

Université de Québec  
Institut national de la recherche scientifique (INRS)  
Eau, Terre et Environnement (ETE)

**Analyse et modélisation non-déterministe de la superficie  
horaire des refuges thermiques potentiels : étude de cas de la  
rivière Sainte-Marguerite (Québec, Canada)**

Mémoire présenté pour l'obtention du grade de  
Maîtrise ès sciences (M.Sc.)

Par

**Ilias Hani**

**Date du dépôt initial : 02 août 2022**

**Date du dépôt final : 28 septembre 2022**

**Jury d'évaluation**

Examinatrice externe

**Zoe Li**

**Université McMaster**

Examineur interne

**Saeid Homayouni**

**INRS - ETE**

Directeur de recherche

**André St-Hilaire**

**INRS - ETE**

Co-directeur de recherche

**Taha B.M.J. Ouarda**

**INRS - ETE**

**À mes grands-parents...**

## REMERCIEMENTS

J'aimerais tout d'abord remercier mon directeur de recherche, André St-Hilaire à qui je dois dans un premier temps mon passage à la maîtrise recherche, merci de m'avoir encadré, enseigné et dirigé tout au long de ce parcours et de m'avoir appris à être moins « bon élève » et plus autonome. Je voudrais aussi le remercier pour sa patience et sa disponibilité malgré les conditions difficiles des deux dernières années (pandémie). Merci pour ton engagement, ta passion et surtout tes judicieux conseils, qui ont contribué à alimenter ma réflexion.

Des remerciements à mon co-directeur Taha Ouarda pour son support continu, son encadrement dans la rédaction des articles et l'important bagage scientifique qu'il m'a légué. C'est en grande partie grâce à ses précieux enseignements et conseils que j'ai pu continuer jusqu'au bout. Je vous suis très reconnaissant, merci pour votre collaboration.

Une mention spéciale à Nassir El-Jabi et à Daniel Caissie de l'université de Moncton qui m'ont orienté pour le choix de mon programme de maîtrise et m'ont procuré les outils nécessaires durant mes cours de baccalauréat et mes stages Coop pour mener à bien ce projet.

J'adresse aussi mes remerciements à André Boivin qui m'a accompagné et aidé quotidiennement durant les campagnes de terrain à la rivière Ste-Marguerite au cours des étés 2020 et 2021. Un grand merci à tous les intervenants de près ou de loin qui œuvrent au sein du CIRSA. Ce fût une expérience enrichissante dans un cadre exceptionnel.

Je tiens également à remercier toute l'équipe St-Hilaire. Ce fût un privilège de travailler avec vous. Une mention particulière à Amandine St-Hilaire et Gaëlle Ricaud pour leur aide durant les campagnes de terrain, à mes collègues Olfa Abidi, Freddy Houndekindo, Mostapha Khorsandi et Eisenhower Vargas Rincon.

Finalement, merci à mes parents. Leur soutien inconditionnel, leur présence et leurs encouragements sont pour moi les piliers fondateurs de ce que je suis et de ce que je fais.

## RÉSUMÉ

En raison des températures élevées des rivières en été, de nombreux efforts sont déployés pour améliorer les conditions environnementales des espèces poïkilothermes, tels que les salmonidés. Pour éviter le stress causé par l'augmentation de la température des cours d'eau, ces poissons s'abritent dans des zones d'eau froide, appelées refuges thermiques. Les panaches de confluence des tributaires froids avec un cours d'eau plus chaud représentent des écosystèmes où la dynamique des cours d'eau est très complexe. Ces affluents ont tendance à être considérés comme une zone d'eau froide relativement constante au fil du temps, ce qui les rend très propice à la formation de potentiels refuges thermiques. Dans la présente étude, deux confluences de tributaires de la rivière Ste-Marguerite Québec, (Canada) ont été instrumentées durant deux étés consécutifs, afin d'estimer la superficie horaire des refuges thermiques potentiels et analyser sa variabilité temporelle à l'aide de méthodes non déterministes. L'aire des refuges thermiques potentiels (ARTP) est estimée par une méthode d'interpolation déterministe nommée la Pondération Inverse à la Distance (PID) appliquée aux températures maximales horaires de l'eau. Par la suite, différents modèles d'apprentissage automatique sont proposés pour prédire la superficie maximale horaire des refuges soit :

- (i) Le modèle additif général (GAM),
- (ii) Le modèle de régression multivariée par splines adaptatifs (MARS),
- (iii) La régression par machine à vecteur de support (SVM),
- (iv) La régression par forêt aléatoire (RF).

Les modèles ont été développés en utilisant entre trois et cinq variables hydrométéorologiques explicatives. Les critères utilisés pour évaluer la performance des algorithmes proposés dans la présente étude sont l'erreur moyenne quadratique relative (rRMSE) et le critère de Nash-Sutcliffe (NASH). Le modèle RF a indiqué les meilleures performances et a démontré une grande précision aux deux stations étudiées ( $rRMSE \leq 15\%$  et  $NASH \approx 93\%$ ). Les modèles SVM et RF ont démontré une forte capacité à fournir des estimations précises des ARTP comparés aux deux modèles de référence GAM et MARS; une étude antérieure connexe traitant des prévisions quotidiennes moyennes de l'ARTP. D'autre part, la variabilité diurne et intra-saisonnière de l'ARTP a été étudiée en

calculant le coefficient de variation (CV) sur différentes périodes de la journée avec une emphase sur l'étiage extrême de l'été 2021. Les résultats montrent un inversement de tendance dans la variabilité des superficies de refuges entre les mois de juillet et août à travers les deux étés. Les niveaux d'eau du mois d'août 2021 ont atteint des seuils critiques et les résultats montrent une distribution de données homogène ( $CV \leq 13\%$ ) avec des valeurs élevées et quasi constantes de l'ARTP. Les outils d'interaction du modèle GAM ont été utilisés dans le but de modéliser le cycle diurne de l'ARTP en fonction des composantes horaires et journalière (jour julien). Durant les conditions de faible débit, les résultats illustrent des fenêtres temporelles spécifiques pour les valeurs extrêmes quotidiennes des aires de refuges thermiques potentiels. Les modèles considérés dans cette étude s'avèrent d'une grande utilité pour améliorer les stratégies de gestion de pêches et aussi pour la conservation les habitats de certaines espèces de salmonidés qui subissent des pressions de plus en plus critiques.

**Mots clés:** L'aire des refuges thermiques potentiels; Température horaire de l'eau; Variabilité temporelle; modèles de régressions; apprentissage automatique.

# TABLE DES MATIÈRES

<b>REMERCIEMENTS</b> .....	<b>iii</b>
<b>RÉSUMÉ</b> .....	<b>iv</b>
<b>LISTE DES TABLEAUX</b> .....	<b>viii</b>
<b>LISTE DES FIGURES</b> .....	<b>ix</b>
<b>1- INTRODUCTION</b> .....	<b>1</b>
1.1- MISE EN CONTEXTE .....	2
1.2- VARIABLE HYDROMÉTÉOROLOGIQUES ET OUTILS DE MODELISATION .....	3
1.3- APPROCHE NON-PARAMETRIQUES .....	4
1.4- OBJECTIFS .....	6
1.5- STRUCTURE DU MEMOIRE.....	7
<b>2- SYNTHÈSE DE L'ARTICLE EN FRANÇAIS</b> .....	<b>8</b>
2.1- METHODOLOGIE ET SITE D'ÉTUDE .....	9
2.1.1- Estimation de l'ARTP maximale horaire .....	9
2.1.2- Algorithmes d'apprentissage automatique et modélisation horaire de l'ARTP.....	12
2.1.3- Procédure d'évaluation et mesures de la performance .....	17
2.2- RÉSULTATS .....	18
2.2.1- L'ARTP et les conditions de débits en 2020 .....	18
2.2.2- L'ARTP et les conditions d'étiages en 2021 .....	19
2.2.3- Variabilité diurne et intra-saisonnière .....	20
2.2.4- Performance de la prédiction avec les modèles d'apprentissage automatique .....	21
2.3- DISCUSSION ET CONCLUSION .....	23
2.3.1- Variabilité temporelle de l'ARTP .....	23
2.3.2- Performance de la modélisation statistique .....	24
<b>3- ARTICLE VERSION ANGLAISE INTÉGRALE</b> .....	<b>27</b>
3.1- ABSTRACT .....	28
3.2- INTRODUCTION.....	29
3.3- STUDY SITE.....	32
3.4- MATERIAL AND METHODS .....	34
3.4.1- Estimation of hourly PTRA.....	34
3.4.2- Machine learning modeling.....	36
3.4.3- Evaluation procedures and performance metrics .....	40
3.5- RESULTS .....	41
3.5.1- PTRA & mainstem discharge in 2020 .....	41

3.5.2-	PTRA & low flow conditions in 2021 .....	44
3.5.3-	Diel and intra-seasonal variability .....	46
3.5.4-	Machine learning performance .....	49
3.6-	DISCUSSION AND CONCLUSION .....	51
3.6.1-	PTRA characteristics and temporal variability .....	51
3.6.2-	Machine learning performance .....	53
	<b>Acknowledgments .....</b>	<b>55</b>
	<b>Credit author statement.....</b>	<b>55</b>
	<b>Declaration of Competing Interest .....</b>	<b>55</b>
<b>4-</b>	<b>BIBLIOGRAPHIE.....</b>	<b>56</b>

**LISTE DES TABLEAUX**

**Table 1:** Monitoring period and the number of deployed water temperature sensors for stations 1 and 2 .34  
**Table 2:** Hourly mainstream water temperature (°C) of the Ste-Marguerite River during summertime’  
warmest months (July & August 2020-2021)..... 44  
**Table 3:** Best models’ performance evaluated using the rRMSE, RMSE, Bias, rBias and NASH. .... 50



## LISTE DES FIGURES

<b>Figure 1:</b> (a) Location of river mainstems and tributaries, sites, and hydrometeorological stations, and (b) example of deployed water temperature sensors (in yellow) arrays at Ste-Marguerite northeast (Station 2). .....	33
<b>Figure 2:</b> The time series of potential predictors and PTRA at Station 1 (top left & top right) and Station 2 (bottom left and bottom right) (2020-2021). .....	42
<b>Figure 3:</b> Interpolated PTRA and corresponding time series of potential predictors at Station 1 during July 11 <sup>th</sup> & 12 <sup>th</sup> 2020 at 11 p.m.(Figures (3a), (3b) & (3e)) and July 30 <sup>th</sup> at noon & 11 p.m. (Figures (3c), (3d) & (3f)). The blue arrow indicates the mainstem flow direction. ....	43
<b>Figure 4:</b> Interpolated PTRA at station 2 on the 16 <sup>th</sup> July at 9 a.m. (4a) & on the 17 <sup>th</sup> July at 9 a.m. (4b), and the corresponding time series of potential predictors (4c), in green the 24 hour period between (4a) and (4b). .....	45
<b>Figure 5:</b> Coefficient variation for the warmest months (2020-2021) based on three different day periods. S1 and S2 refer to stations 1 & 2, respectively. ....	46
<b>Figure 6:</b> Contour (2020-2021) and perspective plot (2021) using GAM tensor product smoother for stations 1 & 2 (S1 & S2): PTRA~Te (JD, Hourly). PTRA in m <sup>2</sup> , Julian days (JD) & the hour components. ....	48
<b>Figure 7:</b> PTRA (m <sup>2</sup> ) and factor-smooth interaction using the summer' warmest months (2021) fitted with GAM cyclic cubic regression splines for stations 1 & 2 (S1 & S2): PTRA ~ Month + s(Hourly, by=Month).....	49
<b>Figure 8:</b> Evaluation procedures comparison based on rRMSE: K-fold vs. LooCV vs. Split Sample, for stations 1 & 2 (S1 & S2) .....	50
<b>Figure 9:</b> Observed vs. Predicted PTRA (m <sup>2</sup> ). Red and blue dots represent hourly PTRA estimates for summer 2020 and 2021, respectively.....	53

# **1- INTRODUCTION**

## 1.1- MISE EN CONTEXTE

À l'échelle globale, le changement climatique devrait modifier considérablement la distribution et les processus biophysiques des organismes aquatiques dans les écosystèmes lotiques (Dugdale et al., 2018; Isaak & Rieman, 2013). Une augmentation de la température de l'eau est prévue dans la plupart des scénarios de changement climatique pour de nombreuses rivières de l'est du Canada (Dugdale et al., 2018; Morrill et al., 2005; van Vliet et al., 2013). Christensen & Lettenmaier. (2007) ont montré que la température moyenne annuelle entre 1980-1999 et 2080-2099 augmentera de 3,6 °C dans l'est de l'Amérique du Nord (25-50° N et 50-85°O) selon certains scénarios. Une température élevée de l'eau peut avoir des effets néfastes sur les ressources halieutiques en limitant l'habitat des espèces poïkilothermes telles que les salmonidés (Caissie, 2006; Caissie et al., 2004; Lund et al., 2002). Durant les périodes de haute température, ces espèces réduisent leur température corporelle en se déplaçant vers des zones d'eau froide, appelées refuges thermiques (Armstrong et al., 2016; Breau et al., 2011; Dugdale et al., 2015; Torgersen et al., 1999). Il existe un consensus selon lequel les températures entre 20 °C et 23 °C créent un stress thermique chez le saumon Atlantique adulte (*Salmo salar*) (Breau & Caissie, 2013; Shepard, 1995). Pour les saumons juvéniles, la plage de croissance optimale est entre 15 °C et 19 °C, tandis que les seuils limites d'alimentation varient entre 21 °C et 22 °C et les conditions létales surviennent entre 25 °C et 28 °C (Elliott, 1991; Guillemette et al., 2011) .

La disparité thermique dans une rivière est largement reconnue comme un modérateur des effets du changement climatique qui rendent certaines espèces ectothermes, dont les plages de tolérance sont restreintes, vulnérables aux variations de température (Corey et al., 2020; Isaak et al., 2015). Durant la période estivale, les tributaires de taille moyenne (ordre de Strahler : 3-4) peuvent créer des panaches d'eau froide à la confluence (Dugdale et al., 2013; Torgersen et al., 2012). Ces intrusions latérales d'eau froide provenant d'affluents pérennes abaissent souvent la température du cours d'eau principal sur une partie du chenal (généralement près des berges) et ont un débit plus important que celui de la résurgence hyporhéique ou des sources d'eau souterraines (Gendron, 2013; Poole & Berman, 2001; Sutton et al., 2007). Cependant, toute confluence de tributaire ne représente pas forcément un refuge thermique, car cela dépend principalement des fluctuations saisonnières et de la variation spatiale et temporelle de la température (Dugdale et al., 2013).

Ces panaches d'eau froide permettent aux individus de survivre au réchauffement de l'eau et sont souvent une priorité pour la conservation (Davis et al., 2013). Sullivan et al. (2021) ont élaboré une typologie écohydrologique pour ces habitats qui permet de clarifier les différentes nuances et d'unifier le vocabulaire utilisé dans la littérature :

- (i) Un panache d'eau froide, désigne une zone où la température de l'eau est plus froide que la température ambiante du cours d'eau immédiatement en amont ( $\Delta T = 2-10^{\circ}\text{C}$ ) (par ex, Ebersole et al., 2001; Kurylyk et al., 2015).
- (ii) Un refuge thermique est un panache d'eau froide utilisé par les espèces poïkilothermes.
- (iii) Un refuge physiologique désigne une parcelle d'eau dont la température est inférieure à un seuil biologiquement critique (par ex, Breau et al., 2007).

Les poissons poïkilothermes peuvent souvent percevoir des changements de température de moins de  $0.5^{\circ}\text{C}$  (Murray, 1971). La truite arc-en-ciel (*Oncorhynchus mykiss*), en particulier, réagirait à des variations aussi faibles que  $0.1^{\circ}\text{C}$  (Bardach & Bjorklund, 1957).

## **1.2- VARIABLE HYDROMÉTÉOROLOGIQUES ET OUTILS DE MODELISATION**

La température des cours d'eau dépend principalement des conditions météorologiques (température de l'air, les apports en précipitations, rayonnement solaire, la vitesse du vent, etc.) et des caractéristiques physiques du bassin (localisation - longitude, latitude et altitude - utilisation des terres, la couverture végétale, etc.). L'hydrologie de surface spécifique à chaque tronçon de rivière influence fortement la température de l'eau (Caissie, 2006; Gardner et al., 2003) et dépend principalement du régime d'écoulement. Le fait de passer d'un régime nival, i.e., dominé par la fonte des neiges au printemps à un régime pluvial dominé par les précipitations durant l'été, contribue à l'augmentation de la température de l'eau dans la rivière (Barnett et al., 2005). Nuhfer et al. (2017) ont observé que les réductions du débit n'avaient pas d'effet significatif sur la densité d'omble de fontaine (*salvelinus fontinalis*), mais que la croissance des poissons du printemps à l'automne diminuait de manière significative en cas de réduction du débit de 75 % ou plus. D'après Webb et al. (2003), la variable hydrologique la plus importante reste le débit et intègre à son tour d'autres variables hydrologiques, tels le degré de turbulence ou les dimensions des sections de l'écoulement qui influencent à leur tour la température des cours d'eau (Caissie, 2006; Moatar & Gailhard, 2006).

D'un point de vue physique, un grand nombre d'études déterministes ont été menées sur la dynamique des fluides des confluences et du mélange dans les canaux ouverts, cependant la plupart des études ont porté soit sur l'examen de la configuration de la turbulence à la confluence des grands fleuves (Biron et al., 2004; Rhoads & Sukhodolov, 2001) ou ont été motivées par la nécessité de quantifier le mélange de masses d'eau ayant des températures différentes en aval de l'exutoire d'une industrie qui utilise l'eau de rivière comme agent de refroidissement (Zavarsky & Duester, 2020). Peu d'études ont tenté spécifiquement de décrire la dynamique des refuges thermiques et leur modélisation hydrodynamique, ce qui permettrait de répondre aux questions de disponibilité des refuges thermiques, d'atténuation de la température et de survie estivale des poissons dans un cadre naturel de rivières à saumons. La nécessité de caractériser et de quantifier l'évolution temporelle possible de ces habitats, a suscité un grand intérêt pour les techniques non déterministes (stochastiques) qui ont fourni des résultats prometteurs concernant l'utilisation ou la distribution des espèces de poissons d'eau froide (Frechette et al., 2018; Jeong et al., 2013; T. Wang et al., 2020; Wilbur et al., 2020). Par exemple, Jeong et al. (2013) ont modélisé la température de l'eau en un seul point du refuge thermique en utilisant des modèles stochastiques de température de l'eau à la rivière Ouelle (Canada) dans un contexte de changements climatiques. Plus récemment, Saadi et al., (2021) ont développé un modèle statistique à l'échelle des confluences d'affluents à la rivière Ste-Marguerite (Canada), qui utilise un nombre limité de variables hydrométéorologiques afin d'estimer la moyenne journalière de l'aire des refuges thermiques potentiels (ARTP). Ils ont également suggéré que ces panaches d'eau froide présentaient une importante variabilité diurne dans certains cas.

### **1.3- APPROCHE NON-PARAMETRIQUES**

Des études axées sur la température de l'eau ont montré que les modèles stochastiques non paramétriques capturent mieux les relations non linéaires entre les variables et offrent de meilleures performances comparativement à de nombreux modèles paramétriques (Adamowski & Labatiuk, 1987; Benyahya et al., 2007; Caissie et al., 2001; Chenard & Caissie, 2008; Daigle et al., 2010; St-Hilaire et al., 2012). La régression non paramétrique et les réseaux de neurones artificiels (ANN) sont les deux structures principales d'apprentissage automatique utilisées à cette fin. Nous n'avons pas inclus de modèles ANN dans cette étude, mais avons plutôt comparé le pouvoir prédictif de divers algorithmes de régressions non paramétriques. Bien que les ANN aient un nombre fixe de

paramètres, ils sont généralement considérés comme des modèles non paramétriques en raison du grand nombre de paramètres "cachés" (Lee et al., 2017).

Dans le prolongement des travaux de Saadi et al. (2021), la résolution temporelle a été augmentée dans la présente étude. Nous avons estimé l'ARTP au pas de temps horaire en utilisant la température maximale horaire des cours d'eau et une méthode d'interpolation spatiale nommée la Pondération Inverse à la Distance (PID). Nous avons également évalué les performances de prévision de quatre modèles de régression non paramétrique pour la prévision de l'ARTP maximale horaire. Les techniques d'apprentissage automatique proposées comprennent le modèle additif généralisé (GAM; Hastie & Tibshirani, 1987), le modèle de régression multivariée par splines adaptatifs (MARS; Friedman, 1991) la régression par machine à vecteurs de support (SVM; Vapnik, 1998) et la régression par forêts aléatoires (RF; Breiman, 2001). Ces modèles peuvent être classés en fonction de leurs algorithmes d'apprentissage, où GAM et MARS sont des extensions d'algorithmes linéaires, SVM utilise des méthodes à noyaux alors que RF se base sur l'apprentissage par des arbres de décision (Zhang et al., 2020).

Des études antérieures ont montré que le modèle GAM permet de modéliser de façon satisfaisante la température de l'eau en fonction de la température moyenne quotidienne de l'air et du débit de la rivière Ste-Marguerite (Laanaya et al., 2017). Frechette et al. (2019) ont appliqué le même modèle pour évaluer la relation entre l'abondance des poissons, la température de l'eau et le débit. D'autre part, l'approche MARS a amélioré le pouvoir prédictif de l'analyse fréquentielle régionale des crues (RFA), comparée au modèle GAM (Msilini et al., 2020). Récemment, le modèle RF a été introduit pour la première fois dans la RFA et a été combiné avec l'analyse de corrélation canonique (CCA) pour obtenir une qualité prédictive élevée (Desai & Ouarda, 2021). Saadi et al. (2021) ont montré que les modèles GAM et MARS possèdent un fort potentiel pour simuler adéquatement l'ARTP moyenne journalière. Les méthodes à noyaux telles que le modèle SVM ont montré leur capacité à modéliser un large éventail de processus hydrologiques (Deka, 2014). Weierbach et al. (2021) ont comparé la régression linéaire multiple (MLR), SVM et RF afin de prédire la température mensuelle des cours d'eau. Ils ont souligné la surestimation des valeurs extrêmes par le modèle RF et la bonne performance globale du modèle SVM. Allahbakhshian-Farsani et al. (2020) ont appliqué SVM, MARS et le modèle arbre de régression boosté (BRT) dans l'analyse fréquentielle régionale des crues. Les résultats ont montré que les meilleures performances ont été obtenues avec

le modèle SVM en utilisant la fonction de base radiale (RBF) comme noyau. Quan et al. (2020) ont utilisé un modèle optimisé du SVM en utilisant un modèle de réseaux de neurones connu sous le nom d'algorithme génétique (GA) pour prédire la température de l'eau des grands réservoirs à haute altitude dans l'ouest de la Chine. Leur résultat fournit des indications utiles sur la prédiction de la température de l'eau à différentes profondeurs des réservoirs.

Les modèles de régression basés sur les arbres de décisions se sont avérés compétitifs par rapport aux réseaux de neurones artificiels et ont montré un fort potentiel pour la modélisation de la température de l'eau. Ebersole et al. (2015) ont utilisé RF pour modéliser l'occurrence de parcelles d'eau froide aux confluences des affluents en fonction du bassin-versant et caractéristiques climatiques. Ferchichi et al. (2021) ont utilisé le modèle RF et un réseau de neurones à rétropropagation (BPNN) afin d'étudier l'impact des futurs scénarios de température des eaux côtières sur le risque de croissance microbienne marine. Ils ont obtenu des performances très similaires entre les deux modèles. Feigl et al., (2021) ont prédit la température des cours d'eau en Europe à l'aide de différentes méthodes d'apprentissage automatique et ont constaté que le modèle basé sur les arbres de décisions appelé '*extreme gradient boosting*' (XGBoost) avait des performances comparables à celles du réseau de neurones à propagation avant (FNN).

## 1.4- OBJECTIFS

Les objectifs de cette présente étude sont :

- (i) Estimer l'ARTP maximale horaire à la confluence de deux tributaires de la rivière Ste-Marguerite (Canada) au cours des deux étés 2020 et 2021.
- (ii) Étudier la variabilité diurne de l'ARTP en mettant l'accent sur l'évènement d'étiage de l'été 2021.
- (iii) Valider les modèles GAM et MARS sur un pas de temps horaire; qui ont précédemment été utilisés pour la prévision de l'ARTP moyenne journalière (Saadi et al., 2021).
- (iv) Présenter les algorithmes de régression SVM et RF pour l'estimation de l'ARTP maximale horaire en utilisant des variables hydrométéorologiques relativement faciles à mesurer.

## 1.5- STRUCTURE DU MEMOIRE

L'ensemble des travaux effectués dans le cadre de cette maîtrise est présenté sous forme d'un mémoire par article. Le chapitre deux représente une synthèse de l'article en français, suivie par la version anglaise de l'article (chapitre 3) qui s'intitule « *Machine learning analysis and modeling of hourly potential thermal refuge area: case study of the Ste-Marguerite River* » et a été soumis à la revue à la revue *Science of the Total Environment* (STOTEN). Les références relatives à l'article sont reportées dans la bibliographie présentée à la dernière section de ce mémoire.

L'encadrement, l'orientation de la méthodologie de recherche et la validation des analyses et résultats obtenus durant ces travaux ont été assurés et effectués par le directeur de recherche en collaboration avec le co-directeur de recherche. Durant la réalisation de l'article, l'étudiant avait pour rôle de participer à l'élaboration de la méthodologie de recherche, réaliser les campagnes de terrain, de compiler et extraire les données échantillonnées, et de procéder à la calibration et l'analyse statistiques des résultats obtenus. L'étudiant a aussi écrit l'article avec l'aide des suggestions et révisions des co-auteurs.



## **2- SYNTHÈSE DE L'ARTICLE EN FRANÇAIS**

## 2.1- METHODOLOGIE ET SITE D'ÉTUDE

Deux sites potentiels pour la présence de refuges thermiques ont été identifiés sur la rivière Sainte-Marguerite dans la région du Saguenay, entre Chicoutimi et Sacré-Cœur, (Québec, Canada). La rivière se divise en trois branches distinctes, dont deux sont considérées dans cette étude (**Figure 1a**). La branche principale coule parallèlement à la rive nord du fjord du Saguenay et se déverse dans la baie Sainte-Marguerite. Les branches nord-est et nord-ouest sont plus difficiles d'accès. La branche nord-est se joint au cours d'eau principal à environ 2 km de l'exutoire. Dans notre étude de cas, la première station est située sur la branche principale, tandis que la deuxième station est localisée sur la branche nord-est, avec une superficie des bassins versants de 1 097 km<sup>2</sup> et de 980 km<sup>2</sup>, respectivement. Les deux stations ont été choisies en raison de la présence d'affluents permanents et froids. La Station 1 est caractérisée par une forte pente et un écoulement torrentiel, avec une série de petites fosses en cascades (*step pools*), un angle à la confluence entre le tributaire et le cours d'eau principal d'environ 90 degrés et un tributaire plus froid que celui de la Station 2. Cette dernière est caractérisée par un lit de faible pente pour le cours d'eau principal et un angle d'entrée de l'affluent d'environ 60 degrés. Le substrat des deux sites est dominé par du gravier et du galet.

### 2.1.1- Estimation de l'ARTP maximale horaire

Les stations 1 et 2 ont été instrumentées pendant les étés 2020 et 2021 (**Table 1**). Comme le montre la **Figure 1b** (Station 2), un réseau de thermographes (en jaune) a été placé à la confluence de l'affluent et du cours d'eau principal. Plusieurs thermographes (capteurs Pendant Hobo, précision =  $\pm 0.5$  °C) ont été fixés sur des barres d'armature ancrées dans le lit du cours d'eau principal pour former des rangées transversales (transects) le long de la confluence. Ces rangées ont été conçues pour mesurer la température dans les panaches d'eau froide qui ont été observés premièrement à l'aide d'une caméra thermique infrarouge (Saadi et al., 2021; Wang et al., 2020). Des thermographes supplémentaires ont été déployés en amont de chaque confluence pour enregistrer la température de l'eau qui n'était pas influencée par le processus de mélange à la confluence. La distance latérale entre les points de mesure de température (capteurs) variait en fonction de la morphologie du site (la bathymétrie, la taille de l'affluent, l'angle à la confluence). Les thermographes ont enregistré les données à des intervalles de 15 minutes, à partir desquels nous avons calculé les températures maximales horaires de l'eau.

Kurylyk et al., (2015) suggèrent que les tacons du saumon atlantique utilisent des refuges thermiques dont la différence de température avec le cours d'eau en amont est inférieure à 2°C, ce qui est conforme à la définition de Saadi et al. (2021). En conséquence, nous avons défini l'ARTP par la région dont la température de l'eau est inférieure d'au moins un degré Celsius à celle du cours d'eau principal situé directement en amont ( $\Delta T \geq 1$  °C). Ce critère a été choisi non pas sur la base d'observations locales, mais afin de disposer de séries temporelles horaires des ARTP suffisamment longues pour permettre d'appliquer des techniques d'apprentissage automatique et pour que la différence de température minimale qui permet l'observation d'une ARTP ( $\Delta T = 1$ °C) soit supérieure à la précision des thermographes utilisés durant cette étude ( $\pm 0.5$  °C). L'ARTP maximale horaire est estimée à partir d'une interpolation spatiale et la limite du panache a été déterminée par la zone dont la température de l'eau est au moins 1°C inférieur à celle du cours d'eau principal ( $\Delta T \geq 1$  °C).

#### 2.1.1.1- Pondération inverse à la distance (PID)

La pondération inverse à la distance est une méthode déterministe qui utilise une moyenne pondérée des valeurs connues (sites jaugés) pour estimer une valeur inconnue à un endroit donné (sites non jaugés). Ce processus nécessite l'estimation des pondérations  $w_{ij}$  comme une fonction inverse de la distance entre le site d'intérêt et les sites jaugés. L'équation 1 montre la base de l'interpolation PID :

$$Z_j = \frac{\sum_{i=1}^n x_i * w_{ij}}{\sum_{i=1}^n w_{ij}} \quad \text{Équation 1}$$

Où:

$$w_{ij} = \frac{1}{D_{ij}^p}$$

Dans notre cas d'étude,  $Z_j$  est la valeur estimée de la température de l'eau à l'emplacement non jaugé  $j$ ,  $x_i$  est la température de l'eau du site  $i$  voisin,  $w_{ij}$  est le poids attribué aux sites jaugés  $i$ ,  $D_{ij}$  est la distance entre le site jaugé  $i$  et le site non jaugé d'intérêt  $j$ ,  $n$  est le nombre de sites jaugés, et  $p$  est l'exposant qui varie entre 3 et 4 dans cette étude.

L'erreur d'extrapolation des températures de l'eau est évaluée à l'aide de l'approche de validation croisée 'leave-one-out' (LOOCV) pour chaque station. Cette méthode est étroitement liée à la

méthode d'estimation statistique appelée *jack-knife* (Efron, 1982). La racine de l'erreur quadratique moyenne (RMSE) est utilisée pour évaluer l'erreur par validation croisée ( $RMSE_{station\ 1} = 0,3\ ^\circ C$  et  $RMSE_{station\ 2} = 0,27\ ^\circ C$ ). Nous voulions nous assurer que l'erreur globale, qui comprend à la fois la précision des thermographes ( $\pm 0.5\ ^\circ C$ ) et l'erreur de la validation croisée, restait inférieure à  $1\ ^\circ C$ ; le critère utilisé pour établir les limites spatiales minimales du refuge :  $(|\pm 0.5|^\circ C + RMSE_{station\ i} \leq 1^\circ C)$ . Les erreurs moyennes globales des stations 1 et 2 étaient de  $0.8^\circ C$  et  $0.77^\circ C$ , respectivement. Toute la programmation et les calculs concernant la PID ont été effectués à l'aide du logiciel *Python* 3.8, la librairie *Scipy.Spatial* et la fonction *cKDTree* ([docs.scipy.org/doc/scipy/reference/spatial.html](https://docs.scipy.org/doc/scipy/reference/spatial.html)).

#### 2.1.1.2- Les variables potentielles prédictives

Afin de prédire l'ARTP maximale horaire, nous avons évalué quatre modèles d'apprentissage automatique qui ont utilisé jusqu'à cinq variables explicatives potentielles à un pas de temps horaire:

- (i) Le débit du cours principal ( $Q_m$ ),
- (ii) La température de l'air ( $T_a$ ),
- (iii) La température maximale horaire du cours principal directement en amont de la confluence ( $T_m$ ),
- (iv) La température du tributaire associé ( $T_t$ ),
- (v) La différence de température entre  $T_m$  et  $T_t$  ( $T_\Delta$ ).

Dans cette étude, le débit horaire est obtenu à partir de la station hydrométrique #062803 du ministère provincial de l'environnement et de la lutte contre le changement climatique (*MELCC*), située sur la branche nord-est de Ste-Marguerite ( $48^\circ 16' 5'' N$ ,  $69^\circ 54' 33'' W$ ). La température de l'air est mesurée à la station météorologique *CIRSA*, indiquée sur la **Figure 1(a)**. Une méthode de présélection de variables appelée "*l'élimination récursive des variables (Recursive Feature Elimination; RFE)*" est utilisée pour identifier le meilleur sous-ensemble de variables pour chaque modèle.

#### 2.1.1.3- L'élimination récursive des variables

La RFE est un algorithme d'apprentissage automatique de sélection de variables explicatives qui peut s'adapter à n'importe quel modèle. Il produit le meilleur ensemble possible de

variables explicatives qui donne les meilleures performances. L'objectif de la RFE est d'ajuster le modèle d'apprentissage automatique, de classer les variables par importance, d'écarter les variables les moins importantes et de réajuster le modèle. Le processus est répété jusqu'à ce que le critère de performance soit adéquat (Kuhn & Johnson, 2013). Dans cette étude, l'erreur quadratique moyenne a été utilisée comme critère de performance pour sélectionner le meilleur sous-ensemble de variables. Pour évaluer la performance de prédiction des sous-ensembles de variables possibles, la validation croisée avec  $k$  répétitions ( $k$ -fold) est menée avec  $k=10$ . On ajuste ensuite le modèle en utilisant les  $k-1$  ( $10 - 1 = 9$ ), par la suite le modèle est validé en utilisant le sous-ensemble ( $fold$ ) restant. Les scores et les erreurs  $RMSE_j$  sont notés (Équation 3) et le processus est répété jusqu'à ce que chaque sous-ensemble serve au sein de l'ensemble d'entraînement. La moyenne des erreurs est considérée comme la métrique de performance du modèle ( $RMSE_{cv}$ ) et s'écrit comme suit :

$$RMSE_{cv} = \frac{1}{k} \sum_{j=1}^k RMSE_j \quad \text{Équation 2}$$

Pour chaque  $fold$   $j$ , la RMSE est calculée comme suit :

$$RMSE_j = \sqrt{\frac{1}{N} \sum_{i=1}^N (o_i - s_i)^2} \quad \text{Équation 3}$$

Où  $k$  est le nombre de *sous-ensembles* ( $k=10$ ),  $N$  est la taille de l'échantillon du *sous-ensemble*  $j$ ,  $s_i$  sont les données simulées et  $o_i$  sont les données observées.

## 2.1.2- Algorithmes d'apprentissage automatique et modélisation horaire de l'ARTP

### 2.1.2.1- Le modèle additif généralisé (GAM) et la variabilité temporelle de l'ARTP.

Le GAM est un modèle non paramétrique et une version étendue du modèle linéaire généralisé GLM (McCullagh & Nelder, 1989). Ce modèle remplace le prédicteur linéaire dans le GLM par un prédicteur additif. Les GAM modélisent les données continues avec une fonction de lissage non linéaire (Splines de régression) qui peut prendre diverses formes (Dominici et al., 2002). De plus, une fonction de liaison adaptée à la distribution de la variable réponse doit également être définie. Dans le processus de prévision, nous avons utilisé le GAM pénalisé, où les fonctions de lissage utilisées, sont des splines cubiques afin d'éviter le surajustement (*Penalized splines*, Wood, 2006). Un des grands avantages d'utiliser le modèle GAM est que la forme optimale de la non linéarité, autrement appelé le degré de lissage de la fonction, est contrôlé en utilisant une

régression pénalisée qui est déterminée automatiquement à l'aide d'une méthode de validation croisée généralisée (GCV). Une famille de distribution gaussienne est utilisée pour la variable réponse : l'ARTP maximale horaire (Marra & Wood, 2011). Le modèle GAM s'exprime comme suit :

$$g(E(Y)) = \beta_0 + s_1(x_1) + s_2(x_2) \dots + s_p(x_p) + \varepsilon \quad \text{Équation 4}$$

La fonction de liaison  $g$  est une fonction paramétrique qui lie la moyenne de la variable dépendante à un ensemble de variables explicatives ;  $E(y)$  est l'espérance de la variable de réponse prédite.  $\beta_0$  est l'ordonnée à l'origine,  $s_i$  est la fonction lisse de la  $i^{\text{ème}}$  variable explicative,  $x_i$  sont les variables indépendantes utilisées pour estimer l'aire potentiel des refuges thermiques horaire,  $\varepsilon$  est un terme d'erreur qui est normalement distribué avec la variance  $\sigma_\varepsilon$ . GAM s'est révélé être un outil puissant qui peut fournir une estimation de l'ARTP moyenne journalière à l'aide d'un nombre relativement faible de variables explicatives qui sont généralement faciles à mesurer (Saadi et al., 2021).

Afin d'analyser la variation spatiotemporelle des aires de refuge, nous avons organisé l'heure de la journée en trois périodes distinctes, et avons calculé le coefficient de variation ( $CV = \frac{\text{Standard deviation}}{\text{Mean}}$ ) de l'ARTP estimée pour chacune de ces périodes tout au long des mois de juillet et août (2021/2021) (Heure Normale de l'Est, HNE) :

- (i) Tôt (00:00 - 07:59),
- (ii) Mi-journée (08:00 - 15:59),
- (iii) Tard (16:00-23:59).

D'autre part, nous avons également étudié la variabilité diurne à travers les composantes horaires et journalières du calendrier Julien en utilisant les outils d'interactions de GAM. Ainsi, les données collectées en continu durant deux étés consécutifs ont été analysées pour mieux comprendre la variabilité intra-saisonnière de l'ARTP. Une base de lissage par produit tensoriel (Te) a été implémentée pour simuler l'effet d'interaction de l'heure de la journée et le jour Julien sur l'ARTP en utilisant l'équation 5. Cet outil apporte une très grande flexibilité au GAM, puisqu'il permet d'expliquer une variable réponse par des fonctions impliquant plusieurs variables explicatives. L'idée générale est que l'on va utiliser des bases de lissage marginales et les combiner de telle sorte que l'on va construire une fonction multivariée (Boucher et al., 2017; Mahardja et al., 2021).

Pour modéliser le cycle diurne de l'ARTP durant les mois les plus chaud de l'année, nous avons utilisé la variable 'mois' comme un facteur à deux niveaux (juillet/août) et on a lissé la variable réponse (ARTP) selon les différentes heures de la journée (Équation 6). Dans le modèle GAM, lorsque l'on modélise des données cycliques, on souhaite généralement que le prédicteur soit identique aux deux bouts des phases. Pour y parvenir, nous devons modifier la fonction de base. Dans notre cas d'étude on va utiliser la fonction de lissage cubique cyclique pour modéliser les effets de l'heure de la journée et la fonction cubique simple pour les effets du jour Julien :

$$Area \sim Te(j, h, bs = c('cr', 'cc')) \quad \text{Équation 5}$$

$$Area \sim m + s(h, by = m, bs = 'cc') \quad \text{Équation 6}$$

Où *Area* représente la variable réponse (ARTP) au pas de temps horaire, le jour julien est représenté par *j*, l'heure de la journée par *h*, et le mois correspondant par '*m*'. '*cr*' représente la fonction de lissage cubique simple et '*cc*' la fonction de lissage cubique cyclique.

#### 2.1.2.2- La régression multivariées par splines adaptatifs (MARS)

Le modèle MARS peut être considéré comme un cas spécifique flexible du modèle GAM, exprimé comme une combinaison linéaire de fonctions de base et de leurs interactions (Msilini et al., 2020). Il s'agit d'une approche de régression non paramétrique permettant de traiter des données à haute dimension. Le modèle MARS construit une suite de régressions linéaires en subdivisant l'ensemble des variables explicatives en plusieurs régions (Conoscenti et al., 2015). Les valeurs de rupture entre les régions sont appelées "nœuds" et les intervalles sont conçus comme des fonctions de bases linéaires (Conoscenti et al., 2016). Chaque nœud marque la fin d'une région de données et le début d'une autre. Les fonctions de bases linéaires sont générées par une recherche par étapes et la comparaison des sous-modèles est établie par validation croisée généralisée (GCV) (Roy et al., 2018). MARS s'écrit comme suit :

$$Y = \beta_0 + \sum_{k=1}^N \beta_i f_i(x) \quad \text{Équation 7}$$

Où  $\beta_0$  est l'intercept, et  $\beta_i$  sont les coefficients de régression des fonctions de base  $f_i(x)$ .

Dans une étude récente, le modèle MARS a montré une meilleur précision et a utilisé moins de variables explicatives que GAM dans la modélisation de l'ARTP moyenne journalière (Saadi et

al., 2021). Par conséquent, les deux modèles ont été sélectionnés pour valider la modélisation de l'ARTP au pas de temps horaire.

### 2.1.2.3- La régression par machine à vecteur de support (SVM)

Le SVM a été introduit par Vapnik, (1998) comme un outil robuste d'apprentissage automatique. Le but est de faire correspondre les ensembles de données originaux de l'espace d'entrée à un espace de caractéristiques à haute dimension ou même à dimension infinie, de sorte que les problèmes de classification ou de régression deviennent plus faciles dans cet espace (Deka, 2014). Cette approche a pour principal avantage de minimiser la complexité du modèle et l'erreur de prédiction, en utilisant des fonctions noyaux. Le succès du SVM dépend principalement de la détermination de la fonction noyau appropriée et de ses hyper-paramètres associés. Dans une certaine mesure, on peut considérer que le choix de la fonction noyau est équivalent au choix de la structure d'un réseau de neurones artificiels. Dans cette étude, les variables explicatives ont un comportement non linéaire. À cet égard, les noyaux non linéaires donnent au SVM la capacité de modéliser de hyperplans complexes avec des frontières non linéaires. En régression, SVM a des avantages dans les espace de haute dimensionnalité car l'optimisation de la régression des vecteurs de support ne dépend pas de la dimensionnalité de l'espace d'entrée (Drucker et al., 1996). SVM utilise une fonction de perte  $L_\epsilon(v, g(u))$  qui décrit la déviation de la fonction estimée par rapport à la fonction originale. Dans le présent contexte, on utilise la fonction de perte insensible standard de Vapnik - *epsilon* ( $\epsilon$ ) qui est définie comme suit :

$$L_\epsilon(v, g(u)) = \begin{cases} 0 & \text{for } |v - g(u)| \leq \epsilon \\ |v - g(u)| - \epsilon & \text{Sinon} \end{cases} \quad \text{Équation 8}$$

En utilisant la fonction de perte insensible ( $\epsilon$ ), on peut trouver  $g(u)$  qui peut mieux estimer le vecteur de sortie original  $v$ .

Nous avons utilisé le modèle de régression  $\epsilon$ -SVM avec un noyau non linéaire dénommé fonction de base radiale (Radial Basis Function - RBA), comme le montre l'équation 9. Il s'agit de la forme la plus généralisée de noyau et l'une des plus utilisées en raison de sa similitude avec la distribution gaussienne.

$$K(u_i, u_j) = \exp\left(-\gamma(u_i - u_j)^2\right), \gamma \geq 0 \quad \text{Équation 9}$$

La fonction  $K$  calcule la proximité entre deux points  $u_i$  et  $u_j$ , et " $\gamma$ " est le paramètre qui contrôle la portée radiale. Enfin, la fonction de décision du SVM non linéaire est exprimée comme suit :



$$g(u) = \sum_{i=1}^l (\alpha_i - \bar{\alpha}_i) K(u_i, u) + b \quad \text{Équation 10}$$

$\alpha_i$  et  $\bar{\alpha}_i$  représentent un couple de variables pour les contraintes correspondantes et  $b$  le biais.

#### 2.1.2.4- La régression par forêt aléatoire (Random Forest (RF))

Le RF a été introduit par Breiman, (2001) comme un modèle d'apprentissage automatique basé sur l'idée de mise en sac ou '*bagging*' (agrégation bootstrap) (Breiman, 1996). Les variables explicatives du *bagging* font la moyenne des estimations de plusieurs modèles, où chaque modèle est entraîné sur un échantillon '*bootstrap*' au lieu de la totalité de l'échantillon observé. Le caractère aléatoire introduit par le *bootstrap* augmente la capacité du modèle à généraliser et à produire des résultats de prédiction stables. Les modèles RF sont des prédicteurs qui utilisent les arbres de classification et de régression CART comme apprenant de base. Les CART appliquent récursivement des divisions binaires aux données afin de minimiser l'entropie (une mesure de la dissimilitude) dans les nœuds de l'arbre. Cette opération est effectuée jusqu'à ce que chaque nœud atteigne une taille minimale, ou qu'une profondeur maximale de l'arbre définie au préalable soit atteinte. Breiman, (2001) a montré que l'ajout d'un caractère aléatoire supplémentaire à la méthode de mise en sac améliore la précision des prédictions. Dans le modèle RF, on y parvient en sélectionnant uniquement un sous-ensemble aléatoire de variables disponibles pour le fractionnement à chaque nœud. L'estimation de la variable dépendante est donnée par :

$$Y = \frac{1}{M} \sum_{m=1}^M f_m(X) \quad \text{Équation 11}$$

Où  $f_m$  Désigne un seul arbre de régression (RT),  $M$  le nombre de RT utilisés,  $X$  la matrice des variables indépendantes et  $Y$  le vecteur de la variable dépendante.

Tous les modèles d'apprentissage automatique ont été mis en application sur R (R Core Team, 2020). Le développement des modèles s'est appuyé sur la librairie "*Caret*" (Kuhn et al., 2020), pour l'élimination récursive des variables (RFE) et pour l'optimisation des hyper-paramètres à l'aide des méthodes de grilles de recherche et de validation-croisée. L'implémentation de GAM, MARS, SVM et RF a été complétée avec les librairies "*mgcv*" (Wood, 2006), "*earth*" (Milborrow, 2018), "*e1071*" (Meyer & Wien, 2001) et "*random forest*" (Breiman, 2001), respectivement. Les visualisations ont été réalisées à l'aide la librairie "*ggplot2*" (Wickham & Grolemund, 2016).

### 2.1.3- Procédure d'évaluation et mesures de la performance

Nous avons comparé trois différentes procédures de validation :

- (i) *Train-Test Split* qui consiste à décomposer de manière aléatoire un ensemble de données. Une partie servira à l'entraînement du modèle de d'apprentissage automatique, l'autre partie permettra de le tester pour la validation. En général, on réserve 70% à 80% des données pour l'entraînement. Les 20 à 30% restants seront exploités pour la validation-croisée. Cette technique est efficace, sauf si la taille d'échantillon est relativement courte. Il peut alors manquer certaines informations sur les données qui n'ont pas été utilisées pour l'entraînement, et les résultats peuvent donc être hautement biaisés. En revanche, si l'ensemble de données est vaste et que la distribution est égale entre les deux échantillons, cette approche convient tout à fait.
- (ii) La validation croisée *k-fold* : Par rapport aux autres approches de validation, elle résulte généralement sur un modèle moins biaisé. Elle permet d'assurer que toutes les observations de l'ensemble de données original apparaissent dans l'ensemble d'entraînement et dans l'ensemble de test. En cas de séries courtes, il s'agit donc de l'une des meilleures approches. La procédure a un paramètre unique appelé ' *k* ' faisant référence au nombre de sous-groupes dans lesquels l'échantillon sera divisé. On commence tout d'abord par séparer l'ensemble de données de manière aléatoire en *k* différentes *folds*. La valeur de *k* ne doit être ni trop basse ni trop haute, et on choisit généralement une valeur comprise entre 5 et 10 en fonction de l'envergure de la base de données.
- (iii) La validation croisée *leave one-out* (LOOCV) : il s'agit d'un cas particulier de la validation croisée *k-fold* où  $k=N$ , avec *N* l'ensemble de la base de données. C'est-à-dire qu'à chaque itération, l'apprentissage se fait sur  $N-1$  observations et la validation sur l'unique observation restante.

Les métriques suivantes sont utilisées pour évaluer la qualité de nos modèles de régression :

1- Le coefficient de Nash (NASH) 
$$= 1 - \frac{\sum_{i=1}^n (o_i - s_i)^2}{\sum_{i=1}^n (o_i - o_m)^2}$$
 Équation 12

2- La racine carrée de l'erreur quadratique (Root Mean Square Error (RMSE)) 
$$= \sqrt{\frac{1}{n} \sum_{i=1}^n (o_i - s_i)^2}$$
 Équation 13

3- La racine carrée de l'erreur quadratique relative (Relative Root Mean Square Error (rRMSE)) 
$$= \sqrt{\frac{1}{n} \sum_{i=1}^n \frac{(o_i - s_i)^2}{o_i^2}}$$
 Équation 14

4- Le biais 
$$= \frac{1}{n} \sum_{i=1}^n (o_i - s_i)$$
 Équation 15

5- Le biais relatif (rBias) 
$$= \frac{1}{n} \sum_{i=1}^n \frac{(o_i - s_i)}{(o_i)}$$
 Équation 16

Où  $o_i$  est la  $i^{\text{ième}}$  valeur observée,  $s_i$  est la  $i^{\text{ième}}$  valeur simulée à l'aide du modèle,  $o_m$  est la valeur moyenne observée et  $n$  est la taille de l'échantillon du site.

## 2.2- RÉSULTATS

### 2.2.1- L'ARTP et les conditions de débits en 2020

La **Figure 2** montre que la superficie des refuges thermiques potentiels (en rouge) aux stations 1 et 2 était nulle ( $0 \text{ m}^2$ ) pour un nombre total de 20 et 10 jours, respectivement. Ces absences d'ARTP (c.-à-d.  $= 0 \text{ m}^2$ ) ont été associées à une augmentation importante du débit du cours d'eau principal ( $Q_m$  jusqu'à  $128 \text{ m}^3/\text{s}$ ) (fin juillet, début août et mi-septembre). Pour un débit inférieur à  $50 \text{ m}^3/\text{s}$  (juin, juillet et mi/fin août), les aires estimées à la Station 2 en 2020 sont relativement constantes ( $\approx 400 \text{ m}^2$ ) pendant une longue période, tandis qu'à la Station 1, l'ARTP a exprimé une plus grande variabilité, avec une chute significative observée pour des valeurs de débit aussi faibles que  $25 \text{ m}^3/\text{s}$  (le 11 juillet et vers la mi-août). La **Figure 3** illustre la délimitation spatiale de l'ARTP (frontière en rouge) et montre l'effet de l'augmentation du débit sur la variabilité temporelle du panache à la station 1. Entre la période du 11 au 12 juillet 2020, la température maximale du cours d'eau principal a atteint  $26 \text{ °C}$  ( $T_m \geq 22 \text{ °C}$  pendant plus de 10 jours consécutifs) et le débit  $Q_m$  est passé de  $12,5$  à  $25 \text{ m}^3/\text{s}$ . Par conséquent, l'aire du refuge thermique a connu sa plus grande chute ( $210.6 \text{ m}^2$ ) en 24 heures tout au long de la saison estivale (de  $260.7 \text{ m}^2$  à  $50.1 \text{ m}^2$ ). Par la suite, les valeurs de débit ont continué à augmenter ( $Q_m \geq 40 \text{ m}^3/\text{s}$ ) et le 30 juillet, la superficie a atteint  $29 \text{ m}^2$ . Les observations de l'été 2020 ont révélé que lorsque le débit du cours

d'eau principal dépassait 45 m<sup>3</sup>/s et 60 m<sup>3</sup>/s respectivement à la station 1 et 2, l'ARTP était réduite à zéro, et ceci, peu importe la valeur des autres variables explicatives (e.g. la station 1 entre le 3 et le 8 août & début/mi-septembre).

### 2.2.2- L'ARTP et les conditions d'étiages en 2021.

Le débit horaire moyen enregistré à la rivière Ste-Marguerite en août 2021 (7.92 m<sup>3</sup>/s) est quatre fois plus faible que celui observé durant la même période en 2020 (35.01 m<sup>3</sup>/s). Les températures estivales de l'eau dans cette rivière sont généralement plus élevées en juillet qu'en août, ce qui n'était pas le cas en 2021 (**Table 2**). Par conséquent, la variabilité diurne de l'ARTP au cours des deux étés 2020 et 2021 était très contrastée. Durant le mois d'août 2021, la température du cours d'eau principal dans les deux stations a atteint des niveaux dangereux pour les espèces de poisson d'eau froide ( $T_m \approx 27$  °C pendant 10 jours consécutifs; du 15 au 25 Août). Durant ce mois, le débit a atteint des valeurs historiquement basses ( $Q_m \leq 10$  m<sup>3</sup>/s), ce qui soulève des préoccupations de conservation concernant du saumon atlantique (*Salmo Salar*) dans la rivière Ste-Marguerite. Les quantiles et les périodes de retour ont été estimés pour la période estivale à partir des valeurs journalières de débit entre juin et septembre (période estivale) de 1998 à 2021 en utilisant une moyenne mobile de sept jours. Pour un débit d'étiage sur 7 jours consécutifs, une période de retour de 43 ans est estimée pour l'été 2021 ( $Q_{43,7}$ ) en utilisant la distribution de Weibull ajustée avec la méthode du maximum de vraisemblance, ce qui démontre que cet étiage est relativement extrême pour cette rivière.

Au cours de l'été 2021, l'aire du refuge thermique à la Station 1 a été réduite à 0 m<sup>2</sup> seulement pendant quelques heures durant la journée du 28 juin, tandis qu'à la Station 2, seules des valeurs élevées (entre 300 m<sup>2</sup> et 700 m<sup>2</sup>) ont été enregistrées pendant toute la période d'échantillonnage. Ceci est dû aux faibles débits du cours d'eau principal et à la décharge d'eau froide relativement constante de ses affluents. Le débit maximum observé durant toute la période estivale en 2021 est signalée le 17 juillet (32 m<sup>3</sup>/s). Ce jour-là, l'ARTP à la Station 2 a chuté de 681 m<sup>2</sup> à 385 m<sup>2</sup> (**Figure 4**). Par la suite, la superficie des refuges à la même station a augmenté pour se stabiliser dès le début du mois d'août ( $\approx 550$  m<sup>2</sup>). Le débit du cours d'eau principal a continué de diminuer de manière drastique pour atteindre des valeurs très faibles. Ces superficies constantes observées à partir du début du mois d'août sont dues au manque de précipitation (faible niveau d'eau et un lit de rivière peu profond à la Station 2) et à la différence croissante de la température de l'eau entre

le tributaire et le chenal principal ( $T_{\Delta}$  jusqu'à 17 °C à la Station 1). Les effets de l'étiage historique en 2021 sont évidents sur l'ARTP, où certains des thermographes situés sur le bord de la rive se sont avérés être hors de l'eau dès le début du mois d'août et ont donc enregistré la température de l'air au lieu de la température de l'eau. Par conséquent, lorsque ces capteurs ont été retirés de la procédure d'interpolation (PID), les aires de refuge à la Station 2 sont passées de 700 à 570 m<sup>2</sup> (le 3 août 2021), suivie trois jours plus tard par une baisse de 70 m<sup>2</sup> à la Station 1 (de 190 à 120 m<sup>2</sup>).

### **2.2.3- Variabilité diurne et intra-saisonnière**

#### *2.2.3.1- L'effet de l'étiage sur le coefficient de variation*

Selon la **Figure 5**, le coefficient de variation (CV) des ARTP a montré une plus grande variabilité durant le mois d'août que durant le mois de juillet 2020, avec des valeurs allant jusqu'à 75% à la Station 1 et 55% à la Station 2 pour le mois d'août. La rivière Ste-Marguerite a enregistré des débits inhabituellement bas et une température du cours d'eau principal relativement élevée pendant le mois d'août 2021. L'étiage historique survenu durant ce mois a eu un impact significatif sur les valeurs du coefficient de variation de l'ARTP, qui a chuté à environ 13 % et 10 % pour les stations 1 et 2, respectivement. La variation du CV en juillet n'était pas aussi importante que celle enregistrée en août durant les deux étés. Sur la même figure on constate que la Station 1 présentait les valeurs de CV les plus élevées, et que la plus grande variabilité se produisait dans la dernière partie de la journée (Tard : 16:00-23:59). Au cours de l'été 2020, la fin d'après-midi était la période qui représentait la plus grande variabilité à la Station 2, mais cette tendance a changé durant l'été 2021 pour se déplacer vers la mi-journée en juillet et tôt le matin en août.

#### *2.2.3.2- GAM et la variabilité temporelle de l'ARTP*

Tout d'abord, nous avons estimé les interactions entre l'ARTP, les jours juliens et les heures en utilisant une base de lissage par produit tensoriel du modèle GAM (Équation 5). La région bleu foncé dans la **Figure 6** fait référence aux valeurs relativement faibles des aires de refuge, tandis que le rouge foncé identifie les valeurs élevées. Les interactions verticales se sont avérées relativement constantes en 2020, ce qui signifie qu'il n'y avait pas de distinctions dans le moment d'occurrence des aires de refuge maximales et minimales durant la journée en raison des grandes valeurs de débit enregistrées durant cet été. Néanmoins, l'ARTP le long de la période estivale était

fortement corrélée (inversement) aux valeurs de débit du cours d'eau principal (début août et début septembre 2020).

Le modèle GAM a montré des fenêtres précises pour les valeurs maximales de l'ARTP et une variabilité diurne significative durant le mois de Juillet 2021. L'aire du panache à la Station 1 augmente généralement à partir de minuit et atteint le maximum vers 10h00, ce qui est désigné par la région de contour en rouge foncé et/ou bleu clair (blanc) (**Figure 6**). À la même station, l'ARTP diminue progressivement pour atteindre le minimum entre 16h00 et 19h00, désigné par la région de contour en bleu foncé et/ou rouge clair (blanc). La Station 2 a présenté de faibles valeurs de superficie des refuges entre 13h00 et 16h00, tandis que les valeurs les plus élevées ont été observées tôt le matin juste avant 8h00. Un décalage horaire d'environ deux à trois heures est observé du moment d'occurrence des extrêmes entre les stations 1 et 2. La **Figure 7** illustre la variabilité de l'ARTP en juillet et août 2021 en utilisant une fonction de lissage cubique cyclique présentée à l'équation 6. La plus grande variabilité journalière relative a été observée en juillet 2021, avec des valeurs entre  $\approx 40 \text{ m}^2$  &  $\approx 80 \text{ m}^2$  pour les stations 1 et 2, respectivement. Au cours de cette période, l'ARTP à la station 2 était presque constante pour la première partie de la journée (00:00 - 07:59), entraînant de petites fluctuations. La variabilité journalière durant le mois d'août 2021 était moins importante à cause des faibles valeurs débits du cours d'eau principale. Les valeurs extrêmes journalières de l'ARTP ont été capturées avec plus de précision à la station 1, affichant une courbe cyclique et des créneaux horaires précis.

## 2.2.4- Performance de la prédiction avec les modèles d'apprentissage automatique

### 2.2.4.1- Procédure d'évaluation

Les trois procédures de validation décrites plus haut ont été testées pour chacun des deux sites et les résultats sont présentés à la **Figure 8**. En utilisant la méthode *LOOCV*, les stations 1 et 2 ont présenté des valeurs moyennes de rRMSE (moyenne sur les quatre modèles) de 20.53 % et 17.82 %, respectivement. Pour la méthode *k-fold* les rRMSE étaient de 21.73 % et 19.17 % et pour la méthode Train-Test Split les valeurs étaient de 21.45 % et 19.29 %. Dans notre étude de cas, la méthode *LOOCV* a donné les meilleures performances en termes de rRMSE et a été considérée pour la validation des résultats obtenues par les modèles d'apprentissage automatique. D'autre part, le biais de l'estimation de l'erreur est faible avec la procédure *LOOCV* puisque cette méthode utilise presque toutes les données pour l'entraînement (Alpaydin, 2020; Desai & Ouarda, 2021).

#### 2.2.4.2- Les modèles de références GAM et MARS

Dans la **Table 3** les aires de refuge interpolées avec la PID sont comparées aux prévisions horaires des modèles d'apprentissage automatique. GAM et MARS ont présenté en moyenne des coefficients de Nash plus élevés à la Station 1 (0.85) qu'à la Station 2 (0.78). Cependant, la Station 2 avait un rRMSE plus faible (20%) comparé à celui de la Station 1 (23%). En général, les deux modèles de référence avaient une performance relativement similaire pour la même station et aucun biais systématique n'est observé dans les prévisions horaires des deux modèles, en grande partie grâce à la méthode de validation-croisée *LOOCV*.

Le modèle GAM a utilisé cinq variables explicatives alors que le modèle MARS en a utilisé trois et quatre pour les stations 2 et 1 respectivement. Les résultats ont montré que ces deux techniques non paramétriques avaient un potentiel significatif pour décrire des estimations précises de la superficie de panache horaire et expliquer la plupart de la variation temporelle. Néanmoins, aucun des deux modèles de référence n'a dépassé les modèles de régression SVM ou RF, en termes de qualité de performance prédictive (rRMSE et Nash).

#### 2.2.4.3- La régression par machine à vecteur de support (SVM) et forêt aléatoire (RF)

Le modèle SVM a montré de meilleures performances si on le compare aux modèles GAM et MARS. Les stations 1 et 2 affichaient respectivement des valeurs de rRMSE entre 20.83 % et 18.57 %, des valeurs de Nash entre 89% et 83% et un biais relatif entre 1.8 % et 2.4 %. Le modèle SVM était un peu plus biaisé que les autres modèles. SVM et RF ont utilisé des ensembles identiques de variables prédictives pour les stations 1 et 2 ( $Q_m$ ,  $T_t$ ,  $T_a$  et  $T_\Delta$ ). Néanmoins, RF a surpassé tous les autres modèles avec en moyenne (sur les deux sites) des valeurs de rRMSE et de coefficient de Nash de 13,14% et 93%. D'après la même **Table 3**, les performances obtenues à la Station 2 étaient supérieures à celles de la Station 1. Ceci est dû au fait qu'il y avait moins d'observations disponibles à la Station 1 (4438 heures) comparé à la Station 2 (4846 heures) et aussi la surface globale couverte par le réseau de thermographes à la Station 2, qui était trois fois plus supérieure, en raison de l'importance du panache d'eau froide à cette station.

## 2.3- DISCUSSION ET CONCLUSION

### 2.3.1- Variabilité temporelle de l'ARTP

Nous avons estimé l'aire horaire maximale des refuges thermiques potentiels au cours des étés 2020-2021 et mis en évidence la variabilité selon différentes composantes temporelles (heure et jour Julien). Il est possible que le réseau de thermographes déployé eût parfois une couverture spatiale insuffisante pour couvrir l'intégralité du panache d'eau froide, mais la variation temporelle de l'ARTP était assez évidente. Les valeurs mesurées et estimées permettent une analyse comparative, mais il se peut qu'en absolu, certaines valeurs d'ARTP soient sous-estimées. Pour capturer l'étendue complète du panache, la résolution spatiale et la surface couverte par les capteurs devraient être augmentées en utilisant les technologies adaptées telles que l'imagerie aérienne infrarouge thermique (TIR) et les capteurs de température distribués (DTS) (Dzara et al., 2019). D'autre part, une importante variabilité des débits a été observée à travers les deux étés ou des valeurs élevées ont été enregistrées au cours de l'été 2020 et des niveaux critiques observés durant l'été 2021. Dans les confluences de tributaires de la rivière Ste-marguerite, Les débits dépassant le seuil  $50 \text{ m}^3/\text{s}$  sont souvent associés à un mélange important des masses d'eau. Par conséquent, la différence de température entre la rivière principale et le tributaire associé diminue au fur et à mesure que le débit augmente. Il pourrait en résulter une zone plus homogène aux confluences en termes de températures d'eau, empêchant ainsi la formation d'un refuge thermique à cet endroit.

En 2020, la Station 1 a connu le plus grand nombre de jours consécutifs sans observation d'ARTP. Cette particularité est principalement attribuable aux caractéristiques morphologiques et environnementales de ce site. Comme indiqué précédemment, la Station 1 est caractérisée par une pente très forte entre la rive et le thalweg. Un mélange homogène de température entre la rivière principale et l'affluent associé est donc plus probable à la Station 1 qu'à la Station 2. Cette dernière est caractérisée par un panache relativement isolé du cours d'eau principal car il repose sur un lit de gravier et de cailloux peu profond. Durant l'été 2020, les aires de refuge interpolées à la Station 2 étaient relativement constantes pour des valeurs de débit inférieures à  $50 \text{ m}^3/\text{s}$ , alors qu'à la Station 1, la moitié de cette quantité ( $25 \text{ m}^3/\text{s}$ ) réduisait l'ARTP de plus de 50 %.

Les conditions d'étiage en 2021 du cours d'eau principal sont caractérisées par une récurrence de 43 ans sur 7 jours consécutifs. Il en résulte un comportement très similaire des superficies du refuge



dans les deux stations durant cette période. Par exemple, dès le début du mois d'août, de nombreux capteurs ancrés proches de la rive ont été exondés, ce qui a conduit une chute considérable du coefficient de variation (CV) entre les deux étés. Cela s'explique par la baisse continue des valeurs de débit qui a entraîné un niveau d'eau historiquement bas et un cours d'eau principal à température élevée ( $T_m \approx 27 \text{ }^\circ\text{C}$ ). Lorsque les thermographes situés proche de la rive se sont retrouvés hors de l'eau, la température de l'ARTP est calculée avec un nombre plus restreint de thermographes et la variabilité diminue de manière significative. En effet, dans une confluence de tributaire, la variabilité latérale le long de la rive jusqu'au thalweg peut être très importante. En période de faibles débits et de températures élevées les superficies des refuges demeurent plus constantes qu'en d'autres circonstances. La **Figure 6** confirme cette conclusion, avec un gradient estimé constant (couleur homogène) à partir du début du mois d'août (jour Julien  $\geq 220$ ). En dépit de cette stabilité dans les superficies, la Station 2 a vu ses aires de panache se consolider autour de  $500 \text{ m}^2$ . Quant à la Station 1, une aire constante d'environ  $130 \text{ m}^2$  est observée durant la période d'étiage, mais qui reste relativement importante étant donnée la taille de la confluence à cette station.

Cela dit, alors que nos analyses ont porté sur les superficies, les volumes des refuges thermiques potentiels (VRTP en  $\text{m}^3$ ) n'ont pas été mesurés et pourraient être plus important à la Station 1 en raison de la pente raide et de la profondeur des bassins à cet endroit. Benda et al., (2004) ont montré que la géométrie locale d'un réseau fluvial affecte la morphologie de la confluence. Plus précisément, l'angle auquel les affluents entrent dans le cours principal est généralement aigu (moins de  $90^\circ$ ), mais lorsque l'angle approche  $90^\circ$  (par exemple, Station 1), la morphologie du lit de la rivière à la confluence est fortement influencée, ce qui peut avoir des effets sur l'étendue de l'ARTP/VRTP.

### **2.3.2- Performance de la modélisation statistique**

GAM était le modèle qui utilisait le plus de variables explicatives pour la prédiction horaire des aires de refuge thermique potentiels. Le modèle MARS n'a pas inclus la température du tributaire ( $T_1$ ) pour les deux stations, car une partie de l'information contenue dans  $T_1$  est incluse dans la variable  $T_\Delta$ . D'autre part, MARS ne considère pas la température de l'air  $T_a$  comme une variable explicative à la Station 2, ce qui peut être expliqué par la distance qui la sépare de la station météorologique (CIRSA), où la température de l'air est mesurée. SVM et RF n'ont pas inclus la température du cours d'eau principal ( $T_m$ ), mais ont considéré la différence de température  $T_\Delta$  à la

place. La **Figure 9** présente les aires de refuge estimées avec la PID versus celles prédites par les modèles de régressions. Les points rouges et bleus représentent respectivement les données de 2020 et 2021. À la Station 1, aucune différence significative n'a été observée entre les deux étés à cause des bassins profonds qui caractérisent cette station. Cependant, à la Station 2, l'effet de l'événement d'étiage survenu en 2021 est clairement visible, avec deux populations distinctes (2020 et 2021). Les modèles ont surestimé l'ARTP pendant la période de juin et juillet 2020 (points rouges  $\approx 390 \text{ m}^2$ ) et ont sous-estimé celle d'août 2021 ( $500 \text{ m}^2 \leq$  points bleus  $\leq 700 \text{ m}^2$ ). Pour les valeurs supérieures à  $390 \text{ m}^2$ , la performance de prédiction est améliorée comparé aux modèles de références (GAM et MARS) en utilisant les modèles SVM et RF, qui ont rapproché les points de la ligne diagonale ( $y=x$ ). Les modèles développés avec un pas de temps horaire ont donné de bon résultats et parfois même meilleurs que ceux construits par Saadi et al. (2021) sur une base journalière, en utilisant les modèles de régression SVM et RF.

Dans cette étude, les approches non paramétriques se sont avérées utiles pour décrire les estimations horaires des aires de refuge thermique. Les forêts d'arbres de décision (RF) ont capturé la majorité de la variation des aires de refuge au cours de la saison estivale, ou des valeurs élevées ont été observées pendant une longue période dues aux températures élevées et au faibles niveaux d'eau observés. Les résultats montrent que les deux stations contiennent des refuges thermiques potentiels et que la superficie peut être prédite avec une erreur inférieure à 15% (rRMSE) en utilisant la régression RF. D'autres variantes des modèles utilisant les arbres de décision méritent d'être étudiées, comme les arbres extrêmement aléatoires (*Extremely Randomized Tree* - ERT), l'*eXtreme gradient boosting* (XGBoost) et le modèle *M5Tree* qui ont démontré d'excellentes performances dans la modélisation de la température de l'eau des rivières des bassins versants Autrichiens. (Feigl et al., 2021; Heddam et al., 2020). Nous n'avons pas inclus les modèles ANN, qui devraient également être considérés pour les prochaines études comparatives, en raison de leur performance exceptionnelle dans la prédiction à long terme dans plusieurs applications hydrologiques récentes (Qiu et al., 2021; L. Wang et al., 2022). La capacité de prédire la surface d'un refuge thermique dans une confluence de tributaire à un pas de temps horaire et de pouvoir étudier sa variabilité temporelle en utilisant des outils non déterministes et des approches non paramétriques, peut-être d'une grande aide pour les gestionnaires de ces habitats naturels. Les études futures devraient également :

- (i) Explorer la possibilité de construire des modèles multi-sites à une échelle régionale
- (ii) Inclure les prévisions hydrométéorologiques et le débit des tributaires comme variables explicatives supplémentaires, pour améliorer les performances des modèles d'apprentissage automatique
- (iii) Analyser la variabilité diurne dans un contexte de changement climatique.

Enfin, bien que cette étude ait été principalement motivée par le fait que les refuges thermiques sont des habitats clés pour l'ichtyofaune, les modèles développés dans cette étude peuvent avoir d'autres applications. Des panaches thermiques artificiels existent dans de nombreux environnements fluviaux et côtiers. Le cas le plus courant est probablement le rejet d'effluents d'eau chaude d'une centrale électrique (par exemple, Penk & Williams, 2019; Zavorsky & Duester, 2020) mais d'autres applications pourraient inclure d'autres phénomènes naturels, tels que des rivières se jetant dans des lacs (Smith & Simpkins, 2018), les affluents affectant les régimes thermiques des lacs (Råman Vinnå et al., 2018) ou les panaches d'eau douce dans les estuaires (Huang et al., 2020).

### **3- ARTICLE VERSION ANGLAISE INTÉGRALE**

### 1 3.1- ABSTRACT

2 To avoid heat stress during high summer river temperatures, aquatic ectotherm species such  
3 as the Atlantic salmon (*Salmo salar*) seek refuge in discrete areas of cold water called thermal  
4 refuges. Small perennial tributary inflows in warm conditions have been found to create reach-  
5 scale cold-water plumes at confluences, potentially vital for cold-water taxa survival. This paper  
6 investigated the potential thermal refuge area (PTRA) during two consecutive summers (2020-  
7 2021) at a sub-daily time-step in two tributary confluences of the Ste-Marguerite River (Canada).  
8 We first delineated the hourly PTRA at both confluences using a spatial interpolation technique  
9 termed Inverse Weighted Distance (IWD). Then, the relative variation of PTRA was investigated  
10 at a sub-daily time step, highlighting the impact of the low flow conditions (2021) on the diel cycle  
11 of PTRA extremes. Furthermore, we used four different supervised machine learning regression  
12 models and a few hydrometeorological variables to assess hourly PTRA estimates. The proposed  
13 algorithms were (i) The multivariate adaptive splines regression (MARS), (ii) The generalized  
14 additive model (GAM), (iii) The support vector machine regression (SVM), and (iv) The random  
15 forest regression (RF). The results showed that GAM and MARS were outperformed by tree-based  
16 and kernel-based regression models SVM and RF. The latter worked best with high accuracy at  
17 both sites and showed superior performance with mean relative root mean square error (rRMSE)  
18 and mean Nash–Sutcliffe efficiency coefficient (Nash) of 13% and 93%. Due to the decreasing  
19 availability of thermal refuge for salmonids, monitoring stream temperatures at small spatial and  
20 temporal scales, in conjunction with data-driven techniques, is an essential first step toward fully  
21 understanding stream temperature heterogeneity at tributary confluences.

22 **Keywords:** Hourly potential thermal refuge area (PTRA); Hourly water temperature; Tributary  
23 confluences; Diel variability; Machine learning; Regression model.

24

25

26

## 27 3.2- INTRODUCTION

28 On a global scale, climate change is expected to significantly change aquatic organisms'  
29 distribution and biophysical processes across lotic ecosystems (Dugdale et al., 2018; Isaak &  
30 Rieman, 2013). Numerous rivers around the world are likely to experience an overall increase in  
31 water temperature according to most climate change scenarios (Dugdale et al., 2018; Morrill et al.,  
32 2005; van Vliet et al., 2013). High stream water temperature can adversely affect fisheries'  
33 resources by limiting fish habitat and generating sublethal heat stress or even mortality (Caissie,  
34 2006; Caissie et al., 2004; Lund et al., 2002). Under those circumstances, thermal patchiness has  
35 been widely recognized as a moderator for the harmful effects of climate change on a variety of  
36 ectotherm species, many of which reside in environments that are on the edge of their thermal  
37 tolerance range, such as salmonids (Corey et al., 2020; Isaak et al., 2015). Such species reduce their  
38 body temperature by moving to habitat patches that remain cool ( $< 21\text{ }^{\circ}\text{C}$ ) during periods of warm  
39 weather (Armstrong et al., 2016; Breau et al., 2011; Dugdale et al., 2015; Torgersen et al., 1999).  
40 These cold habitats allow individuals to survive regardless of the warming and are often high  
41 priorities for conservation (Davis et al., 2013; Dugdale et al., 2013). There have been many  
42 different definitions of salmonid fish thermal refuges, including areas with temperatures above  
43  $23^{\circ}\text{C}$  (Sutton et al., 2007), temperatures below the mean mainstem temperature (Baird & Krueger,  
44 2003), and temperatures  $>3^{\circ}\text{C}$  cooler than surrounding areas (Brewitt & Danner, 2014). Sullivan  
45 et al. (2021) came up with an eco-hydrological typology that classifies thermal habitat for cold-  
46 water fish species as follows: (i) A cold-water patch which refers to an area of water temperature  
47 cooler ( $2\text{--}10^{\circ}\text{C}$ ) than ambient streamflow immediately upstream (Ebersole et al., 2001, 2015), (ii)  
48 A thermal refuge is a cold-water patch used by poikilotherms, (iii) Physiological refuge refers to a  
49 patch of water with temperatures below a biologically critical threshold (e.g., Breau et al., 2007).  
50 Nonetheless, Greer et al., (2019) showed that different threshold-based definitions can lead to  
51 different conclusions about thermal refuge status and may oversimplify different aspects of a  
52 refuge.

53 Poikilotherms fish could often perceive temperature changes of less than  $0.5^{\circ}\text{C}$  (Murray, 1971).  
54 However, Rainbow trout (*Oncorhynchus mykiss*) has been shown to respond to temperature  
55 differences as low as  $0.1^{\circ}\text{C}$  under hourly decreasing water temperature (Bardach & Bjorklund,  
56 1957). During summer, medium-sized tributaries (Strahler order: 3-4) with a constant cool flow

57 may create reach-scale cold-water plumes at confluences where the tributaries enter the mainstem  
58 (Dugdale et al., 2013; Torgersen et al., 2012). Such lateral, point-source intrusions of cold water  
59 from perennial tributaries often lower the mainstem water temperature on a portion of the channel  
60 (usually near the banks) through a mixing effect and have greater discharge than in-stream  
61 hyporheic or groundwater upwelling (Gendron, 2013; Poole & Berman, 2001; Sutton et al., 2007).

62 By identifying tributary confluences as potentially thermal refuges, target species could be  
63 encouraged to utilize such shelters more extensively. As an example, Biron et al., (2004) installed  
64 channel deflectors upstream of tributary confluence in order to extend the cold-water plume into  
65 the river channel, which would prevent mixing between tributaries and mainstems. However,  
66 juvenile steelhead prefers thermally mixed regions that allow them to thermoregulate while  
67 feeding (Brewitt et al., 2017). Alternatively, management strategies such as riparian shading can  
68 decrease water temperatures by as much as 4°C (Ebersole et al., 2015; Marteau et al., 2022) and  
69 provide cover from predatory birds (Kurylyk et al., 2015), all while keeping flow relatively  
70 undisturbed.

71 The need to characterize, quantify and investigate the possible temporal evolution of these  
72 microhabitats sparked a lot of interest in statistical techniques using supervised machine learning,  
73 which delivered promising results regarding the use or distribution of cold-water fish species  
74 (Frechette et al., 2018; Jeong et al., 2013; T. Wang et al., 2020; Wilbur et al., 2020). Ebersole et al.  
75 (2015) modeled the occurrence of cold-water patches at tributary confluences as a function of the  
76 watershed and climatic characteristics. More recently, in a case study of the Ste-Marguerite River  
77 (Canada), Saadi et al. (2021) developed a dependable tributary-river scale statistical model that  
78 uses a limited number of hydrometeorological predictors to estimate the daily mean potential  
79 thermal refuge area (PTRA). They also suggested that these cold-water patches exhibited  
80 significant diel variability in some cases. In water temperature studies, nonparametric regression  
81 has been shown to capture the nonlinear relationship between variables better and provide higher  
82 performance than many parametric models (Adamowski & Labatiuk, 1987; Benyahya et al., 2007;  
83 Caissie et al., 2001; Chenard & Caissie, 2008; Daigle et al., 2010; St-Hilaire et al., 2012).  
84 Nonparametric models and artificial neural networks (ANN) are the main supervised machine  
85 learning structures used for this purpose. In this study, we did not include ANN models but rather  
86 investigated the predictive power of different nonparametric regression algorithms. ANN has a

87 fixed number of parameters but a large number of "hidden" parameters; hence they are usually  
88 categorized as nonparametric (Lee et al., 2017).

89 As a follow-up to the work of Saadi et al. (2021), temporal resolution is increased in the present  
90 study, as we estimated hourly PTRA using hourly maximum water temperature and the Inverse  
91 Weighted Distance (IWD) interpolation method. We also assessed the performance of four  
92 nonparametric regression models for hourly PTRA forecasting. The proposed machine learning  
93 models include the Generalized Additive Model (GAM; Hastie & Tibshirani, 1987), the  
94 Multivariate Adaptive Regression Splines (MARS; Friedman, 1991), the Support Vector Machine  
95 regression (SVM; Vapnik, 1998), and the Random Forest regression (RF; Breiman, 2001). These  
96 models could be classified based on their learning algorithms, where GAM and MARS are  
97 extensions to linear-based learners, SVM is a kernel-based learner, and RF is a tree-based learner  
98 (Zhang et al., 2020).

99 In previous studies, GAM successfully modeled water temperature using daily average air  
100 temperature and discharge as main predictors at the Sainte-Marguerite River, Northern Canada  
101 (Laanaya et al., 2017). Frechette et al. (2019) applied GAM to assess the relationship between fish  
102 count and key environmental variables such as discharge and temperature. The MARS approach  
103 improved regional flood frequency analysis (RFA) and slightly outperformed GAM (Msilini et al.,  
104 2020). Saadi et al. (2021) simulated daily mean PTRA using GAM and MARS. The latter had  
105 better accuracy and used fewer predictors, but both models showed strong potential. SVM has been  
106 successfully used to model a range of hydrological processes (Deka, 2014). Weierbach et al. (2021)  
107 compared multiple linear regression (MLR), SVM, and RF to predict monthly water stream  
108 temperature. They pointed out the weakness of RF in predicting extreme values but the overall  
109 good performance of SVM. Allahbakhshian-Farsani et al. (2020) applied SVM, MARS, and  
110 boosted regression tree (BRT) in regional frequency analysis and showed that the SVM model  
111 based on the radial basis function (RBF) kernel resulted in the best performance. Quan et al. (2020)  
112 used SVM and the genetic algorithm (GA) to predict the water temperature of large high-altitude  
113 reservoirs in western China. The result provides valuable insights into predicting vertical water  
114 temperature at different depths of reservoirs. Tree-based regression models proved to be  
115 competitive with ANN and showed great potential for water temperature modeling. Desai &  
116 Ouarda. (2021) introduced random forest regression (RF) combined with the Canonical correlation



117 analysis (CCA) in RFA and resulted in high prediction quality. Ferchichi et al. (2021) used RF and  
118 backpropagation neural network (BPNN) to study the possible impact of future coastal water  
119 temperature scenarios on the risk of potential growth of marine bacteria. They observed very  
120 similar performance between the two models. Feigl et al. (2021) predicted stream temperature in  
121 Europe using different machine learning methods. They showed that the tree-based model termed  
122 extreme gradient boosting (XGBoost) had comparable performance to the feed-forward neural  
123 network (FNN).

124 Kurylyk et al. (2015) suggested that Atlantic salmon parr use thermal refuges with temperature  
125 differences with the immediate upstream less than 2°C, which is consistent with Saadi et al. (2021)'  
126 definition. Consequently, we defined PTRA as water temperature regions at least one-degree  
127 Celsius colder than the immediate upstream. This criterion was selected to have a sufficiently long  
128 hourly PTRA time series and ensure that the minimum temperature difference (1°C) is greater than  
129 the sensors' precision (0.5°C) (Saadi et al., 2021).

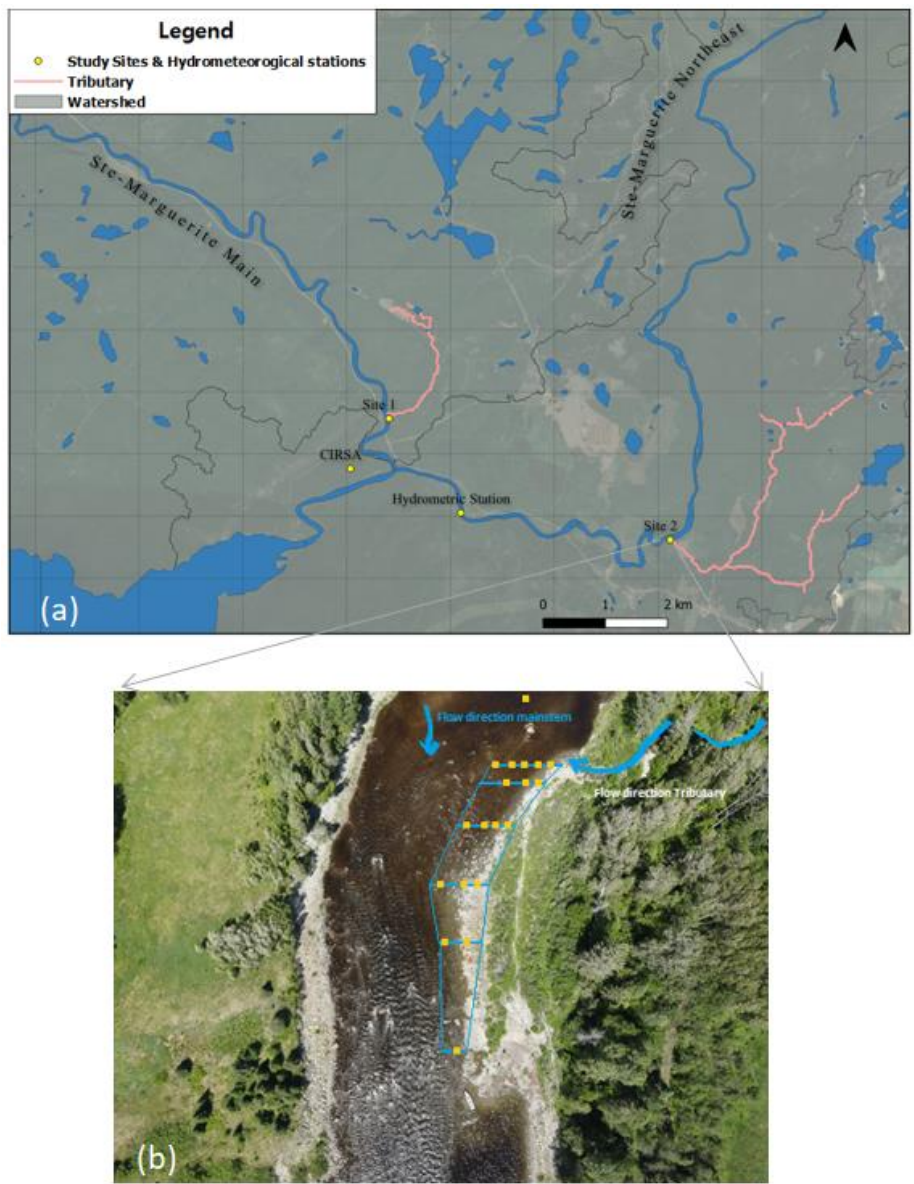
130 The purpose of the current study is to:

- 131 (i) Investigate the PTRA' diel variability across two consecutive summers (2020-2021)  
132 with an emphasis on the low flow period (2021),
- 133 (ii) Compare the performance of four regression models for hourly PTRA estimates using  
134 a few hydrometeorological variables.

### 135 **3.3- STUDY SITE**

136 Two sites were identified in the Sainte Marguerite River on Quebec's North Shore, between  
137 Chicoutimi and Sacré-Coeur, Canada. The river splits into three distinct branches, two of which  
138 are considered in this study (**Figure 1a**). The main river runs from the west along Provincial Road  
139 172 parallel to the Saguenay Fjord and flows into Sainte-Marguerite Bay. The northeast and  
140 northwest branches are more difficult to access, mainly served by gravel roads. In this study, the  
141 first site is located on the main branch, and the second is on the northeast branch. The drainage  
142 area is 1,097 km<sup>2</sup> and 980 km<sup>2</sup> for the main and the northeast branches, respectively, out of 2,077  
143 km<sup>2</sup>. Both stations were selected because of known cold affluents that are not ephemeral. Station 1  
144 has the coldest tributary and is characterized by a lower reach and series of steep step pools. The  
145 angle between the tributary and the mainstem at the confluence is approximately 90 degrees.

146 Station 2 has a downward slope and a tributary entry angle of roughly 60 degrees. The dominant  
147 substrate categories are gravel and cobble for both sites.



148  
149 **Figure 1:** (a) Location of river mainstems and tributaries, sites, and hydrometeorological stations,  
150 and (b) example of deployed water temperature sensors (in yellow) arrays at Ste-Marguerite  
151 northeast (Station 2).

## 152 3.4- MATERIAL AND METHODS

### 153 3.4.1- Estimation of hourly PTRA

154 Both stations 1 and 2 were monitored during the summers of 2020 and 2021 (**Table 1**).  
 155 Parallel arrays of thermographs were placed at the confluence of tributaries and river mainstems  
 156 (**Figure 1b**). We created sensor grids by anchoring rebar into the streambed, on which  
 157 thermographs (Pendant Hobo Temperature Logger, accuracy =  $\pm 0.5$  °C) were attached to form  
 158 cross-stream rows and along-stream columns. Grids were designed to capture the cold-water  
 159 plumes of the incoming tributaries that were first observed using a thermal camera (Saadi et al.,  
 160 2021; Wang et al., 2020). An additional sensor was deployed upstream of each mainstem and  
 161 tributary to record water temperatures ( $T_m$  and  $T_t$ ) that were not influenced by the mixing process  
 162 at the confluence. The lateral distance between thermographs in each transect varied (up to 15 m)  
 163 according to the morphology of the sites (e.g., bathymetry, size of the tributary, angle at the  
 164 confluence between tributary and associated mainstem). Sensors recorded data at 15 min time  
 165 steps, from which we computed maximum hourly water temperatures. We estimated hourly PTRA  
 166 using the water temperature grid and a spatial interpolation technique called Inverse Weighted  
 167 Distance (IWD). The area determined the plume boundaries with water temperature at least one-  
 168 degree Celsius colder than the main river stem.

169 **Table 1:** Monitoring period and the number of deployed water temperature sensors for stations 1  
 170 and 2

Site-Year	Deployed sensors		Monitoring period		Observations (Hours)		
	2020	2021	2020	2021	2020	2021	Total
Station 1	18	18	June 20 <sup>th</sup> to September 23 <sup>rd</sup>	June 9 <sup>th</sup> to September 5 <sup>th</sup>	2303	2135	4438
Station 2	17	18	June 20 <sup>th</sup> to September 23 <sup>rd</sup>	June 9 <sup>th</sup> to September 23 <sup>rd</sup>	2303	2543	4846

171

#### 172 3.4.1.1- *Inverse weighted distance interpolation method (IWD)*

173 The IWD is a deterministic spatial interpolation method that uses known values (gauged  
 174 sites) and corresponding weighted values to estimate an unknown value at a given location  
 175 (ungauged sites). This process requires the estimation of weights  $w_i$  as an inverse function of the  
 176 distance between the site of interest and gauged sites. Equation 1 shows the basis of IWD  
 177 interpolation:

178 
$$Z_j = \frac{\sum_{i=1}^n x_i * w_i}{\sum_{i=1}^n w_i}$$
 Equation 1

179 Where:

180 
$$w_i = \frac{1}{D_{ij}^p}$$

181  $Z_j$  is the estimated value of water temperature at the ungauged site  $j$ ,  $x_i$  is the water temperature  
182 value of the neighboring gauged site,  $w_i$  is the weight assigned to the gauged sites  $i$ ,  $D_{ij}$  is the  
183 distance between the gauged site  $i$  and ungauged site  $j$  of interest,  $n$  is the number of gauged sites,  
184 and  $p$  is the exponent.

185 The extrapolation error is evaluated using the leave-one-out cross-validation (LOOCV) approach.  
186 This process is closely related to the statistical method of Jackknife estimation (Efron, 1982). Root  
187 Mean Square Error (RMSE) is used to assess the cross-validation error (RMSE<sub>station 1</sub> = 0.3 °C and  
188 RMSE<sub>station 2</sub> = 0.27 °C). We wanted to make sure that the overall error, which includes both sensor  
189 accuracy and cross-validation error, remained less than the one-degree Celsius criteria used to  
190 establish the minimal plume boundaries ( $|\pm 0.5|^\circ\text{C} + \text{RMSE}_{\text{station } i} \leq 1^\circ\text{C}$ ). The average overall  
191 errors for stations 1 and 2 were 0.8°C and 0.77°C, respectively. All programming and computing  
192 for IWD were completed using Python 3.8, and cKDTree function class from  
193 "Scipy.Spatial" package (docs.scipy.org/doc/scipy/reference/spatial.html).

194 3.4.1.2- *Potential predictive variables*

195 To predict maximum hourly PTRA at both stations, we evaluated four machine learning  
196 models and used up to five potential predictors on an hourly time step, as shown in **Figure 2**:

- 197 (i) Discharge of mainstream ( $Q_m$ ),  
198 (ii) Air temperature ( $T_a$ ),  
199 (iii) Upstream mainstem water temperature ( $T_m$ ),  
200 (iv) Upstream tributary water temperature ( $T_t$ ),  
201 (v) The temperature difference between  $T_m$  and  $T_t$  ( $T_\Delta$ ).

202 Hourly  $Q_m$  is obtained from the hydrometric station #062803, located on the northeast branch of  
203 Ste-Marguerite (48° 16' 5" N, 69° 54' 33" W); records are downloaded from the provincial ministry  
204 of the environment and fight against climate change (MELCC). Hourly  $T_a$  is measured from a

205 pressure-temperature sensor (Onset Hobo-Level Logger, accuracy =  $\pm 0.44^{\circ}\text{C}$ ) located at CIRSA  
206 meteorological station (**Figure 1a**). All potential predictive variables are used to identify each  
207 model's best subset of inputs using a backward selection method known as Recursive Feature  
208 Elimination (RFE).

#### 209 3.4.1.3- *Feature selection (Recursive feature elimination)*

210 RFE is a feature selection algorithm that searches for a subset of features by starting with  
211 all features in the training dataset and successively removing them based on the lowest contribution  
212 to the model's accuracy. The goal is to fit the given machine learning algorithm, rank features by  
213 importance, discard the least important features, and re-fit the model (Kuhn & Johnson, 2013). The  
214 best-selected subset is the one that optimizes most of the performance criteria. RMSE was used as  
215 the performance criterion for selecting the best subset of features. We use the k-fold cross-  
216 validation (k=10) in the performance prediction assessment of the possible feature subsets. First,  
217 the goodness of model fitting is evaluated by computing the RMSE for each feature subset  
218 ( $\text{RMSE}_j$ ). Then, the average  $\text{RMSE}_{cv}$  is calculated, and the equations are shown below:

$$219 \text{RMSE}_{cv} = \frac{1}{k} \sum_{j=1}^k \text{RMSE}_j \quad \text{Equation 2}$$

220 For each fold  $j$ , RMSE is computed as follows:

$$221 \text{RMSE}_j = \sqrt{\frac{1}{N} \sum_{i=1}^N (o_i - s_i)^2} \quad \text{Equation 3}$$

222 Where  $k$  is the number of folds (10),  $j$  is one of the  $k$  folds;  $N$  is the sample size of the fold  $j$ ,  $s_i$  are  
223 simulated/predicted PTRAs, and  $o_i$  represents the observed PTRAs.

### 224 3.4.2- Machine learning modeling

#### 225 3.4.2.1- *Generalized additive model (GAM)*

226 GAM is a nonparametric and extended version of the generalized linear model GLM  
227 (McCullagh & Nelder, 1989). This model replaces the linear predictor in GLM with an additive  
228 one, which can use categorical and continuous data. GAMs model the continuous data as a  
229 nonlinear smoothing function by regression splines that can take various forms (Dominici et al.,  
230 2002). In order to avoid overfitting, we employed the penalized GAM in this study, where  
231 smoothing functions or nonlinear "basis functions" (BFs) are simple cubic splines (also referred to

232 as penalized splines) (Wood, 2006). In addition, a Gaussian distribution family is used for the  
233 response variable (Marra & Wood, 2011). GAM is expressed as follows:

$$234 \quad g(E(Y)) = \beta_0 + s_1(x_1) + s_2(x_2) \dots + s_p(x_p) + \varepsilon \quad \text{Equation 4}$$

235 The link function  $g$  is a parametric function that links the dependent variable mean to a set of  
236 explanatory variables;  $E(y)$  is the expectation of predicted response variables.  $\beta_0$  is the intercept,  
237  $s_i$  is the  $i^{\text{th}}$  explanatory variable's smooth function,  $x_i$  represent the independent variables used to  
238 forecast PTRAs,  $\varepsilon$  is an error term. GAM was a powerful approach for estimating daily PTRA  
239 using a small number of easy-to-measure hydrometeorological predictors (Saadi et al., 2021).

240 Diel variability was also investigated across Julian days using tensor product smooth  $Te$  (equation  
241 5) to analyze the combined effect of the two temporal variables (Julian days & hours) on the  
242 response variable (PTRA). In addition, factor-smooth interaction was used to create a separate  
243 smoother for each month and cyclic cubic splines to reflect the periodic diel pattern of PTRA during  
244 the warmest months of summer (July & August) (equation 6):

$$245 \quad PTRA \sim te(j, h) \quad \text{Equation 5}$$

$$246 \quad PTRA \sim m + s(h, by = 'm')$$

247 Where  $PTRA$  represents the response variable, Julian days are represented by  $j$ ,  $h$  represents the 24  
248 hours of the day. Months are represented by  $m$  and refer to July and August. Cyclic cubic splines  
249 are selected as smoothing functions for the variable ' $h$ '.

250 Finally, we organized the time of day across different periods to assess PTRA variability across a  
251 24-hour cycle (Eastern Standard Time) (Mahardja et al., 2021):

- 252 • Early (00:00-07:59),
- 253 • middle (08:00- 15:59),
- 254 • Late (16:00-23:59).

255 The coefficient of variation of estimated PTRA ( $CV = \frac{\text{Standard deviation}}{\text{Mean}}$ ) were then calculated over  
256 the two warmest summer months (July & August) for each of the three periods described above.  
257 CV measures the extent of variation in reference to the population mean. The larger the dispersion,

258 the bigger the CV. The coefficient of variation is useful because it is dimensionless (i.e.,  
259 independent of the unit of measurement) and thus comparable across data sets with various units  
260 or widely differing means.

#### 261 3.4.2.2- *Multivariate adaptive regression splines (MARS)*

262 The MARS model can be seen as a flexible extension of GAM, expressed as a linear  
263 combination of basis functions and their interactions (Msilini et al., 2020). It is a flexible  
264 nonparametric regression approach that can deal with high-dimensional data. The MARS model  
265 builds a suite of linear regression models by subdividing all explanatory variables into several  
266 regions (Conoscenti et al., 2015). The break values between the regions are called "knots", and the  
267 intervals are designed as linear basis functions (BFs) (Conoscenti et al., 2016). The knots mark the  
268 end of one region of data and the beginning of another. The BFs are generated by stepwise research,  
269 and the comparison of sub-models is made based on the generalized cross-validation (GCV) (Roy  
270 et al., 2018). The MARS nonlinear model is a combination of several BFs as follows:

$$271 Y = \beta_0 + \sum_{k=1}^N \beta_i f_i(x) \quad \text{Equation 7}$$

272 Where  $\beta_0$  is the intercept, and  $\beta_i$  are the regression coefficients of the basis functions  $f_i(x)$ , N is  
273 the number of basis functions.

274 In a previous study (Saadi et al., 2021), MARS predicted daily mean PTRA more accurately than  
275 GAM and employed fewer predictors. The two models were chosen for this study to test their  
276 adaptability for hourly PTRA modeling.

#### 277 3.4.2.3- *Support vector machine regression (SVM)*

278 SVM was introduced by (Vapnik, 1998) as a robust machine learning tool. SVM maps the  
279 original data sets from the input space to a high-dimensional or even infinite-dimensional feature  
280 space to make regression problems easier in that feature space (Deka, 2014). This approach  
281 simultaneously minimizes model complexity and prediction error, using Kernel functions and soft  
282 margin principles. Because support vector regression optimization does not depend on the size of  
283 the input space, it will have advantages in high dimensionality space (Drucker et al., 1996). In this  
284 regard, nonlinear kernels allow the SVM to model complex separating hyperplanes (nonlinear  
285 boundaries) (Smits & Jordaan, 2002). In the present context, we used the Vapnik's *epsilon* ( $\epsilon$ ) SVM,  
286 where the value of  $\epsilon$  defines a margin of tolerance where no penalty is given to errors. The larger  $\epsilon$

287 is, the larger errors are admitted in the solution. By contrast, if  $\epsilon$  is very low, every error is  
 288 penalized. SVM uses a loss function  $L_\epsilon(y, g(x_i, w))$  (equation 8), which ensures the existence of  
 289 global minimum and, at the same time, the optimization of reliable generalization bound :

$$290 \quad L_\epsilon(y, g(x_i, w)) = \begin{cases} 0 & \text{for } |y - g(x_i, w)| \leq \epsilon \\ |y - g(x_i, w)| - \epsilon & \text{otherwise} \end{cases} \quad \text{Equation 8}$$

291  $x_i$  represents the vector of input features, and  $w$  is the set of weights. Using the  $\epsilon$ -insensitive loss  
 292 function, one can find  $g(x_i, w)$  that can better approximate the actual output vector  $y$ . The support  
 293 vectors are the instances across the margin.

294 We used the  $\epsilon$ -SVM with a nonlinear Radial Basis Function kernel (equation 9) because of its  
 295 similarity to the Gaussian distribution :

$$296 \quad K(x_i, x) = \exp(-\gamma(x_i - x)^2), \gamma \geq 0 \quad \text{Equation 9}$$

297  $(x_i - x)^2$  is defined as the squared Euclidean distance between the two feature vectors, " $\gamma$ " is the  
 298 parameter in RBF, controlling the radial range.

299 SVM tries to reduce model complexity by minimizing  $\|w^2\|$ . This can be described by introducing  
 300 slack variables or the concept of Soft margin ( $\xi_i, \xi_i^* \geq 0 \quad i = 1, \dots, n$ ) to measure the deviation of  
 301 training samples outside the  $\epsilon$ -insensitive zone, and a penalty error term is added using the  
 302 regularization parameter  $C$  :

$$303 \quad \text{Min } \frac{1}{2} \|w^2\| + C \sum_1^n (\xi_i + \xi_i^*)$$

$$304 \quad \begin{cases} y_i - g(x_i, w) \leq \epsilon + \xi_i^* \\ -y_i + g(x_i, w) \leq \epsilon + \xi_i \end{cases}$$

305 The solution to this minimization problem is the decision function of nonlinear SVM expressed as  
 306 follows:

$$307 \quad g(x) = \sum_{i=1}^l (\alpha_i - \bar{\alpha}_i) K(x_i, x) + b \quad \text{Equation 10}$$

308  $l$  is the number of support vectors and the dual set of variables  $\alpha_i$  and  $\bar{\alpha}_i$  represent the  
 309 corresponding constraints, and  $b$  is the bias term.

#### 310 3.4.2.4- Random Forest regression (RF)

311 RF was introduced by (Breiman, 2001) as an ensemble learning model based on the idea of  
 312 bagging (bootstrap aggregating) (Breiman, 1996). Bagging predictors average multiple model  
 313 predictions, where each model is trained on a bootstrapped sample instead of the full-observed



314 sample. This randomness introduced by bootstrapping increases the model's ability to generalize  
315 and produce stable prediction results. RF models are bagging predictors using classification and  
316 regression trees (CARTs) as a base learner. RF recursively applies binary splits to the data to  
317 minimize entropy in the tree nodes (a measure of dissimilarity). This is done until each node reaches  
318 a minimum node size or a previously defined maximum tree depth is reached. Breiman (2001)  
319 showed that adding further randomness to the bagging method improves prediction accuracy. In  
320 RF, this is achieved by only selecting a random subset of available variables for the split at each  
321 node. The estimate for the dependent variable is given by:

$$322 \quad Y = \frac{1}{M} \sum_{m=1}^M f_m(X) \quad \text{Equation 11}$$

323 Where  $f_m$  denotes a single fitted regression tree, X the matrix of predictors, and Y the dependent  
324 variable vector (PTRA).

325 All machine learning modeling was done in R (R Core Team, 2020). The models' development  
326 relied heavily on 'Caret' (Kuhn et al., 2020) for recursive feature selection (RFE) and hyper-  
327 parameter optimization using grid search and cross-validation techniques. Implementation of  
328 GAM, MARS, SVM, & RF was carried respectively under the following packages : 'mgcv'  
329 (Wood, 2006), 'earth' (Milborrow, 2018), 'e1071' (Meyer & Wien, 2001), and 'random Forest'  
330 (Breiman, 2001). The visualizations were completed using ggplot2 (Wickham & Grolemund,  
331 2016).

### 332 **3.4.3- Evaluation procedures and performance metrics**

333 In machine learning modeling, a robustness test is usually conducted by calibrating a model  
334 using one part of the data and then validating it with another. This is a widely used and accepted  
335 procedure known as the split sample approach. Nevertheless, Shen et al., (2022) suggest this  
336 approach has some drawbacks and should be abandoned. The researchers contend that all data  
337 should be used for model development and calibration before the model is integrated into decision-  
338 making processes, based on an impressive empirical study of two hydrological models and 463  
339 catchments in the United States. However, in our study case, the models' robustness is assessed  
340 using three different validation procedures due to the novelty of the nondeterministic approach for  
341 hourly maximum PTRA estimation. The proposed procedures are : (i) Split Sample, (ii) K-fold  
342 cross-validation, and (iii) Leave-one-out cross-validation (LOOCV). We used rRMSE (Equation

343 14) to assess the performance of the best evaluation procedure. The following metrics are used to  
 344 evaluate the quality and the accuracy of the forecasts:

345 1- Nash-Sutcliffe efficiency coefficient (Nash)  $= 1 - \frac{\sum_{i=1}^n (o_i - s_i)^2}{\sum_{i=1}^n (o_i - o_m)^2}$  Equation 12

346 2- Root Mean Squared Error (RMSE)  $= \sqrt{\frac{1}{n} \sum_{i=1}^n (o_i - s_i)^2}$  Equation 13

347 3- Relative RMSE (rRMSE)  $= \sqrt{\frac{1}{n} \sum_{i=1}^n \frac{(o_i - s_i)^2}{o_i^2}}$  Equation 14

348 4- Bias  $= \frac{1}{n} \sum_{i=1}^n (o_i - s_i)$  Equation 15

349 5- Relative Bias (rBias)  $= \frac{1}{n} \sum_{i=1}^n \frac{(o_i - s_i)}{(o_i)}$  Equation 16

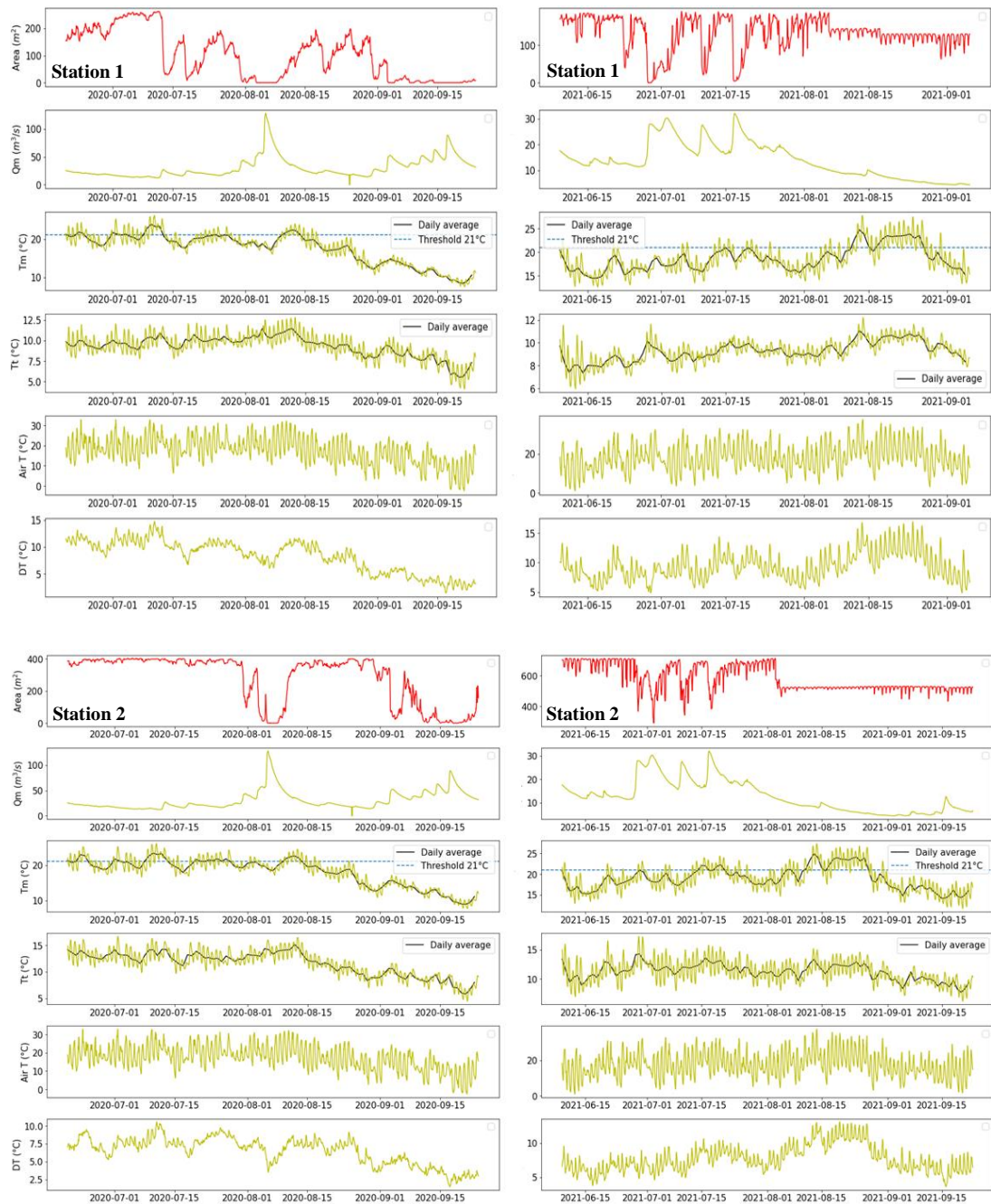
350 Where  $o_i$  is the observed value at site  $i$ ,  $s_i$  is the simulated value using the model for the same site  
 351  $i$ ,  $o_m$  is the mean observed value, and  $n$  is the sample size site  $i$ .

### 352 3.5- RESULTS

#### 353 3.5.1- PTR A & mainstem discharge in 2020

354 **Figure 2** shows that the cold-water plume area in 2020 reached 0 m<sup>2</sup> for 20 and 10 days,  
 355 respectively, at stations 1 and 2. These absences of PTR A (i.e., PTR A = 0 m<sup>2</sup>) were associated with  
 356 a significant increase in main river discharge (up to 128 m<sup>3</sup>/s) (e.g., end of July, early August, and  
 357 mid-September). Estimated areas in 2020 remained relatively constant ( $\approx$  400 m<sup>2</sup>) for an extended  
 358 period at Station 2 for mainstem discharge below 50 m<sup>3</sup>/s (e.g., June, July & Mid/Late August).  
 359 While at Station 1, the PTR A expressed greater variability, with a significant drop observed for  
 360 discharge values as low as 25 m<sup>3</sup>/s (e.g., July 11<sup>th</sup>, Mid August).

361  
362

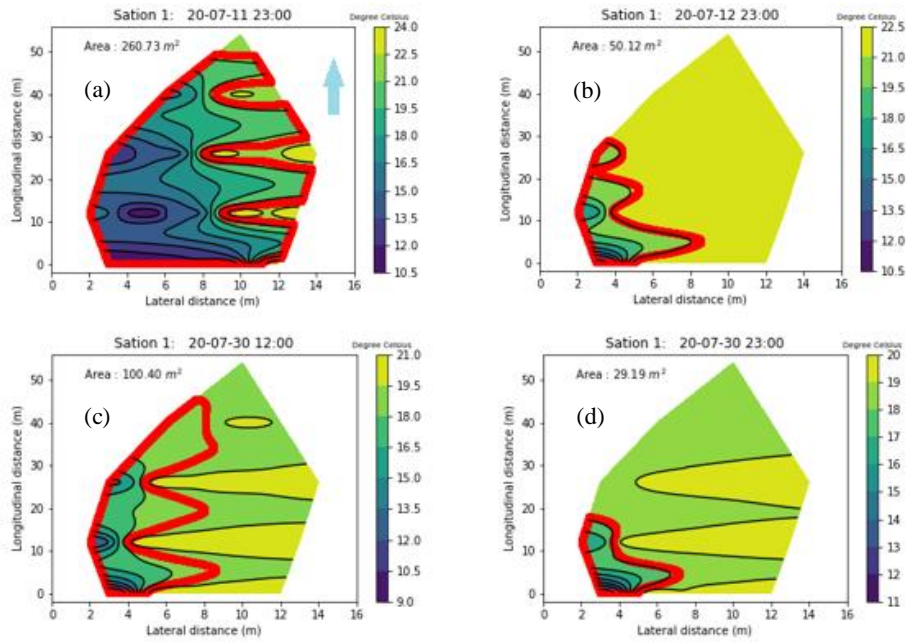


363

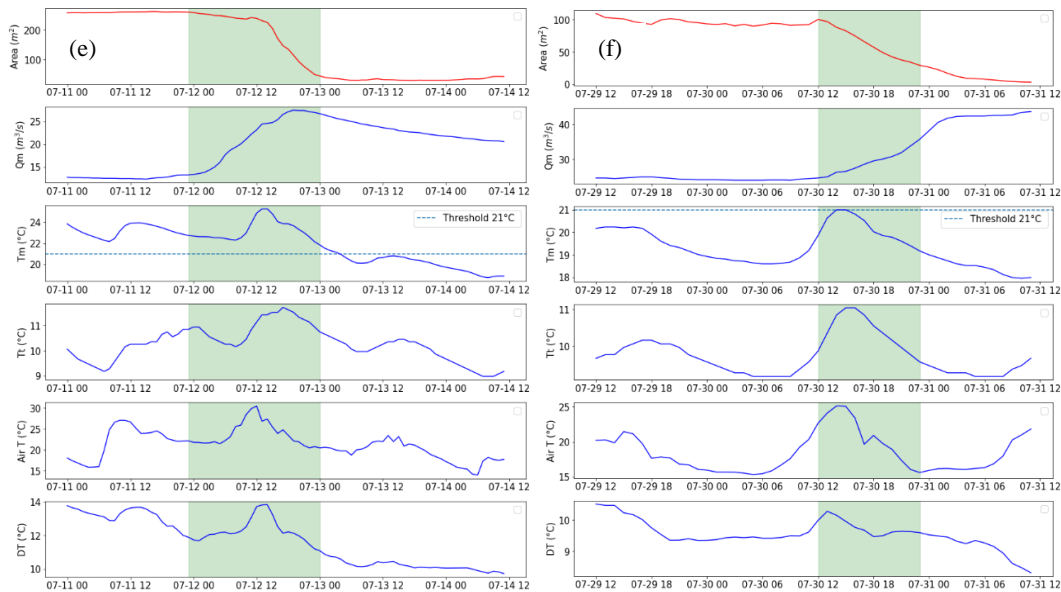
364 **Figure 2:** The time series of potential predictors and PTRA at Station 1 (top left & top right) and  
365 Station 2 (bottom left and bottom right) (2020-2021).

366 **Figure 3** illustrates the spatial delineation of PTRA at Station 1 (red) and shows the effect of  
367 discharge increase on the plume diel variability. During July 11-12<sup>th</sup>, 2020, the maximum  
368 mainstem water temperature reached 26 °C as the flow increased from 12.5 m<sup>3</sup>/s to 25 m<sup>3</sup>/s ( $T_m \geq$   
369 22 °C for more than ten consecutive days). Therefore, the thermal refuge area showed the most  
370 significant drop in 24 hours (from 260.7 m<sup>2</sup> to 50.1 m<sup>2</sup>). On July 30<sup>th</sup>, PTRA reached 29 m<sup>2</sup> and

371 flow values continued to increase (up to 40 m<sup>3</sup>/s). Observations throughout the summer of 2020  
 372 revealed that when mainstem discharge exceeded 45 m<sup>3</sup>/s (Station 1) and 60 m<sup>3</sup>/s (Station 2), the  
 373 cold-water plume area was reduced to zero, regardless of other predictors values (e.g., Station 1  
 374 from August 3<sup>rd</sup> to 8<sup>th</sup> and early/mid-September).



375



376

377 **Figure 3:** Interpolated PTRA and corresponding time series of potential predictors at Station 1  
 378 during July 11<sup>th</sup> & 12<sup>th</sup> 2020 at 11 p.m.(Figures (3a), (3b) & (3e)) and July 30<sup>th</sup> at noon & 11 p.m.  
 379 (Figures (3c), (3d) & (3f)). The blue arrow indicates the mainstem flow direction.

380

### 381 3.5.2- PTRA & low flow conditions in 2021

382 The hourly average discharge recorded in August 2021 (7.92 m<sup>3</sup>/s) is four times lower than  
383 that observed over the same period in 2020 (35.01 m<sup>3</sup>/s). Ste-Marguerite's summer water  
384 temperatures are typically higher during July than August, but that was not the case in 2021 (**Table**  
385 **2**).

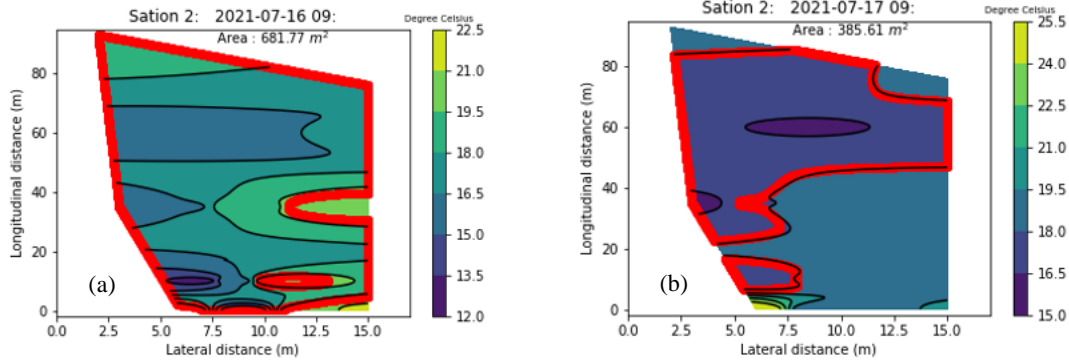
386 **Table 2:** Hourly mainstream water temperature (°C) of the Ste-Marguerite River during  
387 summertime' warmest months (July & August 2020-2021).

	2020			2021		
Months	Min	Max	Mean	Min	Max	Mean
July	16.63	26.01	20.81	14.78	24.33	19.43
August	10.88	24.46	18.68	14.63	27.31	21.35

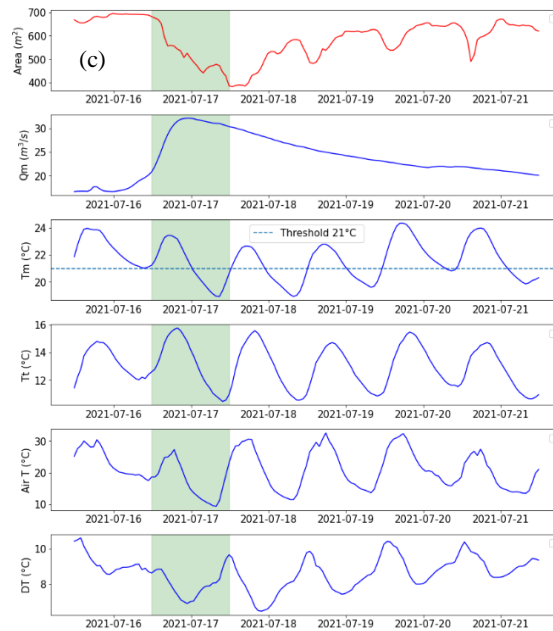
388

389 Therefore, diel PTRA variability across the summers of 2020 and 2021 was highly  
390 contrasted. From mid to end of August 2021, the mainstem water temperature reached 27 °C at  
391 both stations for more than ten consecutive days. Also, discharge reported historically low levels  
392 ( $Q_m \leq 10$  m<sup>3</sup>/s), raising conservation concerns about the survival of Atlantic salmon (*Salmo Salar*)  
393 in the Ste-Marguerite River. Quantiles and return periods were estimated from hourly discharge  
394 values between June and September from 1998 to 2021 using a one-week moving average. For a  
395 low-flow duration of 7 days, a return period of 43 years is computed for summer 2021 using the  
396 Weibull distribution and fitted with the maximum likelihood method.

397 During summer 2021, PTRA at Station 1 was reduced to 0 m<sup>2</sup> for only a few hours (June 28th),  
398 whereas, at Station 2, only high values were recorded for the whole sampling period (between 300  
399 m<sup>2</sup> and 700 m<sup>2</sup>). This was due to low flows in the main river and relatively constant cold water  
400 discharge from its tributaries. The largest mainstem discharge value is reported on July 17<sup>th</sup> (32  
401 m<sup>3</sup>/s). On this day, PTRA at Station 2 dropped from 681 m<sup>2</sup> to 385 m<sup>2</sup> (**Figure 4**). Subsequently,  
402 constant values are observed during August, and discharge drastically decreased to reach critical  
403 values ( $Q_m \leq 10$  m<sup>3</sup>/s).



404



405

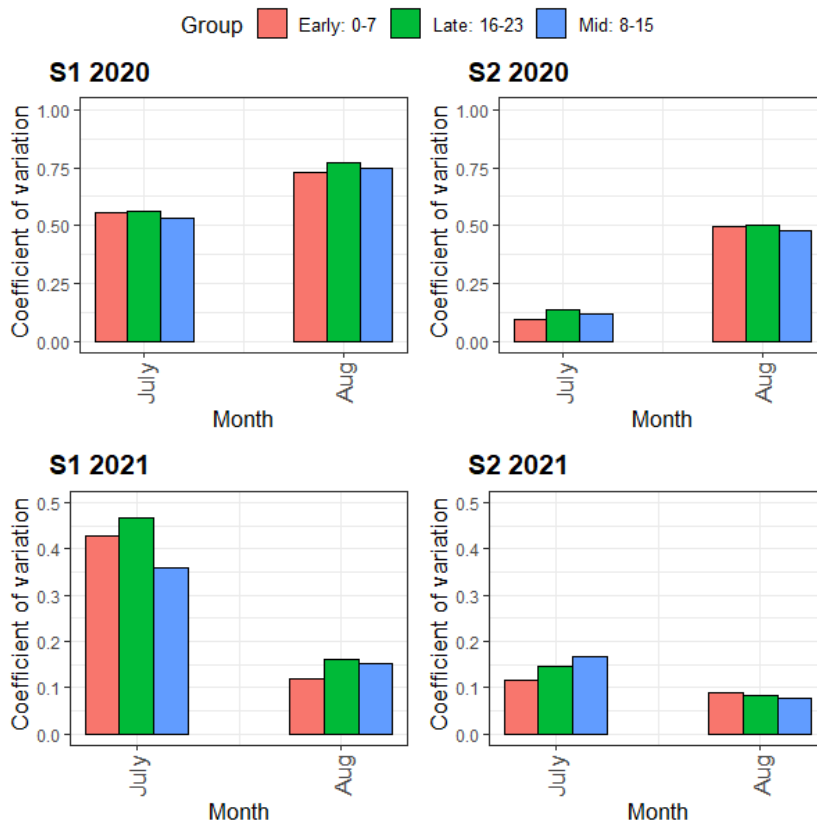
406 **Figure 4:** Interpolated PTRA at station 2 on the 16<sup>th</sup> July at 9 a.m. (4a) & on the 17<sup>th</sup> July at 9  
 407 a.m. (4b), and the corresponding time series of potential predictors (4c), in green the 24 hour  
 408 period between (4a) and (4b).

409 Those persistent cold-water plumes were due to low water levels (shallow riverbed at Station 2)  
 410 and increasing difference in water temperatures between tributaries and mainstem channels ( $T_{\Delta}$  up  
 411 to 17 °C at Station 1). The effects of the historic low flow are evident in the estimated PTRA, where  
 412 some of the edge/shore-located sensors were found to be out of the water as soon as early August  
 413 at both stations; hence recorded air temperatures instead. Consequently, when we removed these  
 414 sensors from the IWD interpolation procedure, the observed PTRA decreased. On August 3<sup>rd</sup>, the  
 415 thermal refuge area at Station 2 dropped from 700 to 570 m<sup>2</sup>, followed three days later by a drop  
 416 of 70 m<sup>2</sup> at Station 1 (From 190 to 120 m<sup>2</sup>).

417 **3.5.3- Diel and intra-seasonal variability**

418 *3.5.3.1- The change in CV under low flow conditions*

419 According to **Figure 5**, August 2020 showed the most variability in the coefficient of  
420 variation (CV), reaching values up to 75% & 55% for stations 1 & 2, respectively. In 2021, the Ste-  
421 Marguerite River recorded low flow rates and relatively high mainstem water temperature. These  
422 low flow conditions significantly impacted the PTRA's coefficient of variation in August 2021,  
423 where values plummeted to 13% and 10% for stations 1 and 2, respectively. The drop in variability  
424 across both summers during July was not as significant as in August. However, a slight decrease is  
425 observed in July's coefficient of variation in 2021. Station 1 had the greatest CV values for both  
426 summers, and the most variation occurred during the latter part of the day. Station 2 showed a shift  
427 in the peak period variability across both summers. Late afternoon in 2020 exhibited the highest  
428 CV values, then shifted to the middle of the day and early morning in 2021 for July and August,  
429 respectively.



430

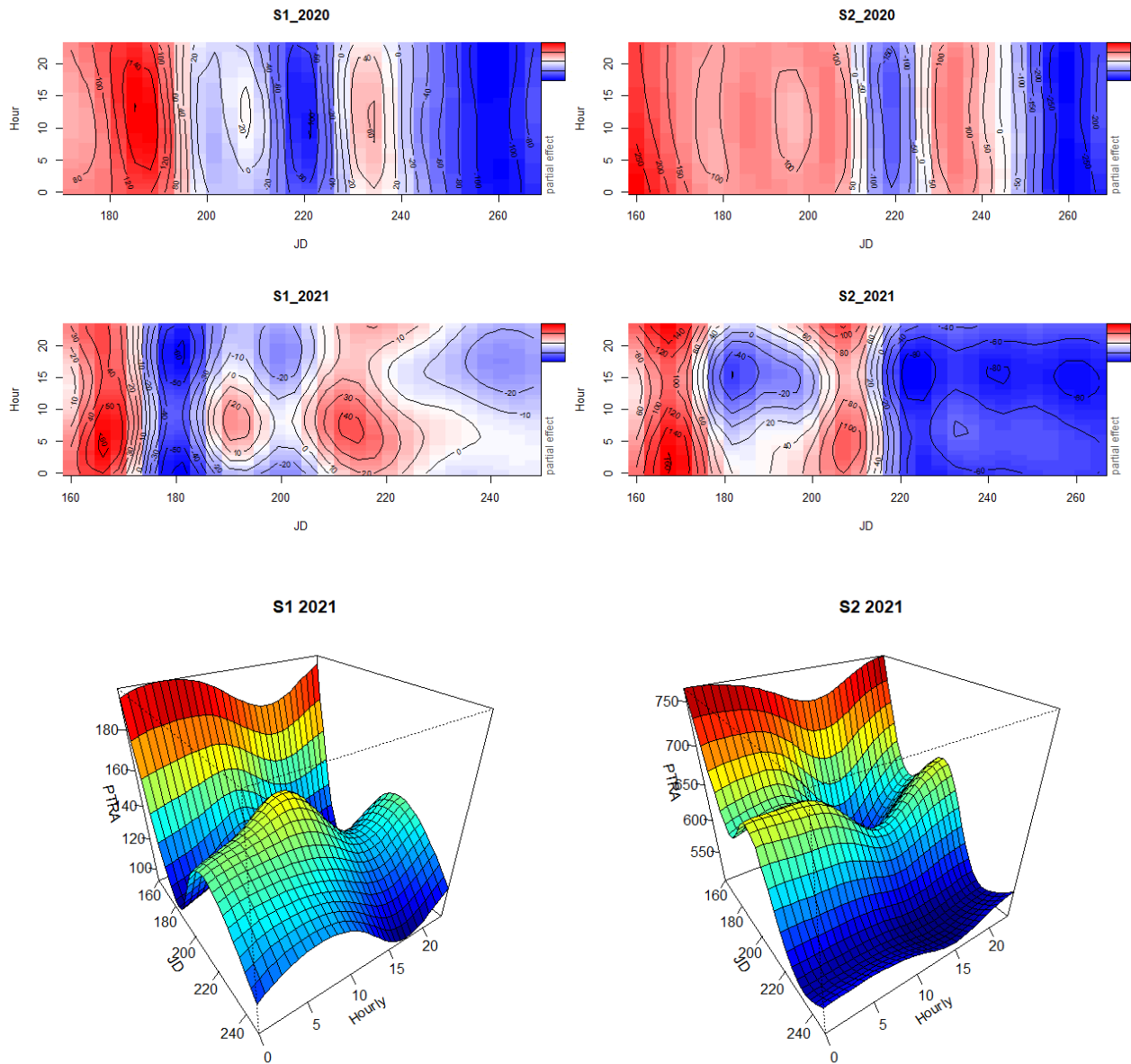
431 **Figure 5:** Coefficient variation for the warmest months (2020-2021) based on three different day  
432 periods. S1 and S2 refer to stations 1 & 2, respectively.

433

434 *3.5.3.2- GAM and PTRAs diel variability*

435           We computed interactions across Julian days, hours, and months using GAM's smoothers  
436 and tensor product smooth (Equations 5 & 6). The dark blue region in **Figure 6** (Contour line)  
437 refers to low values of PTRAs, while dark red identifies the high PTRAs values. GAM product  
438 smooth was shown to be relatively constant vertically in 2020, meaning there was no discernible  
439 difference in the time of occurrence of low and high peaks of PTRAs throughout the day.  
440 Nevertheless, PTRAs along Julian days (horizontally) was highly correlated with main river high  
441 discharge values (e.g., early August & September 2020).





442

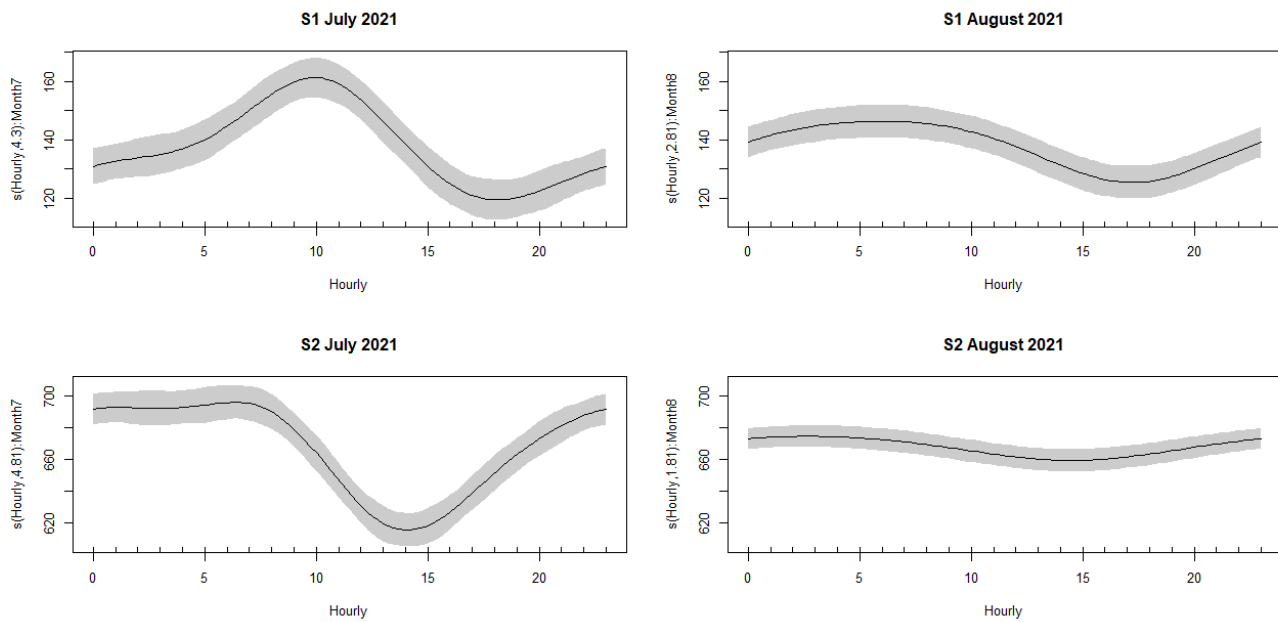
443

444 **Figure 6:** Contour (2020-2021) and perspective plot (2021) using GAM tensor product smoother  
 445 for stations 1 & 2 (S1 & S2): PTRA~Te (JD, Hourly). PTRA in m<sup>2</sup>, Julian days (JD) & the hour  
 446 components.

447

448 In 2021, GAM showed precise timing for PTRA extremes and significant diel variability. The area  
 449 at Station 1 typically increased from midnight to peak just before 10:00 a.m. (dark red and light  
 450 blue region), then gradually decreased to reach the lowest values between 4:00 p.m. and 7:00 p.m.  
 451 (dark blue and light red region). Station 2 exhibited the lowest PTRA values between 1:00 p.m.

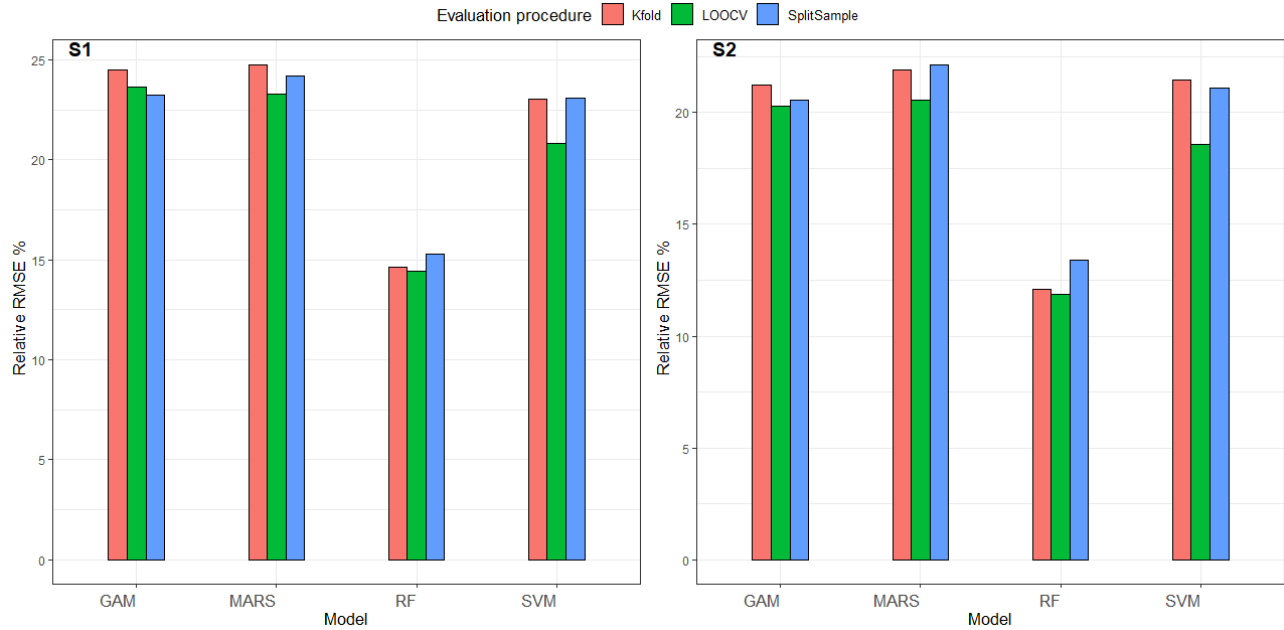
452 and 4:00 p.m., whereas the highest values were recorded late at night and early morning, just before  
 453 8:00 a.m. A time delay of about three hours is observed between PTRA extremes at the two  
 454 locations. **Figure 7** illustrates PTRA diel variability during July and August 2021. The greatest  
 455 relative diel variability was observed in July 2021, with values of  $\approx 40 \text{ m}^2$  &  $\approx 80 \text{ m}^2$  for stations  
 456 1 & 2, respectively. During this period, PTRA at Station 2 was near-constant for the early part of  
 457 the day (00:00-07:59), resulting in small fluctuations. PTRA extremes were captured more  
 458 accurately in Station 1, displaying a clear cyclic pattern and precise time slots.



459  
 460 **Figure 7:** PTRA ( $\text{m}^2$ ) and factor-smooth interaction using the summer's warmest months (2021)  
 461 fitted with GAM cyclic cubic regression splines for stations 1 & 2 (S1 & S2):  $\text{PTRA} \sim \text{Month} +$   
 462  $s(\text{Hourly}, \text{by}=\text{Month})$ .

### 463 3.5.4- Machine learning performance

464 LOOCV assessed the best performance based on rRMSE (**Figure 8**) across all the proposed  
 465 machine-learning models. As such, only this method is presented in the results (**Table 3**).



466

467 **Figure 8:** Evaluation procedures comparison based on rRMSE: K-fold vs. LooCV vs. Split  
 468 Sample, for stations 1 & 2 (S1 & S2)

469 **Table 3:** Best models' performance evaluated using the rRMSE, RMSE, Bias, rBias and NASH.

Model	Site	Predictors	rRMSE	RMSE (m <sup>2</sup> )	rBias	Bias (m <sup>2</sup> )	Nash
GAM	Station 1	<u>Qm, TΔ, Tt, Ta</u>	23.63%	27.02	0%	0	0.85
	Station 2	<u>Qm, TΔ, Tt, Ta</u>	20.27%	88.84	0%	0	0.78
	<b>Mean</b>	<b>NA</b>	<b>21.95%</b>	<b>57.93</b>	<b>0%</b>	<b>0</b>	<b>0.81</b>
MARS	Station 1	<u>Qm, TΔ, Tm, Ta</u>	23.27%	26.65	0%	0.001	0.86
	Station 2	<u>Qm, TΔ, Tm</u>	20.57%	90.16	0%	0	0.78
	<b>Mean</b>	<b>NA</b>	<b>21.92%</b>	<b>58.4</b>	<b>0%</b>	<b>0</b>	<b>0.82</b>
SVM	Station 1	<u>Qm, TΔ, Tt, Ta</u>	20.83%	23.86	-1.8%	-2	0.89
	Station 2	<u>Qm, TΔ, Tt, Ta</u>	18.57%	81.39	-2.4%	-10.41	0.83
	<b>Mean</b>	<b>NA</b>	<b>19.70%</b>	<b>52.62</b>	<b>-2.1%</b>	<b>-6.2</b>	<b>0.86</b>
RF	Station 1	<u>Qm, TΔ, Ta, Tt</u>	14.41%	16.5	0%	-0.01	0.94
	Station 2	<u>Qm, TΔ, Ta, Tt</u>	11.88%	52.11	-0.2%	-0.68	0.92
	<b>Mean</b>	<b>NA</b>	<b>13.14%</b>	<b>34.30</b>	<b>-0.2%</b>	<b>-0.68</b>	<b>0.93</b>

470 The bias of the test error estimate will be low with the LOOCV procedure since it uses almost all  
471 the data to train (Alpaydin, 2020; Desai & Ouarda, 2021). Modeling hourly PTRA was done by  
472 pooling data from both summers (2020-2021). We compared the simulated and interpolated PTRA  
473 using the performance metrics described above in the *Material and Methods* section (Equations  
474 12-16).

#### 475 *3.5.4.1- Benchmark models (GAM & MARS)*

476 GAM and MARS resulted in Nash scores of 85% at Station 1 and 78% at Station 2. However,  
477 Station 2 exhibited the lowest rRMSE (20%). The performances of both models were similar, and  
478 neither showed systematic bias. GAM used all five predictors for both stations, whereas MARS  
479 used fewer predictors (four at Station 1 & three at Station 2). The results showed that these  
480 nonparametric techniques have significant potential for describing accurate estimates of PTRA and  
481 account for most of the variation. Based on rRMSE and Nash score, neither the benchmark machine  
482 learning models outperformed the kernel-based SVM nor the tree-based RF.

#### 483 *3.5.4.2- Support vector machine & Random forest regression*

484 SVM showed better Nash scores and rRMSEs across both stations than benchmark models GAM  
485 & MARS. However, it expressed a slightly higher rBias than the other machine learning models,  
486 with values between -1.8% and -2.1%. Stations 1 and 2 had rRMSEs of 20.83% and 18.57%,  
487 respectively, and Nash scores of 89% and 83%. SVM and RF employed identical predictor sets for  
488 both stations ( $Q_m$ ,  $T_t$ ,  $T_a$ ,  $T_\Delta$ ). Nonetheless, RF outperformed all proposed models in terms of  
489 rRMSE and Nash, with mean values of 93% and 13.14%, respectively. According to **Table 3**, the  
490 regression quality performed at Station 2 was superior to those obtained at Station 1. One  
491 explanation would be that Station 1 had fewer observations available (4438 hours) than Station 2  
492 (4846 hours), and Station 2's thermograph arrays covered a larger overall area.

## 493 **3.6- DISCUSSION AND CONCLUSION**

### 494 **3.6.1- PTRA characteristics and temporal variability**

495 High flow values are often associated with increased mixing of water masses and low  
496 mainstem temperature; consequently, the difference in water temperature between the mainstem  
497 and the tributary decreases ( $T_\Delta$ ). This could result in a more homogenous zone of water

498 temperatures at the confluences, thus preventing PTRA from forming. In 2020, Station 1 had the  
499 most consecutive days of PTRA absence, primarily attributable to the early August high flow  
500 period and the geomorphic peculiarities of this site. A steeper cross-river slope from the bank to  
501 the main river thalweg characterizes station 1. Therefore, mixing between the main river and  
502 tributary water is more likely than at Station 2, where the plume is somewhat separated from the  
503 river's mainstem as it sits on a shallow gravel/cobble bed. During summer 2020, interpolated areas  
504 at Station 2 remain relatively constant for discharge values inferior to  $50 \text{ m}^3/\text{s}$ ; however, half this  
505 amount ( $25 \text{ m}^3/\text{s}$ ) can lower PTRA at Station 1 by more than 50%.

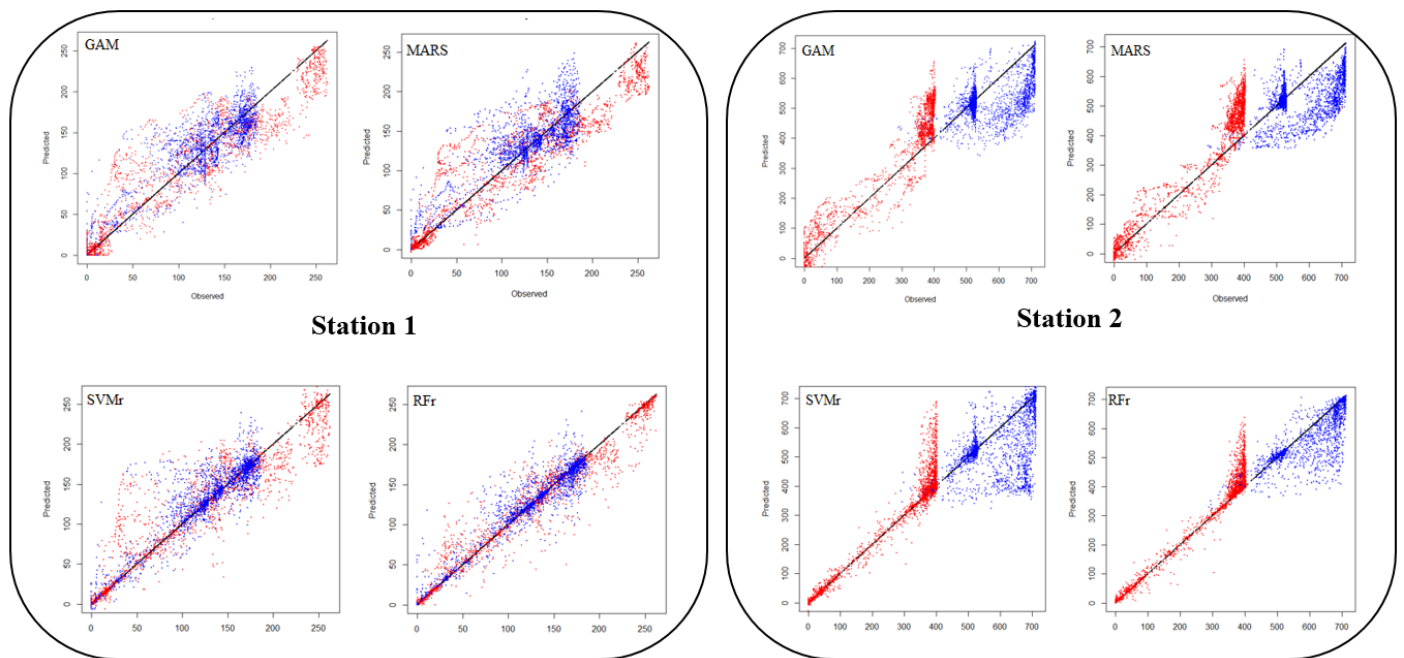
506 Main river low flow condition in 2021 is characterized by a return period of 43 years. According  
507 to Daigle et al., (2015), sites that do not currently function as thermal refuges but do supply cold  
508 water output are to become essential during hot weather. Continuously decreasing discharge values,  
509 low water levels, and high mainstem water temperature ( $T_m \geq 26 \text{ }^\circ\text{C}$ ) resulted in a similar PTRA  
510 behavior across the two stations. The estimated PTRA in 2021 showed specific time windows for  
511 daily extremes in July and a nearly constant gradient (low diel variability) in August (**Figure 6**).  
512 Additionally, several shore-located sensors were found out of the water starting in early August.  
513 Dugdale et al., (2013) suggested that lateral temperature heterogeneity (along the bank to the  
514 thalweg) in the mainstem river is likely to result from tributaries that supply a steady volume of  
515 cold water throughout the summer. In August 2021, PTRA diel variability decreased significantly  
516 (**Figure 6 & Figure 7**) after these sensors were removed from the interpolation procedure, resulting  
517 in relatively persistent PTRA. Nevertheless, the extent of PTRA was well captured at both stations  
518 and was significantly greater at Station 2. Still, the potential thermal refuge volume (PTRV in  $\text{m}^3$ )  
519 may be more significant at Station 1 due to its series of steep step pools. In addition to the cross-  
520 river slope, Benda et al., (2004) found that a river network's local geometry affects confluence  
521 morphology. More specifically, the angle at which tributaries enter the mainstem, where in most  
522 situations, the tributary-mainstem intersection angle is acute (less than  $90^\circ$ ), but as the angle  
523 approaches  $90^\circ$  (e.g., Station 1); the extent of the PTRA at the confluence is strongly influenced  
524 by the tributary's effects on riverbed morphology.

525 In some cases, the array of thermographs may have been insufficient to cover the complete cold-  
526 water plume. Therefore, spatial resolution and areal coverage should be increased to capture the  
527 full extent of the plume using thermal infrared aerial imagery (TIR) or/and distributed temperature

528 sensors (DTS) (Dzara et al., 2019). As part of an ecosystem assessment, thermal refuges should be  
529 evaluated according to the complex spatial (riverbed slope, depth, tributary-mainstem angle) and  
530 temporal patterns (diel and intra-seasonal variability), as well as considering other factors, such as  
531 food availability (Brewitt et al., 2017).

### 532 3.6.2- Machine learning performance

533 GAM was the least parsimonious model among the four suggested models and used all five  
534 potential predictors. The MARS model did not include the tributary temperature ( $T_t$ ), while SVM  
535 and RF did not include the mainstem temperature ( $T_m$ ). A part of the information contained in  $T_t$   
536 and  $T_m$  is included in  $T_\Delta$ , which is kept in the final best-performing models. MARS is not  
537 considering air temperature ( $T_a$ ) as a predictor at Station 2, being the farthest from the  
538 meteorological station where  $T_a$  is recorded. Observed versus predicted plots are presented in  
539 **Figure 9** for stations 1 and 2, where red and blue dots represent data for 2020 and 2021,  
540 respectively.



541  
542 **Figure 9:** Observed vs. Predicted PTRA ( $m^2$ ). Red and blue dots represent hourly PTRA  
543 estimates for summer 2020 and 2021, respectively.

544 At Station 1, no discernible difference was observed between the two summers. However, the  
545 effects of low discharge and high mainstem water temperature are more pronounced at Station 2,  
546 where two distinct populations are easily distinguished (2020 & 2021). The models overestimated

547 the PTRA at Station 2 during June and July 2020 (red dots  $\approx 390 \text{ m}^2$ ) and underestimated it in  
548 August 2021 ( $500 \text{ m}^2 \leq$  blue dots  $\leq 700 \text{ m}^2$ ). For those high estimated values ( $\geq 390 \text{ m}^2$ ), prediction  
549 performance was improved using SVM and RF regressions, as points got closer to the diagonal  
550 line than those obtained using GAM & MARS. The hourly models had overall good performance  
551 compared to those built on a daily time step by Saadi et al. (2021) and sometimes performed even  
552 better. In addition, the results showed that both stations contain potential thermal refuges, and their  
553 area can be estimated with rRMSE of less than 15%. Furthermore, by using smoothers and product-  
554 smooth GAM tools, it is possible to gain a better understanding of the intra-seasonal and diel  
555 behaviour of PTRA and incorporate it into salmon river fisheries management plans, particularly  
556 during low flow events when mainstem temperatures approach the lethal upper limits temperature  
557 of Atlantic salmon.

558 Nonparametric approaches have helped describe hourly PTRA estimates, and the excellent  
559 performance of RF makes other variants of tree-based models worth the investigation. For example,  
560 the extremely randomized trees (ERT), extreme gradient boosting (XGBoost), and M5Tree  
561 demonstrated outstanding performances in water temperature modeling (Feigl et al., 2021; Heddam  
562 et al., 2020). We did not include neural network models, which should also be considered for  
563 comparative studies where long datasets are available due to their outstanding performance in long-  
564 lead-time prediction in recent hydrological applications (Qiu et al., 2021; L. Wang et al., 2022).  
565 Future studies should also explore the possibility of constructing multi-site models at a regional  
566 scale and investigate the diel variability in a climate change context. Moreover, including  
567 forecasted hydrometeorological variables, morphological features (cross-river slope & angle of  
568 entry), and tributary discharge as additional predictors may improve the models' performances.  
569 Given these points, assessing hourly PTRA estimates and investigating the diel variability using  
570 nondeterministic approaches may be beneficial to improve strategies for the survival of cold-water  
571 taxa.

572 Finally, although this study was primarily motivated by the fact that thermal refuges are key  
573 habitats for ichthyofauna, the models developed herein may have other applications. The  
574 implications of PTRA changes, could affect the survival rate and the distribution of Atlantic salmon  
575 during summer in-river residence in the Sainte-Marguerite River. Salmon migration may be  
576 physically impeded by low river discharge and hot water temperatures, making it difficult for them

577 to reach adequate refuges. Man-made thermal plumes exist in many river and coastal environments  
578 and the most common case is likely the discharge of warm-water effluent from a power plant (e.g.,  
579 Penk & Williams, 2019; Zavarisky & Duester, 2020). Nevertheless, other applications could include  
580 other natural phenomena, such as rivers emptying into lakes (Smith & Simpkins, 2018), tributaries  
581 affecting thermal regimes of lakes (Råman Vinnå et al., 2018) or freshwater plumes in estuaries  
582 (Huang et al., 2020).

### 583 **Acknowledgments**

584 This work was partly funded by the Natural Sciences and Engineering Research Council (NSERC).  
585 The authors would like to thank André Boivin for his assistance during the fieldwork related to this  
586 study, and the Ministry of Sustainable Development, Environment, and Fight Against Climate  
587 Change of the Province of Quebec (MDDELCC) for the employed datasets.

### 588 **Credit author statement**

589 **Ilias Hani:** Conceptualization, Methodology, Formal analysis, Software, Investigation, Writing -  
590 Original Draft, Field-related work.

591 **André St-Hilaire:** Conceptualization, Methodology, Formal analysis, Supervision, Writing -  
592 Review & Editing, Field-related work.

593 **Taha B.M.J. Ouarda:** Conceptualization, Methodology, Formal analysis, Review & Editing.

### 594 **Declaration of Competing Interest**

595 The authors declare that they have no known competing financial interests or personal relationships that  
596 could have appeared to influence the work reported in this paper.

597

598

599

600

601

602



603

604

605

#### 606 **4- BIBLIOGRAPHIE**

607 Adamowski, K., & Labatiuk, C. (1987). Estimation of flood frequencies by a nonparametric  
608 density procedure. In *Hydrologic Frequency Modeling* (pp. 97–106). Springer.

609 Allahbakhshian-Farsani, P., Vafakhah, M., Khosravi-Farsani, H., & Hertig, E. (2020). Regional  
610 flood frequency analysis through some machine learning models in semi-arid regions.  
611 *Water Resources Management*, 34(9), 2887–2909.

612 Alpaydin, E. (2020). *Introduction to machine learning*. MIT press.

613 Armstrong, J. B., Ward, E. J., Schindler, D. E., & Lisi, P. J. (2016). Adaptive capacity at the  
614 northern front: Sockeye salmon behaviourally thermoregulate during novel exposure to  
615 warm temperatures. *Conservation Physiology*, 4(1), cow039.

616 Baird, O. E., & Krueger, C. C. (2003). Behavioral thermoregulation of brook and rainbow trout:  
617 Comparison of summer habitat use in an Adirondack River, New York. *Transactions of*  
618 *the American Fisheries Society*, 132(6), 1194–1206.

619 Bardach, J. E., & Bjorklund, R. G. (1957). The temperature sensitivity of some American  
620 freshwater fishes. *The American Naturalist*, 91(859), 233–251.

621 Barnett, T. P., Adam, J. C., & Lettenmaier, D. P. (2005). Potential impacts of a warming climate  
622 on water availability in snow-dominated regions. *Nature*, 438(7066), 303–309.

623 Benda, L., Andras, K., Miller, D., & Bigelow, P. (2004). Confluence effects in rivers:  
624 Interactions of basin scale, network geometry, and disturbance regimes. *Water Resources*  
625 *Research*, 40(5).

626 Benyahya, L., Caissie, D., St-Hilaire, A., Ouarda, T. B., & Bobée, B. (2007). A review of  
627 statistical water temperature models. *Canadian Water Resources Journal*, 32(3), 179–  
628 192.

629 Biron, P. M., Ramamurthy, A. S., & Han, S. (2004). Three-dimensional numerical modeling of  
630 mixing at river confluences. *Journal of Hydraulic Engineering*, 130(3), 243–253.

631 Boucher, J.-P., Côté, S., & Guillen, M. (2017). Exposure as duration and distance in telematics  
632 motor insurance using generalized additive models. *Risks*, 5(4), 54.

633 Breau, C., & Caissie, D. (2013). *Adaptive management strategies to protect Atlantic salmon*  
634 *(Salmo salar) under environmentally stressful conditions*. Canadian Science Advisory  
635 Secretariat= Secrétariat canadien de consultation ....

636 Breau, C., Cunjak, R. A., & Bremset, G. (2007). Age-specific aggregation of wild juvenile  
637 Atlantic salmon *Salmo salar* at cool water sources during high temperature events.  
638 *Journal of Fish Biology*, 71(4), 1179–1191.

639 Breau, C., Cunjak, R. A., & Peake, S. J. (2011). Behaviour during elevated water temperatures:  
640 Can physiology explain movement of juvenile Atlantic salmon to cool water? *Journal of*  
641 *Animal Ecology*, 80(4), 844–853.

642 Breiman, L. (1996). Bagging predictors. *Machine Learning*, 24(2), 123–140.

643 Breiman, L. (2001). Random Forests. *Machine Learning*, 45(1), 5–32.  
644 <https://doi.org/10.1023/A:1010933404324>

645 Brewitt, K. S., & Danner, E. M. (2014). Spatio-temporal temperature variation influences  
646 juvenile steelhead (*Oncorhynchus mykiss*) use of thermal refuges. *Ecosphere*, 5(7), 1–26.

647 Brewitt, K. S., Danner, E. M., & Moore, J. W. (2017). Hot eats and cool creeks: Juvenile Pacific  
648 salmonids use mainstem prey while in thermal refuges. *Canadian Journal of Fisheries*  
649 *and Aquatic Sciences*, 74(10), 1588–1602.

650 Caissie, D. (2006). The thermal regime of rivers: A review. *Freshwater Biology*, 51(8), 1389–  
651 1406.

652 Caissie, D., El-Jabi, N., & Satish, M. G. (2001). Modelling of maximum daily water temperatures  
653 in a small stream using air temperatures. *Journal of Hydrology*, 251(1–2), 14–28.

654 Caissie, D., St-Hilaire, A., & El-Jabi, N. (2004). Prediction of water temperatures using  
655 regression and stochastic models. *57th Canadian Water Resources Association Annual  
656 Congress., Montreal*, 6.

657 Chenard, J.-F., & Caissie, D. (2008). Stream temperature modelling using artificial neural  
658 networks: Application on Catamaran Brook, New Brunswick, Canada. *Hydrological  
659 Processes: An International Journal*, 22(17), 3361–3372.

660 Christensen, N. S., & Lettenmaier, D. P. (2007). A multimodel ensemble approach to assessment  
661 of climate change impacts on the hydrology and water resources of the Colorado River  
662 Basin. *Hydrology and Earth System Sciences*, 11(4), 1417–1434.

663 Conoscenti, C., Ciaccio, M., Caraballo-Arias, N. A., Gómez-Gutiérrez, Á., Rotigliano, E., &  
664 Agnesi, V. (2015). Assessment of susceptibility to earth-flow landslide using logistic  
665 regression and multivariate adaptive regression splines: A case of the Belice River basin  
666 (western Sicily, Italy). *Geomorphology*, 242, 49–64.

667 Conoscenti, C., Rotigliano, E., Cama, M., Caraballo-Arias, N. A., Lombardo, L., & Agnesi, V.  
668 (2016). Exploring the effect of absence selection on landslide susceptibility models: A  
669 case study in Sicily, Italy. *Geomorphology*, 261, 222–235.

670 Corey, E., Linnansaari, T., Dugdale, S. J., Bergeron, N., Gendron, J.-F., Lapointe, M., & Cunjak,  
671 R. A. (2020). Comparing the behavioural thermoregulation response to heat stress by  
672 Atlantic salmon parr (*Salmo salar*) in two rivers. *Ecology of Freshwater Fish*, 29(1), 50–  
673 62.

674 Daigle, A., Jeong, D. I., & Lapointe, M. F. (2015). Climate change and resilience of tributary  
675 thermal refugia for salmonids in eastern Canadian rivers. *Hydrological Sciences Journal*,  
676 60(6), 1044–1063.

677 Daigle, A., Ouarda, T. B., & Bilodeau, L. (2010). Comparison of parametric and non-parametric  
678 estimations of the annual date of positive water temperature onset. *Journal of Hydrology*,  
679 390(1–2), 75–84.

680 Davis, J., Pavlova, A., Thompson, R., & Sunnucks, P. (2013). Evolutionary refugia and  
681 ecological refuges: Key concepts for conserving Australian arid zone freshwater  
682 biodiversity under climate change. *Global Change Biology*, 19(7), 1970–1984.

683 Deka, P. C. (2014). Support vector machine applications in the field of hydrology: A review.  
684 *Applied Soft Computing*, 19, 372–386.

685 Desai, S., & Ouarda, T. B. (2021). Regional hydrological frequency analysis at ungauged sites  
686 with random forest regression. *Journal of Hydrology*, 594, 125861.

687 Dominici, F., McDermott, A., Zeger, S. L., & Samet, J. M. (2002). On the use of generalized  
688 additive models in time-series studies of air pollution and health. *American Journal of*  
689 *Epidemiology*, 156(3), 193–203.

690 Drucker, H., Burges, C. J., Kaufman, L., Smola, A., & Vapnik, V. (1996). Support vector  
691 regression machines. *Advances in Neural Information Processing Systems*, 9.

692 Dugdale, S. J., Bergeron, N. E., & St-Hilaire, A. (2013). Temporal variability of thermal refuges  
693 and water temperature patterns in an Atlantic salmon river. *Remote Sensing of*  
694 *Environment*, 136, 358–373.

695 Dugdale, S. J., Bergeron, N. E., & St-Hilaire, A. (2015). Spatial distribution of thermal refuges  
696 analysed in relation to riverscape hydromorphology using airborne thermal infrared  
697 imagery. *Remote Sensing of Environment*, 160, 43–55.

698 Dugdale, S. J., Curry, R. A., St-Hilaire, A., & Andrews, S. N. (2018). Impact of future climate  
699 change on water temperature and thermal habitat for keystone fishes in the lower Saint  
700 John River, Canada. *Water Resources Management*, 32(15), 4853–4878.

701 Dzara, J. R., Neilson, B. T., & Null, S. E. (2019). Quantifying thermal refugia connectivity by  
702 combining temperature modeling, distributed temperature sensing, and thermal infrared  
703 imaging. *Hydrology and Earth System Sciences*, 23(7), 2965–2982.

704 Ebersole, J. L., Liss, W. J., & Frissell, C. A. (2001). Relationship between stream temperature,  
705 thermal refugia and rainbow trout *Oncorhynchus mykiss* abundance in arid-land streams  
706 in the northwestern United States. *Ecology of Freshwater Fish*, 10(1), 1–10.

707 Ebersole, J. L., Wigington Jr, P. J., Leibowitz, S. G., Comeleo, R. L., & Sickie, J. V. (2015).  
708 Predicting the occurrence of cold-water patches at intermittent and ephemeral tributary  
709 confluences with warm rivers. *Freshwater Science*, 34(1), 111–124.

710 Efron, B. (1982). *The jackknife, the bootstrap and other resampling plans*. SIAM.

711 Elliott, J. (1991). Tolerance and resistance to thermal stress in juvenile Atlantic salmon, *Salmo*  
712 *salar*. *Freshwater Biology*, 25(1), 61–70.

713 Feigl, M., Lebedzinski, K., Herrnegger, M., & Schulz, K. (2021). Machine-learning methods for  
714 stream water temperature prediction. *Hydrology and Earth System Sciences*, 25(5), 2951–  
715 2977.

716 Ferchichi, H., St-Hilaire, A., Ouarda, T. B., & Lévesque, B. (2021). Impact of the future coastal  
717 water temperature scenarios on the risk of potential growth of pathogenic *Vibrio* marine  
718 bacteria. *Estuarine, Coastal and Shelf Science*, 250, 107094.

719 Frechette, D. M., Dugdale, S. J., Dodson, J. J., & Bergeron, N. E. (2018). Understanding  
720 summertime thermal refuge use by adult Atlantic salmon using remote sensing, river

721 temperature monitoring, and acoustic telemetry. *Canadian Journal of Fisheries and*  
722 *Aquatic Sciences*, 75(11), 1999–2010.

723 Frechette, D. M., St-Hilaire, A., & Bergeron, N. (2019). *Statistical analysis of fish ladder*  
724 *attractivity on the Nord-Est Sainte-Marguerite River*.

725 Friedman, J. H. (1991). *Estimating functions of mixed ordinal and categorical variables using*  
726 *adaptive splines*. Stanford Univ CA Lab for Computational Statistics.

727 Gardner, B., Sullivan, P. J., & Lembo, J., Arthur J. (2003). Predicting stream temperatures:  
728 Geostatistical model comparison using alternative distance metrics. *Canadian Journal of*  
729 *Fisheries and Aquatic Sciences*, 60(3), 344–351.

730 Gendron, J.-F. (2013). *Physical controls on summer thermal refuges for salmonids in two gravel-*  
731 *cobble salmon rivers with contrasting thermal regimes: The Ouelle and Ste. Marguerite*  
732 *rivers*.

733 Greer, G., Carlson, S., & Thompson, S. (2019). Evaluating definitions of salmonid thermal  
734 refugia using in situ measurements in the Eel River, Northern California. *Ecohydrology*,  
735 12(5), e2101.

736 Guillemette, N., St-Hilaire, A., Ouarda, T. B., & Bergeron, N. (2011). Statistical tools for thermal  
737 regime characterization at segment river scale: Case study of the Ste-Marguerite River.  
738 *River Research and Applications*, 27(8), 1058–1071.

739 Hastie, T., & Tibshirani, R. (1987). Generalized additive models: Some applications. *Journal of*  
740 *the American Statistical Association*, 82(398), 371–386.

741 Heddam, S., Ptak, M., & Zhu, S. (2020). Modelling of daily lake surface water temperature from  
742 air temperature: Extremely randomized trees (ERT) versus Air2Water, MARS, M5Tree,  
743 RF and MLPNN. *Journal of Hydrology*, 588, 125130.

744 Huang, W., Li, C., White, J. R., Bargu, S., Milan, B., & Bentley, S. (2020). Numerical  
745 experiments on variation of freshwater plume and leakage effect From Mississippi River  
746 Diversion in the Lake Pontchartrain Estuary. *Journal of Geophysical Research: Oceans*,  
747 *125*(2), e2019JC015282.

748 Isaak, D. J., & Rieman, B. E. (2013). Stream isotherm shifts from climate change and  
749 implications for distributions of ectothermic organisms. *Global Change Biology*, *19*(3),  
750 742–751.

751 Isaak, D. J., Young, M. K., Nagel, D. E., Horan, D. L., & Groce, M. C. (2015). The cold-water  
752 climate shield: Delineating refugia for preserving salmonid fishes through the 21st  
753 century. *Global Change Biology*, *21*(7), 2540–2553.

754 Jeong, D. I., Daigle, A., & St-Hilaire, A. (2013). Development of a stochastic water temperature  
755 model and projection of future water temperature and extreme events in the Ouelle River  
756 basin in Québec, Canada. *River Research and Applications*, *29*(7), 805–821.

757 Kuhn, M., & Johnson, K. (2013). *Applied predictive modeling* (Vol. 26). Springer.

758 Kuhn, M., Wing, J., Weston, S., Williams, A., Keefer, C., Engelhardt, A., Cooper, T., Mayer, Z.,  
759 Kenkel, B., & Team, R. C. (2020). Package ‘caret.’ *The R Journal*, *22*(3), 7.

760 Kurylyk, B. L., MacQuarrie, K. T., Linnansaari, T., Cunjak, R. A., & Curry, R. A. (2015).  
761 Preserving, augmenting, and creating cold-water thermal refugia in rivers: Concepts  
762 derived from research on the Miramichi River, New Brunswick (Canada). *Ecohydrology*,  
763 *8*(6), 1095–1108.

764 Laanaya, F., St-Hilaire, A., & Gloaguen, E. (2017). Water temperature modelling: Comparison  
765 between the generalized additive model, logistic, residuals regression and linear  
766 regression models. *Hydrological Sciences Journal*, *62*(7), 1078–1093.

767 Lee, J., Bahri, Y., Novak, R., Schoenholz, S. S., Pennington, J., & Sohl-Dickstein, J. (2017).  
768 Deep neural networks as gaussian processes. *ArXiv Preprint ArXiv:1711.00165*.

769 Lund, S. G., Caissie, D., Cunjak, R. A., Vijayan, M. M., & Tufts, B. L. (2002). The effects of  
770 environmental heat stress on heat-shock mRNA and protein expression in Miramichi  
771 Atlantic salmon (*Salmo salar*) parr. *Canadian Journal of Fisheries and Aquatic Sciences*,  
772 59(9), 1553–1562.

773 Mahardja, B., Bashevkin, S., Pien, C., Nelson, M., Davis, B., & Hartman, R. (2021). *Escape*  
774 *From the Heat: Thermal Stratification in a Well-Mixed Estuary and Implications for Fish*  
775 *Species Facing a Changing Climate*.

776 Marra, G., & Wood, S. N. (2011). Practical variable selection for generalized additive models.  
777 *Computational Statistics & Data Analysis*, 55(7), 2372–2387.

778 Marteau, B., Piégay, H., Chandesris, A., Michel, K., & Vaudor, L. (2022). Riparian shading  
779 mitigates warming but cannot revert thermal alteration by impoundments in lowland  
780 rivers. *Earth Surface Processes and Landforms*.

781 McCullagh, P., & Nelder, J. A. (1989). *Generalized Linear Models II*. Chapman and Hall,  
782 London.

783 Meyer, D., & Wien, F. T. (2001). Support vector machines. *R News*, 1(3), 23–26.

784 Milborrow, S. (2018). *Earth: Multivariate Adaptive Regression Splines derived from mda: Mars*  
785 *by T. Hastie and R. Tibshirani*.

786 Moatar, F., & Gailhard, J. (2006). Water temperature behaviour in the River Loire since 1976 and  
787 1881. *Comptes Rendus Geoscience*, 338(5), 319–328.

788 Morrill, J. C., Bales, R. C., & Conklin, M. H. (2005). Estimating stream temperature from air  
789 temperature: Implications for future water quality. *Journal of Environmental Engineering*,  
790 131(1), 139–146.



791 Msilini, A., Masselot, P., & Ouarda, T. B. (2020). Regional Frequency Analysis at Ungauged  
792 Sites with Multivariate Adaptive Regression Splines. *Journal of Hydrometeorology*,  
793 21(12), 2777–2792.

794 Murray, R. W. (1971). 5 Temperature Receptors. In *Fish physiology* (Vol. 5, pp. 121–133).  
795 Elsevier.

796 Nuhfer, A. J., Zorn, T. G., & Wills, T. C. (2017). Effects of reduced summer flows on the brook  
797 trout population and temperatures of a groundwater-influenced stream. *Ecology of*  
798 *Freshwater Fish*, 26(1), 108–119.

799 Penk, M. R., & Williams, M. A. (2019). Thermal effluents from power plants boost performance  
800 of the invasive clam *Corbicula fluminea* in Ireland’s largest river. *Science of the Total*  
801 *Environment*, 693, 133546.

802 Poole, G. C., & Berman, C. H. (2001). An ecological perspective on in-stream temperature:  
803 Natural heat dynamics and mechanisms of human-caused thermal degradation.  
804 *Environmental Management*, 27(6), 787–802.

805 Qiu, R., Wang, Y., Rhoads, B., Wang, D., Qiu, W., Tao, Y., & Wu, J. (2021). River water  
806 temperature forecasting using a deep learning method. *Journal of Hydrology*, 595,  
807 126016.

808 Quan, Q., Hao, Z., Xifeng, H., & Jingchun, L. (2020). Research on water temperature prediction  
809 based on improved support vector regression. *Neural Computing and Applications*, 1–10.

810 Råman Vinnå, L., Wüest, A., Zappa, M., Fink, G., & Bouffard, D. (2018). Tributaries affect the  
811 thermal response of lakes to climate change. *Hydrology and Earth System Sciences*, 22(1),  
812 31–51.

813 Rhoads, B. L., & Sukhodolov, A. N. (2001). Field investigation of three-dimensional flow  
814 structure at stream confluences: 1. Thermal mixing and time-averaged velocities. *Water*  
815 *Resources Research*, 37(9), 2393–2410.

816 Roy, S. S., Roy, R., & Balas, V. E. (2018). Estimating heating load in buildings using  
817 multivariate adaptive regression splines, extreme learning machine, a hybrid model of  
818 MARS and ELM. *Renewable and Sustainable Energy Reviews*, 82, 4256–4268.

819 Saadi, A. M., Msilini, A., Charron, C., St-Hilaire, A., & Ouarda, T. B. (2021). Estimation of the  
820 area of potential thermal refuges using generalized additive models and multivariate  
821 adaptive regression splines: A case study from the Ste-Marguerite River. *River Research*  
822 *and Applications*.

823 Shen, H., Tolson, B. A., & Mai, J. (2022). Time to Update the Split-Sample Approach in  
824 Hydrological Model Calibration. *Water Resources Research*, 58(3), e2021WR031523.

825 Shepard, S. L. (1995). *Atlantic salmon spawning migrations in the Penobscot River, Maine:*  
826 *Fishways, flows and high temperatures*.

827 Smith, B. J., & Simpkins, D. G. (2018). Influence of river plumes on the distribution and  
828 composition of nearshore Lake Michigan fishes. *Journal of Great Lakes Research*, 44(6),  
829 1351–1361.

830 Smits, G. F., & Jordaan, E. M. (2002). *Improved SVM regression using mixtures of kernels*. 3,  
831 2785–2790.

832 St-Hilaire, A., Ouarda, T. B., Bargaoui, Z., Daigle, A., & Bilodeau, L. (2012). Daily river water  
833 temperature forecast model with ak-nearest neighbour approach. *Hydrological Processes*,  
834 26(9), 1302–1310.

835 Sullivan, C. J., Vokoun, J. C., Helton, A. M., Briggs, M. A., & Kurylyk, B. L. (2021). An  
836 ecohydrological typology for thermal refuges in streams and rivers. *Ecohydrology*, e2295.

837 Sutton, R. J., Deas, M. L., Tanaka, S. K., Soto, T., & Corum, R. A. (2007). Salmonid  
838 observations at a Klamath River thermal refuge under various hydrological and  
839 meteorological conditions. *River Research and Applications*, 23(7), 775–785.

840 Torgersen, C. E., Ebersole, J. L., & Keenan, D. M. (2012). *Primer for identifying cold-water*  
841 *refuges to protect and restore thermal diversity in riverine landscapes.*

842 Torgersen, C. E., Price, D. M., Li, H. W., & McIntosh, B. A. (1999). Multiscale thermal refugia  
843 and stream habitat associations of chinook salmon in northeastern Oregon. *Ecological*  
844 *Applications*, 9(1), 301–319.

845 van Vliet, M. T., Franssen, W. H., Yearsley, J. R., Ludwig, F., Haddeland, I., Lettenmaier, D. P.,  
846 & Kabat, P. (2013). Global river discharge and water temperature under climate change.  
847 *Global Environmental Change*, 23(2), 450–464.

848 Vapnik, V. (1998). The support vector method of function estimation. In *Nonlinear modeling* (pp.  
849 55–85). Springer.

850 Wang, L., Xu, B., Zhang, C., Fu, G., Chen, X., Zheng, Y., & Zhang, J. (2022). Surface water  
851 temperature prediction in large-deep reservoirs using a long short-term memory model.  
852 *Ecological Indicators*, 134, 108491.

853 Wang, T., Kelson, S. J., Greer, G., Thompson, S. E., & Carlson, S. M. (2020). Tributary  
854 confluences are dynamic thermal refuges for a juvenile salmonid in a warming river  
855 network. *River Research and Applications*, 36(7), 1076–1086.

856 Webb, R. M., Peters, N. J., Aulenbach, B. T., & Shanley, J. B. (2003). *Relations between*  
857 *hydrology and solute fluxes at the five water, energy, and biogeochemical budget (WEBB)*  
858 *watersheds of the United States Geological Survey.* 332–339.

859 Weierbach, H., Lima, A. R., Christianson, D., Faybishenko, B., Hendrix, V., & Varadharajan, C.  
860 (n.d.). *A Comparison of Data-Driven Models for Predicting Stream Water Temperature.*

861 Wickham, H., & Grolemund, G. (2016). *R for data science: Import, tidy, transform, visualize,*  
862 *and model data.* O'Reilly Media, Inc.

863 Wilbur, N. M., O'Sullivan, A. M., MacQuarrie, K. T., Linnansaari, T., & Curry, R. A. (2020).  
864 Characterizing physical habitat preferences and thermal refuge occupancy of brook trout  
865 (*Salvelinus fontinalis*) and Atlantic salmon (*Salmo salar*) at high river temperatures. *River*  
866 *Research and Applications*, 36(5), 769–783.

867 Wood, S. N. (2006). Low-rank scale-invariant tensor product smooths for generalized additive  
868 mixed models. *Biometrics*, 62(4), 1025–1036.

869 Zavarisky, A., & Duester, L. (2020). Anthropogenic influence on the Rhine water temperatures.  
870 *Hydrology and Earth System Sciences*, 24(10), 5027–5041.

871 Zhang, Y., Ma, J., Liang, S., Li, X., & Li, M. (2020). An evaluation of eight machine learning  
872 regression algorithms for forest aboveground biomass estimation from multiple satellite  
873 data products. *Remote Sensing*, 12(24), 4015.

874