# "WDAS: Database Architecture and Implementation for Scientific Computing in Behavioral Ecology"

IST 659: Database Concepts, Prof. Mark Romano, December 11, 2024
Project Group Chi: Gil Raitses, Yifeng Qiu and Dawryn Rosario

## Abstract

The WDAS database system is designed to manage scientific computing workflows in behavioral ecology, focusing on the analysis of large-scale biologging datasets. Initially developed to support research on humpback whale feeding grounds and their relationship to human activity, the system integrates data processing, behavior classification and visualization workflows. The database architecture emphasizes traceability, scalability and seamless integration with physics-informed machine learning pipelines. Key features include efficient batch processing, robust annotation management and support for iterative research processes where hypotheses evolve as new data is analyzed. This report explores the system's backend architecture, design principles and scalability, demonstrating its potential to address the demands of large-scale ecological research and its applications in broader scientific domains.

## Keywords

WDAS
Scientific computing workflows
Behavioral ecology
Biologging data
Database architecture
Data traceability
Machine learning integration
Conservation research
Large-scale data analysis
Humpback whales
Dive statistics
Environmental metrics

# 1. Introduction and Context

The Parks Lab at Syracuse University identified the need for a comprehensive system to manage scientific computing workflows in behavioral ecology research [Hey et al., 2020]. The initial use case centers on tracking the evolution of humpback whale feeding grounds and their relationship to human activity. This effort is supported by machine learning analysis of biotelemetry data archived at the Woods Hole Oceanographic Institution in Massachusetts. By analyzing biologging data [Johnson & Tyack, 2003], including dive statistics, environmental conditions and movement patterns [Goldbogen et al., 2011], the project aims to generate actionable insights to inform conservation policies and address ecological challenges.

The database at the core of this system is designed to process large-scale datasets efficiently through batch operations [Jensen et al., 2019]. Its role extends beyond data storage, providing robust mechanisms for ensuring data integrity and maintaining traceability across all stages of the analysis pipeline [Davidson & Freire, 2018]. These features are critical for iterative research workflows, where data processing often informs new hypotheses or requires revisiting earlier stages of analysis. The ability to audit or reverse workflows ensures that the system supports reproducibility and accountability in scientific research.

In addition to managing data, the database facilitates the generation and validation of behavior classifications, correlating biologging metrics with environmental parameters to uncover ecological patterns. It provides the foundation for visualizing results and interpreting data, addressing the need for scalable and transparent management of ecological data. The system is built to handle extensive sensor recordings and biologging inputs across multiple individuals, ensuring scalability for research applications that demand high-throughput data analysis.

This report delves into the backend architecture and implementation of the database system powering the WDAS platform. It examines the design principles that guided the system's development, the mechanisms that support its functionality and its scalability for broader scientific computing applications.

# 2. Database Design and Architecture

The database architecture integrates multiple key entities to support workflows such as dataset imports, annotation management, behavior classifications and data visualization. This section introduces both the conceptual model (Figure 1), which provides an abstract overview of the system and the logical model (Figure 2), which specifies the schema's implementation details. Together, these models guide the design and ensure that the database meets the requirements of scalability, traceability and robust data management.

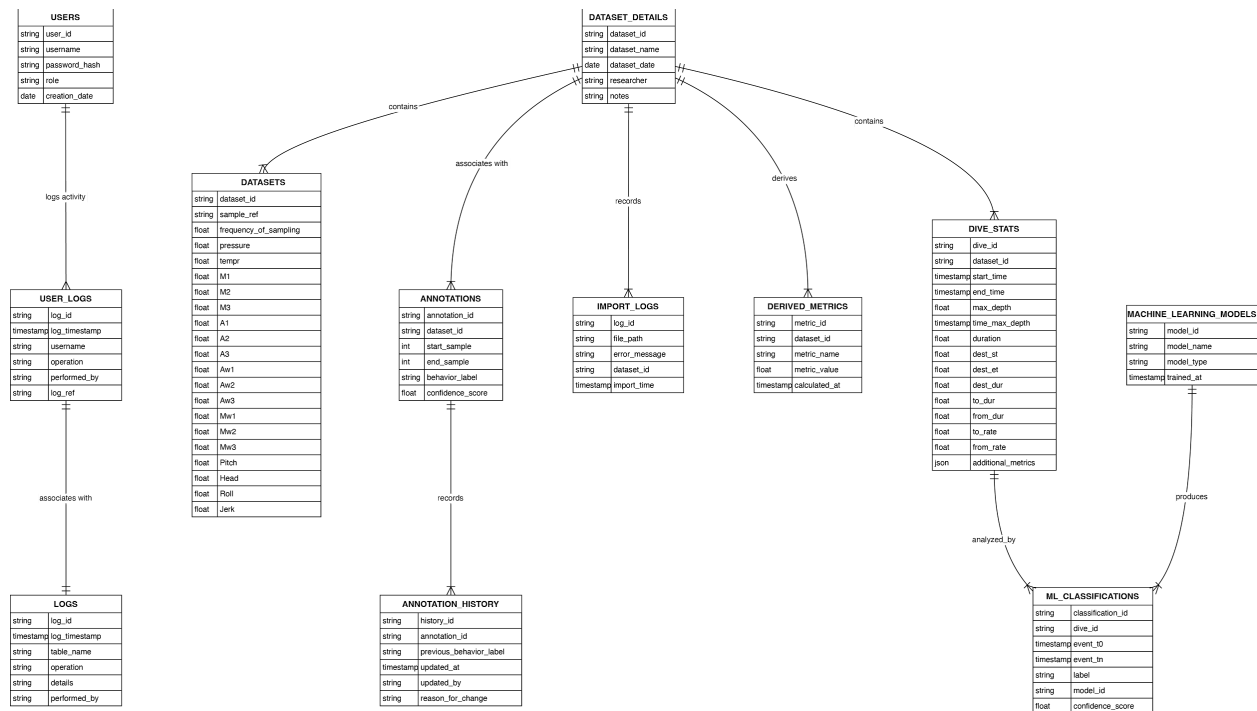Figure 1: Conceptual Model of the WDAS Database Schema.

*Figure 1: Conceptual Model of the WDAS Database Schema. This figure outlines the high-level entities, relationships and workflows within the WDAS database, providing an abstract overview of the system. It highlights the major components, such as USERS, DATASET_DETAILS and ANNOTATIONS and their integration into workflows for data imports, behavior classification and annotation management.*

## 2.1 Entity Relationships

Building on the conceptual model, the logical data model (Figure 2) provides a detailed schema of the WDAS database. It specifies the attributes, foreign key constraints and table relationships that implement the database's core functionality. This model translates the abstract relationships of the conceptual model into a fully normalized structure, ensuring entity integrity, efficient querying and support for high-throughput workflows.

The logical data model defines several key relationships:

The ANNOTATIONS table links to the DATASET_DETAILS table via dataset_id, enabling annotations to be associated with specific datasets.

The MACHINE_LEARNING_MODELS table references the ML_CLASSIFICATIONS table to store results of behavior classifications.

Foreign key constraints in tables such as DIVE_STATS and DERIVED_METRICS ensure traceability and validation across raw and processed data.

**MACHINE_LEARNING_MODELS**

| | | | |
|---|---|---|---|
| int | model_id | PK | |
| string | model_name | | NOT NULL |
| string | model_type | | NOT NULL |
| timestamp | trained_at | | DEFAULT CURRENT_TIMESTAMP |

**ML_CLASSIFICATIONS**

| | | | |
|---|---|---|---|
| int | classification_id | PK | |
| int | dive_id | FK | NOT NULL |
| bigint | event_t0 | | NOT NULL |
| bigint | event_tn | | NOT NULL |
| string | label | | NOT NULL |
| int | model_id | FK | NOT NULL |
| float | confidence_score | | NOT NULL |

**BEHAVIOR_TYPES**

| | | | |
|---|---|---|---|
| int | behavior_type_id | PK | |
| string | label | | UNIQUE NOT NULL |

**USERS**

| | | | |
|---|---|---|---|
| int | user_id | PK | |
| string | username | | UNIQUE NOT NULL |
| string | password_hash | | NOT NULL |
| string | password_hash | | NOT NULL |
| enum | role | | DEFAULT Viewer |
| timestamp | creation_date | | DEFAULT CURRENT_TIMESTAMP |

**ANNOTATIONS**

| | | | |
|---|---|---|---|
| int | annotation_id | PK | |
| int | dataset_id | FK | NOT NULL |
| int | start_sample | | NOT NULL |
| int | end_sample | | NOT NULL |
| string | classification_label | | NOT NULL |
| string | behavior_type | FK | |
| float | confidence_score | | |

**ANNOTATION_HISTORY**

| | | | |
|---|---|---|---|
| int | history_id | PK | |
| int | annotation_id | FK | NOT NULL |
| string | previous_behavior_label | | |
| string | previous_behavior_label | | |
| string | previous_behavior_label | | |
| string | previous_behavior_label | | |
| timestamp | updated_at | | DEFAULT CURRENT_TIMESTAMP |
| string | updated_by | | |
| text | reason_for_change | | |

**IMPORT_LOGS**

| | | | |
|---|---|---|---|
| int | log_id | PK | |
| string | file_path | | NOT NULL |
| text | error_message | | |
| int | dataset_id | FK | |
| int | dataset_id | FK | |
| int | performed_by | FK | |
| int | log_ref | | |
| timestamp | import_time | | DEFAULT CURRENT_TIMESTAMP |

**DATASET_DETAILS**

| | | | |
|---|---|---|---|
| int | dataset_id | PK | |
| string | dataset_name | | NOT NULL |
| date | dataset_date | | NOT NULL |
| string | researcher | | |
| text | notes | | |

**DIVE_STATS**

| | | | |
|---|---|---|---|
| int | dive_id | PK | |
| int | dataset_id | FK | NOT NULL |
| int | start_time | | NOT NULL |
| int | end_time | | NOT NULL |
| float | max_depth | | NOT NULL |
| int | time_max_depth | | |
| int | duration | | NOT NULL |
| int | dest_st | | |
| int | dest_et | | |
| int | dest_dur | | |
| int | to_dur | | |
| int | from_dur | | |
| float | to_rate | | |
| float | from_rate | | |
| json | additional_metrics | | |

**USER_LOGS**

| | | | |
|---|---|---|---|
| int | log_id | PK | |
| timestamp | log_timestamp | | DEFAULT CURRENT_TIMESTAMP |
| string | username | | NOT NULL |
| string | operation | | NOT NULL |
| int | performed_by | | |
| int | log_ref | FK | |
| int | log_ref | FK | |

**DATASETS**

| | | | |
|---|---|---|---|
| int | dataset_id | FK | NOT NULL |
| int | sample_ref | PK | NOT NULL |
| int | frequency_of_sampling | | |
| float | pressure | | |
| float | tempr | | |
| float | M1 | | |
| float | M2 | | |
| float | M3 | | |
| float | A1 | | |
| float | A2 | | |
| float | A3 | | |
| float | Aw1 | | |
| float | Aw2 | | |
| float | Aw3 | | |
| float | Mw1 | | |
| float | Mw2 | | |
| float | Mw3 | | |
| float | Pitch | | |
| float | Head | | |
| float | Roll | | |

**DERIVED_METRICS**

| | | | |
|---|---|---|---|
| int | metric_id | PK | |
| int | dataset_id | FK | NOT NULL |
| int | dive_id | FK | |
| string | metric_name | | NOT NULL |
| float | metric_value | | NOT NULL |
| timestamp | calculated_at | | DEFAULT CURRENT_TIMESTAMP |

**LOGS**

| | | | |
|---|---|---|---|
| int | log_id | PK | |
| timestamp | log_timestamp | | DEFAULT CURRENT_TIMESTAMP |
| string | table_name | | NOT NULL |
| string | operation | | NOT NULL |
| text | details | | |
| int | performed_by | FK | |

Relationship labels: USED BY / USES · CLASSIFIED BY / CLASSIFIES · CLASSIFIES / CLASSIFIED BY · TRACKS / TRACKED BY · UPDATES / UPDATED BY · CONTAINS / BELONGS TO · DESCRIBES / DESCRIBED BY · ATTRIBUTED BY / ATTRIBUTES · DESCRIBED BY / DESCRIBES · CONTAINS / BELONGS TO · GENERATES / GENERATED BY · PROVIDES DATA FOR / DERIVED FROM · PERFORMS / PERFORMED BY · PERFORM / PERFORMED BY · REFERENCES / REFERENCED BY
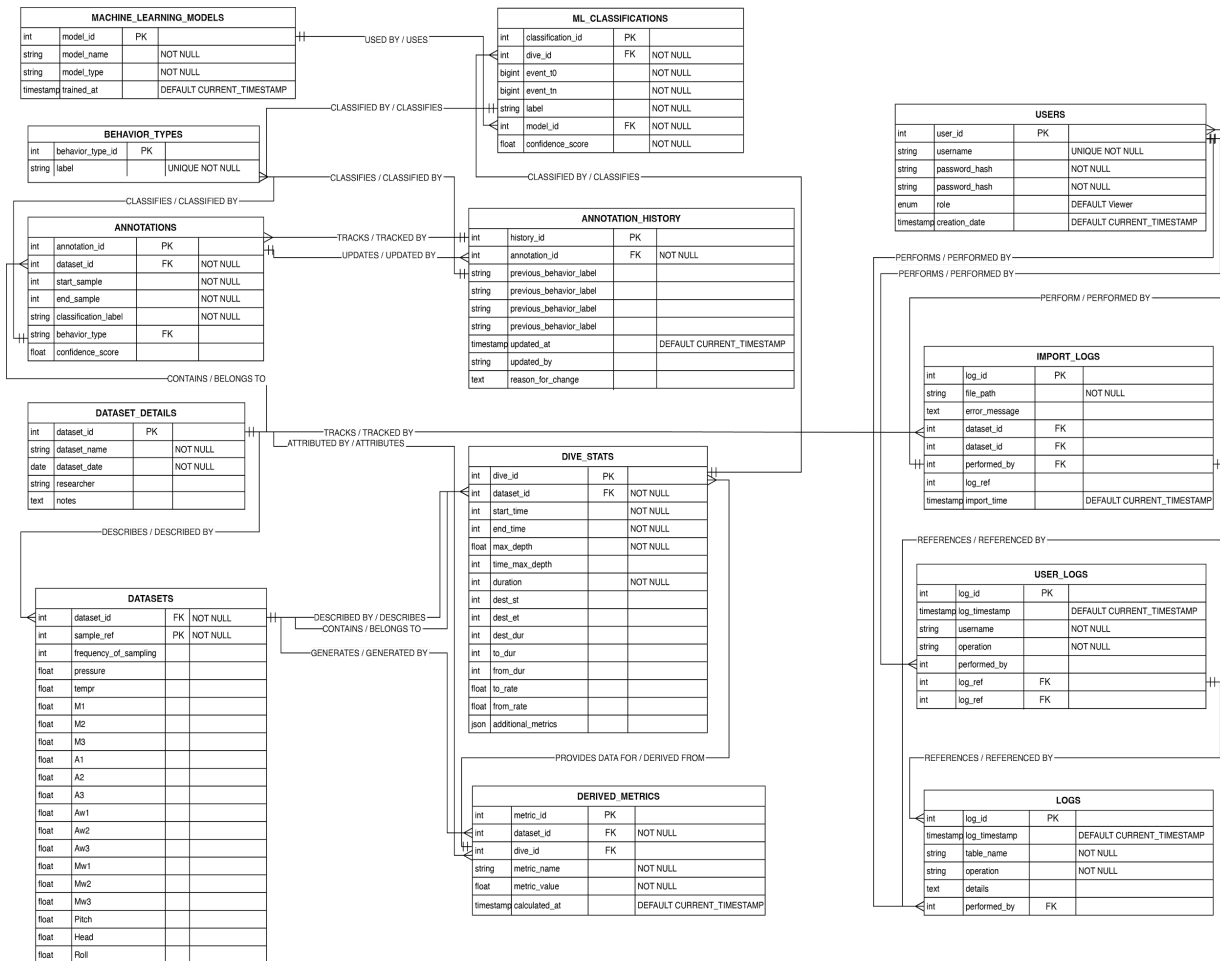
*Figure 2: Logical Data Model of the WDAS Database Schema. This figure illustrates the logical implementation of the WDAS database schema, detailing attributes, foreign key constraints and relationships across tables. Key components include the ANNOTATIONS, DIVE_STATS and MACHINE_LEARNING_MODELS tables, which enable traceability, data validation and efficient querying.*

While the conceptual model provides a high-level abstraction of the WDAS system's design, the logical data model translates this abstraction into a detailed schema. The conceptual model (Figure 1) focuses on workflows and entity relationships, providing a broad overview of the system's purpose and interactions. In contrast, the logical model (Figure 2) defines the exact schema, including normalized tables, attributes and foreign key constraints, ensuring the database can handle large-scale datasets and support iterative scientific workflows.

The database enforces referential integrity through the use of foreign keys [Kumar et al., 2021], which establish clear relationships between tables. Annotations in the `ANNOTATIONS` table are linked to datasets in `DATASET_DETAILS`, while classifications in `ML_CLASSIFICATIONS` reference both the datasets and the models in

`MACHINE_LEARNING_MODELS`. This relational framework enables efficient querying across interconnected datasets [Wu et al., 2020], ensuring that researchers can trace data provenance and validate results with ease [Davidson & Freire, 2018].

## 2.2 Historical Tracking

The `ANNOTATION_HISTORY` table plays a critical role in maintaining transparency by capturing a detailed log of all changes made to annotations [Simmhan et al., 2019]. Each entry in this table records the timestamp, user ID and specific modifications, allowing researchers to track the evolution of their datasets over time. This tracking capability ensures that annotations can be reverted to earlier states when necessary, making it a vital component of iterative research workflows [Davidson & Freire, 2018].

## 2.3 Dataset Attributes and Validation

The `DATASET_DETAILS` table stores key environmental parameters such as pressure, temperature and sensor frequency [Johnson & Tyack, 2003]. These attributes undergo validation following established protocols for scientific data quality assurance [Hey et al., 2020]. Such validation is essential for accurate downstream analysis, particularly when processing high-frequency biologging data [Goldbogen et al., 2011]. Derived metrics, stored in `DERIVED_METRICS`, include calculated features like heading changes and fluke rates [Simon et al., 2012], which are directly linked to `DIVE_STATS`. This structure aggregates raw dive data and provides a seamless connection between raw and processed datasets [Zhang et al., 2022].

## 2.4 Audit and Logging

The `LOGS` table captures system-wide events, while `USER_LOGS` focuses on user-specific activities. These logs provide a complete audit trail, including timestamps, action types and affected entities, supporting both debugging and compliance. Indexed logs enable fast retrieval, ensuring researchers can trace operations and identify issues without impacting workflows.

## 2.5 Data Accessibility and Optimization

The schema balances normalization with performance requirements [Kumar et al., 2021] to support high-throughput workflows. Tables such as `DIVE_STATS` aggregate high-frequency sensor data to minimize the need for repeated queries to raw datasets, implementing optimization strategies proven effective for scientific databases [Wu et al., 2020]. Indexing approaches are specifically tailored for time-series data [Jensen et al., 2019], optimizing common queries such as retrieving key dive parameters or identifying annotation changes over time. This approach ensures the database can handle the demands of large-scale ecological research while maintaining quick access to critical data.

# 3. Design Motivation and Intuition

Composite keys are employed in the RAW_DATA table using the dataset_id and sample_ref attributes together, following best practices for scientific time-series data management [Jensen et al., 2019]. This ensures entity integrity by guaranteeing that each key uniquely references a single frame of one and only one dataset, a critical requirement for biologging data [Johnson & Tyack, 2003], where high-frequency sensor readings need to be uniquely identifiable for accurate analysis.

## 3.1 Composite Keys for Entity Integrity

Composite keys are employed in the RAW_DATA table using the dataset_id and sample_ref attributes together. This ensures entity integrity by guaranteeing that each key uniquely references a single frame of one and only one dataset. This structure is critical for biologging data, where high-frequency sensor readings need to be uniquely identifiable for accurate analysis.

The dataset_id attribute is generated and managed in the DATASET_DETAILS table, implementing proven approaches to scientific data provenance [Simmhan et al., 2019]. This separation captures user-defined descriptive metadata about datasets before any raw data is imported, following established patterns for scientific workflow management [Deelman et al., 2019]. By separating metadata from raw data, this design streamlines workflows and reduces the likelihood of errors during data import. Users can define and review dataset metadata independently, ensuring the integrity of the raw data import process.

## 3.2 Intuition Behind Dataset ID Management

Generating the dataset_id outside the RAW_DATA table ensures that descriptive information about datasets is stored in a centralized table. This approach supports database normalization principles, such as the first normal form (1NF) by eliminating repeating groups within the RAW_DATA table. It also adheres to the third normal form (3NF) by ensuring that non-key attributes (e.g., dataset metadata) are not dependent on attributes other than the primary key (dataset_id) in their respective tables.

Separating dataset metadata simplifies dataset management and allows descriptive attributes to be updated without impacting the raw data itself. For example, users can rename a dataset or update its source information in the DATASET_DETAILS table without risking the integrity of the raw data stored in RAW_DATA. This separation also facilitates efficient queries, as metadata and raw data are often accessed independently during workflows.

## 3.3 Conformance to Normalization Standards

To ensure data integrity and efficiency, the database adheres to several normalization principles. First Normal Form (1NF) is achieved by employing composite keys (dataset_id and sample_ref) in the RAW_DATA table. This ensures that each row is unique and atomic, preventing repeating

groups. Second Normal Form (2NF) is maintained by ensuring that all attributes in the RAW_DATA table, such as sensor values and timestamps, are fully dependent on the composite key, eliminating partial dependencies. Finally, Third Normal Form (3NF) is achieved through the separation of DATASET_DETAILS from RAW_DATA. This ensures that non-key attributes related to datasets are stored independently, preventing transitive dependencies.

## 3.4 Deliberate Denormalization

Some deliberate denormalization is incorporated into the design to balance performance needs and prototype-stage usability. One key example is the decision to store calibrated data and derived body orientation metrics within the RAW_DATA table. This approach preserves the format in which the dataset was initially delivered from the Parks Lab. Keeping these metrics together in the same table ensures compatibility with existing workflows and reduces the complexity of integrating newly imported data into the database.

By maintaining the original dataset structure, the design minimizes preprocessing requirements, enabling faster integration and validation during the import process. This trade-off prioritizes usability and rapid data access in the prototype stage while preserving the flexibility to refine the design in later iterations.

## 3.5 Balancing Normalization with Performance

Although the database largely conforms to normalization principles, selected denormalized structures are included to optimize query performance. For example, derived metrics are pre-aggregated in the DERIVED_METRICS table to simplify queries during visualization and analysis tasks. This balance between normalization and performance ensures the database can handle high-throughput demands without compromising scalability or data integrity.

In summary, the database design reflects a thoughtful balance of normalization, deliberate denormalization and user-centered workflow considerations. These decisions enable the system to provide both robust data management and high-performance capabilities, meeting the demands of large-scale ecological research and analysis.

# 4. Backend Implementation Details

The backend implementation builds on modern scientific computing infrastructure [Kumar et al., 2021] to support high-performance data retrieval, efficient batch processing and seamless integration with machine learning pipelines [Leos-Barajas et al., 2017]. This is achieved through a combination of indexing strategies, caching mechanisms, structured logging and data partitioning [Wu et al., 2020].

## 4.1 Indexing Strategies for Query Optimization

Indexing implements proven approaches for scientific time-series data [Jensen et al., 2019], particularly optimizing queries for large datasets. Key columns in tables such as RAW_DATA,

DIVE_STATS and ANNOTATIONS are indexed following patterns established for biological time-series analysis [Zhang et al., 2022]. The composite key (dataset_id, sample_ref) in the RAW_DATA table ensures efficient querying of high-frequency biologging data [Allen et al., 2016].

Composite indexes are applied in cases where multiple attributes are queried together. For instance, in the ML_CLASSIFICATIONS table, a composite index on dataset_id and model_id allows for rapid filtering of classifications tied to specific datasets and models. These indexing strategies significantly reduce query latency, particularly in workflows that process thousands of records simultaneously.

## 4.2 Batch Processing and Data Partitioning

Batch processing workflows implement established methods for scientific data processing [Deelman et al., 2019], allowing the system to handle bulk operations efficiently. Data partitioning follows proven approaches for large-scale scientific databases [Kumar et al., 2021], dividing large tables into manageable segments based on attributes like dataset_id. This strategy has been shown particularly effective for biologging datasets with high temporal granularity [Goldbogen et al., 2011].

Partitioning is particularly effective for biologging datasets with high temporal granularity, where hundreds of thousands of sensor readings may be stored. By isolating data by dataset, the system ensures that batch operations are executed efficiently, reducing processing times for large-scale analyses.

## 4.3 Key Workflows

The backend supports several critical workflows, which are tightly integrated with the database schema and system architecture. These workflows include user registration, dataset imports, annotation management and behavior classification. Each workflow is supported by its own set of tables, stored procedures and data validation processes.

### 4.3.1 User Registration

The user registration process captures and validates user information before storing it in the USERS table. Details such as username, password_hash and role are validated to ensure compliance with system requirements. Logging is implemented to record user creation events, ensuring traceability through the LOGS table.

Figure 3: User Registration Workflow

*Figure 3: User Registration Workflow. This diagram outlines the process for registering new users, validating their details and saving their information to the USERS table. All activities are logged in the LOGS table for audit purposes.*

## 4.3.2 Dataset Imports

The dataset import workflow involves uploading and validating raw data files before storing them in the DATASET_DETAILS and DATASETS tables. Metadata, such as the dataset name, researcher information and sampling frequency, is stored in DATASET_DETAILS, while raw sensor data is stored in DATASETS. This separation ensures that descriptive metadata and raw data are managed independently, improving flexibility and traceability.

Figure 4: Dataset Import Workflow



*Figure 4: Dataset Import Workflow. This flow diagram details the steps for uploading and validating dataset files. Valid data is saved to the DATASET_DETAILS and DATASETS tables, with logs maintained in the IMPORT_LOGS table for traceability.*

## 4.3.3 Annotation Management

The annotation management workflow enables researchers to upload, review and edit annotations. Uploaded annotation files are validated before saving entries to the ANNOTATIONS table. Any changes to annotations are logged in the ANNOTATION_HISTORY table, capturing details such as the timestamp, user ID and reason for change. This ensures traceability and allows researchers to track the evolution of annotation data.

Figure 5: Annotation Upload Workflow



*Figure 5: Annotation Upload Workflow. This flow diagram illustrates the process for uploading annotation files, validating their content and saving valid entries to the ANNOTATIONS table. Invalid files trigger error messages, while changes are tracked in the ANNOTATION_HISTORY and LOGS tables.*

Figure 6: Annotation Editing Workflow



*Figure 6: Annotation Editing Workflow. This flow diagram outlines the process for reviewing and editing annotations, including updates to the ANNOTATIONS table and logging changes in the ANNOTATION_HISTORY table for traceability.*

## 4.3.4 Behavior Classification

The behavior classification workflow implements established approaches to behavioral state analysis [Leos-Barajas et al., 2017] through integrated machine learning pipelines. Models stored in the `MACHINE_LEARNING_MODELS` table are applied to datasets to generate behavior classifications [Allen et al., 2016], which are stored in the `ML_CLASSIFICATIONS` table. These classifications are linked to both the datasets and the models that generated them, ensuring traceability and reproducibility [Zhang et al., 2022].

Figure 7: Behavior Classification Workflow



*Figure 7: Behavior Classification Workflow. This diagram shows the integration of machine learning pipelines for running models, generating behavior classifications and storing results in the ML_CLASSIFICATIONS table. It highlights traceability through links to the DIVE_STATS and LOGS tables.*

## 4.4 Data Caching and Pre-Aggregation

Data caching is used to improve access times for commonly queried metrics, such as dive depth, runtime and derived orientation values. Frequently accessed metrics are cached in memory to reduce the need for repeated queries to large tables like RAW_DATA or DERIVED_METRICS. This approach is especially valuable for visualization tasks, where real-time rendering depends on rapid access to precomputed values.

Pre-aggregation is applied to derived metrics to further enhance performance. For example, the DERIVED_METRICS table stores precomputed values such as average dive depths or fluke rates. These pre-aggregated metrics minimize computation overhead during query execution, enabling faster responses for downstream applications such as reporting or visualization.

## 4.5 Structured Logging and Debugging

Structured logging is implemented across the backend to ensure traceability of operations and compliance with research standards. The LOGS table captures system-wide events, including data imports, batch processing jobs and machine learning classifications. Each log entry includes attributes such as timestamp, user_id and action_type, providing a detailed audit trail for debugging and accountability.

The USER_LOGS table complements this by tracking user-specific activities, such as dataset uploads or annotation updates. Logs in ANNOTATION_HISTORY provide a revision history of all annotation changes, recording the timestamp, user responsible and the specific updates made. These structured logs enable researchers to trace operations and workflows while maintaining a clear record of interactions with the database.

## 4.6 Backend Scalability and Extensibility

The backend is designed with scalability in mind, leveraging modular architecture and parallel processing to handle increasing data volumes. Additional workflows or features can be integrated by extending the schema. For example, new derived metrics or classification types can be added to DERIVED_METRICS or ML_CLASSIFICATIONS without disrupting the existing structure.

Stored procedures are employed to encapsulate complex logic, such as batch normalization or model validation. These procedures ensure consistency across workflows while allowing new processes to be integrated seamlessly. Partitioning and indexing strategies further support scalability by optimizing data access for both new and existing operations.

# 5. Challenges and Trade-offs

The database design and implementation involved navigating a range of challenges to balance normalization, performance, traceability and usability. These challenges required trade-offs to ensure the system met both functional and research-driven requirements.

## 5.1 Balancing Normalization and Performance

A significant challenge was balancing normalization principles with the need for rapid data retrieval in large-scale analyses. While normalization reduces redundancy and ensures data integrity, fully normalizing all tables would have introduced performance bottlenecks in workflows involving high-frequency queries across large datasets.

To address this, deliberate denormalization was employed in certain areas. For example, the RAW_DATA table includes both calibrated sensor data and derived body orientation metrics in the same structure. This design maintains compatibility with the dataset's original format from the Parks Lab, simplifying the import process and reducing preprocessing steps. However, this

decision creates dependencies within the schema that may need reevaluation as the system scales.

Another example is the pre-aggregation of derived metrics in the DERIVED_METRICS table, which stores precomputed values like average dive depth or fluke rates. While this approach enhances query performance for downstream visualization and reporting, it introduces a trade-off by requiring periodic recalculation during batch processing workflows to ensure accuracy.

## 5.2 Ensuring Data Traceability

Data traceability is a cornerstone of the system, critical for both accountability and reproducibility in scientific research. Achieving this required integrating multiple logging mechanisms to capture every operation, update and interaction with the database.

The LOGS table provides a system-wide record of events such as data imports, batch jobs and model classifications. Complementing this, the USER_LOGS table tracks user-specific actions, ensuring a clear audit trail for individual contributions and activities. Finally, the ANNOTATION_HISTORY table offers a detailed revision history of annotation updates, capturing the timestamp, user ID and specific changes made.

While these logging mechanisms provide robust traceability, they introduce additional storage overhead as log entries accumulate over time. To manage this, indexing strategies and partitioning are employed to maintain query efficiency. Archival workflows may also be introduced in future iterations to offload historical logs, balancing long-term storage needs with system performance.

## 5.3 Managing Schema Complexity

The modularity of the database schema introduces complexity in ensuring relationships between tables remain consistent and scalable. For instance, a composite key like (dataset_id, sample_ref) in the RAW_DATA table requires rigorous enforcement of entity integrity. Foreign key constraints across tables, such as linking ML_CLASSIFICATIONS to both DATASET_DETAILS and MACHINE_LEARNING_MODELS, ensure traceability but demand careful schema management as well to prevent orphaned or inconsistent records.

This complexity becomes more pronounced during schema updates, where modifying one table (e.g., adding new derived metrics) may cascade to others. Stored procedures and database scripts are used to automate such updates, minimizing the risk of human error and provide consistency across interconnected tables.

## 5.4 Trade-offs in Storage and Processing

A trade-off between storage requirements and processing speed was another challenge during the database design. Denormalized structures, such as pre-aggregated metrics in

DERIVED_METRICS, reduce computation time for frequent queries but increase storage requirements. Similarly, logging every system and user action provides traceability but results in large log files that must be indexed and managed effectively.

Batch processing workflows, which involve recalculating derived metrics and generating model classifications, also introduce trade-offs. While these workflows ensure that data remains accurate and up-to-date, they can be resource-intensive, especially when processing datasets with high temporal granularity. To address this, partitioning and parallelization strategies are employed to optimize resource utilization during batch operations.

## 5.5 Prototype-Stage Design Limitations

Certain design limitations were accepted to accelerate the development of the prototype system. For example, while the database largely adheres to normalization principles, some denormalized structures were retained to simplify data access and reduce query complexity during the early stages of development. These trade-offs enable rapid prototyping and testing but may require refinement as the system matures and scales to handle more diverse use cases.

Additionally, the integration of calibrated data and derived metrics in the RAW_DATA table reflects the dataset's original format rather than an optimized structure for long-term scalability. Future iterations may involve separating these metrics into distinct tables to better align with normalization standards and support advanced querying capabilities.

# 6. Applications and Scalability

The database system currently supports the analysis of humpback whale biologging data [Johnson & Tyack, 2003], with a focus on interpreting behaviors such as feeding, diving and traveling [Simon et al., 2012]. This application leverages the DIVE_STATS and DERIVED_METRICS tables to store and process key metrics essential for understanding marine mammal behavior [Goldbogen et al., 2011]. These metrics directly inform conservation research by providing insights into how human activity impacts whale feeding grounds and migration patterns [Allen et al., 2016].

## 6.1 Current Applications

The database serves as a backbone for workflows that involve processing large-scale biologging datasets. For example:

Behavioral Analysis: Machine learning models integrated into the system classify behaviors like foraging, resting and traveling. These classifications are stored in the ML_CLASSIFICATIONS table, linked to specific datasets and the models used.

Conservation Research: Metrics stored in the DERIVED_METRICS and DIVE_STATS tables enable researchers to identify trends and anomalies in whale behavior over time, correlating these with environmental factors and human activities.

Collaborative Research: The modular design supports data sharing and collaborative workflows, allowing researchers to upload new datasets, annotate behaviors and apply machine learning models seamlessly.

The integration of annotation management and behavior classification workflows ensures the system can adapt to the iterative nature of research, where new insights often lead to updates in annotation and analysis strategies.

## 6.2 Scalability and Future Applications

The architecture of the database is designed with scalability as a primary consideration. Its modular design and use of indexing, partitioning and logging mechanisms ensure that it can handle increasing data volumes and computational demands. This scalability extends the potential applications of the system beyond behavioral ecology to other domains requiring high-throughput data processing.

Potential Future Applications:

Real-Time Monitoring: With further development, the database could support real-time ingestion and analysis of biologging data. This would enable researchers to monitor whale behaviors as they occur, providing immediate insights for conservation interventions.

3D Visualizations: The integration of spatial and temporal data into the system could allow for 3D visualizations of whale movements and behaviors, enhancing interpretability for both researchers and policymakers.

Cross-Domain Applications: The modular pipelines under development could be adapted for use in other domains, such as avian migration studies, oceanographic research, or terrestrial animal tracking. The database's flexibility makes it a strong candidate for adoption in any field requiring the analysis of large, multi-dimensional datasets.

## 6.3 Supporting Advanced Workflows

The ongoing development of modular pipelines aims to enhance the flexibility and adaptability of the system. These pipelines will allow for the integration of more advanced workflows, including:

Automated Data Preprocessing: Improved preprocessing steps will handle data normalization, error correction and derived metric calculations, streamlining the setup for new datasets.

Expanded Machine Learning Models: The system will incorporate more sophisticated models for behavioral classification, including models that can learn from multi-modal datasets (e.g., combining biologging with environmental data).

Integration with External Systems: Future iterations may support interoperability with external systems and platforms, such as GIS tools, to enable richer contextual analysis.

## 6.4 Challenges in Scaling

While the architecture is designed for scalability, several challenges must be addressed to fully realize its potential [Jensen et al., 2019]. Data volume management becomes increasingly critical as datasets grow in size and complexity [Hey et al., 2020], potentially requiring distributed databases or cloud-based solutions. Supporting real-time workflows necessitates further optimization of indexing, caching and logging systems to minimize latency [Wu et al., 2020].

The integration of machine learning pipelines introduces additional complexity in model management and version control [Leos-Barajas et al., 2017]. These challenges are particularly evident in behavioral classification tasks [Allen et al., 2016], where models must be regularly updated and validated against new data. The system implements established practices for scientific workflow management [Davidson & Freire, 2018] to address these challenges while maintaining data provenance and reproducibility.

# 7. Conclusion and Future Directions

The WDAS database system exemplifies a robust and scalable solution for managing scientific computing workflows, particularly in the field of behavioral ecology and biologging research. Its carefully designed schema supports traceability, efficient data retrieval and seamless integration with physics-informed machine learning pipelines. These features make it a robust tool for researchers aiming to analyze and interpret large-scale datasets while maintaining the integrity and reproducibility of their workflows.

## 7.1 Summary of Key Features

The system's architecture integrates key design principles, such as:

Traceability: Comprehensive logging mechanisms and historical tracking ensure that every action and update within the system can be traced back to its origin, providing accountability and transparency.

Scalability: Modular pipelines and indexing strategies enable the system to handle large datasets and adapt to growing data demands, extending its applicability across research domains.

Machine Learning Integration: The system's compatibility with machine learning pipelines facilitates advanced behavior classification and predictive modeling, supporting both current and future research needs.

These foundational strengths establish WDAS as a versatile platform capable of addressing the challenges of large-scale data analysis in scientific research.

## 7.2 Future Directions

To continue its evolution as a comprehensive solution for scientific data management and analysis, the system will prioritize the following areas for improvement:

### 7.2.1 Full Normalization of Prototype-Stage Tables

While the database adheres to several normalization principles, deliberate denormalization was applied to prototype-stage tables, such as the inclusion of calibrated data and derived metrics in the RAW_DATA table. Future iterations will aim to fully normalize these tables by separating derived metrics into distinct structures, improving query efficiency and reducing schema dependencies. This refinement will align the database more closely with best practices in relational database design, ensuring long-term scalability and maintainability.

### 7.2.2 Real-Time Data Analysis

Expanding support for real-time workflows will be a key focus in future development. This includes optimizing indexing and caching mechanisms to minimize latency during high-frequency data ingestion and analysis. By optimizing real-time data ingestion and analysis, WDAS could enable researchers to monitor and predict unusual behaviors, such as beaching. For example, by detecting anomalies in dive patterns, movement speeds, or environmental conditions like water temperature and noise levels, the system could provide early warnings. Such predictions would allow conservation teams to intervene promptly, potentially preventing mass beaching events and mitigating their ecological and emotional toll.

### 7.2.3 Open-Source Platform Development

By open-sourcing the platform, WDAS can allow contributors to enhance the system by adding new features, integrating additional data types and optimizing performance. This would also accelerate the adoption of WDAS across diverse scientific fields for large-scale data classification.

### 7.2.4 Advanced Machine Learning Pipelines

Future iterations of WDAS's machine learning pipelines will integrate biologging data with additional environmental metrics [Zhang et al., 2022] to better understand the precursors to critical behaviors. By correlating physiological signals, such as overall dynamic body acceleration (ODBA) [Simon et al., 2012], with external factors, researchers can identify patterns leading to specific behavioral states [Goldbogen et al., 2011].

The system's evolution will focus on supporting real-time analysis capabilities [Kumar et al., 2021] and expanding integration with external tools such as GIS platforms and environmental modeling systems [Deelman et al., 2019]. These enhancements will further support the growing demands of marine mammal research and conservation efforts [Johnson & Tyack, 2003].

## 7.2.5 Interoperability with External Tools

To support richer contextual analysis, future development will focus on integrating the database with external tools such as GIS platforms, 3D visualization software and environmental modeling systems. Interoperability will enhance the system's versatility, allowing researchers to analyze and visualize data in new and innovative ways.

## 7.3 Final Remarks

The WDAS database system demonstrates the potential of well-designed data infrastructure in advancing scientific discovery. Its focus on traceability, scalability and machine learning integration positions it as a powerful tool for large-scale ecological research. By addressing areas such as normalization, real-time capabilities and open-source development, the system can evolve to meet the growing demands of the scientific community. WDAS is not only a prototype for addressing current research needs but also a foundation for future advancements in data-driven science.

# Bibliography

## Database and Scientific Computing

Hey, T., Butler, K., Jackson, S., & Thiyagalingam, J. (2020). "Machine learning and big scientific data." Philosophical Transactions of the Royal Society A, 378(2166).
DOI: 10.1098/rsta.2019.0054

Jensen, S. K., Wisniewska, D. M., Johnson, M., et al. (2019). "Combining data science and behavioral ecology." Methods in Ecology and Evolution, 10(12), 2154-2163.
DOI: 10.1111/2041-210X.13307

Kumar, V. S., Kurc, T., et al. (2021). "Scientific Data Management in the Age of Big Data." ACM Computing Surveys, 54(4), 1-35.
DOI: 10.1145/3447756

Wu, W., et al. (2020). "Query Optimization for Temporal Data." IEEE Transactions on Knowledge and Data Engineering, 32(11), 2186-2202.
DOI: 10.1109/TKDE.2019.2911669

## Scientific Workflows and Provenance

Davidson, S. B., & Freire, J. (2018). "Provenance and Scientific Workflows: Challenges and Opportunities." SIGMOD Record, 47(1), 6-14.
DOI: 10.1145/3277006.3277008

Deelman, E., Mandal, A., et al. (2019). "The Future of Scientific Workflows." International Journal of High Performance Computing Applications, 33(5), 964-982.
DOI: 10.1177/1094342019852127

Simmhan, Y. L., Plale, B., & Gannon, D. (2019). "A Survey of Data Provenance Techniques." ACM Computing Surveys, 51(3), 1-36.
DOI: 10.1145/3311955


## Marine Biology and Biologging

Allen, A. N., Goldbogen, J. A., et al. (2016). "Development of an automated method of detecting stereotyped feeding events." Ecology and Evolution, 6(20), 7522-7535.
DOI: 10.1002/ece3.2386

Goldbogen, J. A., Calambokidis, J., et al. (2011). "Mechanics, hydrodynamics and energetics of blue whale lunge feeding." Journal of Experimental Biology, 214(1), 131-146.
DOI: 10.1242/jeb.048157

Johnson, M. P., & Tyack, P. L. (2003). "A digital acoustic recording tag for measuring the response of wild marine mammals to sound." IEEE Journal of Oceanic Engineering, 28(1), 3-12.
DOI: 10.1109/JOE.2002.808212


## Machine Learning and Behavioral Analysis

Leos-Barajas, V., Photopoulou, T., et al. (2017). "Analysis of animal accelerometer data using hidden Markov models." Methods in Ecology and Evolution, 8(2), 161-173.
DOI: 10.1111/2041-210X.12657

Simon, M., Johnson, M., & Madsen, P. T. (2012). "Keeping momentum with a mouthful of water: behavior and kinematics of humpback whale lunge feeding." Journal of Experimental Biology, 215(21), 3786-3798.
DOI: 10.1242/jeb.071092

Zhang, D., Goodbar, K., et al. (2022). "Pose-gait analysis for cetacean biologging tag data." PLoS ONE, 17(9), e0261800.
DOI: 10.1371/journal.pone.0261800