

Database documentation: vessel

D. O. Fisher, C. Sutton, and F. Wei

NIWA Fisheries Data Management
Database Documentation Series

Version 1.0 May 2005

Contents

1	Introduction to the Database Document series	3
2	Vessel Registrations.....	3
2.1	Data sources.....	4
2.2	Data validation.....	4
3	Data Structures.....	5
3.1	Table relationships.....	5
3.2	Database design.....	9
4	Table Summaries	10
5	vessel Tables.....	11
5.1	Table 1: t_meta.....	11
5.2	Table 2: t_vessel.....	12
5.3	Table 3: t_vessel_org.....	13
5.4	Table 4: t_history.....	14
5.5	Table 5: t_lloyds.....	15
5.6	Table 6: t_rule.....	16
6	vessel business rules.....	17
6.1	Introduction to business rules.....	17
6.2	Summary of rules.....	18
Appendix 1	Reference code tables.....	19
Appendix 2	Data Grooming Rules.....	21
Appendix 3	Final Manual Grooming Procedures.....	23

List of Figures

Figure 1: Entity Relationship Diagram (ERD) for the vessel database.	7
---	---

Revision History

Version	Change	Date	Responsible
0.1	Initial version	10 Jan 2005	David Fisher
0.2	Incorporating comments from G Straker, MFish 21Mar2005	March 2005	David Fisher
1.0	First official release	3 May 2005	David Fisher

1 Introduction to the Database Document series

The National Institute of Water and Atmospheric Research (NIWA) is Data Manager and Custodian for the fisheries research data owned by the Ministry of Fisheries (MFish).

The Ministry of Fisheries data sets incorporate historic research data, data collected or held by MAF Fisheries prior to the split in 1995 of policy to the Ministry of Fisheries and research to NIWA, and more recent data collected by NIWA and other research providers for the Ministry of Fisheries.

This document describes the vessels specifications database **vessel**, and is part of the database documentation series produced by NIWA.

All documents in this series include an introduction to the database design, a description of the main data structures accompanied by an Entity Relationship Diagram (ERD), and a listing of all the main tables. The ERD graphically shows how all the tables link together.

This document is intended as a guide for users and administrators of the **vessel** database.

Access to this database is restricted to specific nominated personnel as specified in the current Schedule 6 of the Data Management contract between the Ministry of Fisheries and NIWA. Any requests for data should in the first instance be directed to the Ministry of Fisheries.

2 Vessel Registrations

Since the 1970's it has been mandatory for vessels fishing in New Zealand waters to be registered. Registration includes the completion of a form detailing the specifications of the vessel, including the size of the vessel, year built, and engine(s) make and model. The Ministry of Fisheries and its predecessors have held these registration data.

Fishing vessels may change their name or call sign, or re-power with a new engine with a different kilowatt (or horsepower) rating. Vessels may come and go from NZ waters. In the early years of New Zealand's fisheries foreign licensed vessels came to NZ waters to fish within the EEZ. Keeping track of individual vessel hulls, particularly if they return to New Zealand under a different name can be a difficult process and errors have been made in the past in assigning various vessel identifiers to generate a unique value for each individual vessel hull.

The objective of this database is to create a research version of these vessel registration data, for use by fisheries scientists in their analysis of catch and effort data from various New Zealand fisheries. It is also likely that these data will be useful to fisheries managers and compliance staff from the Ministry of Fisheries. This database provides the best-groomed dataset of individual fishing vessels operating in New Zealand waters to date.

2.1 Data sources

MFish provided NIWA with 3 overlapping vessel datasets.

1. An initial excel spreadsheet that included Lloyds IMO numbers that had been added by a student employed by MFish.
2. Data from the MFish vessel database based on the ‘old form’ used for vessel registrations.
3. Data from the MFish vessel database based on the ‘new form’ used for vessel registrations.

No one dataset contained all attributes or the full time period (i.e., all datasets contained one or more all null columns and/or did not cover the complete time period). The initial datasets included approximately one record per vessel per year per dataset. The scope of the contract was to groom data for vessels over 20 metres overall length and registered between October 1983 and September 2003.

2.2 Data validation

These three datasets referred to above were ‘collapsed’. That is the records were sorted by vessel and then by *spec_from* and *spec_to* attributes. These *spec_from* and *spec_to* attributes contain the start date and end date applicable to each record. Each consecutive record was compared for each vessel and if all attributes were identical except *spec_from* and *spec_to* then the first *spec_from* and second *spec_to* dates were assigned to the ‘collapsed’ or combined record. If any attribute changed then a record with these specification dates was written out and the process continued. This process was repeated for all records for each vessel.

MFish had recorded the *spec_to* as 31 Dec 2999 (or 2099) when the vessel was not known to have left NZ waters for that registration record. If a record had a *spec_to* of 31 Dec 2999 in the middle of a consecutive sequence for one dataset (i.e. one *owner_key*) then the *spec_to* was set to equal the next *spec_from* value.

In creating this research version of these vessel registration data, we did not retain information such as when a vessel was not registered to fish in NZ waters. This decision is based on the assumption that fisheries research scientists want to know the vessel details when it was fishing NZ waters and are not concerned if it was not fishing in NZ waters.

We then loaded these three datasets to the original vessel table, *t_vessel_org* in the database. These data were extensively groomed based on the rules in the *t_rule* table and in Appendix 2. Initially, this grooming process was to identify individual vessel hulls. We used the four main identifying attributes: MFish *vessel_key*, *vessel_id*, *call_sign* and *vessel_name*, plus combinations of other attributes to determine those records that represented the same vessel. A new attribute *vid* (NIWA vessel ID) was created to identify individual vessels. A series of electronic routines and manual checking of marginal cases was undertaken to assign distinct *vid* values to individual vessels.

Where two or more MFish *vessel_key* values had been assigned to the same individual vessel, identified by a distinct *vid*, the *vessel_key* that was the smallest number was adopted. This should equate to the *vessel_key* that was assigned first.

The data grooming process followed the rules in Appendix 2. For each vessel identified by a vid, records were compared electronically and where there were different values for the same attribute these rules were applied. Null values were populated based on not null values in other records for the same attribute for the same vid, as per the rules in Appendix 2. The data was ‘collapsed’ again at the end of this electronic grooming process.

Final manual grooming was then undertaken, which included comparing remaining records for each vid. For those records in the *t_vessel* table where an imo_no value existed these records were compared with the data on the Lloyds Register of Shipping CD-ROM. The Lloyds database contains the full history of each vessel including previous names and changes of engine.

The details of this final manual grooming are included in Appendix 2 – Rules.

Approximately 10%, i.e., 150 vessels were checked visually against the *t_vessel_org* and *t_history* tables. This was to ensure that the grooming process had worked appropriately. No obvious discrepancies were apparent.

Any remaining null values were then populated electronically from the Lloyds data in the table *t_lloyds* where appropriate.

A final ‘collapse’ was undertaken so in some cases the final groomed dataset in the *t_vessel* table contains one record for each vessel hull. Where genuinely different values exist for one vessel hull at different times, such as the vessel was re-engined with a more powerful engine, then 2 or more records will exist for that vessel hull.

3 Data Structures

3.1 Table relationships

The vessel database contains several tables. The ERD for **vessel** (Figure 1) shows the logical structure (i.e. schema) of the database and its entities (each entity is implemented as a database *table*) and the relationships between these tables and tables in other databases. This schema is valid regardless of the database system chosen, and it can remain correct even if the Database Management System (DBMS) is changed. Each table represents an object, event, or concept in the real world that is selected to be represented in the database. Each attribute of a table is a defining property or quality of the table. All of the table’s attributes are shown in the ERD. The underlined attribute represents the table’s primary key¹.

Some of the tables in the **vessel** database have attributes called foreign keys². The foreign keys define the relationships between the tables in **vessel**.

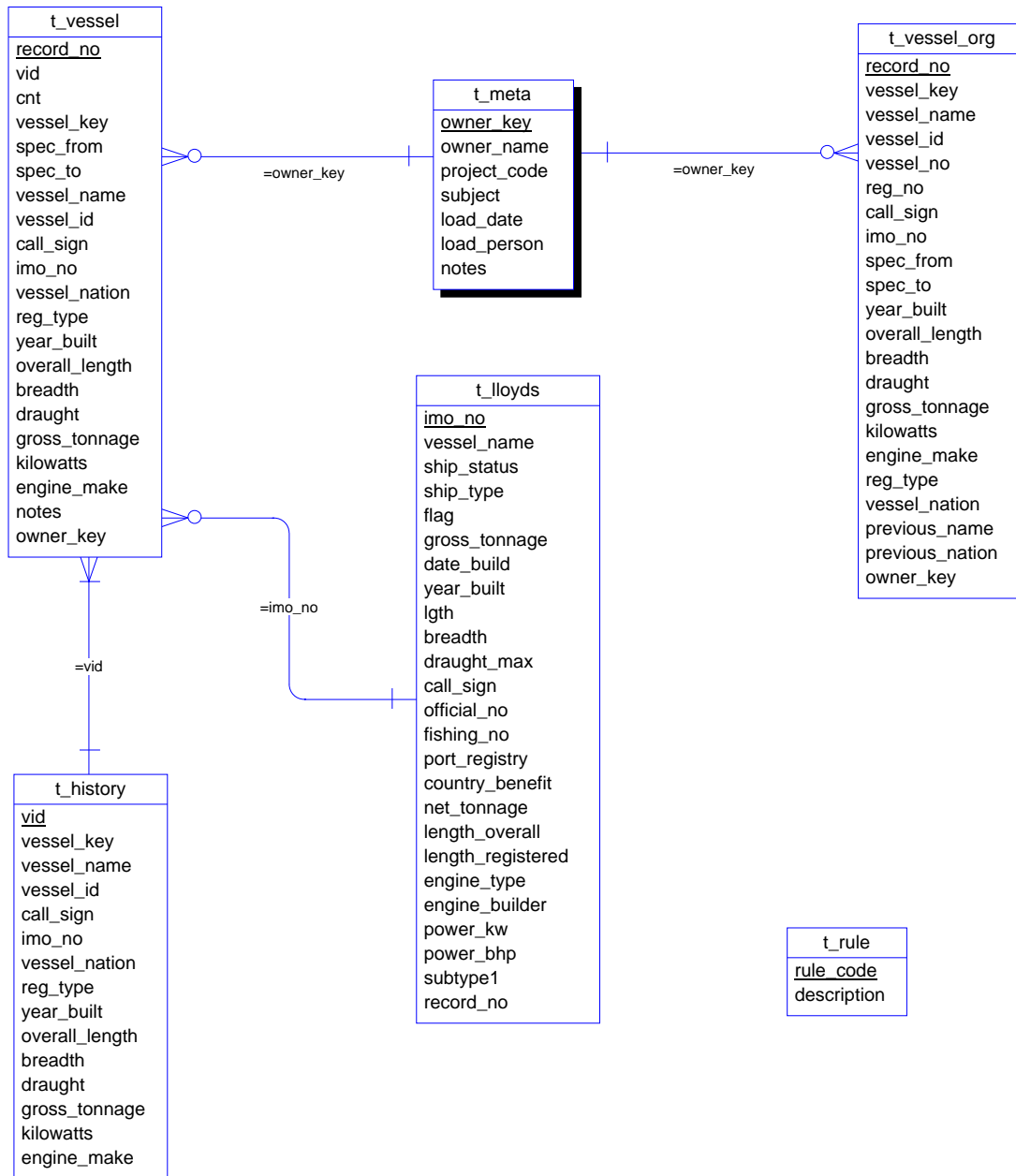
¹ A primary key is an attribute or a combination of attributes that contains an unique value to identify that record.

² A foreign key is any attribute, or a combination of attributes, in a table that is a primary key of another table. Tables are linked together through foreign keys.

The **vessel** database is implemented as a relational database; i.e., each table is a special case of the mathematical construct known as a *relation* and hence elementary relation theory is used to deal with the data within tables and the relationships between them. There are three types of relationships possible between tables, but only one exists in **vessel**: one-to-many³. These relationships can be seen in ERDs by connecting a crow's foot (indicating 'many') from the child table; e.g., *t_vessel_org*, to the parent table; e.g., *t_meta*, with a straight line (indicating 'one') pointing to the parent.

Every relationship has either a mandatory or optional aspect to it. If a relationship is mandatory, then it has to occur at least once, while an optional relationship might not occur at all. For example, in Figure 1, consider the relationship between the table *t_meta* and its child table *t_vessel_org*. The symbol 'o' by the child table *t_vessel_org* means that a metadata record (in table *t_meta*) can have zero or many original vessel records (in table *t_vessel_org*), while the bar by the parent table *t_meta* means that for every original vessel record there must be a matching record in *t_meta*.

³ A one-to-many relationship is where one record (the *parent*) in a table relates to one or many records (the *child*) in another table; e.g., one landing in *t_landing* can have many catches in *t_catch* but one catch can only come from one landing.



Conceptual Data Model		
Project :		
Model : vessel database		
Author : Fred Wei	Version 1	7/01/2005

Figure 1: Entity Relationship Diagram (ERD) for the vessel database.

These links are enforced by referential constraints (or integrity checks). Constraints do not allow *orphans* to exist in any table; i.e., where a child record exists without a related parent record. This may happen when either a parent record is deleted, the parent record is altered so the relationship is lost, or a child record is entered without a parent record

Constraints are shown in the table listings by the following format:

```
Referential:      constraint name (attribute[, attribute]) | INSERT |
                                                           | DELETE |
                  parent table (attribute[, attribute])
```

Note that the typographical convention for the above format is that square brackets [] may contain more than one item or none at all. Items stacked between vertical lines || are options of which one must be chosen.

For example, consider the following constraint found in the table *t_meta* :

```
Referential: (owner_key) REFERRED t_vessel (owner_key)
```

This means that the value of the attribute *owner_key* in the current record must already exist in the parent table *t_meta* or the record will be rejected and the following message will be displayed:

```
*** Error: constraint violation, primary key entry does not exist in table
't_meta'
```

For tables residing in external databases, the parent table name will be prefixed by the name of the database.

Section 5 lists all the **vessel** tables as implemented by the Empress DBMS. As can be seen in the listing of the tables, each table has a primary key that forms a unique index on the attribute. Primary keys are generally listed using the following format:

```
Indices:          PRIMARY KEY BTREE ON (attribute[, attribute])
```

where attribute(s) make up the primary key and the index name (if present) is the primary key name. Primary keys prevent records with duplicate keys from being inserted into the database tables; e.g., a record with vessel_id number that already exists in the table.

The database listing (Tables 1-6) show that the tables also have indices on many attributes. That is, attributes that are most likely to be used as a searching key have like values linked together to speed up searches. These indices are listed using the following format:

```
Indices:          NORMAL (2, 15) index_name ON (attribute[, attribute])
```

Note that indices may be simple (pointing to one attribute) or composite (pointing to more than one attribute). The numbers "...(2, 15)..." in the syntax are Empress DBMS default values relating to the amount of space allocated for the index.

3.2 Database design

The top-level table in **vessel** is *meta* (Table 1). This table holds some summary information (metadata) for this database. The primary key of this table is *owner_key*, which identifies the source of the data to one of the three source datasets.

The *owner_key* provides a link to the *t_vessel_org* table, which contains all the original data from the three MFish datasets.

These data from the *t_vessel_org* table were groomed using a series of electronic and manual processes. The resulting groomed data are contained in the *t_vessel* table and are considered to be the best version to use for research purposes.

The *t_rule* table documents the rules applied to the data between the *t_vessel_org* and *t_vessel* tables. The rules that were applied to groom data as part of this process are documented as a series of rule codes separated by the '|' character in the *notes* attribute in the *t_vessel* table.

The *t_history* table lists all variations recorded for each attribute for each vessel. With only one record for each unique *vid*, multiple values for each attribute are separated by the '/' character. This table allows users to select, for example, all vessels that have had multiple vessel keys assigned by selecting "where *vessel_key* like '%/%'".

The *t_lloyds* table contains the current data extracted from the Lloyds Register of Shipping database in 2004. These Lloyds data have been used in the grooming process as a reference dataset and can be joined with the *t_vessel* table using the *imo_no* attribute. The Lloyds dataset only holds information on vessels with a minimum gross tonnage of 100 tonnes.

4 Table Summaries

The **vessel** database has five tables containing vessel specifications data, plus one table which documents the rules used in grooming the vessel specification data.

The following is a listing and brief outline of the tables contained in **vessel**:

1. **t_meta**: contains data ownership information for the database
2. **t_vessel**: contains the groomed details for each vessel.
3. **vessel_org**: contains the 3 original datasets supplied by MFish.
4. **t_history**: includes all values for each attribute for each vessel hull in 1 record.
5. **t_lloyds**: contains vessel details from Lloyds register of shipping.
6. **t_rule**: documents the rules used to groom the vessel data and the associated codes used in the notes attribute of table *t_vessel*.

5 vessel Tables

The following are listings of the tables in the **vessel** database, including attribute names, data types (any range restrictions), and comments.

5.1 Table 1: t_meta

Comment: This table contains data summary information for the database.

Attributes	Data Type	Null?	Comment
owner_key	integer		Number as primary key.
owner_name	character(32,1)		Name of the dataset owner.
project_code	character(32,1)		The project code associated with the dataset.
subject	character(64,1)		Descriptive text relating to the dataset.
load_date	date(5)		Date when the dataset is loaded into the database.
load_person	character(32,1)		Name of the person who loaded the dataset.
notes	character(512,1)		Any commentary text regarding the dataset.

Creator: dba

Referential: (owner_key) REFERRED t_vessel (owner_key)

(owner_key) REFERRED t_vessel_org (owner_key)

Indices: PRIMARY KEY BTREE ON (owner_key)

5.2 Table 2: t_vessel

Comment: This table contains the groomed vessel dataset.

Attributes	Data Type	Null?	Comment
record_no	longinteger		Unique identification number for each record in this table.
vid	longinteger		Unique identification number assigned by NIWA for each vessel.
cnt	integer		The number of records (count) collapsed into the current record in the grooming process.
vessel_key	longinteger		MFish assigned number to identify a vessel.
spec_from	date(0)		Starting date for a vessel's specifications.
spec_to	date(0)		Finishing date for a vessel's specifications.
vessel_name	character(32,1)		Vessel name.
vessel_id	character(8,1)		MFish assigned alphanumeric code to identify a vessel, typically registration number or call sign.
call_sign	character(10,1)		Signal letters or radio call sign assigned by the relevant Registration (Flag) Authority.
imo_no	longinteger		Unique International Maritime Organisation number assigned by Lloyd's Register to each ship.
vessel_nation	character(3,1)		3 char code for vessel's flag or nationality.
reg_type	character(1,1)		MFish assigned classification code - Charter, Domestic, Foreign, Unknown.
year_built	integer		The year the vessel was built.
overall_length	integer		Overall length in metres
breadth	integer		Breadth in metres
draught	integer		Draught in metres
gross_tonnage	integer		Gross tonnage.
kilowatts	integer		Engine power in kilowatts.
engine_make	character(32,1)		Engine make.
notes	character(256,1)		Rules used in grooming process and other comments separated with ' '
owner_key	integer		Foreign key link to t_meta record.

Creator: dba

Referential:
 (owner_key) REFER t_meta (owner_key)
 (vid) REFER t_history (vid)
 (imo_no) REFER t_lloyds (imo_no)

Indices:
 PRIMARY KEY BTREE ON (record_no)
 NORMAL (2, 15) BTREE ON (vessel_key)
 FOREIGN KEY BTREE ON (owner_key)
 FOREIGN KEY BTREE ON (vid)
 FOREIGN KEY BTREE ON (imo_no)

5.3 Table 3: t_vessel_org

Comment: This table contains the original vessel datasets from MFish

Attributes	Data Type	Null?	Comment
record_no	longinteger		Unique identification number.
vessel_key	longinteger		MFish assigned number to identify a vessel.
vessel_name	character(32,1)		Vessel name.
vessel_id	character(8,1)		MFish assigned alphanumeric code to identify a vessel, typically registration number or call sign.
vessel_no	character(8,1)		MFish assigned number to identify a vessel.
reg_no	character(16,1)		Vessel registration number.
call_sign	character(10,1)		Signal letters or radio call sign assigned by the relevant Registration (Flag) Authority.
imo_no	longinteger		Unique International Maritime Organisation number assigned by Lloyd's Register to each ship.
spec_from	date(1)		Starting date for a vessel's specification.
spec_to	date(1)		Finishing date for a vessel's specification.
year_built	integer		The year a vessel is built.
overall_length	decimal(6,1)		Overall length in metres.
breadth	decimal(4,1)		Breadth in metres.
draught	decimal(4,1)		Draught in metres.
gross_tonnage	decimal(7,1)		Gross tonnage.
kilowatts	decimal(5,1)		Engine power in kilowatts.
engine_make	character(32,1)		Engine make.
reg_type	character(1,1)		MFish assigned classification code - Charter, Domestic, Foreign, Unknown.
vessel_nation	character(3,1)		Vessel flag or nationality.
previous_name	character(32,1)		Previous name of vessel.
previous_nation	character(3,1)		Previous flag or nationality of vessel.
owner_key	integer		Foreign key link to t_meta record.
Creator:	dba		
Referential:	(owner_key) REFER t_meta (owner_key)		
Indices:	FOREIGN KEY BTREE ON (owner_key)		

5.4 Table 4: t_history

Comment: This table lists all values of each attribute separated by '|' for each vessel

Attributes	Data Type	Null?	Comment
vid	longinteger		Unique identification number assigned by NIWA for each vessel.
vessel_key	character(32,1)		MFish assigned number to identify a vessel.
vessel_name	character(128,1)		Vessel name.
vessel_id	character(128,1)		MFish assigned alphanumeric code to identify a vessel, typically registration number or call sign.
call_sign	character(128,1)		Signal letters or radio call sign assigned by the relevant Registration (Flag) Authority.
imo_no	character(32,1)		Unique International Maritime Organisation number assigned by Lloyd's Register to each ship.
vessel_nation	character(32,1)		3 char code for vessel's flag or nationality.
reg_type	character(8,1)		MFish assigned classification code - Charter, Domestic, Foreign, Unknown.
year_built	character(64,1)		The year a vessel is built.
overall_length	character(64,1)		Overall length in metres.
breadth	character(64,1)		Breadth in metres.
draught	character(64,1)		Draught in metres.
gross_tonnage	character(64,1)		Gross tonnage.
kilowatts	character(128,1)		Engine power in kilowatts.
engine_make	character(128,1)		Engine make.
Creator:	dba		
Referential:	(vid) REFERRED t_vessel (vid)		
Indices:	PRIMARY KEY BTREE ON (vid)		

5.5 Table 5: t_lloyds

Comment: This table contains selected information from the Lloyds Register of Shipping record.

Attributes	Data Type	Null?	Comment
imo_no	longinteger		Unique International Maritime Organisation number assigned by Lloyd's Register to each ship.
vessel_name	character(32,1)		The current name of the ship.
ship_status	character(32,1)		In Service Commission(S), Laid Up(L), In Casualty or Repairing(R), Converting / Rebuilding(C), To Be Broken Up(T), Unconfirmed Ships(X).
ship_type	character(16,1)		Bulk Carrier, Container Ship, Dredger, Fishing, General Cargo Ship, Icebreaker, (Offshore) Supply Ship, Passenger Ship, Refrigerated Cargo Ship, RoRo Cargo Ship, Tanker, Vehicles Carrier.
flag	character(32,1)		Indicates the flag country of registry under which the ship normally operates.
gross_tonnage	integer		Gross tonnage.
date_build	character(16,1)		The year and month when the new construction survey process is completed.
year_built	character(4,1)		Reported year of completion of construction.
lgth	decimal(6,2)		Overall Length between perpendiculars else Registered length.
breadth	decimal(5,2)		Extreme Breadth else Moulded Breadth of the vessel.
draught_max	decimal(5,2)		In most cases this is the maximum summer draught amidships.
call_sign	character(16,1)		Signal letters or radio call sign assigned by the relevant Registration (Flag) Authority.
official_no	character(16,1)		The identification number assigned by the national authority.
fishing_no	character(16,1)		The identification number assigned by the national authority to ships engaged in the fishing industry.
port_registry	character(32,1)		Place where the ship is registered. Home port is shown where there is no port of registry.
country_benefit	character(32,1)		The country considered to be the main beneficiary from the earnings generated by the operation of the vessel.
net_tonnage	integer		Represent a measure of the ship's freight earning capacity.
length_overall	decimal(6,2)		The extreme length of the ship.
length_registered	decimal(6,2)		Measured from the extreme fore point of the hull to the after end of the stern post, or if there is no stern post to the fore side of the rudder stock.
engine_type	character(32,1)		Oil, Steam Turbine or Steam Reciprocating.
engine_builder	character(128,1)		Manufacturer of the main engine.
power_kw	integer		Power in kilowatts. Power for LR Class ships is Design Power; for non-LR Class Service Power is usually recorded.
power_bhp	integer		Power in bhp. Power for LR Class ships is Design Power; for non-LR Class Service Power is usually recorded.

subtype1	character(32,1)	Subtypes of vessel, e.g. if fishing or reefer, and types of fishing vessel by method.
record_no	longinteger	Non-lloyds field, unique number assigned while loading the data.
Creator:	dba	
Referential:	(imo_no) REFERRED t_vessel (imo_no)	
Indices:	NORMAL (2, 15) BTREE ON (vessel_name)	
	NORMAL (2, 15) BTREE ON (call_sign)	
	PRIMARY KEY BTREE ON (imo_no)	

5.6 Table 6: t_rule

Comment: This table contains the rules used in the grooming process.

Attributes	Data Type	Null?	Comment
rule_code	character(16,1)		Alphanumeric Code for a rule.
description	character(256,1)		Description of the rule.
Creator:	dba		
Indices:	PRIMARY KEY BTREE ON (rule_code)		

6 vessel business rules

6.1 Introduction to business rules

The following are a list of business rules applying to the **vessel** database. A business rule is a written statement specifying what the information system (i.e., any system that is designed to handle vessel specifications data) must do or how it must be structured.

There are three recognised types of business rules:

Fact	Certainty or an existence in the information system.
Formula	Calculation employed in the information system.
Validation	Constraint on a value in the information system.

Fact rules are shown on the ERD by the cardinality (e.g., one-to-many) of table relationships. Formula and Validation rules are implemented by referential constraints, range checks, and algorithms both in the database and during validation. These rules state that a value **must** meet the specified criteria.

Validation rules may be part of the preloading checks on the data as opposed to constraints or checks imposed by the database. These rules sometimes state that a value should be within a certain range. All such rules containing the word 'should' are conducted by preloading software. The use of the word 'should' in relation to these validation checks means that a warning message is generated when a value falls outside this range and the data are then checked further in relation to this value.

6.2 Summary of rules

Vessel details (*t_vessel*)

<i>record_no</i>	Must be a unique integer.
<i>vid</i>	Must be an integer greater than zero, and must be a unique value for a specific hull.
<i>cnt</i>	Must be an integer greater than zero.
<i>vessel_key</i>	Must be an integer greater than zero, and must be a unique value for a specific hull.
<i>spec_from</i>	Must be a valid date, and should be between 1978 and the current date.
<i>spec_to</i>	Must be a valid date, and should be between 1978 and the current date, or 2099, or 2999.
	Multiple column checks on spec dates: The <i>spec_from</i> date must not be greater than the <i>spec_to</i> date.
<i>vessel_id</i>	Should consist only of the characters A-Z and 1-9.
<i>call_sign</i>	Should be a valid NZ or international radio call sign. Typically 4 to 7 alpha numeric characters.
<i>imo_no</i>	Must be an integer greater than zero.
<i>vessel_nation</i>	Should be a valid MFish 3 character code for vessel nationality.
<i>reg_type</i>	Must be a 1 character code, and should be one of ('C', 'D', 'F', 'U').
<i>year_build</i>	Must be an integer greater than zero, and should be a valid year between 1900 and the current year.
<i>overall_length</i>	Must be an integer greater than zero, and should be between 20 and 150 metres.
<i>breadth</i>	Must be an integer greater than zero, and should be between 3 and 24 metres.
<i>draught</i>	Must be an integer greater than zero, and should be between 1 and 12 metres.
<i>gross_tonnage</i>	Must be an integer greater than zero, and should be between 1 and 9000.
<i>kilowatts</i>	Must be an integer greater than zero, and should be between 40 and 9000.
<i>notes</i>	Should contain codes separated by ' ', as documented in the <i>t_rule</i> table

Metadata (t_meta)

owner_key Must be a unique integer.

load_date Must be a valid date, and should be between 2004 and the current date.

History data (t_history)

vid Must be an integer greater than zero, and must be a unique value for a specific hull.

Rules details (t_rule)

rule_code Rule code must be unique.

For the Lloyds reference data in table *t_lloyds*, as supplied by Lloyds Register of Shipping, business rules are not applied other than enforcing numeric data types. This is because the Lloyds data are a reference data set and were not groomed or checked before loading to this table.

For the original data in table *t_vessel_org*, as supplied by MFish, business rules are not applied other than enforcing numeric and date data types. This is to retain the original data as supplied by MFish.

7 Acknowledgments

The authors would like to thank Peter Shearer for his review and editorial comment for this document.

Appendix 1 Reference code tables

Nationality type	Nationality type description
AUS	Australia
BZE	Belize
CHI	China (People's Republic of)
COO	Cook Islands
FIJ	Fiji
GRE	Greece
JAP	Japan
KOR	Korea
NOR	Norway
NZL	New Zealand
PHI	Philippines
POL	Poland (Republic of)
RUS	Russian Federation
SNG	Singapore (Republic of)
SVG	Saint Vincent
TAI	Taiwan
UKR	Ukraine
USA	U.S.A
VAN	Vanuatu (Republic of)

Codes used for the reg_type attribute, documenting the registration type.

Vessel reg type	Vessel reg type description
C	Charter
D	Domestic
F	Foreign License
U	Unknown

Appendix 2 Data Grooming Rules

Rules 201 to 207 are applied by script, and 211 to 250 are used in manual grooming process.

rule_code description

201	Identify vessels by 4 major id attributes - vessel_key, vessel_id, call_sign and vessel_name, together with overall_length, built_year as major reference parameters, and tonnage, breadth, draught as secondary reference parameters.
202	If two vessels have no common id attribute, they are different vessels unless they have all the same reference parameters and subject to manual id process.
203	Two vessels are the same vessel if both have one common id attribute and either a. the same imo number or b. the same built year and overall length and tonnage and breadth and draught.
204	Two vessels are the same vessel if both have two common id attributes and either a. the same imo number or b. the same year built and overall length and tonnage and breadth and draught.
205	Two vessels are the same vessel if both have three common id attributes and either a. the same imo number or b. the same built year and overall length and tonnage.
206	Two vessels are the same vessel if both have four common id attributes, ie vessel_key, vessel_id, call_sign and vessel_name.
207	Two vessels are the same vessel if both have three common id attributes and one null id attribute and either a. the same imo number or b. two out of the three of the following attributes built year, length, tonnage.
208-210	Rules 208, 209 & 210 do not exist.
211	Same as 212.
212	If 6 attributes agreed then I considered that only one hull existed for the two sets of data.
213	If 5 attributes agreed (provided one of the attributes that did not agree related to vessel_name, call_sign, or a dimension, and was minor) then one hull existed for the two sets of data.
214	If 4 attributes agreed (provided at least 1 of the attributes that did not agree met one of the rules listed in 213) then one hull existed for the two sets of data.
215	If 4 attributes agreed, and it was not possible to compare other attributes (due to the presence of null values) then only one hull existed for the two sets of data, but only if there was compelling data to support this assumption. For example, imo numbers were the same. It was decided that 4 was the minimum number of attributes required for comparison, as any less meant that it was not possible to make a meaningful conclusion.
216	In the instance that the conditions in rule 215 are not met, or fewer than 4 attributes were present it was assumed that 2 hulls existed.
232	Same as 212 except there are three agreed ID attributes (vessel_name, call_sign, vessel_id, vessel_key).
233	Same as 213 except there are three agreed ID attributes (vessel_name, call_sign, vessel_id, vessel_key).
234	Same as 214 except there are three agreed ID attributes (vessel_name, call_sign, vessel_id, vessel_key).
235	Same as 215 except there are three agreed ID attributes (vessel_name, call_sign, vessel_id, vessel_key).
236	Same as 216 except there are three agreed ID attributes (vessel_name, call_sign, vessel_id, vessel_key).
242	Same as 212 except there are three agreed ID attributes and one null ID attribute.
243	Same as 213 except there are three agreed ID attributes and one null ID attribute.

244 Same as 214 except there are three agreed ID attributes and one null ID attribute.

245 Same as 215 except there are three agreed ID attributes and one null ID attribute.

246 Same as 216 except there are three agreed ID attributes and one null ID attribute.

BR rules applied to breadth

CL rules applied to collapse procedure

CS rules applied to call sign

DR rules applied to draught

EM rules applied to engine maker

EMrMAN Manually simplify engine maker.

EMrSTN for engine make string, first char of each word set upper case.

GT rules applied to gross tonnage

GTrOv if gross tonnage > 8000t, set to null.

IM rules applied to imo number

IMrDIS same vid got different IMO numbers, manual fix.

IMrDUP different VIDs have the same IMO number, manual fix.

IMrNU use unique IMO no to fill in all the nulls.

KW rules applied to kilowatts

NT rules applied to vessel nationality

OL rules applied to overall length

RT rules applied to vessel registration type

RTrU code in C/D/F/U, set to U for unknown or null if only one code exists except U, then set all to the code.

VK rules applied to vessel key of MFish

VKrDUP different VIDs have the same vessel key, manual fix.

VKrMin use the lowest vessel_key where a vessel has multiple vessel_keys.

VN rules applied to vessel name

VNrSTN for name string, a. first char of each word set upper case; b. for alpha-numeric names, use "No." in front of the number if "No." appears at least once; c. if name strings contain Arabic and Roman numbers, use Roman number; d. strip off leading F. V.

YB rules applied to year built

YBrIMO for multiple year built values, chose the one that also in Lloyds record, this rule overwrites YBrMF.

YBrMF for multiple year built values, chose the one that occurs most frequently.

rABA if two values in pattern of a*b*a*, then a is adopted.

rConv if two values are within 5% difference by unit converting to metric system, then use the value in metric unit.

rDCI decimals round to nearest integer

rMAN manual fix.

rNEW manually added new value.

rNMF if absolute difference between two numbers are within 5%, then choose the most frequent one.

rNUImo populate a null attribute with Lloyds record.

rNUP populate nulls with values in the following patterns: a. NNNv to vvvv; b. vNNv to vvvv; c. vNNvNu to vvvvNu; d. vNNuNuzNNz to vNNuuuzzzz. Where 'N' represents a null value.

rNUT for a null/unknown value, if there's a value in the same time period, then use this value.

rONE If the difference between two big integers is 1 then merge to the most frequent one.

rSIM If both strings are very similar due to misspelling then assess occurring frequency of each value, and adopt the most frequent one.

rSTU set string to upper case.

Appendix 3 Final Manual Grooming Procedures

Processing of data where IMO numbers were present and the Lloyd's Register could be referred to

The Lloyd's Register was referred to in all cases where IMO numbers existed and different values occurred for a single vid. These data represented about 70% of the final ungroomed dataset.

The Lloyd's Register was of limited use for call sign and registration type because these fields were often null. Only when the discrepancies were obvious were changes made. For example, on one occasion a vessel had been given the call sign '7 Jun' (which appears to be part of a date and not a call sign) and this was changed to the call sign recorded in the Lloyd's Register. It is recognised that this approach does not fully address the issue of differences in call sign and registration type.

The notation used to document amendments using the Lloyd's Register is an extension of that listed in Appendix 2, and includes the following prefixes:

Vessel name (VN)
Callsign (CS)
IMO number (IM)
Nation (NT)
DOB (YB)
OA Length (OL)
Breadth (BR)
Draught (DR)
Gross tonnage (GT)
Kilowatts (KW)
Engine make (EM)

1. RULE: - For static attributes, such as year built

Where different values occurred for a single vid and one of these values agreed with what was recorded in the Lloyd's Register then this value was assumed to be the correct one. The notation **YBrIMO** was used to document this change.

2. RULE: - For non-static attributes, such as vessel name, overall length, kilowatts

The same approach was taken for non-static attributes as for static attributes outlined above in point 1.

However, the history table in the Lloyd's Register was manually accessed in all cases to check for changes in attributes, and/or conversions. If it was unclear whether the difference was legitimate, then it was assumed that it was a legitimate change of a value in *t_vessel*. The notation **VNrIMO|OLrIMO|KW rIMO,** was used to document these changes.

Processing of data where IMO numbers were not present and the Lloyd's Register could not be referred to

The frequency (count) of each attribute was examined in cases where no IMO numbers existed and different values occurred for a single vid. These data represented about 30% of the final ungroomed dataset.

For example, if a vessels breadth was 5m on 20 occasions, and 7m on 2 occasions then it was assumed that 5m was the correct breadth. It would be unusual for a vessel to be widened.

The notation used to document such amendments was **BRrMF** (Breadth changed to most frequent value).

3. RULE: - Most frequent value

It was considered that up to a 10% difference in values, between a given numeric attribute, was an acceptable level to adopt this “most frequent” approach.

However, if differences of more than 10% existed for a specific attribute, but a clear pattern was evident, it was also considered appropriate to adopt a single value. For example, if the kilowatts were:

Count	Kilowatts	Engine Make
10	2800	Caterpillar 3142
8	1750	Caterpillar 3142
9	2800	Caterpillar 3142

In this case 1750 would be changed to 2800 because there is no evidence that Engine make/model has changed. However, if the Engine make attribute was null then it would not have been possible to conclude whether 1750 was legitimate or not. Therefore, the value for kilowatts would not have been changed.

4. RULE: - Null values

In a few cases the null rules were used. These rules include:

rNUT- If the value is null/unknown and there is a value in the same time period, then use this value.

rNUP – nulls were populated with values in the following patterns: a. vNNv to vvvv; b. vNNvNu to vvvvNu; c. vNNuNuzNNz to vNNuuuzzzz. Where ‘N’ represents a null value

Checking of range constraints

Following the grooming process each attribute was examined to identify improbable values. For example, Breadth was restricted to the range 1-24m; Draught to the range 1-12m. There were a number of instances where these attributes had very large improbable values, however it was generally easy to determine what the error was. For example, in one instance a breadth of 103m was due to the fact that the overall length had been placed in the wrong column.

Checking of a data subset

Approximately 10% (i.e., 150) of the vessels were checked visually against the *t_vessel_org* and *t_history* tables. This was to ensure that the grooming process had worked appropriately. No obvious discrepancies were apparent.

General comments

The two main attributes that required grooming were gross tonnage and kilowatts. Gross tonnage was easier to groom than kilowatts. End users of the data should be aware that these two attributes are likely to be the least robust in this dataset because they were often null or inconsistent in the Lloyds dataset. We are confident that all other attributes, including year built, overall length, breadth, and draught are reasonably robust.