

# 601.465/665 — Natural Language Processing

## Homework 7: Finite-State Programming\*

Prof. Jason Eisner — Fall 2022  
Due date: Friday 9 December, 11pm

This homework is a tutorial that walks you through some exercises in constructing finite-state machines. All of the finite-state algorithms are provided for you via the OpenFST toolkit. Hence no conventional programming will be needed, just regular expressions and command-line tools. You'll construct a variety of FSMs, primarily using the Thrax regular expression compiler and the Ngram utility for backoff language modeling.

**Homework goals:** First, to get some experience in writing regular expressions for unweighted and weighted languages and relations. Second, to get a feel for the range of problems that can be easily handled in this declarative framework. After some initial practice, you'll get to try applying FSTs in various NLP situations, some of which are new to you: syntactic chunking, pronunciations from a lexicon, language models, word segmentation, character-level edit distance, and noisy-channel decoding.

**Collaboration:** *You may work in pairs on this homework*, as it is rather long: it may be more fun to work through the exercises with a friend. That is, if you choose, you may collaborate with one partner from the class, handing in a single homework with multiple names on it. However:

1. You are expected to do the work *together*, not divide it up: your solutions should emerge from collaborative real-time discussions with the whole group present.
2. Your README file should describe at the top what each of you contributed, so that we know you shared the work fairly.
3. Your partner for homework 7 can't be the same as your partner from homework 6 (tagging).

In any case, observe **academic integrity** and never claim any work by third parties as your own.

**Reading:** There is no separate reading this time. Instead, we'll give you information and instructions as you work through the homework.



**What to hand in (via Gradescope):** As usual, you should submit a README.pdf file with answers to all the questions in the text. We'll also ask you to submit a .zip archive of all the grammar files you create:

File	Questions	Should contain
binary.grm	2, 3, 9	First, Second, Disagreements, Triplets, NotPillars, Oddlets, WFlip, WeightedMultipath
rewrite.grm	4, 5	Cross, BitFlip1, BitFlip2, Parity1, Parity2, Parity3, UnParity, Split
chunker.grm	6a	NP, MakeNmod, TransformNP, BracketTransform, BracketResults
stress.grm	7	Stress (commented), StressWords
dactyls.grm	8d	Results (modified)
rhyme.grm	8	Ending, WordEnding
noisy.grm	10, 11, 12	CompleteWord, DelSpaces, SpellText, Generate, RandomWord, Spell (revised), PrintText (revised)

\*Many thanks to Frank Ferraro, who co-wrote this homework and wrote the accompanying scripts.

## Software:

- **OpenFST** is a very efficient C++ toolkit for building and manipulating semiring-weighted FSMs. You can use the C++ API to directly specify states, arcs, and weights, or to combine existing FSMs through operations like union and composition. You can also store these FSMs in `.fst` files and manipulate them with command-line utilities like `fstunion` and `fstcompose`.

A symbol table (`.sym` file, or part of some `.fst` files) specifies the integerized internal representation of the upper or lower alphabet. E.g., the integers 1, 2, 3, ... might internally represent the letters a, b, c, ... or perhaps the words `aardvark`, `aback`, `abacus`, .... OpenFST uses these integers to label arcs in its data structures and file format ( $\epsilon$  arcs are labeled with 0). It is only the symbol table that reveals what a given FSM's integer labels are supposed to *mean*.

However, we will not be using OpenFST directly (nor its Python interface, **Pynini**). Instead, we will use two packages that provide a fairly friendly interface to OpenFST:

- **Thrax** is an extended regular expression language that you can use to define collections of finite-state machines. A Thrax grammar can be created with any text editor and is stored in a `.grm` file. The Thrax compiler *compiles* this into a `.far` file—an “fst archive” that contains multiple named OpenFST machines and symbol tables.
- Since regular expressions are not good at specifying the topology of  $n$ -gram models, there's also the **NGram** toolkit, which builds a  $n$ -gram backoff language model from a corpus. It supports many types of smoothing. The resulting language model is represented as a weighted FSA in OpenFST format.

### 1. First, get set up!

- (a) The `nlp-class` conda environment has been updated to include the OpenFST, Thrax, and NGram packages. If you are running on your own machine rather than `ugrad`, then you should get that update by typing

```
conda env update -f nlp-class.yml
```

using the latest version of `nlp-class.yml` from the course website.

- (b) To view drawings of FSMs, download and install `graphviz`:

<http://www.graphviz.org/Download.php>

On Linux systems, you can just do `sudo apt-get install graphviz`.

- (c) Download a copy of the homework directory `hw-ofst`, either from the `ugrad` network<sup>1</sup> or the web.<sup>2</sup>
- (d) Look in the `hw-ofst` directory.<sup>1</sup> Our scripts are in the `bin` subdirectory, which you should probably add to your `PATH` so that you can execute these scripts without saying where they live. Run the following command<sup>3</sup> (and maybe put it in your `~/.bashrc` so that it will be executed automatically next time you log in).

```
export PATH=${PATH}:/usr/local/data/cs465/hw-ofst/bin
```

---

<sup>1</sup>`/usr/local/data/cs465/hw-ofst/`. You can copy this directory to your local machine. Or just use it in place: we suggest making a symlink to it in your working directory, to give yourself easy access to the files.

<sup>2</sup><http://cs.jhu.edu/~jason/465/hw-ofst>

<sup>3</sup>This command assumes you're using the nice `bash` shell. If it doesn't work, switch shells (temporarily by typing `bash` or permanently by typing `chsh`). Or if you really want to keep using `tcsh`, modify the command accordingly.

- (e) We've given you a script `grmtest` to help streamline the compilation and testing of Thrax code. Its usage is:

```
grmtest <grm file> <transducer_name> [max number output lines]
```

This script compiles the specified `.grm` file into a `.far` file (using a makefile produced by `thraxmakedep`), and then passes the standard input through the input through the exported FST named by `<transducer_name>`. You'll get to try it out below.<sup>4</sup>

*Warning:* If the output string is the empty string  $\epsilon$ , then for some reason `grmtest` skips printing it. This seems to be a bug in `thraxrewritetester`, which `grmtest` calls. Just be aware of it.

2. Now get to know Thrax. We highly recommend looking through the online manual<sup>5</sup> and perhaps the commented examples that come with Thrax.<sup>6</sup> The following tutorial leads you through some of the basic FSM operations you can do in Thrax.

- (a) Let's first define some simple FSMs over a binary alphabet. Type the following declarations into a new file `binary.grm`.

```
Zero = "0";
One = "1";
Bit = Zero | One;
export First = Optimize[Zero Zero* Bit* One One One One?];
```

This defines four named FSMs using *Thrax's regular expression syntax*.<sup>7</sup> Each definition ends in a semicolon. The first and second FSMs accept only the strings 0 and 1, respectively. The third defines our entire alphabet, and hence accepts either 0 or 1. The fourth accepts some subset of `Bit*`.

We can compile this collection of named FSMs into a `fst` archive (a `.far` file"). More precisely, the archive provides only the FSMs that have been marked with an `export` declaration; so here `Zero`, `One`, and `Bit` are just intermediate variables that help define the exported FSM `First`.

- i. Try compiling and running it using our `grmtest` script:

```
$ grmtest binary.grm First
[compiler messages appear here]
Input string: [type your input here]
```

The FSA `First` is interpreted as the identity FST on the corresponding language. So entering an input string will transduce it to itself if it is in that language, and otherwise will fail to transduce. Type `Ctrl-D` to quit.

You'll get an error if you try running `grmtest binary.grm Zero`, because `Zero` wasn't exported.

- ii. What language does `First` accept (describe it in English)? Why are 0 and 1 quoted in the

<sup>4</sup>It will help to have a good shell such as `bash` that lets you recall and edit previous command lines. It will also help to have a good terminal program that lets you scroll up to see earlier parts of the input. `gnome-terminal` does this, as does the Emacs shell mode (`ESC x shell`).

<sup>5</sup><http://www.openfst.org/twiki/bin/view/GRM/ThraxQuickTour>

<sup>6</sup>`/usr/local/share/thrax/grammars/` on the ugrad or grad machines.

<sup>7</sup>Whereas XFST defines many special infix operators, Thrax instead writes most operators (and all user-defined functions) using the standard form `Function[arguments]`. Thrax uses square brackets `[]` for these function calls, and parentheses `()` for grouping. Optionality is denoted with `?` and composition with `@`. There is apparently no way to write a wildcard that means "any symbol"—you need to define `Sigma = "a"|"b"|"c"|...` and then you can use `Sigma` within other regular expressions.

.grm file?

- iii. Let's get information about **First**. First, we need to extract the FSA **First** from the FST archive:<sup>8</sup>

```
$ far2fst binary.far First
```

Now use the **fstinfo** shell command<sup>9</sup> to analyze **First.fst**:

```
$ fstinfo First.fst
```

Look over this output: how many states are there? How many arcs?

- iv. Optionally, look at a drawing of **First** (as an identity transducer over the alphabet  $\{0, 1\}$ ):

```
$ fstview First.fst
```

Note that **fstview** is a wrapper script that we are providing for you.<sup>10</sup> The picture will take a few seconds to appear if the graphics pixels are being sent over a remote X connection.<sup>11</sup>

- (b) Now let's look at equivalent ways of describing the same language.

- i. Can you find a more concise way of defining **First**'s language? Add it to **binary.grm** as a new regexp **Second**, of the form

```
export Second = Optimize[ ...fill something in here... ];
```

Run **grmtest** to check that **First** and **Second** seem to behave the same on some inputs.

- ii. Here's how to check that **First** and **Second** really act the same on *all possible inputs*—that they define the same language:

```
export Disagreements = Optimize[ (First - Second) | (Second - First) ];
```

If **First** and **Second** are equivalent, then what strings should **Disagreements** accept?

To check that, run **fstinfo** on **Disagreements.fst**. From the output, can you conclude that **First** and **Second** must be equivalent?

*Note:* The **fstequal** utility is another way to check:

```
if fstequal First.fst Second.fst; then echo same; else echo different; fi
```

One way to program **fstequal** would be to construct the **Disagreements** FSA.

- (c) You might have wondered about those **Optimize[ ... ]** functions. The Thrax documentation notes that **Optimize**

...involves a combination of removing epsilon arcs, summing arc weights, and **determinizing** and **minimizing** the machine ... [Details are [here](#).]

To find out what difference that made, make a new file **binary-unopt.grm** that is a copy of **binary.grm** with the **Optimize** functions *removed*. Then try:

```
grmtest binary-unopt.grm First    # and type Ctrl-D to exit
far2fst binary-unopt.far
fstview First.fst Second.fst Disagreements.fst
```

<sup>8</sup>Use **far2fst binary.far** to extract *all* exported FSTs, or **far2fst binary.far First Second** to extract multiple ones that you specify.

<sup>9</sup><http://man.sourcentral.org/f14/1+fstinfo>

<sup>10</sup>After printing **fstinfo**, it calls **fstdraw** to produce a logical description of the drawing, then makes the drawing using the Graphviz package's **dot** command, and finally displays the drawing using **evince**. Each of these commands has many tweakable options. What if you're running on your own machine and don't have **evince**? Then edit the **fstview** script to use a different PDF viewer such as **xreader**, **atril**, **xpdf**, or **acroread**.

<sup>11</sup>If you start getting "can't open display" errors, then try connecting via **ssh -Y** instead of **ssh -X**. An alternative is to copy the (small) **.pdf** file produced to your local machine and use your local PDF viewer.

Questions:

- i. Although **First** and **Second** may be equal, their unoptimized FSMs have different sizes and different topologies, reflecting the different regular expressions that they were compiled from. How big is each one?
- ii. The drawing of the unoptimized **Disagreements.fst** shows that it immediately branches at the start state into two separate sub-FSAs. Why? (*Hint*: Look back at the regexp that defined **Disagreements**.)
- iii. Now test some sample inputs with  
`grmtest binary-unopt.grm First`  
How are the results different from the optimized version? Why?

- (d) You may not want to call **Optimize** on every machine or regular sub-expression. The documentation offers the following warning:

When using composition, it is often a good idea to call **Optimize[]** on the arguments; some compositions can be massively sped up via argument optimization. However, calling **Optimize[]** across the board (which one can do via the flag `--optimize_all_fsts`) often results in redundant work and can slow down compilation speeds on the whole. Judicious use of optimization is a bit of a black art.

If you optimize **Disagreements** *without* first optimizing **First** and **Second**, what do you get and why?

3. Now try some slightly harder examples. Extend your **binary.grm** to also export FSAs for the following languages. (You are welcome to define helper FSAs beyond these.)

- (a) **Triplets**: Binary strings where 1 only occurs in groups of three or more, such as 000000 and 0011100001110001111111.
- (b) **NotPillars**: All binary strings *except* for even-length strings of 1's:  $\epsilon$ , 11, 1111, 111111, 11111111, ... (These correspond to binary numbers of the form  $2^{2k} - 1$  written in standard form.) Some strings that *are* in this language are 0, 1, 000011, 111, 0101, 011110.
- (c) **Oddlets**: Binary strings where 1's only appear in groups of odd length. Careful testing this one! (Note that 0000 is in this language, because 1's don't appear in it at all.)

You will extend **binary.grm** further in question 9.

4. So far we have only constructed FSAs. But Thrax also allows FSTs. Create a new file **rewrite.grm** in which you will define some FSTs for this question and the next question.

Complicated FSTs can be built up from simpler ones by concatenation, union, composition, and so on. But where do you get some FSTs to start with? You need the built-in `:` operator:

**input : output**

which gives the cross-product of the input and output languages.

In addition, any FSA also behaves as the identity FST on its language.

Place the following definition into **rewrite.grm**:

```
export Cross = "a" (("b":"x")* | ("c"+ : "y"*) | ("":"fric")) "a";
```

Note that "" denotes the empty string  $\epsilon$ .

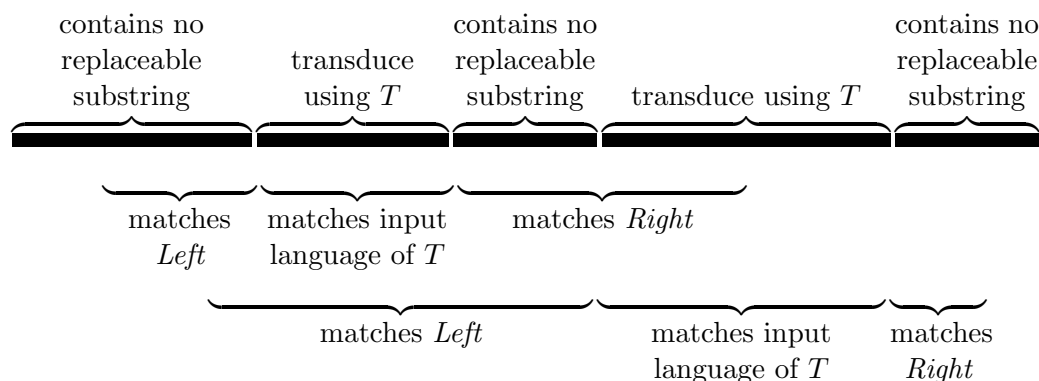
- (a) What is the input language of this relation (answer with an ordinary regexp)?
- (b) Give inputs that are mapped by **Cross** to 0 outputs, 1 output, 2 outputs, and more than 2 outputs.
- (c) How would you describe the **Cross** relation in English? (You do not have to hand in an answer for this, but at least think about it.)
- (d) Make an **Optimized** version of **Cross** and look at it with **fstview**. Is it consistent with your answers above? How many states and arcs?

Check your answers by transducing some inputs, using the following commands:

```
grmtest rewrite.grm Cross
grmtest rewrite.grm Cross 3
```

The second version of the command limits the number of outputs that are printed for each input that you enter.

5. If you have a simple FST,  $T$ , then you can make a more complicated one using **context-dependent rewriting**. Thrax's **CDRewrite** operator is similar to the  $\rightarrow$  operator in XFST. It applies  $T$  “everywhere it can” within the input string, until there are no unreplaced substrings that could have been replaced. The following shows schematically how two substrings of **■■■■■** might be replaced:



The braces underneath the string show that each of the 2 replaced substrings appears in an appropriate context—immediately between some substring that matches *Left* and some substring that matches *Right*. The 2 replaced substrings are not allowed to overlap, but the contexts can overlap with the replaced substrings and with one another, as shown in the picture.<sup>12</sup>

If you want to require that the *maximal* substring to the left matches *Left*, then start *Left* with the special symbol [BOS], which can only match at the beginning of the string. Similarly for *Right* and [EOS] (end of string).

<sup>12</sup>If you are wondering how this is accomplished, have a look at [Mohri & Sproat \(1996\)](#), section 3.1.

The above example shows only one way of dividing up the input string into regions that are transduced by  $T$  and regions that are left alone. If there are other ways of dividing up this input string, then the rewrite transducer will try them too—so it will map this input to multiple outputs.<sup>13</sup>

`CDRewrite[T, Left, Right, Any, Dir, Force]` specifies a rewrite transducer with arguments

- $T$  : any FST
- $Left, Right$  : unweighted FSAs describing the left and right contexts in which  $T$  should be applied. They may contain the special symbols "[BOS]" and "[EOS]", respectively. (These symbols are only to be used when describing contexts, as in these arguments to `CDRewrite`, which interprets them specially. They do *not* appear in the symbol table.)
- $Any$  : a minimized FSA for  $\Sigma^*$ , where  $\Sigma$  is the alphabet of input and output characters. The FST produced by `CDRewrite` will only allow input or output strings that are in  $Any$ , so be careful!
- $Dir$  : the direction of replacement.
  - 'sim' specifies "simultaneous transduction":  $Left$  and  $Right$  are matched against the original input string. So all the substrings to replace are identified first, and then they are all transduced in parallel.
  - 'ltr' says to perform replacements "in left-to-right order." A substring should be replaced if  $Left$  matches its left context *after* any replacements to the left have been done, and  $Right$  matches its right context *before* any replacements to the right have been done.
  - 'rtl' uses "right-to-left" order, the other way around.
- $Force$  : how aggressive to be in replacement?
  - 'obl' ("obligatory," like  $\rightarrow$  in XFST) says that the *unreplaced* regions may not contain any more replaceable substrings, as illustrated above. That is, they may not contain a substring that matches the input language of  $T$  and which falls between substrings that match  $Left$  and  $Right$ .
  - 'opt' ("optional," like  $\rightarrow$  in XFST) says it's okay to leave replaceable substrings unreplaced. Since the rewrite transducer has the freedom to replace them or not, it typically has even more outputs per input.

Define the following FSTs in your `rewrite.grm`, and test them out with `grmtest`:

- (a) **BitFlip1**: Given a string of bits, changes every 1 to 0 and every 0 to 1. This is called the "1's complement" of the input string. Define this *without* using `CDRewrite`.
- (b) **BitFlip2**: Like **BitFlip1**, but now it should work on any string of digits (e.g., transducing 1123401 to 0023410). Define this version using `CDRewrite`.  
*Hint:* The  $Any$  argument in this case should be `Digit*` where

```
Bit = "0" | "1";
Digit = "0" | "1" | "2" | "3" | "4" | "5" | "6" | "7" | "8" | "9";
```

- (c) **Parity1**: Transduces even binary numbers to 0 and odd binary numbers to 1. Write this one *without* using `CDRewrite`.  
 It's always good to think through the unusual cases. Some people think the empty string  $\epsilon$  is a valid binary number (representing 0), while others don't. What does your transducer think?

<sup>13</sup>So there is no notion here of selecting the regions to transduce in some deterministic way, such as the left-to-right longest match used by the XFST  $\textcircled{\rightarrow}$  operator.



(d) **Parity2**: Same thing as **Parity1**, but use the **Reverse** function in your definition. So start by writing a transducer that keeps the *first* bit of its input instead of the *last* bit.

(e) **Parity3**: Same thing as **Parity1**, but this use **CDRewrite** in your definition.

What does this transducer think about  $\epsilon$ ?

*Hint*: You may find it helpful to use **[BOS]** and **[EOS]**, or the composition operator **@**.

(f) **UnParity**: Define this as **Invert[Parity3]**. What does it do?

(g) **Split**: Split up a binary string by nondeterministically inserting spaces into the middle, so that input 0011 maps to the eight outputs

$\{0011, 001\ 1, 00\ 11, 00\ 1\ 1, 0\ 011, 0\ 01\ 1, 0\ 0\ 11, 0\ 0\ 1\ 1\}$

*Hint*: Use **CDRewrite["": " ", ...]** and figure out the other arguments.

(h) **Extra credit: SplitThree**: Similar to **Split**, but always splits the input string into exactly three (nonempty) binary strings. This will produce multiple outputs for strings longer than 3 bits, and no outputs for strings shorter than 3 bits.

*Hint*: Compose **Split** with something else. The composition operator is **@**.

6. Bit strings are great, but let's move on now to natural language. You know from Homework 1 that precisely describing syntax can be challenging, and from Homework 4 that recovering the full parse of a sentence can be slow. So, what if you just want a quick finite-state method for finding simple NPs in a sentence? This could be useful for indexing text for search engines, or as a preprocessing step that speeds up subsequent parsing. Or it could be part of a cascade of FSTs that do information extraction (e.g., if you want to know who manufactures what in this world, you could scan the web for phrases like “*X manufactures Y*” where *X* and *Y* are NPs).

Identifying interesting substrings of a sentence is called “chunking.” It is simpler than parsing because the “chunks” are not nested recursively. This tutorial question will lead you through building an FST that does simple NP chunking.

We will assume that the input sentence has already been tagged (perhaps by a tagging FST, which you could compose with this chunking FST). You'll build the following objects:

- An FSA that accepts simple noun phrases: an optional determiner, followed by zero or more adjectives **Adj**, followed by one or more nouns **Noun**. This will match a “base” NP such as **the ghastly orange tie**, or **Mexico**—though not the *recursive* NP **the ghastly orange tie from Mexico that I never wear**.

The regexp defining this FSA is a kind of simple grammar. To make things slightly more interesting, we suppose that the input has two kinds of determiners: quantifiers (e.g., **every**) are tagged with **Quant** whereas articles (e.g., **the**) are tagged with **Art**

- A transducer that matches exactly the same input as the previous regular expression, and outputs a *transformed* version where non-final **Noun** tags are replaced by **Nmod** (“nominal modifier”) tags. For example, it would map the input **Adj Noun Noun Noun** deterministically to **Adj Nmod Nmod Noun** (as in **delicious peanut butter filling**). It would map the input **Adj** to no outputs at all, since that input is not a noun phrase and therefore does not allow even one accepting path.
- A transducer that reads an arbitrary input string and outputs a single version where all the *maximal* noun phrases (chosen greedily from left to right) have been transformed as above and bracketed.



- (a) You'll be editing the provided file `chunker.grm`:

```
import 'byte.grm' as bytelib;
import 'tags.grm' as tags;
Sigma = (tags.Tags) | (bytelib.kBytes);
SigmaStar = Optimize[Sigma*];
```

Copy this file along with `byte.grm` and `tags.grm` from the `grammars/` directory. Line 2 defines `tags` to be the collection of FSMs that are exported by `tags.grm`. Expressions like `tags.Tags` in line 3 then refer to individual FSMs in that collection. You should look at these other files referenced by lines 1–2. Now:

- i. Define an FSA NP that accepts an optional article (**Art**) or quantifier (**Quant**); followed by an arbitrary number of adjectives (**Adj**); followed by at least one noun (**Noun**). We would like to write:

```
export NP = Optimize[(Art|Quant)? Adj* Noun+];
```

What goes wrong? (*Hint*: look at importable FSMs from `tags`.) Fix the definition in `chunker.grm`, and in your `README`, provide evidence of what you were able to accept.

You will use the fixed `chunker.grm` for the rest of this question.

(*Note*: Really we should be working over the alphabet of tags rather than the default alphabet of ASCII characters. Later in the homework we'll see how to define other alphabets using *symbol tables*.)

- ii. Have a look at NP:

```
far2fst chunker.far NP
fstview NP.fst
```

How many states and arcs are there? Comment on the structure of the machine.

- (b) In a noun-noun compound, such as *the old meat packing district*, the nouns *meat* and *packing* act as *nominal modifiers*. Define and try out a transducer `MakeNmod`, using `CDRewrite`, that replaces `Noun` with `Nmod` immediately before any `Noun`. So `ArtAdjNounNounNoun` as in the example becomes `ArtAdjNmodNmodNoun`.

To define `MakeNmod`, you'll need to figure out what arguments to use to `CDRewrite`.

- (c) Now define an FST

```
export TransformNP = Optimize[NP @ MakeNmod];
```

- i. Describe in words what this composition is doing.  
ii. What are the results on `ArtAdjNounNounNoun` and `AdjNounNounNounVerb`?  
iii. What is the size of `TransformNP` compared to `MakeNmod`?  
iv. How does the topology of `TransformNP` differ from that of `NP`?

- (d) This FST transduces a noun phrase to one that has <angle brackets> around it:

```
export BracketNP = (" : "<") NP (" : ">");
```

Here the `NP` language is being interpreted as the identity relation on that language, and concatenated with two other simple regular relations. So `BracketNP` reads  $\epsilon$ , any `NP`,  $\epsilon$  and writes `<`, the same `NP`, `>`.

What, if anything, is the difference between the following?

```
export Brackets1 = Optimize[SigmaStar (BracketNP SigmaStar)*];
export Brackets2 = CDRewrite[BracketNP, "", "", SigmaStar,'sim','obl'];
```

Try them out on short and long strings, such as `ArtAdj`, `AdjNoun`, and `VerbArtAdjNounNounNounVerbPrepNoun`.

- (e) Now define `BracketTransform` to be like `Brackets2`, except that it should not only bracket noun phrases but also apply the transformation defined by `TransformNP` within each noun phrase. This should be a fairly simple change.
- (f) One interesting thing about FSTs is that you can pass many strings through an FST at once. Define `BracketResults` to be the regular language that you get by applying `BracketTransform` to *all* strings of the form `Quant Noun+ Verb` at once.  
(*Hint*: Check out the `Project` operator in the Thrax documentation.<sup>14</sup> You may want to optimize the result.

You can check the FSA by using `fstview` on `BracketResults` (note that it may be drawn as an identity FST). To print out the strings it accepts (up to a limit), run `grmtest` on the cross-product machine `":BracketResults`, and enter an empty input string to see all the outputs.

- (g) **Extra credit:** To get a sense of how `CDRewrite` might do its work, define your own version of `TransformNP` that does *not* use `CDRewrite`. It should be a composition of three FSTs:
- *Optionally* replace each `Noun` with `Nmod`, without using `CDRewrite`. This will transduce a single input to many outputs.
  - *Check* that no `Noun` is followed by another `Noun` or `Nmod`. This filters outputs where you didn't replace enough nouns.
  - *Check* that every `Nmod` is followed by a `Noun` or another `Nmod`. This filters outputs where you replaced too many nouns. (*Hint*: It is similar to the XFST “restrict” operator that we defined in class.)

Call your version `TransformNP2`.

24

7. Let's learn a little bit about pronunciation, which will require some trickier rewriting. We will define

```
import 'byte.grm' as bytelib; # copy it from grammars/byte.grm
Letter = bytelib.kAlpha;      # kAlpha is defined in byte.grm
Sigma = Letter | "'";        # alphabet also includes stress mark
Vowel = "a" | "e" | "i" | "o" | "u" | "A" | "E" | "I" | "O" | "U";
Consonant = Letter - Vowel;
Nucleus = Vowel+;
```

Syllables have the form `Consonant* Nucleus Consonant*`. The word `roundabout` consists of three syllables, `round-a-bout`, whose nuclei are `ou`, `a`, `ou`.

For this problem, we won't have to worry about the complete rules for dividing a word into syllables (e.g., why did the letters `nd` both belong to the first syllable of `roundabout`?). We will only care about

<sup>14</sup>This operator is used to “project” a relation onto the upper or lower language, like the `.u` and `.l` operators in XFST. Why is that called projection? Consider a set of points on the plane:  $\{(x_1, y_1), (x_2, y_2), \dots\}$ . The projection of this set onto the  $x$  axis is  $\{x_1, x_2, \dots\}$  and the projection onto the  $y$  axis is  $\{y_1, y_2, \dots\}$ . Same thing when projecting a regular relation, except that each  $(x_i, y_i)$  pair is a pair of strings.

the nuclei. Let's assume for this problem that adjacent vowels always belong to the same nucleus, as in **roundabout** or **tree**. (This is an oversimplification since in truth, adjacent vowels are sometimes the nuclei of adjacent syllables, as in **cadmium** or **pancreatic**.)

A conceptual flaw in the definitions above is that vowels are actually *sounds*, not letters. In English, spelling is complicated: a letter does not reliably stand for a particular sound. The rules above will therefore not find enough nuclei in **happy** or **rhythm**, and it will find too many in **borehole**. But let's pretend that they're right.

Most spoken languages use both stressed and unstressed syllables. Stressed (or “accented”) syllables tend to pronounce their nuclei louder and longer than adjacent unstressed syllables. Also, when the intonation system assigns a “melody” to a sentence, it places the most emphatic high and low pitches onto stressed syllables.

Stressed and unstressed syllables tend to roughly alternate with each other, but the exact pattern depends on the language. We explain the pattern in a language by saying that stress is mostly assigned by rules—which we can write as FSTs. A first approximation to the English stress system is this.<sup>15</sup>

Visit the nuclei from left to right. Add stress to a nucleus if (1) it is not stressed yet, and (2) we did not just visit a stressed nucleus, and (3) we are about to visit an unstressed nucleus.

Read those rules carefully! Here are some results of applying them. We indicate stress by placing the mark ' before the nucleus. (Try saying the words aloud.)

(2 syllables)	m'ountain	spl'endid
(3 syllables)	c'annonball	r'oundabout
(4 syllables)	M'issis'ippi	p'arab'olic
(5 syllables)	'oper'ational	'ophthalm'ology
(6 syllables)	t'inntinn'abul'ation	'incomm'unic'ado

A final step makes sure that 1-syllable words are stressed:<sup>16</sup>

(4) Add a stress if there isn't one anywhere in the word.

(1 syllable) j'ump t'oad

Now, many words might not seem to fit the pattern (say them aloud):

(2 syllables)	contr'ol	destr'oy
(3 syllables)	c'abar'et	sem'antics
(4 syllables)	K'alamaz'oo	carn'ivorous
(5 syllables)	p'eripat'etic	ev'apor'ating

But these exceptions still *partly* fit the pattern. We can describe them by saying that the underlined stress marks    were already present in the lexicon—that is, an English speaker had *memorized* that those syllables had to be stressed. Any remaining stress marks were still filled in automatically by

---

<sup>15</sup>This is not the whole story. Things get more complicated because of morphology. In a compound word like **bl'ackb'oard**, **f'antasyl'and**, or **v'ideot'aping**, the rules apply to the two halves separately. Certain prefixes and suffixes similarly act as if they were separate words. Other suffixes cause the syllable *before* them to be stressed. Other prefixes and suffixes are glued onto a word *after* its stress has been assigned, but then the rules are applied again to possibly fill in additional stresses on the new, compound word.

<sup>16</sup>There are exceptions for little function words such as **the**, **a**, **of**, **with**. We'll ignore that here.

applying the rules above. Notice that the nucleus before `_` never gets stressed by the rules, because of condition (3).

- (a) Copy `grammars/stress.grm` into your working directory. Add a definition of an FST `Stress` that adds stress marks according to the rules above. The input may or may not contain some lexical stress marks. For example, it should map `operational` to `'oper'ational` and either `Kalamaz'oo` or `K'alamaz'oo` to `K'alamaz'oo`. Comment your Thrax code to explain how you are solving the problem!

Your `Stress` transducer is really quite supercalifragilisticexpialidocious,<sup>17</sup> and you can have some fun trying it out on additional words. You will also find it interesting to view the FST. The drawing may be frighteningly large at first, but to make it small and comprehensible, use `Optimize` and try shrinking the alphabet to just one vowel and one consonant, by viewing `Optimize[("a"|"b")* @ Stress]` instead.

*Hint* on how to define `Stress`: Use a single `CDRewrite` with the `'ltr'` option to implement (1)–(3). Once you've got that working, compose with a separate step that handles (4). The hardest part is (2), and there are several ways to solve it. You can try to translate the English statement directly into a regular expression, or you can come up with an equivalent description. One strategy is to use an FST to massage the input string into something easier to work with, perhaps by marking the starts of nuclei, or prepending an extra syllable so that you don't have to worry about `[BOS]`; you'd later delete this extra material with another FST.

- (b) Run the transducer backwards using `Invert[Stress]`. What possible inputs from the lexicon could give rise to the observed pronunciation `ev'apor'ating`? How about `'incomm'unic'ado`?
- (c) Now that you can assign stress to a word, try assigning stress to a sentence. The domain of `Stress` is `Sigma*` (where `Sigma = Letter | "'`). Make a new transducer that will replace all words in a sentence with their stressed forms:

```
export StressWords = CDRewrite[Stress, ...];
```

`StressWords` needs to use the larger alphabet `bytelib.kBytes`, which consists of all characters including spaces and punctuation.

*Note:* This problem would be easier with directed replacement, but Thrax doesn't have that yet. It's still not too hard.

- (d) **Extra credit:** Make an improved FST `Stressy` that is like `Stress` but properly handles the letter `y` (including `Y`), which is sometimes a vowel. You will have to look at examples of words containing `y` to figure out a rule for when it is and isn't a vowel. An obscure example is the city of `Ypsilanti` (where it's a vowel).

You may want to transduce the input so that `y` is split into two symbols (perhaps `y` for consonant and `^` for vowel), and similarly `Y`; then use your old `Stress`; then convert the `y`'s back again.

Explain in your `README` what your rule is, and comment your `stress.grm` code, of course.

8. Now let's make a rhyming dictionary. Two words rhyme if they have the same "rhyming endings." The rhyming ending of a word is the suffix that starts with its last stressed nucleus.

In this problem (unlike the `Stress` problem) we are really going to work with sound rather than spelling. `flight`, `contrite`, and perhaps `leukocyte` all rhyme because they all have the same ending,

<sup>17</sup>Although for reasons discussed earlier, maybe you should spell that as `"...expiyala..."`

which sounds like 'ite. However, they do not rhyme with **bauxite** or **Semite**, whose endings sounds like 'auxite and 'emite. The difference is that they have stress on the next-to-last nucleus rather than the last nucleus.

- 27 (a) The *flight* vs. *site* examples above demonstrated that words can rhyme even when they spell their endings differently. Conversely, can you give an example of words that *don't* rhyme even when they spell their endings the same way?

We will represent the pronunciation of a word as a sequence of phonemes, where a phoneme represents an indivisible unit of “meaningful” sound. The government agency ARPA<sup>18</sup> defined an ASCII text notation for English phonemes, known as the “ARPAbet.” You can look at the ARPAbet symbol table in `symbol_tables/arpabet.sym`. Each symbol in the left column is designed to represent an English phoneme.<sup>19</sup>

But how do we map from English orthography (spelling) to phonemes? Rather than write a complicated FST, we'll just look it up in the CMU pronunciation dictionary: `data/cmudict.txt`.<sup>20</sup> Take a look at this two-column, tab-separated file:

```
academic      AE2 K AH0 D EH1 M IH0 K
academy       AH0 K AE1 D AH0 M IY0
...
```

Below we'll use the `StringFile` function, which creates a finite relation from a file like this one. Each line of the file has the format

*input string* <tab character> *output string*

and the relation (FST) is the union of all these *input* : *output* pairs. In this case, each input is an English word and the output is its pronunciation. Notice again that there's not a simple mapping between letters and phonemes. There can be more letters than phonemes, one letter can map to many phonemes, and many letters can map to one phoneme. There are some errors in the file, but it's pretty good and let's assume it's correct.

Each vowel includes a stress level: 0 for unstressed syllables, and 1 or 2 for stressed syllables, where 1 indicates the most strongly stressed syllable in the word.

- (b) The words **academic** and **academy** contain four tokens of the letter **a**. How many of these are pronounced differently? You do not have to hand in a written answer for this, but think carefully about it. (Say the words repeatedly, focusing on what your tongue and mouth are doing. Is the tip of the tongue near the top or bottom of your mouth? How open is your mouth? What's happening with your lips? Are they pulled back, or relaxed?)
- (c) Warmup: The following Thrax grammar, `grammars/dactyls.grm`, illustrates a very useful way to define regular *languages* with the help of regular *relations*. What regular language is defined by **Results** below? Why?

<sup>18</sup>The Defense Advanced Projects Research Agency, now known as DARPA.

<sup>19</sup>Technically speaking, the ARPAbet vowel symbols represent entire nuclei, which may consist of 1 or 2 phonemes in English. For example, the nucleus of **bite** is technically a pair of vowels (roughly, “ah” morphing into “ee”—try saying **bite** very slowly if you don't believe this!), but the ARPAbet represents this *diphthong* with the single symbol **AY**.

<sup>20</sup>The smaller version in `data/cmumini.txt` may be faster for getting your code working.

```

import 'byte.grm' as bytelib;
Sigma = bytelib.kGraph | bytelib.kSpace;

Pronounce = StringFile['data/cmudict.txt'];
StressPattern = CDRewrite[(Sigma-bytelib.kDigit) : "", "", "", Sigma*, 'sim', 'obl'];
Dacytl = ("1" | "2") "0" "0"; # what's a dacytl? look it up!

export Results = Optimize[Project[Pronounce @ StressPattern @ (Dacytl*), 'input']];
export ListResults = "" : Results;

```

*Hint:* You can see lots of strings from the `Results` language by running

```
grmtest dactyls.grm ListResults 100000
```

and giving the empty input string.

*Hint:* The general technique illustrated here is to say that  $x$  is in our language  $L$  iff a conveniently modified version of  $x$ —as transduced by  $T$ —would fall in some other language  $L'$  that is easier to define. We then throw away that result of the transduction as well as  $T$  and  $L'$ . We don't care about them for their own sake, only as helpers in defining  $L$ :

$$L \stackrel{\text{def}}{=} \{x : T(x) \cap L' \neq \emptyset\}$$

This type of definition is not available in standard Unix regular expression syntax, which doesn't support transduction.

The previous grammar was a bit silly (and slow to compile), because it treated the pronunciation AE2 K AH0 D EH1 M IH0 K as a sequence of 23 *bytes* (ASCII characters), including spaces. We really want to treat it as a sequence of 8 *phonemes*.

This should be possible because regular relations are just subsets of  $\Sigma^* \times \Delta^*$ , where  $\Sigma$  and  $\Delta$  are the input and output alphabets. (This generalizes functions from  $\Sigma^* \rightarrow \Delta^*$ .) Here, we would like to define a new  $\Delta$ , namely the ARPAbet.

OpenFST defines a new alphabet using a “symbol table.” The ARPAbet is defined in `symbol_tables/arpabet.sym`, which defines the ARPAbet phonemes to be the integers 1 through 70 (internally, just as ASCII characters are 1 through 255) and gives them string names for I/O purposes.

When Thrax encounters a string like AE2 K AH0 D EH1 M IH0 K, you need to tell it whether to interpret that string as describing a sequence of ASCII characters (“byte mode”), a sequence of Unicode characters (“utf8 mode”), or a sequence of symbols from some other specified alphabet (“symbol mode”). It may encounter such strings in three settings:

- in a `.txt` file (like `cmudict.txt`)      Specify mode with extra arguments to `StringFile`.
- in a `.grm` file (between double quotes)      Specify mode using a modifier on the quoted string.
- in the input you enter into `grmtest`      Use `grmtest-with-symbols` and specify the mode via extra command-line arguments. These are passed on to `thraxrewrite-tester`.

- (d) Modify your copy of `dactyls.grm` to accomplish the same thing as before, but using the actual ARPAbet alphabet. You will need to do

```
arpa = SymbolTable['arpabet.sym'];
PronounceArpa = StringFile['data/cmudict.txt', byte, arpa];
```

so that the second column of `cmudict.txt` will be parsed using `arpabet.sym`. You will also find it useful to do

```
import 'arpabet.grm' as arpabet;
```

so that you can use regexps defined therein, such as `arpabet.Stressed`. You'll notice when reading `arpabet.grm` that quoted strings over that alphabet must be followed by `.arpa`. You will have to do the same in your grammar, even in the case of the empty string of phonemes, which is written as `"".arpa` and has a different type (in the sense of programming language types) from the empty string of bytes `""`.

Now it's time to make our rhyming dictionary.

- (e) Create a new file `rhyme.grm`. Start by include some definitions from the previous question, such as `arpa`, `Pronounce` (the second version), and `arpabet`.

Define an FST `Ending` that maps a phoneme string to its rhyming ending. Both input and output are strings over the ARPAbet alphabet—so to test out this FST, you will have to use a variant of `grmtest`:

```
grmtest-with-symbols <grm file> <transducer>
                    <input symbol file>
                    <output symbol file>
                    [max number output lines]
```

This has two extra arguments—the input and output symbol files. So you will do

```
grmtest-with-symbols rhyme.grm Ending arpabet.sym arpabet.sym
```

- (f) Now define an FST `WordEnding` that maps an ASCII word to the rhyming ending of its pronunciation. It should map `academic` (8 letters) to `EH1 M IH0 K` (4 phonemes). Probably a good idea to `Optimize` this machine. What are its domain and range?

You can test it with

```
grmtest-with-symbols rhyme.grm WordEnding byte arpabet.sym
```

where the special argument `byte` says to interpret the input string in byte mode, rather than using a symbol table from a file. This means that the input symbols are individual characters and do not have to be separated by spaces.

- (g) What does `WordEnding @ Invert[WordEnding]` describe?  
What are its input and output alphabets?
- (h) What happens when you try to build the FSM for `WordEnding @ Invert[WordEnding]`? Why? Be specific about what is happening in the machine composition.
- (i) Luckily, we're not stuck. Instead of using `grmtest` to pass an input word `w` through the composition above, i.e.,

```
w @ (WordEnding @ Invert[WordEnding])
```

you can pass it through one transducer at a time, i.e.,

29

30

31

32



`(w @ WordEnding) @ Invert[WordEnding]`

which is equivalent since composition is an associative operator. Why is this more efficient?

You can do this via `grmtest`, by listing a pipeline of FSTs that the input string will pass through in sequence:

`grmtest rhyme.grm WordEnding,InvWordEnding`

where you have defined `InvWordEnding` appropriately. Try it out and give some interesting input/output examples. Does the dictionary have a rhyme for `orange`? For `adventureland`?

- (j) **Extra credit:** The word `academic` is considered to be a *trivial* rhyme for `academic`, `nonacademic`, or `pandemic` because *too much* of the endings match: Most poets avoid such rhymes because they sound too repetitive (whereas a non-trivial rhyme is just repetitive enough to be pleasing!).

Come up with a good rule for detecting trivial rhymes; describe and justify it in your `README`.

Can you give an FST that will transduce its input word to all nontrivial rhymes for that word? This is harder than it might seem because the difference operator `-` is not defined on regular relations. (The difference of two regular relations is not always regular.)

9. In OpenFST, you can define weighted FSMs. By default, OpenFST, NGram and Thrax all use the tropical semiring,  $\langle \mathbb{W}, \oplus, \otimes, \odot, \oslash \rangle = \langle \mathbb{R} \cup \{\pm\infty\}, \min, +, \infty, 0 \rangle$ .<sup>21</sup> Thus, weights can be interpreted as costs (e.g., negated log-probabilities). Concatenating two paths or regexps will *sum* their costs, whereas if a string is accepted along two parallel paths or by a union of regexps, then it gets the *minimum* of their costs.

Augment an FSM's definition by appending a weight  $w$  in angle brackets, `<` and `>`, and wrapping the entire definition in parentheses.

- (a) i. What is the minimum-weight string accepted by this FSA, and what is the weight of that string? (Remember that parentheses in Thrax just represent grouping, not optionality.)  
`(Zero <1>) (Bit+ <0.2>) (One <0.5>)`
- ii. What is the minimum-weight pair of strings accepted by this FST, and what is the weight of that pair?  
`(Zero : One <1>) (Bit+ <0.2>) (One : One One <0.75>)`
- (b) In your old `binary.grm` file, define a weighted transducer `WFlip` that accepts the language of the above FSA and, reading left to right:

- Nondeterministically flips the leftmost bit. Flipping has weight 2, while not flipping has weight 1.
- In the `Bit+` portion, replaces every 0 with 01 (at cost 0.5), and replaces every 1 with 0 (at cost 0.4).
- Accepts the final 1 bit with weight 0.5.

Don't use `CDRewrite` here.

For example, `WFlip` should produce the following:

Input string: 0011

Output string: 00101 <cost=2.4>

Output string: 10101 <cost=3.4>

<sup>21</sup>And currently a portion of the Thrax we're using supports only this semiring.

(c) Now let's consider cases where we aggregate the weights of multiple paths using  $\oplus$ .

- i. In your `README`, name any two binary strings  $x, y$ : for example,  $(x, y) = (00, 1)$ . In `binary.grm`, define `WeightedMultipath` to be a simple *weighted* FST of your choice, such that the particular pair  $(x, y)$  that you named will be accepted along at least two different paths, of different weights.

To confirm that these two accepting paths exist, view a drawing of the machine, and use `grmtest` to find out what  $x$  maps to. What are the weights of these paths?

- ii. Now define `WeightedMultipathOpt = Optimize[WeightedMultipath]`. In this FST, how many paths accept  $(x, y)$  and what are their weights? Why?
- iii. Suppose  $T$  is an FST with weights in some semiring, and  $x$  and  $y$  are strings. So  $T$  accepts the pair  $(x, y)$  along 0 or more weighted paths.

Describe, in English, what the following weighted languages or relations tell you about  $T$ :

```
T_out      = Project[      T,      'output']; # erases input from T
xT_out     = Project[ x @ T,      'output']; # erases input x from x @ T
Ty_in      = Project[      T @ y, 'input'];  # erases output y from T @ y
xTy        =      x @ T @ y;
exTy       = ("":x) @ T @ (y:""); # erases input x & output y from x @ T @ y
```

```
xT_out_opt = Optimize[xT_out];
Ty_in_opt  = Optimize[Ty_in];
exTy_opt   = Optimize[exTy];
```

How big is the last FSM, in general? Why do the last three FSMs have practical importance?<sup>22</sup>

You can try these all out for the case where  $T$  is `WeightedMultipath` and  $x$  and  $y$  denote the strings  $(x, y)$  you named above.

- (d) **Extra credit:** Define an FSM `NoDet` that has no deterministic equivalent. (Unweighted FSAs can always be determinized, but it turns out that either outputs (FSTs) or weights can make determinization impossible in some cases that have cycles.)

How will you know you've succeeded? Because the Thrax compiler will run forever on the line `Determinize[RmEpsilon[NoDet]]`—the determinization step will be building up an infinitely large machine. (`RmEpsilon` eliminates  $\epsilon$  arcs from the FSM, which is the first step of determinization. Then `Determinize` handles the rest of the construction.)

10. Throughout the remainder of this homework, we'll be focused on building noisy-channel decoders, where weighted FSTs really shine. Your observed data  $y$  is assumed to be a distorted version of the "true" data  $x$ . We would like to "invert" the distortion as best we can, using Bayes' Theorem. The most likely value of  $x$  is

$$x^* \stackrel{\text{def}}{=} \underset{x}{\operatorname{argmax}} \Pr(x \mid y) = \underset{x}{\operatorname{argmax}} \Pr(x, y) = \underset{x}{\operatorname{argmax}} \underbrace{\Pr(x)}_{\text{"language model"}} \underbrace{\Pr(y \mid x)}_{\text{"channel model"}} \quad (1)$$

<sup>22</sup>In general, one might want to do one of these computations for many different  $x$  (or  $y$ ) values. Rather than compiling a new Thrax file for each  $x$  value, you could use other means to create  $x$  at runtime and combine it with  $T$ . For example, to work with `BracketTransform` from question 6e, try typing this pipeline at the command line: `echo "ArtNounNounNoun" | fstcompilestring | fstcompose - BracketTransform.fst | fstproject --project_output | fstoptimize | fstview`. (Other handy utilities discussed in this handout are `fstrandgen` for getting random paths, `fstshortestpath` for getting the lowest-cost paths, `farprintstrings` for printing strings from these paths, our `fstprintstring` script for printing the output of a single random path, and our `grmfilter` script for transducing a file.) Or you could just use OpenFST's C++ API.

In finite-state land, both language and channel models should be easy to represent. A channel is a weighted FST that can be defined with Thrax, while a language model is a weighted FSA that can be straightforwardly built with the NGram toolkit.<sup>23</sup> To make even easier to build a language model, we've given you a wrapper script, `make-lm`:

```
make-lm corpus.txt
```

By default, this will create a Kneser-Ney back-off trigram language model called `corpus.fst`. Every sentence of `corpus.txt` should be on its own line.<sup>24</sup>

- (a) Create the default language model for the provided file `data/entrain.txt`. Each line of this file represents an observed sentence from the English training data from the HMM homework.<sup>25</sup> Notice that the sentences have already been tokenized for you (this could be done by another FST, of course).

You should now have a language model `entrain.fst` in your working directory, along with two other files you'll need later: `entrain.alpha` is the alphabet of word types and `entrain.sym` is a symbol table that assigns internal numbers to them.

Look at the files, and try this:

```
wc -w entrain.txt      # number of training tokens
wc -l entrain.alpha    # number of training types
fstinfo entrain.fst    # size of the FSA language model
```

- (b) The NGram package contains a number of useful shell commands for using FST-based  $n$ -gram models. Three that you may find particularly interesting are `ngramprint`, `ngramrandgen`, and `ngramperplexity`. You can read up on all three if you like (use the `--help` flag).

Sampling several sentences from the language model is easy:

```
ngramrandgen --max_sents=5 entrain.fst | farprintstrings
```

Each `<epsilon>` represents a backoff decision in the FSM. You can see that there is a *lot* of backoff from 2 to 1 to 0 words of context. That's because this corpus is very small by NLP standards: the smoothing method realizes that very few of the possible words in each context have been seen yet.

To read the sentence more easily, use the flag `--remove_epsilon` to `ngramrandgen` (this prevents `<epsilon>`s from being printed). Alternatively, just use our `fstprintstring` script to print a random output from any FST:

```
fstprintstring entrain.fst
```

If you want to know the most probable path in the language model, you can use `fstshortestpath`, which runs the Viterbi algorithm.<sup>26</sup>

```
fstshortestpath entrain.fst | fstprintstring
```

How long are the strings in both cases, and why? What do you notice about backoff?

<sup>23</sup>Language models can't be concisely described with regular expressions: at least, not the standard ones.

<sup>24</sup>Otherwise you may get weird non-normalization errors for Kneser-Ney smoothing.

<sup>25</sup>Because the NGram toolkit assumes that sentence boundaries are marked by newlines, we've omitted `###`.

<sup>26</sup>In general `fstprintstring` will print a randomly chosen string, but in the output of `fstshortestpath`, there is only one string to choose from.

- (c) **Recommended:** Although it's a small language model, `entrain.fst` is far too large to view in its entirety. To see what a language model FSA looks like, try making a *tiny* corpus, `tiny-corpus.txt`: just 2–3 sentences of perhaps 5 words each. Try to reuse some words across sentences. Build a language model `tiny-corpus.fst` and then look at it with `fstview`. There is nothing to hand in for this question.
- (d) Now, let's actually use the `entrain.fst` language model. Copy `noisy.grm` from the `grammars/` directory:

```
import 'byte.grm' as bytelib;          # load a simple grammar (.grm)
export LM = LoadFst['entrain.fst'];    # load trigram language model (.fst)
vocab = SymbolTable['entrain.sym'];     # load model's symbol table (.sym)
```

Each line loads a different kind of external resource. In particular, the second line loads the trigram language model FSA, and the third line loads the symbol table used by that FSA. The symbol table consists of the vocabulary of the language model, as well as the OOV symbol `<unk>` (“unknown”).<sup>27</sup>

You can therefore use LM to transduce some strings:

```
grmtest-with-symbols noisy.grm LM entrain.sym entrain.sym
```

What is the result of transducing the following? Explain your answers. What are the domain and range of the relation LM?

- Andy cherished the barrels each house made .
- If only the reporters had been nice .
- Thank you

We now want to compose LM with a noisy channel FST. Because the language model is nothing more than an FSA, we can use it in Thrax. Of course, we're going to have to be careful about symbol tables: the noisy channel's input alphabet must be the same as LM's output alphabet.

Remember to be careful when creating FSTs over a nonstandard alphabet. If you write

```
("barrels barrels" : (" | "ship"))*;
```

then the input to this FST must be a multiple of 15 symbols in the default `byte` alphabet. But if you write

```
("barrels barrels".vocab : ("".vocab | "ship".vocab))*;
```

then Thrax will parse the quoted strings using the `vocab` symbol table. So here, the input must be an even number of symbols in the `vocab` alphabet. Writing `"Thank".vocab` will give an error because that word is not in the symbol table (it's not in the file `entrain.sym` from which `vocab` was loaded).

A noisy channel that “noises up” some text might modify the sequence of words (over the `vocab` alphabet) or the sequence of letters (over the `byte` alphabet). If you want to do the latter, you'll need

---

<sup>27</sup>The `make-lm` script takes the vocabulary to be all the words that appear in training data, except for a random subset of the singletons. These singletons are removed from the vocabulary to ensure that some training words will be OOV, allowing the smoother to estimate the probability of OOV in different contexts. Ideally the smoothing method would figure this out on its own.

to convert words to letters. Recall from question 8a that the `StringFile` function interprets its given tab-separated text file as an FST, with the domain and range as the second and third parameters, respectively. So we'll add the following line to `noisy.grm`:

```
Spell = Optimize[StringFile['entrain.alpha', vocab, byte]];
```

This maps a word to its spelling, just as `Pronounce` in 8a mapped a word to its pronunciation.

- (e) In `noisy.grm`, define a transducer called `CompleteWord` that could be used to help people enter text more quickly. The input should be the first few characters in a word, such as `barr`. Each output should be a word that starts with those characters, such as `barrel` or `barrage`.

Use LM to assign a cost to each word, so that each completed word is printed together with its cost. Is a word's cost equal to the unigram probability of the word, or something else?

*Hint:* Be careful to think about the input and output alphabets, and to pass them as arguments to `grmtest-with-symbols`. The input alphabet should be `byte` (not `byte.sym`), as explained in question 8f.

- (f) **Extra credit:** Now define `CompleteWordInContext` similarly. Here the input is a sequence—separated by spaces—of 0 or more complete words followed by a partial word. Each output is a single word that completes the partial word, as before. But this time the cost depends on the context: that's what language models are for.

Try it out and give a few example inputs that illustrate how the context can affect the ranking of the different completions.

*Hint:* You might not be able to get away with exporting `CompleteWordInContext` as a single transducer—it's rather big because it spells out the words in the language model. It will be more efficient to use the pipelining trick from question 8i. In your README, tell the graders what command to enter in order to try out your pipeline, and give the original `CompleteWordInContext` definition that your pipeline was derived from.

11. Question 10 defined our language model,  $\text{Pr}(\mathbf{x})$ . Now let's compose it with some channel models  $\text{Pr}(\mathbf{y} \mid \mathbf{x})$  that we'll define. In this question, we'll practice by working through a simple deterministic noisy channel.

- (a) Still working in `noisy.grm`, define a deterministic transducer `DelSpaces` that deletes all spaces. Define this using `CDRewrite`, and use the alphabet `bytelib.kGraph | bytelib.kSpace`. Using `grmtest` you should be able to replicate the following:

Input string: If only the reporter had been nice .

Output string: Ifonlythereporterhadbeennice.

Input string: Thank you .

Output string: Thankyou.

Input string: The reporter said to the city that everyone is killed .

Output string: Thereportersaidtothecitythateveryoneiskilled.

- (b) `grmtest` will transduce each string that you type in, providing multiple outputs when they exist. To transduce a whole file to a single output, once you've tested your transducer, we've provided another wrapper script `grmfilter`:

```
$ grmfilter
Usage:
    cat input.txt | grmfilter [-opts] <grammar file> <name1>,<name2>,...
-r: select a random path
-s: find shortest path (default)
-h: print this help message (and exit)
```

Just like `grmtest`, it takes two required arguments, a `.grm` file and a comma-separated list of FST names defined in that file. It reads strings from the standard input, one per line, and writes their transductions to the standard output. The output string comes from *one* of the paths that accept the input. The default (which can be signaled explicitly with the `-s` flag) is to choose a maximum-probability path. The alternative (the `-r` flag) is to select a path randomly in proportion to its probability. We know each path's probability because its total cost gives the negative log of its probability.

Try running `DelSpaces` on the text file `data/entest.txt`, which contains the first 50 sentences of the English test data `entest` from the HMM homework. Save the result as `entest-noisy.txt`. In general, you should use the `-r` flag to pass text through a noisy channel, so that it will randomly noise up the output (although in this introductory case the channel happens to be deterministic):

```
grmfilter -r noisy.grm DelSpaces <entest.txt >entest-noisy.txt
```

Uh-oh! Someone got into your files and used your own `DelSpaces` against you! Now how will you ever read any of your files?

After despairing for a while, you realize that you can just reverse `DelSpaces`'s actions. So you try `Invert[DelSpaces]`, but unfortunately that turns `Ifonlythereporterhadbeennice.` back into all kinds of things like

```
I fon lyt      he reporterh adbeenni  ce.
```

The correct solution is somewhere in that list of outputs, but you need to find it. What a perfect opportunity to use your language model LM and the Viterbi algorithm for finding the most probable path!

The idea is that the text actually came from the generative process (1), which can be represented as the composition

```
Generate = LM @ DelSpaces;    # almost right!
```

Unfortunately the output of LM is words, but the input to `DelSpaces` is characters. So they won't compose. You will need to stick a transducer `SpellText` in between. This transducer represents another deterministic step in the generative process that resulted in the noisy sequence of characters.

- (c) Define `SpellText` in `noisy.grm`. It should spell the first input word, output a space, spell the second input word, output another space, and so on. This yields the kind of text that actually appeared in `entrain.txt` (there is a space after each word in a sentence, including the last).

Now revise your definition of `Generate` to use `SpellText`.

- (d) Now you should be able to decode noisy text via

```
Decode = Invert[Generate];
```

Unfortunately, this machine will be too large (and slow to compile). So you should use the same approach as in question 8i, and ask `grmtest` to pass the noisy text through a sequence of inverted machines.

*Important:* At the end of your sequence of machines, you should add `PrintText`, which you can define for now to be equal to `SpellText`. This has the effect of pretty-printing the decoded result. It will turn the recovered sequence of words back into characters, and put spaces between the words.

Using `grmtest` in this way, try decoding each of the following. Note that the lowest-cost results are shown first. Discuss the pattern of results, and their costs, in your `README`:

- `Ifonlythereporterhadbeennice.`
- `If only.`
- `ThereportersaidtothecitythatEveryoneIskilled.`
- `Thankyou.`

- (e) The reason `Thankyou` failed is because we didn't account for OOVs. The vocabulary has an OOV symbol `<unk>`, but it is treated like any other word in the vocabulary.<sup>28</sup> So LM will accept phrases like `<unk> you`, but not `Thank you`.

So just as we described how to spell in the above questions, we'll now describe how to spell OOV words. We'll say that `<unk>` can rewrite as an arbitrarily long sequence of non-space text characters (`bytelib.kGraph`):

```
RandomChar = bytelib.kGraph <4.54>;
RandomWord = Optimize[(RandomChar (RandomChar <w1>)* ) <w2>];
SpellOOV = "<unk>".vocab : RandomWord;
```

The weight in `RandomChar` is saying that each of the 94 characters in `bytelib.kGraph` has the same probability, namely  $\frac{1}{94}$ , since  $-\log \frac{1}{94} \approx 4.54$ .

How about `RandomWord`? When you define it in `noisy.grm`, you'll have to give actual numbers for the numeric weights  $w_1$  and  $w_2$ . Try setting  $w_1 = 0.1$  and  $w_2 = 2.3$ . To check out the results, try these commands:

```
grmtest noisy.grm RandomWord    # evaluate cost of some strings
far2fst noisy.far RandomWord    # (get the FSA for commands below)
fstprintstring RandomWord.fst   # generate a random string
fstview RandomWord.fst         # look at the FSA
```

- i. What do  $w_1$  and  $w_2$  represent? Hint: the costs 0.1 and 2.3 are the negative logs of 0.9 and 0.1.
- ii. For each  $n \geq 0$ , what is the probability  $p_n$  that the string generated by `RandomWord` will have length  $n$ ?
- iii. What is the sum of those probabilities,  $\sum_{n=0}^{\infty} p_n$ ?
- iv. How would you change  $w_1$  and  $w_2$  to get longer random words on average?
- v. If you decreased *both*  $w_1$  and  $w_2$ , then what would happen to the probabilities of the random words? How would this affect the behavior of your decoder? Why?



56

vi. How could you improve the probability models `RandomChar` and `RandomWord`?

Once you've answered those questions, **reset  $w_1 = 0.1$  and  $w_2 = 2.3$  and proceed.**

- (f) Now, revise `Spell` so that it is not limited to spelling words in the dictionary, but can also randomly spell `<unk>`. (*Hint:* Use `Spell100V`.)

Also revise `PrintText` so that if your decoder finds an unknown word `<unk>`, you will be able to print that as the 5-character string `"<unk>."`

To check your updated decoder, try running the sentences from question 11d through it. Again discuss the pattern of results. Remember that if you want, you can add an extra argument to `grmtest` to limit the number of outputs printed per input.

- (g) Remember that your goal was to de-noise your corrupted files, whose spaces were removed by `DelSpaces`. Just run `grmfilter` again, but with three differences:

- Before, you were converting `entest.txt` to `entest-noisy.txt`. Now you should convert `entest-noisy.txt` to `entest-recovered.txt`.
- Instead of running the noisy channel forward, run it backward, using your pipeline from 11d. You can leave out the `PrintText` step of the pipeline since `grmfilter` is a bit smarter than `grmtest` about how it prints outputs.
- Since you want the most likely decoding and not a random decoding, don't use the `-r` flag this time.

57

Look at the results in `entest-recovered.txt`. What kinds of errors can you spot? Does this *qualitative* error analysis give you any ideas how to improve your decoder?

- (h) Suppose you'd like to *quantify* your performance. The metric we'll consider is the **edit distance** between `entest.txt` and `entest-recovered.txt`.

Edit distance counts the minimum number of edits needed to transduce one string ( $x$ ) into another ( $y$ ). The possible edits are

- **substitute** one letter for another;
- **insert** a letter;
- **delete** a letter;
- **copy** a letter unchanged.

Each of these operations has a cost associated with it. We'll stick with the standard unweighted edit distance metric in which substitutions, insertions and deletions all have cost 1; copying a character unchanged has cost 0. For simplicity we will treat the unknown word symbol as if really were the 5-character word `<unk>`, which must be edited into the true word.

As you know, edit distance can easily be calculated using weighted finite-state machines:

```
Sigma = bytelib.kBytes;
export Edit = (Sigma | (("|Sigma) : ("|Sigma)  <1>) )*;
```

The `Edit` machine transduces an input string  $x$  one byte at a time: at each step, it either passes an input character through with cost 0, or does an insert, delete or substitute with cost 1. That gives an edited version  $y$ . The cheapest way to get from  $x$  to a given  $y$  corresponds to the shortest path through  $x @ \text{Edit} @ y$ . As we saw in class, that machine has the form of an  $|x + 1| \times |y + 1|$  grid

<sup>28</sup>Except by some of the `ngram` utilities that we're not using.

with horizontal, vertical, and diagonal transitions. It has exponentially many paths, of various total cost, that represent different sequences of edit operations for turning  $x$  into  $y$ .

**We’ve given you an edit distance script** to calculate the edit distances between the corresponding lines of two files:

```
editdist entest.txt entest-recovered.txt
```

This will compare each recovered sentence to the original. Do the scores match your intuitive, qualitative results from 11g?

Please look at `grammars/editdist.grm`, the Thrax grammar used by the `editdist` script. You’ll see that it’s more complicated than `Edit`, but this construction *reduces* the size of the overall machine by a couple of orders of magnitude. While it still computes  $x \text{ Edit } y$ , it splits `Edit` up into two separate machines, `Edit1` and `Edit2`. We still find the shortest path, but now through

$$(x @ \text{Edit1}) @ (\text{Edit2} @ y);$$

By doing the composition this way, both  $x$  and  $y$  are able to impose their own constraints (what letters actually appear) on `Edit1` and `Edit2`, thus reducing the size of the intermediate machines. The resulting FST can be built quite quickly, though as mentioned before, it does have  $|x + 1| \times |y + 1|$  states and a similar number of arcs.

58

- (i) **Extra credit:** How can you modify your pipeline so that it recovers an appropriate spelling of each unknown word, rather than `<unk>`? For example, decoding `Thankyou` should give `Thank you` rather than `<unk> you`.<sup>29</sup>

59

12. **Extra credit (but maybe the real point of this homework, so read the question even if you don’t do it):** Finally, it’s time to have some fun. We just set up a noisy-channel decoder to handle a simple deterministic noisy channel. Now try it for some other noisy channels! The framework is nearly identical—just replace `DelSpaces` with some other FST. For each type of noisy channel,
- Define your channel in `noisy.grm` as a weighted FST.
  - Explain in `README` what you implemented and why, and how you chose the weights.
  - Use your channel to corrupt `entest.txt` into `entest-noisy.txt`.
  - Use your inverted channel, the language model, and `SpellText` to decode `entest-noisy.txt` back into `entest-recovered.txt`.
  - Look at the files and describe in your `README` what happened, with some examples.
  - Report the edit distance between `entest.txt` and `entest-recovered.txt`.

Have fun designing some of the channels below. Each converts a string of bytes into a string of bytes. In general make them non-deterministic (in contrast to `DelSpaces`), and play with the weights.

- DelSomeSpaces:** Nondeterministically delete none, some, or all spaces from an input string.
- DelSuffixes:** Delete various word endings. You may find <http://grammar.about.com/od/words/a/comsuffixes.htm> helpful.

<sup>29</sup>The recovered spelling is determined by the language model and the channel model. It won’t always match the noisy spelling. E.g., if the noisy channel tends to change letters into lowercase, then decoding `Thank you` might yield `THANK you`.

- (c) **Typos:** Introduce common typos or misspellings. You may get some inspiration from [http://en.wikipedia.org/wiki/Wikipedia:Lists\\_of\\_common\\_misspellings](http://en.wikipedia.org/wiki/Wikipedia:Lists_of_common_misspellings) or <http://en.wikipedia.org/wiki/File:Qwerty.svg>. Some real-world typos are due to the fact that some words *sound* similar, so if you're ambitious, you might be able to make some use of `data/cmudict.txt` or the **Pronounce** transducer.
- (d) **Telephone:** Deterministically convert (lower-case) letters to digits from the phone keypad. For example, `a` rewrites as `2`.
- (e) **Tinyphone:** Compose your **Telephone** FST with another FST that allows common cellphone keypad typos. For example, there should be a small chance of deleting a digit, doubling a digit, or substituting one of the adjacent digits for it.
- (f) Try composing some of these machines in various orders. As usual, give examples of what happens, and discuss the interactions.

Feel free to try additional noisy channels for more extra credit. You could consider capitalization, punctuation, or something crazy out of your imagination.<sup>30</sup>

---

<sup>30</sup>It might be fun to replace each word deterministically with its rhyming ending, using your **WordEnding** FST from question 8f (composed with something that transduces ARPAbet characters to the byte alphabet). Then your noisy channel decoder will find the highest-probability string that rhymes word-by-word with your original input text. Should be amusing.