

# Comp 40 Homework 1: Files, Pictures and Interfaces

*Please read the entire assignment before starting work.*

## Table of Contents

1. Summary of the Assignment
2. Background and preparation
  1. Purpose of this assignment
  2. [C vs. C++ and Hanson's Interfaces and Implementations](#)
  3. [Setting up your environment](#)
  4. [Getting started with Hanson's Interfaces and Implementations](#)
  5. [Other preparation](#)
3. Part A: Brightness of a grayscale image
  1. [Specification](#)
  2. [Examples of Brightness in Use](#)
  3. [Help with image files](#)
  4. [Getting images](#)
  5. [Problem analysis and advice](#)
4. Part B: Read a line
  1. [readaline specification](#)
  2. [Partial credit](#)
  3. [Hints](#)
5. Part C: Similarities in files
  1. [simlines specification overview](#)
  2. [Definition of line similarity](#)
  3. [Simlines output specification](#)
  4. [Hints](#)
  5. [Performance target](#)
  6. [Problem analysis and advice](#)
6. Part C (DESIGN): Simlines Design
  1. [Design overview](#)
  2. [Design document specifics](#)
  3. [Submitting your design document](#)
  4. [Before submitting your code](#)
  5. [Submitting your completed work](#)
7. [Suggested order of work](#)
  1. [Deadlines and tokens](#)
8. [General advice for new C programmers](#)

## Summary of the Assignment

In this assignment you will design, build and test two application programs and one supporting file input routine. The first program reports the average brightness of an image file. The second program reports similarities in file data, and it uses an input routine called `readaline` that you will implement. That input routine must conform to an interface that we provide.

The primary specifications for these programs are contained in the sections [Part A \(Brightness\)](#), [Part B \(Read a Line\)](#) and [Part C \(Similarities in Files\)](#) below.

Note that for Part B, you get partial credit for a limited implementation that only handles short input lines. *We strongly urge you to implement that limited version first, then go on and complete Part C, and*

only if you have time go back to enhance Part B for full credit. You will lose much more credit for doing a poor job on Part C than for failing to handle long lines in Part B.

## Background and preparation

The sections below give information about the purpose of this assignment, as well as other background you should understand before starting work.

### ***Purpose of this assignment***

This assignment has several goals:

1. To teach you to work in pairs on assignments that are more challenging than you have encountered before
2. To help you make the transition from C++ programming to the kind of C programming we expect in COMP 40
3. To start you thinking about the *interface* as a unit of *design*
4. To give you practice in identifying interfaces and existing implementations that can help you solve problems
5. To give you experience reading and conforming exactly to *specifications*, such as the those contained in this document
6. To start to convince you that writing good test cases, not just for correct or obvious input but also for edge cases and error cases, is as important as writing good program code
7. To give you experience *teaching yourself* about languages like C, systems like Linux, and system features like `stdin`, `fopen`, etc.
8. To give you experience working with a partner to design, document, implement and test a computer program
9. To introduce you to multiple representations of numbers

We understand that assignments like this will take you out of your comfort zone. When you read these instructions there will be many things that at first you won't understand. *Not everything you need to know will be explained to you in detail.*

These are the challenges that professional programmers and computer scientists face every day. Stick with it, figure things out, use the system documentation, get help. You may not succeed completely at everything we ask of you even by the end of the assignment, although it's certainly possible and many students do. Experience shows that over time almost all of you will learn not just to build programs like these, but also to teach yourself about the language features and tools that you need.

### ***C vs. C++, and Hanson's Interfaces and Implementations***

The C language is for the most part a proper subset of C++. Indeed, C came first and was only many years later extended to add object orientation, parameterized types and other higher-level features features. Note that [both languages are very widely used to this day](#).

So, why C in comp 40? A primary goal of COMP 40 is to teach you how computers work, and to show you how modern software uses the hardware on which it runs. You will find that C, being a smaller and simpler language, translates much more directly to hardware primitives. Furthermore, working in C allows us to build ourselves many of the higher-level features whose implementation is hidden in the runtimes of more complex languages. So, using C we learn more about the internal workings of languages like C++ (and Java and Python, etc.)

Although C itself does not include the sorts of high-level structures like Lists, Sets, or hash-based Tables/Maps you might find in C++ or Java, we in COMP 40 give you [implementations by Dave Hanson](#). We expect you to use these rather than building your own, and we offer them as examples of interesting, well-written C code from which you can learn. There is an online [Quick Reference guide](http://ciibook.webhop.net/pdf/quickref.pdf) at <http://ciibook.webhop.net/pdf/quickref.pdf> **Note: the Hanson Array class is not available for use in COMP 40, and is not needed for this assignment.** Of course, you are also encouraged where appropriate to use C language `structs`, `arrays`, etc.

You will find some useful information for new C programmers in the section below titled: [General advice for new C programmers](#) Although we will give you some hints like these about the differences

between C and C++, we also expect to teach yourselves many of the details. Other good sources include the books that we have recommended, online resources, etc. Help each other learn! In COMP 40 you must not share work on your solutions, but you are welcome to work with your fellow students to learn new languages and technologies.

## Setting up your environment

- To add the course binaries to your execution path, run

```
use comp40
```

You may want to put this command in your `.cshrc` or your `.profile`:

```
use -q comp40
```

Without the `-q`, you may have difficulties with `scp`, `ssh`, [git](#), and [rsync](#).

- To get started:

```
git clone /comp/40/git/filesnpix
```

You will get two files: `compile` and `Makefile`. You should leave the `compile` script alone: it's a program that just runs `make`. It's there for legacy reasons as we transition the course from `compile` scripts to `makefiles`. Use `make` and the `Makefile` to build your programs. Depending on how you structure your code, you may have to make minor modifications to the `Makefile` so that it will compile and link together the correct source modules for your projects. The handout [A simple introduction to Compile Scripts and Makefiles](#) (concentrate on the `Makefile` part) should give you the information you need. In later projects you will likely have to make more extensive modifications to `Makefiles`, so you should read and learn about them so you understand what they're doing. The `Makefile` for Lab 0 and this one are intended as gentle introductions: we will gradually make use of more sophisticated features of `make` as the semester progresses. (You may want to compare this one to the Lab 0 `Makefile` to see what I mean.)

## Getting started with Hanson's Interfaces and Implementations

For [Part C: Similarities in files](#) you will need several of the *C interfaces and implementations* David Hanson. You'll also need a general understanding of Hansons conventions, exceptions for assertions in error handling, etc. for [Part A: Brightness](#).

- To get started, read Chapters 1 and 2 (pages 1–31) of Hanson's book.
- To learn what Hanson has built for you, skim the beginnings of the relevant chapters: pages 45–52, pages 33–34, pages 103–107, pages 115–118 (pages 118–125 recommended), pages 137–140, pages 161–164, pages 171–173, and the first sections of Chapters 15 and 16, which I don't have the page numbers for.

## Other preparation

Be sure you have read the [course policies](#), especially notes on the [general expectations for COMP 40 homework](#) (including all subsections!), [collaboration](#), guidelines for [pair programming](#) and the course [coding standards](#). *Note that terminology definitions (e.g. for the term Checked Runtime Error), and the expectations set in the sections on [Homework Grades](#) and [Errors, Exceptions, Output, Valgrind and Grades](#) etc. are implicitly part of the specifications for all COMP 40 homework, including this one.*

We emphasize again that pair programming is **required** for all COMP 40 assignments. **It is your responsibility to understand and abide by the COMP 40 [policies on pair programming](#).** *If you have not been assigned a partner or do not have permission to work separately (such permission is very rare), please contact the instructor immediately!* You will get no credit if you work alone unless you have permission.

## Part A: Brightness of a grayscale image

Please write a C program `brightness` that prints the average brightness of a grayscale image. Every pixel in a grayscale image has a brightness between 0 and 1, where 0 is black and 1 is as bright as possible.

## Specification

Print to standard output a single newline-terminated line containing the *average* brightness of the supplied image. The brightness should be printed in decimal notation with exactly one digit before the decimal point and three digits after the decimal point.

The program takes at most one argument:

- If an argument is given, it should be the name of a portable graymap file (in `pgm` format).
- If no argument is given, `brightness` reads from standard input, which should contain a portable graymap.
- If more than one argument is given, `brightness` halts with an error (see below).
- If a portable graymap is promised but not delivered, `brightness` halts with an error (see below).
- Upon successful completion, your program must terminate with an exit code of `EXIT_SUCCESS` (from `stdlib.h`). This is true of all programs you write in this course unless otherwise specified.

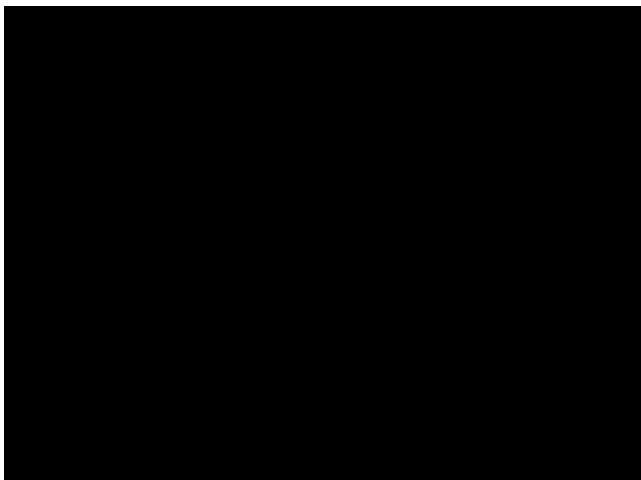
Where the specification above requires that you halt with an error, you have two choices:

1. You may print an error message on `stderr` and exit your program using `exit(EXIT_FAILURE)` from `stdlib.h`.
2. You may exit with a failed assertion or other Hanson-style Checked Runtime Exception (see explanation of CREs in the Hanson book)

When exiting due to such an error you must not produce any output on `stdout`.

## Examples of Brightness in Use

- Here are two photos:



Black cat in coal cellar



Polar bear in snowstorm

If `cellar.pgm` is a picture of a black cat in a coal cellar at midnight, and if `bear.jpg` is a picture of a polar bear in an snowstorm, then output should look something like this:

```
sunfire33{noah}: brightness cellar.pgm
0.000
sunfire33{noah}: djpeg -grayscale bear.jpg | brightness
0.972
sunfire33{noah}:
```

The first example takes its input from a file named on the command line, and the second example takes its input from standard input, as part of a Unix *pipeline*.

My solution to this problem takes fewer than 35 lines of C code.

## Help with image files

We provide [code to help you read image files](#); you will find the `Pnmrdr` interface in `/comp/40/include` and the implementation in `/comp/40/lib64`. If you use the supplied `Makefile`, these should be found automatically when your code needs them. Creating a `Pnmrdr_T` will read the graymap header for you, and from the header you can compute how many pixels are in the image. (You should read exactly as many pixels as are there—no more, no fewer.) Don't forget that the brightness of each pixel is represented as a *scaled integer*, as described in the `Pnmrdr` interface.



## Getting images

You can get images to play with by using one or more of the following programs:

- `djpeg` (use the `-grayscale` option)
- `pngtopnm`
- `pstopnm`
- `ppmtopgm`

## Problem analysis and advice

The main issues here are:

- In place of much of the C++ technique you already know, you have new C techniques to learn. The ideas are all similar, like old wine; C is a new bottle.
- You will have to read and understand the interface for [Pnm\\_rdr](#), and you will have to understand a little bit about the [pgm image format](#).
- You will have to do something appropriate if somebody hands you input that is *not* a portable graymap.

There is a subtlety here: we are asking you to use `Pnm_rdr` to read image files, and we are also sending a strong signal that your programs should be careful to check for and handle erroneous input. What if `Pnm_rdr` fails to detect an error in what is supposed to be an image file? Answer: for this assignment and others in COMP 40, don't worry about it. The whole purpose of using `Pnm_rdr` is to hide from you details of how the image file is represented on disk. Therefore, it's inappropriate to ask you to do more checking than `Pnm_rdr` does.

Of course, if this were a commercial program that might not be good enough: you would have to do due diligence to check that `Pnm_rdr` was indeed detecting all significant errors. For our purposes, just trust that it does.

- You are still responsible for making sure that when `Pnm_rdr` does detect an error in the image file, your program reports the error in an appropriate way. Maybe or maybe not this will involve writing nontrivial code.
- You will have to understand the compile-time and link-time options that gcc needs to work with the `Pnm_rdr` interface (`-I/comp/40/include` and `-L/comp/40/lib64 -lpnmrdr`). A good start would be to understand the `Makefile` you received for this assignment (and you can also look in [A simple introduction to Compile Scripts and Makefiles](#), which is for a previous version of this assignment, but is carefully explained there).
- To get the `Makefile` for *this* program, you must (if you didn't already):

```
git clone /comp/40/git/filesnpix
```

This will create a `filesnpix` directory with the compile script in it.

- You'll have to deal with **two representations** of real numbers between 0 and 1.

There is also an important issue of style:

- When using an enumeration literal such as `Pnm_rdr_gray`, refer to it by *name*, **not** by number.

## Part B: Read a line

Please create a single source file named `readaline.c`. Within that file you must:

- Include the header file `readaline.h` using `#include`. (The header file itself is provided for you in `/comp/40/include`, which is in the include path if you use the provided `Makefile`.)
- Implement a single function named `readaline`, conforming to the function declaration in the header file, which is:

```
size_t readaline(FILE *inputfp, char **datap);
```

We will separately test and grade the correctness of your implementation of this function. You will also use the function in Part C below, so problems with this function can also affect your grade on that. Test

carefully (and then test again!)

*Do not make a copy of `readaline.h` in your project directory!* Your submission will not be accepted if you do. If your compiles fail when including that header file there is almost surely something wrong with your `Makefile`.

## ***readaline specification***

The purpose of this function is to read a single line of input from file `inputfp`, which is presumed to have been opened for reading. As is common in specifications for computer programs and interfaces, we carefully define some terms, and then use those to specify the behavior of `readaline`:

- The term *character* refers to any of the 256 characters of [ISO Latin-1 extended ASCII](#). The bytes in the input file are interpreted as such *characters*.
- The characters comprising each file are grouped into zero or more *lines* as follows:
  - Each line contains at least one character
  - New lines begin at the start of the file, and after each newline character (`'\n'`)
  - Each newline character is included in the line that it terminates
- Each invocation of `readaline` retrieves the next unread line in the file. The characters comprising the line are placed into a contiguous array of bytes, and `*datap` is set to the address of the first byte. `readaline` returns the number of bytes in the line.
- The array of bytes is allocated by `readaline` using `malloc` (or the related allocation functions described in `man 3 malloc`.) It is the responsibility of the caller of `readaline` to free the array using `free`.
- `readaline` leaves the file seek pointer at the first (i. e., unread) character of the following line (if any) or at EOF
- If `readaline` is called when there are no more lines to be read, it sets `*datap` to `NULL` and returns 0.
- `readaline` terminates with a Checked Runtime Error in any of the following situations:
  1. Either or both of the supplied arguments is `NULL`
  2. An error is encountered reading from the file
  3. Memory allocation fails
- `readaline` must not cause memory leaks. That is, it must not leave allocated any dynamically acquired memory other than that returned to the caller through `*datap`.

For handling input lines you **MUST NOT** use the system library routines `getline` or `getdelim`, and you must not consult the man pages or other documentation relating to them. Reason: we want you to learn to do the work of reading input lines yourself.

## ***Partial credit***

For full credit, your `readaline` implementation must support input lines of any size. Significant partial credit is available if you support input files in which no line exceeds 200 characters in length, including any new line character. If you read a line that is longer than your implementation can handle, your `readaline` **MUST** write the message:

```
readaline: input line too long\n
```

to `stderr` and must immediately cause the program to exit with status code = 4 by calling the system function `exit(4)` (You learned to call `exit()` in Part A.)

## ***Hints***

A few hints:

- The `datap` parameter to `readaline` is a *call by reference (CBR)* parameter, i.e., it's used for function *output*.

Analogy: “Please slip the address of the party under my door.” In this case, there are two locations involved: the location of the party, and the location where we want location of the party to be stored (under the door).

Application: For this function, the caller wants access to the data in the next line. The `readaline` function will collect the data and store it somewhere (that's where the party is). The caller is asking `readaline` to store the location of the data in a location it has access to. That's what `*datap` refers to.

Draw a stack diagram. Stack diagrams are not analogies or “the way you think about how functions work”; they are precise descriptions of what *actually happens in the computer*. Therefore, you want to get these right, and there is a definite right way to draw these. Please ask the course staff for help, but try to draw it out first. Then you can explain your thinking to a member of the course staff, and we can applaud your insight or correct misunderstandings as appropriate.

- Under the covers, Hanson's `ALLOC` and `NEW` macros use `malloc`. So, you have a choice: if you are more comfortable using `malloc` and friends directly, go ahead. If you prefer the extra error-checking provided by Hanson's macros, you may use them. (FWIW: Norman Ramsey prefers Hanson's versions, Noah Mendelsohn is generally more comfortable with `malloc`. Either will get you full credit if used properly.)
- You will want to come up with some strategy for carefully testing your `readaline` function. Of course, you could just use it in your Part C `simlines` program and hope everything works, but we think you'll find that debugging is actually much harder that way. When something goes wrong (as it almost surely will), you will have to look everywhere to find the problem. Therefore: *we strongly urge you to come up with a strategy for carefully testing `readaline` by itself*.
- Question: according to the specification, is there *any* file you cannot read in its entirety by repeatedly calling your `readaline` function? Think about it. Why or why not? (You do not need to submit your answer.) Make sure your implementation can handle every file that the specification requires. Design your test cases accordingly.
- Handling lines of arbitrary size may be trickier than you think. DO NOT spend all your time trying to implement that at the cost of not getting to Part C. As noted above, you will get significant partial credit for both parts B and C if you can handle lines of at least 200 characters (you will use your `readaline` implementation in part C). We suggest you plan for longer lines, but for a start support only the 200 character minimum. Go on to part C and get that running. If you get all of that working, go back and extend `readaline` to handle longer lines. In principle, your Part C program should immediately start working with longer lines too!
- The `size_t` return type is an integer large enough to hold the number of bytes in the largest supported file on our Linux system. It is a standard type defined in `stddef.h` and used for this purpose by system library routines such as `fread`.

## Part C: Similarities in files

You are to write a program named `simlines`, the purpose of which is to read one or more data files to detect lines in the files that are similar to each other. Below we give you a detailed specification of what the program must do. Informally, it reads through one or more files looking for cases in which two or more lines contain the same words in the same order. It produces on standard output a report with a section for each such list of words, indicating which line numbers in which files contain those words.

### ***simlines specification overview***

The specification for `simlines` is as follows:

- The `simlines` program accepts zero or more command line arguments, each of which is the name of a file
- The program identifies all cases in which two or more lines in any of the input files are "similar" to each other, where similarity is determined by the definition below
- Similarities are reported if they occur within any single file and/or if lines in more than one of the files are similar
- The program writes its results to standard output in exactly the form specified for output below
- Upon successful completion, your program must terminate with an exit code of `EXIT_SUCCESS` (from `stdlib.h`). This is true of all programs you write in this course unless otherwise specified.
- The `simlines` program raises a Checked Runtime Error if any of the following occur:
  - Any one or more of the named input files cannot be opened
  - An error is encountered reading from any of the input files

- Memory cannot be allocated using `malloc`

If any of these situations arise, the program produces no output on `stdout`. The output on `stderr` must be only what is produced by the default Checked Runtime Error exception handler.

- If the same file is named more than once on the command line, then it is read and processed repeatedly, once for each command line reference. Note that this is very likely to result in similarities being detected and reported, since files are mostly similar to themselves!
- You **MUST** use your implementation of `readaline` from part B to read the data. If you are using a partial credit version of `readaline` that supports input lines of limited length, then you must rely on the error handling specified in part B to exit from `simlines` if an unsupported long input line is encountered.

## Definition of line similarity

The following gives the rules for determining whether two lines are *similar*:

- The terms *character* and *line* have the same definition as in the specifications for `readaline`
- A *word character* is any of lowercase 'a' &ndash; 'z'; uppercase 'A' &ndash; 'Z'; the digit characters '0' &ndash; '9', and the underscore character '\_' (which is ASCII code 95 decimal). All other characters are *non-word characters*
- Any contiguous grouping of word characters is a *word*. As common sense would suggest the term always refers to the largest possible groupings, thus the line:

```
'abc  def '
```

contains the two words `abc` and `def`, but not the words `ab`, `ef` or `a`.

- Thus any line contains zero or more words, optionally preceded, followed by and/or separated by non-word characters.
- **Two lines are similar if and only if they contain the same words in the same order.** Note that non-word characters are significant only as separators. Thus, the following lines are all similar to each other:

```
'abc  def '
'abc    def '
'  abc    def '
'abc,def '
```

- Lines containing no words (including empty lines) are not similar to other lines containing no words

## Simlines output specification

We define the term *match group* to refer to any set of two or more lines in the input file(s) that are *similar*. Such lines contain the same words in the same order, and we refer to that ordered list of words as the *matching words* list.

The output written to `stdout` by `simlines` must conform *exactly* to the following specification:

- There is one *output section* for each match group.
- If there is more than one match group, the corresponding output sections are separated from each other by a single additional newline (`\n`). The sections may appear in *any* order.
- If a match group contains  $n$  matches, then the corresponding output section consists of  $n + 1$  newline-terminated lines:
  - The first line of each output group contains the list of match words, separated from each other (if there is more than one) by a single space ' ' character.
  - For each match, a line consisting of:
    - The name of the file in which the match occurred, exactly as specified in the corresponding commandline argument. The name is printed left justified (i.e. starting at the beginning of the line), and padded to the right with blanks to fill a total of 20 characters. If the filename is longer than 20 characters, then it is written in its entirety, but with no additional padding.



- A single additional space (typically the 21st character, unless the file name was long)
  - A right justified seven digit integer giving the line number in which the match occurred. Line numbering is 1-based, i.e. the first line in the file is line 1, the second line 2, etc. No leading zeros are included in the line numbers.
- Within a match group the lines corresponding to each match are written in order. That is all matches in the file named by the first command line argument appear ahead of those in the files that follow. If there are multiple matches in a single file, they are written in order of increasing line number.
- As noted above, the same file named repeatedly on the command line results in no special processing: the file is read repeatedly, almost surely resulting in similar lines in the output, and those lines must be reported in order according to where the file references were in the command line. E.g. the command

```
simlines testfile otherfile testfile
```

might well result in an output group like this:

```
hello world
testfile-----2
otherfile-----7
testfile-----2
```

If `testfile` contained the words "hello world" on line 2, and `otherfile` contain the same words on line 7. (The light gray dashes in the sample output above are spaces in the actual output of `simlines`; the dashes are shown above to make the spaces easier to count.)

A consequence of the above rules is that `simlines` produces no output if there are no matches.

## Hints

Here are some hints that may help you in writing your `simlines` program:

- Writing that left-justified 20 character filename is easy if you know how to use `printf`. Look online for hints about printing fixed width left-justified strings or ask for help. If you do this right, all the formatting including correct handling of very long filenames, the proper right justification of the line number, etc. can easily be handled by a single `printf`.
- If you read carefully you will note that, except for allowing the output groups to be written in any order, the specification determines character by character exactly what your output must be
- **Be careful! If you do not exactly match the specification, you will not get credit.** Note that formatting errors tend to impact all of your output, and can easily result in failure of multiple tests or even of all tests ... read the specification and check your output carefully!
- There is no *additional* deduction for Part C if your `readaline` supports only lines of limited length, *as long as your `simlines` correctly exits when confronted with a long line that it cannot handle*. If your `simlines` does not exit in that way, then `simlines` will be graded on its handling of long as well as short lines.

## Performance target

Your `simlines` program should perform well on inputs resulting in at least 10000 match groups with up to hundreds of thousands of lines of input. (Due to limitations in the implementations of Hansons libraries, if you use the data structures we think likely, you may find performance slows dramatically if the number of match groups grows much larger than 10000.) By perform well, we mean completing in under 20 seconds or so on an unloaded server, if the output is redirected to a file or to `/dev/null` (If you look up `/dev/null` you'll find that it is a pseudo-file that throws away whatever is written to it... writing to that won't let you check your output, but it's a good way to time your program without waiting for hundreds of thousands of lines of output to be written to your display or even to a real file.) Of course, for shorter inputs of hundreds of lines with tens of matches, your program should respond more or less instantaneously.

## Problem analysis and advice

This problem boils down to simple string processing and standard data structures.

- The key to solving the `simlines` problem is to think very carefully about the data structures you will be creating and about *how* those data structures will result in a solution. Ask yourself questions like the following. Note that these overlap a lot with the information we ask you to provide us in your [Design Document](#)
  - What data structure(s) will you create to represent a match group?
  - How will you efficiently find out whether a given line from an input file belongs in a match group?
  - If it does, what information will you retain about the match?
  - For which of these are Hanson's datatype implementations such as List, Table, Set, Sequence, etc. useful, or when is it more appropriate to use C-language arrays, structs, etc.?
  - Which data, if any, should be converted to Hanson Atoms, and why?
  - When the time comes to write out a match group, how will you write the matches in order? (Hint: you should *not* try to sort them using a sort routine...there are much easier and more efficient ways if you use the right data structures. Be very careful to check the ordering rules for the various data structures you might be tempted to use!)
  - How is the memory for each of your data structures allocated? Do you have a way to free it to avoid memory leaks (remember, there is no way to free the memory for Atoms, and you will not be penalized for using Atoms... all other memory leaks must be avoided.)

Draw pictures! Take some interesting but small test cases and dry out in detail a picture of all your data structures and how they connect to each other. Once you have this picture absolutely clear and correct, organizing your code and then writing and debugging it will become much easier.

- We've said it before, and you'll hear us say it again and again in COMP 40: you will be tempted to put the majority of your time into designing and coding your program. You'll be tempted to build all or most of it, test the whole thing, and then try to find where the bugs are. *This will almost surely waste a lot of your time!* When you test all of your work together, the bugs could be almost anywhere. A bug in one routine might not cause an immediate crash, but might produce bad results because your program fails after executing hundreds of thousands of more lines of code.

The way to avoid this is to put as much energy into your test plan as you get into the design and coding of your solution. Find ways to test individual pieces of your code. Create test cases explore not just the obvious paths, but the unusual ones.

- C strings are different from C++ strings, and C's string library can be tricky to use properly. Make absolutely certain you understand the role of the NULL character '`\0`' in the termination of C strings. If you just code without understanding this, you are in for hours of frustration. Hanson offers some string processing libraries that you may (or might not) find helpful...my solution does not use them, but Norman Ramsey did use them in the solution to a somewhat similar problem given in years past.
- When it comes to string comparison, **what you know is wrong**. In C, writing

```
if (s1 == s2) { ... }
```

does *not* compare equality of two strings—it compares equality of *pointers*. To compare two strings for equality, you must write

```
if (strcmp(s1, s2) == 0) { ... }
```

Many old hands write this code more briefly:

```
if (!strcmp(s1, s2)) { ... }
```

but the briefer style requires a sharp eye for the exclamation mark and is based on a pun: `strcmp()` returns 0 when the strings match, and an `if` will regard that as false.

We prefer to make tests involving integers that don't actually stand for booleans explicit, so use the first form.

- Hanson's `Atom` interface maps equal strings to identical pointers, so pointer equality is OK on strings created with `Atom_string` or `Atom_new`. To use strings as keys or for similar purposes with Hanson's data structures, you **must** use the `Atom` interface. (When using a string as a

*value* (as opposed to a key), then use of Atom's is optional.)

- Note: **Hansons's Array data structure is *not* available for use in COMP 40.** As you will see in the next assignment, we have developed an improved (or at least more interesting) version that we will introduce then. You should not need to use Hansons's array for this first homework. (Of course, C character strings are C arrays, so you will definitely use C arrays when working with those.)
- Don't forget to run `valgrind` on your code

Repeat: **the data structures are already built for you**; your job is to figure out which ones will be useful. We are looking for a clean, straightforward design.

## Part C (DESIGN): Simlines Design

**DUE EARLY!** (see [course calendar](#))

### *Design overview*

The key to doing a good job on any program of significant complexity is to think thoroughly *in advance* two related questions:

1. How will data in the program be represented and interconnected?
2. How will the program logic be organized, and how will the computation be done?

For many programs, the first question is the more fundamental: you will find that if you have a sound and well considered approach to representing the data in your program, the logic and the program structure will easily follow. Thus, the heart of a software design lies in its representation of data.

There are several reasons why it's useful not just to think through, but also to write down the answers to the above questions *before* writing your program logic. These include:

- Writing down your design will tend to make you think it through more carefully. You will do a better job of discovering whether you have anticipated all the possibilities.
- Once your design is written down, it will be easy for you to have others review it (in the case of COMP 40, you will develop your design with your partner, and you may show it to teaching assistants and instructors, but to no one else)
- In COMP 40, we ask you to submit a design document early enough that we can give you some comments while your program is still under development. Often, we can warn you of design problems before they derail you.

It is crucial to cover not just the obvious main path considerations for your program; you need to consider all the different inputs and circumstances that your program has to handle properly.

### *Design document specifics*

Because documentation is such an essential part of COMP 40, we offer a number detailed guides on the course [referenced page](#). Included are some general guidance to preparing documentation, as well as specific guides for design documents for entire programs and for new abstract data types.

**For this assignment, we are making things simpler and asking you only for the following subset of the information we might normally expect for a new program. In the design document you submit you must:**

- Identify what data structures you will need to compute simlines and **what each data structure will contain**
- Hanson's data structures are *polymorphic*, so you will have to **explain what each `void *` pointer will point to**
- Please **describe the invariant** that will hold when your `simlines` program is partway through reading input lines.
- Based on these two explanations, **explain briefly how you will determine the simlines groups once all input has been read and explain how you will use your match group data structures to write the required output sections** .
- **Explain in detail your testing plan for both part B and part C** You'll get essentially no credit

for saying "I plan to try it with lots of inputs and see if it works". We need to know how you plan to test (e.g. `readaline` is a function, so you'll have to explain who calls it), and you should give us some sense of the range of test cases you'll be trying.

## Suggested order of work

This is a big assignment, and it's your first assignment. You will need to balance working in an orderly way with looking ahead so you can plan your time and overlap thinking about some of the harder problems in Part C with some of the development work for Part A and Part B. As described below, you will need to get far enough on the design of Part C to submit your design document on time.

The exact order in which you will want to do all this will depend a bit on what you already know and on your working style, but something like this will work for most of you:

- Read all of this instruction document several times. Try to make sure you understand what is being asked of you, and if you don't, ask on Piazza or talk to a TA in the labs.
- Make an annotated copy of these instructions that you and your partner can share (either on paper or online), noting things you don't understand, and indicating specific areas that will require research, design work, etc. **Note in your annotations everything that you are required to submit, and every specific case that you must cover.** We are not asking you to hand in your annotated copy, but you can use it as a valuable checklist as you work through the assignment.
- Make a tentative schedule and update it as things change. You can base it on the following if you like. Plan for trouble! Things will almost surely go wrong. If you aim to have your work done a day or two early, then the pressure will not be so great if you have a nasty bug.
- Your main early activities should include:
  - Learning all the technologies you will need to solve these problems so that you can make good design decisions. If necessary, experiment! Almost surely, you will want to write little test programs that open files, read data, and maybe create a Hanson data structure or two, just to make sure you know how. Your annotated instructions will be a good guide to what to learn first (e.g. image file formats and `Pnm_rdr` for your brightness program)
  - Setting up your development environment, as described above
- Then, get started designing [Part A: Brightness](#) (remember: both partners must be working together on all aspects of the design and implementation and testing ... you must not split the work!)
- *Before you implement brightness*, make a detailed test plan for it. Include in your schedule the time to develop the test cases and to do the actual testing.
- Start working on the implementation of [Part A: Brightness](#). Try to finish brightness in the first half of the time available for this assignment. Part B and especially Part C are more time-consuming, so leave plenty of time for them!
- As you start to turn the corner on brightness, or if you are getting short on time, begin preparing your design document for [Part C: Similarities in files](#). You probably won't make good progress on your first try, as you'll realize that you need to look more deeply into some of the Hanson interfaces and what they can do for you.
- Make a detailed test plan for [Part B \(Read a Line\)](#) and for [Part C: Similarities in files](#) (we need an outline of both test plans in your design document.)
- Be sure to submit your design document by the deadline (or up to two days late if you are using tokens — **if you use tokens for your design document you may not get feedback in time to be useful for your implementation**)
- Design, implement and thoroughly test *the partial credit version* of [Part B \(Read a Line\)](#). Do not rely on your `simlines` program as the only test framework for `readaline`.
- With luck, around this time you'll be getting back our comments on your Part C design document. Either way, as soon as you finish Part B, start refine your design for and then code `simlines` (part C).
- Test, test, and test some more.
- If you have time, go back and extend your `readaline` to handle long lines. Thoroughly unit test `readaline` and, only once it works reliably, relink your `simlines` with it. **Before you start modifying `readaline` to handle long lines, be sure to save a copy of the source for your simple implementation!** That way, if you get in trouble with the enhancement you'll still have



something that works. Without that, you will have neither a working Part B nor Part C!

- As you get ready to submit, read carefully the section on [Organizing and submitting your solutions](#).

To review: the [design document](#) covers the design of [part C](#) only and also testing plans for both [part B](#) and [part C](#). Your design document is due *before* the rest of the assignment. Parts A and B do not require a written design. The full assignment (parts A, B, and C) is due a few days later. Having working versions of all parts with only short line support in Part B will earn you more credit than a full version of Part B but missing or buggy Part C. See the [course calendar](#) for specific due dates.

## ***Deadlines and tokens***

Like all programming assignments for this class, the programming parts of this assignment are due one minute before midnight on the day indicated in the course calendar. You may turn in assignments *up to 48 hours after the due date*, which will cost you one or two [extension tokens](#). If you wish not to spend an extension token, then when midnight arrives submit whatever you have. We are very willing to give partial credit.

**If you spend an extension token on any part of the assignment, it automatically applies to *all* parts of the assignment.** I.e. submit either your design or your code a day late you use a token; submit them both a day late and it's still only one token total. Of course, if either is two days late, that's two tokens.

You may resubmit until the original deadline (i.e. not counting token extensions) and we will grade the latest version available at the deadline. *What you may not do is to submit before the deadline, then decide to use tokens and resubmit.* Reason: at the deadline we gather your work and start grading. If you resubmit, we will likely not notice, and would then have to regrade. If you are planning to use tokens, please do not submit. If you have an unusual problem, please email [comp40-staff@cs.tufts.edu](mailto:comp40-staff@cs.tufts.edu).

## **Organizing and submitting your solutions**

You will make two submissions to complete this assignment. Several days before the final submission, you will submit a design document. At the end of your work, you will make a final submission that includes your code.

Only one partner makes each submission, and the same partner should submit both the design and the final code submission. When you submit, you will identify your partner.

### ***Submitting your design document***

By the design deadline, submit a [design document](#) that explains what data structures you will use to write `simlines` and how you will use them.

Your document should be a plain ASCII document called `DESIGN` and formatted to fit in 80 columns. If you prefer to use a word processor you can submit `design.pdf`.

To submit your design, change to the directory that contains the file and run the `submit40-filesnpix-design` command. For this to work, you will have had to run `use comp40` either by hand or in your `.cshrc` or your `.profile`.

### ***Before submitting your code***

The purpose of these exercises is not just to teach you some tricky programming, it's to familiarize you with certain techniques, technologies, and approaches to programming.

**Before submitting, *each* member of your team should go back and reread the [purpose statement](#) for this assignment.**

Has *each of you* truly achieved these goals? Pair programming is terrific, but you need to stop once in awhile to make sure that every member of your team is getting the experience needed to move ahead. Your next steps in COMP 40 depend on your being solid in these basics. It's a lot to learn in a week: if you're still unsure, talk to a TA or schedule time with the instructor.



## Submitting your completed work

In your final submission, don't forget to include a README file which

- Identifies you and your programming partner by name and
- Acknowledges help you may have received from or collaborative work you may have undertaken with classmates, programming partners, course staff, or others
- Identifies what has been correctly implemented and what has not
- Says **approximately how many hours you have spent** completing the assignment

Your final submission should include at least these files:

```
README
brightness.c
simlines.c
```

A carefully designed, modular solution for `simlines` will probably include at least two other files.

When you get everything working, `cd` into the directory you are submitting and type `submit40-filesnpix` to submit your work.

## General advice for new C programmers

You will find some very helpful introductions to C on the [COMP 40 course references page](#). Some of the primary differences between C and C++ include:

- C has structures but not classes: a C programmer must find good ways to manage together the definition of data structures and the code that manipulates those structures, but C will not make the connection for you.
- I/O in C is different from C++, and generally a somewhat lower level. Although C has a few simple facilities for reading in numbers etc., most of the input and output you do will be at the level of individual characters or bytes. Instead of doing input and output using the `<iostream>` C++ library with its convenient `<<` operator and the `cout` and `cin` streams, you will use the `<stdio.h>` C library, with its less convenient `printf`, `fgetc` and `fgets` functions, etc. and the `stdout` and `stdin` file handles. (In C++, the operators `<<` and `>>` are overloaded at multiple types, but C does not support operator overloading, so another approach is needed.)
- Don't forget to initialize each C pointer variable, and if necessary, the data to which it points. In C++, an object may be initialized implicitly by a constructor, so a construction like:

```
Car *myCar_p = new Car;
```

might allocate space for a `Car` object, *initialize all the car member data*, and set the `myCar_p` pointer.

In C, you should use Hanson's `NEW` or `malloc` from the C library to allocate space for your data and to return a pointer:

```
struct Car *myCar_p;

/* Hanson */
NEW(myCar_p);      /* allocate space for a Car struct and set myCar_p

/* C library: one of the following */
myCar_p = malloc(sizeof(struct Car)); /* obvious but less foolproof
myCar_p = malloc(sizeof(*myCar_p));   /* better in most cases
```

Crucially, *none of the C constructions shown above actually initializes members of the `Car` structure*; they just allocate the needed space. You will see in Hanson's book good techniques (e.g. functions like `Table_new`) for writing code that combines allocation with initialization to achieve roughly what C++ gives you.

Of course, if you want to avoid memory leaks, you will eventually want to match each Hanson `NEW` with a `FREE`, and each `malloc` with a `free`. *One excellent technique that professionals use is to write the cleanup code at the same time as the allocation code.* So, typically, each time you

code a `malloc` you code a matching `free` somewhere else. Very often the latter will be in the termination code that runs just before your program exits.

One more detail to keep in mind: C++ and Hanson's library each have support for *Exceptions*, and those can result when C++ `new` or Hanson `NEW` fail to find the needed space. C itself does not have exceptions; if space can't be found then `malloc` returns `NULL`. *You should always check after each `malloc` to ensure that it worked!*

- You might wish to take three minutes to view, or review, [Pointer Fun with Binky](#), and compare the code you see there with Hanson's `NEW` macro.
- **Use the [C programming idioms](#)** that we suggest for all good C programs!