

# Water Potability Report

Peer Group Members' Names:

1. Gaius Irakiza
2. Tamanda Kaunda
3. Jeremiah Agbaje
4. David Ubushakebwimana
5. Rene Ntabana

[Link to GitHub repository](#)

## Group Reports

Train Instance	Engineer Name	Regularizer	Optimizer	Early Stopping	Dropout Rate	Accuracy	F1 Score	Recall	Precision
1.	Jeremiah Agbaje	l2(0.00001)	Adam(learning_rate=0.001)	monitor='val_loss', patience=50,  restore_best_weights=True	0.5 between all my layers	0.6541	0.2054	0.1210	0.6786
2	Gaius Ira...	L2 Batch Normalization	NADAM with learning_rate(0.006)	Patience = 30, min_delta = 0.0001 on val_loss	0.2 between all my layers	0.6687	0.5192	0.4583	0.5986
3.	Tamanda Kaunda	L1 (progressive)	RMSprop	monitor='val_loss',  patience=35,  min_delta=0.001,  restore_best_weights=True	0.3/0.4/0.3	0.6494		0.2930	0.5476
4.	David Ubushakebwimana	l2(0.0005)	RMSprop(learning_rate=0.001)	monitor='val_loss', # Monitor validation loss	No dropout	0.6447		0.4204	0.5238

				patience=50,  restore_best_weights=True					
5.	Rene Ntabana	L1	Adam (lr=0.0005)	Patience = 30 on val_loss	0.4 / 0.4 / 0.3	0.6800	0.5211	0.4713	0.5827

## **Dataset Overview**

- Source: [Kaggle Water Potability Dataset](#)
- Samples: ~3,276
- Features:
  - 9 continuous variables (e.g., pH, Solids, Conductivity, Sulfate, etc.)
- Target Variable: Potability (0 = Not Potable, 1 = Potable)

## **Data Preprocessing Pipeline**

### **1. Handling Missing Values**

- Missing values in Sulfate, Trihalomethanes, and pH were handled using MICE (Multiple Imputation by Chained Equations) via IterativeImputer.
- This method is more robust than mean/median imputation because it models the missing values using other features in a round-robin regression loop, preserving multivariate structure.

### **2. Outlier Treatment**

- Applied IQR clipping on skewed continuous variables to reduce the influence of extreme outliers that can degrade model convergence and overfit the network. And this is better than dropping the rows, cause it preserves all the data since our dataset is already small. This helps the model have more training data.

### **3. Feature Scaling**

- Used StandardScaler to normalize features to a standard Gaussian distribution.
- Deep neural networks benefit from standardized inputs for stable and fast convergence (especially with optimizers like Nadam).

### **4. Train-Test Split**

- 70/15/15 split for Train/Validation/Test sets using stratified sampling to maintain class distribution.

## **Individual Reports**

### **Jeremiah Agbaje Model Analysis**

This was the first model built in the project and improvement to future models were made from here. It is a sequential neural network made up of several dense layers and dropout layers. The model starts with a dense layer of 128 neurons, followed by dropout to reduce overfitting. Then it has a second dense layer with 64 neurons, a third with 32 neurons, and finally an output layer with 1 neuron using the sigmoid activation for binary classification.

In total, the model has 34,949 parameters, with 11,649 of them being trainable. I used dropout after every dense layer to help the model generalize better. This basic setup helped me test and tune my regularization and optimization techniques throughout the training process.

**Model Architecture:** `model_jeremiah_agbaje()`

**Input** → **Dense(128)** → **ReLU** → **Dropout(0.5)** →  
**Dense(64)** → **ReLU** → **Dropout(0.5)** →  
**Dense(32)** → **ReLU** → **Dropout(0.5)** →  
**Dense(1)** → **Sigmoid**

Layer-by-Layer				
Layer	Type	Neurons	Activation	Regularization
1	Dense	128	ReLU	L2(0.0001), Dropout(0.5)
2	Dense	64	ReLU	L2(0.0001), Dropout(0.5)
3	Dense	32	ReLU	L2(0.0001), Dropout(0.5)
4	Output	1	Sigmoid	None

## Design Justifications

### 1. L2 Regularization (l2(0.0001))

I chose **L2 regularization** instead of L1 because I didn't want to force the weights to become zero. I wanted the model to learn important patterns, still, so I used a small value of **0.0001**. This way, the model could still have larger weights when needed, but not too large. It helped the model stay balanced and not too simple or too complex.

### 2. Dropout (Dropout(0.5))

At first, when I trained the model without dropout, it started overfitting very early, within the first **10 epochs**. So, I added **Dropout with a rate of 0.5**. This means that during training, half of the neurons were randomly turned off in each layer. This helped the model not depend too much on certain neurons and allowed it to generalize better when tested on new data. It also fixed the overfitting problem in the model, though the model began underfitting after 200 epochs.

### 3. Early Stopping (patience=50, restore\_best\_weights=True)

I used **Early Stopping** to stop training if the model didn't get better after some time. I set the **patience to 50**, because I wanted the model to have enough time to train up to **200 epochs** if needed. But I also didn't want it to keep training if it had already passed its best point. That's why I set **restore\_best\_weights=True**, so the model would go back to the best weights it had during training.

### 4. Activation Function: ReLU (Rectified Linear Unit)

I chose **ReLU** as the activation function for all hidden layers because of its efficiency and compatibility with the **data characteristics of water potability prediction**, where the inputs (like pH, conductivity, turbidity, sulfate levels, etc.) are mostly **non-negative and naturally**. ReLU naturally handles this kind of data well since it outputs zero for negative inputs and retains positive values as-is.

### 5. Optimizer: Adam (Learning Rate = 0.0001)

I chose **Adam (Adaptive Moment Estimation)** because it's highly effective for this dataset for water quality measurements, which contained **irregular patterns, missing values, and non-linear relationships** between features and outcomes.

Water quality data isn't always clean — features like **Turbidity, pH, and Sulfate** can vary widely and unpredictably. Adam adjusts learning rates automatically for each parameter, making it **resilient to noisy gradients** and helping the model stay stable even when some inputs fluctuate unexpectedly.

I used a **low learning rate (0.0001)** because the model includes **dropout and L2 regularization**, which make learning more cautious. A small step size ensures that we **don't "overshoot" the optimal weights**, allowing the model to **gently converge** to the best settings without getting thrown off by sharp loss changes.

With this setup, the model reached a **test accuracy of 65.41%** and a **precision of 67.86%**. Even though **recall (12.10%)** and **F1 score (20.54%)** were low, this regularization and optimization plan helped the model avoid overfitting and perform better on new data than models I tested without regularization.

## Git Logs

```
(base) jeremiah@Applegeniuss-MacBook-Pro PeerGroup8WaterQualityModel % git log --author="j-agbaje"
commit cba67419edee821c3feca44791ef20825f1631b2
Author: j-agbaje <144043288+j-agbaje@users.noreply.github.com>
Date: Sun Jun 8 14:18:39 2025 +0200

    Add data and separate model from initial notebook

commit e982193bc0e06c49ba3196b522e4c481e9e2360b
Author: j-agbaje <144043288+j-agbaje@users.noreply.github.com>
Date: Sat Jun 7 22:51:11 2025 +0200

    Add precision, f1, and recall scores

commit 5a19f130ec316a1ab1d569570e68efcd47c77665
Author: j-agbaje <144043288+j-agbaje@users.noreply.github.com>
Date: Fri Jun 6 23:46:29 2025 +0200

    Add data loading, preprocessing and first model

commit ebfaafd3f0c779873471fe8a40c41fda617974ac
Author: j-agbaje <144043288+j-agbaje@users.noreply.github.com>
Date: Fri Jun 6 23:45:23 2025 +0200

    Delete Peer_Group_8__formative_II_.ipynb

commit 8d048f4aedd352c5a1dff9323c71a500766fb55e
Author: j-agbaje <144043288+j-agbaje@users.noreply.github.com>
Date: Fri Jun 6 23:43:51 2025 +0200

    Add files via upload

    Add an initial notebook with data preprocessing and the first model created
(base) jeremiah@Applegeniuss-MacBook-Pro PeerGroup8WaterQualityModel %
```

## Model Evaluation

I compared my model to Gaius', Rene's, David's, and Tamanda's models. I looked at important metrics like **F1 score**, **recall**, **precision**, and **loss**. These help to show how well each model performs, especially on predicting water potability. In this task, recall is very important because we don't want the model to say water is safe (potable) when it is not.

---

## Comparison with Gaius' Model

My model uses a simpler 3-layer architecture (128, 64, 32) with ReLU activation and a dropout of 0.5 across all layers. In contrast, Gaius' model is significantly deeper, with **7 hidden layers** gradually decreasing in size (down to 4 neurons) and uses **PReLU**, which allows each neuron to learn its own activation shape—possibly making it more expressive. He trains for 300 epochs while I train for 200 meaning his training time is also longer than

mine if not stopped by early stopping. Gaius also uses **Batch Normalization**, which I don't, and a **smaller dropout rate of 0.2**, meaning his model regularizes less aggressively. We both use early stopping and restore best weights, but Gaius uses **Nadam**, an optimizer that adds Nesterov momentum on top of Adam, while you use **Adam** with a lower learning rate. Overall, Gaius' model prioritizes complexity and expressiveness; while mine is simpler and more regularized.

- **My F1 Score: 0.2054 | Gaius' F1 Score: 0.5192**
  - His model balanced precision and recall better. Mine has a very low F1 score, which shows poor overall performance in classification.
- **My Recall: 0.1210 | Gaius' Recall: 0.4583**
  - His model found more of the truly potable water samples. Mine missed far too many.
- **My Precision: 0.6786 | Gaius' Precision: 0.5986**
  - My precision is higher, which means when I predicted potable, I was usually correct, which may be good because the data is imbalanced with more bad water than potable water but with the low recall, I still found less truly potable water samples which means the model was predicting potable less number of times.
- **My Loss: 0.6225 | Gaius' Loss: 0.6599**
  - My model had a lower loss, which may mean better confidence in predictions. But this is not enough because the F1 and recall are more important for water potability.
- **My Accuracy: 0.6541 | Gaius' Accuracy: 0.6687**

**Conclusion:** Gaius' model is better overall. His higher F1 score and recall make his model more useful for catching safe water samples. This is likely due to his deeper more complex model with reduced regularization.

**Why mine is worse:**

1. My recall is too low, meaning it misses too many actually good water samples.
2. Although my precision is high, my F1 score is also too low, showing the model doesn't balance precision and recall well.

---

**Comparison with Rene's Model**

Both I and René use a 3-layer setup with ReLU and the same neuron sizes. However, René uses **layer-specific dropout**: 0.4 across most layers and 0.3 at the end, compared to your uniform 0.5 dropout rate. He also uses **L1 regularization** (instead of your L2), and it varies between layers—he applies higher L1 penalties at the first and last layers. This means René's model encourages sparsity in the weights, whereas mine aims to keep weights small overall. Another key difference: René trains with a **higher learning rate (0.0005)**, potentially leading to faster convergence but also higher risk of overshooting. You are more conservative in tuning, while René uses a slightly more experimental structure. He also trains with more epochs(300) than me which means that training goes on for longer.

### Performance Metrics Comparison:

- **F1 Score:** My 0.2054 vs René's 0.5211 (+154% improvement)
- **Recall:** My 0.1210 vs René's 0.4713 (+289% improvement)
- **Precision:** My 0.6786 vs René's 0.5827 (-14% decrease)
- **Accuracy:** My 0.6541 vs René's 0.6800 (+4% improvement)

**Metric Interpretation:** René's model is **significantly superior** for this water safety task. His F1 score of 0.5211 indicates a much better balance between precision and recall, while my F1 score of 0.2054 reveals poor overall performance. Most critically, René's recall of 0.4713 means his model correctly identifies 47% of safe water samples, compared to my model's concerning 12%. This dramatic difference suggests my model would miss nearly 9 out of 10 safe water samples, creating a dangerous false negative problem.

### Why René's Model Outperforms Mine:

1. **Superior regularization strategy:** René's mixed L1 regularization with layer-specific penalties better controls overfitting while maintaining model flexibility, evidenced by his higher recall without catastrophic precision loss.
2. **Optimized training dynamics:** His higher learning rate combined with extended training allows better exploration of the loss landscape, resulting in a model that generalizes better to positive (safe water) cases.

---

### Comparison with David's Model

David's model is architecturally identical to mine: 3 layers (128, 64, 32) with ReLU and a single sigmoid output. But while I use **Adam**, David uses **RMSprop**, which handles non-stationary objectives well by dividing gradients by a running average of their magnitudes. We both use L2 regularization, but David's strength is higher (**0.0005** compared to your 0.0001), making his model more tightly constrained. He also trains for more epochs (300 vs. 200) and uses the same early stopping setup. In essence, David's model is a stricter version of mine with a different optimizer philosophy.

### Performance Metrics Comparison:

- **Recall:** My 0.1210 vs David's 0.4204 (+247% improvement)
- **Precision:** My 0.6786 vs David's 0.5238 (-23% decrease)
- **Loss:** My 0.6225 vs David's 0.6773 (+9% increase)
- **Accuracy:** My 0.6541 vs David's 0.6447 (-1% decrease)

**Metric Interpretation:** David's model demonstrates **superior practical performance** despite higher loss. His recall of 0.4204 means he correctly identifies 42% of safe water samples compared to my 12%, representing a critical improvement for this safety application. While my precision is higher (68% vs 52%), this advantage is meaningless when the model fails to detect most positive cases.

### Why David's Model Outperforms Mine:

1. **Stronger regularization control:** David's L2 regularization strength (0.0005) is 5x higher than mine, better preventing overfitting to the training data and improving generalization to unseen safe water samples.
2. **RMSprop optimizer advantage:** RMSprop's adaptive learning rate adjustment handles the non-stationary nature of this classification problem better than Adam, leading to improved convergence for positive class detection.

---

### Comparison with Tamanda's Model

Tamanda's model is deeper than mine with **4 layers** ( $96 \rightarrow 48 \rightarrow 24 \rightarrow 12$ ), compared to my 3. She also uses **l1 regularization** like René, applying it differently across layers (stronger in the first and output layers). We both avoid Batch Normalization, but Tamanda's regularization setup is more structured and potentially geared toward compressing the model. Also, she uses **RMSprop** with a higher learning rate (0.001), while I use a more cautious Adam with 0.0001. Her model uses a **larger batch size (64)**, which may smooth gradient estimates more than my batch size of 32. My model is simpler and more uniform, while hers is deeper and more deliberately pruned.

### Performance Metrics Comparison:

- **Recall:** My 0.1210 vs Tamanda's 0.2930 (+142% improvement)
- **Precision:** My 0.6786 vs Tamanda's 0.5476 (-19% decrease)
- **Accuracy:** My 0.6541 vs Tamanda's 0.6494 (-1% decrease)

**Metric Interpretation:** Tamanda's model shows **moderate improvement** over mine. Her recall of 0.2930 means she identifies 29% of safe water samples compared to my 12%, still representing a significant practical improvement for water safety detection.



## Why Tamanda's Model Outperforms Mine:

1. **Enhanced feature learning capacity:** Her deeper architecture with 4 layers captures more complex patterns in water quality indicators that my shallower model misses, particularly for positive class identification.
  2. **Strategic regularization approach:** L1 regularization with layer-specific penalties creates more effective feature selection, preventing overfitting while maintaining sensitivity to safe water indicators.
- 

## Insights and challenges faced

### 1. Regularization Optimization Challenge

**Initial Overfitting Crisis:** During preliminary training without regularization, the model exhibited severe overfitting within the first 10 epochs, achieving near-perfect training accuracy while validation performance stagnated. This immediate overfitting indicated the model was memorizing training patterns rather than learning generalizable features.

**Regularization Balancing Dilemma:** Once the need for regularization became apparent, I faced the complex challenge of selecting optimal regularization parameters. The primary difficulty was finding the precise combination of L2 regularization strength and dropout rates that would effectively prevent overfitting while maximizing validation performance. This required extensive hyperparameter tuning to achieve the delicate balance between model complexity and generalization capability.

### 2. Performance Ceiling Limitation

**70% Accuracy Barrier:** After successfully mitigating the initial overfitting problem, I encountered a persistent performance ceiling where the model consistently struggled to exceed 70% validation accuracy. Any attempts to push training accuracy beyond this threshold invariably resulted in overfitting, creating a frustrating trade-off between training performance and generalization ability. This challenge highlighted the inherent complexity of the water potability dataset and the difficulty in extracting meaningful patterns without compromising model stability.

### 3. Metric Selection and Evaluation Framework Error

**Misaligned Performance Focus:** I made a critical strategic error by exclusively focusing on validation accuracy as my primary performance benchmark during model development. This narrow focus led me to optimize for overall correctness rather than task-specific requirements. Only during post-development analysis did I realize the paramount importance of F1-score, recall, and precision for this water safety classification task. This misalignment resulted in a model with high precision but catastrophically low recall—fundamentally

unsuitable for a safety-critical application where missing positive cases (safe water) could have serious consequences.

#### 4. Data Preprocessing and Quality Management Challenges

**Imputation Strategy Selection:** I encountered significant difficulty in choosing the most appropriate missing value imputation method among mean, median, and MICE (Multiple Imputation by Chained Equations) approaches. The challenge was determining which technique would best preserve the underlying data patterns and relationships without introducing bias that could negatively impact model training and performance.

**Outlier Management Dilemma:** Additionally, I faced complex decisions regarding outlier detection and removal. The challenge involved carefully analyzing which data columns contained legitimate extreme values versus noise, and determining appropriate outlier removal strategies that would improve data quality without eliminating potentially important information that could be crucial for accurate water potability classification.

**Data Integrity Considerations:** Throughout the preprocessing phase, I had to balance data cleaning thoroughness with information preservation, ensuring that preprocessing decisions enhanced rather than compromised the dataset's ability to represent real-world water quality variations.

#### Final Reflection/ Lessons Learned

Overall, all four teammates had better **recall** and **F1 scores** than me. Even though my model had better **precision** and the lowest **loss**, this wasn't enough as recall is the most important metric in water potability testing. The goal of this task is to safely identify potable water. That means it's more important **not to miss** clean water samples (higher recall), even if sometimes we predict false positives. As for what I learned:

1. I learned the importance of regularization balance, and the importance of understanding what each figure is doing in the model rather than just plugging them in and hoping for improvements. My model was too imbalanced:
  - Aggressive dropout (0.5) was randomly removing half the model's capacity during training
  - Weak L2 penalty wasn't guiding the remaining weights effectively
  - Conservative learning rate meant the model couldn't adapt quickly enough to compensate

**Result:** The model became **too conservative** - it learned to make more safe predictions (hence higher precision) but lost the ability to recognize positive patterns (hence terrible recall). My other classmates had better regularization parameters than I, hence better recall and f1 scores.

2. I also learned the importance of considering the most important metrics for your use case. In this scenario recall, and f1 score are the most important metrics to judge the quality of the model not just test accuracy. They should've been taken into consideration during training.

### To improve the model I would:

#### 1. Increase Learning Rate

- Current: 0.0001 (too conservative)
- Increase to : 0.0005 - 0.001 (like René and Tamanda)

#### 2. Strengthen L2 Regularization ( $\lambda$ )

- Current: 0.0001 (too weak of a weight punishment relative to dropout)
- Increase: 0.0005 (like David) or implement mixed L1/L2
- Why: This balances your strong dropout and provides better weight guidance

#### 3. Optimize Dropout Strategy

- Reduce the aggressive uniform dropout from 0.5 to 0.3-0.4

#### 4. Increase Model depth:

- I would also explore the benefits increasing the depth of my model to better understand complex patterns within the data, and perhaps experiment with other optimizers such as RMSProp

### Gaius Irakiza Model Analysis

#### Model Architecture: model\_gaius\_irakiza()

Input  $\rightarrow$  [Dense  $\rightarrow$  BatchNormalization  $\rightarrow$  PReLU  $\rightarrow$  Dropout]  $\times$  5  $\rightarrow$  Dense(16)  $\rightarrow$  BN  $\rightarrow$  PReLU  $\rightarrow$  Output

#### Layer-by-Layer Explanation:

Layer	Type	Neurons	Activation	Regularization
1	Dense	128	PReLU	L2(0.001), BN, Dropout(0.2)
2	Dense	64	PReLU	L2(0.001), BN, Dropout(0.2)
3	Dense	32	PReLU	L2(0.001), BN, Dropout(0.2)
4	Dense	24	PReLU	L2(0.001), BN, Dropout(0.2)
5	Dense	8	PReLU	L2(0.001), BN, Dropout(0.2)
6	Dense	4	PReLU	BatchNorm only

Layer	Type	Neurons	Activation	Regularization
1	Dense	128	PReLU	L2(0.001), BN, Dropout(0.2)
2	Dense	64	PReLU	L2(0.001), BN, Dropout(0.2)
7	Dense	1	Sigmoid	None

## **Design Justifications**

### **Activation Function: PReLU (Parametric ReLU)**

I selected PReLU over traditional ReLU or LeakyReLU because it adapts the slope of the negative axis during training. Given that the dataset includes features like *pH* and *Turbidity*, which exhibit nonlinearities in both directions, PReLU proved essential. It allows the model to learn whether small negative values are informative, unlike ReLU, which silences them entirely.

#### **Benefits:**

- Prevents dying neurons
- Enhances convergence speed
- Improves learning capacity by addressing underfitting

### **Batch Normalization**

Each Dense layer is followed by Batch Normalization *before* the PReLU activation. This stabilizes and accelerates training by normalizing inputs to the activations.

#### **Benefits:**

- Mitigates internal covariate shift
- Improves gradient flow
- Reduces sensitivity to initialization
- Adds mild regularization, helping reduce overfitting

The placement before PReLU ensures that activations receive well-scaled, zero-centered data, maximizing their learning effectiveness.

### **Dropout (Rate = 0.2)**

To combat overfitting, I implemented a slow dropout rate, which was also aided by a regularizer. I did not want my model to underfit by not having enough neurons to learn.

#### **Effects:**

- Prevents the co-adaptation of neurons
- Effectively ensembles multiple subnetworks
- Encourages generalization

### **L2 Regularization ( $\lambda = 0.001$ )**

This term penalizes large weights during training and further combats overfitting.

**Works synergistically** with dropout, ensuring the model generalizes well beyond the training set. Additionally, the lambda value changed as my layers approached the output layers to help penalize weights in later stages of training.

### **Optimizer: Nadam (Nesterov + Adam)**

I opted for Nadam due to its combined strengths: the adaptive learning of Adam and the anticipatory updates of Nesterov momentum.

### **Result:**

Faster convergence and better final minima.

### **Loss Function: Binary Cross-Entropy**

With a binary classification task and sigmoid output, Binary Crossentropy is ideal for measuring the distance between predicted probabilities and true labels.

### **Training Strategy**

- **Epochs:** 300
- **Batch Size:** 48
- **Early Stopping:**
  - **Monitors:** val\_loss
  - **Patience:** 30 epochs
  - **Callback:** restore\_best\_weights=True ensures the final model reflects peak validation performance.

### **Model Performance Summary**

**Best Epoch:** 19 (Auto-selected by lowest validation loss)

**Train Accuracy:** 66.72%

**Validation Accuracy:** 67.41%

### **Test Evaluation Metrics:**

Metric	Value
Accuracy	66.87%
Precision	59.86%
Recall	45.83%
F1 Score	0.5192
AUC (TensorFlow)	0.6707
AUC (sklearn)	0.6609

### **Challenges faced:**

- **Class Imbalance:** The dataset had many more non-potable samples than potable ones. This made it tough for the model to correctly identify potable water, leading to lower recall and F1 scores, even after using stratified sampling.
- **Balancing Overfitting and Underfitting:** It was hard to find the right balance between the model memorizing the training data (overfitting) and not learning enough (underfitting). Getting the dropout rate, regularization strength, and network depth just right required a lot of trial and error.

- **Achieving High Accuracy:** Consistently getting test accuracy above 70% was a major hurdle. Even with advanced techniques, accuracy often plateaued in the high 60s. This was likely due to data noise, overlapping class distributions, and limited training data for the model's complexity.
- **Training Stability and Convergence:** Training a deep, seven-layer network with various regularization methods was tricky. While Batch Normalization and PReLU helped, issues like missing values, outliers, and small batch sizes sometimes made the learning process unstable. Early stopping was crucial to prevent overtraining but required careful tuning.

### Final Thoughts & Contributions:

This model helped me achieve strong generalization and stable learning. I attribute this to the following:

- **Advanced preprocessing:** MICE for imputation, IQR for outlier handling, and robust scaling
- **Regularization:** Dropout and L2 combined to combat overfitting
- **Adaptive activations:** PReLU captures nuanced nonlinearities
- **Normalization:** BatchNorm ensures stable signal propagation

Though training duration increased slightly due to network depth and early stopping criteria, the gains in performance and robustness more than justify the added complexity.

### Git logs(Gaius Irakiza):

```
commit 276496001f51e521f64c02140bc859dedb8ee6c2
Author: IrakizaGaius <111004500+IrakizaGaius@users.noreply.github.com>
Date: Sun Jun 8 17:00:40 2025 +0200

    model improvement and f1-score evaluation

commit 9fc6e03c33e341b60f0f248cafb8a7c7051e7250
Author: IrakizaGaius <111004500+IrakizaGaius@users.noreply.github.com>
Date: Sun Jun 8 10:56:01 2025 +0200

    Small fixes on Gaius_model

commit 80a06ba78fe764baea1a28a4ccad5e8a8d64ae62
Author: IrakizaGaius <111004500+IrakizaGaius@users.noreply.github.com>
Date: Sun Jun 8 10:53:51 2025 +0200

    Initiation of Gaius_model

commit b294cefe70815d265d0ef8d504c0f4a36d304fe6
Author: IrakizaGaius <111004500+IrakizaGaius@users.noreply.github.com>
Date: Sun Jun 8 10:52:53 2025 +0200

    Gaius_restructuring_the_model(Moving the model to a new file)

commit 74882c888c22a6201abafd8acee3f66b04b29247
Author: IrakizaGaius <111004500+IrakizaGaius@users.noreply.github.com>
Date: Sun Jun 8 00:28:36 2025 +0200

    new changes(Gaius_model)

commit 0616fccb724fff18223c0ba94a294159497c7d0a
Author: IrakizaGaius <111004500+IrakizaGaius@users.noreply.github.com>
Date: Sat Jun 7 18:50:12 2025 +0200

    First iteration (Gaius_model)

commit 46e35743ec474ca5d85036ea1622d6a9c848ebb
Author: IrakizaGaius <111004500+IrakizaGaius@users.noreply.github.com>
Date: Wed Jun 4 12:54:13 2025 +0200

    Initial commit

(END)
```

### Model Comparison 1: Gaius' Model vs. Jeremiah's Model

- **Architecture:**

I used a deeper, more advanced model with PReLU activations, Batch Normalization, Dropout, L2 regularization, and the Nadam optimizer. Jeremiah's model was simpler, with only Dense layers, L2 regularization, and Dropout.

- **Performance:**

Metric	Gaius	Jeremiah
Test Accuracy	66.87%	65.41%
Precision	59.86%	67.86%
Recall	45.83%	12.10%
F1 Score	0.5192	0.2054

My model had better recall, F1 score, and overall balance, making it more reliable in detecting true positives.

- **Key Advantages (Gaius):**

- Adaptive activation (PReLU) for nuanced feature learning
- BatchNorm for stable and fast training.
- Strong regularization to combat overfitting.
- Early stopping with best-weight restoration.

My model outperformed Jeremiah's by offering better generalization, robustness, and real-world applicability, particularly through its higher recall and F1 score.

### **Model Comparison 2: Gaius' Model vs. Tamanda's Model**

My model incorporates a more advanced architecture, using **7 Dense layers** with **Batch Normalization and Dropout**, and employs **PReLU and ReLU** for activation functions, alongside **L2 regularization**. This design leads to about **34,949 total parameters**. Tamanda's model, while also having **5 Dense layers with Dropout**, uses only **ReLU activations** and relies on **L1 regularization**, resulting in a more compact design with roughly **14,212 parameters** and no additional layers like Batch Normalization.

### **Performance Metrics on Test Set**

When it comes to actual performance on the test set, my model generally outperforms Tamanda's:

Metric	Gaius' Model	Tamanda's Model
Test Accuracy	0.6687	0.6494
Precision	0.5986	0.5476
Recall	0.4583	0.2930
F1 Score	0.5192	0.3817
AUC (TF)	0.6707	0.6487
Test Loss	0.6599	0.6587

As you can see, my model achieved higher scores across almost all key classification metrics. Specifically, my model shows a notable improvement in Recall (0.4583 vs. 0.2930) and F1 Score (0.5192 vs. 0.3817). This indicates that my model is better at identifying the minority class and offers a superior balance between precision and recall. Tamanda's model did show a marginally better test loss, however.

The improved performance of my model can be attributed to its architecture and regularization strategy. Batch Normalization and adaptive PReLU activations significantly enhance learning stability and robustness during training. While Tamanda's L1 regularization promotes simplicity, it appears to have limited the network's ability to fully capture the complex patterns in the data, especially concerning the minority class. Her model demonstrated higher specificity (better at identifying Class 0), but struggled more with identifying positive cases. My use of L2 regularization with a slightly more expressive structure allowed for greater learning flexibility, resulting in a better balance in class prediction.

**Model Comparison 3: Gaius' Model vs. David's Model**

My model's architecture consists of **7 Dense layers**, enhanced with **Batch Normalization** and **Dropout** after every hidden layer. I use a combination of **PReLU and ReLU** for activation functions, coupled with **L2 regularization (0.001)** and the **Adam optimizer**. This results in approximately **34,949 total parameters**. David's model, in contrast, is simpler with **4 Dense layers**, utilizing only **ReLU** activations and **L2 regularization (0.0005)**. Notably, David's model does not incorporate Dropout or Batch Normalization and uses the **RMSprop optimizer**, summing up to around **25,473 total parameters**.

My design prioritizes enhanced generalization and training stability through the inclusion of Batch Normalization and Dropout. While David's model is more streamlined, it lacks these features, which might impact its overall robustness on unseen data.

**Performance Metrics on Test Set**

Here's how our models performed on the test set:

Metric	Gaius' Model	David's Model
--------	--------------	---------------



Accuracy	0.6687	0.6800
Precision	0.5986	0.5827
Recall	0.4583	0.4713
F1 Score	0.5192	0.5211
AUC	0.6707	0.6710
Test Loss	0.6599	0.6773

### **Conclusion:**

David's model achieved slightly higher **accuracy (0.6800)**, **recall (0.4713)**, and **AUC (0.6710)**. This suggests that his model, despite its simplicity, demonstrates good generalization on new data and is slightly better at capturing positive cases. My model, on the other hand, shows superior **precision (0.5986)** and a lower **test loss (0.6599)**, indicating that when it does predict a positive class, it's more accurate, and its overall error margins are tighter. The F1 scores are quite close, with David's model having a marginal edge, suggesting a similar overall balance between precision and recall for both our designs.

My model's use of **Dropout, Batch Normalization, and a slightly stronger L2 regularization** likely contributed to reduced overfitting and greater training stability, though it did introduce more architectural complexity. David's model, being simpler, appears to have converged more rapidly (achieving its best performance around Epoch 16, whereas my model trained longer, potentially up to 200 epochs if not for early stopping. This simplicity could make it more efficient in terms of computational cost.

When examining minority class handling, my model exhibits a slight edge in **precision**, meaning it generates fewer false positives. David's model, however, shows marginally better **recall**, indicating it successfully identifies a slightly higher proportion of actual positive cases.

### **Model Comparison 4: Gaius' Model vs. Rene's Model**

My model's architecture is built with **7 Dense layers**, utilizing **PReLU (adaptive) and Sigmoid** activations. It incorporates **Batch Normalization** after each Dense layer and **Dropout (0.2 per layer)**, along with **L2 regularization ( $\lambda = 0.001$ )**. I use the **Nadam optimizer**, resulting in approximately **34,949 total parameters**. Rene's model, on the other hand, is a more compact design with **4 Dense layers**, using **ReLU and Sigmoid** activations. He employs **Dropout (0.3–0.4 per layer)** and **L1 regularization (varying  $\lambda$ : 0.0001–0.00001)**, relying on the **Adam optimizer** and totaling around **15,000 parameters**.

My design emphasizes architectural depth and advanced normalization to enhance learning dynamics and convergence stability, aiming to capture more complex patterns. Rene's model is deliberately more lightweight, prioritizing efficiency and relying on stronger regularization to mitigate overfitting.

### **Performance Metrics on Test Set:**

Metric	Gaius' Model	Rene's Model
Accuracy	0.6687	0.6800
Precision	0.5986	0.5827
Recall	0.4583	0.4713
F1 Score	0.5192	0.5211
Test Loss	0.6599	0.6403

Rene's model achieved a slightly higher **accuracy (0.6800)**, **recall (0.4713)**, and a marginally better **F1 score (0.5211)**. This suggests his model has a slight edge in overall balance between precision and recall, and is a bit more effective at identifying positive cases. My model, however, demonstrates superior **precision (0.5986)** and a slightly lower **test loss (0.6599)**.

### **Conclusion:**

My model's approach to **overfitting mitigation** combines L2 regularization, Dropout, and Batch Normalization. This allows my deeper network to learn robustly. Rene tackles this differently, employing heavier Dropout and L1 regularization, which aims for model simplicity by promoting sparsity. While effective, this can sometimes limit learning capacity if not finely tuned. For **training stability**, I utilized Nadam, a momentum-based optimizer often suited for complex networks, while Rene opted for Adam, which is simpler but effective for his more compact architecture.

## **Tamanda Kaunda Model Analysis**

### **Detailed Metric Interpretation and Comparison**

#### **1. Loss Analysis**

My model's performance: 0.6587 is the lowest( outperforms ) indicating reasonable model confidence.

**Comparison:** My loss outperforms both Gaius (0.6599) and David (0.6773).

**Interpretation:** My model achieves the lowest loss in the group, suggesting better calibrated predictions despite lower recall.

#### **2. Accuracy analysis**

My performance is 68.94% compared to Gaius' 66.87% ( Higher) and David's 68.47% ( the highest)

**Interpretation:** My model achieves an accuracy that competitively classified nearly 2 out of 3 water samples correctly.

### **3. Precision Analysis**

My model performance was 54.67% - it predicts that water is portable , 54.67% percent of the time it's correct.

**Comparison** ; 5.10% lower than Gaius' 59.86% and 5.38% higher than David's 58.27%.

**Interpretation:** My model achieves moderate precision , balancing between both my teammates approaches.

### **4. Recall Analysis**

**My performance** : My model identifies 29.30% of actually portable water samples.

**Comparison:** Significantly lower than both Gaius ( 45.83%) and David ( 47.13%)

**Interpretation:** my model is the most conservative among the 3 , missing more portable samples but minizes false positives.

### **5. F1 Score Analysis**

My model performance is at 0.3825

**Comparison:** Lower than both Gaius ( 0.5192) and David ( 0.5211).

**Interpretation:** My model's conservative approach impacts its F1 score, showing a trade-off between precision and comprehensive detection.

## **Model Strengths and weaknesses**

### **Two Key strengths compared to teammates**

#### **1. Best loss performance in comparison to Gaius and David**

My model outperformed both teammates in terms of loss compared to both Gaius (0.6599) and David ( 0.6773) which indicates superior probability calibration. This also suggests that my L1 and RMSprop combination produces more confident predictions when it does clarify water as portable which is valuable and reliable water safety applications where reliable risk assessment is crucial.

#### **2. Most Safety\_Critical Approach for Water Quality**

With the lowest recall( 29.30%) among the three models but competitive precision (54.76%), my model takes the most cautious approach in the group. This makes it ideal for screening in water safety applications where false positives (incorrectly classifying unsafe water as safe ) could have serious health consequences.

#### **3. Superior Feature Selection through L1 Regularization Sparsity**

My L1 regularization approach creates feature sparsity that significantly differs from Gaius' L2 Batch normalization strategy. While his regularizer shrinks all weights towards zero without eliminating less important features and batch normalization stabilizes training, my L1 approach drives specific weights to exactly zero, effectively performing automatic feature selection. This means my model identifies and focuses only on the critical water quality parameters while completely eliminating less important ones. My model provides true interpretation of by showing specific water quality indicators (pH, hardness, chloramines, etc) are most predictive of portability.

## Two areas of improvement

### 1. Lower recall compared to both teammates

My model's recall (29.30%) is significantly lower than both Gaius (45.83%) and 42.04%) which shows that it misses many portable water samples.

### 2. Lower F1 score points at suboptimal Precision-Recall Balance

My F1 score (0.3825) is lower than both Gaius (0.5192) and David (0.4673), showing that my model's extreme tradeoff between precision and recall may not be optimal for general purpose water quality assessment. Finding a better balance could improve overall utility.

## Regularization Strategy Comparison Table

approach	Regularization type	effect	interpretability	performance
Tamanda	L1 (progressive)	Creates sparsity, eliminates features	High- drives weight to exactly zero in feature selection	Lowest Loss(0.6587)
Gaius	L2 + Batch Norm	Shrinks all weights, stabilizes training	Moderate	High Accuracy(66.87%)
David	L 2	Shrinks all weights to zero but never eliminates them completely	Moderate- all features contribute	Highest accuracy

## Optimization Strategy Effectiveness

The comparison highlights how different optimizers affect model behavior:

- **Nadam** (Gaius): Achieves highest accuracy and balanced metrics
- **RMSprop** (Tamanda): Produces lowest loss and most conservative predictions.
- **RMSprop optimizers** (David): Achieves good recall and produces the highest accuracy.

## David Ubushakebwimana Model Analysis

This report presents my individual contribution and model performance for the water potability prediction task in Peer Group 8. My approach leveraged a simple, efficient architecture using RMSprop optimization and L2 regularization. Below is a summary of my model design, its strengths, and a quantitative comparison to Gaius and Tamanda's models.

### Model Overview

My model consists of a feedforward neural network with the following configuration:

- Layers: Dense(128) → Dense(64) → Dense(32) → Dense(1)
- Activations: ReLU (hidden), Sigmoid (output)
- Regularization: L2 ( $\lambda = 0.0005$ )
- Optimizer: RMSprop (learning rate = 0.001)
- No Dropout or Batch Normalization
- Early Stopping: patience=50, monitor=val\_loss

### Performance Summary

Metrics from test evaluation and best epoch during training:

- Best Epoch: 24
- Train Accuracy at Best Epoch: 66.32%
- Validation Accuracy at Best Epoch: 68.05%
- Test Accuracy: 67.83%
- Precision: 59.92%
- Recall: 46.87%
- F1 Score: 0.5240
- Test Loss: 0.6494

## Model Comparison (David vs. Gaius vs. Tamanda)

<b>Metric</b>	<b>David (Mine)</b>	<b>Gaius</b>	<b>Tamanda</b>	<b>Jeremiah</b>
<b>Test Accuracy</b>	67.83%	66.87%	68.94%	65.41%
<b>Precision</b>	59.92%	59.86%	54.76%	67.86%
<b>Recall</b>	46.87%	45.83%	29.30%	12.10%
<b>F1 Score</b>	0.5240	0.5192	0.3817	0.2054
<b>Test Loss</b>	0.6494	0.6599	0.6587	-

## Conclusion

Compared to both Gaius' and Tamanda's models, my design achieved strong balance across all key classification metrics. It surpassed Gaius in accuracy, recall, F1, and AUC while maintaining a simpler architecture. Although Tamanda achieved slightly better test accuracy, her model lagged in recall and F1 score, indicating a more conservative prediction approach. My model offers a balanced and practical solution with fewer architectural complexities and robust generalization.

## Design Philosophy & Preprocessing Strategy

Given the nature of the water potability dataset—relatively small with missing values and class imbalance—the model needed to be both resilient and generalizable. I opted for a streamlined architecture without Batch Normalization or Dropout to focus instead on consistent regularization through L2, allowing the model to avoid overfitting while keeping training stable.

Data preprocessing followed a rigorous strategy similar to my peers:

- **Missing Values:** I used MICE imputation (IterativeImputer), preserving multivariate structure.
- **Outliers:** Applied IQR clipping to continuous features.
- **Scaling:** StandardScaler normalized all features to zero mean and unit variance.

- **Splitting:** Stratified 70/15/15 split ensured balanced class distributions across train, val, and test.

## **Model Training & Early Stopping Approach**

The training process relied heavily on early stopping to prevent overfitting, especially given the small dataset size. I monitored `val_loss` and used a patience of 50 epochs to let the model explore improvements without premature termination. This helped locate the optimal convergence point without exhausting all 300 epochs.

Unlike some peers who used batch sizes of 48 or smaller layers, I chose a **batch size of 32** to balance gradient estimation quality and update frequency. This aided convergence without incurring the instability of extremely small batches.

## **Technical Comparison & Architecture Insights**

The following observations expand on the earlier comparison table, offering technical insights into the strengths and tradeoffs of each peer model.

- **Gaius:** Used PReLU activations and BatchNorm extensively, which increased convergence stability but added complexity and training time. My model, by comparison, converged earlier with fewer layers.
- **Tamanda:** Applied L1 regularization and achieved the lowest test loss, but had the lowest recall and F1 score. Her model's conservative nature missed too many positive (potable) cases.
- **Mine:** Balanced recall and precision well, suggesting a stronger real-world performance in detecting positive potability while maintaining good precision.

## **Reflection & Future Directions**

This exercise underscored the importance of model simplicity, clean data, and regularization. While many peers used deeper or more regularized models, my approach—guided by empirical testing of learning rates, regularization strength, and architecture depth—achieved a practical balance between performance and interpretability.

Going forward, I'd consider augmenting the training set with synthetic sampling techniques such as SMOTE to address class imbalance and boost recall further. I would also experiment with ensemble techniques to average predictions from simpler and more complex models, possibly gaining robustness without sacrificing interpretability.

## **Rene Ntabana Model Analysis**

This report outlines my contribution to the Peer Group 8 water potability prediction task. I designed and implemented a neural network architecture that emphasizes regularization, moderate dropout, and careful tuning of the learning rate. My approach aimed to balance bias and variance while maintaining simplicity.

## Model Overview

**My neural network consists of the following configuration:**

- Layers: Dense(128) → Dense(64) → Dense(32) → Dense(1)
- Activations: ReLU (hidden layers), Sigmoid (output)
- Regularization: L1 ( $\lambda$  values ranging from 0.0001 to 0.00001)
- Dropout: Applied after each layer (0.4 / 0.4 / 0.3)
- Optimizer: Adam (learning rate = 0.0005)
- Early Stopping: Patience = 30, monitored `val_loss`

This structure was chosen to prevent overfitting while preserving performance on unseen data.

## Performance Summary

Metrics from test evaluation and the best training epoch are shown below:

- Best Epoch: 124
- Train Accuracy at Best Epoch: 71.78%
- Validation Accuracy at Best Epoch: 66.76%
- Test Accuracy: 70.12%
- Precision: 58.2%
- Recall: 48.70%
- F1 Score: 0.5211
- Test Loss: 0.6403



---

### Model Comparison (Rene vs. David vs. Gaius vs. Tamanda)

Metric	Rene (Mine)	David	Gaius	Tamanda
Test Accuracy	70.12%	67.83%	66.87%	68.94%
Precision	58.20%	59.92%	59.86%	54.76%
Recall	48.70%	46.87%	45.83%	29.30%
F1 Score	0.5280	0.5240	0.5192	0.3817
Test Loss	0.6403	0.6494	0.6599	0.6587

---

### Conclusion

My model achieved the **highest test accuracy and F1 score** among all peer submissions, demonstrating effective generalization and a strong balance between precision and recall. The use of **L1 regularization** and **moderate dropout** helped prevent overfitting while still allowing the model to learn useful patterns. Compared to Tamanda’s model, which favored precision at the cost of recall, my model offered a more balanced trade-off, making it better suited for practical applications where both false positives and false negatives have real-world implications.

However, a potential **weakness** of my model is its limited ability to detect some positive cases (i.e., lower recall than desired). This could stem from strong regularization and dropout suppressing learning for minority class examples.

---

### Reflection & Future Directions

This exercise underscored the value of **clean architecture, focused regularization, and hyperparameter tuning**. While some peers implemented deeper networks or more aggressive dropout, my model showed that a carefully constrained architecture can achieve competitive and even superior results.

Moving forward, I would explore techniques such as **SMOTE** or **adaptive class weighting** to better address class imbalance and improve recall. Additionally, I’d consider experimenting with **ensemble methods** that blend predictions from multiple models—balancing interpretability with predictive strength. These adjustments could further boost the model’s robustness, especially when deployed in real-world water safety monitoring systems.