

# SAMUEL JOHNSTON

johnston.samuelj@gmail.com · <https://sjjohnst.github.io>  
Montréal, H2S 3G5

*Je suis ingénieur en données géospatiales spécialisé dans le développement de pipelines de données d'observation de la Terre (Earth Observation), pour les analyses et l'apprentissage automatique. Mon expérience pratique en tant que scientifique des données m'a appris la valeur de la qualité, de la fiabilité et de la scalabilité, que je mets à profit dans mon travail d'ingénieur.*

## PROFESSIONAL EXPERIENCE

### Scientifique/Ingénieur de données *H2O Geomatics*

Septembre 2024 - présent  
*Kitchener, ON*

- Optimisation du pipeline de données EO qui traite 1.5 PB de données de réflectance MODIS, permettant ainsi un gain de vitesse de 12 fois et une réduction de 16 fois de la consommation de mémoire.
- Développement d'un pipeline de données pour la production d'échantillons de formation et de test de haute qualité pour le développement d'algorithmes pour supprimer la couverture nuageuse, intégrant de la qualité et les tests unitaires, et s'adaptant à des flux de travail parallèles.
- Utilisation d'un environnement informatique haute performance (HPC) distribué pour déployer un 3D Convolutional Neural Network (CNN) pour supprimer la couverture nueuseuse, ce qui a permis de produire un ensemble de données quotidien sans lacunes sur plus de 20 ans concernant la couverture de glace lacustre de 1,391 lacs à travers le monde.
- Harmonisation des données géospatiales relatives à la couverture de glace des lacs à partir de plusieurs types de capteurs et formats, afin de permettre la validation d'algorithmes à grande échelle.
- **Technologies utilisées :** Dask, Xarray, Linux, Python, Zarr

### Data Scientist *Université of Waterloo*

Mai 2024 - Août 2024  
*Waterloo, ON*

- Découverte de modèles et de tendances dans les choix agricoles à travers les Prairies canadiennes à l'aide de l'ensemble de données Annual Crop Inventory, telles que la popularité croissante de la rotation des cultures céréalières et légumineuses.
- **Technologies utilisées :** Google Earth Engine, Python, GeoPandas

### Chercheur Invité *ECMWF*

Octobre 2023 - Décembre 2023  
*Reading, Royaume-Uni*

- Utilisation d'un environnement informatique haute performance (HPC) pour mener des expériences sur la paramétrisation d'un modèle numérique de prévision météorologique haute résolution.
- Amélioration du modèle lacustre dans ECLand (la composante terrestre de l'IFS) grâce à la paramétrisation de la neige sur la glace, afin de tenir compte de son effet isolant, ce qui a permis de corriger avec succès la surestimation de l'épaisseur de la glace.
- **Technologies utilisées :** Fortran, Python, HPC

### Chercheur en apprentissage automatique *H2O Geomatics*

Janvier 2023 - Septembre 2023  
*Kitchener, ON*

- Formation de modèles distribués sur des clusters HPC à l'aide d'un paradigme de parallélisme de données distribué afin d'accélérer le développement de plus de 60/
- Développement d'un pipeline d'ingestion et d'harmonisation des données, qui combine ERA5 et IMS (ensemble de données spatiales sur la couverture de neige) dans un format NetCDF harmonisé, avant de les partitionner en exemplaires individuel d'entraînement et de test pour l'apprentissage automatique.
- Suivi et contrôle du développement des modèles et des expériences à l'aide de Weights and Biases (WandB), permettant une allocation optimisée des ressources et une réponse rapide aux défaillances
- Développement d'un nouvel algorithme de prévision de la glace sur lacs basé sur des réseaux de transformateurs spatiotemporels, surpassant le modèle physique prédominant FLake.
- **Technologies utilisées :** Python, Linux, HPC

**Ingénieur de données**  
*Lakes Environmental Software*

Sep 2022 - Dec 2022  
Waterloo, ON

- Conception d'un pipeline de données d'ingestion et d'harmonisation pour traiter les étiquettes d'experts Dynamic World V1 et les observations SAR Sentinel-1, qui ont produit des échantillons pour l'entraînement de modèles d'apprentissage automatique en aval afin d'effectuer la classification de l'utilisation des sols à partir du SAR
- **Technologies utilisées :** Google Earth Engine, Rasterio, Python, Dask

**Chercheur en apprentissage profond**  
*Lakes Environmental Software*

Septembre 2021 - Avril 2022  
Waterloo, ON

- Évaluation de diverse architecture d'apprentissage profond (FNO, ConvLSTM, CNN-LSTM) sur la simulation d'un problème classique d'écoulement dans une cavité à couvercle, obtenant d'excellentes performances (RMSE  $0.0061 m \cdot s^{-1}$ ) et donnant lieu à une publication.
- Intégration du bilan massique dans l'entraînement du modèle, ce qui a permis d'obtenir des valeurs de bilan massique comparables ( $10^{-5}$  à celles de la solution numérique directe produite par OpenFOAM).
- **Technologies utilisées :** Python, PyTorch, Keras, Dask

**Assistant de recherche**  
*Université of Waterloo*

Janvier 2021 - Août 2021  
Waterloo, ON

- Développement d'un modèle de forêt aléatoire pour prédire les charges et nutriments à partir des caractéristiques d'utilisation de sols, obtenant des scores  $R^2$  entre 0.6 et 0.8, et donnant lieu à une publication.
- **Technologies utilisées :** Python, Scikit-Learn, Pandas

## ÉDUCATION

**Licence en informatique**  
Diplôme en développement durable  
Université of Waterloo

2019 - 2024

## CERTIFICATIONS

**Certificat professionnel en ingénierie des données (DeepLearning.AI)** Septembre 2025 - aujourd'hui  
Cours complet sur le cadre de l'ingénierie des données, dispensé par Joe Reis, auteur du livre 'Fundamentals of Data Engineering'. Apprentissage des principes d'une bonne architecture des données et application de ces principes dans des laboratoires pratiques pour créer des systèmes de données sur le cloud AWS.  
**Concepts clés appris :** Architecture des données, pipelines par lots et en continu, infrastructure en tant que code, orchestration, DataOps [Automatisation, Observabilité, Réponse aux incidents]  
**Technologies Applied :** S3, Airflow, Terraform, Kinesis Data Streams, PostgreSQL, DynamoDB, REST API

**Land in Focus – Basics of Remote Sensing**  
EO College

2025

## COMPÉTENCES TECHNIQUES

<b>Cloud &amp; Orchestration</b>	S3, Airflow, Terraform
<b>Bases de données</b>	PostgreSQL, Amazon RDS, DynamoDB
<b>Big Data &amp; Traitement distribué</b>	Dask, Xarray, Apache Spark
<b>Formats de données</b>	Zarr, NetCDF COG, HDFS
<b>Langages</b>	Python, C, C++, Bash, SQL, JavaScript, R
<b>Géospatial</b>	QGIS, STAC, Zarr, GDAL
<b>CI/CD</b>	GitHub Actions
<b>Apprentissage automatique</b>	PyTorch, PyTorch Lightning, Keras, Scikit-Learn

## COMPÉTENCES TRANSVERSALES

- Excellentes compétences en communication en français et en anglais
- Capacité à travailler de manière autonome avec un minimum d'encadrement
- Priorité aux exigences des clients

## PROJECTS

---

### **Plateform automatisée de détection des rochers et falaise (en cours de développement)**

Plateforme qui soutient le développement local de l'escalade en permettant aux utilisateurs de détecter les rochers et les falais potentiels. Elle fournit des raster d'une résolution de 1 mètre à partir de données DEM prétraitées, en utilisant la détection des contours et l'apprentissage automatique pour mettre en évidence les caractéristiques périentiques.

**Technologies utilisées :** Titiler, Cloud-Optimized GeoTiff, GDAL, Vite, Python.

## PUBLICATIONS ET PRÉSENTATIONS

---

**Samuel J. Johnston**, Justin Murfitt, Claude Duguay, and Clément Albergel. 2025. «Addressing Cloud Contamination in Satellite Derived Lake Ice Cover Products.» *Remote Sensing of Environment*, soumis.

**Samuel J. Johnston**, Justin Murfitt, and Claude Duguay. 2025. «A Deep Learning Approach to Lake Ice Cover Forecasting» *GMD*, soumis.

Costa Rocha, Paulo A., **Samuel J. Johnston**, Victor Oliveira Santos, Amir A. Aliabadi, Jesse V.G. Thé, and Bahram Gharabaghi. 2023. «Deep Neural Network Modeling for CFD Simulations: Benchmarking the Fourier Neural Operator on the Lid-Driven Cavity Case.» *Applied Sciences*, 13, (5): 3165-65. <https://doi.org/10.3390/app13053165>

Basu, Nandita B., J. Dony, K.J. Van Meter, **Samuel J. Johnston**, and Anita T. Layton. 2023. «A Random Forest in the Great Lakes: Stream Nutrient Concentrations across the Transboundary Lake Basin.» *Earth's Future*, 11, (4). <https://doi.org/10.1029/2021EF002571>