1	
2	
3	
4	
5	Nonlinear post-processing of numerical seasonal climate
5	i toninical post processing of numerical seasonal chinate
6	forecasts
7	
8	
9	
10	
11	
12	
13	Joel Finnis*, Dept. of Geography, Memorial University of Newfoundland
14	William W. Hsieh, Dept. of Earth & Ocean Sciences, University of British Columbia
15	Hai Lin, Environment Canada
16	William Merryfield, Canadian Centre for Climate Modeling and Analysis, Environment
17	Canada
18	
19	
20	Submitted to Atmosphere-Ocean (20 May 2011)
21	
22	
23	
24	
25	Corresponding Author: jfinnis@mun.ca; Phone: 709-864-8987
26	Dept. of Geography, Memorial University of Newfoundland, St. John's NL, A1B 3X9
27	

#### 3 Abstract

4 Although numerical models are increasingly being used to generate operational seasonal 5 forecasts, the reliability of these products remains relatively low. Regression-based post-6 processing methods have proven useful in increasing forecast skill, but as yet efforts have 7 focused on linear regression. Given the nonlinear nature of the climate system and 8 sources of model error, nonlinear analogues of these post-processing methods may offer 9 considerable improvements. The current study tests this hypothesis, applying both linear 10 and nonlinear regression to the correction of climate hindcasts produced with general 11 circulation models. Results indicate that nonlinear support vector regression is better able 12 to extract indices of the Pacific-North American pattern and the North Atlantic 13 Oscillation from coupled model output, while linear approaches are better suited to 14 atmosphere-only model output. Statistically significant predictions are produced at lead 15 times of up to nine months, and can be obtained from model output with no forecast skill 16 prior to processing.

# 17 **1. Introduction**

Numerical seasonal climate forecasts are becoming increasingly common, as the research community experiments with new tools and techniques while various user groups explore their potential value. Building on improvements in climate models and computing power realized in recent years, several prominent forecasting centers have developed operational systems to predict the climate state at lead times of several months to several seasons (Derome et al. 2001; Livezey and Timofeyeva 2008; Saha et al. 2006). A common

1	approach has been to force an atmospheric general circulation model (GCM) with sea
2	surface conditions observed at the beginning of a forecast period, recognizing that the
3	ocean boundary captures various slow-varying processes with strong atmospheric
4	influences. Several of these atmosphere-only (AGCM) forecasting systems have
5	demonstrated modest skill at lead times of one month to a season, during which sea
6	surface conditions are unlikely to change significantly. These early successes have
7	stoked interest in systems based on coupled atmosphere-ocean GCMs (AOGCMs). The
8	addition of an evolving ocean state has the potential to extend the useful lead-time of a
9	forecast, if the coupled model responds to initial conditions in a realistic manner. In
10	Canada, an AOGCM-based forecasting system has been developed through the Coupled
11	Historical Forecasting Project (CHFP). This effort is a natural extension of the first and
12	second Historical Forecasting Projects (HFP & HFP2), which formed the basis for
13	Environment Canada's current AGCM-based forecasting system (Merryfield et al. 2010).
14	The finite inherent predictability of the climate system presents one of the primary
15	challenges to seasonal forecasting. The influence of the slow-varying climate processes
16	used in dynamical forecasting is often much smaller than the climate's internal variability;
17	i.e. predictable signals are obscured by climate noise. The magnitude of this problem
18	varies by region; for example, predictability is high in the tropics relative to the
19	extratropics (Rowell 1998). Skillful extratropical forecasts are therefore often associated
20	with accurate prediction of processes with the highest signal-to-noise ratios. Much of the
21	seasonal forecast skill in North America can be attributed to the influence of the El Nino-
22	Southern Oscillation (ENSO; e.g. Derome et al. 2001; Livezey and Timofeyeva 2008;
23	Shukla et al. 2000), which generates a highly predictable signal over much of the

1 continent. Other teleconnection patterns have similar, if smaller, value as North

2 American climate predictors, including the North Atlantic Oscillation (NAO; e.g. Hurrell

3 and van Loon 1997) and the Pacific-North American Pattern (PNA; e.g. Leathers et al.

4 1991), among others.

5 Other key challenges in numerical forecasts involve the quality of the models being used. 6 All models produce systematic error in their representation of the climate system and its 7 integration from initial conditions. Limitations in resolution and model physics push the 8 model integrations towards a biased interpretation of the true climate, and this drift 9 obscures the simulated evolution from initial conditions. Furthermore, the response itself 10 may deviate considerably from that of the true climate. Although a model may produce a 11 predictable, systematic response to prescribed forcings, the spatial characteristics of this 12 response may differ from that of the true climate. That is, a GCM's response may be 13 statistically related to the true outcome, yet produce forecast fields with limited accuracy. 14 These problems are exacerbated in coupled GCMs, as nonlinear interactions between 15 component models can amplify any existing biases.

16 Much of the effort to improve numerical forecasts focuses on improving the models used. 17 This work is necessary, but also expensive and time intensive. A complementary, and 18 less resource intensive, approach involves the development of post-processing methods 19 that identify and correct systematic model biases. Post-processing techniques typically 20 use computationally efficient linear statistical analyses, such as regression, to relate GCM 21 projections with true climate outcomes. This process has proven useful in improving 22 temperature, precipitation, and teleconnection forecasts, (Jia et al. 2010; Lin et al. 2005; 23 Lin et al. 2008). However, given that both the climate system and sources of GCM bias

1	involve nonlinear processes, it is possible that the use of linear statistics has limited the
2	effectiveness of this approach. If this is the case, nonlinear analogues of the statistical
3	methods applied in past studies should give superior results. In many cases, appropriate
4	nonlinear methods exist, and have already proven useful in climate and weather
5	applications. In particular, nonlinear regression techniques are well developed, and can
6	be readily adapted to numerical climate forecast correction. Past studies have
7	successfully applied nonlinear regression to seasonal forecasting of El Nino events (Wu
8	et al. 2006), extreme precipitation (Zeng et al. 2011), and air temperatures (Wu et al.
9	2005), among other climate and weather applications.
10	The current study assesses the potential value of nonlinear regression in climate forecast
11	post-processing, applying the method to HFP2 and CHFP hindcast data. Classical
11 12	post-processing, applying the method to HFP2 and CHFP hindcast data. Classical multivariate linear regression is compared against support vector regression, a machine
11 12 13	post-processing, applying the method to HFP2 and CHFP hindcast data. Classical multivariate linear regression is compared against support vector regression, a machine learning method with nonlinear capabilities. The present study focuses on prediction of
11 12 13 14	<ul> <li>post-processing, applying the method to HFP2 and CHFP hindcast data. Classical</li> <li>multivariate linear regression is compared against support vector regression, a machine</li> <li>learning method with nonlinear capabilities. The present study focuses on prediction of</li> <li>key teleconnection patterns influencing North American climate; these predictions are</li> </ul>
11 12 13 14 15	<ul> <li>post-processing, applying the method to HFP2 and CHFP hindcast data. Classical</li> <li>multivariate linear regression is compared against support vector regression, a machine</li> <li>learning method with nonlinear capabilities. The present study focuses on prediction of</li> <li>key teleconnection patterns influencing North American climate; these predictions are</li> <li>themselves of limited operational value, but are strongly related to forecast variables of</li> </ul>
11 12 13 14 15 16	<ul> <li>post-processing, applying the method to HFP2 and CHFP hindcast data. Classical</li> <li>multivariate linear regression is compared against support vector regression, a machine</li> <li>learning method with nonlinear capabilities. The present study focuses on prediction of</li> <li>key teleconnection patterns influencing North American climate; these predictions are</li> <li>themselves of limited operational value, but are strongly related to forecast variables of</li> <li>value and provide a simple illustration of the post-processing methods' ability to extract a</li> </ul>
11 12 13 14 15 16 17	post-processing, applying the method to HFP2 and CHFP hindcast data. Classical multivariate linear regression is compared against support vector regression, a machine learning method with nonlinear capabilities. The present study focuses on prediction of key teleconnection patterns influencing North American climate; these predictions are themselves of limited operational value, but are strongly related to forecast variables of value and provide a simple illustration of the post-processing methods' ability to extract a useful signal from forecast noise.

**2. Data** 

Post-processing methods are tested here using numerical climate hindcasts produced for
the first phase of Environment Canada's Coupled Historical Forecasting Project (CHFP;
Merryfield et al. 2010) and second Historical Forecasting Project (HFP2; Kharin et al.
2009), downloaded from the Canadian Centre for Climate Modelling and Analysis

(CCCMA). The simulations examined were generated with version three of the CCCMA
 atmospheric GCM (AGCM3), a global spectral model run at T63 horizontal resolution
 with 32 vertical levels (Scinocca et al. 2008). The two prediction systems both generate
 ensembles of 10 parallel integrations for each hindcast period, facilitating comparison of
 their predictive skill. Brief descriptions of the initialization and ensemble generation
 schemes follow; a more detailed discussion is available in Merryfield et al. (2010).

7 The complete set of HFP2 forecasts includes output from four atmospheric models. For 8 the purposes of the current student, only forecasts generated with AGCM3 have been 9 used, in order to ease comparison with CHFP1. These hindcasts were generated by 10 forcing AGCM3 with static sea surface temperature (SST) anomalies, observed during 11 the month preceding a simulation's start date. Sea ice extents were also initialized to the 12 previous month's observations, then relaxed to climatology over the first fifteen days of 13 the simulation. Hindcasts extend four months from the start date, a period over which 14 SST anomalies are expected to change relatively little. The ensemble is generated using 15 different atmospheric initial conditions, taken from the NCEP-NCAR reanalysis (Kalnay 16 et al. 1996) at 12 hr intervals preceding the simulation start date. Twelve hindcasts were 17 generated for each year, beginning at the first of every month. Data cover the period 18 1969-2001 (Kharin et al. 2009).

CHFP1 uses the version 3 of the CCCMA Coupled GCM (CGCM3), which couples the
AGCM3 with the National Center for Atmospheric Research Community Ocean Model
(NCOM; Gent et al. 1998). The inclusion of a dynamic ocean boundary is intended to
improve predictive skill at long lead times, and CHFP1 hindcasts have a twelve month
duration. The 10 member ensemble was produced by combining two ocean

1	initializations, taken one and five days prior to the start date, and five atmospheric
2	initializations, taken at 00:00Z from the five days preceding the simulation. These states
3	were produced by a coupled model run in which SSTs were relaxed strongly to observed
4	values leading up to the forecast, so that the atmospheric states contain the influence of
5	observed SSTs. (Unlike the HFP2 initial states they do not contain details of the
6	atmospheric state as represented by the NCEP-NCAR reanalysis at those times; this
7	difference tends to degrade CHFP1 skills relative to HFP2 in the DJF period.) Four
8	hindcasts were produced for each year, beginning on the first of December, March, June,
9	and September respectively (Merryfield et al. 2010). Data cover 1972-2001.
10	The present study examines December (winter) runs valid for 1972-2001, covering the
11	largest coincident period of the HFP2 and CHFP1 experiments. Winter runs were chosen
12	due to the greater influence of synoptic weather systems and teleconnection indices on
13	winter climate relative to other seasons. Stronger upper level westerlies at this time of
14	year favor Rossby wave propagation, which in turn is influenced by teleconnections such
15	as ENSO. This generally results in more skillful seasonal forecasts relative to other
16	seasons, and renders it particularly well suited to numerical climate prediction. Results
17	are evaluated against data from the NCEP-NCAR reanalysis and time series of prominent
18	teleconnection patterns from the National Oceanographic and Atmospheric
19	Administration's Climate Prediction Center (NOAA-CPC;
20	http://www.cpc.ncep.noaa.gov). Loading patterns in the 500 hPa geopotential height
21	(Z500) field for the North Atlantic Oscillation (NAO) and Pacific-North America (PNA)
22	pattern were generated by regressing NAO and PNA time series from NOAA-CPC onto
23	gridded anomaly fields from NCEP-NCAR reanalysis, producing regressed anomaly

("loading") maps that were subsequently normalized to a unit vector (Figures 1a and 2a).
Separate loading patterns were calculated for overlapping three-month periods (e.g. DJF,
JFM, FMA, etc) in order to reflect seasonal variations in the structure and intensity of the
two teleconnections. CHFP1 and HFP2 predictions of the teleconnection indices were
then calculated as the projection of the loading patterns onto the ensemble mean Z500
anomaly hindcast for each three-month period.

# 7 **3. Methods**

# 8 *3.a Support Vector Regression:*

9 Support vector regression (SVR) is an extension of the support vector machine method
10 (developed originally for classification problems) to regression problems (Vapnik 1997).
11 A brief description of the technique is described here; more detailed explanation can be
12 found in Hsieh (2009).

Let x be a vector of predictor variables and y be a desired predictand. Linear
relationships between y and x can be readily identified through classical linear regression,
while nonlinear relationships cannot. However, a nonlinear mapping function, Φ, can be
used to convert the nonlinear relationship into a linear regression problem, i.e.

17 
$$f(\mathbf{x},\mathbf{w}) = \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle + b,$$
 (1)

18 where  $\langle , \rangle$  denotes the inner product, and **w** and *b* are the regression coefficients that 19 minimize the error between predictions *f* and observations *y*. SVR calculates an  $\varepsilon$ -20 insensitive error norm as:

1 
$$|y - f(\mathbf{x}, \mathbf{w})|_{\varepsilon} = \begin{cases} 0, & \text{if } |y - f(\mathbf{x}, \mathbf{w})| < \varepsilon \\ |y - f(\mathbf{x}, \mathbf{w})| - \varepsilon, & \text{otherwise} \end{cases}$$
, (2)

2 i.e., the error is ignored when smaller than  $\varepsilon$ , and approximates the mean absolute error 3 when large. For a given training data set with sample size N, SVR calculates w and b 4 coefficients that minimize the following error function R,

5 
$$R = \frac{C}{N} \sum_{i=1}^{N} \left| f(\mathbf{x}_{i}, \mathbf{w}) - y_{i} \right|_{\varepsilon} + \frac{1}{2} \left\| \mathbf{w}^{2} \right\| , \qquad (3)$$

6

13

7 where C is a parameter that controls the relative influence of the two terms in the error 8 function. Large values of C emphasize the first ( $\varepsilon$ -insensitive error) term at the expense 9 of the second (regularization) term. Conversely, low values of C lead to a more 10 prominent regularization term, penalizing (i.e. suppressing) complex models which tend 11 to need weight vectors w with larger magnitude. Both C and  $\varepsilon$  are user-selected 12 parameters, commonly referred to as hyperparameters.

By linearizing a nonlinear regression problem via (1), SVR eliminates the need for 14 nonlinear optimization and problems associated with local minima of a nonlinear error 15 function. In this regard, SVR presents an advantage over the neural network approach to 16 nonlinear regression. The remaining concern is that the mapping function  $\Phi(x)$  may 17 transform the data into a high dimensional space, and thereby become computationally 18 expensive to solve. This problem is addressed through the use of a kernel trick, which replaces the inner product  $\langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle$  with a scalar kernel function  $K(\mathbf{x}, \mathbf{x}')$ . For the 19

purposes of this study, a linear and a nonlinear kernel have been used for linear and
nonlinear SVR, respectively. The linear kernel is simply the inner product *K*(**x**, **x**<sub>i</sub>) = <**x**, **x**<sub>i</sub>>. While the linear kernel does not allow the SVR to model nonlinear relations, linear
SVR is not the same as the classic linear regression, as the ε-insensitive error norm
renders it a robust linear regression method. Nonlinear SVR was performed with the
following Gaussian kernel:

7 
$$K(\mathbf{x},\mathbf{x}_{i}) = \exp(\frac{-\|\mathbf{x}-\mathbf{x}_{i}\|^{2}}{2\sigma^{2}}).$$
 (4)

8

9 The use of this kernel adds  $\sigma$  as an additional hyperparameter, describing the width of the 10 kernel.

The selection of hyperparameters has a significant impact on the performance of the SVR model. Although objective methods for selecting some hyperparameters have been proposed, in practice model optimization requires additional tuning. Here, the approach of Cherkassky and Ma (2002) has been used to generate initial objective estimates of the hyperparameters, which are subsequently tuned using a fine grid search.

#### 16 3.b Post-processing Model Testing

17 Hindcasts of the NAO and PNA indices were post-processed with three regression

- 18 algorithms; multivariate linear regression (LR), linear support vector regression (LSVR),
- 19 and nonlinear support vector regression (NLSVR). Post-processing models use five
- 20 predictors: the teleconnection hindcast from the unprocessed model output, and the time

series of four leading patterns from one of the following decompositions of GCM output: 1 2 i) principal components (PCs) of the Z500 anomaly field over the northern hemisphere 3 (45-90°N), ii) PCs of tropical Pacific SST anomalies (20°S-20°N, 120-270°E), or iii) a 4 singular value decomposition (SVD) (i.e. maximum covariance analysis; von Storch and 5 Zwiers 2001) of these two fields, giving spatial patterns for both. In the present paper, 6 anomaly refers to deviations from climatology, with monthly mean fields subtracted from 7 the raw data. As HFP2 does not predict the ocean state, SST predictors for HFP2 were 8 derived from the static SST anomaly driving the AGCM3. CHFP1 predictors were 9 derived from the forecast state of both SSTs and Z500. The predictor sets were chosen 10 for their similarity to sets used in previous post-processing studies (e.g. Lin et al. 2008), 11 and for their relative simplicity. The first step in model training selects one of the four 12 PC or SVD time series as the best predictor set, which is then used in subsequent training 13 steps. The result is that individual post-processing models use the unprocessed 14 teleconnection index with additional spatial data related to a single field (SST or Z500), 15 summarized with a single decomposition method (PC or SVD). Although it is possible 16 that combining various predictors from the three GCM decompositions could generate 17 better forecasts, there are concerns that the skill of forecasts produced with screened 18 predictor sets can be overstated (DelSole and Shukla 2009). The approach used here 19 represents a compromise between limiting the complexity of predictor screening, while 20 maintaining some flexibility to account for the possibility that the relative value of the 21 ocean and atmospheric component models may change as a coupled run progresses. 22 For CHFP1 data, separate regression models were built for each of ten overlapping three-

23 month periods. Hindcasts used begin in December, and run through to the end of the

1 following November. The run is then split into ten forecast periods;

December/January/February (DJF), JFM, etc. through to SON. Predictors are derived
from hindcast data valid for a given forecast period; i.e. the DJF post-processing model
uses the hindcast output averaged over DJF. DJF forecasts represent prediction at a zeromonth lead, JFM at one-month lead, etc. For each lead time, the best predictor set was
first identified, using a coarse 5-fold cross-validation linear regression procedure to
compare their relative value.

8 The selected predictor set was then used by the LR, LSVR and NLSVR methods for post-9 processing. The use of linear regression to select predictor sets may give LR and LSVR a 10 slight advantage over NLSVR, hence any improved skill by NLSVR over the linear 11 methods is likely to be understated. The three methods were each evaluated using a 12 rigorous double cross-validation scheme, illustrated in Figure 3. Thirty-fold cross-13 validation produced predictions for each year, using a model trained with the majority of 14 the remaining years in the study period (outer round, or CV1 in Figure 3). Years 15 immediately preceding and following the year being predicted were excluded from the 16 training period, in order to reduce information leakage from autocorrelation. 17 In the case of LSVR and NLSVR, hyperparameters were selected using a second nine-18 fold cross-validation applied to the training data (inner round, or CV2 in Figure 3). 19 Values identified as optimal for the training data were then used in an SVR model to 20 predict the testing year in CV1. This double cross-validation scheme is employed to 21 avoid validating forecast skills with data already used for model training and 22 hyperparameter selection.

#### 1 4. Results

# 2 *4.a Teleconnection Loading Patterns & Relationship to Predictors*

3 It is helpful to begin with an examination of the physical structure of the teleconnections 4 being predicted, and considering the ability of the forecast models to reproduce them. 5 Figures 1a) and 2a) respectively show the Z500 loading patterns of the PNA and NAO 6 used in the current study, for the DJF period. The PNA features two opposing centres of 7 action, located over the North Pacific and Western North America. Similarly, the NAO 8 places opposing centres over the climatological Iceland Low in the North Atlantic and the 9 Azores High in the subtropical Atlantic. The spatial characteristics of these correlation-10 derived patterns are essentially identical to those used by the NCDC. Time series of the 11 amplitudes of these patterns also agree very closely to the original NAO and PNA time 12 series used to define them, giving correlation scores above 0.9. 13 Past studies have shown that an SVD of the Z500 and tropical SST fields used to generate 14 predictor sets in the current work produce patterns related to the teleconnections of 15 interest. Specifically, the first Z500 SVD (SVD1) shares spatial characteristics with the 16 PNA, while the second (SVD2) typically resembles the NAO (Lin & Derome, 2005). 17 This emphasizes the influence of coupled atmosphere-ocean processes on these 18 teleconnections, and provides insight into the physical mechanisms behind these 19 phenomena. In the context of the present study, it follows that an SVD analysis of GCM 20 output can provide an indication of whether the model accurately simulates these 21 mechanisms.

1	Results of SVD analyses of NCEP data (Figures 1b & 2b), HFP2 output (Figures 1c &
2	2c), and CHFP1 output (Figures 1d & 2d) for DJF are presented alongside the
3	teleconnection loadings. The NCEP SVDs both feature centres of action similar to the
4	associated loading patterns. However, while the teleconnection patterns de-emphasize
5	regions outside of these centres, the NCEP SVDs also place high loadings in other areas
6	of the Northern Hemisphere. The two teleconnections appear to be somewhat convoluted
7	in the SVD analysis, with SVD1 emphasizing the PNA centres but also introducing weak
8	loadings in the vicinity of the Iceland Low and Azores High. SVD2 emphasizes the
9	Iceland Low and Azores High, but shifts their location relative to the teleconnections. It
10	also introduces a dipole between western and eastern North America, resembling a
11	weakened and eastwards shifted PNA pattern.
12	SVD results for HFP2 output closely resemble NCEP results, particularly for SVD1.
13	There is again a tendency to convolve PNA and NAO-like patterns, although the NAO
14	dipole is less well-defined. However, CHFP1 output gives very different results. CHFP1
15	SVDs feature weaker centres of action, and less defined features in general than NCEP or
16	HFP2. This could result from greater variability in atmosphere-ocean interactions,
17	perhaps related to the adjustment process of the component models during the early
18	stages of the coupled runs. Despite this, CHFP1 SVD1 closely resembles the PNA
19	loading pattern, with a well-defined Pacific/North America dipole. The NAO is not
20	captured well in either SVD, with SVD2 showing only an Azores High combined with a
21	weaker PNA dipole. It can be inferred that the model captures PNA-like processes well,
22	but may be misrepresenting those related to the NAO

1	Table 1 summarizes relationships between the SVDs and the teleconnection indices,
2	showing correlation scores between SVD time series and observed PNA and NAO
3	indices. The strong relationship between NCEP SVD1 and the PNA is confirmed with
4	high values in both DJF and JFM. The PNA-like aspects of NCEP SVD2 also leads to
5	reasonably strong correlations. However, neither SVD is strongly related to the observed
6	NAO. HFP2 SVD1 also correlates reasonably well with the observed PNA, suggesting
7	this is a relatively useful predictor. This is not the case for CHFP1; although the CHFP1
8	SVD1 spatial pattern resembles the PNA, the temporal amplitude of this pattern in the
9	CHFP1 output does not provide a useful estimate of the teleconnection's state.
10	Surprisingly, SVD2 does show a statistically significant, if weak, relationship with the
11	observed NAO, suggesting there is a signal related to the NAO embedded in a
12	structurally dissimilar pattern.
13	Table 2 gives the number of times the four predictor sets described in the previous section
14	were identified as the best basis for a post-processing model, using the linear regression-
15	based screening process. Z500 SVDs were selected most often, likely due to their
16	marked relationship to the PNA and NAO.
17	4.b Post-processing Results
18	Performance of unprocessed and processed model simulations are primarily evaluated in

19 terms of correlation between observed and projected teleconnection indices, with results

- 20 presented in Tables 3 (HFP2) and 4 (CHFP1). Post-processing of CHFP1 data are also
- 21 evaluated using correlation, mean absolute error (MAE), and index of agreement

(Willmott et al. 1985), reported as skill scores (SS) relative to the unprocessed CHFP1
 output (Figures 4 & 5), with SS defined by

3 
$$SS = \frac{A - A_{ref}}{A_{perfect} - A_{ref}},$$
 (5)

where A is a particular accuracy measure, A<sub>ref</sub> is the value of A from a reference forecast
set (here the unprocessed CHFP1 output), and A<sub>perfect</sub> is the value for a set of perfect
forecasts. As presented here, SS illustrates the improvements of post-processing methods
relative to raw CHFP1 output.

8 It is useful to first compare performance of the raw model projections from the two 9 forecasting systems, both to provide a baseline for post-processing assessment and to 10 evaluate the relative strengths of the AOGCM and AGCM approaches. HFP2 performs 11 reasonably well for both teleconnections at all lead times (Table 3). Correlation scores 12 larger than 0.33 indicate statistical significance at the 95% confidence level, and this 13 value is exceeded by all unprocessed HFP2 forecasts. However, in absolute terms these 14 correlations are relatively low, with projections capturing only 9-25% of teleconnection 15 variability. CHFP1 output is considerably less reliable (Table 4). Although capable of 16 predicting the PNA index relatively well at short lead times (>0.5 for JFM & FMA), 17 CHFP1 scores decrease rapidly at longer leads. NAO simulation is uniformly poor, with 18 the model producing no statistically significant positive correlations. These results 19 suggest the CHFP1 offers a negligible improvement over the simpler HFP2 forecasting 20 system, extending useful PNA forecasts by one forecast period but degrading simulation 21 of the NAO. Although not shown, MAE values are comparable for the two forecasting

systems, ranging from 0.5-1 of the observed teleconnection's standard deviation for all
 forecast periods.

3 The post-processing methods used here vary considerably in their effectiveness when 4 applied to HFP2 output (Table 3). Post-processing reduces correlation scores to 5 statistically insignificant values for the NAO, but all methods tend to improve PNA 6 forecasts. This may be due to the fact that the additional predictors (the SVD and PC 7 time series from Z500 and tropical Pacific SST) are less directly related to NAO than to 8 PNA, and extra irrelevant predictors tend to decrease skills. The two linear methods give 9 comparable results and outperform NLSVR, which fails to improve the DJF PNA 10 prediction. The post-processing of CHFP1 output is considerably more effective (Table 11 4; Figures 4 & 5). Although the methods used here have difficulty improving the few 12 statistically significant correlations found in the raw forecasts (e.g. PNA at JFM & FMA), 13 they frequently improve on the poorer CHFP1 results. In the case of the NAO, the only 14 significant results are generated through post-processing which yields relatively high 15 correlations even at a nine-month lead (SON). Similarly, relatively high PNA 16 correlations are generated for most lead times by one or more of the methods. Although 17 no post-processing method consistently generates the highest correlations, NLSVR has 18 the greatest success rate, giving the best scores for five three-month periods for the NAO 19 and four for the PNA, accounting for almost half of the twenty post-processing periods. 20 By contrast, linear regression gives the best correlation in four instances (FMA, MJJ, JJA, 21 and ASO for the NAO), and LSVR performs best in another four cases (NAO at DJF, and 22 PNA at MAM, AMJ & SON). Unprocessed CHFP1 output gives the best results in only 23 three instances (PNA at JFM, FMA and MJJ). It should be noted that the statistically

significant post-processed correlations still capture only 10-30% of teleconnection
 variability. With the 30 year sample available, it is difficult to assess the source of this
 skill. However, the ability of all methods to generate significant correlations at extended
 leads suggests a useful signal is embedded in the model output.

5 Figure 4 also presents mean absolute error (MAE) skill, relative to unprocessed CHFP1

6 forecasts. As with correlation, MAE results indicate post-processing generally improves

7 CHFP1 forecasts, but the relative advantage of any single method is less clear. LR

8 marginally outperforms NLSVR, giving the lowest MAE for eight prediction periods to

9 NLSVR's six. LSVR and the unprocessed model both give the best results for only three

10 periods.

11 To better illustrate the relative advantages of various post-processing methods,

12 Willmott's index of agreement (IA) has also been used to assess the skill of CHFP1-

13 derived forecasts (Willmott et al. 1985), with

15 where

16 
$$\begin{aligned} f' &= f - \overline{y} \\ y' &= y - \overline{y} \end{aligned}$$
 (7)

with *f* the predicted value, and *y* the observed. IA incorporates elements of MSE andcorrelation into a single value between zero and one, with one being a perfect score.

Results are given in Figure 5, again presented as skill scores relative to the raw model
forecast. This measure emphasizes the value of nonlinear post-processing, with NLSVR
forecasts improving CHFP1 output in all but five cases, and giving the highest positive
skill in nine. By contrast, LR and LSVR reduce raw CHFP1 skill in a substantial number
of cases (seven each), and combined provide the highest positive skill in only seven
instances.

7

# 8 5. Summary & Discussion

9 Results suggest regression-based post-processing is most effective when applied to 10 coupled model output, improving CHFP1 output considerably but only improving PNA 11 forecasts when applied to HFP2. The inference is that CHFP1 output contains a useful 12 signal that is obscured by biases in the model's integration from a prescribed initial 13 condition. Although CHFP1 may not faithfully reproduce loading patterns associated 14 with the NAO and PNA, it does produce a predictable evolution related to the observed 15 teleconnections. This relationship appears to be nonlinear, allowing NLSVR to 16 outperform linear post-processing methods in the majority of cases. The lower 17 performance of LSVR relative to NLSVR confirms it is the nonlinear capabilities that 18 give this method an advantage, rather than the SVR methodology itself. 19 Results indicate that HFP2 captures the structure of the NAO and PNA better than 20 CHFP1, and recreates teleconnection time series relatively well prior to post-processing. 21 This is related to the use of an atmosphere-only GCM, which removes the influence of 22 ocean model and coupled atmosphere/ocean error present in the CHFP1 simulations.

1	However, regression-based post-processing improved only PNA forecasts from HFP2
2	output, a result not seen in past studies using similar methods (Lin et al. 2005). Reduced
3	performance in the current study may be related to differences in predictor selection, the
4	time period examined, and the double cross-validation methodology employed here. It is
5	also possible that the HFP2 predictors supplied in this study are not sufficiently related to
6	the NAO; with the addition of irrelevant predictors, it is expected that forecast skills will
7	drop. However, similar predictor sets have produced better skill in previous studies (e.g.
8	Lin et al. 2005). Additional sensitivity experiments indicate results are most sensitive to
9	variations in the cross-validation approach, such as changes in the number of cross-
10	validation steps and relaxing the auto-correlation buffer in the outer CV round (Figure 3).
11	In contrast, using predictors and periods identical to those in past studies did not change
12	results appreciably. As the LR-processed skill is lower in the present study than
13	previously reported suggests the approach used here gives a conservative assessment of
14	performance. It follows that performance of processed CHFP1 output and SVR-
15	processed results may also be underestimated, adding weight to the argument that
16	nonlinear post-processing offers valuable improvements to this data set.
17	Much of the skill in current North American climate forecasts has been attributed to the
18	El Nino/Southern Oscillation (Livezey and Timofeyeva 2008), which produces a
19	persistent, predictable response in large-scale climate. ENSO also exerts a strong
20	influence on the PNA, and has been tentatively associated with the NAO (Hoerling et al.
21	2001). In order to assess the influence of ENSO on post-processed forecasts, correlation
22	scores for the CHFP1 were recalculated without years initialized with a strong ENSO
23	signal, defined as a NINO3 index of magnitude 1.5 or higher; results are presented in

1	Figure 5. Removing strong ENSO events reduces raw CHFP1 skill somewhat, leaving
2	only a single significant correlation (PNA at FMA). Post-processing scores are also
3	generally reduced, particularly the LR and LSVR results. However, NLSVR continues to
4	generate statistically significant correlations for almost half of the forecasts. A
5	substantial portion of the skill generated by nonlinear post-processing is apparently
6	independent of the persistent, predictable influence of strong ENSO events. This
7	approach therefore appears to generate the most versatile and robust forecasts, while the
8	raw model, LSVR, and LR forecasts are less suitable for non-ENSO years.
9	Although the present study demonstrates the value of SVR as post-processing tool in
10	climate forecasting, significant obstacles remain if it is prove useful operationally. Of
11	particular concern is the method's greater complexity relative to traditional linear
12	regression, and a lack of familiarity with this method in the meteorological community.
13	For SVR post-processing to be adopted by forecasting centers, a clear advantage over
14	linear regression must be displayed. The results presented here suggest the additional
15	complexity yields valuable improvements, even when tested on a relatively short set of
16	hindcast experiments. Although PNA and NAO forecasts are of little operational value
17	themselves, these teleconnections are strongly connected with surface temperature and
18	precipitation over North America. The implication is that SVR post-processing should
19	also be capable of improving climate forecasts of these more valuable fields. A cell-by-
20	cell post-processing approach, similar to that presented by Lin et al. (2008), could be used
21	to test this hypothesis. However, such an exercise with SVR is computationally
22	intensive, and beyond the scope of the current paper. Still, the current work suggests the
23	method has some value and could be used to improve seasonal forecasts considerably.

# 2 Acknowledgements

- 3 J. Finnis and W. Hsieh are grateful for the support from the Canadian
- 4 Foundation for Climate and Atmospheric Sciences via project GOAPP. W.
- 5 Hsieh is supported by a Discovery Grant from the Natural Sciences and
- 6 Engineering Research Council of Canada.

		DJF		JFM	
	_	SVD1	SVD2	SVD1	SVD2
	Obs.	0.65	-0.48	0.88	-0.40
PNA	HFP2	0.41	-0.27	0.57	-0.16
	CHFP1	0.24	-0.31	0.29	-0.19
	Obs.	0.21	0.32	0.28	0.06
NAO	HFP2	-0.23	-0.13	0.32	-0.17
	CHFP1	0.26	0.33	0.16	0.02

3 Table 1: Pearson correlation scores between teleconnection time series (PNA and NAO),

4 from observed, HFP2 and CHFP1 data, and the first two Z500 SVD time series from

5 observed (NCEP) data. Correlations significant at the 95% are in bold type.

-	Decomposition Times Selected		Selected
	Method	NAO	PNA
SST	PC	4	1
	SVD	1	1
Z500	PC	0	2
-	SVD	5	6

Table 2: Frequency with which the four predictor sets were chosen for model training, as a best set is chosen for each of the ten lead times and the two teleconnections.

	NAO Projections					
	HFP2	LR	LSVR	NLSVR		
DJF	0.361	0.112	0.185	0.063		
JFM	0.361	0.184	0.033	-0.093		
		PNA I	Projections			
-	HFP2	PNA I LR	Projections LSVR	NLSVR		
DJF	HFP2 <b>0.396</b>	PNA I LR <b>0.440</b>	Projections LSVR 0.473	NLSVR 0.354		
DJF JFM	HFP2 0.396 0.487	PNA I LR 0.440 0.621	Projections LSVR 0.473 0.607	NLSVR 0.354 0.577		

3 Table 3: Pearson correlation scores between the observed NAO and PNA time series and

4 projections from unprocessed and processed HFP2 data. All three month periods covered

5 by the HFP2 data set are shown. Correlations significant at the 95% level are in bold

6 type.

1			

		NAO Pr	ojections	
	CHFP	LR	LSVR	NLSVR
DJF	0.234	0.313	0.402	0.326
JFM	0.191	0.056	0.179	0.566
FMA	0.150	0.202	0.084	0.013
MAM	0.076	0.200	0.313	0.466
AMJ	0.169	0.298	0.243	0.477
MJJ	-0.230	0.241	0.130	-0.163
JJA	-0.214	0.417	0.163	0.267
JAS	-0.260	0.337	0.183	0.409
ASO	-0.371	0.393	0.292	0.358
SON	-0.178	0.140	0.266	0.332
Mean	-0.043	0.260	0.225	0.305
	PNA Projections			
	CHFP	LR	LSVR	NLSVR
DJF	0.242	0.279	0.309	0.340
JFM	0.522	0.390	0.277	0.344
FMA	0.573	0.520	0.529	0.541
MAM	0.374	0.388	0.403	0.277
AMJ	0.084	0.287	0.442	0.376
MJJ	0.126	-0.123	-0.214	-0.105

JJA

JAS

ASO

SON

Mean

0.320

0.264

-0.274

0.023

0.225

3

4 Table 4: Pearson correlation scores between the observed NAO and PNA time series and

0.298

0.326

-0.197

-0.060

0.211

0.257

0.381

-0.291

0.149

0.224

0.364 0.402

0.236

0.118

0.289

5 projections from unprocessed and processed CHFP1 data. All three month periods

covered by the CHFP1 data set are shown. Correlations significant at the 95% level are 6 7 in bold type.









- 3 Figure 2: (a) The DJF Z500 NAO loading pattern (unitless), with the second Z500 SVD
- 4 pattern from (b) NCEP, (c) HFP2 and (d) CHFP1, in units of meters.



7 Figure 3: Illustration of the double cross-validation approach to model training and

8 testing. An outer round of validation (CV1) is used to generate hindcast time series,

9 using models with hyperparameters selected through an inner round of cross-validation

10 (CV2). In CV1, the training data are shown in grey and the 1-year validation data

shaded. The 1-year data segments (shown in white) bridging the training data and the
validation data are not used, to avoid autocorrelation leaking

13 information from the training data to the adjacent validation data. The validation data

14 segment is moved repeatedly in 1-year increments from the start of the data record to the

end in this cross-validation loop, so forecast performance is validated over the whole

16 record. Meanwhile in CV2, the training data from CV1 are assembled and divided into 9

17 segments, with 8 used for training and one (shaded) for validation. Again the training and

18 validation segments are rotated in the loop so all segments are eventually used for

19 validation to determine the optimal model values/hyperparameters. The optimal model

20 determined from CV2 is then used to forecast over the 1-year validation segment in CV1.









3 Figure 5: Same as Figure 4, but showing index of agreement skill scores.



2 Figure 6: Correlation scores of raw and post-processed CHFP1 predictions, calculated

3 without years initialized during moderate-to-strong ENSO events.

4

1 2	
3	
4	References
5 6 7	Cherkassky, V., and Y. Q. Ma, 2002: Selection of meta-parameters for support vector regression. <i>Lect Notes Comput Sc</i> , <b>2415</b> , 687-693.
8 9	DelSole, T., and J. Shukla, 2009: Artificial Skill due to Predictor Screening. <i>J Climate</i> , <b>22</b> , 331-345, doi: 10.1175/2008jcli2414.1.
10 11	Derome, J., and Coauthors, 2001: Seasonal predictions based on two dynamical models. <i>Atmos Ocean</i> , <b>39</b> , 485-501.
12 13 14	Gent, P. R., F. O. Bryan, G. Danabasoglu, S. C. Doney, W. R. Holland, W. G. Large, and J. C. McWilliams, 1998: The NCAR Climate System Model global ocean component. <i>J Climate</i> , <b>11</b> , 1287-1306.
15 16	Hoerling, M. P., J. W. Hurrell, and T. Y. Xu, 2001: Tropical origins for recent North Atlantic climate change. <i>Science</i> , <b>292</b> , 90-92.
17 18	Hsieh, W., 2009: <i>Machine Learning Methods in the Environmental Sciences</i> . Cambridge University Press.
19 20	Hurrell, J. W., and H. van Loon, 1997: Decadal variations in climate associated with the north Atlantic oscillation. <i>Climatic Change</i> , <b>36</b> , 301-326.
21 22 23	Jia, X. J., H. Lin, and J. Derome, 2010: Improving Seasonal Forecast Skill of North American Surface Air Temperature in Fall Using a Postprocessing Method. <i>Mon Weather Rev</i> , <b>138</b> , 1843-1857, doi: 10.1175/2009mwr3154.1.
24 25	Kalnay, E., and Coauthors, 1996: The NCEP/NCAR 40-year reanalysis project. <i>B Am Meteorol Soc</i> , <b>77</b> , 437-471.
26 27 28	Kharin, V. V., Q. B. Teng, F. W. Zwiers, G. J. Boer, J. Derome, and J. S. Fontecilla, 2009: Skill Assessment of Seasonal Hindcasts from the Canadian Historical Forecast Project. <i>Atmos Ocean</i> , <b>47</b> , 204-223, doi: 10.3137/Ao1101.2009.
29 30 31	Leathers, D. J., B. Yarnal, and M. A. Palecki, 1991: The Pacific North-American Teleconnection Pattern and United-States Climate .1. Regional Temperature and Precipitation Associations. <i>J Climate</i> , <b>4</b> , 517-528.
32 33 34	Lin, H., J. Derome, and G. Brunet, 2005: Correction of atmospheric dynamical seasonal forecasts using the leading ocean-forced spatial patterns. <i>Geophys Res Lett</i> , <b>32</b> , -, doi: 10.1029/2005gl023060.

- 1 Lin, H., G. Brunet, and J. Derome, 2008: Seasonal forecasts of Canadian winter
- 2 precipitation by postprocessing GCM integrations. *Mon Weather Rev*, **136**, 769-783, doi:
- 3 10.1175/2007mwr2232.1.
- 4 Livezey, R. E., and M. M. Timofeyeva, 2008: The first decade of long-lead US seasonal
- 5 forecasts Insights from a skill analysis. *B Am Meteorol Soc*, **89**, 843-854, doi:
- 6 10.1175/2008bams2488.1.
- 7 Merryfield, W. J., W. S. Lee, G. J. Boer, V. V. Kharin, B. Pal, J. F. Scinocca, and G. M.
- 8 Flato, 2010: The First Coupled Historical Forecasting Project (CHFP1). *Atmos Ocean*,
- 9 **48,** 263-283, doi: 10.3137/Ao1008.2010.
- 10 Rowell, D. P., 1998: Assessing Potential Seasonal Predictability with an Ensemble of
- 11 Multidecadal GCM Simulations. J Climate, 11, 109-120, doi: 10.1175/1520-
- 12 0442(1998)011<0109:APSPWA>2.0.CO;2.

- 15 Scinocca, J. F., N. A. McFarlane, M. Lazare, J. Li, and D. Plummer, 2008: Technical
- Note: The CCCma third generation AGCM and its extension into the middle atmosphere.
   *Atmos Chem Phys*, 8, 7055-7074.
- Shukla, J., and Coauthors, 2000: Dynamical seasonal prediction. *B Am Meteorol Soc*, 81,
  2593-2606.
- 20 Vapnik, V., Golowich, S., Smola, A., 1997: Support vector method for function
- 21 approximation, regression estimation, and signal processing. *Advances in Neural*
- *Information Processing Systems*, M. Mozer, Jordan, M., Petsche, T., Ed., MIT Press, 281 287.
- von Storch, H., and F. W. Zwiers, 2001: *Statistical analysis in climate research*. 1st pbk.
  ed. Cambridge University Press, x, 484 p. pp.
- Willmott, C. J., and Coauthors, 1985: Statistics for the Evaluation and Comparison of
  Models. *J Geophys Res-Oceans*, **90**, 8995-9005.
- 28 Wu, A. M., W. W. Hsieh, and A. Shabbar, 2005: The nonlinear patterns of North
- American winter temperature and precipitation associated with ENSO. *J Climate*, 18, 1736-1752.
- 31 Wu, A. M., W. W. Hsieh, and B. Y. Tang, 2006: Neural network forecasts of the tropical
- 32 Pacific sea surface temperatures. *Neural Networks*, **19**, 145-154, doi:
- 33 10.1016/J.Neunet.2006.01.004.
- Zeng, Z., W. W. Hsieh, A. Shabbar, and W. R. Burrows, 2011: Seasonal prediction of
- winter extreme precipitation over Canada by support vector regression. *Hydrol Earth Syst Sc*, 15, 65-74, doi: 10.5194/Hess-15-65-011.

<sup>Saha, S., and Coauthors, 2006: The NCEP Climate Forecast System.</sup> *J Climate*, 19, 34833517, doi: 10.1175/JCLI3812.1.