

abaraw: A shorthand notation for bird records

Zoological
Data Processing

Toward more efficient data entry

John W. Shipman

2013-02-07 13:30

Abstract

Describes a system for the rapid and efficient encoding of ornithological field records.

This publication is available in Web form¹ and also as a PDF document². Please forward any comments to john@nmt.edu.

Table of Contents

1. Introduction	1
2. Installation	2
3. Operation of the <i>abaraw</i> script	2
4. The input syntax	3
4.1. Syntax of the day line	3
4.2. The locality definition line	4
4.3. The locality back-reference line	4
4.4. The census line	4
4.5. A small, complete example file	6

1. Introduction

In the document *A system for encoding bird field notes*³, we describe a technique for using an XML document type to represent field records of wild birds. The author has used this system for some time to publish his field notes⁴.

Initially, the author used the `emacs` text editor and the related `nxml-emacs` package to enter and maintain the XML notes files; see the documentation for `nxml-emacs`⁵.

However, for most notes, that system proved to be somewhat slow and cumbersome. The current document describes a “shorthanding” system that allows rapid creation of most of the XML from a terse textual notation that greatly speeds up the entry of *most* of the information.

The shorthand notation is not intended to create every type of information in the XML field notes schema. This schema provides for a wide variety of types of information: notes on breeding status, vocalization, and many other kinds of specialized content. However, the great bulk of the records are pretty simple—what kind of bird, how many.

¹ <http://www.nmt.edu/~shipman/aba/raw/doc/>

² <http://www.nmt.edu/~shipman/aba/raw/doc/abaraw.pdf>

³ <http://www.nmt.edu/~shipman/aba/doc/>

⁴ <http://www.nmt.edu/~shipman/aba/field.html>

⁵ <http://www.nmt.edu/tcc/help/pubs/nxml/>

The author decided that a mixed strategy might be best: develop a terse shorthand notation that represents most of the data, and use that to create the bulk of the XML notes file; then use the full XML editor to add in anything that isn't captured in the shorthand notation.

Hence, the workflow goes like this:

1. Transcribe the field notes into the shorthand notation as a file whose name ends with the extension “.in”. Omit the kinds of information that cannot be represented in the shorthand notation (breeding, vocalization, etc.).
2. Use the *abaraw* script to process that file, creating an XML file whose name ends with “.out”.
3. Use *emacs* with *nxml -mode* to add any information that was omitted in the shorthand notation.

At this writing, the author has been using the shorthand system for several months, and finds that it greatly reduces the time required for data entry of field notes.

2. Installation

The directory in which this script is run must have the following files available, as either copies or soft links:

- `abaraw`⁶: The *abaraw* script itself.
- `abbr.py`⁷: Python module to process bird abbreviations.
- `aou.xml`⁸: XML file defining the latest AOU Check-List taxonomy.
- `rnc.py`⁹: Python constants for names from the schema.
- `rnc_txny.py`¹⁰: Used by `txny.py`.
- `txny.py`¹¹: Module to read the taxonomy authority file.

Additionally, the Python `lxml` package must be installed; see *Python XML processing with lxml*¹².

3. Operation of the *abaraw* script

An input file must have a name of this form:

```
yyyy-mm.in
```

The `yyyy-mm` part gives the year and month of the notes. For example, a notes file for July 2008 must be named “2008-07.in”.

To process one file, use this command:

```
abaraw yyyy-mm.in
```

Output is sent to the standard output stream. Typically you will want to redirect it to a file named `yyyy-mm.out`. For example, to process the July 2008 file:

⁶ <http://www.nmt.edu/~shipman/aba/raw/doc/ims/abaraw>
⁷ <http://www.nmt.edu/~shipman/xnomo/ims2/abbr.py>
⁸ <http://www.nmt.edu/~shipman/xnomo/aou/aou.xml>
⁹ <http://www.nmt.edu/~shipman/aba/doc/pyims/rnc.py>
¹⁰ http://www.nmt.edu/~shipman/xnomo/ims2/rnc_txny.py
¹¹ <http://www.nmt.edu/~shipman/xnomo/ims2/txny.py>
¹² <http://www.nmt.edu/tcc/help/pubs/pylxml/>

```
abaraw 2008-07.in >2008-07.out
```

4. The input syntax

Each input file is a mixture of these kinds of lines:

- A *day line* starts with a bang (“!”) character and defines the date, state code, and default locality for the lines that follow it, up to the next day line or to the end of the file. See Section 4.1, “Syntax of the day line” (p. 3).
- A *locality definition line* starts with an at-sign (“@”) character and defines a new locality code that applies to following lines up to the next locality line, day line, or end of file. See Section 4.2, “The locality definition line” (p. 4).
- If the observer leaves a particular locality but then returns to that locality later in the same day, a *locality back reference line* specifies the locality that applies to following lines up to the next locality line, day line, or to the end of the file, whichever comes first. See Section 4.3, “The locality back-reference line” (p. 4).
- A *census line* describes the sighting of one or more kinds of birds. See Section 4.4, “The census line” (p. 4).

The reader may wish at this time to review the schema¹³, since the output of the shorthanding system is destined for an XML file conforming to that schema. We will use names from the schema to indicate where shorthand items will go in the output.

4.1. Syntax of the day line

Here is the general form of a day line:

```
!state yyyy-mm-dd locality-def
```

state

The state or region code, to be copied to the `state` attribute of `day-notes`.

yyyy-mm-dd

The date of the following records; this will be copied to the `date` attribute of the `day-notes` element.

locality-def

The rest of the line has the same format as a locality definition line, including the initial “@”. This locality’s code will be copied to the `day-loc` attribute of the `day-notes` element, and its definition will become a `loc` element. See Section 4.2, “The locality definition line” (p. 4).

Here are some examples of day lines.

```
!nm 2008-06-15 @SocCo Socorro County
!nm 2008-07-22 @BdA Bosque del Apache NWR
!ca 2009-07-10 @PAB Palo Alto Baylands
```

¹³ <http://www.nmt.edu/~shipman/aba/doc/schema.html>

4.2. The locality definition line

To define the location of one or more sightings, a locality definition line defines a new locality code and the corresponding full description. This locality is applied to all the census lines that follow, up to the next locality line or day line (or to the end of the file).

Here is the format of this form of locality line:

```
@code loc-name
```

In the output file, the *code* is copied to the `code` attribute of the `loc` element, and the `loc-name` becomes the `name` attribute of that `loc` element.

Examples:

```
@WCRd Magdalenas: Water Canyon Road
@ABQ Albuquerque
@TerryL Socorro: New Mexico Tech: Terry Lake
```

4.3. The locality back-reference line

If the observer leaves a certain locality and then returns to it later, it is not necessary to repeat the entire locality definition line. Just use a locality back-reference line, which has this general form:

```
@code
```

The supplied locality *code* is applied to all following census lines up to the next locality line, day line, or end of file. The code must have been defined on a locality definition line used earlier in the same day.

Examples:

```
@SocCo
@BdA
@PAB
```

The back-reference form may not be used in a day line.

4.4. The census line

Any line that does not start with “!” or “@” is a census line. A census line consists of one or more *bird groups* separated by spaces.

Each bird group describes the sighting of one particular kind of bird. Each group results in the output of one `form` element in the XML output, and has this syntax:

- A *bird ID group* that describes the taxonomic position of the bird or birds seen. Generally this will be a single six-letter bird code, but it may be a pair of codes describing a hybrid or species pair.

The bird ID group for a hybrid has this form.

```
ab6^alt
```

The *ab6* and *alt* codes are the six-letter codes for the probable parents of the hybrids. For example, this group describes a probable hybrid of American Wigeon and European Wigeon:

```
amewig^eurwig
```

The bird ID group for a species pair has this form:

`ab6|alt`

The `ab6` and `alt` codes are the six-letter codes for the two forms that the observer could not distinguish. For example, this group describes a bird that might have been either a Dusky Flycatcher or a Hammond's Flycatcher:

`dusfly|hamfly`

For a single code, the code will appear as the `ab6` attribute of the `taxon-group`. For hybrids and species pair records, the two codes will appear as the `ab6` and `alt` attributes, and the `^` or `|` will appear as the `rel` attribute of the `taxon-group`.

- If the record is considered notable, an exclamation point `!` follows the bird ID group.
- Next are zero or more *suffix groups*, as defined below. If there are multiple suffix groups, each group becomes an `floc` child element of the record's `form` element.
- If there is information to be added later while editing the XML output, a suffix `*`, called the *long flag*, is added. This forces any output `form` or `floc` elements to be represented as a pair of tags, rather than as an empty tag. This makes it easier to add content to the record later during direct editing of the XML.

Here is the general form of a suffix group, where square brackets indicate an optional element.

`count-group age-group sex-group status-group`

count-group

This group indicates the number of individuals seen. For permissible values, refer to the definition of the `count` attribute of the `age-sex-group` pattern in the schema. Examples: `1`; `57`; `8-10`; `6+`; `40-;` `#`.

age-group

Any of the permissible `age` attributes of the output: `"a"`, `"i"`, or `"p"` (for "female or immature").

sex-group

An optional sex code, `"m"` or `"f"`.

status-group

A code indicating the status of the sighting.

<code>?</code>	Questionable identity; outputs an attribute <code>"q='?'"</code> .
<code>=</code>	Not countable under American Birding Association rules; outputs an attribute <code>"q='-'"</code> .

Here are some examples of complete bird groups and their interpretation.

<code>rinduc</code>	One or more Ring-necked Ducks.
<code>cangoo12</code>	Twelve Canada Geese.
<code>whfibi#</code>	Two or more White-faced Ibis.
<code>heptan!</code>	Hepatic Tanager, notable record.
<code>westan1m</code>	One male Western Tanager.
<code>haiwool-2</code>	One or two Hairy Woodpeckers.


```

<!--=====-->
<day-notes date="2008-04-04" day-loc="Soc" state="nm">
  <day-summary default-loc="Soc">
    <loc code="Soc" name="Socorro"/>
    <loc code="home" name="Socorro: 507 Fitch"/>
    <loc code="Terry" name="Socorro: New Mexico Tech: Terry Lake"/>
    <loc code="NMT" name="Socorro: New Mexico Tech"/>
    <loc code="Speare" name="Socorro: New Mexico Tech: Speare Hall"/>
  </day-summary>
  <form ab6="rinduc">
    <floc loc="Terry" sex="f" count="1"/>
    <floc loc="Terry" sex="m" count="1"/>
  </form>
  <form ab6="rudduc" loc="Terry" sex="f" count="1"/>
  <form ab6="amecoo">
    <floc loc="Terry" count="10"/>
    <floc loc="NMT" count="2"/>
  </form>
  <form ab6="belkin" loc="NMT" count="1"/>
  <form ab6="saypho" loc="Speare"/>
  <form ab6="barswa" loc="Terry"/>
  <form ab6="grtgra" loc="NMT" sex="m" count="1"/>
  <form ab6="casfin">
    <floc loc="home" sex="m" count="1"/>
    <floc loc="home" sex="f" count="2"/>
  </form>
  <form ab6="houfin" loc="Speare">
  </form>
</day-notes>
</note-set>

```

