# Enhancing ML Models for Solar Weather Forecasting using Clustering and Adversarial Anomaly Detection

MSc. Dissertation by Ivo Saavedra for MEIC

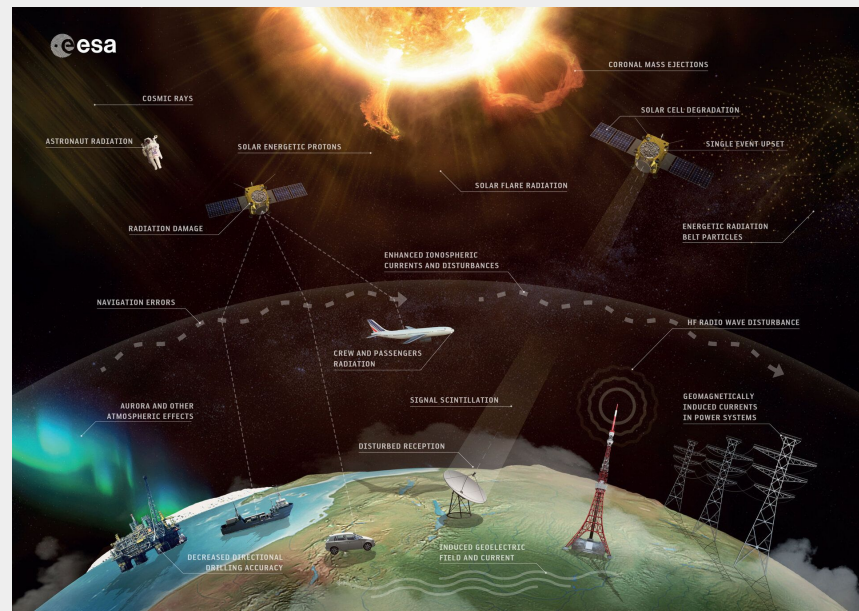Supervised by André Restivo and Filipa Barros

**U.PORTO**

FEUP **FACULDADE DE ENGENHARIA**
UNIVERSIDADE DO PORTO

# 01.

# INTRODUCTION

# CONTEXT - SOLAR WIND

- **Space Weather Science** is a field of research that aims to understand and predict solar phenomena

- An example of these is the **solar wind**

- **Coronal Mass Ejections** (CMEs) are known to affect:
  - Electrical Grids
  - Geolocation Systems
  - Radio-Communication Systems
  - Spacecraft and human in orbit



Space Weather Effects
Taken from: https://www.esa.int/ESA_Multimedia/Images/2018/01/Space_weather_effects

# MOTIVATION

- Most solar phenomena are still not fully understood
  - □ Prediction remains difficult

- **Magnetohydrodynamic** (MHD) simulators have been developed to extrapolate the conditions that lead to these events

- An example of this is **MULTI-VP**[1]
  - □ Simulates the 3D structures of solar wind
  - □ Calculates many 1D solar wind solutions from flux-tube geometries and heating functions
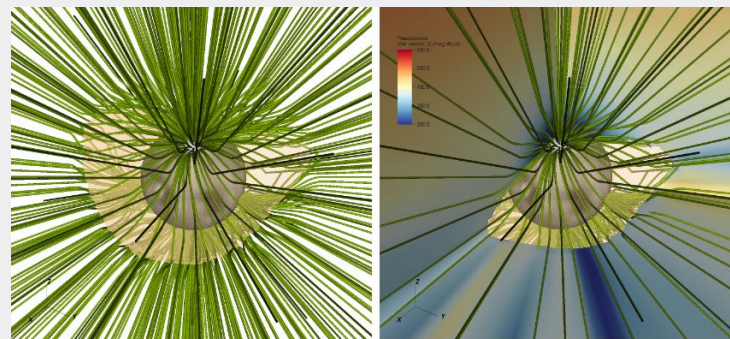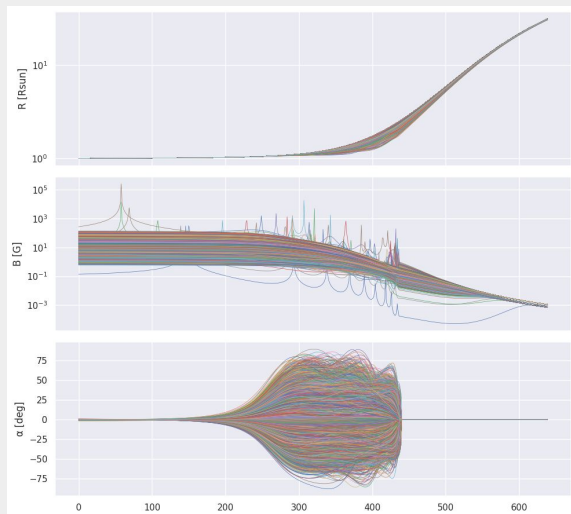


Illustration of the operation of MULTI-VP
Taken from: Rui F. Pinto and Alexis P. Rouillard [1]
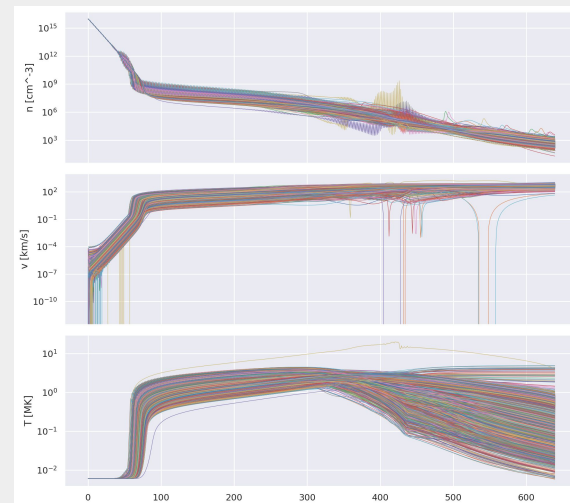
# MAGNETOGRAM DATA AND PREDICTIONS

**Input (Partial Flows):**
- **R[Rsun]** - radial coordinate radius
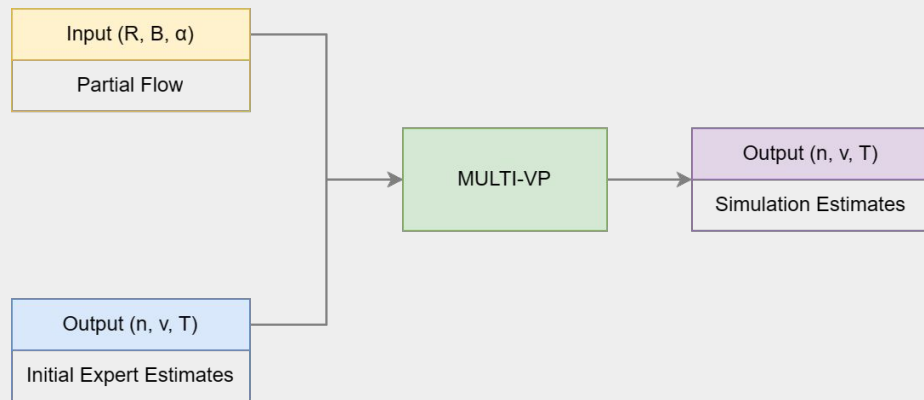- **B[G]** - magnetic field
- **ɑ[degree]** - flux tube inclination

**Outputs (Predicted Flows):**
- **n[cm$^{-3}$]** - number of protons per unit volume
- **v[km/s]** - speed-oriented through the line
- **T[MK]** - temperature at a point in space
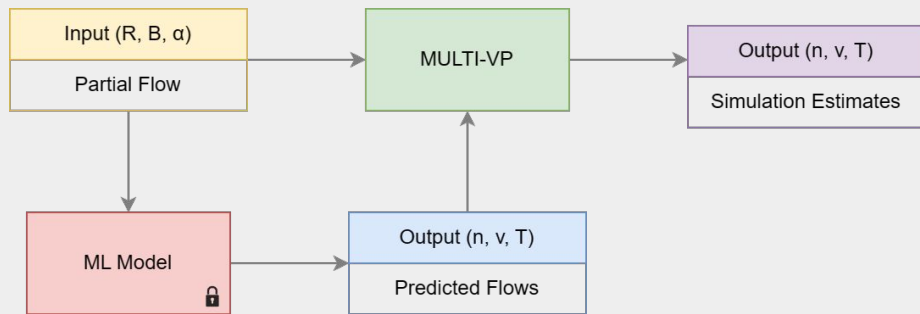
# INITIAL APPROACH: MULTI-VP



MULTI-VP Methodology

- MHD that simulates 3D structures of the solar wind
- Takes magnetogram data as input
- Predicts solar wind flows based on the initial partial flows and expert estimates

**Problems:**
- Takes a long time to converge
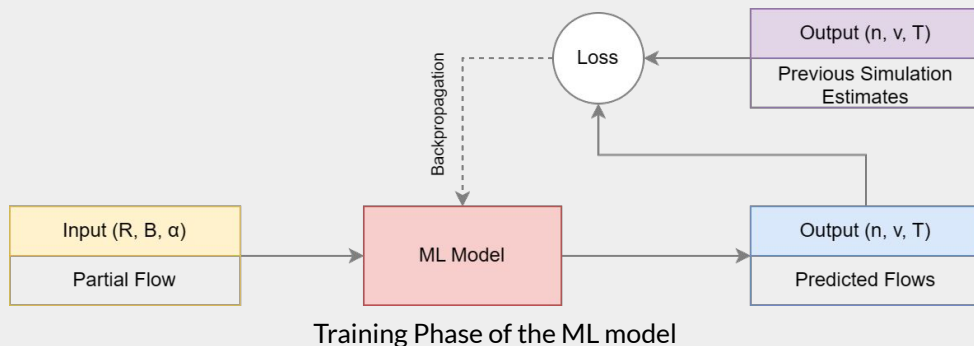- Requires expert initial guesses

# BASELINE APPROACH: ML FOR INITIAL CONDITION ESTIMATION



| Input (R, B, α) | | MULTI-VP | | Output (n, v, T) |
| --- | --- | --- | --- | --- |
| Partial Flow | | | | Simulation Estimates |

| ML Model | | Output (n, v, T) |
| --- | --- | --- |
| | | Predicted Flows |

ML for Initial Condition Estimation Methodology

- Introduces ML model to predict initial flows from partial flows
  - These are passed to MULTI-VP

- Estimations don't need to be made by hand

- Known to reduce the total simulation time of MULTI-VP

# BASELINE APPROACH: ML FOR INITIAL CONDITION ESTIMATION



Training Phase of the ML model

**Training:**
- Model takes initial flows and produces estimations
- The predictions are compared with initial simulation results from MULTI-VP
- Parameters are updated based on the calculated Loss

**Evaluation:**
- Comparing MULTI-VP's execution time
- MSE between real estimations and predicted ones

**Problem:**
- Doesn't capture periphery values
- Sensitive to anomalies in the training dataset

# 03.
# RESEARCH STATEMENT

# HYPOTHESIS

"By integrating clustering and adversarial anomaly detection techniques, the initial conditions predicted by RNNs for the MULTI-VP simulator will be closer to the final simulation results and contribute to faster executions."

# RESEARCH QUESTIONS

**RQ1 -** Are clustering methods capable of detecting characteristics in the dataset that were overlooked by the original RNN and would help with the prediction task?

**RQ2 -** Do the estimates obtained with clustering-based training significantly improve the simulation's performance?

**RQ3 -** Can adversarial learning methods detect anomalies in solar wind profiles?

**RQ4 -** Does the resulting dataset significantly improve the predictive ability of the RNN?

**RQ5 -** Does the improved predictive ability of the RNN result in a reduction of execution time for MULTI-VP?

# 03.

# CLUSTERING

# STATE-OF-THE-ART

**Focus:** Clustering techniques to enhance ML models
**Platforms:** Scopus, Google Scholar

| Paper | Year | Type | Topic | Dimensionality Reduction |
|---|---|---|---|---|
| JAIN[3] | 2010 | Survey | KMeans Analysis | No |
| FAHAD EA.[4] | 2014 | Survey | Large Data Clustering | Yes |
| FAHIMAN EA.[5] | 2017 | Primary | Improving ML | No |

# EXPERIMENTS

**Methods:**
- TimeSeriesKMeans
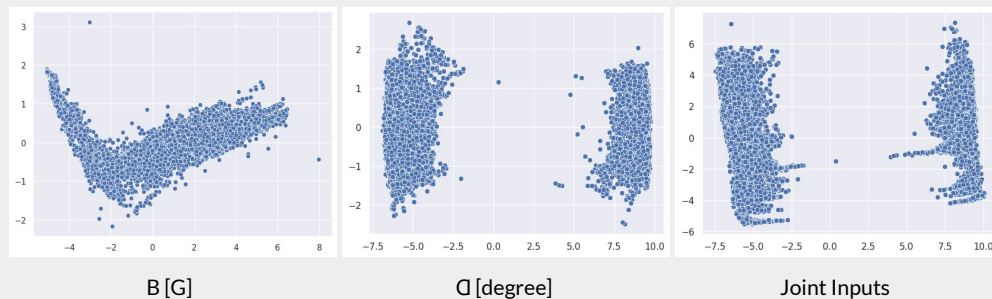- SOM
- KMeans
- AgglomerativeClustering
- DBSCAN

**Validity:**
- Elbow Test (KMeans)
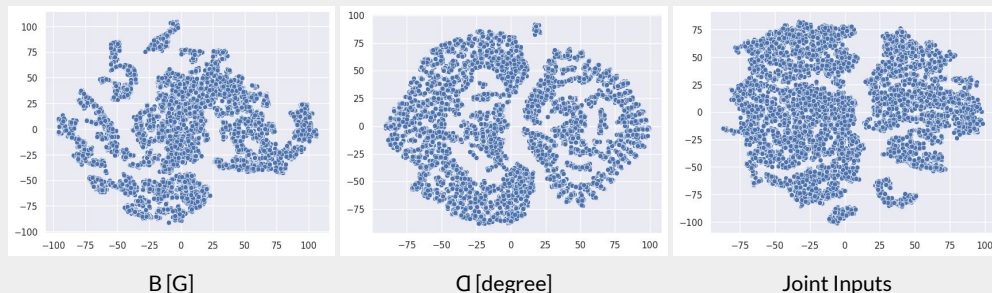- Silhouette Score
- Calinski-Harabasz Index
- Davies-Bouldin Index

**ML Evaluation:**
- Train a model for each cluster
- Measure Predictions' MSE

## Principal Component Analysis



B [G]  　　　  Cl [degree]  　　　  Joint Inputs

## t-distributed Stochastic Neighbor Embedding



B [G]  　　　  Cl [degree]  　　　  Joint Inputs

# EXPERIMENTAL RESULTS



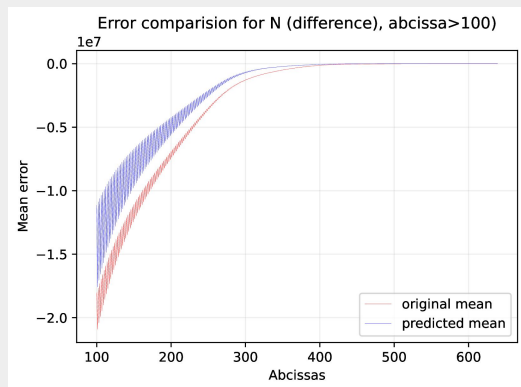KMeans of the PCA of the joint input variables

**Selection Criteria:**
- Predictions' MSE (after model training with clustering approach)
- Validity metrics
- Cluster Distribution
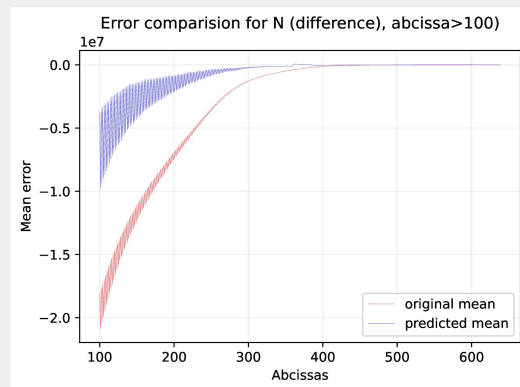
**Chosen methodology:**
1. Apply **PCA** on the **joint input** variables
2. **KMeans** clustering of the representation
3. Train prediction model for each cluster
4. Generate **validation predictions** and use them as **initial conditions** to MULTI-VP

# MULTI-VP EVALUATION: N [CM$^{-3}$]
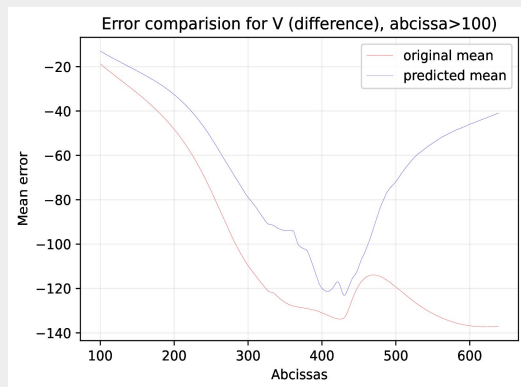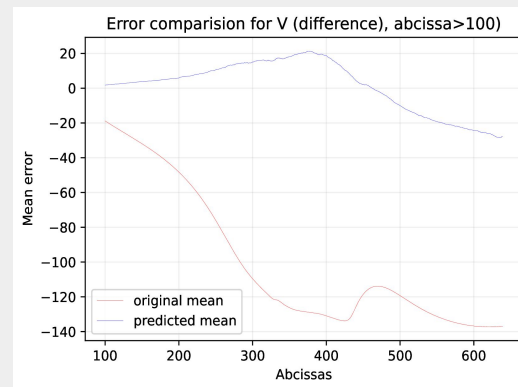
Baseline Results

Clustering RNN Results



- MSE between the initial conditions and simulation outputs for the *n [cm$^{-3}$]* variable
- Significant reduction in the distance between initial conditions and final outputs in the clustering approach

# MULTI-VP EVALUATION: V [KM/S]

Baseline Results

Clustering RNN Results



- MSE between the initial conditions and simulation outputs for the *v [Km/s]* variable
- Significant reduction in the distance between initial conditions and final outputs in the clustering approach

# MULTI-VP EVALUATION: T [MK]
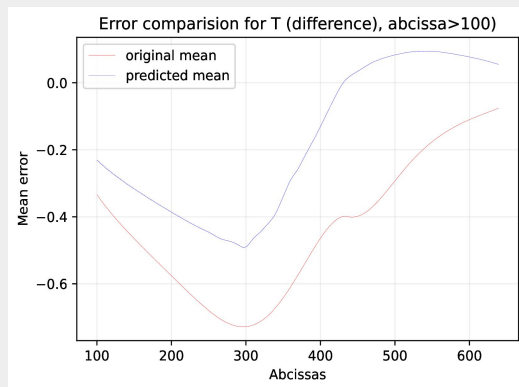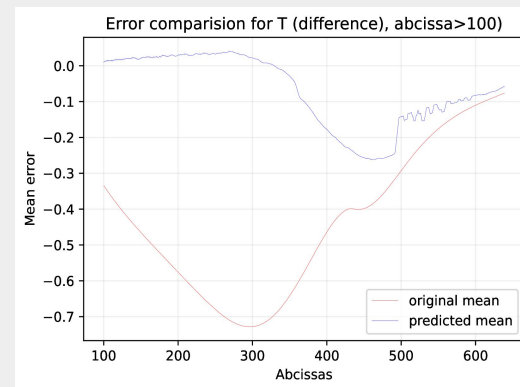
Baseline Results



Clustering RNN Results



- MSE between the initial conditions and simulation outputs for the *T [MK]* variable
- Significant reduction in the distance between initial conditions and final outputs in the clustering approach

# MULTI-VP EVALUATION: MEAN SPEED-UP

1.06 → 1.05

Baseline        Clustering Models

# 03.

# ADVERSARIAL ANOMALY DETECTION

# STATE-OF-THE-ART

**Focus:** Adversarial anomaly detection in tabular data
**Platform:** Scopus

| Paper | Year | Training Method | Anomaly Detection | Architecture | Application |
|-------|------|-----------------|-------------------|--------------|-------------|
| MAD-GAN[6] | 2019 | Normal | Reconstruction and D Loss | Vanilla | Time Series |
| TANO-GAN[7] | 2020 | Normal | Reconstruction | Vanilla | N/A |
| FGAN[8] | 2019 | Normal and Division Boundary | Adapted G and D Loss | Vanilla | Time Series |

# EXPERIMENTS

**Architectures:**
- Linear GAN
- MAD-GAN (adapted from [4])
- Adversarial AE

**Anomaly Scores (AS):**
- Discriminator Score
- MSE Reconstruction (Generator)
- MSE-Discriminator Reconstruction (Generator and Discriminator)

**General Approach:**
1. Learn the normal distribution of the data in the training phase
2. Use either/both generator and discriminator to determine anomaly score for each profile
3. Filter profiles with anomaly scores above a predefined threshold (hyperparameter)

(Each architecture was designed and optimized around the input variables and later applied to the output variables from MULTI-VP)
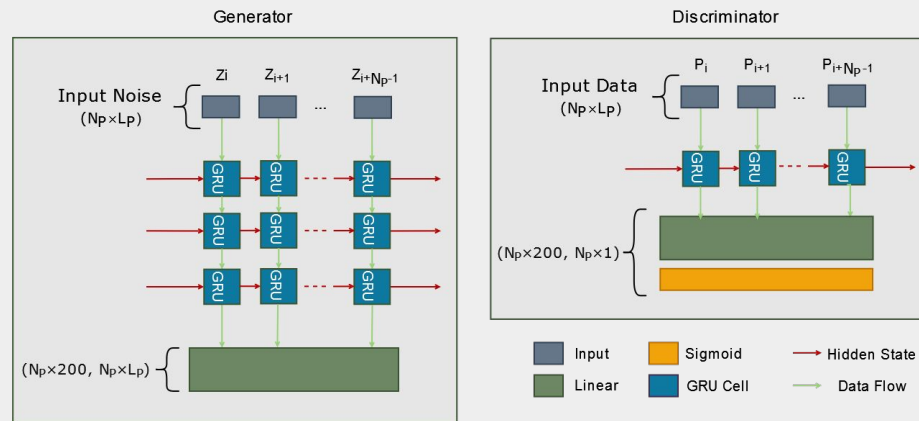
# EXPERIMENTS: ARCHITECTURE SELECTION

- Select best AS for each architecture (based on visual inspection of the datasets after filtering)
- Compare architecture+AS performances and select best

| Architecture | Function | Inputs | | Outputs | |
|---|---|---|---|---|---|
| | | Thresh (%) | #Profiles | Thresh (%) | #Profiles |
| LINEAR GAN | Rerr | 10 | 1177 | 10 | 1177 |
| MAD-GAN | **Rerr** | **3** | **352** | **3** | **352** |
| AAE | RDerr | 10 | 1177 | 10 | 1177 |

Anomalies for each of the architecture+AS combination
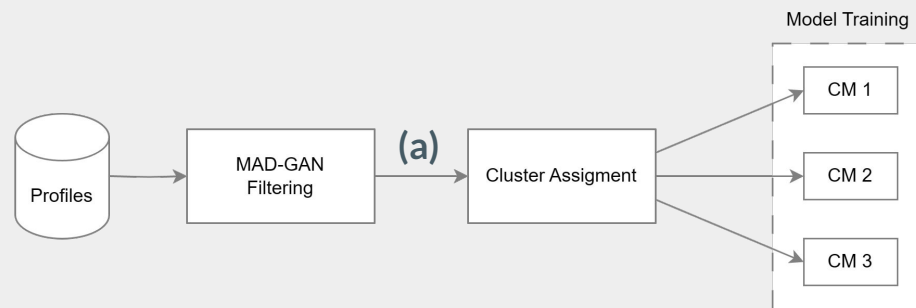
# MAD-GAN

- Takes dataset windows of consecutive observations

- Determines anomaly score based on the abnormality of samples within the window and the entire dataset

- Adapted to accept consecutive profiles instead of time-series:
  - Dimensionality reduction of the data features
  - Anomaly detection in the magnetic field for the input model
  - Anomaly detection in all the output variables for the output model
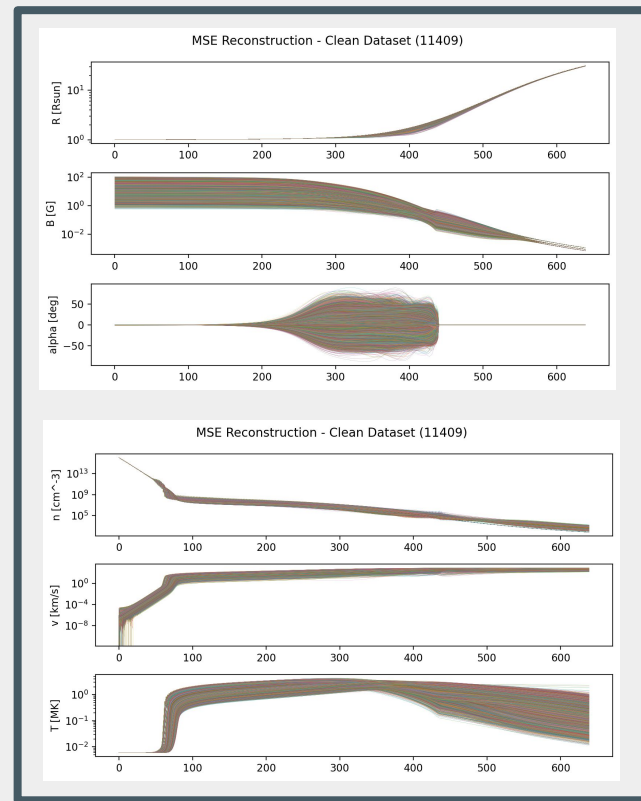


MAD-GAN Architecture

# MULTI-VP EXPERIMENTS - SETUP



Model Training

Anomaly Detection Clustering Dataflow
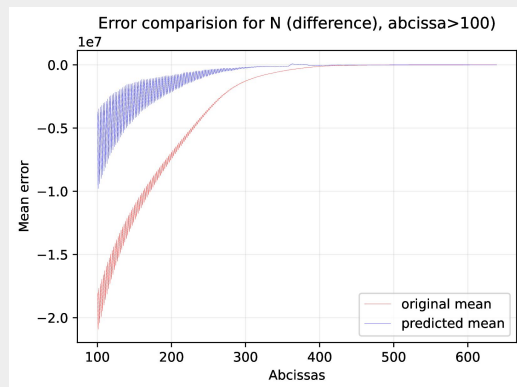
**(a) Filtered Dataset**



1. Detect and filter profiles with **anomalous input** or **output variables**
2. Separate the dataset into clusters
3. Retrain the clustering models from the previous experiments with the new datasets
4. Generate **validation predictions** and use them as **initial conditions to MULTI-VP**

# MULTI-VP EVALUATION: N [CM⁻³]

Clustering Models' Results

Anomaly+Clustering Models' Results



- MSE between the initial conditions and simulation outputs for the *n [cm⁻³]* variable
- Slight improvement in some of the initial abscissas when compared to the method without anomaly detection
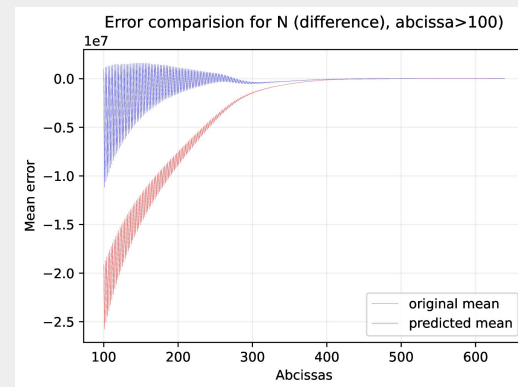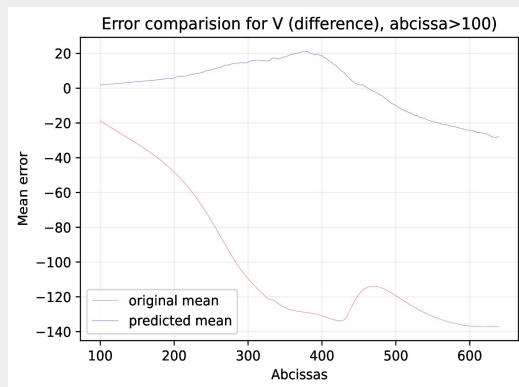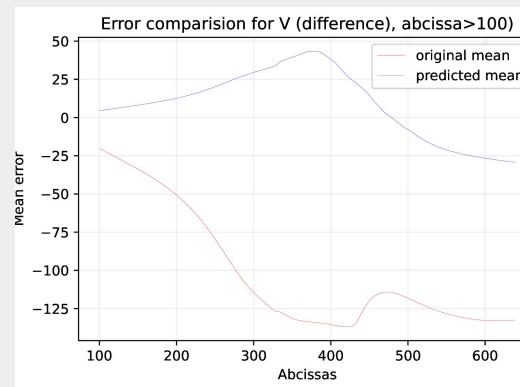
# MULTI-VP EVALUATION: V [KM/S]

Clustering Models' Results

Anomaly+Clustering Models' Results





- MSE between the initial conditions and simulation outputs for the *v [Km/s]* variable
- Significant increase in the abscissa distance in comparison with the previous results
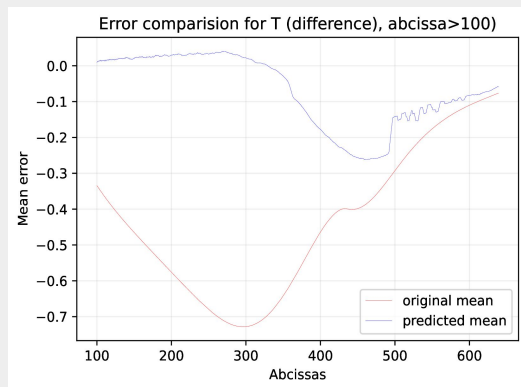
# MULTI-VP EVALUATION: T [MK]

Clustering Models' Results

Anomaly+Clustering Models' Results



- MSE between the initial conditions and simulation outputs for the *T [MK]* variable
- Worse distance to the simulation predictions when compared to the initial expert estimates and the previous experiments

# MULTI-VP EVALUATION: MEAN SPEED-UP

1.06  →  1.05  →  1.06

Baseline        Clustering Models        Anomaly
+
Clustering Models

# RESEARCH QUESTION ANALYSIS I

**RQ1 - *Are clustering methods capable of detecting characteristics in the dataset that were overlooked by the original RNN and would help with the prediction task?*** Yes. We managed to produce closer initial conditions to the final simulation outputs. Dividing the data into clusters of approximately the same size, made it possible for the RNN to capture previously unseen/ignored features.

**RQ2 - Do the estimates obtained with clustering-based training significantly improve the simulation's performance?** No. The overall computation time didn't improve over the baseline model. A mean speedup of 1.05 was obtained with the clustering method, in contrast to the 1.06 speedup from the baseline.

# RESEARCH QUESTION ANALYSIS II

**RQ3 - Can adversarial learning methods detect anomalies in solar wind profiles?** Yes. Every implemented architecture was able to detect anomalous profiles in the dataset, with MAD-GAN outperforming the others. We managed to produce an apparently cleaner version of the original dataset.

**RQ4 - Does the resulting dataset significantly improve the predictive ability of the RNN?** No. The removal of anomalous profiles hindered the quality of the predictions from the RNN model. By excluding entire profiles from training, we might have removed important features.

**RQ5 - Does the improved predictive ability of the RNN result in a further reduction of execution time for MULTI-VP?** Even with worse initial conditions, the simulation took less time to reach a viable solution. This might indicate that the computation time is not directly linked to the proximity of the initial conditions to the simulation outputs.

# CONCLUSIONS

- ■ Managed to produce initial conditions closer to simulation outputs
- ■ Didn't reduce computation time of the simulation

**Future Work:**

- ■ Explore if the speedup is being well calculated
- ■ Test the physical validity of the estimates
- ■ Surrogate model
- ■ Test on other MHD simulators

# REFERENCES

[1] Rui F. Pinto and Alexis P. Rouillard. A Multiple Flux-tube Solar Wind Model. The Astrophysical Journal, 838(2):89, 2017.

[2] Ana Filipa Sousa Barros. Initial Condition Estimation in Flux Tube Simulations using Machine Learning, 2021.

[3] Anil K. Jain. Data clustering: 50 years beyond K-means. Pattern Recognition Letters, 31(8):651–666, 2010

[4] Adil Fahad, Najlaa Alshatri, Zahir Tari, Abdullah Alamri, Ibrahim Khalil, Albert Y. Zomaya, Sebti Foufou, and Abdelaziz Bouras. A Survey of Clustering Algorithms for Big Data: Taxonomy and Empirical Analysis. IEEE Transactions on Emerging Topics in Computing, 2(3):267–279, 2014.

[5] Fateme Fahiman, Sarah M. Erfani, Sutharshan Rajasegarar, Marimuthu Palaniswami, and Christopher Leckie. Improving load forecasting based on deep learning and K-shape clustering. In 2017 International Joint Conference on Neural Networks (IJCNN), pages 4134–4141, 2017.

[6] D. Li, D. Chen, B. Jin, L. Shi, J. Goh, and S.-K. Ng. MAD-GAN: Multivariate Anomaly Detection for Time Series Data with Generative Adversarial Networks. Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics), 11730 LNCS:703–716, 2019.

[7] Md Abul Bashar and Richi Nayak. TAnoGAN: Time Series Anomaly Detection with Generative Adversarial Networks. In 2020 IEEE Symposium Series on Computational Intelligence (SSCI), pages 1778–1785, 2020.

[8] Phuc Cuong Ngo, Amadeus Aristo Winarto, Connie Khor Li Kou, Sojeong Park, Farhan Akram, and Hwee Kuan Lee. Fence GAN: Towards Better Anomaly Detection. In 2019 IEEE 31st International Conference on Tools with Artificial Intelligence (ICTAI), pages 141–148, 2019.

# THANK YOU FOR YOUR TIME