

Accessing *Starling* Databases

Rob Verhoeven

April 3, 2014

Abstract

This document contains information over the way information is stored in and can be extracted from *Starling* databases (`.dbf`). The information has been compiled using existing documentation on *dBase III* databases and reverse engineering efforts by mainly prof.dr. Andries Brouwer.

Document status

revision	date	reason for change
1.0	28-11-2013	initial document release
1.1	18-02-2014	update of special character table and addition of note on packing

Contents

1	<i>Starling</i>	3
2	Database header	3
3	Records	3
3.1	Record header	3
3.2	Record contents	4
4	References to associated .var file	4
5	Character encodings	5
5.1	Single byte mode	5
5.2	Double byte mode	7
6	Layout codes	10
7	Auxiliary files	10
7.1	.inf file	10
7.2	.prt file	10

1 *Starling*

Starling is a tool for comparative linguists, which allows them to catalogue linguistic data in a flat database and apply a diverse range of procedures upon that data. The tool was created by Sergei Starostin, who passed away in 2005, and is currently maintained by his son George.¹

Much credit goes to prof.dr. Andries Brouwer,² who, through personal correspondence, shared with me some of his code for extracting data from *Starling* databases, allowing me to profit from his reverse engineering efforts on the database structure and encoding schemes.

2 Database header

Starling databases (.dbf) are in an adapted *dBase III* database format. Going by the *dBase III* header as documented on the *DBF Viewer 2000* website³, the header format seems to be as follows:

bytes	description
0	.dbf file type 0x03 dBase III
1	Year of last update (offset from 1900)
2	Month of last update
3	Day of last update
4–7	Number of data records
8–9	Offset of first data record (i.e. header size (HDR))
10–11	Record size (including delete flag)
12–27	Reserved bytes (unknown content)
28	Table bit flags (zero by default?) 0x01 has a database index (.cdx) 0x02 has a memo field 0x04 is a database container (.dbc)
29	Code page mark (zero by default)
30–31	Reserved bytes (0 and 1 by default, respectively)
32–HDR	Record structure information (see Sec. 3.1)

3 Records

To avoid complications, it is assumed that we are only dealing with packed databases, i.e. databases on which the operation **Pack** from the menu **Assist** has been invoked. This avoids having to deal with fragmentation in *Starling*'s variable length fields (which are discussed in Sec. 4).

3.1 Record header

Record structure information is included in the *Starling* database header, where each field is specified in a block of 32 bytes.

¹<http://starling.rinet.ru/>

²<http://www.win.tue.nl/~aeb/>

³<http://www.dbf2002.com/dbf-file-format.html>

bytes	description
0–10	Name (padded with null characters)
11	Field type (ASCII)
	C character
	Y currency
	N numeric (can be a real number)
	F float
	D date
	T datetime
	B double
	I integer
	L logical
	M memo
	G general
	P picture
12–15	Offset of the field, within the record. (Do not trust this value.)
16	Field length
17	Number of decimal places
18	Field bit flags
	0x01 system field (invisible to user)
	0x02 nullable
	0x04 binary (for field types <i>C</i> and <i>M</i> only)
19–31	Reserved bytes

3.2 Record contents

Each records start with a ‘delete byte’ with the following possible values.

0x20	normal record
0x2A	deleted record

Numeric fields have their value stored in text (each digit is represented by its ASCII encoded value). Furthermore, fields containing character data are padded with null characters (0x00).

4 References to associated .var file

A **.dbf** file is typically accompanied by a homonymous **.var** file. Due to the way records are specified in *dBase III*, a field (containing character data) has a maximum length of 255 characters. To accommodate longer texts, *Starling* databases treat a character field with a length of 6 bytes as a reference to a chunk of data, stored within the associated **.var** file.

bytes	description
0–3	data offset
4–5	data length

When a field that normally contains such a reference should instead be empty, its 6 bytes are all ASCII encoded spaces (0x20).

5 Character encodings

Stemming from a time before Unicode, *Starling* works with its own character encoding. Unfortunately, this encoding is not straightforward. Much work has been done on reverse engineering this encoding by prof.dr. Andries Brouwer and the following information relies heavily on the information compiled on his website.⁴

Characters are encoded in either one or two bytes. By default, upcoming data is to be read in ‘single byte mode’. Byte flags are used to indicate when upcoming text is to be interpreted in ‘double byte mode’ or when to return to ‘single byte mode’.

0x01	Enter ‘double byte mode’
0x7F	Return to ‘single byte mode’

There are, however, complications. When in double byte mode, encountering a byte in the ASCII range (below **0x7F**) also indicates a return to single byte mode (with the exception of the **0x01** flag for double byte mode). However, combining characters that are encoded in a single byte can seemingly occur in both single and double byte mode without affecting the mode.

5.1 Single byte mode

Among the byte values up to **0x1F** (the range for ASCII control characters) are those for switching byte mode (see above). In addition to those, the following bytes have their own semantics:

0x09	Tab character.
0x0A	Newline.
0x0D	Occasionally follows a 0x0A , together forming one newline.
0x15	New paragraph.
0x1D	Marks the next byte as ‘special’ (see below).

The characters encoded by the remainder of the code points (**0x20–0xFF**) is shown in Table 1. The encoding is a customization of the Cyrillic code page CP866.

In case a byte is marked as special, by a preceding byte **0x1D**, the byte encodes the characters displayed in Table 2.

⁴<http://www.win.tue.nl/~aeb/natlang/charsets/starling-charset.html>

	_0	_1	_2	_3	_4	_5	_6	_7	_8	_9	_A	_B	_C	_D	_E	_F
2_	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	
3_	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
4_	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
5_	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	◌̂	_
6_	‘	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
7_	p	q	r	s	t	u	v	w	x	y	z	{		}	◌̃	
8_	А	Б	В	Г	Д	Е	Ж	З	И	Й	К	Л	М	Н	О	П
9_	Р	С	Т	У	Ф	Х	Ц	Ч	Ш	Щ	Ъ	Ы	Ь	Э	Ю	Я
A_	а	б	в	г	д	е	ж	з	и	й	к	л	м	н	о	п
B_	ā	◌̇	ă	◌̈	ǎ	ç	č	ċ	ō	ē		◌̊	ö	ε	ƒ	◌̋
C_	ç	γ	ѡ	ħ	◌̄	ı	ı̇	ı̈	◌̌	◌̍	ķ	λ	λ̇	-	λ̈	†
D_	Ł	η	ō	ö	ō̄	ɔ	ō	ř	q̇	ß	~	◌̊	◌̋	š	ţ	◌̌
E_	р	с	т	у	ф	х	ц	ч	ш	щ	ъ	ы	ь	э	ю	я
F_	ø	ū	ü	ũ	ə	ē	◌̇	ʷ	ħ	χ	Ʒ	ž	ž	ʔ	Ɔ	Λ

Table 1: The range of characters encoded by a single byte, each along with their Unicode code point(s). A red background denotes the ASCII code space, a green background denotes a character from the CP866 code page, and a blue background denotes characters from the custom *Starling* encoding. Red characters indicate combining characters.

	_0	_1	_2	_3	_4	_5	_6	_7	_8	_9	_A	_B	_C	_D	_E	_F
4_		Æ 04D4			đ 0111								† 026C			
5_					ț 0167										◌̈ 0311	
6_		æ 00E6	ð 0180	ƒ 0255	ð 00F0	œ 0153		ɡ 01E5	ĥ 1E2B	ı 0131	ĵ 0237					ø 00F8
7_			ŗ 027E	ſ 0283	þ 00FE				ž 1D9A							
8_																
9_																
A_				ġ 0260	À 0467	Š 0455						Љ 0459		Њ 045A	Ћ 046B	
B_		◌̈ 030B														◌̈ 032E
C_																
D_													◌̈ 030F			
E_	Ѕ 0281													€ 0454		Ž 017E
F_								ћ 0195	ĥ 0452							

Table 2: Characters encoded by a ‘special’ byte in single byte mode.

5.2 Double byte mode

A sequence of two bytes encodes a character using the Chinese Big5 encoding,⁵ of which the first byte always has a value greater than 0x7F. *Starling* makes use of code points in the space reserved for user-defined characters (0x81–0xA0) for its own character encoding. The following provides an overview of *Starling*’s user defined characters.

1st byte	2nd byte	character
0x83		(refer to Table 3)
0x85	0xAF	ƒ (U+03DD)
0x87		(refer to Table 4)
0x88	0x81	ŗ (U+0475)
	0x83	À (U+0467)

⁵Cf. Microsoft’s code page 950.

	_0	_1	_2	_3	_4	_5	_6	_7	_8	_9	_A	_B	_C	_D	_E	_F
9_	ᾱ	ᾰ	ᾡ	ᾠ							ᾶ	ᾷ	Ᾱ	Ὰ	·	
	0344	0314 0301	0313 0301	0301							0308 0300	0314 0300	0313 0300	0314 0342	0387	
	ᾱ	ᾰ		ᾲ	A	B	X	Δ	E	Φ	Γ	H	I	Ϝ	K	Λ
	0314	0313		0308	0391	0392	03A7	0394	0395	03A6	0393	0397	0399	1FF3	039A	039B
	M	N	O	Π	Θ	P	Σ	T	Υ	η	Ω	Ξ	Ψ	Z		
B_	039C	039D	039F	03A0	0398	03A1	03A3	03A4	03A5	1FC3	03A9	039E	03A8	0396		
C_	ᾶ	ᾷ	α	β	χ	δ	ε	φ	γ	η	ι	ς	κ	λ	μ	ν
	0313 0342	0300	03B1	03B2	03C7	03B4	03B5	03C6	03B3	03B7	03B9	03C2	03BA	03BB	03BC	03BD
D_	ο	π	θ	ρ	σ	τ	υ	ϙ	ω	ξ	ψ	ζ	ᾰ			
	03BF	03C0	03B8	03C1	03C3	03C4	03C5	1FB3	03C9	03BE	03C8	03B6	0342			

Table 3: The Greek characters encoded by the second byte in double byte mode, when the first byte is 0x83. The ‘α’ that is present with the combining characters (shown in red) is merely there for reference.

	_0	_1	_2	_3	_4	_5	_6	_7	_8	_9	_A	_B	_C	_D	_E	_F
8_				ѓ 0433 0311	б 0431		ю 044E									
9_				ж 0436			Ю 042E				И 0418	С 0421			А 0410	П 041F
A_	Р 0420		О 041E	Л 041B	Д 0414				З 0417		К 041A		Е 0415	Г 0413	М 041C	
B_		Н 041D		х 0445					ъ 044A	ф 0444	и 0438	с 0441	в 0432	у 0443	а 0430	п 043F
C_	р 0440	ш 0448	о 043E	л 043B	д 0434	ь 044C	т 0442		з 0437		к 043A	ы 044B	е 0435	г 0433	м 043C	ц 0446
D_	ч 0447	н 043D	я 044F	Х 0425		ѣ 0463			ѧ A657		Ѧ 0467					
E_								Ѫ 046B	Е 0415	е 0065						
F_			Ѡ 0461	Ѳ 046C	Ѵ 046D	Ѷ 0464	Ѹ 0465									

Table 4: The Cyrillic characters encoded by the second byte in double byte mode, when the first byte is 0x87.

6 Layout codes

Starling allows character data to be formatted. For this, it uses layout tags that are inserted into the character data. These tags are the following:

tag	formatting style
\B...\b	bold text
\I...\i	italic text
\U...\u	underline text
\H...\h	superscript text
\L...\l	subscript text
\C...\c	condensed text

In addition to these layout tags, there are also tags that produce links. The exact format of these links is somewhat unclear, but seems to be as follows:

tag	formatting style
\X0...\x	Link to an entry containing the tagged text. The text itself also makes up the clickable link. (Perhaps the number 0 is variable and refers to the column number.)
\X<0.n>...\x	The clickable link is again the text between the tags, but the target is now entry <i>n</i> . (Perhaps the number 0 is variable and is an index to a list of different databases.)

Starling does not enforce the contents of a cell to have correctly nested tags. Even closing tags without a preceding opening tag can occur.

7 Auxiliary files

A *Starling* database file (**.dbf**) can be accompanied by multiple homonymous files with a different extension. One of these is the **.var** file, already mentioned in Sec. 4, but beside that one, there can also be an **.inf** file and/or a **.prt** file.

7.1 .inf file

This is a file containing metadata for the database, such as multilingual aliases for the field names (because they can only be at most 10 characters long) or additional information describing the contents of the database. The file is encoded in the same *Starling* encoding as set forth in Sec. 5.

7.2 .prt file

This file contains settings for the *Starling* print functionality, such as column delineators etc. It is itself a *Starling* database.