

Massachusetts Institute of Technology
Department of Electrical Engineering and Computer Science

6.035, Fall 2006

Handout 3 — Project Overview

Wednesday, September 6

This is an overview of the course project and how we'll grade it. You should not expect to understand all the technical terms, since we haven't yet covered them in class. We're handing it out today to give you some idea of the kind of project we're assigning, and to let you know the various due dates. Handout 5 describes the technical details of the project.

The class will be partitioned into groups of three students. You will be allowed to choose your own partners as much as possible. Each group is to write, in Java, a compiler for a simple programming language. We expect all groups to complete all phases successfully. The start of the class is very fast-paced: do not fall behind!

Important Project Dates

Thursday, September 7	Scanner and Parser Project assigned
Tuesday, September 19	Scanner and Parser Project due Semantic Checker assigned
Thursday, October 5	Semantic Checker due Code Generation assigned
Thursday, October 19	Code Generation Design and Checkpoint due
Thursday, October 26	Code Generation due Data Flow Analysis assigned
Thursday, November 2	Data Flow Analysis Design due
Thursday, November 9	Data Flow Analysis due Optimizer Project Assigned
Thursday, November 30	Optimizer Design due
Monday, December 11	Optimizer Project due
Wednesday, December 13	Compiler Derby

The Project Segments

Descriptions of the five parts of the compiler follow in the order that you will build them.

Scanner and Parser

This part scans the input stream (the program), and encodes it in a form suitable for the remainder of the compiler. You will need to decide exactly what you want the set of tokens to be and create the regular expressions for the scanner generator. The convention for this partitioning is quite standard in practice.

We'll supply a scanner generator (JFlex). This will consist of a program that takes a specification of the set of token types and outputs a Java program, the scanner. This specification uses regular definitions to describe which lexical tokens, i.e., character sequences, are mapped to which token types. The resulting scanner processes the input source code by interpreting the DFA (Deterministic Finite Automaton) that corresponds to the regular definition.

The parser checks the syntactic correctness of the token stream generated by the scanner. We'll supply a parser generator (Java CUP). This will consist of a primitive table-driven parser and a program that converts an LALR(1) grammar into a parse table suitable for use by that parser. You will have to transform the reference grammar into an LALR(1) grammar.

Semantic Checker

This part checks that various non-context free constraints, e.g., type compatibility, are observed. We'll supply a complete list of the checks. It also builds a symbol table in which the type and location of each identifier is kept. The experience from past years suggests that many groups underestimate the time required to complete the static semantic checker, so you should pay special attention to this deadline.

It is important that you build the symbol table, since you won't be able to build the code generator without it. However, the completeness of the checking will not have a major impact on subsequent stages of the project. At the end of this project the front-end of your compiler is complete and you have designed the intermediate representation (IR) that will be used by the rest of the compiler.

Code Generation

THIS IS A VERY TIME-CONSUMING PROJECT. For example, in previous years, it has been done in two parts. **YOU NEED TO START EARLY!**

In this assignment you will create a working compiler by generating unoptimized x86-64 assembly code from the intermediate format you generated in the previous assignment. Because you have relatively little time for this project you should concentrate on correctness and leave any optimization hacks out, no matter how simple.

The steps of code generation are as follows: first, the rich semantics of Decaf are broken-down into a simple intermediate representation. For example, constructs such as loops and conditionals are expanded to code segments with simple comparison and branch instructions. Next, the intermediate representation is matched with the Application Binary Interface, i.e., the calling convention and register usage. Then, the corresponding x86-64 machine code is generated. Finally, the code, data structures, and storage are laid-out in the assembly format. We will provide a description of the object language. The object code created using this interface will then be run on a testing machine (more on the testing machines soon).

This phase requires a Project Design Document and a Project Checkpoint, which is due one week prior to the deadline (see the Important Project Dates Section). The group has to provide two parts. First, a design document describing your design. This will be reviewed by the TA and feedback will be provided. This document will also count towards the project grade. Second, the group has to submit a checkpoint of the implementation. The checkpoint exists to strongly encourage you to start working on the project early. If you get your project working at the end, the checkpoint will have little effect. However, if your group is unable to complete the project, the checkpoint submission has a critical role in your grade. If we determine that your group did not do a substantial amount of work before the checkpoint, you will be severely penalized.

Data Flow Analysis

This assignment phase consists of building a data-flow framework to help optimize the code generated by your compiler. For this phase, you are required to implement the data-flow framework and a single data-flow optimization pass to test the framework. This framework will be used in the Optimizer project to build data-flow optimization passes.

We will provide a description of the framework and the required optimization to be implemented in a later handout.

This assignment also has a project design requirement (see the Important Project Dates Section).

Optimizer

The final project is a substantial open-ended project. In this project your team's task is to generate optimized code for programs so that they will be correctly executed in the shortest possible time.

There are multitude of different optimizations you can implement to improve the generated code. You can perform data-flow optimizations such as constant propagation, common sub-expression elimination, copy propagation, loop invariant code motion, unreachable code elimination, dead code elimination and many others using the framework created in the previous segment. You can also implement instruction scheduling, register allocation, peephole optimizations and even parallelization across the cores of the target architecture.

In order to identify and prioritize optimizations, you will be provided with a benchmark suite of three simple applications. Your task is to analyze these programs, perhaps hand optimizing these programs, to identify which optimizations will have the highest performance impact. Your write-up needs to clearly describe the process you went through to identify the optimizations you implemented and justify them.

The last class will be the "Compiler Derby" at which your group will compete against other groups to identify the compiler that produces the fastest code. The application used for the Derby will be provided to the groups one day before the Derby. This is done in order for your group to debug the compiler and get it working on this program. However, you are forbidden from adding any application-specific hacks to make this specific program run faster.

Grading

Make sure you understand this section so you won't be penalized for trivial oversights. The entire project is worth 70% of your 6.035 grade. The grade is divided between the segments in the following breakdown:

Scanner-Parser	5%
Semantic Checker	10%
Code Generator	15%
Data-flow Analyzer	10%
Optimizer	30%

The remaining 30% comes from three quizzes, each worth 10%.

The first 4 phases of the project (Scanner-Parser, Semantic Checking, Code Generation, and Data-flow Analysis) will be graded as follows:

- (25%) Design and Documentation (subjective). Your score will be based on the quality of your design, clarity of your design documentation, and incisiveness of your discussion on design alternatives and issues. Some parts of the project require additional documentation. Always read the *What to Hand In* section. For the segments that require a Project Design Document, 10% will be assigned to the design document and 15% to the final write-up.
- (75%) Implementation (objective). Points will be awarded for passing specific test cases. Each project will include specific instructions for how your program should execute and what the output should be. If you have good reasons for doing something differently, consult your TA first.
 - Public Tests (25%)
 - Hidden Tests (50%)

The Optimizer project phase will be graded differently:

- (20%) Design and Documentation, with particular attention given to your description of the optimization selection process.
- (40%) Implementation. As each group implements different optimizations, the only public test is the generation of correct results for the benchmark suite and the Derby program (10%). The hidden tests will check for overtly optimistic optimizations and incorrect handling of programs (30%).
- (40%) Derby Performance. The formula for translating the running time of the program compiled by your compiler into a grade will be announced later.

All members of a group will receive the same grade on each part of the project unless a problem arises, in which case you should contact your TA as soon as possible.

What To Hand In

For each part of the project, you will have to provide us with two things: online sources and hardcopy documentation.

<code>-o <outname></code>	Write output to <code><outname></code>
<code>-target <stage></code>	<code><stage></code> is one of scan, parse, inter, or assembly. Compilation should proceed to the given stage.
<code>-opt [optimization...]</code>	Perform the listed optimizations. <code>all</code> stands for all supported optimizations. <code>-<optimization></code> removes optimizations from the list.
<code>-debug</code>	Print debugging information. If this option is not given, there should be no output to the screen on successful compilation.

Table 1: Compiler Command-line Arguments

Online Source

Create a new directory with the project name (e.g., scanner-parser) at the top-level of your group locker. Copy all relevant source files to this directory. You are required to provide an Apache Ant build-file that compiles your sources and packages the resulting *class* files into a java archive. The *jar* file will be used by the TAs to run the various tests.

Files in your submission directory should not be modified after the deadline for submitting the project (5:00 p.m. on the due date). Please consult additional project handouts for phase-specific requirements.

Command-line Interface

Your compiler should have the following command line interface.

```
javac -jar Compiler.jar [option | filename...]
```

The command line arguments you must implement are listed in table 1. Exactly one filename should be provided, and it should not begin with a '-'. The filename must not be listed after the -opt flag, since it will be assumed to be an optimization.

The default behavior is to compile as far as the current assignment of the project and produce a file with the extension changed based on either the target or ".out" if the target is unspecified.

By default, no optimizations are performed. The list of optimization names will be provided in the optimization assignments.

We have provided a class, CLI, which is sufficient to implement this interface. It also returns a Vector of arguments it did not understand which can be used to add features. The TAs will not use any extra features you add for grading. However, you can tell us which, if any, to use for the compiler derby. You may wish to provide a "-O" flag, which turns on the optimizations you like.

Hardcopy Documentation

Documentation should be turned in to the course secretary by 5:00 p.m. on the due date. It should be clear, concise and readable. Fancy formatting is not necessary; plain text is perfectly acceptable. You are welcome to do something more extravagant, but it will not help your grade.

Your documentation must include the following parts:

1. A brief description of how your group divided the work. This will not affect your grade; it will be used to alert the TAs to any possible problems.
2. A list of any clarifications, assumptions, or additions to the problem assigned. The project specifications are fairly broad and leave many of the decisions to you. This is an aspect of real software engineering. If you think major clarifications are necessary, consult your TA.
3. An overview of your design, an analysis of design alternatives you considered, and key design decisions. Be sure to document and justify all design decisions you make. Any decision accompanied by a convincing argument will be accepted. If you realize there are flaws or deficiencies in your design late in the implementation process, discuss those flaws and how you would have done things differently. Also include any changes you made to previous parts and why they were necessary.
4. A brief description of interesting implementation issues. This should include any non-trivial algorithms, techniques, and data structures. It should also include any insights you discovered during this phase of the project.
5. Only for scanner and parser: A listing of commented source code relevant to this part of the project. For later stages, the TAs can print out source code if needed. Do not resubmit source code from other parts of the project, even if minor changes have been made.
6. A list of known problems with your project, and as much as you know about the cause. If your project fails a provided test case, but you are unable to fix the problem, describe your understanding of the problem. If you discover problems in your project in your own testing that you are unable to fix, but are not exposed by the provided test cases, describe the problem as specifically as possible and as much as you can about its cause. If this causes your project to fail hidden test cases, you may still be able to receive some credit for considering the problem. If this problem is not revealed by the hidden test cases, then you will not be penalized for it. It is to your advantage to describe any known problems with your project; of course, it is even better to fix them.

It is entirely up to you to determine how to test your project. The thoroughness of your testing will be reflected in your performance on the hidden test cases.