

O101001010

1 0 1001001010100

Алексей Евгеньевич Хмельнов,

101010101010101010100

И.В. Бычков,

А.А. Михайлов

Институт динамики систем и теории управления СО РАН

{hmelnov,bychkov,mikhailov}@icc.ru

http://FlexT.icc.ru

Выбор: текстовые или бинарные 201010101010101010 файлы

Бинарные файлы:

0

- более компактные
- код для чтения/записи пишется проще
- возможен быстрый переход к любому месту файла (Seek) => чтение фрагмента

Текстовые файлы:

- более прозрачные
- легче редактируются
- не зависят от порядка байт платформы

Новое поколение выбирает текст (XML)

- сейчас часто жертвуют эффективностью ради упрощения контроля за правильностью структуры файла
- часто текстовые форматы основываются на синтаксисе XML, поскольку для этого синтаксиса существуют готовые библиотеки и инструментальные средства, облегчающие разработку алгоритмов чтения/записи информации,
- для XML можно автоматизировать контроль правильности структуры файла (используя XML-схемы)

Позиция авторов

- В небольших текстовых файлах удобно хранить информацию, которую необходимо часто редактировать (например, настройки программы). Иначе придётся писать специальные формы редактирования для таких файлов.
- Лучше избегать использования текстовых форматов для обмена большими объёмами информации.
- Для обеспечения прозрачности бинарных файлов и контроля правильности их создания используются спецификации форматов бинарных файлов

Цели разработки языков спецификации бинарных форматов данных

- Документирование форматов данных
- Проверка данных на соответствие спецификации с диагностикой ошибок
- Генерация кода чтения данных
- Метаинформация для библиотек обработки/передачи данных

Обзор средств спецификации бинарных форматов данных

_	Средство	Цель разработки	Вид	Описанные форматы	Битовые типы	Указатели	Варианты	Для человека	Дополнительно	
	BFF (Binary File Format Definition)	обратное проектирование формата AutoCAD DWG	Грамматика	DWG	_	?	?	?		
	DFDL (Data Format Description Language)	описания текстовых и бинарных данных, используемых в GRID-системах, Open Grid, IBM	XML, простые записи	таблицы в специальных обменных форматах	-	_	?	_		
	MFL (Message Format Language)	Описание обменных форматов в WebLogic Integration фирмы Oracle	XML, потоки данных	данные в специальных обменных форматах	-	-	+, необяза- тельные поля	-		
	NetPDL	описание сетевых протоколов, документирования форматов пакетов	XML, потоки данных	сетевые пакеты	+, прими- тивные	_	+	_)
0	BinPAC	описание сетевых протоколов, генерация кода чтения	Модифицирован-ные файлы Bison, потоки данных	сетевые пакеты	-	-	+	+		1
	EAST (Enhanced Ada SubseT)	разработан CCSDS для облегчения обмена информацией между космическими системами	Типы данных на базе Ada, потоки данных	последова-тельно передаваемые данные космические данные	+	-	+, выбор по внешнему значению	+	только выражения по одному значению	-
	HUDDL	описание форматов гидрографических данных, генерация кода	XML, потоки данных	гидрографические форматы	?	_	?	_		1
	Advanced Language Processing Technology Applied to Digital Records	интеграции разнородной информации, используемой для принятия военных решений, военный НИИ США, генерация кода чтения	атрибутивные грамматики	Форматы с последова- тельной записью	?	-	?	+		
	DataScript	генерации библиотек для чтения данных на языке Java	определения типов	Java CLASS (без маш. команд), DVI, ELF	+, прими-	—, есть вычисляе мые метки	+, по условиям на поля в union	+	с 2003 г. исходные тексты проекта не обновлялись	
	Miraplacid Binary DOM	реализация универсальной библиотеки для доступа к бинарным и текстовым данным в стиле DOM	определения типов	Описан 21 формат, часть из них - текстовые	-	+	-	+	Просмотра дерева структур, выделение текущей в дампе	
	Kaitai Struct	создание парсеров бинарных данных по спецификациям	~JSON, последовательность с необязательными и повторяющимися фрагментами	Описано 36 форматов	-	+	+	-	Появился Switch	
	Synalyze It! Hexinator	Просмотр бинарных данных	XML, грамматики, Python	80 форматов, в основном для macOS	+	+	+	_	Просмотра дерева структур, раскраска дампа	B
	FlexT	Просмотр бинарных данных	определения типов	100 форматов	+	+	+	+	U/ 111	

Обзор: выводы

- во многих предметных областях осознана необходимость использования спецификаций форматов данных
- иногда ставится задача не описания произвольных форматов, а лишь облегчения работы с некоторым их подмножеством
- иногда задача описания произвольных форматов ставится, но реализуются не все конструкции, необходимые для её решения
- часто для спецификаций выбирается представление XML, что затрудняет восприятие текста спецификации человеком
- в некоторых проектах явно ставится задача лёгкости восприятия текста человеком
- спецификацию бинарного формата рассматривают как грамматику, как структуру последовательности, или как набор типов данных
- многие языки спецификаций были созданы для облегчения написания кода для работы с данными
- разработчики очень редко изучают аналоги
- задача описания бинарных данных востребована

Виды спецификаций форматов 0101010101010101 данных

	10/100		21/1/0	
1	Специф	икация формата	данных	4
1 1 1	Спецификация	я интерпретации	Спецификация изменения	100
1	Идентификация типов данных	Графика, таблица, текст,		1

10101001010 101010

■■ Язык FlexТ

- декларативные спецификации форматов бинарных данных
- FlexT Flexible Types
- Основными конструкциями языка FlexT являются определения типов данных, сложность которых может существенно превышать сложность аналогичных определений в традиционных языках программирования

Дамп файла LSKAT.DBF и содержащаяся в нём информация

			110
LSKAT.DBF		DOS 241	Col 0 100%
00000000:	03 60 0B 1C 05 00 00 00	82 00 16 00 00 00 00 00	♥`♂∟ ⊕ B =
00000010:	00 00 00 00 00 00 00 00	00 00 00 00 00 00 00 00	
00000020:	4E 41 4D 45 00 69 B1 04	94 1D 00 43 38 64 0A 00	NAME i∰+Φ+ C8d© -
00000030:	12 00 D4 6F 34 69 B2 21	F1 52 0A 00 04 6F 44 69	‡ ⊧o4i‱!ëR⊠ ⊧oDi
00000040:	47 43 4F 4C 00 69 B1 04	94 1D 00 4E 38 64 0A 00	GCOL i ◆Φ+ N8d□
X 00000050:	02 00 04 6F 34 69 B2 21	F1 52 0A 00 04 6F 44 69	⊕ ⊧o4i‱!ëR⊡ ⊧oDi //
00000060:	54 4C 49 4E 45 00 B1 04	94 1D 00 4E 38 64 0A 00	TLINE • 🗣 👭 N8d 🗖 🐰
00000070:	01 00 04 6F 34 69 B2 21	F1 52 0A 00 04 6F 44 69	⊜ ⊧o4i∭!ëR⊡ ⊧oDi 💢
00000080:	0D 00 20 87 A5 AC AF AE	AB AE E2 AD AE 20 20 20	♪ Земполотно
00000090:	20 20 20 20 20 31 30 20	20 88 91 91 8E 20 20 20	10 ИССО 🗥
000000AO:	20 20 20 20 20 20 20 20 20	20 20 20 20 39 20 20 82	9 B
000000080 :	91 8F 20 20 20 20 20 20	20 20 20 20 20 20 20 20	CN ' "
7 000000000:	20 31 35 30 20 8F EO AE	E7 A8 A5 20 20 20 20 20	
7 00000000:	20 20 20 20 20 20 20 31	34 20 20 90 A5 AC AE AD	14 Pemon
000000E0:	E2 20 20 20 20 20 20 20	20 20 20 20 20 31 32 20	т 12
000000F0:	1A	20 20 20 20 20 31 32 20	· · · · · · · · · · · · · · · · · · ·
DUUUUU .	TVIUINI III	1/6//	

NAME	GCOL	TLINE				
Земполотно	10		[•]	Структура	D:ODES\L	SKAT.DBF —
NCCO	.9		l	_		
ВСП	15	0	<u>Имя</u>	Тип	Дпина	Десятки
Прочие	14		NAME	Сим	вол 18	
Ремонт	12		GCOL	Чис	по 2	0
			TLINE	Чис	по 1	0

Спецификация на FlexT формата DBF

```
TBinDate LastChangeDate
type
                                                    ulong RecCnt //Число записей
  TBinDate array[3] of Byte //Дата в
                   двоичном формате (ГГММДД)
                                                    PDataArray HdrLen //Длина заголовка
 TDBF3FldKind enum Char (
                                                                       в байтах
    fkChar='C', fkNumeric='N', fkLog='L',
                                                    Word RecLen
                                                                      //Плина записи в
                                                                       байтах
    fkDate='D', fkMemo='M'
                                                    (array[20] of Byte) Reserved
                                                  ends
  TDBF3FldDsc struc
    array[11] of Char, 0; Name //Имя -
                                                  TDBF3HdrWithFields struc
                   ASCIIZ строка
                                                    TDBF3Hdr H
    TDBF3FldKind hType
   ulong DataP //like Delphi Tag
                                                    array[(@.H.HdrLen-@.H:Size-1) div
                                                      32] of TDBF3FldDsc Fields
   Byte Len
                 //Длина поля в байтах
                                                  ends
   Byte DecNum
                 //Число знаков после
                   ТОЧКИ
                                                data
   Word MUsrRsrv1//Зарезерв. для
                                                  0x0000 TDBF3HdrWithFields Hdr
                   многопольз. систем
                                                type
    Byte WorkID //ID рабочей области
                                                  TFieldData array[Hdr.Fields[#].Len]of
   Word MUsrRsrv2//Зарезерв. для
                                                                  Char
                   многопольз. систем
                                                  TFieldsData array of TFieldData:
    Byte SetFldData//Используется командой
                                                                [@:Size=Hdr.H.RecLen-1]
                   SET FIELDS
                                                  TRecData struc
   array[8] of Byte Reserved
                                                    Char F
  ends
                                                    TFieldsData D
 PDataArray ^TDataArray near
                                                  ends
  TDBF3Hdr struc
                                                  TDataArray array[Hdr.H.RecCnt] of
   Byte Ver //$02-dBase II,$03-dBase III
                                                                       TRecData
             //$83-dBase III с Мето-полями
```

Результат разбора по спецификации файла LSKAT.DBF

```
00:Hdr: TDBF3HdrWithFields = (
 H: (Ver:dBase III { 03 }; LastChangeDate: (Y:96; M:11; D:28); RecCnt: 00000005; HdrLen: 0082;
    RecLen:0016; Reserved:
   (0:00,1:00,2:00,3:00,4:00,5:00,6:00,7:00,8:00,9:00,10:00,11:00,12:00,13:00,14:00,
    15:00,16:00,17:00,18:00,19:00));
  Fields: (
    0: (Name: 'NAME'; hType:fkChar{'C'}; DataP:000A6438; Len:12; DecNum:00;
       MUsrRsrv1:6FD4; WorkID:34; MUsrRsrv2:B269; SetFldData:21;
       Reserved: (0:F1,1:52,2:0A,3:00,4:D4,5:6F,6:44); InMDX:69),
    1: (Name: 'GCOL'; hType:fkNumeric{'N'}; DataP: 000A6438; Len: 02; DecNum: 00;
       MUsrRsrv1:6FD4; WorkID:34; MUsrRsrv2:B269; SetFldData:21;
       Reserved: (0:F1,1:52,2:0A,3:00,4:D4,5:6F,6:44); InMDX:69),
    2: (Name: 'TLINE'; hType:fkNumeric{'N'}; DataP:000A6438; Len:01; DecNum:00;
       MUsrRsrv1:6FD4; WorkID:34; MUsrRsrv2:B269; SetFldData:21;
       Reserved: (0:F1,1:52,2:0A,3:00,4:D4,5:6F,6:44); InMDX:69),
    3:0D))
    81:00 |...
82:Hdr.H.HdrLen^: TDataArray = (
                                   ',1:'10',2:' ')),
  0: (F: ''; D: (0: 'Земполотно
                                     ',1:' 9',2:' ')),
  1: (F: ' '; D: (0: 'MCCO
  2: (F: ''; D: (0: 'ΒCΠ
                                    ',1:'15',2:'0')),
  3: (F:''; D: (0:'Прочие
                                     ',1:'14',2:' ')),
  4: (F:' '; D: (0:' Ремонт
                                 ',1:'12',2:' ')))
   F0:1A |..
```

Особенности динамических типов данных

		110
Тип	Статический	Динамический
Размер	Фиксирован	Может зависеть от данных
Смещение составляющих	Фиксировано	Может зависеть от данных
Назначение	Тип переменных	Тип констант
Примеры на Паскале	<pre>TName=array[031] of Char; TId = string[32]</pre>	<pre>const CP: PChar = 'Hello'; var S: String; S := 'Hello';</pre>

Пример динамического типа данных: представление RTTI

пришер Н		
П	редставление RTTI	10101
1010010		1010101010
Определение типа	Функция для получения	Описание этого
данных из модуля	адреса поля TypeData	типа на FlexT
TypInfo в Delphi		
PTypeInfo = ^TTypeInfo;	<pre>function GetTypeData(TypeInfo:</pre>	TTypeInfo struc
TTypeInfo = record	PTypeInfo): PTypeData;	TTypeKind Kind
Kind: TTypeKind;	asm	Str Name
Name: ShortString;	XOR EDX, EDX	TTypeData TypeData
{TypeData: TTypeData}	MOV DL, [EAX].TTypeInfo.	ends
end;	Name.Byte[0]	
PTypeData=^TTypeData;	LEA EAX, [EAX].TTypeInfo.	
TTypeData= packed record	Name[EDX+1]	
case TTypeKind of	end;	

■ Особенности языка FlexT (1/2)

1. Составляющие переменного размера

```
FldNameTbl: array[@.numFields] of pchar;
```

2. Параметры и свойства типов

```
THdrData(Kind, Cnt) case THdrValType(@:Kind) of CHAR: array[@:Cnt] of Char
```

endc

3. Ссылки на свойства и составляющие в выражениях

```
V[i]- элемент массива, V.A- поле записи, V.0x15- случай варианта. V:Size- параметр или свойство. @ - переменная определяемого типа. V@- родительский тип, V:@- содержащая переменная, V:\#- порядковый номер переменной
```

■ Особенности языка FlexT (2/2)

- 4. Дополнительные блоки определений типов:
- Блок утверждений : [@.offset:Cnt=@.count]
- Блок условий корректности

:assert[@.Op>=0x80]

- Блок отображения :displ=(INT(2*@))
- Блок именования : **autoname**= (@.tag)
- Блок доп. свойства типа :let Val=(@.0)exc(@.1)
- 5. Вызовы типов: PSubSecData (Kind=@.subsec) lfc

Основные конструкторы типов языка FlexT(1/2)

		A
Целые числа	num -(6)	
Пустой тип	void	
Символьные типы	char, wchar	
Перечислимые типы		31
Собственно перечисление	enum byte (A=1,B,C)	01
Отображение выражением	<pre>enum byte = Tbl[@].Name;</pre>	11
Перечисление термов	<pre>enum TBit8 fields (R0: TReg @0.3,) of(rts(R0) = 000020_,)</pre>)
Множество	<pre>set 8 of (OLD_SECURED ^ 0x02,)</pre>	
Запись	struc Byte Len array[@.Len]of Char S ends	17

Основные конструкторы типов языка FlexT(2/2)

		4
Вариант	<pre>case @.Kind of vkByte: Byte</pre>	l
7	else ulong	
	endc	
Проверка	try	
1 1	FN: TFntNum	
	Op: TDVIOp	6
4	endt	0
Массив	array[@.Len]of str	0
	array of str ?@[0]=0!byte;	1
	<pre>array of str:[@:Size=@:ASize]</pre>) 7
Сырые данные	raw[@.S]	1
Сырые данные		
Выравнивание	align 16 at &0;	Ĭ
Beipubilibuilit		
Указатель	^TTable near=DWORD, ref=@:Base+@;	
Предварительное объявление	forward	
Машинные команды	codes of TOpPDP ?(@.Op>=	7
тугашинные команды	TWOpCode.br) and;	1
	<u> </u>	

Порядок байтов в слове и битовые типы данных 01010101010101010101010

type bit

TBitRec struc

num+(6)

num+(8)

num+ (2) C

ends

1	01 01	01	0	00	101010	0101
	1	0	-		. 5101	011

№ бита в байте Адрес байта № бита в слове

0		7	0		7	
	0			1		
0		7	8		15	

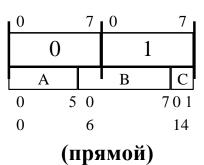
(прямой)

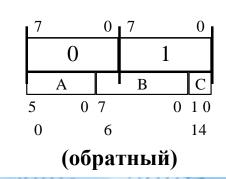
7		0	7		0
	0			1	
15		8	7		0

(обратный)

0101010100/ Порядок байтов в слове

№ бита в байте ● 0Адрес байта Поле TBitRec № бита в значении поля Смещение поля в записи





Побитовое размещение полей записи **TBitRec**

Типы данных ArcView Shape (общие для файлов SHP и SHX)

```
set byteorder rev
                                       TShpHeader struc
                                         longrev Magic
type
  longrev num-(4)
                                         longrev Rsrv0
 ulongrev num+(4)
                                         longrev Rsrv1
set byteorder norm
                                         longrev Rsrv2
include Float.rfi
                                         longrev Rsrv3
                                         longrev Rsrv4
                                         longrev FileLength
type
TShapeType enum long (
                                         long Ver
  Point=1, Arc=3,
                                         TShapeType ShapeType
  Polygon=5, MultiPoint=8)
                                         TBBox BBox
                                         TDouble Zmin
TBBox struc
                                         TDouble Zmax
  TDouble Xmin
                                         TDouble Mmin
  TDouble Ymin
                                         TDouble Mmax
  TDouble Xmax
                                       ends:assert[@.Magic=9994,
  TDouble Ymax
                                           @.FileLength*2=FileSize,
                                           @.Ver=1000]
ends
                                       data
TXPoint struc
                                       O TShpHeader Hdr
  TDouble X
  TDouble Y
                                       assert Hdr:assert;
ends
```

01

0

Формат файла SHP

```
include shp.rfi
                                      TShapeRecData(Len) struc
                                        TShapeType ST
                                        case @.ST of
type
                                          Point: TXPoint
TMultyPointData struc
                                          Arc, Polygon: TArcData
                                          MultiPoint:
  TBBox BBox
  long NumPoints
                                            TMultyPointData
  array[@.NumPoints] of TXPoint
                                        endc SD
    Points
                                        (raw[]at &@;) rest
ends
                                      ends:[@:Size=@:Len*2]:
                                         assert[@.ST=Hdr.ShapeType]
TArcData struc
                                      TShapeRec struc
  TBBox BBox
                                        longrev RecNo
  long NumParts
                                        ulongrev Len
  long NumPoints
                                        TShapeRecData (@.Len) Data
  array[@.NumParts] of long Parts
                                      ends
  array[@.NumPoints] of TXPoint
     Points
                                      data
ends
                                      Hdr:Size (array of TShapeRec:
                                         [@:Size=FileSize-&@]) Tbl
```

Float.rfi – представление действительных чисел

```
type bit
  TBit1 num+(1)
  TBit4 num+(4)
  TBit8 num+(8)
  TBit11 num+(11)
  TBit20 num+(20)
  TBit23 num+(23)
  TDoubleV struc
    array[13]of TBit4 M
   TBit11 P
    TBit1 S //Sign 1->Neg
  ends
  TSingle struc
    TBit23 M
    TBit8 P
    TBit1 S //Sign 1->Neg
  ends: displ=(FLOAT(@.M,@.P-0x7F,@.S))
type
  TDouble struc
    TDoubleV V
  ends: displ=(FLOAT(@.V.M,@.V.P-0x3FF,@.V.S))
```

Дамп небольшого Shape-файла

10700														216	170.		
err2.shp									DC)S			372	Co 1	0	100%	
00000000:	00 00	27 OA	00 (00 00	00	00	00	00	00	00	00	00	00	, o			
00000010:	00 00	00 00	00 (00 00	00	00	00	00	BA	E8	03	00	00		ll l	ш♥	
00000020:	05 00	00 00	EE I	EB CO	39	23	0A -	3B	40	CF	D5	56	EC	👲 юы	-9#0;Œ	≐ե≬ь ։	
00000030:	2F 1B	3F 40	C6 1	DC B5	84	70	10	3B	40	DF	BE	0E	9C	/+?@ <u>-</u> -	Д:►; €	引ル	17
00000040:	33 22	3F 40	00 (00 00	00	00	00	00	00	00	00	00	00	3"?@			1
00000050:	00 00	00 00	00 (00 00	00	00	00	00	00	00	00	00	00				h
@00000060:	00 00	00 00	00 (00 00	01	00	00	00	40	05	00	00	00			•	Ľ
00000070:	EE EB	CO 39		0A 3B	40	CF	D5	56	EC	2F	1B	3F	40	юы [∟] 9# <mark>⊙</mark> ;	(С≐ р∪ь.	/+? @	C
<u> </u>	8F E4	F2 1F		OF 3B	40	DF	BE	OE :	9C	33	22	3F	40	-Πφ€ ∀ π * έ;	Gard MP:	3"?@	1
00000090:	01 00	00 00		00 00	00	00	00	00	00	ΑE	69	DE	71	Θ 👲	1	oi q	U
_000000A0:	8A OE	3B 40		BE OE	9C	33	22	3F	40	8F	E4	F2	1F	кл;е⊌л			١
000000B0:	D2 OF	3B 40		FO 85	C9		21	3F	40	5E	DC	46	03	π#; CoEF			r
000000co:	78 OB	3B 40		DS 56	EC			3F	40	EE	EB	CO.	39	×₫;0≐린	b/ +?@	юы <u>-</u> 9	1
000000D0:	23 OA	3B 40		A4 DF	BE			3F	40	AΕ	69	DE	71	#0; CCA	₽ <i>∏</i> ∟?@	oi q	h
000000E0:	8A OE	3B 40		BE OE	9C	33		3F	40	00	00	00	02	K∏; e=4 f			П
000000F0:	00 00	00 40		00 00	00		47	E1		14	0E	3B	40		o Gc z'		Г
00000100:	E8 6A	2B F6		1D 3F	40			B5	84	7C	10	3B	40	_ <u>wj</u> +94+?		-	L
00000110:	E4 14	1D C9	E5 :	1F 3F	40	01	<u> 00 </u>	00	00	<u>05</u>	00	00	00	⊸φ¶₩ _{[P} κ♥?			
00000120:	00 00	00 00		B5 15	FB	CB	OE	3B	40	E4		1D	C9		N### @		L
00000130:	E5 1F	3F 40		DC B5	84			3B	40	91	OF.	78	36	×▼?@ -			L
00000140:	AB 1E	3F 40		6F 81	04		OF.	3B	40	E8	<u>6A</u>	2B	F6	л≜?@@o I			ı
00000150:	97 1D	3F 40		47 E1	<u>78</u>	14	OE .	3B	40	<u>C9</u>	E5	3F	A4	<u>4+?@oGc</u>	24M; F	IPX? μ	L
00000160:	DF 1E	3F 40		B5 15	FB	CB	0E	3B	40	E4	14	1D	C9	- - 4?@t= 8	Manti G	фЧ₩Г	
00000170:	E5 1F	3F 40												x₹?@			
														THE RESERVE TO SERVE THE PERSON OF THE PERSO			_

Результат разбора Shape-файла

```
000:Hdr: TShpHeader = (Magic:9994; Rsrv0:0; Rsrv1:0; Rsrv2:0; Rsrv3:0; Rsrv4:0;
  FileLength:186; Ver:1000; ShapeType:Polygon{5};
  BBox: (Xmin:27.0396; Ymin:31.1062; Xmax:27.0644; Ymax:31.1336); Zmin:0;
  Zmax:0; Mmin:0; Mmax:0)
064:Tbl: #Tf1 144 c10 =
  0: (RecNo:1; Len:00000040;
    Data: (ST:Polygon { 5 };
      SD: (BBox: (Xmin: 27.0396; Ymin: 31.1062; Xmax: 27.0618; Ymax: 31.1336);
        NumParts:1; NumPoints:5; Parts: (0:0);
        Points: (0: (X:27.0568; Y:31.1336),1: (X:27.0618; Y:31.1302),
           2: (\mathbf{X}: 27.0448; \ \mathbf{Y}: 31.1062), 3: (\mathbf{X}: 27.0396; \ \mathbf{Y}: 31.1096),
           4: (X:27.0568; Y:31.1336))); rest:[])),
  1: (RecNo:2; Len:00000040;
    Data: (ST:Polygon{5};
      SD: (BBox: (Xmin: 27.055; Ymin: 31.1156; Xmax: 27.0644; Ymax: 31.1246);
        NumParts:1; NumPoints:5; Parts: (0:0);
        Points: (0: (X:27.0578; Y:31.1246),1: (X:27.0644; Y:31.1198),
           2: (X:27.0616; Y:31.1156),3: (X:27.055; Y:31.1206),
           4:(X:27.0578; Y:31.1246))); rest:[])))
```

Источники информации о формате данных

- Описание на естественном языке неточность, неполнота, неоднозначность
- Исходные тексты программ чтения/записи информация о формате смешана с файловыми операциями, и другими командами обработки данных
- Примеры данных (reverse engineering) сначала спецификация отсутствует

Преимущества формальных спецификаций

- компактность, отсутствие не относящихся к способам хранения данных сведений, что облегчает их восприятие человеком;
- верифицируемость посредством их применения для анализа относящихся к формату корректных данных;
- возможность использования для локализации ошибок в сгенерированных файлах

Проверка соответствия данных спецификации

Результат разбора данных по описанию может отразить значение каждого бита исходного файла.

Разбор данных по спецификации может использоваться для:

- проверки данных на соответствие спецификации (как это делается для XML с использованием XML-схем)
- проверки спецификации на соответствие примеру данных

Спецификация, предназначенная для проверки, может допускать частично соответствующие полной спецификации данные.

Уточнение спецификации в ходе её записи на FlexT

- 1. Описать фрагмент формата в соответствии с тем, как понята спецификация
- 2. Применить спецификацию для разбора файла
- 3. При неудачном разборе данных исправить спецификацию, go to 2
- 4. Если не закончено, go to 1

Спецификация по примерам данных (reverse engineering)

- Является крайним случаем уточнения спецификации
- Лучше всего применять при наличии программыгенератора примеров данных с задаваемыми свойствами, например, компилятора
- Может применяться просто при наличии достаточного числа примеров

■■II Команды блока displ

```
Строковая константа
                                         вывод строки
'NL'[<Целочисл. выражение>]
                                         перевод строки с отступом по уровню вложенности
                                         конструкции
'NL'
                                         перевод строки без отступа
'SHOWARRAY' '('<Адресное выражение> ',' '(' <Список команд>,','* ')' ')'
                                         вывод массива
'COND' '(' <Логическое выражение> ',' '(' <Список команд>,','* ')' [ ',' '(' <Список команд>,','* ')' ] ')'
                                         условный оператор
'FLOAT' '(' M:<Целочисл. выражение> ',' Pwr:<Целочисл. выражение> ','
   Sign:<Целочисл.выражение> ')'
                                         вывод действительного числа, заданного мантиссой,
                                         порядком, и знаком
'INT' <Целочисл. выражение>
                                         вывод целого числа
'HEX' '(' V:<Целочисл. выражение> [',' Prec:<Целочисл. выражение>] ')'
                                         вывод целого числа в шестнадцатеричном
                                         представлении
'CHR' <Целочисл. выражение>
                                         вывод символа, заданного кодом ASCII
'WCHR' <Целочисл. выражение>
                                         вывод символа, заданного кодом Unicode
<Адресное выражение>
                                         вывод элемента данных
'ADDRNAME' <Целочисл. выражение>
                                         вывод имени переменной с указанным адресом
```

Перевод из Shape в MIF при помощи запроса на отображение Shp2Mif.displ

```
Hdr ('Version 300',NL,'Charset "WindowsCyrillic"',NL,'Delimiter ","',NL,'CoordSys NonEarth Units "km" Bounds
    (',@.BBox.Xmin,', ',@.BBox.Ymin,') ','(',@.BBox.Xmax,', ',@.BBox.Ymax,')',NL,
    'Columns 1', NL,' ID Integer', NL, 'Data', NL)
Tbl (ShowArray(@,(
  COND(@.Data.ST=TShapeType.Polygon,(
   'REGION',INT(@.Data.SD.Polygon.NumParts),NL ,
   ShowArray(@.Data.SD.Polygon.Points,
    (INT(@.T:N),NL
     ShowArray(@.T,(' ',@.X,' ',@.Y,NL ))
    ),((
  COND(@.Data.ST=TShapeType.Arc,(
   COND(@.Data.SD.Arc.NumParts=1, (
     'PLINE', INT(@.Data.SD.Arc.NumPoints), NL ,
     ShowArray(@.Data.SD.Arc.Points[0].T,(' ',@.X,' ',@.Y,NL ))
    ),('PLINE MULTIPLE ',INT(@.Data.SD.Arc.NumParts),NL ,
     ShowArray(@.Data.SD.Arc.Points, (INT(@.T:N),NL , ShowArray(@.T,(' ',@.X,' ',@.Y,NL )))))),
    (COND(@.Data.ST=TShapeType.Point,(
     'POINT',@.Data.SD.Point.X,'',@.Data.SD.Point.Y,NL
   ),(
     'Other', NL
   )))))))))
```

Использование спецификаций для анализа машинного кода

- Машинные команды используются в:
 - исполняемых и объектных файлах
 - коде виртуальных машин
 - сценариях в: ТТГ, файлах инсталляторов
- исследование файлов многих форматов будет неполным без отображения содержащихся там машинных инструкций
- представление машинной команды состоит из постоянной части, задающей вид команды, и переменной части, задающей её аргументы, если таковые имеются
- Постоянная часть для байт-кодов виртуальных машин, как правило, является перечислимым типом данных
- Для машинных команд реальных процессоров используется более плотное кодирование, поэтому в постоянной части могут выделяться битовые поля, содержащие часть аргументов, например, номера регистров
- Будем называть сочетания таких битовых полей с кодом инструкции термами
- Кодовый тип данных позволяет использовать статический дизассемблер для любых спецификаций кодирования машинных команд
- Использование спецификаций кодирования машинных команд позволяет применять все методы анализа, отображения и поиска данных для анализа кода

Тип данных перечисление термов

- Задаётся список всех битовых полей, применяемых для кодирования команд
- Каждое поле характеризуется битовыми смещением, размером и типом
- Затем описываются сами термы с указанием того, какие из этих полей они используют
- Вместе с термом указывается его база число, задающее значения не входящих в поля битов
- В записи базы терма допускается использовать подчерки для обозначения переменной части битов, входящих в состав полей (такая маска проверяется на соответствие упоминаемым в терме полям)
- В реализации типа данных «перечисление термов» используется дерево решений, позволяющее быстро найти терм по его коду
- В выражениях можно обращаться к любым полям перечисления термов, как к полям записи, однако, если у выбранного терма это поле не используется, то такое обращение вызывает ошибку вычисления этой части выражения
- Реализована разновидность типа данных «вариант» с выбором содержимого по значению перечисления термов для описания переменной части команд. Для выбора варианта в этом случае также строится дерево решений

В І І Кодовый тип данных

- Реализуется на базе массива со стоп-условием и записывается аналогичным образом
- Используется идентификатор codes вместо array
- Стоп-условие используется для выделения команд, завершающих последовательность, например, команд безусловного перехода или выхода из процедуры
- В описаниях команд, вызывающих переход, должны использоваться указатели на определяемый кодовый тип
- Может быть определено вычисляемое свойство isCall, которое в этом случае используется дизассемблером для различения команд перехода и вызовов процедур

Фрагменты спецификации кодирования команд процессора ARM

```
type bit
TBit num+(1)
TBit2 num+(2)
TsBit24 num-(24)
TRN enum TBit4 (R0,R1,... R12,SP,LR,PC)
TShiftOp enum TBit2 (SHL,SHR,ASR,ROR)
TBrOfs24 ^TOpARMSeq NIL- =TsBit24 REF=
 (&(@:@)+4+@*4);
type
TARMOpCode enum TBit32 fields (
 Cond: TCond @28.4, //Op.execution condition
 Rn: TRN @16.4, //Source register
 Rd: TRN @12.4, //Destination register
 Shift: TShiftOp @5.2,
 SM: TBit @20.1
) of (
 //Data processing rotate right with extend
  DPRRX(cond,ArOp,S,Rn,Rd,Rm) =
   0b 000 I I 0000I0110 ,
```

```
//Move status register to register
  MSR(cond,R,Rd) =
   0b 0001I0 001111I 0000I00000000,
 //Branch and Branch with link
  Branch(cond,BL,BrOfs) =
   0b 101 I I
 //Software interrupts
 SWI (cond, SWIOfs) =
   0b 11111 I I
):let IsStop=(((@.BL=0) and (@.cond>=
  TCond.AL)) exc 0
);:let isCall=((@.BL=1) exc 0);
TOpARMSeq codes of TARMOpCode ?@:isStop;:
  displ=('(', ShowArray(@,(NL,HEX(&@,4),':
  @)),NL,')')
```

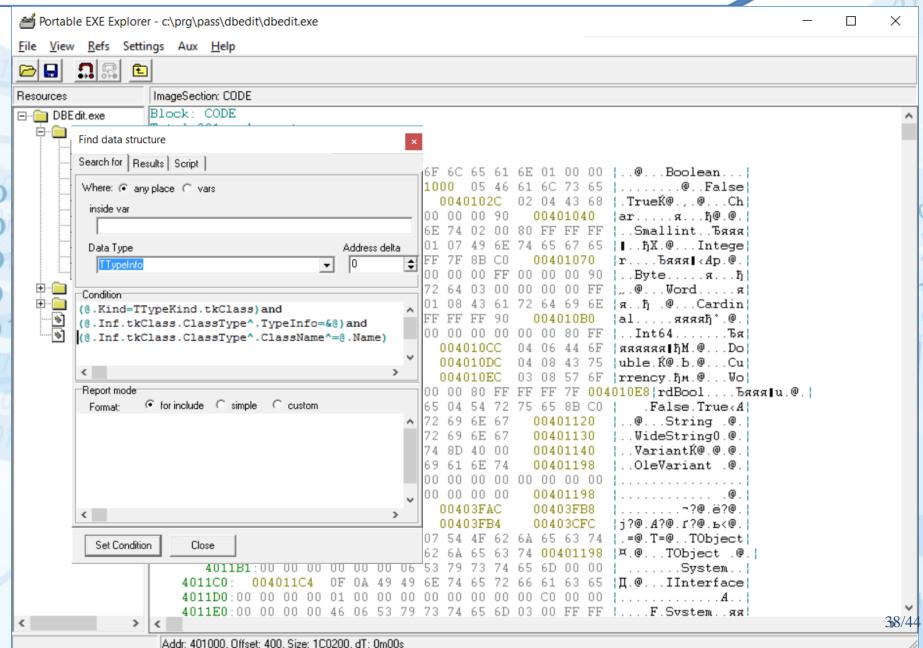
Разбор файла класса Java

```
Binary file explorer - C:\PRG\Java\ExprTreeOk\ExprTree.class
File Show Refs Aux Help
descrNdx: (LExprItem;)Ljava/awt/Component; (0077); attr_count:0001; attributes:
      (0: (name_index: 'Code' {005F}; len: 00000037;
          info:(max_stack:0001; max_locals:0002; CodeLen:0000000E; code:(
              0000: aload 1
              0001: invokevirtual 'isLeaf' {0048}
              0004: ifne 000C
              0007: aload 0
              0008: getfield 'edInPlaceOp' {0005}
              000B: areturn
              000C: aconst_null
              000D: areturn
            ); ExcTblLen:0000; AttrCnt:0002;
            attributes: (0:(name_index:'LineNumberTable'{0060}; len:0000000E;
                info: (Len: 0003; LinNumTbl: (0: (pc: 0000; 1: 00C7), 1: (pc: 0007; 1: 00C8),
                    2:(pc:000C; 1:00C9)))),1:(name_index:'StackMapTable'{0062}; len:00000003; info:
   000:00 01 0c |...|
                )))))).
   10: (access_flags: [ACC_PUBLIC]; nameNdx: 'StartEditExprItem' {0078};
     descrNdx:'(LExprItem;)Z'{0079}; attr_count:0001;
     attributes: (
       0: (name_index: 'Code' {005F}; len: 000000A3;
          info:(max_stack:0002; max_locals:0002; CodeLen:0000005A; code:(
              0000: aload 1
              0001: invokevirtual 'isLeaf' {0048}
              0004: ifeq 0035
              0007: aload 1
              0008: getfield 'h'{0040}
              000B: ifnonnull 0014
              000E: aload 1
              000F: ldc (Tag:C_String{08}; info:'?'{00C3}){49}
              0011: putfield 'h'{0040}
              0014: aload 0
              0015: getfield 'edVal'{002D}
              0018: aload_1
              0019: getfield 'h'{0040}
              001C: checkcast 'java/lang/String' {0041}
              001F: invokevirtual 'EditText' {004A}
              0022: ifne 0027
```

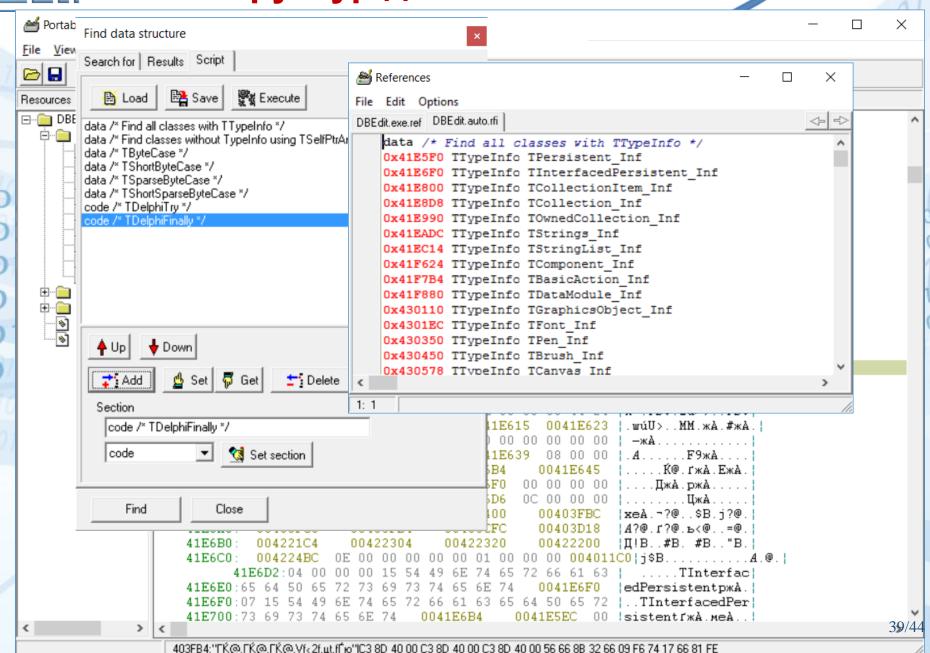
Просмотр результатов разбора и дампа

```
Binary file explorer - C:\PRG\Java\ExprTreeOk\ExprTree.class
File Show Refs Aux Help
0008: return
              ); ExcTblLen:0000; AttrCnt:0001;
              attributes: (0:(name_index:'LineNumberTable'{0060}; len:0000000A;
                   info: (Len: 0002; LinNumTbl: (0: (pc: 0000; 1: 00C2), 1: (pc: 0008; 1: 00C3)))))))),
    9:(access_flags:[ACC_PUBLIC]; nameNdx:'getInPlaceEditor'{0076};
       descrNdx:'(LExprItem;)Ljava/awt/Component;'{0077}; attr_count:0001; attributes:
        (0: (name_index: 'Code' {005F}; len: 00000037;
            info:(max_stack:0001; max_locals:0002; CodeLen:0000000E; code:(
                 0000: aload_1
                0001: invokevirtual 'isLeaf' {0048}
                 0004: ifne 000C
                 0007: aload_0
                0008: getfield 'edInPlaceOp' {0005}
                 000B: areturn
                 000C: aconst_null
                000D: areturn
              ); ExcTblLen:0000; AttrCnt:0002;
              attributes: (0:(name_index:'LineNumberTable'{0060}; len:0000000E;
                   info: (Len: 0003; LinNumTbl: (0: (pc: 0000; 1: 00C7), 1: (pc: 0007; 1: 00C8),
                        2:(pc:000C; 1:00C9)))),1:(name_index:'StackMapTable'{0062}; len:00000003; info:
    000:00 01 0c |...|
                   )))))),
    10: (access_flags: [ACC_PUBLIC]; nameNdx: 'StartEditExprItem' {0078};
         72 01 01 07 00 73 07 00 6E 00 00 FC 00 21 01 FC r...... s.. n b !..b
|00000F90:|00 14 07 00|6E F9 00 0B|00 01 00 74|00 75 00 01| ..... nw ... t u .
00000FA0: 00 5F 00 00 00 25 00 02 00 02 00 00 00 09 2A B4
|00000FB0:|00 14 2B B6|00 47 B1 00|00 00 01 00|60 00 00 00| ...+¶ G±
00000FC0: 0A 00 02 00 00 00 C2 00 08 00 C3 00 01 00 76 00 ... ...
|00000FD0:|77 00 01 00|5F 00 00 00|37 00 01 00|02 00 00 00|w ... _
|00000FE0:|0E||2B||B6 00||48 9A 00 08||2A||B4 00 05||B0 01 B0 00||...<mark>+¶</mark> H.b...*r ...*
00000FF0: 00 00 02 00 60 00 00 0E 00 03 00 00 00 C7 00
|00001000:|07 00 C8 00|0C 00 C9 00|62 00 00 00|03 00 01 0C|... И ... Й Ь
| 00001010: | 00 01 00 78 00 79 00 01 00 5F 00 00 00 A3 00 02 | ... x y ...
|00001020:|00 02 00 00 00 5A 2B B6 00 48 99 00 31 2B B4 00 ... Z+¶ H<sup>™</sup> 1+r
|00001030:|40 C7 00 09|2B 12 49 B5|00 40 2A B4|00 2D 2B B4|@3 ...+...Iu @*f --+r
|00001040:|00 40 C0 00 41 B6 00 4A 9A 00 05 03 AC 2B 2A B4| @A A¶ Jљ .....¬+*r
@s &+r @3 راع 00 40 E6 00 40 B6 00 40 B5 00 40 A7 00 26 28 B4 00 40 C7 🖣 لاي @s
. $£ با±1 ...+*r ...¶ لب# 1 00001060: 00 11 كا £1 28 كا £4 00 05 B6 00 4C B5 00 40 A7 00 11 ...+*r
00001070: 2A B4 00 05 2B B4 00 40 C0 00 41 B6 00 4D 04 AC *r ...+r @A A¶ M...-
|00001080:|00 00 00 02|00 60 00 00|00 2A 00 0A|00 00 00 CE|
00001090 00 07 00 CE 00 0E 00 D0 00 14 00 D1 00 25 00 D2
```

■ Поиск структур данных 1/4

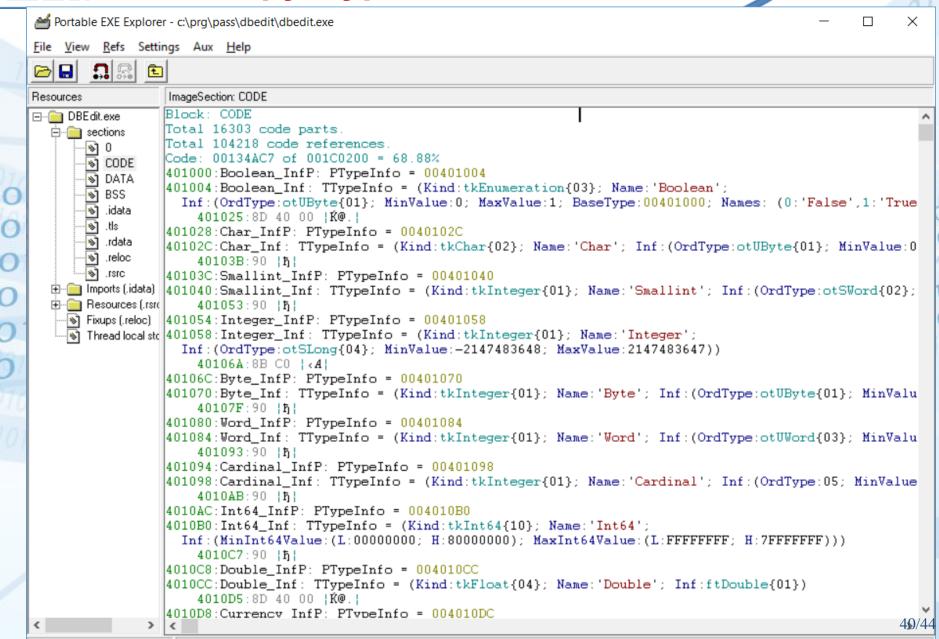


Поиск структур данных 2/4



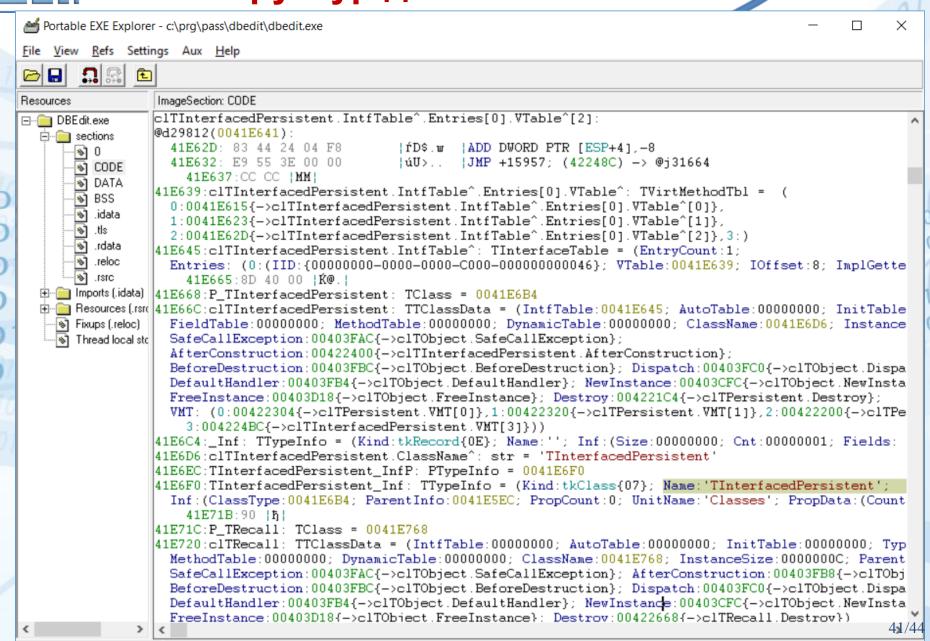
Поиск структур данных 3/4

Addr: 401000, Offset: 400, Size: 1C0200, dT: 0m00s



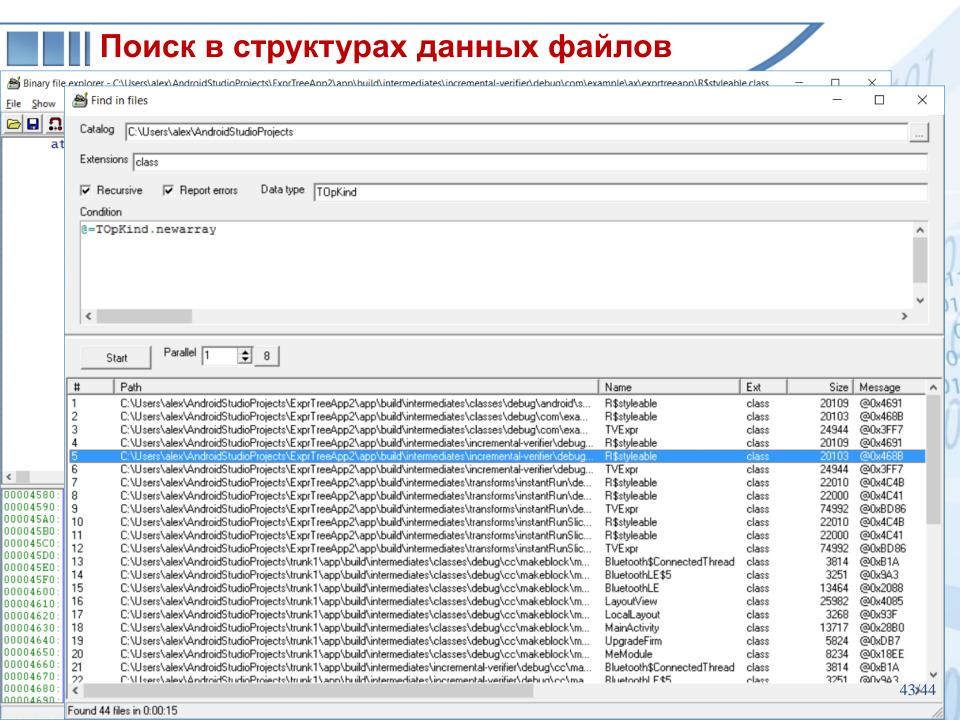


Addr: 401000, Offset: 400, Size: 1C0200, dT: 0m00s



Пример применения поиска структур: отчёт по анализу кода интерпретатора => описание типов данных для конвертера

			1/11
TPushHParentOrZ try		DocWin Obj	Project Doc
P: TPushHParent		Owner 11	VisGUIWidth 0
Z: TPushZ		Open 18	Flags 3
Fail: TFail		X 0	Doc 12
endt		Y 0	ActiveDoc 11
		W O	Buttons 12
TRegClassCall struc	11/	Н 0	Dependencies 12
TPushBOrD pushSize	110	1.	Available 12
TSkipToPushConst Skip1		Doc Obj	Scripts 11
TPushFieldTblOrZ pushF		Name 4	Prefs 11
TPushMethodTblOrZ pushM TPushHParentOrZ pushP		Creator 5	GUI 12
byte bPushClName		·	
PPChar ClName	10	CreationDate 5	WorkDir 11
byte bCallReg		Comments 5	WinX 0
long CallRegAddr	entra di Parisi	Open 5	WinY 0
byte bAdd	de phi manelon	Update 5	WinW O
word wAddArg	(COSTILI	Close 5	WinH O
TSkipToA3 Other	UI	GUIName 4	StartUp 5
byte bMov		GroupGUIName 4	ShutDown 5
ulong hClAddr		Win 11	SerialNumber 5
<pre>ends:assert[@.pushSize.*:assert, @.pushF.*:assert,@.pushM.*:assert,</pre>		CSMgr 11	SelColor 11
@.pushP.*:assert, @.bPushClName=0x68, IsFixup @.ClName,		DocExts 12	FN 11
@.bCallReg=0xE8, @.CallRegAddr+&@.bCallReg+5=0x50ED90,		ObjectTag 11	Windows 13
@.bAdd=0x83,@.wAddArg=0x14C4, @.bMov=0xA3]	111	IsPrinting 18	Modified 18
	14	,	Codepage 2
fldlist.txt:	04	ExtObj Obj	GUINames 12
@.ClName^,	JI		LegendLineFlatness 2
ShowArray(@.pushF.F.FT^, (NL,' ',@.Name^,' ',Int(@.hDT))),NL	4	Two-Dag Lob-	-
27/27	0.	InvRec Obj	TocWidthScale 6
AVC1.txt	UI	Object 11	ODBVersion 2
@.ClName^,' ',AddrName(@.pushP.P.CallPar^.Addr),		Client 11	ODBGUIs 12
ShowArray(@.pushF.F.FT^, (NL,@.Name^,' ',Int(@.hDT))),NL			42

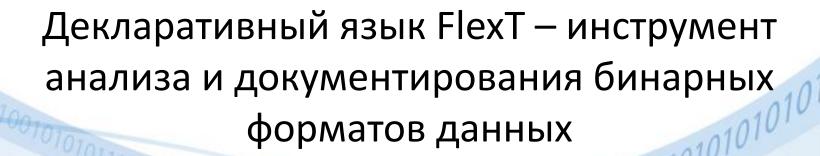


Ограничения текущей версии FlexТ

- спецификации на FlexT являются чисто декларативными, поэтому они не предназначены для описания таких форматов, при интерпретации которых необходимо строить сложные вспомогательные структуры данных, непосредственно не представленные в файле
- взаимосвязь нескольких файлов можно описать только через генерацию REF-файлов

Заключение

- при наличии спецификации на FlexT бинарный файл становится не менее прозрачным, чем текстовый
- спецификации форматов бинарных данных предоставляют альтернативу использованию текстовых форматов для упрощения контроля за правильностью работы алгоритмов чтения/записи
- возможность автоматического контроля правильности структуры файла делает спецификации на языке FlexT альтернативой использованию форматов на основе XML



O101001010

1 0 1001001010100

Алексей Евгеньевич Хмельнов,

101010101010101010100

И.В. Бычков,

А.А. Михайлов

Институт динамики систем и теории управления СО РАН

{hmelnov,bychkov,mikhailov}@icc.ru

http://FlexT.icc.ru

Спецификация формата BCRL

```
const
  BitCntBits = 4;
  LogMaxVal = 3;
  MinIHdrSz=BitCntBits+LogMaxVal+1;
type bit
TBit num+(1)
TBitCnt num+(BitCntBits):displ=(int(@))
TLPart num+(LogMaxVal)
TILen (Cur) struc pas
  Cont: TBit
  Part: TLPart.
  Rest: case @.Cont of
   1: TILen (@@:Cur shl LogMaxVal+@@.Part)
  endc
ends:let Val=(@.Rest.1:Val exc(@:Cur shl
   LogMaxVal+@.Part));:displ=(Int(@:val))
```

```
TInterval (Lim) struc pas
  BitCnt: TBitCnt
  Len: TILen(0)
  D: array[@.Len:val+1]of
       num-(@@.BitCnt)
ends
TBCRLData (PkSz, Sz) struc pas
  D: array of TInterval (&@+
       @@:PkSz)?&@>=@:Lim!void;
  Al: align 8;
ends: [@:Size=@:PkSz*8]
```



Спецификация кодировки UTF-8

```
type
T1ByteChar char():assert[(@>0)and(@<0x80)]
T2ByteChar array[2]of Byte:assert[((@[0]and 0xE0)=0xC0)
  and ((@[1]) and 0xC0) = 0x80) : displ=(HEX)((@[0]) and 0x1F) shl 6+(@[1]) and
   0x3F))
T3ByteChar array[3]of Byte:assert[((@[0]and 0xF0)=0xE0)
  and ((@[1] and 0xC0) = 0x80) and ((@[2] and
   0xC0) = 0x80) : displ=(HEX((@[0]and 0x0F)shl 12
  +(@[1]and 0x3F)shl 6+(@[2]and 0x3F)))
TUtf8Char try
  B1: T1ByteChar
  B2: T2ByteChar
  B3: T3ByteChar
  X: Byte
endt
TUtf8Array(Sz) array of TUtf8Char: [@:Size=@:Sz]:displ=(
  '"', ShowArray (@, (WCHR (@.B1 exc(
    ((@.B2[0]and 0x1F)shl 6+(@.B2[1]and 0x3F)) exc(
    ((@.B3[0]) and 0x0F) shl 12+(@.B3[1]) and 0x3F) shl 6+(@.B3[2]) and 0x3F))
   exc (0.X)))),'"')
```