious model. The results indicated that considering streamflow indices in stochastic models improves the overall performance; however, the improvement is modest and the rise in Nash–Sutcliffe coefficient gained by adding flow variables is small. The best model that included streamflow indices (SM5) has an RMSE of 0.507 °C, whereas the best model including only air temperatures (SM2) had an RMSE of 0.527 °C. The performance of these two models was also compared on an annual basis (Table VIII). Model SM5 had the minimum annual RMSE and maximum annual NSC and outperformed SM2 for every year; however, the associated reduction in RMSE was modest and varied between 0.015 and 0.03 °C. Both models have high Nash coefficients (>0.9) for all years. These results are consistent with findings from a previous study that reported modest improvements by addition of flow (Webb *et al.*, 2003). Webb *et al.* (2003) developed regression models to predict water temperatures from air temperature and flow for four catchments in southwest England ( $R^2$  values between

0.857 and 0.919) and reported that the improvement by adding flow in the regression equation was modest.

The relative flow change was the most important streamflow variable in the best model including streamflow indices (SM5, Equation (21)) and was inversely related to stream temperature. Using the best subsets regression guarantees that the best models for each group having an equal number of variables are obtained. As flow and lagged flows were compared within the same groups, it can be concluded that the partial correlation between water temperature residuals  $R_w$  and relative flow change was higher than that of  $R_w$  and flow itself.

Relative flow change is probably a better indicator of change in depth, which modifies the local heat budget. The impact of discharge on water temperature occurs primarily through the associated increase/decrease in depth and the increased/decreased thermal inertia of the river. Thus, a change in discharge indicates associated changes in depth and thermal inertia of the system. For instance, a positive value of relative discharge (i.e. if flow increases from the day before) indicates an increased thermal inertia of the water body, which results in decreased heat exchange between the system and the environment. Higher discharges also result in reduced travel time of flow through the channel system, and, in turn, decreased exposure of stream to solar radiation. For instance, during the summer, precipitation from storms may decrease stream temperature if the volume of discharge increases at a faster rate (i.e. high relative flow change) than it could be heated in a stream channel (Smith and Lavis, 1975).

The only other water temperature model previously used to analyse Moisie River water temperatures is CEQUEAU, a conceptual deterministic model. The water temperature model in CEQUEAU is based on a general heat budget methodology (Morin and Couillard, 1990). Outputs from the hydrological model are used as inputs in the temperature model to calculate local temperature. Detailed equations and methods for calculating the terms of the surface heat budget can be found in Morin and Couillard (1990). Morin and Sochanski (1990) used CEQUEAU to estimate the water temperatures of the Moisie River at another station, for the period of 1973 to 1989 (16 years), and they obtained an RMSE of about 1 °C. The prediction accuracy of the stochastic models for water temperatures presented here for 1993–1998 (6 years) is greater than that of the deterministic approach. However, the comparison between model performances is difficult to make because, although both models were applied on the same river, different stations and time periods were used.

When applying regression models to water temperatures, care should be taken to avoid collinearity between independent variables, as this results in unstable regression coefficients with large standard errors. To alleviate this problem, application of ridge regression was proposed and the ridge regression was successfully applied to the best model obtained in previous sections, which was the model including relative flow change (SM5). This

Table VIII. Performance measures (RMSE, RB, NSC) of water temperature from models SM2 and SM5 for every year

Year	RMSE (°C)		RB		NSC	
	SM2	SM5	SM2	SM5	SM2	SM5
1993	0.548	0.528	−0.017	−0.021	0.938	0.943
1994	0.530	0.501	0.0150	0.014	0.936	0.943
1995	0.511	0.489	−0.018	−0.018	0.926	0.933
1996	0.515	0.501	0.033	0.033	0.930	0.933
1997	0.529	0.499	0.033	0.034	0.899	0.910
1998	0.558	0.543	0.002	0.001	0.918	0.922

resulted in an increase in RMSE of about 0.15 °C, which is largely compensated by the stability of the regression coefficients. The empirical models presented in this article present a simple way of predicting water temperatures by linking water temperatures to both air temperatures and flow variables. The results from the present study suggest that stochastic models are capable of predicting daily mean water temperatures with an RMSE of 0.51 °C and offer a simple means of successfully predicting water temperature time-series. The stochastic models provided comparable or better results than those reported in the literature for both stochastic models (Kothandaraman, 1971; Cluis, 1972; Caissie *et al.*, 1998) and deterministic models (Morin and Sochanski, 1990).

River water temperature is a relatively inexpensive variable to monitor. As long time-series of water temperature and flow become increasingly available, models such as the ones developed in this study will be of greater use. In practice, these models could be used successfully in many applications, including (but not limited to): predicting water temperatures and calculating many biotic indices, such as growth rate of aquatic organisms (Crisp and Howson, 1982); the effect of climate change on aquatic habitat (Mohseni *et al.*, 2003); the implication of climate change in water quality (Morrill *et al.*, 2005); filling the gaps and estimating missing water temperatures in biological studies (Swansburg *et al.*, 2002).

The wider implications of this study are: (1) The family of time-series or stochastic models originally developed on smaller systems are adequate for larger rivers. (2) Our preliminary study shows that considering stream-flow indices improves the performance of stochastic models. The inclusion of flow as an independent variable opens the way to applications of such models on regulated rivers. (3) Application of ridge regression was also proposed for collinear data, which may be used as a useful tool. Given that in many water temperature and hydrological studies the modeller is faced with the challenge of modelling multicollinear data, ridge regression may be used as an additional predictive tool.

The results obtained from the models presented in this study are only directly applicable to the specific river. However, the methodology presented in this paper can be applied in other rivers to predict water temperatures by considering different streamflow indices along with lagged water and air temperature residuals and performing the best subsets regression. It may, therefore, become an interesting alternative or complement to more complex deterministic approaches for which the required data are not always available.

In this study, linear stochastic approaches were used to model daily water temperature. There are still many methods and tools that have potential applications in water temperature modelling, among which are nonlinear methods such as artificial neural networks and nonparametric methods like  $K$  nearest neighbours. These methods have been used successfully in many forecasting applications in engineering and science; however, they have

not yet been fully explored in water temperature modelling. Further research is needed to explore and refine these methods and to find suitable applications in water temperature modelling. Further research is also needed to develop stochastic models at a sub-daily scale (i.e. hourly). Stochastic modelling of hourly stream temperatures is more challenging because it requires modelling of the diurnal cycle, which may vary with the season.

#### ACKNOWLEDGEMENTS

We are indebted to two anonymous reviewers and the editor for their constructive comments on this paper and to Zeljka Ristic Rudolf for her assistance with some figures. This work was funded by the Natural Sciences and Engineering Research Council (Canada), Hydro-Québec, and the Institut national de la recherche scientifique (University of Québec, Canada).

#### REFERENCES

- Affi AA, Clark V. 1996. *Computer-Aided Multivariate Analysis*. Chapman and Hall: New York.
- Akaike H. 1974. A new look at the statistical model identification. *IEEE Transactions on Automatic Control* **19**: 716–723.
- Bartholow JM. 1989. *Stream temperature investigations: field and analytic methods*. Instream Flow Information Paper No. 13, US Fish and Wildlife Service.
- Bartholow JM. 1999. *SSTEMP for Windows: the stream segment temperature model* (Version 1-1-3). US Geological Survey computer model and help file.
- Bovee KD. 1982. *A guide to stream habitat analysis using the instream flow incremental methodology*. National Biological Service, Fort Collins, CO.
- Brown GW, Krygier JT. 1970. Effects of clear-cutting on stream temperature. *Water Resources Research* **6**(4): 1133–1139.
- Caissie D, El-Jabi N, St-Hilaire A. 1998. Stochastic modelling of water temperatures in a small stream using air to water relations. *Canadian Journal of Civil Engineering* **25**: 250–260.
- Caissie D, El-Jabi N, Satish M. 2001. Modeling of maximum daily water temperatures in a small stream using air temperatures. *Journal of Hydrology* **251**: 14–28.
- Chatfield C. 2000. *Time Series Forecasting*. Chapman and Hall: London.
- Chatterjee S, Price B. 1991. *Regression Analysis by Example*, 2nd edn. Wiley: New York.
- Cluis D. 1972. Relationship between stream water temperature and ambient air temperature—a simple autoregressive model for mean daily stream water temperature fluctuations. *Nordic Hydrology* **3**(2): 65–71.
- Crisp DT, Howson G. 1982. Effect of air temperature upon mean water temperature in streams in the north Pennines and English Lake District. *Freshwater Biology* **12**: 359–367.
- Edwards RW, Densem JW, Russell PA. 1979. An assessment of the importance of temperature as a factor controlling the growth rate of brown trout in streams. *Journal of Animal Ecology* **48**: 501–507.
- Elliott JM, Hurley MA, Fryer RJ. 1995. A new, improved growth model for brown trout, *Salmo trutta*. *Journal of Functional Ecology* **9**: 290–298.
- Erickson TR, Stefan HG. 2000. Linear air/water temperature correlations for streams during open water periods. *ASCE, Journal of Hydrologic Engineering* **5**(3): 317–321.
- Gu R, Li Y. 2002. River temperature sensitivity to hydraulic and meteorological parameters. *Journal of Environmental Management* **66**: 43–56.
- Gu R, Montgomery S, Austin T. 1998. Quantifying the effects of stream discharge of summer river temperature. *Hydrological Science Journal* **43**(6): 885–894.
- Hodgson S, Quinn TP. 2002. The timing of adult sockeye salmon migration into fresh water: adaptations by populations to prevailing thermal regimes. *Canadian Journal of Zoology* **80**: 542–555.

- Hockey JB, Owens IF, Tapper NJ. 1982. Empirical and theoretical models to isolate the effect of discharge on summer water temperatures in the Hurunui River. *Journal of Hydrology (New Zealand)* **21**(1): 1–12.
- Hoerl A, Kennard R. 1970a. Ridge regression: biased estimation for nonorthogonal problems. *Technometrics* **12**(1): 55–67.
- Hoerl A, Kennard R. 1970b. Ridge regression: applications to nonorthogonal problems. *Technometrics* **12**(1): 69–82.
- Janssen PHM, Heuberger PSC. 1995. Calibration of process-oriented models. *Ecological Modelling* **83**: 55–66.
- Johnson SL, Jones JA. 2000. Stream temperature response to forest harvest and debris flows in western Cascades, Oregon. *Canadian Journal of Fisheries and Aquatic Science* **57**(suppl. 2): 30–39.
- Kothandaraman V. 1971. Analysis of water temperature variations in large river. ASCE, *Journal of the Sanitary Engineering Division* **97**(SA1): 19–31.
- Ljung L. 1997. *System Identification Toolbox User's Guide*. Mathworks Inc.: Natick, MA.
- Miller AJ. 1990. *Subset Selection in Regression*. Chapman and Hall: London.
- Mitchell S. 1999. A simple model for estimating mean monthly stream temperatures after riparian canopy removal. *Journal of Environmental Management* **24**(1): 77–83.
- Mohseni O, Stefan HG, Erickson TR. 1998. A nonlinear regression model for weekly stream temperatures. *Water Resources Research* **34**(10): 2685–2692.
- Mohseni O, Erickson TR, Stefan HG. 1999. Sensitivity of stream temperatures in the United States to air temperatures projected under a global warming scenario. *Water Resources Research* **35**(12): 3723–3733.
- Mohseni O, Stefan HG, Eaton JG. 2003. Global warming and potential changes in fish habitat in US streams. *Climate Change* **59**: 389–409.
- Montgomery DC, Peck EA. 1992. *Introduction to Linear Regression Analysis*. 2nd edn. Wiley: New York.
- Morin G, Couillard D. 1990. Predicting river temperatures with a hydrological model. In *Encyclopedia of Fluid Mechanics: Surface and Groundwater Flow Phenomena*. Volk Gulf Publishing Company: Houston, TX; 171–209.
- Morin G, Sochanski W. 1990. *Rivière Moisie—étude de milieu physique—volume 5: régimes thermiques de la rivière Moisie avant et après détournement de la rivière aux Pékans*. Rapport scientifique No. 296.
- Morrill JC, Bales RC, Conklin MH. 2005. Estimating stream temperature from air temperature: implications for future water quality. ASCE, *Journal of Environmental Engineering* **131**(1): 139–146.
- Nash JE, Sutcliffe JV. 1970. River flow forecasting through conceptual models. Part A. Discussion of principles. *Journal of Hydrology* **10**: 282–290.
- Neumann DW, Zagona EA, Rajagopalan B. 2003. A regression model for daily maximum stream temperature. *Journal of Environmental Engineering* **129**(7): 667–674.
- Petts GE. 2000. A perspective on the abiotic processes sustaining the ecological integrity of running waters. *Hydrobiologia* **422**: 15–27.
- Pilgrim JM, Fang X, Stefan HG. 1998. Stream temperature correlations with air temperatures in Minnesota: implications for climate warming. *Journal of the American Water Resources Association* **34**(5): 1109–1121.
- Ryan TP. 1997. *Modern Regression Methods*. Wiley: New York.
- Schwartz G. 1978. Estimating the dimension of a model. *The Annals of Statistics* **6**: 461–464.
- Sinokrot BA, Stefan HG. 1993. Stream temperature dynamics: measurements and modeling. *Water Resources Research* **29**(7): 2299–2312.
- Smith K, Lavis ME. 1975. Environmental influences on the temperature of a small upland stream. *Oikos* **26**: 228–236.
- Stefan HG, Preud'homme EB. 1993. Stream temperature estimation from air temperature. *Water Resources Bulletin* **29**(1): 27–45.
- St-Hilaire A, Morin G, El-Jabi N, Caissie D. 2000. Water temperature modelling in a small forested stream: implication of forest canopy and soil temperature. *Canadian Journal of Civil Engineering* **27**: 1095–1108.
- St-Hilaire A, Morin G, El-Jabi N, Caissie D. 2003. Sensitivity analysis of a deterministic water temperature model to forest canopy and soil temperature in Catamaran Brook (New Brunswick, Canada). *Hydrological Processes* **17**: 2033–2047.
- Swansburg E, Chaput G, Moore D, Caissie D, El-Jabi N. 2002. Size variability of juvenile Atlantic salmon: links to environmental conditions. *Journal of Fish Biology* **61**: 661–683.
- Webb BW, Nobilis F. 1997. Long-term perspectives on the nature of the air–water temperature relationship: a case study. *Hydrological Processes* **11**: 137–147.
- Webb BW, Walling DE. 1993. Temporal variability in the impact of river regulation on thermal regime and some biological implications. *Freshwater Biology* **29**: 167–182.
- Webb BW, Clack PD, Walling DE. 2003. Water–air temperature relationships in a Devon river system and the role of flow. *Hydrological Processes* **17**: 3069–3084.