# Source Code: Querying and Serving $N$-gram Language Models with Python

**Nitin Madnani**
Laboratory for Computational Linguistics and Information Processing
Institute for Advanced Computer Studies
University of Maryland, College Park
`nmadnani@umiacs.umd.edu`

## Abstract

Statistical $n$-gram language modeling is a very important technique in Natural Language Processing (NLP) and Computational Linguistics used to assess the fluency of an utterance in any given language. It is widely employed in several important NLP applications such as Machine Translation and Automatic Speech Recognition. However, the most commonly used toolkit (SRILM) to build such language models on a large scale is written entirely in C++ which presents a challenge to an NLP developer or researcher whose primary language of choice is Python. This article first provides a gentle introduction to statistical language modeling. It then describes how to build a native and efficient Python interface (using SWIG) to the SRILM toolkit such that language models can be queried and used directly in Python code. Finally, it also demonstrates an effective use case of this interface by showing how to leverage it to build a Python language model server. Such a server can prove to be extremely useful when the language model needs to be queried by multiple clients over a network: the language model must only be loaded into memory once by the server and can then satisfy multiple requests. This article supplements (Madnani, 2009) and provides the entire set of source code listings along with appropriate technical comments where necessary. Some of the listings may already be included with the primary article (in complete or excerpted form) but are reproduced here for the sake of completeness.

## 1 Overview

Natural Language Processing (NLP) is a very active area of research and development in Computer Science and with good reason. NLP applications such as machine translation and automatic speech recognition pervade almost every area of our lives. An important part of such applications is a statistical $n$-gram language model which is used to assess the fluency of an utterance in the chosen language. These models are estimated by using large amounts of monolingual text. Given the easy availability of such text today, $n$-gram language models are now used widely in NLP research. However, the toolkit that is most commonly used to construct these models is written in C++ and does not provide bindings in any other language. Given the rising popularity of dynamic languages as the primary language for building NLP applications and teaching NLP courses around the world—especially Python (Madnani, 2007; Madnani and Dorr, 2008; Bird et al., 2008)—the lack of such bindings represents a challenge. To address this issue, we develop a native Python interface to the toolkit using popular the Simplified Wrapper and Interface Genera-

tor (SWIG) tool. A more complete description of language modeling and the process undertaken to build the interface and the server may be found in the accompanying article (Madnani, 2009).

## 2 Requirements

The work presented in this article and (Madnani, 2009) is based on creating a native Python interface to the SRI Language Modeling (SRILM) toolkit that is most commonly used to build large-scale $n$-gram language models (Stolcke, 2002). The interface is constructed using the Simplified Wrapper and Interface Generator or SWIG (Beazley, 1996). Before delving into the source code, it is important to list the salient requirements for this project:

1. **SRILM Toolkit**: The source for the SRILM toolkit was obtained from the web. The version used in this work is the latest stable release (v1.5.8). It was compiled and installed on both 32-bit and 64-bit machines. It requires gcc version 3.4.3 or higher along with GNU make. Although the platform used in this paper is Linux-x86, SRILM can be compiled on a variety of platforms including Mac OS X and Win32. The reader is referred to the installation instructions included with the source for more details.

2. **SWIG**: The source for SWIG was also obtained from the web and compiled manually. The version used in this work is the latest stable release (v1.3.39). It was also compiled and installed on both 32-bit and 64-bit machines. SWIG is easy to install from source on almost any flavor of UNIX (including Mac OS X). In addition, pre-compiled binaries of the latest version for the Win32 platform are available from the SWIG web site.

3. **Python**: Python version 2.5 was used for this work and was installed on both 32-bit and 64-bit machines. Since the end-product is a Python interface, the Python header files must be installed. On most Linux distributions, these can be installed by either building Python from source or installing the `python-devel` package in addition to the standard `python` package.

## 3 Code

This section provides the complete listings of source code that are necessary for the reader of the primary article (Madnani, 2009) to build the Python interface on her own. Listings are grouped by subsections that correspond directly to the sections in the primary article. Only brief technical descriptions are included for each listing. For more details on creating the Python interface, the reader should refer to the primary article.

**Section 1: Introduction**

In this section, statistical $n$-gram language models are introduced and the reader is shown how to build a simple unsmoothed unigram language model using tools that are very easily available on any machine. Listing 1 shows how to find the most frequent words from Jane Austen's *Persuasion*. Listing 2 shows how to write a Python script that uses this corpus to build a very simple unigram language model. Note that we ignore all casing information when computing the unigram counts to build the model. Finally, Listing 3 shows how to use this unigram language model to evaluate some given test sentences. Note that the script requires the input file to contain one sentence per line.

However, this model is unsmoothed and will assign a zero probability to any sentence that contains words that were not seen in the training text. Listing 4 shows how to modify `unigram.py`

to incorporate Laplace smoothing. Since Laplace smoothing requires a $1$ to be added to each unigram count, the sum of *all* counts (which forms the denominator for the maximum likelihood estimation of unigram probabilities) increases by $1 * N$ where $N$ is the number of unique words in the training corpus. The results of using this smoothed model to evaluate the given test sentences are shown in Listing 5.

Trying to train language models by writing our own scripts is not ideal since it is very difficult to scale up to higher order $n$-grams and use more intelligent smoothing techniques. It is much easier to use a toolkit written specifically for this purpose: the SRI Language Modeling toolkit. Listing 6 shows how to build and use a smoothed trigram language model with this toolkit. The model is trained on Leo Tolstoy's *War and Peace* and can compute both probability and perplexity values for a file containing multiple sentences as well as for each individual sentence.

### Section 2: A Python Interface for Language Models

This section describes how to use SWIG to build a Python interface for querying language models built with the SRILM toolkit. The first step is to learn how SWIG works. Listings 7 and 8 illustrate how to use SWIG for generating a Python interface for a simple C program that computes the volume and surface area of a cylinder, given the radius $r$ and height $h$. Listing 9 shows how to compile and use the Python interface based on the wrapper files that are generated by SWIG. Note that the files `_cylinder.so` and `cylinder.py` together constitute the `cylinder` Python module.

Once we know how to write a SWIG interface file for a given C program, we can use the same procedure to create the SRILM Python module. We first create a C program containing all the functionality that the Python module should have. For the purpose of this article, this set of functionalities includes **introspection** (how many $n$-grams are in the model) and **evaluation** (score any sentence or a file containing multiple sentences in terms of both probability and perplexity values). Listing 10 shows this C program. Salient points regarding each function are provided as comments. This file defines functions to read in an SRILM file and map it to an internal data structure. Once this mapping is complete, the structure can be queried for various pieces of useful information. The code only includes very minimal error checking due to space considerations. In situations where perplexity cannot be calculated (all the words in a sentence are unseen, for example), this code returns a value of $-1.0$. Listing 11 shows the corresponding header file `srilm.h`. The SRILM toolkit is written in C++ which means a C++ compiler must be used. However, since our interface functions are written in C, we need to wrap our function declarations in an `extern C` block.

Given the C program and the header file declaring the functions, a corresponding SWIG interface file can now be created. This file must include the appropriate header files from the SRILM toolkit so that SWIG has knowledge of the data structures and other supporting functions from SRILM. In addition, it is important to note that the functions `getSentenceProbPpl()` and `getFileProbPpl()` require a pre-allocated array of size 3 as input which are then populated with the probability and perplexity values. SWIG provides the `carrays.i` library file which makes this quite easy. The `%array_class(<type>, <name>)` macro creates wrappers for an array object that can be passed around between C and Python easily. Listing 12 shows this interface definition file.

With the C program, the header file and the Python interface definition in hand, the interface itself can now be compiled using SWIG as shown in Listing 13. While compiling the Python shared library module (line 4), it is necessary to pass three SRILM static libraries (`liboolm.a`, `libdstruct.a` and `libmisc.a`) to the linker. These libraries contain the basic classes such as `Ngram`

(defined in `$SRILM/lm/src/LM.cc`) and `TextStats` (defined in `$SRILM/lm/src/TextStats.cc`) that come bundled with the toolkit for third-party development. Note that the environment variable `$SRILM` here refers to the top level directory where the SRILM toolkit was installed.

At this point the Python SRILM module is compiled and ready to use. Listing 14 shows a Python script that outputs information similar to the output of the SRILM program `ngram` that we looked at earlier. The script is fairly self-explanatory with the provided comments. The script also generates 5 random sentences using the trigram language model and writes them out to a file. The output of the script is shown in Listing 15. From manual inspection, we can easily verify that the output of this script is identical to `ngram`. This confirms that the module has been implemented correctly and works as intended.

**Section 3: Serving Language Models with Python**

This section details using the above SRILM Python module to build a language model server that can service multiple clients. To build such a server, we rely on the XML-RPC server functionality that comes bundled with Python in the `SimpleXMLRPCServer` module.

Listing 16 shows an example script that creates a very simple XML-RPC server. This particular server has just one function that is callable by remote clients. All this method does is takes in a list of numbers and returns a list of their squares. The same listing also shows how to register an instance of a class that reverses the string passed to it. Note that in this example, we are starting the server on port 8585 on the local machine itself. In a production environment, we would use the hostname of the machine on which we want the server to run.

Once the server is up and running, writing a client to talk to the server is also very simple. Listing 17 shows a simple client script: a proxy object is created that manages the communication with the actual LM server and all the remote methods are called directly on this proxy object. Listing 18 runs both the client and the server, and also shows the output of the client. We see that both methods were called successfully and produced the expected output. Since remote introspection was also turned on, the server can also be asked to list all the methods that are available for remote invocation. Note also the lines that are interspersed between the various outputs. These are logs of the various requests that were sent to the server. The number 200 represents a status code definition in the HTTP protocol and indicates that the POST request succeeded. These logs are useful for debugging purposes but can be turned off in a production setting by setting the `logRequests` parameter to `False` when instantiating the server. We will do this for the language model server.

Now let's see how we can combine the `srilm` module and the `SimpleXMLRPCServer` module to create a language model server. Since we are writing the server in Python, we need to be able to interface to an SRILM language model in our Python code. For this purpose, we will reuse the SWIG-generated Python module from our toolbox. In order to use it, we need to install the `srilm` module somewhere in the Python search path. There are several ways to do this:

(a) Keep the relevant module files (`srilm.py` and `_srilm.so`) in the same directory as the script using them.

(b) Install the module files in an appropriate place and then add that path to the search path by using the `sys.path.append()` method.

(c) Add the module installation path to the PYTHONPATH environment variable.

Listing 19 shows the code for the language model server. This server assumes that the `srilm` module generated by SWIG is already in the search path. The server is structured such that instead of taking a bunch of input arguments on the command line, the arguments can be specified in a configuration file, which is the only input to the server. Another neat thing about this server is that it's remotely stoppable, i.e., it can be cleanly shut down from inside a remote client. The server code can be dissected to highlight some important points:

(a) The parameter `allow_reuse_address` is set to `True` since we want to allow the server socket to be reused after the particular server instance using it has finished running.

(b) Rather than registering each function from the `srilm` module individually with the server, it is much easier to create a wrapper class and register an instance of this class with the server. Doing this makes all the instance's methods remote-callable.

(c) To allow a client to stop the server remotely, the server's `serve_forever()` method is modified such that before serving any request, it first checks to see if a shutdown request has been made by a remote client.

(d) If there is indeed such a request (as indicated by the instance variable `quit`), the server calls its `stop_server()` method which performs a clean shutdown including freeing the memory allocated to the language model data structure (which is why the `StoppableServer` class is designed to contain a pointer to this structure).

(e) Finally, note that the methods `getSentenceProbPpl()` and `getFileProbPpl()` can't simply return the `floatArray` instances since the XML-RPC protocol does not know how to marshal non-builtin datatypes into XML. We get around this limitation by unpacking the individual array elements and returning those instead.

An example configuration file for this server is shown in Listing 20. A configuration file should specify four things:

- The hostname of the machine on which the server will run.

- The port number which the server should use for communication.

- The path to the actual SRILM language model file that needs to be served.

- The order of the language model that's contained in this file.

A corresponding client for the LM server is shown in Listing 21. This client can run on *any* machine that's on the same network as the server machine. For the purpose of illustration, both the server and the client are being run on the same machine here. This can obviously be changed for a networked scenario by setting the appropriate hostnames in the server configuration and the client script. Listing 22 shows how to run the server and the client and also shows the output of running the client. Note that the server needs a couple of minutes after the command is issued to load the LM into memory server. It does not accept any connections until the LM is loaded. However, the server log should indicate when the LM has been successfully loaded.

For the purpose of convenience and easy access, the index below lists all source code listings in this manuscript and indicates the respective page on which it appears.

**Listings**

**References**

Beazley, David M. 1996. SWIG : An Easy to Use Tool for Integrating Scripting Languages with C and C++. In *Proceedings of the Fourth Annual USENIX Tcl/Tk Workshop*. `http://www.swig.org`.

Bird, Steven, Ewan Klein, Edward Loper, , and Jason Baldridge. 2008. Multidisciplinary instruction with the Natural Language Toolkit. In *Proceedings of the Third Workshop on Issues in Teaching Computational Linguistics (TeachCL)*.

Madnani, Nitin. 2007. Getting Started on Natural Language Processing with Python. *ACM Crossroads*, 13(4).

Madnani, Nitin. 2009. Querying and Serving N-gram Language Models with Python. *The Python Papers*, 4(2).

Madnani, Nitin and Bonnie J. Dorr. 2008. Combining Open Source with Research to Re-engineer a Hands-on Introductory NLP Course. In *Proceedings of TeachCL*.

Stolcke, Andreas. 2002. SRILM - An Extensible Language Modeling Toolkit. In *Proc. Intl. Conf. Spoken Language Processing*. `http://www.speech.sri.com/projects/srilm/`.

# Source Code Listings

Listing 1: Using unix tools to find unigram frequencies (Church 1989).

```
1  $ tr −sc 'A−Za−z' '\012' < persuasion.txt | sort \
2       | uniq −c | sort −nr > persuasion.1grams
3
4  $ head persuasion.1grams
5     3277 the
6     2849 to
7     2805 and
8     2669 of
9     1586 a
10    1401 in
11    1331 was
12    1177 had
13    1159 her
14    1124 I
```

Listing 2: Computing sentence probabilities using a unigram model (`unigram.py`).

```
1  inputfile = sys.argv[1]
2  counts = {}
3  for line in open('persuasion.txt','r'):
4      for word in line.strip().lower().split():
5          counts[word] = counts.get(word, 0) + 1
6
7  total = sum(counts.values())
8
9  for sent in open(inputfile, 'r'):
10     prob = 1.0
11     for w in sent.strip().lower().split():
12         prob = prob * (float(counts.get(w, 0))/total)
13     print prob
```

Listing 3: Evaluating sentences with the unigram model.

```
1  $ cat input.sentences
2    this is the first sentence
3    this is the second one
4    this sentence gets a zero score
5
6  $ python unigram.py input.sentences
7    1.8318e−14
8    1.8656e−13
9    0.0
```

Listing 4: Incorporating Laplace smoothing into the unigram model (`unigram_smooth.py`).

```
1   inputfile = sys.argv[1]
2   counts = {}
3   for line in open('persuasion.txt','r'):
4       for word in line.strip().lower().split():
5           counts[word] = counts.get(word, 0) + 1
6
7   total = sum(counts.values()) + len(counts)
8
9   for sent in open(inputfile, 'r'):
10      prob = 1.0
11      for w in sent.strip().lower().split():
12          prob = prob * (float(counts.get(w, 0)+1)/total)
13      print prob
```

Listing 5: Evaluating sentences with the smoothed unigram model.

```
1  $ python unigram_smooth.py input.sentences
2    1.7138e−14
3    1.4629e−13
4    3.1679e−24
```

Listing 6: Building and using a smoothed trigram language model using the SRILM toolkit.

```
1  $ ngram-count -order 3 -text warandpeace.txt -tolower -lm warpeace.lm
2
3  $ cat input.sentences
4    this is the first sentence
5    this is the second one
6    this sentence gets a zero score
7
8  $ ngram -order 3 -lm warpeace.lm -ppl input.sentences
9    file input.sentences: 3 sentences, 16 words, 0 OOVs
10   0 zeroprobs, logprob= -51.5919 ppl= 519.232 ppl1= 1676.84
11
12 $ ngram -lm warpeace.lm -order 3 -ppl input.sentences -debug 1
13   reading 20200 1-grams
14   reading 185625 2-grams
15   reading 67713 3-grams
16
17   this is the first sentence
18   1 sentences, 5 words, 0 OOVs
19   0 zeroprobs, logprob= -12.6611 ppl= 128.879 ppl1= 340.579
20
21   this is the second one
22   1 sentences, 5 words, 0 OOVs
23   0 zeroprobs, logprob= -11.6708 ppl= 88.1321 ppl1= 215.855
24
25   this sentence gets a zero score
26   1 sentences, 6 words, 0 OOVs
27   0 zeroprobs, logprob= -27.26 ppl= 7839.4 ppl1= 34940.6
28
29   file input.sentences: 3 sentences, 16 words, 0 OOVs
30   0 zeroprobs, logprob= -51.5919 ppl= 519.232 ppl1= 1676.84
```

Listing 7: A C program that computes the volume and surface area of a cylinder (`cylinder.c`).

```
1  /* Compute the volume and surface area of a cylinder of radius r and height h */
2  #include <math.h>
3
4  /* Define pi */
5  float pi = 3.1415;
6
7  float volume(float r, float h) { return pi*r*r*h; }
8
9  float surface(float r, float h) { return 2*pi*r*(r+h); }
```

Listing 8: A SWIG interface file for `cylinder.c` (`cylinder.i`)

```
1  %module cylinder
2  %{
3    #include <math.h>
4    extern float volume(float r, float h);
5    extern float surface(float r, float h);
6  %}
7
8  extern float volume(float r, float h);
9  extern float surface(float r, float h);
```

Listing 9: Compiling and Using the Python interface to `cylinder.c`.

```
1   # Tell SWIG that we need a Python module. This will generate a wrapper cylinder_wrap.c
2   # and a Python module file called cylinder.py
3   $ swig -python cylinder.i
4
5   # Compile both the original program cylinder.c and the wrapper cylinder_wrap.c.
6   # Make sure to tell the compiler where the Python header files are.
7   $ gcc -c -fpic cylinder.c cylinder_wrap.c -I/usr/local/include \
8     -I/usr/local/include/python2.5
9
10  # Link the two object files together into _cylinder.so, the companion to cylinder.py.
11  $ gcc -lm -shared cylinder.o cylinder_wrap.o -o _cylinder.so
12
13  # Try out the newly compiled cylinder module in Python (interactive mode)
14  $ python
15  >>> import cylinder
16  >>> cylinder.volume(2,2)
17  25.131999969482422
18  >>> cylinder.surface(2,2)
19  50.26544189453125
```

Listing 10: C program with the desired LM query functions using the SRILM API (`srilm.c`).

```c
1  #include "Prob.h"
2  #include "Ngram.h"
3  #include "Vocab.h"
4  #include "srilm.h"
5  #include <cstdio>
6  #include <cstring>
7  #include <cmath>
8
9  Vocab *swig_srilm_vocab;
10
11 /* Initialize the ngram model */
12 Ngram* initLM(int order) {
13   swig_srilm_vocab = new Vocab;
14   return new Ngram(*swig_srilm_vocab, order);
15 }
16
17 /* Delete the ngram model */
18 void deleteLM(Ngram* ngram) {
19   delete swig_srilm_vocab;
20   delete ngram;
21 }
22
23 /* Read the given LM file into the model */
24 int readLM(Ngram* ngram, const char* filename) {
25   File file(filename, "r");
26   if (!file) {
27       fprintf(stderr,"Error:: Could not open file %s\n", filename);
28       return 0;
29   }
30   else return ngram->read(file, 0);
31 }
32
33 /* How many ngrams are there? */
34 int howManyNgrams(Ngram* ngram, unsigned order) {
35   return ngram->numNgrams(order);
36 }
```

```
37  /* Get the probability and perplexity values for a given sentence */
38  void getSentenceProbPpl(Ngram* ngram, char* sentence, float ans[3]) {
39    char* words[80];
40    unsigned numparsed;
41    TextStats stats;
42    char* scp;
43
44    /* Create a copy of the input string to be safe */
45    scp = strdupa(sentence);
46
47    /* Parse the sentence into its constitutent words */
48    numparsed = Vocab::parseWords(scp, (VocabString *)words, 80);
49    if (numparsed > 80) {
50      fprintf(stderr, "Error: Number of words in sentence should be <= 80.\n");
51    }
52
53    /* Calculate the sentence probability and store it in stats */
54    ans[0] = ngram->sentenceProb(words, stats);
55
56    /* Now calculate the denominator for perplexity */
57    if (stats.numWords + stats.numSentences > 0) {
58      int denom = stats.numWords - stats.numOOVs - stats.zeroProbs + 1;
59
60      /* Calculate ppl value */
61      if (denom > 0)
62        ans[1] = LogPtoPPL(stats.prob/denom);
63      else
64        ans[1] = -1.0;
65
66      /* calculate ppl1 value */
67      denom -= 1;
68      if (denom > 0)
69        ans[2] = LogPtoPPL(stats.prob/denom);
70      else
71        ans[2] = -1.0;
72    }
73  }
```

```
74  /* Run the file through SRILM and store its statistics into a TextStats instance */
75  unsigned fileStats(Ngram* ngram, const char* filename, TextStats &stats) {
76    File givenfile(filename, "r");
77    if (!givenfile) {
78      fprintf(stderr,"Error:: Could not open file %s\n", filename);
79        return 1;
80    }
81    else {
82      ngram->pplFile(givenfile, stats, 0);
83      return 0;
84    }
85  }
86
87  /* Compute probability and perplexity over entire file */
88  void getFileProbPpl(Ngram* ngram, const char* filename, float ans[3]) {
89    TextStats stats;
90
91    if (!fileStats(ngram, filename, stats)) {
92
93      /* calculate the file logprob */
94      ans[0] = stats.prob;
95
96      /* calculate the perplexity values if we can */
97      if (stats.numWords + stats.numSentences > 0) {
98        int denom = stats.numWords - stats.numOOVs - stats.zeroProbs + stats.numSentences;
99
100       /* calculate ppl value */
101       if (denom > 0)
102         ans[1] = LogPtoPPL(stats.prob/denom);
103       else
104         ans[1] = -1.0;
105
106       /* calculate ppl1 value */
107       denom -= stats.numSentences;
108       if (denom > 0)
109         ans[2] = LogPtoPPL(stats.prob/denom);
110       else
111         ans[2] = -1.0;
112     }
113   }
114 }
```

```
115  /* Generate random sentences using the language model */
116  void randomSentences(Ngram* ngram, unsigned numSentences, const char* filename) {
117      VocabString* sent;
118      File outFile(filename, "w");
119      unsigned i;
120
121      /* set seed so that new sentences are generated each time */
122      srand48(time(NULL) + getpid());
123
124      for (i = 0; i < numSentences; i++) {
125          /* call the built-in SRILM method to generate a random sentence */
126          sent = ngram->generateSentence(50000, (VocabString *) 0);
127
128          /* write the generated sentence to the file */
129          swig_srilm_vocab->write(outFile, sent);
130          fprintf(outFile, "\n");
131      }
132      outFile.close();
133  }
```

Listing 11: The accompanying header file for the SRILM C program (`srilm.h`)

```
1   #ifndef SRILMWRAP_H
2   #define SRILMWRAP_H
3
4   extern "C" {
5     Ngram* initLM(int order);
6     void deleteLM(Ngram* ngram);
7     int readLM(Ngram* ngram, const char* filename);
8     int howManyNgrams(Ngram* ngram, unsigned order);
9     void getSentenceProbPpl(Ngram* ngram, char* sentence, float ans[3]);
10    unsigned fileStats(Ngram* ngram, const char* filename, TextStats &stats);
11    void getFileProbPpl(Ngram* ngram, const char* filename, float ans[3]);
12    void randomSentences(Ngram* ngram, unsigned numSentences, const char* filename);
13  }
14
15  #endif
```

Listing 12: The SWIG interface definition file (`srilm.i`)

```
1   /* We want our Python module to be called 'srilm' */
2   %module srilm
3
4   /* Include all needed header files here */
5   %{
6    #include "Ngram.h"
7    #include "Vocab.h"
8    #include "Prob.h"
9    #include "srilm.h"
10  %}
11
12  /* We need to generate interfaces for all functions defined in srilm.h */
13  %include srilm.h
14
15  /* We need a C to Python array bridge and SWIG comes with one */
16  %include "carrays.i"
17  %array_class(float, floatArray);
18
19  /* Since we are allocating memory when creating random sentences */
20  /* make sure that we tell SWIG so that it can free it */
21  %newobject randomSentences;
```

Listing 13: Compiling a Python interface to `srilm.c`.

```
1   $ swig -python srilm.i
2   $ g++ -c -fpic srilm.c srilm_wrap.c -I/usr/local/include \
3       -I/usr/local/include/python2.5 -I$SRILM/include
4   $ g++ -shared srilm.o srilm_wrap.o -loolm -ldstruct -lmisc \
5       -L$SRILM/lib/i686 -o _srilm.so
```

Listing 14: Using the Python interface to SRILM (`test_srilm.py`).

```python
1   # Use the srilm module
2   from srilm import *
3
4   # Initialize a variable that can hold the data for a 3-gram language model
5   n = initLM(3)
6
7   # Read the model we built into this variable
8   readLM(n, "warpeace.lm")
9
10  # How many n-grams of different order are there ?
11  print "There are", howManyNgrams(n, 1), "unigrams in this model."
12  print "There are", howManyNgrams(n, 2), "bigrams in this model."
13  print "There are", howManyNgrams(n, 3), "trigrams in this model."
14  print
15
16  # First go through our test file and get the individual sentence probabilities
17  # and perplexity values
18  infile = 'input.sentences'
19  for sentence in file(infile):
20      svals = floatArray(3)
21      print sentence,
22      getSentenceProbPpl(n, sentence, svals)
23      print 'logprob = %.6f, ppl= %.6f, ppl1= %.6f\n' % (svals[0], svals[1], svals[2])
24
25  # Now get overall values
26  fvals = floatArray(3)
27  getFileProbPpl(n, infile, fvalues)
28  print 'Overall:'
29  print 'logprob = %.6f, ppl= %.6f, ppl1= %.6f\n' % (fvals[0], fvals[1], fvals[2])
30
31  # Generate 5 random sentences
32  randomSentences(n, 5, 'faux-tolstoy.txt')
33
34  # Free model variable (to avoid a memory leak)
35  deleteLM(n);
```

Listing 15: Output for `test_srilm.py`.

```
 1   $ python test_srilm.py
 2
 3     There are 20200 unigrams in this model.
 4     There are 185625 bigrams in this model.
 5     There are 67713 trigrams in this model.
 6
 7     this is the first sentence
 8     logprob = −12.661089, ppl= 128.878830, ppl1= 340.578888
 9
10     this is the second one
11     logprob = −11.670806, ppl= 88.132133, ppl1= 215.854584
12
13     this sentence gets a zero score
14     logprob = −27.259981, ppl= 7839.404297, ppl1= 34940.585938
15
16     Overall:
17     logprob = −51.591873, ppl= 519.232483, ppl1= 1676.841675
18
19   $ cat faux-tolstoy.txt
20
21     it was about to go to the movement of his soldiers , dolokhov .
22     he had not seen looking at the movement of a dangers vasili suddenly jumped up .
23     " why are you ill ?   " and suddenly surrounded added , as soon as the fifteenth .
24     found the rooms opened she said , indeed make a moment , asked rostov .
25     he went to the chief .
```

Listing 16: Creating a Python XML-RPC server (`example-server.py`).

```python
from SimpleXMLRPCServer import SimpleXMLRPCServer

# Instantiate a server on the local machine on port 8585
# Note that the input is a single argument: a tuple
server = SimpleXMLRPCServer(('localhost', 8585))

# Tell the server to support introspection by clients
server.register_introspection_functions()

# Define a method that takes in a list of numbers and returns
# another list with the square of each corresponding number
def squarer(l):
    return [x**2 for x in l]

# Register this function with the server
server.register_function(squarer)

# Define a class with a single method
class reverser:
    # Define a reverse method which returns the reverse of the input string
    def reverse(self, s):
        return s[::-1]

# Register an instance of the above class with the server
server.register_instance(reverser())

# Start the server's main loop
server.serve_forever()
```

Listing 17: A client for `example-server.py` (`example-client.py`).

```
1  import xmlrpclib
2
3  # Connect to the server
4  s = xmlrpclib.ServerProxy('http://localhost:8585')
5
6  # Call the squarer function
7  print s.squarer([1,2,3,4,5])
8
9  # Call the reverse method
10 print s.reverse('reverse this')
11
12 # Get a list of all the methods that are available on this server
13 print s.system.listMethods()
```

Listing 18: Output for `example-client.py`.

```
1  # Run the server forever
2  $ python example-server.py &
3
4  # Run the client
5  $ python example-client.py
6
7  localhost - - [28/May/2009 22:40:55] "POST /RPC2 HTTP/1.0" 200 -
8  [1, 4, 9, 16, 25]
9  localhost - - [28/May/2009 22:40:55] "POST /RPC2 HTTP/1.0" 200 -
10 siht esrever
11 localhost - - [28/May/2009 22:40:55] "POST /RPC2 HTTP/1.0" 200 -
12 ['reverse', 'squarer', 'system.listMethods', 'system.methodHelp', 'system.methodSignature']
```

Listing 19: A language model server (`lmserver.py`).

```python
1  from SimpleXMLRPCServer import SimpleXMLRPCServer
2  import srilm, sys
3
4  class StoppableServer(SimpleXMLRPCServer):
5    allow_reuse_address = True
6
7    def __init__(self, addr, lm, *args, **kwds):
8      self.myhost, self.myport = addr
9      self.lm = lm
10     SimpleXMLRPCServer.__init__(self, addr, *args, **kwds)
11     self.register_function(self.stop_server)
12     self.quit = False
13
14   def serve_forever(self):
15     while not self.quit:
16       try:
17         self.handle_request()
18       except KeyboardInterrupt:
19         break
20     self.server_close()
21
22   def stop_server(self):
23     self.quit = True
24     srilm.deleteLM(self.lm)
25     return 0, "Server terminated on host %r, port %r" % (self.myhost, self.myport)
26
27 class LM:
28   def __init__(self, lm):
29     self.lm = lm
30     self.tmparray = srilm.floatArray(3)
31
32   def howManyNgrams(self, type):
33     return srilm.howManyNgrams(self.lm, type)
34
35   def getSentenceProbPpl(self, s):
36     srilm.getSentenceProbPpl(self.lm, s, self.tmparray)
37     return self.tmparray[0], self.tmparray[1], self.tmparray[2]
38
39   def getFileProbPpl(self, filename):
40     srilm.getFileProbPpl(self.lm, filename, self.tmparray)
41     return self.tmparray[0], self.tmparray[1], self.tmparray[2]
```

```
42   # Get command line arguments
43   args = sys.argv[1:]
44   if len(args) != 1:
45     sys.stderr.write('Usage: lmserver.py <configfile>\n')
46     sys.exit(1)
47   else:
48     configfile = open(args[0],'r')
49
50   # Initialize all needed variables to None
51   address = None
52   port = None
53   lmfilename = None
54   lmorder = None
55
56   # Source the config file
57   exec(configfile)
58
59   # Make sure all needed variables got set properly
60   missing = []
61   for var in ['address', 'port', 'lmfilename', 'lmorder']:
62     if eval(var) is None:
63       missing.append(var)
64   if missing:
65     sys.stderr.write('The following options are missing in \
66                       the configuration file: %s\n' % ', '.join(missing))
67     sys.exit(1)
68
69   # Initialize the LM data structure
70   lmstruct = srilm.initLM(lmorder)
71
72   # Read the LM file into this data structure
73   try:
74     sys.stderr.write('Loading language model ...')
75     srilm.readLM(lmstruct, lmfilename)
76   except:
77     sys.stderr.write('Error: Could not read LM file\n')
78     sys.exit(1)
79   finally:
80     sys.stderr.write(' done.\n')
```

```
81   # Create a wrapper around this data structure that exposes all
82   # the functions as methods
83   lm = LM(lmstruct)
84
85   # Create an instance of the stoppable server with a pointer to
86   # this LM data structure
87   try:
88     server = StoppableServer((address, port), lmstruct, logRequests = False)
89   except:
90     sys.stderr.write('Error: Could not create server\n')
91     sys.exit(1)
92
93   # Tell server to support remote introspection
94   server.register_introspection_functions()
95
96   # Register the LM wrapper instance with the server
97   server.register_instance(lm)
98
99   # Start the server
100  server.serve_forever()
```

Listing 20: An example configuration file for `lmserver.py` (`server.conf`).

```
1    # The address of the server. If you want to run the server and the client
2    # on a local machine that's not on a network, set this to 'localhost'.
3    address = "localhost"
4
5    # Port number. You might want to make sure what ports you have access to
6    port = 8585
7
8    # The LM file you wish to use. You can also use gzipped LMs.
9    lmfilename = "warpeace.lm"
10
11   # The order of the LM
12   lmorder = 3
```

Listing 21: An example client for the LM server (`lmclient.py`).

```python
1   import xmlrpclib, socket, sys
2
3   # Connect to the server (use hostname of server if on network)
4   s = xmlrpclib.ServerProxy("http://localhost:8585")
5
6   # Make the remote procedure calls on the server
7   try:
8     # How many n-grams of different order are there in this LM ?
9     print "There are", s.howManyNgrams(1), "unigrams in this model."
10    print "There are", s.howManyNgrams(2), "bigrams in this model."
11    print "There are", s.howManyNgrams(3), "trigrams in this model."
12    print
13
14    # First go through our test file and get the individual sentence
15    # probabilities and perplexity values
16    infile = 'input.sentences'
17    for sentence in file(infile):
18      print sentence,
19      svals = s.getSentenceProbPpl(sentence)
20      print 'logprob = %.6f, ppl= %.6f, ppl1= %.6f\n' % (svals[0], svals[1], svals[2])
21
22    # Now get overall values
23    fvals = s.getFileProbPpl(infile)
24    print 'Overall:'
25    print 'logprob = %.6f, ppl= %.6f, ppl1= %.6f\n' % (fvals[0], fvals[1], fvals[2])
26
27  except socket.error:
28    sys.stderr.write('Error: could not connect to the server.\n')
29    sys.exit(1)
30  except xmlrpclib.Fault, fault:
31    sys.stderr.write('Error %d: %s\n' % (fault.faultCode, fault.faultString))
32    sys.exit(2)
33  except xmlrpclib.ProtocolError, err:
34    sys.stderr.write('Error %d: %s\n' % (err.errcode, err.errmsg))
35    sys.exit(3)
36
37  # Stop the server. NOTE: You will probably do this in the last client
38  # (if you know which one that is) or in a clean-up script when you are
39  # absolutely sure that all clients are finished.
40  s.stop_server()
```

Listing 22: Output for `lmclient.py`.

```
1   # Run the server and background it
2   $ python lmserver.py server.conf > server.log &
3
4   # Run the client
5   $ python lmclient.py
6
7     There are 20200 unigrams in this model.
8     There are 185625 bigrams in this model.
9     There are 67713 trigrams in this model.
10
11    this is the first sentence
12    logprob = -12.661089, ppl= 128.878830, ppl1= 340.578888
13
14    this is the second one
15    logprob = -11.670806, ppl= 88.132133, ppl1= 215.854584
16
17    this sentence gets a zero score
18    logprob = -27.259981, ppl= 7839.404297, ppl1= 34940.585938
19
20    Overall:
21    logprob = -51.591873, ppl= 519.232483, ppl1= 1676.841675
```