

Assignment 1: Six Degrees of Kevin Bacon

Craving a little Oscar trivia? Try your hand in an Internet parlor game about Kevin Bacon's acting career. He's never been nominated for an Oscar, but he's certainly achieved immortality—based on the premise that he is the hub of the entertainment universe. Mike Ginelli, Craig Fass and Brian Turtle invented the game while students at Albright College in 1993, and their Bacon bit spread rapidly after convincing then TV talk-show host Jon Stewart to demonstrate the game to all those who tuned in. From these humble beginnings, a website was built, a book was published and a nationwide cult-fad was born.

When you think about Hollywood heavyweights, you don't immediately think of Kevin Bacon. But his career spans almost 20 years through films such as *Flatliners*, *The Air Up There*, *Footloose*, *The River Wild*, *JFK* and *Animal House*. So brush up on your Bacon lore. To play an Internet version, visit <http://oracleofbacon.org>.

This assignment is first and foremost a low-level systems programming assignment, but it's also an opportunity to review your C++ while simultaneously exercising your software engineering and low-level memory manipulation skills. You'll also get to see that low-level C coding and high-level C++ data structuring can coexist in the same application.

Due: Wednesday, April 12th at 11:59 pm

The game takes the form of a trivia challenge: Propose two names, and your friend/opponent has to come up with a sequence of movies and mutual co-stars connecting the two. In this case, your opponent takes on the form of an executable, and that executable is shockingly good.

Jack Nicholson and Meryl Streep? That's easy:

```
poohbear@myth10$ ./search "Meryl Streep" "Jack Nicholson (I)"
Meryl Streep was in "Heartburn" (1986) with Jack Nicholson (I).
```

Mary Tyler Moore (rest her soul) and Red Buttons? Not so obvious:

```
poohbear@myth10$ ./search "Mary Tyler Moore" "Red Buttons"
Mary Tyler Moore was in "Change of Habit" (1969) with Regis Toomey.
Regis Toomey was in "C.H.O.M.P.S" (1979) with Red Buttons.
```

Barry Manilow and Lou Rawls? Yes!

```
poohbear@myth10$ ./search "Barry Manilow" "Lou Rawls"
Barry Manilow was in "Bitter Jester" (2003) with Dom Irrera.
Dom Irrera was in "A Man Is Mostly Water" (2000) with Lou Rawls.
```

It's the people you've never heard of that are far away from each other:

```
poohbear@myth10$ ./search "Danzel Muzingo" "Liseli Mutti"
Danzel Muzingo was in "My Day in the Barrel" (1998) with Damian Brown.
Damian Brown was in "Bad Chemistry" (1997) with Dick Welsbacher.
Dick Welsbacher was in "The Attic" (1980) with Carrie Snodgress.
Carrie Snodgress was in "Chill Factor" (1989) with Nathaniel Lees (I).
Nathaniel Lees (I) was in "Rapa Nui" (1994) with Liseli Mutti.
```

Is it true? Your buffoon of a lecturer has a Bacon number of 3?

```
poohbear@myth10$ ./search "Jerry Cain" "Kevin Bacon (I)"
Jerry Cain was in "No Rules" (2005) with Romeo Antonio.
Romeo Antonio was in "Doesn't Texas Ever End" (2009) with Irwin Keyes.
Irwin Keyes was in "Friday the 13th" (1980) with Kevin Bacon (I).
```

I have no idea who this particular Jerry Cain is, but it's certainly not me.

Overview

There are two major components to this assignment:

- You need to provide the implementation for **an imdb class**, which allows you to quickly look up all of the films an actor or actress has appeared in and all of the people starring in any given film. We **could** layer our imdb class over two STL maps—one mapping people to movies and another mapping movies to people—but that would require we read in several megabytes of data from **flat text files**. That type of configuration takes several minutes, even on fast machines, and it's the opposite of fun if you have to sit that long before you play. Instead, you'll tap your sophisticated understanding of data representation and learn about something called memory mapping in order to look up movie and actor information from a prepared data structure that's been saved to disk in its binary form (and I'll describe the binary image format in the pages to come). This is the meatier part of the assignment. (By the way, imdb is short for Internet Movie Database. Our name is a gesture to the company that provides all of the data for the hundreds of thousands of movies and movie stars.)

- You also need to implement a **breadth-first search algorithm** that consults your super-clever `imdb` class to find the shortest path connecting any two actor/actresses. If the search goes on for so long that you can tell it'll be of length 7 or more, then you can be reasonably confident (and pretend that you know for sure that) there's no path connecting them. This part of the assignment is more CS106B-like, and it's a chance to get a little more experience with the STL (using vectors, sets, and lists) and to see a legitimate scenario where a complex program benefits from coding in two different paradigms: high-level, object-oriented C++ (with its STL template containers and template algorithms) and low-level, imperative C (with its exposed memory, brought to you by CS107, *, &, [], and ->).

Task I: The `imdb` class

First off, I want you to complete the implementation of the `imdb` class. Here's the interface:

```
class imdb {
public:
    imdb(const string& directory);
    bool good() const;
    bool getCredits(const string& player, vector<film>& films) const;
    bool getCast(const film& movie, vector<string>& players) const;
    ~imdb();

private:
    const void *actorFile;
    const void *movieFile;
};struct film {
    string title;
    int year;
};

class imdb {
public:
    imdb(const string& directory);
    bool good() const;
    bool getCredits(const string& player, vector<film>& films) const;
```

```

    bool getCast(const film& movie, vector<string>& players) const;

    ~imdb();

private:
    const void *actorFile;
    const void *movieFile;
};

```

The constructor and destructor have already been implemented for you. All the constructor does is initialize actorFile and movieFile fields to point to on-disk data structures using the mmap routine you'll learn about later on in the course. Since you haven't used mmap before, I implemented the constructor and destructor for you. (The destructor just munmaps what was mmap-ed at construction time).

You'll need to implement the getCredits and getCast methods by manually crawling over these binary images in order to produce vectors of movies and actor names. When properly implemented, they provide lightning-speed access to a gargantuan amount of information, because the information is already compactly formatted in a preprepared data structure that permanently lives on the myth machines.

Understand up front that you are implementing these two methods to crawl over two arrays of bytes in order to synthesize data structures for the client. What appears below is a description of how that memory is laid out. You aren't responsible for creating the data files in any way, but you are just responsible for understanding how everything is encoded so that you can rehydrate information from their byte-level representations.

The Raw Data Files

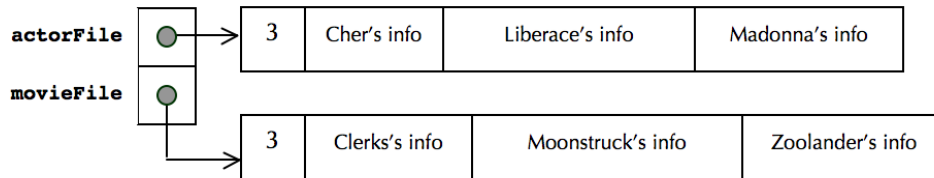
The private actorFile and movieFile fields each address gigantic blocks of memory. Each configured to point to mutually referent database images, and the format of each is described below. The imdb constructor sets these pointers up for you, so you can proceed as if everything is initialized for getCast and getCredits to just work.

For the purposes of illustration, let's assume that Hollywood has produced a mere three movies, and that they've always rotated through the same three actors whenever the time came to cast their three films. Let's pretend those three films are as follows:

- Clerks, released in 1993, starring Cher and Liberace.
- Moonstruck, released in 1988, starring Cher, Liberace, and Madonna.
- Zoolander, released in 1999, starring Liberace and Madonna.

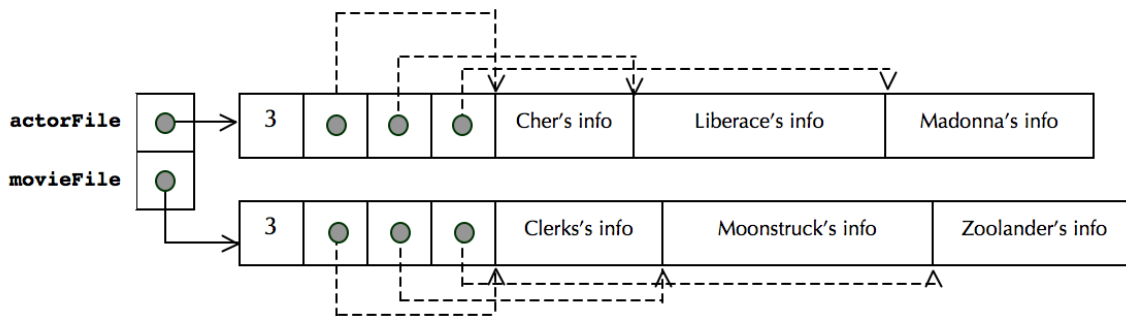
Remember, we're pretending.

If an `imdb` instance is configured to store the above information, you might imagine its `actorFile` and `movieFile` fields being initialized (by the constructor I already wrote for you) as follows:

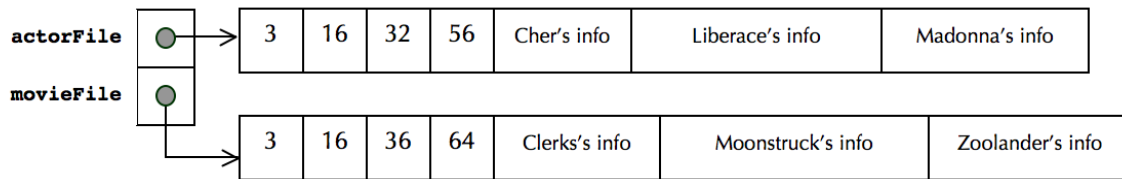


Each of the records for the actors and movies will be of variable size. Some movie titles are longer than others; some films feature 75 actors, while others star only one or two. Some people have prolific careers, while several people are one-hit wonders. Defining a `struct` or `class` to overlay the blocks of data is a fine idea, except that doing so would constrain all records to be the same size. We don't want that, because we'd be wasting a good chunk of memory when storing information on actors who appeared in just one or two films (and for films that feature just a handful of actors).

However, by allowing the individual records to be of variable size, we lose our ability to binary search (hint: via the STL `lower_bound` algorithm) a sorted array of records. The number of actors and actresses is just south of a million, and the number of movies is in the hundreds of thousands, so a linear search would be all turtle-like. All of the actors and movies are sorted by name (and then by year if two movies have the same name), so binary search is still within reach. The strong desire to search quickly motivated my decision to format the data files like this:



Spliced in between the number of records and the records themselves is an array of integer offsets. They're drawn as pointers, but they really aren't stored that way. We want the data images to be *relocatable*—that is, we want the information stored in the data images pointed to by `actorFile` and `movieFile` to be useful, regardless of what addresses get stored there. By storing integer offsets, we can manually compute the location of Cher's record, Madonna's record, or Clerk's record, etc, by adding the corresponding offsets to whatever `actorFile` or `movieFile` turn out to be. A more accurate picture of what gets stored (and this is really what the file format is) is presented here.



Because the numbers are what they are, we would expect Cher's 16-byte record to sit 16 bytes from the front of `actorFile`, Liberace's 24-byte record to sit 32 bytes within the `actorFile` image, and so forth. Looking for Moonstruck? Its 28-byte record can be found 36 bytes ahead of whatever address is stored in `movieFile`. Note that the actual offsets tell me where records are relative to the base address, and the **deltas** between offsets tell me how large the actual records are.

Because all of the offsets are stored as **four-byte integers** (and ints are four bytes, even on 64-bit systems like the `myths`), and because they are in a sense sorted if the records they reference are sorted, we can use binary search. Woo!

To summarize:

- `actorFile` points to a large mass of memory packing all of the information about all of the actors into one big blob. The first four bytes store the number of actors (as an `int`); the next four bytes store the offset to the zeroth actor, the next four bytes store the offset to the first actor, and so forth. The last offset is followed by the zeroth record, then the first record, and so forth. The records themselves are sorted by name. Pinky swear they are.
- `movieFile` also points to a large mass of memory, but this one packs the information about all films ever made. The first four bytes store the number of movies (again, as an `int`); the next `*(int *)movieFile * sizeof(int)` bytes store all of the `int` offsets, and everything beyond the offsets is real movie data. The movies are sorted by title, and those sharing the same title are sorted by year.
- The above description above generalizes to files with 900,000 actors and 300,000 movies.

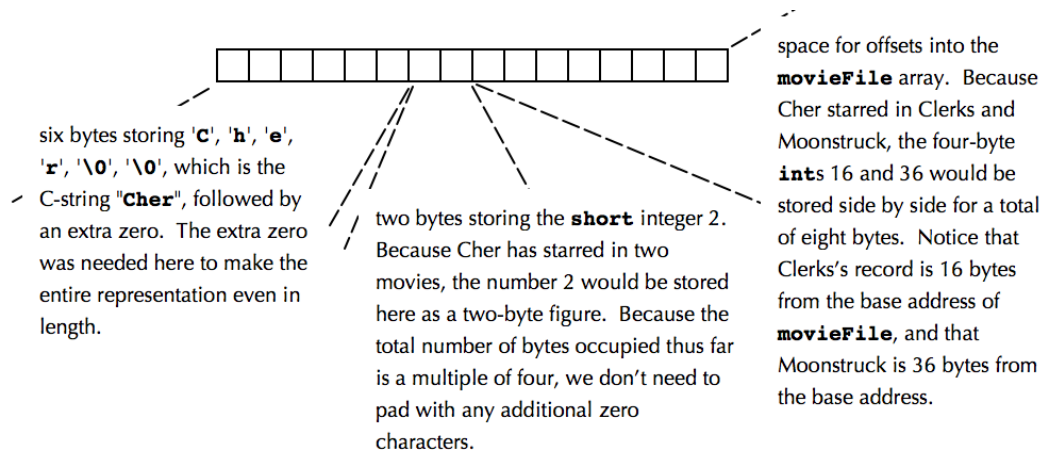
The Actor Record

The actor record is a packed set of bytes collecting information about an actor and the movies he or she's appeared in. We don't use a `struct` or a `class` to overlay the memory associated with an actor, because doing so would constrain the record size to be constant for all actors. Instead, we lay out the relevant information in a series of bytes, the number of which depends on the length of the actor's name and the number of films he's appeared in. Here's what gets manually placed within each entry:

1. The name of the actor is laid out character by character, as a normal null-terminated C-string. If the length of the actor's name is even, then the string is padded with an extra `'\0'` so that the total number of bytes dedicated to the name is always an even number. The information that follows the name is most easily interpreted as a `short`, and the `myths` might constrain addresses manipulated as `short *s` to be even.

2. The number of movies in which the actor has appeared, expressed as a two-byte short. (Some people have been in more than 255 movies, so a single byte isn't always enough). If the number of bytes dedicated to the actor's name (always even) and the short (always 2) isn't a multiple of four, then two additional '\0's appear **after** the two bytes storing the number of movies. This padding is conditionally done so that the four-byte integers that follow sit at addresses that are multiples of four (again, because the 64-bit `myth's` might be configured to require this).
3. An array of offsets into the `movieFile` image, where each offset identifies one of the actor's films.

Here's what Cher's record would look like:



The Movie Record

The movie record is only slightly more complicated. The information that needs to be compressed is as follows:

1. The title of the movie, terminated by a '\0', so the character array behaves as a normal C-string incidentally wedged into a larger binary data figure.
2. The year the film was released, expressed as a single byte. This byte stores the year, minus 1900. Since Hollywood is less than 256 years old, it was fine to just store the year as a delta from 1900. If the total number of bytes used to encode **the name and year** of the movie is odd, then an extra '\0' sits in between the one-byte year and the data that follows.
3. A two-byte **short** storing the number of actors appearing in the film, padding with two additional bytes of zeroes if needed.
4. An array of four-byte integer offsets, where each integer offset identifies one of the actors accessible via `actorFile`. The number of offsets here is, of course, equal to the **short** read during step 3.

One major gotcha: Some movies share the same title even though they are different. (The `Manchurian Candidate`, for instance, was first released in 1962, and then remade in 2004. They're

two different films with two different casts.) If you look in the `imdb-utils.h` file, you'll see that the `film` struct provides `operator<` and `operator==` methods. That means that two `film`s know how to compare themselves to each other using infix `==` and `<` (though not using `!=`, `>`, `>=`, or `<=`). You can just rely on the `<` and `==` to compare two `film` records. In fact, you **have to**, because the movies in the `movieData` binary image are sorted to respect `film::operator<`.

It's best to work on the implementation of the `imdb` class in isolation, not worrying about the details of the search algorithm you'll eventually need to write. I've provided a test harness to exercise the `imdb` all by itself, and that code sits in `imdbtest.cc`. The `make` system generates a test application called `imdbtest` which you can use to verify that your `imdb` implementation is solid. I provide my own version in `/usr/class/cs110/samples/assign1/imdbtest_soln`, so you can run your version and my version side by side and make sure they match character for character.

Note: Your implementation will be—and in fact is intended to be—an interesting mix of C and C++. You'll be relying on your mad C skillz to crawl over these binary images, and you'll be leveraging your C++ mastery to lift that data up into C++ objects. As part of your implementation, you'll need to binary search over the actor and movie offsets to find the actor or movie of interest.

I am **requiring** that you use the STL `lower_bound` algorithm to perform these binary searches, and that you use C++11 lambdas (also known as anonymous functions with capture clauses) to provide nameless comparison functions that `lower_bound` can use to guide its search.

Task II: Implementing Search

You're back in pure C++ mode. At this point, I'm assuming your `imdb` class just works, and the fact that there's some spectacularly shrewd pointer gymnastics going on in the `imdb.cc` file is completely disguised by the delightfully simple `imdb` interface. Use the services of your `imdb` and my `path` class (discussed below) to implement a breadth-first search for the shortest possible path. Leverage the STL containers as much as possible to get this done. Here are the STL classes I used in my solution:

- `vector`: there's no escaping this one, because the `imdb` requires we pull `films` and `actors` out of the binary images as `vectors`.
- `list`: The `list` is a doubly-linked list that provides `O(1)` `push_back`, `front`, and `pop_front` operations. There's also a `queue` template, and you can use that if you want, but I'm so bugged that the STL `queue` calls its methods `push` and `pop` instead of `enqueue` and `dequeue` that I boycotted and used the `list` instead.
- `set`: I used two `sets` to keep track of previously used actors and films. If you're implementing a breadth-first search and you encounter a movie or actor that you've seen before, there's no reason to use it/him/her a second time. You shouldn't need to use anything other than `set::insert`.

Getting Code

Go ahead and clone the mercurial repository that we've set up for you by typing:

```
poohbear@myth10$ hg clone /usr/class/cs110/repos/assign1/$USER assign1
```

Compile often, test incrementally and almost as often as you blink, `hg commit` a bunch so you don't lose your work in some bizarre `rm` disaster, and run `/usr/class/cs110/tools/submit` when you're done. As you make progress, you can invoke `/usr/class/cs110/tools/sanitycheck` to run a collection of tests that compare your solutions to my own. And be sure your solutions are free of memory leaks and errors, since we'll be running your code through `valgrind`. Note: there's a bug in `valgrind` itself that surfaces when virtually *any* C++ program is run through it. You can suppress this error by copying `/usr/class/cs110/tools/config/.valgrindrc` file into your home directory (or copying its contents into an existing `~/.valgrindrc` file):

```
poohbear@myth10$ more /usr/class/cs110/tools/config/.valgrindrc
--memcheck:leak-check=full
--show-reachable=yes
--suppressions=/usr/class/cs110/tools/config/cs110.supp
poohbear@myth10$ cp /usr/class/cs110/tools/config/.valgrindrc ~
```

Assignment 1 Files

Here's the subset of all the files that pertain to just the first of the two tasks:

`imdb-utils.h`

The definition of the `film` struct, and an inlined function that finds the data files for you. *You shouldn't need to change this file.*

`imdb.h`

The interface for the `imdb` class. *You shouldn't change the public interface of this file, though you're free to change the private section if it makes sense to.*

`imdb.cc`

The implementation of the `imdb` constructor, destructor, and methods. This is where your code for `getCast` and `getCredits` belongs.

`imdbtest.cc`

The unit test code we've provided to help you exercise your `imdb`. *You shouldn't have to change this file.*

`Makefile`

By typing `make imdbtest`, you'll compile just the files needed to build `imdbtest`. *You shouldn't need to change this file at all.*

Everything from Task I (except `imdbtest.cc`) contributes to the overall `search` application. Type `make search` to build the `search` executable without building the `imdbtest` application (or you can just type `make` and build both.) There's a sample executable at `/usr/class/cs110/samples/assign1/search_soln` for you to play with. Understand that my sample application and yours aren't obligated to publish the same exact shortest path, but you should be sure that the path lengths themselves actually match. In addition to the files used for Task I, there are these:

`search.cc`

The file where most if not all of your Task II changes should be made.

`path.h`

The definition of the `path` class, which is a custom class useful for building paths between two actors. *You're free to add methods if you think it's sensible to do so.*

`path.cc`

The implementation of the `path` class. *Again, you can add stuff here if you think it makes sense to.*