# Type Encoding/Decoding rules: from C to Protocol Buffer

Shen Liu

(Version 1.5, last modified 07/05/2016)

**Abstract**

Type conversion between C and protocol buffer is an important issue in our project. When the separation is done and the RPC tool begins working, we must automatically restore all the arguments for each RPC function in the receiver process, which means the function parameter types need to be transmitted between two processes fully and exactly. Unfortunately, the type system that protocol buffer supports is quite weak. To make our project automatically run in the end, we have to design a type conversion protocol to let protocol buffer automatically convert some advanced C types(e.g. pointer) into protocol buffer types.

## 1 Background

We use gRPC, which is fully based on google protocol buffer, to deal with RPC issues in our project. In gRPC, a C type must be packed to protocol buffer "`message`" type in a .proto file(IDL file) for further transmission. For example, if you have a C function `int foo(int x)` which needs to be called remotely, then in your .proto file, the argument type `int` can be packed in protocol buffer as follows:

```
1  message M{                        message M{
2    int64 x=1;        or             int32 x=1
3  }                                  }
```

Next, protocol buffer will automatically generate a group of read/write APIs for each message. Here is an API for x's value assignment:

```
1   void set_x(::google::protobuf::int32 value);
```

and an API for getting the value of x:

```
1   inline ::google::protobuf::int32 M::x() const {
2     return x_;
3   }
```

Here is a more complex C-protobuf type conversion sample:

```
1  typedef struct{                    message Circle{
2    int x;                             message Point{
3    int y;                               int64 x=1;
4  }Point;                                int64 y=2;
5                        --->           }
6  typedef struct{                      double radius
       =1;
7    Point center;                    }
8    double radius;
9  }Circle;
```

Our project can automatically finish this conversion for all scalar types and simple composite types as "Circle". However, when parameter types become more and more complex, especially for those structures with multi-level pointers, generating a correct .proto file automatically as before will be a real challenge. To achieve this goal, a possible way is designing a type-conversion protocol to make our project work more intelligently. Simply speaking, for any C type input, first we use such a protocol to convert it into an integer array(encoding), and then construct the "message" type in .proto. On the receiver side, we do array parsing to restore the original C types(decoding) instead of parsing the complex .proto file.

## 2    Type system and encoding/decoding rules

In this draft we only use a small subset of C type system to show how the encoding/decoding idea works. Here is how our toy type system looks like:

```
1  Type t := int | t* | struct {t1; t2; ... ;tn}
2            | tname S
```

Any pair of form (`type`,`value`) based on this type system will be encoded as an array of bytes(`bytes[] lst`), and the first byte(see table 1) in this array denotes what type this array corresponds to.

| lst[0] | type |
|--------|------|
| 0 | int |
| 1 | pointer |
| 2 | struct $\{t_1; t_2; ...; t_n\}$ |
| 3 | tname S |

Table 1: type mapping rules

As we can see from Table 1, any encoding/decoding operation related to type `tname S` requires knowing the associative type `struct{...}` of `tname S`. In our framework, we use a name-type mapping table(e.g. table 2) to map each name string which represents a struct to its corresponding `struct` type.

2

| name | struct |
|------|--------|
| S1 | struct {int;int;} |
| S2 | struct {int;int*;} |
| S3 | struct {int;int*; S1*;} |
| ... | ... |

Table 2: A name–type mapping table example

Besides, in each round for encoding, we also use an auxiliary table called pointer table to record each pointer value that ever appeared. By doing this we can identify some complex function arguments(e.g. circular linked list).

Once we have such auxiliary tables, we can easily construct the encoding/decoding rules for our type system as follows:

```
1  Basic value conversion functions:
2
3  intToBytes(int): convert an integer to a byte string.
4
5  symbolToBytes(S): convert a symbol S to a byte string.
6
7  ptrToBytes(int): convert a pointer address(int) to a
       byte string.
8
9  bytesToInt(bytes[]): convert a byte string to an
       integer.
10
11 bytesToSymbol(bytes[]): convert a byte string to a
       symbol.
12
13 bytesToPtr(bytes[]) convert a byte string to a
       hexadecimal integer.
14
15 dereference(int): return the value that a pointer
       points to.
16
17 getTypeFromTable(S): look up S in the mapping table
       and return its
18                    associative type.
19
20 Definition encode: (type,value) (t,v) -> bytes[]
21   match t with
22   | int     => 0::intToBytes(v)
23   | t*      => 1::ptrToBytes(v)::encode(t,dereference(v
         ))
```

```
24    | struct {(t1,v1);...;(tn,vn)}
25            => 2::intToBytes(n)::encode(t1,v1)::...::
               encode(tn,vn)
26    | tname S
27            => 3::symbolToBytes(S)::encode(
               getTypeFromTable(S),v)
28    end.
29
30
31 Definition decode bytes[] lst =
32    match lst[0] with
33      | 0 => ((int, bytesToInt(lst[1...4])), lst+5)
34
35      | 1 => let ((t1,dereference(v)), l1) = decode (lst
            +5) in ((t1*, v), l1)
36            //v = bytesToPtr(lst[1...4])
37
38      | 2 => let n = bytesToInt(lst[1]) in
39            let ((t1,v1),l1) = decode (lst+5) in
40          let ((t2,v2),l2) = decode l1 in
41            ...
42            let ((tn,vn),ln) = decode l_{n-1} in
43            (struct {(t1,v1);(t2,v2);...;(tn,vn)}, ln)
44
45      | 3 => let S = bytesToSymbol(lst[5...5+length(S)
            -1]) in
46                              /*length(S) =
                                 bytesToInt(lst
                                 [1...4])*/
47            decode(lst+offset)  /* offset = 1+4+length(S
                )*/
48    end
```

Now consider a circular linked list example:

```
1  typedef struct Node{
2    int val;
3    Node* next;
4  }Node_t;
5
6  Node_t *head = (Node_t*) malloc(sizeof(Node_t));
        //head: 0x0004
7  Node_t *tail = (Node_t*) malloc(sizeof(Node_t));
        //tail: 0x0008
8
9  head->val = 10;
10 head->next = tail;
```

```
11
12  tail->val = 20;
13  tail->next = head; // circular linked list
```

Assume that we want to send this circular linked list from sender to receiver, then the encode/decode process is as follow:

```
1  encode(Node_t *head, 0x0004)
2
3  = 1::ptrToBytes(0x0004)
4      ::encode(Node_t, dereference(0x0004))
5      /* dereference(0x0004) = {10,0x0008} */
6
7      (pointer table: {0x0004})
8
9  = 1::ptrToBytes(0x0004)
10      ::3::symbolToBytes(Node_t)
11          ::encode(getTypeFromTable(Node_t),{10,0x0008})
12
13      (pointer table: {0x0004})
14
15  = 1::ptrToBytes(0x0004)
16      ::3::symbolToBytes(Node_t)
17      ::encode(struct Node{int, struct Node*},{10,0x0008
            })
18
19      (pointer table: {0x0004})
20
21  = 1::ptrToBytes(0x0004)
22      ::3::symbolToBytes(Node_t)
23          ::2::intToBytes(2) /*two fields*/
24              ::encode(int,10)
25              ::encode(Node_t*,0x0008)
26
27      (pointer table: {0x0004})
28
29
30  = 1::ptrToBytes(0x0004)
31      ::3::symbolToBytes(Node_t)
32          ::2::intToBytes(2) /*two fields*/
33              ::0::intToBytes(10)
34              ::1::ptrToBytes(0x0008)
35                  ::3::symbolToBytes(Node_t)
36                      ::encode(getTypeFromTable(Node_t),
                          dereference(0x0008))
37
38      (pointer table: {0x0004, 0x0008})
```

```
39
40
41  = 1::ptrToBytes(0x0004)
42      ::3::symbolToBytes(Node_t)
43        ::2::intToBytes(2) /*two fields*/
44          ::0::intToBytes(10)
45          ::1::ptrToBytes(0x0008)
46            ::3::symbolToBytes(Node_t)
47                ::encode(struct Node{int, struct Node
                       *}, {20,0x0004})
48
49    (pointer table: {0x0004, 0x0008})
50
51
52
53
54  = 1::ptrToBytes(0x0004)
55      ::3::symbolToBytes(Node_t)
56        ::2::intToBytes(2) /*two fields*/
57          ::0::intToBytes(10)
58          ::1::ptrToBytes(0x0008)
59            ::3::symbolToBytes(Node_t)
60                ::2::intToBytes(2) /*two fields*/
61                  ::0::intToBytes(20)
62                  ::1::ptrToBytes(0x0004)
63                  ...
64                  (0x0004 is in pointer table already,
                       stop here)
65
66      (pointer table: {0x0004, 0x0008})
```

0x0004 appears again, which means there must be a circle, to remember all pointer values we need an extra data structure for pointer storage and comparison.

The decoding process for the generated `bytes[] lst` can be illustrated as follows:

```
1   decode(bytes[] lst)
2   = decode(1::ptrToBytes(0x0004)
3              ::3::symbolToBytes(Node_t)
4                ::2::intToBytes(2)
5                  ::0::intToBytes(10)
6                  ::1::ptrToBytes(0x0008)
7                    ::3::symbolToBytes(Node_t)
8                      ::2::intToBytes(2)
9                        ::0::intToBytes(20)
10                       ::1::ptrToBytes(0x0004))
```

6

```
11
12  = decode(3::symbolToBytes(Node_t)
13           ::2::intToBytes(2)
14              ::0::intToBytes(10)
15              ::1::ptrToBytes(0x0008)
16                 ::3::symbolToBytes(Node_t)
17                    ::2::intToBytes(2)
18                       ::0::intToBytes(20)
19                       ::1::ptrToBytes(0x0004)) in ((t1
                               *, 0x0004), null)
20
21  = decode(2::intToBytes(2)
22           ::0::intToBytes(10)
23           ::1::ptrToBytes(0x0008)
24              ::3::symbolToBytes(Node_t)
25                 ::2::intToBytes(2) /*two fields*/
26                    ::0::intToBytes(20)
27                    ::1::ptrToBytes(0x0004)) in ((
                            Node_t*, 0x0004), null)
```

Once we have the name Node_t, we can directly look up and retrieve its associative struct type in the mapping table, and then restore a new list on the receiver side like:

```
1  Node_t* head = (Node_t*) malloc(sizeof(Node_t));
```

During the left decoding process, we set up a table(see table 3) which records the values on both sender and receiver sides for each pointer, to help us conveniently restore the sender side point-to relationships in the receiver side.

| pointer | value in sender | value in receiver |
|---------|-----------------|-------------------|
| head | 0x0004 | 0x0012 |
| tail | 0x0008 | 0x0016 |
| head→next | 0x0008 | 0x0016 |
| tail→next | 0x0004 | 0x0012 |
| ... | ... | ... |

Table 3: Pointer values in both sender and receiver

For example, assume in the receiver side pointer head equals to 0x0012, and its sender counterpart equals to 0x0004, then we have an entry like head–0x0004–0x0012 in our table. After we finished the restoration, we can also use this table to check whether the old pointer-to relationships are maintained correctly.

Let's continue decoding the bytes above, we have:

```
1  = decode(2::intToBytes(2)
```

```
2                ::0:: intToBytes (10)
3                ::1:: ptrToBytes (0 x0008)
4                  ::3:: symbolToBytes ( Node_t )
5                    ::2:: intToBytes (2)
6                        ::0:: intToBytes (20)
7                        ::1:: ptrToBytes (0 x0004 ))
8
9  = decode (3:: symbolToBytes ( Node_t )
10             ::2:: intToBytes (2)
11                ::0:: intToBytes (20)
12                ::1:: ptrToBytes (0 x0004 ))
13    head ->val = 10;
14    head ->next = ( Node_t *) malloc ( sizeof ( Node_t )); //
          assume head ->next = 0 x0016
15    (" Node_t " can be directly retrieved from the name -
          type mapping table )
16
17    Pointer table :
18    Pointer        sender    receiver
19    head           0 x0004     0 x0012
20    head ->next    0 x0008     0 x0016
21
22  = decode ( empty )
23
24    head ->next ->val = 20;
25    head ->next ->next = ( Node_t *) malloc ( sizeof ( Node_t ))
          ;// assume 0 x0020
26
27    Pointer table :
28    Pointer           sender    receiver
29    head              0 x0004     0 x0012
30    head ->next       0 x0008     0 x0016
31    head ->next ->next  0 x0004     0 x0020 ( wrong value !)
```

By checking the pointer table we know pointer "head→next→next" should be an alias of pointer "head". So the new allocated value for "head→next→next" should be updated immediately from 0x0020 to 0x0012 as follows:

```
1  = decode ( empty )
2
3    head ->next ->val = 20;
4    head ->next ->next = 0 x0012 ;
5    or
6    head ->next ->next = head ;
7
8    Pointer table :
9    Pointer              sender    receiver
```

```
10    head                  0x0004    0x0012
11    head->next            0x0008    0x0016
12    head->next->next  0x0004    0x0012
```

# 3   Separation Framework

## 3.1   Program dependence graph

Program Dependence Graph(PDG) plays a key role in our separation framework. In our framework, we establish a program dependence graph for each procedure first, and then use some connection techniques to connect the callers and the callees to model the inter-procedural dependences. For global data structures, since they belong to no special procedure at all, we simply use isolated nodes to represent them. Whenever a global data structure is used by a procedure, which means a non-intraprocedural data dependence occurs, we connect the corresponding global nodes with the nodes inside that procedure.

```
1  #include <stdio.h>
2  #include <string.h>
3
4  static char password[20] ;
5
6  void foo(){
7     printf("function foo should be separated!\n");
8  }
9
10 int encrypt(char* str,int key){
11    if(str == NULL) return 1;
12    unsigned int i;
13    for(i=0;i < strlen(str);++i){
14       str[i] = str[i] - key;
15    }
16    return 0;
17 }
18
19 int main(){
20    foo();
21    strncpy(password,"007",20);
22    encrypt(password,5);
23    printf("Encrypted value = %s\n",password);
24    return 0;
25 }
```

Fig 1. a C program to encrypt a password

Let us use the toy program in Fig.1 to show our basic working flow for the PDG construction.

9

We use solid lines to represent the control dependences, and dotted lines to represent the data dependences.

The real difficulty is how to represent the inter-procedure data dependency. Out solution is using parameter trees, which

One problem when using parameter trees is elegantly deal with parameters with recursive data structures(e.g. linked list), which can lead to parameter trees of infinite depth. On this issue, we borrow the idea of J. Graf's work[1]. ...

illustrate how cyclic(recursive) structures are accommodated by ...

linked list, tree...

# References

[1] J. Graf, *Speeding up context-, object- and field-sensitive SDG generation.*, 3rd ed. Harlow, England: Addison-Wesley, 1999.

Graf. In Proc. 9th IEEE International Working Conference on Source Code Analysis and Manipulation . IEEE Computer Society, Sept. 2010