

Type Encoding/Decoding rules: from C to Protocol Buffer

Shen Liu

(Version 1.3, last modified 06/15/2016)

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:

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

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

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

and an API for getting the value of x:

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

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

```

typedef struct{
    int x;
    int y;
}Point;

typedef struct{
    Point center;
    double radius;
}Circle;

message Circle{
    message Point{
        int64 x=1;
        int64 y=2;
    }
    double radius=1;
}

```

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:

```

Type t := int | char | float | unsigned long | t*
        | char[] | struct S {t1; t2; ... ; tn}
        | typedef-name (struct alias declared by "typedef")

```

Any type in this type system will be encoded as a list of integers(int[] lst), and the first 4 bytes(see the following table) in this list denotes what C type this list stands for.

lst[0]	type
0	int
1	unsigned long
2	float
3	pointer
4	char[]
5	struct
6	typedef-name

Once we have such a type system, we can easily construct the encoding/decoding rules as follows:

Basic value conversion functions:

```
intToBytes(int): convert an integer to a byte string.
ulongToBytes(unsigned long): convert an unsigned long integer
                           to a byte string.
floatToBytes(float): convert an float number to a byte string.
strToBytes(char[]): convert a C string to a byte string.
ptrToBytes(unsigned long): convert a pointer address(unsigned long)
                           to a byte string.

bytesToInt(char bytes[]): convert a byte string to an integer.
bytesToUlong(char bytes[]): convert a bytes string to an unsigned
                           long integer.
bytesToFloat(char bytes[]): convert a byte string to a float number.
bytesToStr(char bytes[]): convert a byte string to a char[].
bytesToPtr(char bytes[]) convert a byte string to a hexadecimal number.
```

Encoding:

```
unfold(typedef-name): return the original struct type before
                      redeclaration with "typedef"
```

```
fold(struct S): return a struct type with the name of "typedef-name"
```

```
Definition encode: (type,value) (t,v) -> char bytes[]
match t with
| int      => [0]::intToBytes(v)
| unsigned long
            => [1]::ulongToBytes(v)
| float    => [2]::floatToBytes(v)
| t*       => [3]::ptrToBytes(v)::encode(t,*v)
| char[]   => [4]::intToBytes(strlen(v))::strToBytes(v)
| struct S ((t1,v1);...;(tn,vn))
            => [5]::encode(char[],"S")
                ::intToBytes(n)
                ::encode(t1,v1)::...::encode(tn,vn)
| typedef-name
            => [6]::encode(char[],"typedef-name")
                ::encode(unfold(typedef-name),v)
end.
```

Decoding:

```

Definition decode bytes[] lst =
  match lst[0] with
  | 0 => ((int, bytesToInt([lst+1,lst+4])), tail(lst))

  | 1 => ((unsigned long, bytesToUlong([lst+1,lst+4])), tail(lst))

  | 2 => ((float, bytesToFloat([lst+1,lst+4])), tail(lst))

  | 3 => let ((t,*v), l1) = decode (tail(lst)) in ((t*, v), l1)
        //v = bytesToPtr([lst+1,lst+4])

  | 4 => ((char[],bytesToStr([lst+2,lst+offset-1])),tail(lst))
        /*strlen = lst[1], offset = 2+strlen*/

  | 5 =>
    let (char[],S) = decode([lst+1, lst+offset-1]) in
    /*offset = 2+4+bytesToInt([lst+2,lst+5])*/

    let n = lst[offset] in
    let ((t1,v1),l1) = decode (tail(tail(lst))) in
    let ((t2,v2),l2) = decode l1 in
    ...
    let ((tn,vn),ln) = decode l_{n-1} in
    (struct S {(t1,v1);(t2,v2);...;(tn,vn)}, ln)

  | 6 =>
    let (char[],"typedef-name") = decode([lst+1,lst+offset]) in
    fold(decode(tail(lst)))
    /*offset = 2+4+bytesToInt([lst+2,lst+5])*/
end

```

Now consider a circular linked list example:

```

typedef struct Node{
  int val;
  struct Node* next;
}Node_t;

Node_t *head = (struct Node*) malloc(sizeof(Node_t)); //head: 0x8000
Node_t *tail = (struct Node*) malloc(sizeof(Node_t)); //tail: 0x8008

head->val = 10;
head->next = tail;

tail->val = 20;
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:

```

encode(Node_t *head, 0x8000)

= [3]::ptrToBytes(0x8000)
  :: encode(Node_t, MemRead[0x8000]) //head->next = tail = 0x8008
  /*MemRead returns *head, which equals to {10,0x8008} */

= [3]::ptrToBytes(0x8000)
  :: [6]::encode(char[], "Node_t")
  :: encode(unfold(Node_t), {10, 0x8008})

= [3]::ptrToBytes(0x8000)
  :: [6]::[4]::intToBytes(strlen("Node_t"))
  :: strToBytes("Node_t")
  :: encode(struct Node{int, struct Node*}, {10, 0x8008})

= [3]::ptrToBytes(0x8000)
  :: [6]::[4]::intToBytes(6)
  :: strToBytes("Node_t")
  :: [5]::[4]::intToBytes(4)
  :: strToBytes("Node") // struct name "Node"
  :: [0]::intToBytes(2) //two fields
  :: [0]::intToBytes(10) //head->val = 10
  :: [3]::ptrToBytes(0x8008)
  :: encode(struct Node, MemRead[0x8008])

/* MemRead returns *(head->next), which equals to {20,0x8000}
   since tail->next = head */

= [3]::ptrToBytes(0x8000)
  :: [6]::[4]::intToBytes(6)
  :: strToBytes("Node_t")
  :: [5]::[4]::intToBytes(4)
  :: strToBytes("Node")
  :: [0]::intToBytes(2) //two fields
  :: [0]::intToBytes(10) //head->val = 10
  :: [3]::ptrToBytes(0x8008)
  :: [5]::[4]::intToBytes(4)
  :: strToBytes("Node")
  :: [0]::intToBytes(2)
  :: [0]::intToBytes(20)
  :: [3]::ptrToBytes(0x8000)
  ...
  (stop here...)

```

0x8000 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.