

# Lecture Notes on Generic Data Structures

15-122: Principles of Imperative Computation  
Frank Pfenning, Penny Anderson

Lecture 22  
November 13, 2014

## 1 Introduction

Using `void*` to represent pointers to values of arbitrary type, we were able to implement *generic* stacks in that the types of the elements were arbitrary. The main remaining restriction was that they had to be pointers. Generic queues or unbounded arrays can be implemented in an analogous fashion. However, when considering, say, hash tables or binary search trees, we run into difficulties because implementations of these data structures require operations on data provided by the client. For example, a hash table implementation requires a hash function and an equality function on keys. Similarly, binary search trees require a comparison function on keys with respect to an order. In this lecture we show how to overcome this limitation using function pointers.

## 2 The Hash Table Interface Revisited

Recall the client-side interface for hash tables, online [here](#). The client must provide a type `elem` (which must be a pointer), a type `key` (which was arbitrary), a hash function on keys, an equality function on keys, and a function to extract a key from an element. We write `___` while a concrete type must be supplied there in the actual file.

```

/*****
/* Hash sets, client-side interface */
*****/

// typedef _____* elem;

// f(x, y) returns true if x and y contain equal keys
bool elem_equal(elem x, elem y)
/*@requires x != NULL && y != NULL; @*/ ;

// f(x) returns the hash value of x's key
int elem_hash(elem x)
/*@ requires x != NULL; @*/ ;

```

We were careful to write the implementation so that it did not need to know what these types and functions were. But due to limitations in C0, we could not obtain multiple implementations of hash tables to be used in the same application, because once we fix `elem`, the above two client functions, they cannot be changed.

Given an implementation of the above client interface, the library provides a type `hset` of hash tables and means to create, insert, and search through a hashtable-backed set.

```

/*****
/* Hash table library side interface */
*****/

// typedef _____ hset;
typedef struct hset_header* hset;

hset hset_new(int capacity)
/*@requires capacity > 0; @*/ ;

// Returns the element of H with the equivalent to x, if it exists.
// Otherwise returns NULL.
elem hset_lookup(hset H, elem x)
/*@requires x != NULL; @*/ ;

void hset_insert(hset H, elem x)
/*@requires x != NULL; @*/ ;

```

### 3 Generic Types

Since both keys and elements are defined by the clients, they turn into generic pointer types when we implement a truly generic structure in C. We might try the following in a file `hset.h`.

```
#include <stdbool.h>
#include <stdlib.h>

#ifndef _HSET_H_
#define _HSET_H_

typedef void *elem;

/* Hash table interface */
typedef struct hset_header *hset;

hset hset_new (size_t capacity);
void hset_insert(hset H, elem x);
elem hset_lookup(hset H, elem x);

#endif
```

We use type definitions instead of writing `void*` in this interface so the role of the arguments as keys or elements is made explicit (even if the compiler is blissfully unaware of this distinction).

However, this does not yet work. Before you read on, try to think about why not, and how we might solve it

## 4 Function Pointers

The problem with the approach in the previous section is that the implementation of hash tables must call the functions `elem_equal` and `elem_hash`. Their types would now involve `void*` but in the environment in which the hash table implementation is compiled, there can still only be one of each of these functions. This means the implementation cannot be truly generic. We could not even use two hash tables with different element types simultaneously this way, without copying code and renaming things. (This actually happened with stacks in the Clac programming assignment, if you recall: we had the type `istack`, stacks of ints, as well as the type `qstack`, stacks of queues of strings.)

The underlying issue that we are trying to solve in this lecture is a deep one: how can a language support *generic* implementations of data structures that accommodate data elements of different types. The name *polymorphism* derives from the fact that data take on different forms for different uses of the same data structure. Sophisticated mechanisms to support polymorphism have been developed for modern high-level languages like Java and ML. Here we will look at a simple mechanism, the function pointer. In combination with void pointers and header files, function pointers give us the ability to write generic implementations of data structures. We use void pointers to pass around generic references to data, and function pointers to allow the client to specify to the library how to handle that data.

Because the client knows what these functions should be, it can define them, but must somehow communicate the definitions to the library. The way the client does this is by passing the *address* of a defined function to the library, taking advantage of the fact that the implementation of a function is stored in memory like everything else in C, and therefore a function has an address. These addresses are passed from client to library as pointers to functions.

Leaving generic hash tables aside for a moment, we will use a simple example of sorting to demonstrate this. In C, we can write a string sorting function that takes an array of strings, a lower bound, and an upper bound:

```
void sort(char **A, int lower, int upper);
```

We cannot make this generic by simply changing `char**` to `void**` (an array of void pointers), because we have to be able to compare array elements to sort them.

A comparison function, as we have seen, takes two elements and returns a negative number if the first element is smaller, zero if they are equal,

and a positive number if the first element is bigger. So the comparison function for generic `void*` elements has the following signature:

```
int compare(void* x, void* y);
```

If we want to compare strings (which have C type `char*`), we can use the `strcmp` function from the string library `<string.h>`:

```
#include <string.h>
int string_compare(void* s1, void* s2) {
    return strcmp((char*)s1, (char*)s2);
}
```

We can get a pointer to this function with the *address-of* operator by writing `&string_compare`. If `cmp` is a pointer obtained in this way, we can use it to compare two strings by writing `(*cmp)((void*)"hi", (void*)"yo")`. Note that when we write `(*cmp)`, we are dereferencing the function pointer to get at the actual function!

Generic client functions like this comparison function must be used carefully – if `x` and `y` are pointers to integers, then the result of calling `string_compare((void*)x, (void*)y)` is undefined. This is an easy mistake to make.

What is the *type* of a pointer to the function `string_compare`? We'll, it's a pointer to a function. It's possible to write down the type of a function directly, but it's simpler if we use a type definition:

```
typedef int compare_fn(void *x, void *y);
```

defines `compare_fn` to be the type of functions taking two void pointers and returning an integer. Now we can refer to a pointer to such a function straightforwardly:

```
compare_fn *cmp = &string_compare;
```

With this type definition, we can declare the generic type of sorting functions in `sort.h`:

```
void sort(elem* A, int lower, int upper, compare_fn *compare);
```

where `elem` is defined as `void*` for readability purposes and we can use an implementation of this sorting function to sort an array of `void*` where the elements are actually strings:

```
void** S = xmalloc(4, sizeof(void*));
S[0] = (void*)"pancake";
S[1] = (void*)"waffle";
S[2] = (void*)"toast";
S[3] = (void*)"juice";
sort(S, 0, 4, &string_compare);
```

The sorting library doesn't know, and doesn't need to know, that the void pointers are actually character arrays (that is, C strings). All it needs to know is that the comparison function we passed to the library knows what these pointers are and is able to compare them.

## 5 Generic Operations via Function Pointers

We now return to our hash table implementation problem. With function pointers, we can make the hash table implementation truly generic by allowing the client to provide pointers to the functions for extracting, comparing, and hashing keys.

But where do we pass them? We could pass all three to `ht_insert` and `ht_lookup`, where they are actually used. However, it is awkward to do this on every call. We notice that for a particular hash table, all three functions should be the same for all calls to insert into and search this table, because a single hash table stores elements of the same type and key. We can therefore pass these functions just once, when we first create the hash table, and store them with the table!

This gives us the following interface (in file `ht.h`):

```
#include <stdbool.h>
#include <stdlib.h>

#ifdef _HASHTABLE_H_
#define _HASHTABLE_H_

typedef void *elem;
typedef bool elem_equal_fn(elem x, elem y);
typedef size_t elem_hash_fn(elem x);
typedef void elem_free_fn(elem x);

typedef struct hset_header* hset;
```

```

hset hset_new(size_t capacity,          // Must be > 0
              elem_equal_fn *elem_equal, // Must be non-NULL
              elem_hash_fn *elem_hash,   // Must be non-NULL
              elem_free_fn *elem_free);  // May be NULL

void hset_insert(hset H, elem x);

elem hset_lookup(hset H, elem x);

void hset_free(hset H);

#endif

```

We have added the function `ht_free` to the interface; we'll come back to that.

We have made some small changes to exploit the presence of unsigned integers (in `elem_hash`) and the `size_t` type (also unsigned) to provide more appropriate types to certain functions.

Storing the function for manipulating the data brings us closer to the realm of object-oriented programming where such functions are called *methods*, and the structures they are stored in are *objects*. We don't pursue this analogy further in this course, but you may see it in follow-up courses, specifically 15-214 *Software System Construction*.

## 6 Using Generic Hashtables

First, we see how the client code works with the above interface. We use here the example of word counts, which we also used to illustrate and test hash tables earlier. The structure contains a string and a count.

```

/* elements */
struct wc {
    char *word;          /* key */
    int count;          /* information */
};

```

As mentioned before, strings are represented as arrays of characters (type `char*`). The C function `strcmp` from library with header `string.h` compares strings. We then define:

```
bool words_equal(elem x, elem y) {
    struct wc *w1 = (struct wc*)x;
    struct wc *w2 = (struct wc*)y;
    return strcmp(w1->word , w2->word) == 0;
}
```

Keep in mind that `elem` is defined to be `void*`. We therefore have to cast it to the appropriate struct type before we dereference it and pass it to `strcmp`, which requires two strings as arguments. We have to do the same thing when we write the hash function:

```
bool words_hash(elem x) {
    return hash_string(((struct wc*)x)->word);
}
```

Here is an example where we insert strings created from integers (function `itoa`) into a hash table and then search for them.

```
int n = (1<<10);
ht H = ht_new(n/5, &elem_key, &key_equal, &key_hash);
for (int i = 0; i < n; i++) {
    struct wc* e = xmalloc(sizeof(struct wc));
    e->word = itoa(i);
    e->count = i;
    hset_insert(H, e);
}
for (int i = 0; i < n; i++) {
    struct wc A;
    A.word = itoa(i);
    struct wc *wcount = (struct wc*)(ht->lookup(H, (void*)&A));
    assert(wcount->count == i);
    free(A.word);
}
```



## 7 Implementing Generic Hash Tables

The hash table structure, defined in file `hset.c` now needs to store the function pointers passed to it.

```
struct hset_header {
    size_t size;
    size_t capacity;           /* 0 < capacity */
    chain **table;            /* \length(table) == capacity */
    elem_equal_fn *elem_equal; /* elem_equal != NULL */
    elem_hash_fn *elem_hash;   /* elem_hash != NULL */
    elem_free_fn *elem_free;   /* can be null */
};
```

The `ht_new` function takes its extra arguments and stores them in these new fields:

```
hset hset_new(size_t capacity,
              elem_equal_fn *elem_equal,
              elem_hash_fn *elem_hash,
              elem_free_fn *elem_free) {
    REQUIRES(capacity > 0 && elem_equal != NULL && elem_hash != NULL);

    hset H = xmalloc(sizeof(struct hset_header));
    H->size = 0;
    H->capacity = capacity;
    H->table = xcalloc(capacity, sizeof(chain*));
    H->elem_equal = elem_equal;
    H->elem_hash = elem_hash;
    H->elem_free = elem_free;

    ENSURES(is_hset(H));
    return H;
}
```

We do not require that `elem_free` be non-NULL. The reason for this is that the client may want to retain *ownership* of the data they have stored in the data structure. If this function pointer is NULL, then we assume the client retains ownership and only free the data that is part of the hashset in `hset_free`, but if the client gives us a function, we use it to free the structure.

```
void hset_free(hset H) {
    REQUIRES(is_hset(H));

    for (size_t i = 0; i < H->capacity; i++) {
        chain* p_next;
        for (chain* p = H->table[i]; p != NULL; p = p_next) {
            if (H->elem_free != NULL) (*H->elem_free)(p->data);
            p_next = p->next;
            free(p);
        }
    }

    free(H->table);
    free(H);
}
```

In this free function, we use the explicit syntax for invoking a function pointer: `(*H->elem_free)(p->data)`. This makes sense, because this is the only place in our code where we call `elem_free`. The other functions we store in the data structure are ones we make constant use of, so we'll create some helper functions that allow us to do so more conveniently.

```
static inline size_t hashindex(hset H, elem x) {
    REQUIRES(H != NULL && H->capacity > 0 && H->elem_hash != NULL);
    return (*H->elem_hash)(x) % H->capacity;
}

static inline bool elemequal(hset H, elem x, elem y) {
    REQUIRES(H != NULL && H->capacity > 0 && H->elem_equal != NULL);
    return (*H->elem_equal)(x, y);
}
```

The `static inline` keywords are optional. We use here the directive `static inline` to instruct the compiler to *inline* the function, which means that wherever a call to this function occurs, the compiler just replaces the call by the function body, for the sake of efficiency. This provides a similar but semantically cleaner and less error-prone alternative to C preprocessor macros.

We use these helper functions in the lookup and insertion functions.

```
elem ht_lookup(ht H, elem x)
{
    REQUIRES(is_ht(H));
    int i = elemhash(H, x);
    chain* p = H->table[i];
    while (p != NULL) {
        ASSERT(p->data != NULL);
        if (elemequal(H, p->data, x))
            return p->data;
        else
            p = p->next;
    }
    /* not in chain */
    return NULL;
}
```

This concludes this short discussion of generic implementations of libraries, exploiting `void*` and function pointers.

In more modern languages such as ML, so-called *parametric polymorphism* can eliminate the need for checks when coercing from `void*`. The corresponding construct in object-oriented languages such as Java is usually called *generics*. We do not discuss these in this course.

## 8 A Subtle Memory Leak

Let's look at the beginning code for insertion into the hash table.

```
void hset_insert(hset H, elem x) {
    REQUIRES(is_hset(H) && x != NULL);

    int i = hashindex(H, x);
    for (chain* p = H->table[i]; p != NULL; p = p->next) {
        if (elemequal(H, p->data, x)) {
            p->data = x;
            return;
        }
    }

    // prepend new element
    chain* p = xmalloc(sizeof(chain));
    p->data = x;
    p->next = H->table[i];
    H->table[i] = p;
    (H->size)++;

    ENSURES(is_hset(H));
    return NULL;
}
```

At the end of the while loop, we know that the key  $k$  is not already in the hash table. But this code fragment has a subtle memory leak. Can you see it?<sup>1</sup>

---

<sup>1</sup>The code author overlooked this in the port of the code from C0 to C, but one of the students noticed.

The problem is that when we overwrite `p->data` with `e`, the element currently stored in that field may be lost and can potentially no longer be freed.

There seem to be two solutions. The first is for the hash table to apply the `elem_free` function it was given. We should guard this with a check that the element we are inserting is indeed new, otherwise we would have a freed element in the hash table, leading to undefined behavior.

```
if (elemequal(H, p->data, x)) {
    /* free existing element, if different from new one */
    if (p->data != x) (*H->elem_free)(p->data);
    /* overwrite existing element */
    p->data = e;
    return;
}
```

The client has to be aware that the element already in the table will be freed when a new one with the same key is added.

In order to avoid this potentially dangerous convention, we can also just *return* the old element if there is one, and `NULL` otherwise. The information that such an element already existed may be useful to the client in other situations, so it seems like the preferable solution. The client could always immediately apply the element free function if that is appropriate. This requires a small change in the interface, but first we show the relevant code.

```
int i = hashindex(H, x);
for (chain* p = H->table[i]; p != NULL; p = p->next) {
    if (elemequal(H, p->data, x)) {
        elem old = p->data;
        p->data = x;
        return old;
    }
}
```

The relevant part of the revised header file `ht.h` now reads:

```
typedef void *elem;
typedef bool elem_equal_fn(elem x, elem y);
typedef size_t elem_hash_fn(elem x);
typedef void elem_free_fn(elem x);

typedef struct hset_header* hset;

hset hset_new(size_t capacity,          // Must be > 0
             elem_equal_fn *elem_equal, // Must be non-NULL
             elem_hash_fn *elem_hash,   // Must be non-NULL
             elem_free_fn *elem_free);  // May be NULL

elem hset_insert(hset H, elem x);

elem hset_lookup(hset H, elem x);

void hset_free(hset H);
```

## 9 Separate Compilation

Although the C language does not provide much support for modularity, convention helps. The convention rests on a distinction between *header files* (with extension `.h`) and *program files* (with extension `.c`).

When we implement a data structure or other code, we provide not only `filename.c` with the code, but also a header file `filename.h` with declarations providing the interface for the code in `filename.c`. The implementation `filename.c` contains `#include "filename.h"` at its top, and client will have the same line. The fact that both implementation and client include the same header file provides a measure of consistency between the two.

A header file `filename.h` should never contain any function definitions (that is, code), only type definitions, structure declarations, macros, and function declarations (so-called function prototypes). In contrast, a program file `filename.c` can contain both declarations and definitions, with the understanding that the definitions are not available to other files.

These header files have *header guards* that prevent the compiler from processing them more than once when compiling several files at the same time (thus they are sometimes called “once-only headers”). The guards are

directives to the C preprocessor, perhaps best explained by example. Here again is the header file for hashtables:

```
#include <stdlib.h>
#include <stdbool.h>

#ifndef _HSET_H_
#define _HSET_H_

typedef void *elem;
typedef bool elem_equal_fn(elem x, elem y);
typedef size_t elem_hash_fn(elem x);
typedef void elem_free_fn(elem x);

typedef struct hset_header* hset;

hset hset_new(size_t capacity,           // Must be > 0
             elem_equal_fn *elem_equal, // Must be non-NULL
             elem_hash_fn *elem_hash,   // Must be non-NULL
             elem_free_fn *elem_free);  // May be NULL

elem hset_insert(hset H, elem x);

elem hset_lookup(hset H, elem x);

void hset_free(hset H);

#endif
```

The presence of `#ifndef ... #endif` causes the preprocessor to check whether it has already defined `_HASHTABLE_H_`. The first time it scans the file, it will not have defined it (note the importance of choosing a name that is unlikely to occur in other headers!), and so it processes everything up to the `#endif`. Any subsequent scans will skip everything between `#ifndef` and `#endif`. In the case of this particular header, no harm is done other than a waste of time in processing it more than once. But unpleasant compiler errors can occur if headers in general are not once-only.

We only ever `#include` header files, never program files, in order to maintain the separation between code and interface.

When `gcc` is invoked with multiple files, it behaves somewhat differently than `cc0`. It compiles each file *separately*, referring only to the included

header files. Those come in two forms, `#include <syslib.h>` where `syslib` is a system library, and `#include "filename.h"`, where `filename.h` is provided in the local directory. Therefore, if the right header files are not included, the program file will not compile correctly. We never pass a header file directly to `gcc`.

The compiler then produces a separate so-called *object file* `filename.o` for each `filename.c` that is compiled. All the object files are then *linked* together to create the executable. By default, that is `a.out`, but a name for the executable can be provided with the `-o` switch.

Let us summarize the most important conventions:

- Every file `filename`, except for the one with the main function, has a header file `filename.h` and a program file `filename.c`.
- The program `filename.c` and any client that would like to use it has a line `#include "filename.h"` at the beginning.
- The header file `filename.h` never contains any code, only macros, type definition, structure definitions, and function headers (prototypes). It has appropriate header guards to avoid problems if it is loaded more than once.
- We never `#include` any program files, only header files (with `.h` extension).
- We only pass program files (with `.c` extension) to `gcc` on the command line.

## Exercises

**Exercise 1** *Convert the interface and implementation for binary search trees from C0 to C and make them generic. Also convert the testing code, and verify that no memory is leaked in your tests. Make sure to adhere to the conventions described in Section 9.*