# A Technical Guide to Machine Learning Models for Tuna Biomass Estimation

**Author**: BenchChem Technical Support Team. **Date**: December 2025

| Compound of Interest | | |
| --- | --- | --- |
| Compound Name: | Tuna AI | |
| Cat. No.: | B1682044 | Get Quote |

This guide provides a comprehensive overview of the application of machine learning (ML) for estimating tuna biomass, a critical component of sustainable fisheries management. Traditional stock assessment methods face numerous challenges, but the advent of AI and ML offers powerful tools to process vast and complex datasets, enhancing the accuracy and predictive capabilities of biomass estimations.[1] This document details the data sources, experimental protocols, and performance of various ML models, intended for researchers and scientists in marine biology, fisheries science, and data science.

## Data Sources and Feature Engineering

The foundation of any successful machine learning model is the data it is trained on. For tuna biomass estimation, data is typically drawn from three primary sources: fishery operations, echosounder buoys, and oceanographic satellites.[2][3]
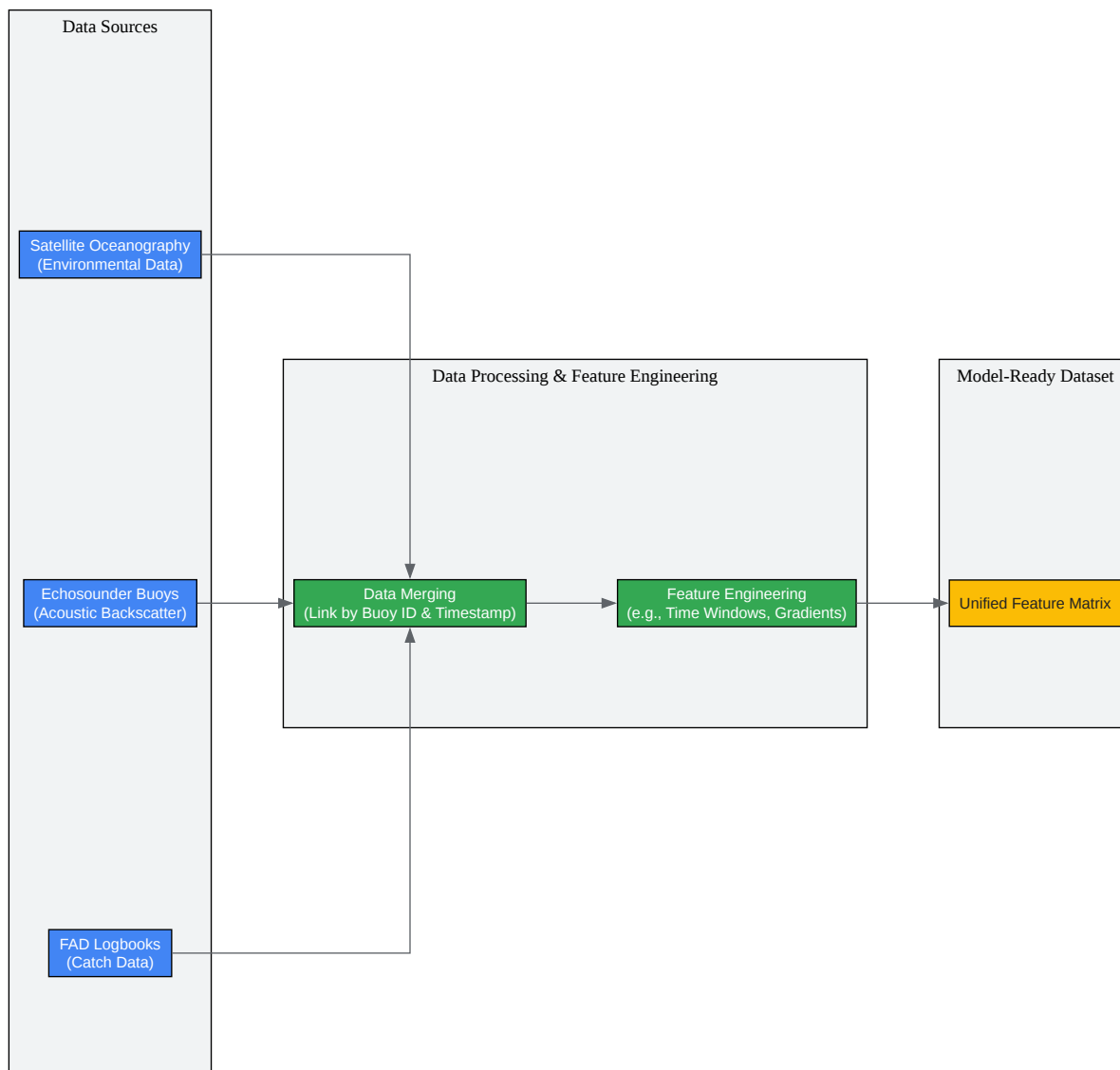
Key Data Sources:

- Fishery-Dependent Data: This includes logbooks from purse-seine fishing vessels and data from observers on board. A crucial component is catch data from fishing "sets" on drifting Fish Aggregating Devices (dFADs), which provides the ground-truth (supervised signal) for training the models.[2][4][5][6]

- Echosounder Buoy Data: Modern dFADs are equipped with satellite-linked echosounder buoys that provide frequent, geo-referenced estimates of fish biomass aggregated beneath

them.[3][4][6] This raw acoustic backscatter data, converted into biomass by manufacturer algorithms, is a primary input for ML models.[2][7]

- Oceanographic Data: Satellite remote sensing provides a wealth of environmental data. These variables are critical as they describe the habitat and ecological conditions that influence tuna distribution and aggregation.[2][8]

The diagram below illustrates the typical workflow for integrating these diverse data sources into a cohesive dataset for model training.

Click to download full resolution via product page

Data Integration Workflow for Tuna Biomass Estimation.

A summary of common features, also known as explanatory variables, used in these models is presented in the table below.

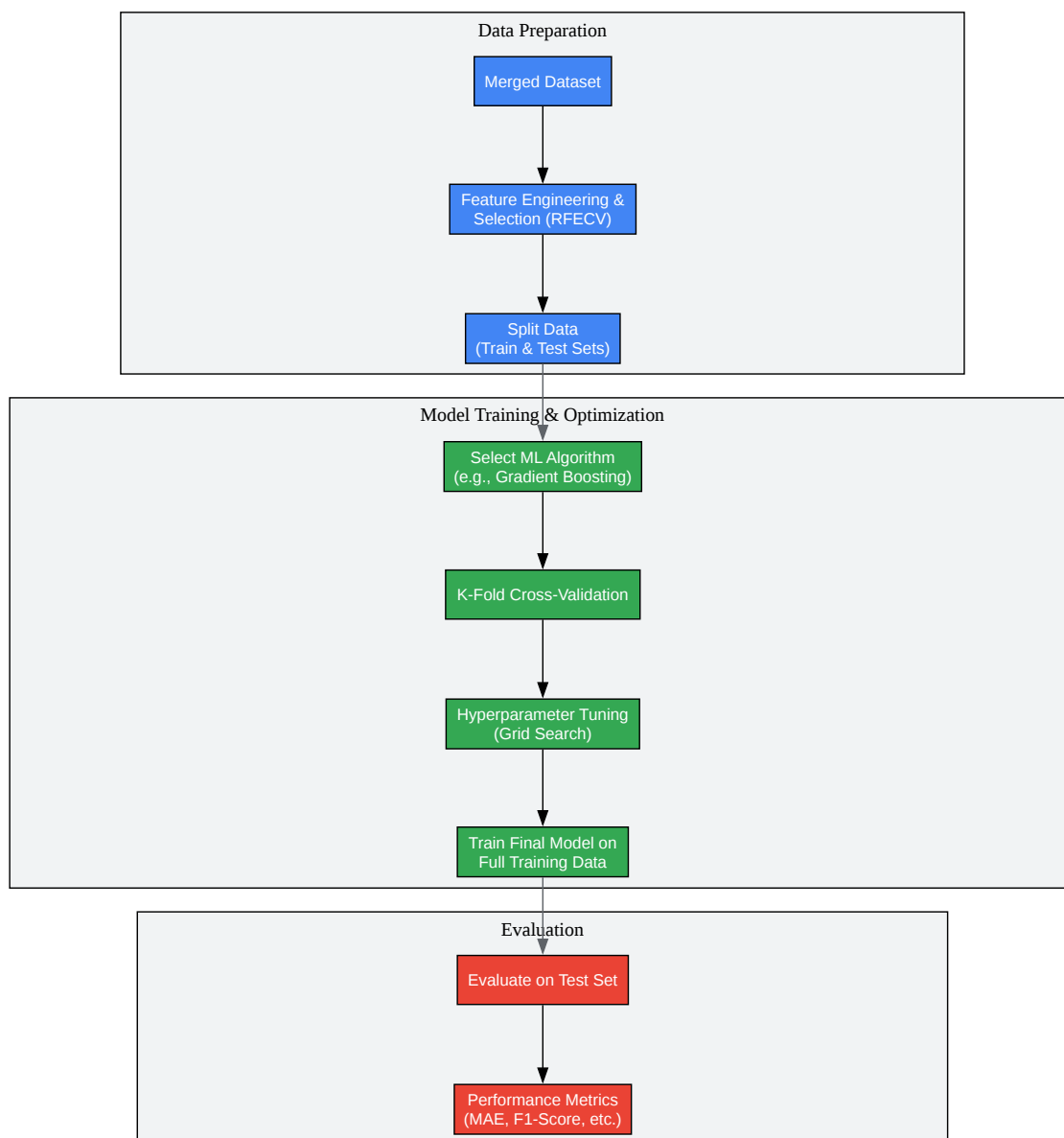| Data Category | Feature | Description |
|---|---|---|
| Echosounder | Biomass Estimates | Time-series of biomass values derived from acoustic backscatter, often aggregated over a 3-day window.[3][4][5][7] |
| Oceanographic | Sea Surface Temperature (SST) | A key environmental factor influencing the metabolic rate and distribution of tuna.[8][9] |
| Chlorophyll-a Concentration (Chl-a) | An indicator of phytoplankton abundance, forming the base of the marine food web.[8][9] | |
| Sea Surface Height / Anomaly (SSH/SLA) | Indicates ocean currents and eddies, which can aggregate nutrients and prey.[8] | |
| Salinity | Affects water density and ocean circulation, influencing tuna habitat suitability.[8][10] | |
| Dissolved Oxygen | Critical for tuna respiration, especially at depth.[8] | |
| Ocean Current Velocity | Influences the drift of dFADs and the movement of tuna and their prey.[10] | |
| Spatiotemporal | Latitude & Longitude | The geographical coordinates of the dFAD buoy at the time of measurement.[2][9] |
| Time-derived Features | Year, month, and day to capture seasonal and long-term patterns.[9] | |
| Climate Indices | ONI, NPGIO, etc. | Large-scale climate patterns (e.g., El Niño) that affect ocean conditions globally.[9][11] |

# Experimental Protocols and Methodologies

A robust and reproducible experimental protocol is essential for developing reliable machine learning models. The process involves several key stages, from data preparation to model evaluation.

Detailed Methodologies:

- Data Preprocessing and Merging: The initial step involves linking records from the different data sources. Echosounder buoy data is merged with FAD logbook events using the unique buoy ID and timestamp.[7] Oceanographic data is then appended based on the GPS coordinates and date of each echosounder record.[2][7]

- Feature Engineering: Raw data is transformed into meaningful features. A critical technique is the use of a time window (e.g., 24, 48, or 72 hours) of echosounder data preceding a fishing event.[7] This captures the temporal dynamics of tuna aggregation. Other engineered features can include the temperature of the previous and subsequent months to capture broader trends.[9]

- Feature Selection: With a high number of potential variables, it's important to select the most informative ones. Techniques like Recursive Feature Elimination with Cross-Validation (RFECV) can systematically identify the optimal combination of variables, improving model performance and interpretability.[8]

- Dataset Splitting: The complete dataset is typically divided into a training set (e.g., 75-80% of the data) and a testing set (20-25%).[6][12] The model learns patterns from the training data, and its performance is evaluated on the unseen testing data.

- Model Training and Validation: The selected ML algorithm is trained on the training dataset. To ensure the model generalizes well and avoids overfitting, k-fold cross-validation (commonly 10-fold) is employed.[12] This involves repeatedly training and validating the model on different subsets of the training data.

- Hyperparameter Tuning: The performance of many ML models is sensitive to their internal settings, known as hyperparameters. Techniques like grid search are used to systematically test various combinations of these settings to find the optimal configuration.[13]

- Model Evaluation: The final, tuned model is evaluated on the held-out test set. The choice of performance metric depends on the specific task.

  - Regression Task (Direct Biomass Estimation): Metrics include Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and Symmetric Mean Absolute Percentage Error (SMAPE).[4][12]

  - Classification Task (e.g., Low vs. High Biomass): Metrics include F1-Score, Accuracy, and Mean Average Precision (mAP).[4][14]

The following diagram outlines this complete experimental workflow.

**Data Preparation**

Merged Dataset

Feature Engineering &
Selection (RFECV)

Split Data
(Train & Test Sets)

**Model Training & Optimization**

Select ML Algorithm
(e.g., Gradient Boosting)

K-Fold Cross-Validation

Hyperparameter Tuning
(Grid Search)

Train Final Model on
Full Training Data

**Evaluation**

Evaluate on Test Set

Performance Metrics
(MAE, F1-Score, etc.)
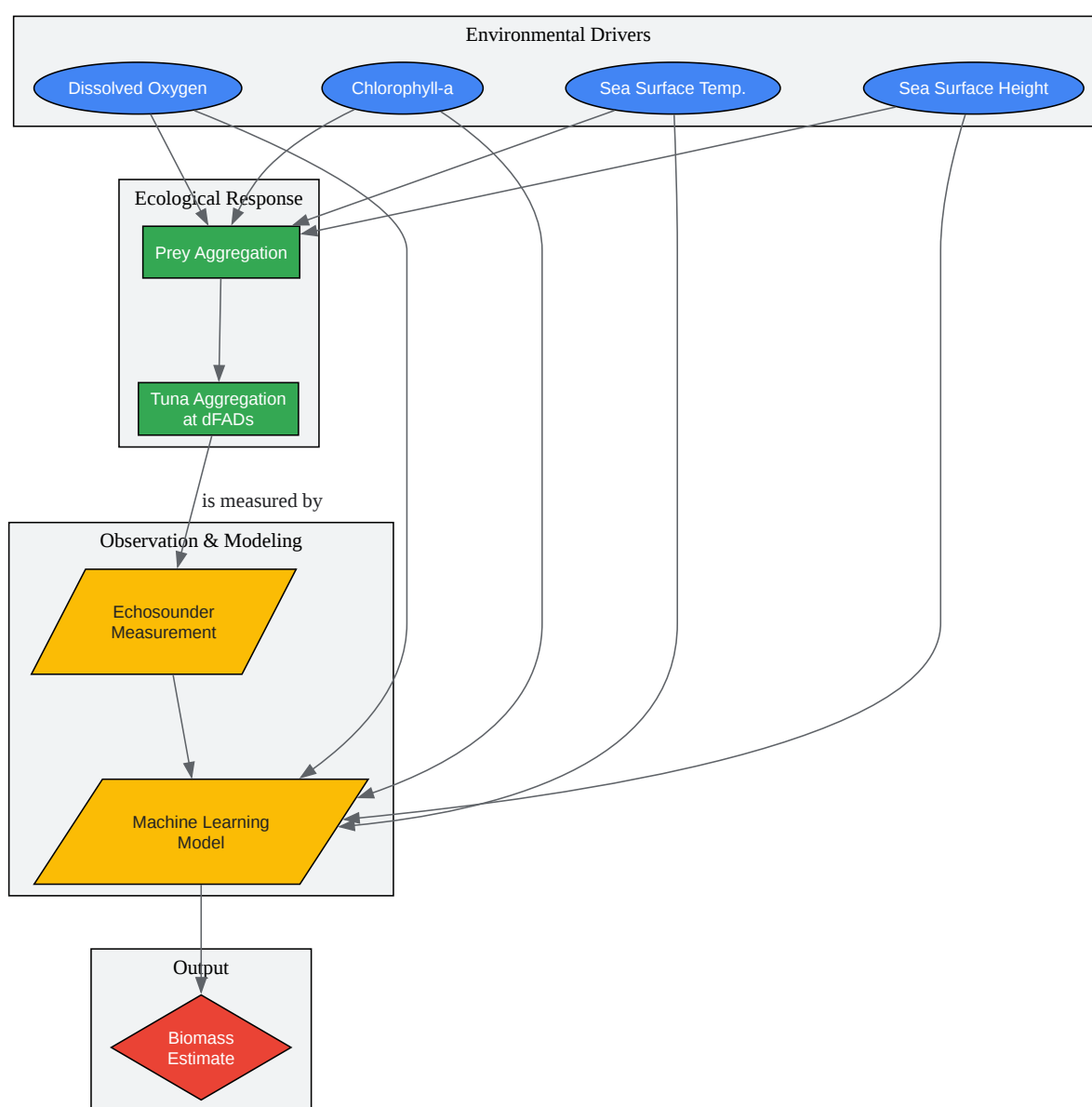
Click to download full resolution via product page

Machine Learning Experimental Protocol Workflow.

# Machine Learning Models and Performance Comparison

Several machine learning algorithms have been successfully applied to the problem of tuna biomass estimation. Tree-based ensemble methods like Gradient Boosting and Random Forest are frequently used and often yield the best results.[2][4][8]

The logical relationship between environmental factors, tuna behavior, and the ML prediction process is visualized below.

Logical Relationships in Tuna Biomass Estimation.

The models are typically configured for one of two primary tasks: regression (predicting the exact biomass) or classification (predicting a biomass category).

| Model | Task Type(s) | Key Features Used | Performance Metrics & Results |
|---|---|---|---|
| Gradient Boosting (GB) | Regression, Classification | 3-day echosounder window, oceanographic data, position/time features.[4] | Regression: MAE = 21.6 t, SMAPE = 29.5%.[4] Binary Classification (>10t): F1-Score = 0.925.[4] |
| Random Forest (RF) | Regression, Classification | Fisheries data, sea temperature, dissolved oxygen, chlorophyll-a, salinity, SSH.[8] | Often used as a robust baseline or primary model.[2][8] In one study, RF was used to explore the changing mechanisms of catch composition. [15] |
| Neural Networks (ANN) | Regression, Classification | Environmental variables (water movement, stream size, water chemistry). [16] | Can achieve high accuracy (e.g., >84% for salmonid abundance) and identify variables with the greatest predictive power.[16] Used to detect complex patterns in fish population dynamics. [1] |
| Support Vector Machine (SVM) | Classification, Regression | Shape features, environmental data. [17] | Used across various fisheries applications, including species classification and predicting potential fishing zones.[1][13] [18] |

| | | | Serve as simpler baseline models to compare against more complex algorithms like Gradient Boosting and Random Forest. [2] |
|---|---|---|---|
| Linear Models (e.g., Elastic Net) | Regression | Echosounder, oceanographic, and positional data.[2] | |

It is consistently noted that models enriched with oceanographic and position-derived features show improved performance over models that use echosounder data alone, highlighting the importance of environmental context.[4]

# Conclusion and Future Directions

Machine learning models, particularly ensemble methods like Gradient Boosting, have demonstrated significant success in estimating tuna biomass around dFADs. By integrating data from echosounder buoys, vessel logbooks, and satellite oceanography, these models provide powerful tools for fisheries monitoring and management.

Future improvements may involve incorporating more diverse data sources, such as information on bycatch or the species composition of tuna schools, which can impact the acoustic properties measured by echosounders.[6] Additionally, hybrid models that combine the mechanistic understanding of fish behavior with data-driven ML approaches hold promise for improving forecast accuracy, especially in the context of a changing climate.[19][20] The continued development and application of these advanced analytical techniques are vital for ensuring the sustainable exploitation of global tuna resources.

> **Need Custom Synthesis?**
>
> *BenchChem offers custom synthesis for rare earth carbides and specific isotopiclabeling.*
>
> *Email: info@benchchem.com or Request Quote Online.*

# References

- 1. taylorfrancis.com [taylorfrancis.com]

- 2. arxiv.org [arxiv.org]

- 3. biorxiv.org [biorxiv.org]

- 4. researchgate.net [researchgate.net]

- 5. [2109.06732] Tuna-AI: tuna biomass estimation with Machine Learning models trained on oceanography and echosounder FAD data [arxiv.org]

- 6. biorxiv.org [biorxiv.org]

- 7. researchgate.net [researchgate.net]

- 8. cdnsciencepub.com [cdnsciencepub.com]

- 9. mdpi.com [mdpi.com]

- 10. scispace.com [scispace.com]

- 11. researchgate.net [researchgate.net]

- 12. mdpi.com [mdpi.com]

- 13. academic.oup.com [academic.oup.com]

- 14. research-repository.griffith.edu.au [research-repository.griffith.edu.au]

- 15. researchgate.net [researchgate.net]

- 16. academic.oup.com [academic.oup.com]

- 17. Fish Classification Using Support Vector Machine | Semantic Scholar [semanticscholar.org]

- 18. tandfonline.com [tandfonline.com]

- 19. [2308.03403] Towards Machine Learning-based Fish Stock Assessment [arxiv.org]

- 20. openreview.net [openreview.net]

- To cite this document: BenchChem. [A Technical Guide to Machine Learning Models for Tuna Biomass Estimation]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b1682044#machine-learning-models-for-tuna-biomass-estimation]

---

**Disclaimer & Data Validity:**

The information provided in this document is for Research Use Only (RUO) and is strictly not intended for diagnostic or therapeutic procedures. While BenchChem strives to provide

Tech Support

accurate protocols, we make no warranties, express or implied, regarding the fitness of this product for every specific experimental setup.

**Technical Support:**The protocols provided are for reference purposes. Unsure if this reagent suits your experiment? [Contact our Ph.D. Support Team for a compatibility check]

**Need Industrial/Bulk Grade?**   Request Custom Synthesis Quote

# BenchChem

Our mission is to be the trusted global source of essential and advanced chemicals, empowering scientists and researchers to drive progress in science and industry.

Contact

Address: 3281 E Guasti Rd

Ontario, CA 91761, United States

Phone: (601) 213-4426

Email: info@benchchem.com