# Indri Project Report

whuang022nccu

May 22, 2019

# 1 Introduction

## 1.1 About Indri

Indri is a sub version of Lemur project , which is designed for information retrieval and made by C++ . It contains apis for makeing a customized full-text search engine thst cover such as stemming ,indexing ,scoring the documents and doing the searching.

## 1.2 Environment Building

We prepared a Linux liked environment for the project. For a cleaning environment to testing and also contains the full GUI system , we choose to install the virtual machine in the original computer which operating system is Ubuntu Desktop 16.04 LTS.

Listing 1: Install VMware Workstation Player

```
1   $ wget https://www.vmware.com/go/getplayer-linux
2   $ sudo bash getplayer-linux
```

After the command we will get the installer of vmplayer , we just need to run install step by step .
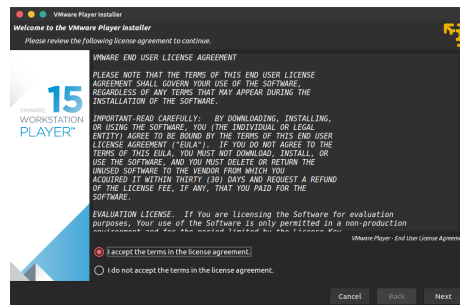


Figure 1: Intital Screenshot of vmplayer installer

Next we need the Ubuntu 16.04 ISO file for the vmplayer.

Listing 2: Get the iso file of Ubuntu 16.04

```
1    $ wget http://releases.ubuntu.com/16.04/ubuntu-16.04.6-desktop-amd64.iso
```

Now we just need to run the vmplayer and install the Ubuntu.We setting our testing environment called 01_Ubuntu_Project_Indri .
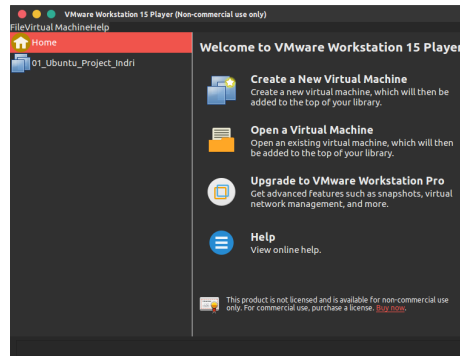


Figure 2: Intital Screenshot of vmplayer

For checking the system , we can type :

Listing 3: Get system info

```
1    $ sudo apt install inxi
2    $ inxi -M
3    $ inxi
```
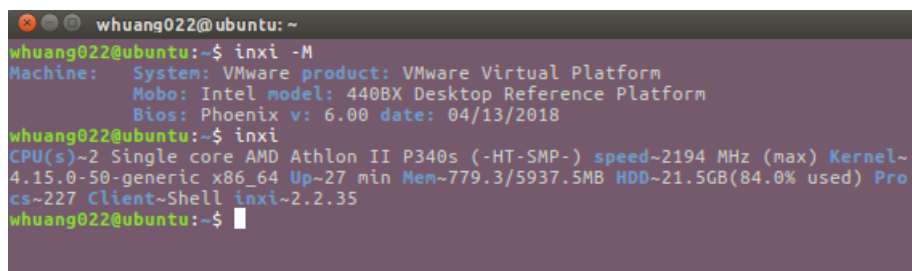


Figure 3: Checking system info

## 1.3   Install Indri

We get Indri 5.8 :

Listing 4: Install Indri in bash

```
1   $ wget https://sourceforge.net/projects/lemur
2   /files/lemur/indri-5.8/indri-5.8.tar.gz/download -O indri-5.8.tar.gz
3   $ tar xvfz indri-5.8.tar.gz
```

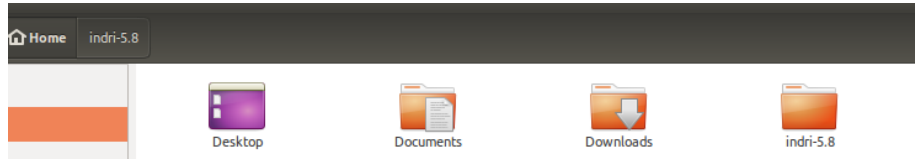After unzip , we can see the Indri 5.8 folder :



Figure 4: Checking system info

### 1.3.1   Installation of dependencies

There is the list [1] that we need to install before compiling the Indri.

1. build-essential:the g++ stuff.

2. git : the version controll system.

3. vim : the bash editor.

4. zlibc : compressed file-system using in Indri .

5. zlib1g : compressed file-system using in Indri .

6. zlib1g-dev : compressed file-system using in Indri

So we just typing :

Listing 5: Get dependencies

```
1   $ apt-get install build-essential
2   git vim wget zlibc zlib1g zlib1g-dev -y
3   $ apt-get update
```

For the zlib , we need to change the Makefile [2] , and add the compiling flag "
-lz " .

```
-include MakeDefns

INSTALLDIRS = $(bindir) $(includedir) $(pkgincludedir) $(includedir)/lemur $(libdir) $(pkgdatadir) $(pkgdatadir)/doc
LDFLAGS = -lz
.PHONY: all dist clean install $(INSTALLDIRS) site-search
all:
	$(MAKE) -C contrib $(LDFLAGS)
	$(MAKE) -C obj -f ../src/Makefile $(LDFLAGS)
ifeq ($(WITH_SWIG), 1)
	$(MAKE) -C swig/src $(LDFLAGS)
endif
	$(MAKE) -C buildindex $(LDFLAGS)
	$(MAKE) -C runquery $(LDFLAGS)
	$(MAKE) -C indri $(LDFLAGS)
	$(MAKE) -C dumpindex $(LDFLAGS)
	$(MAKE) -C harvestlinks $(LDFLAGS)
	$(MAKE) -C pagerank $(LDFLAGS)
	$(MAKE) -C rmodel $(LDFLAGS)
	$(MAKE) -C makeprior $(LDFLAGS)
	$(MAKE) -C site-search $(LDFLAGS)
	$(MAKE) -C modifyfields $(LDFLAGS)
	$(MAKE) -C clarity $(LDFLAGS)
	$(MAKE) -C reformulate $(LDFLAGS)

$(INSTALLDIRS):
	$(INSTALL_DIR) $@

clean:
	$(MAKE) clean -C contrib $(LDFLAGS)
	$(MAKE) clean -C obj -f ../src/Makefile $(LDFLAGS)
ifeq ($(WITH_SWIG), 1)
	$(MAKE) clean -C swig/src $(LDFLAGS)
endif
	$(MAKE) clean -C buildindex $(LDFLAGS)
	$(MAKE) clean -C runquery $(LDFLAGS)
	$(MAKE) clean -C indri $(LDFLAGS)
	$(MAKE) clean -C dumpindex $(LDFLAGS)
	$(MAKE) clean -C harvestlinks $(LDFLAGS)
	$(MAKE) clean -C pagerank $(LDFLAGS)
	$(MAKE) clean -C rmodel $(LDFLAGS)
	$(MAKE) clean -C makeprior $(LDFLAGS)
	$(MAKE) -C site-search clean $(LDFLAGS)
	$(MAKE) -C modifyfields clean $(LDFLAGS)
	$(MAKE) -C clarity clean $(LDFLAGS)
	$(MAKE) -C reformulate clean $(LDFLAGS)
	rm -f depend/*

distclean: clean
	rm -rf autom4te.cache MakeDefns Makefile.app config.log config.status indri.ncb indri.suo
dist: distclean
	$(MAKE) stamp -C obj -f ../src/Makefile
	cd ..; tar zcvf indri-`date +%Y%m%d-%H%M`.tar.gz indri $(LDFLAGS)
```

Figure 5: edit Makefile

### 1.3.2 Installation of Indri

Listing 6: Install Indri in bash

```
1   $ cd ~
2   $ cd indri-5.8
3   $ sudo ./configure
4   $ make
5   $ sudo make install
```

The system will have Indri after the installation.

4

## 2 Using Indri to indexing WT2G

Fist we creat the path of workspace and download the dataset WT2G.And creat
the path for indexing.

Listing 7: Download WT2G File

```
1   $ cd ~
2   $ mkdir workspace
3   $ cd workspace
4   $ wget http://wm5.nccu.edu.tw/base/10001/
5   course/10021115/content/proj02/WT2G.tbz
6   $ tar jxf WT2G.tbz
7   $ mkdir ~/workspace/WT2G_index
8   $ mkdir ~/workspace/WT2G_index_no_stemming
9   $ mkdir ~/workspace/WT2G_index_parameter_file
10  $ mkdir ~/workspace/WT2G_run_shell
```

Next we write the index parameter files ,the format can find at [3] . There will
be two types of index file , one is stemming by porter and the other one do not
have stemming.

Listing 8: project_W2T.xml

```
1   <parameters>
2   <memory>1G</memory>
3   <stemmer>
4     <name>porter</name>
5   </stemmer>
6   <index>/home/whuang022/workspace/WT2G_index</index>
7   <corpus>
8     <path>/home/whuang022/workspace/WT2G</path>
9     <class>trecweb</class>
10  </corpus>
11  </parameters>
```

Listing 9: project_W2T_withoutstemming.xml

```
1   <parameters>
2     <memory>1G</memory>
3     <index>/home/whuang022/workspace/WT2G_index_no_stemming</index>
4     <corpus>
5         <path>/home/whuang022/workspace/WT2G</path>
6         <class>trecweb</class>
7     </corpus>
8   </parameters>
```

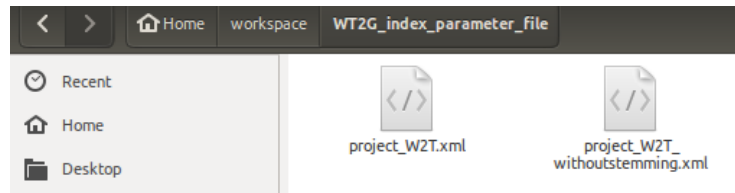The index parameters files will save at path
~/workspace/WT2G_index_parameter_file

Figure 6: The index parameter files

Then we write shell script to do indexing.We put all shell script in the path
~/workspace/WT2G_run_shell

We can use the executable file IndriBuildIndex to build index from bash command . The command format is IndriBuildIndex [*index_parameter_file_path*]

Listing 10: index.sh (use porter stemming)

```bash
1
2  #!/bin/bash
3
4  cd ~/indri-5.8/buildindex
5
6  if [ -e IndriBuildIndex ]
7  then
8      echo "IndriBuildIndex [OK] "
9      ./IndriBuildIndex ~/workspace/WT2G_index_parameter_file
10     /project_W2T.xml
11 else
12     echo "IndriBuildIndex [Error] "
13 fi
```

Listing 11: index_no_stemming.sh (no stemming)

```bash
1
2  #!/bin/bash
3
4  cd ~/indri-5.8/buildindex
5
6  if [ -e IndriBuildIndex ]
7  then
8      echo "IndriBuildIndex [OK] "
9      ./IndriBuildIndex ~/workspace/WT2G_index_parameter_file
10     /project_W2T_withoutstemming.xml
11 else
12     echo "IndriBuildIndex [Error] "
13 fi
```
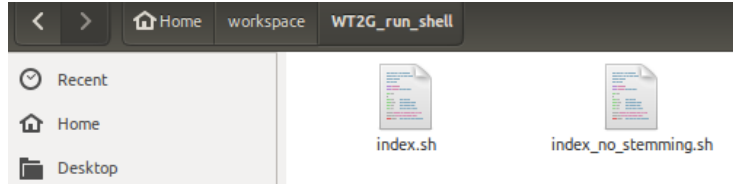
Figure 7: The index shell scripts

Finally , we just run the shell scripts , ande Indri will generate the index files.

Listing 12: Run index shell scripts

```
1   $ cd ~/workspace/WT2G_run_shell
2   $ sudo bash index.sh
3   $ sudo bash index_no_stemming.sh
```

The results is in the index path.



Figure 8: The index files generate from IndriBuildIndexs

# 3 Retrieval Files form Indri index files

## 3.1 Query file and scoring functions

We can use the executable file IndriRunQuery to doing the retrieval from command ,the scoring can write in the query xml file.The detail information of the format about query file , can read at [4]. The important parameters are:

1. <index>:the path of index files

2. <count> : the number of retun files

3. <baseline> : the basic scoring function like Okapi TFIDF

4. <method>: the other scoring functions

5. <query>: the query

### 3.1.1 Vector space model Okapi TFIDF

From the project document there is :

> You will have to use for the weights OKAPI TF x IDF where $OKAPITF = tf/(tf + 0.5 + 1.5 * doclen/avgdoclen)$. For queries, Okapi TF can also be computed in the same way, just use the length of the query to replace doclen. Also note that the definition of OKAPI TF is $tf/tf + k1((1-b) + b*doclen/avgdoclen)$. In the above formula, you can set k1 = 2 and b = 0.75, to end up with:$tf/(tf + 0.5 + 1.5 * doclen/avgdoclen)$.

So we just need to set okapi with k1=2.0, b:0.75. The query file is :

Listing 13: okapi.xml

```
1   <parameters>
2     <baseline>okapi,k1=2.0, b:0.75</baseline>
3     <index>/home/whuang022/workspace/WT2G_index</index>
4     <count>1000</count>
5     <query>
6         <type>indri</type>
7         <number>401</number>
8         <text>
9            foreign minorities germany
10        </text>
11     </query>
12  </parameters>
```

### 3.1.2 Language modeling with Laplace Smoothing

Because Indri dose not have an implement of laplace smoothing , so we need to implement the customized scoring function .In Indri project the scoring functions has an abstract base class named TermScoreFunction.hpp which was designed for define the scoring function class's api.
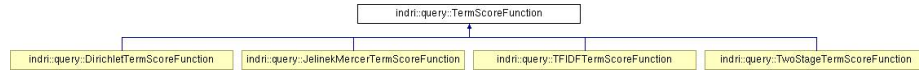


Figure 9: indri::query::TermScoreFunction Class Reference [5]

Listing 14: TermScoreFunction.hpp

```cpp
1   #ifndef INDRI_TERMSCOREFUNCTION_HPP
2   #define INDRI_TERMSCOREFUNCTION_HPP
3   namespace indri
4   {
5     namespace query
6     {
7
8       class TermScoreFunction {
9       public:
10        virtual double scoreOccurrence( double occurrences, int contextLength ) = 0;
11        virtual double scoreOccurrence( double occurrences, int contextLength, double documentOccurrences, int documentLength ) = 0;
12      };
13    }
14  }
15  #endif // INDRI_TERMSCOREFUNCTION_HPP
```

For the customized scoring function , the basic design for the class that extends TermScoreFunction is implement the first function "scoreOccurrence( double occurrences, int contextLength )" , and return the first function's result in the second function. The laplace smoothing mathod is as the following :$p_i = \frac{m_i+1}{n+k}$where m = term frequency, n=number of terms in document (doc length) , k=number of unique terms in corpus. In this implementation , we use log transform for $p_i$ to $log\,(p_i)$. The "occurrences" is term frequency , the "contextLength" is number of terms in document.However,we need the size of unique terms in corpus. So we need to use QueryEnvironment.hpp and it's api termCountUnique () to get k.The following is the code :

Listing 15: LaplaceTermScoreFunction.hpp

```cpp
1
2   #ifndef INDRI_LAPLACETERMSCOREFUNCTION_HPP
3   #define INDRI_LAPLACETERMSCOREFUNCTION_HPP
4   #include <indri/QueryEnvironment.hpp>
5   #include <math.h>
6   #include <iostream>
7   namespace indri
8   {
9     namespace query
10    {
11      class LaplaceTermScoreFunction : public TermScoreFunction {
12      private:
13        double _alpha;
14        indri::api::QueryEnvironment *env;
15      public:
16        LaplaceTermScoreFunction(double alpha , const std::string& indexName) {
17          _alpha = alpha;
18          env = new indri::api::QueryEnvironment();
19          env->addIndex(indexName);
20        }
21        double scoreOccurrence( double occurrences, int contextSize ) {
22          double seen = ( double(occurrences) + _alpha ) / ( double(contextSize) + double(env->termCountUnique ()) );
23          return log( seen );
24        }
25        double scoreOccurrence( double occurrences, int contextSize, double documentOccurrences, int documentLength ) {
26          return scoreOccurrence( occurrences, contextSize );
27        }
28      };
29    }
30  }
```

9

```
31
32   #endif // INDRI_LAPLACETERMSCOREFUNCTION_HPP
```

To add this customized class to Indri wheile calling the IndriRunQuery , we need to add the mapping function in src/TermScoreFunctionFactory.cpp .We add the rule to calling our laplace class :

Listing 16: TermScoreFunctionFactory.cpp

```
1
2   ... (some other codes )
3   else if ( method == "laplace" || method == "l" || method == "lap" ){
4       // laplace −− takes parameter "alpha" & index path
5       double alpha = spec.get( "alpha", 1.0 );
6       return new indri::query::LaplaceTermScoreFunction(alpha,spec.get( "index", "" ));
7   }
8   ... (some other codes )
```

Finally , we need to re compile the code (use "make" and "sudo make install") , the new executable file will contain the customized api design by us ,and the query file is :

Listing 17: laplace.xml

```
1    <parameters>
2      <method>laplace,alpha=1.0,index=/home/whuang022/workspace/WT2G_index</method>
3      <index>/home/whuang022/workspace/WT2G_index</index>
4      <count>1000</count>
5      <query>
6          <type>indri</type>
7          <number>401</number>
8          <text>
9             foreign minorities germany
10         </text>
11     </query>
12   </parameters>
```

### 3.1.3   Jelinek-Mercer smoothing using the corpus

The formula for Jelinek-Mercer smoothing is : $p_i = \lambda P + (1 - \lambda) Q$ and we select $\lambda = 0.8$.

Therefore , the query file is :

Listing 18: jelinek_mercer.xml

```
1    <parameters>
2      <rule>method:linear,collectionLambda:0.8,documentLambda:0.2</rule>
3      <index>/home/whuang022/workspace/WT2G_index</index>
4      <count>1000</count>
5      <query>
6          <type>indri</type>
```

```
 7            <number>401</number>
 8            <text>
 9              foreign minorities germany
10            </text>
11        </query>
12  </parameters>
```

### 3.1.4 Jelinek-Mercer smoothing add Laplace

Listing 19: JelinekMercerTermLaplaceScoreFunction.hpp

```cpp
 1
 2
 3
 4   //
 5   // JelinekMercerTermLaplaceScoreFunction
 6   #ifndef INDRI_JELINEKMERCERLAPLACETERMSCOREFUNCTION_HPP
 7   #define INDRI_JELINEKMERCERLAPLACETERMSCOREFUNCTION_HPP
 8   #include <indri/QueryEnvironment.hpp>
 9   #include <math.h>
10   namespace indri
11   {
12     /// indri query processing and scoring components.
13     namespace query
14     {
15
16       class JelinekMercerLaplaceTermScoreFunction : public TermScoreFunction {
17       private:
18         double _lambda;
19         double _backgroundLambda;
20         double _collectionFrequency;
21         double _collectionComponent;
22         double _oneLevelCollectionComponent;
23         double _contextLambda;
24         double _collectionLambda;
25         double _documentLambda;
26         double _foregroundLambda;
27         indri::api::QueryEnvironment *env;
28       public:
29         JelinekMercerLaplaceTermScoreFunction( double collectionFrequency, double collectionLambda, double documentLambda
   , const std::string& indexName ) {
30           _contextLambda = (1 − collectionLambda − documentLambda);
31           _collectionFrequency = collectionFrequency;
32           _collectionLambda = collectionLambda;
33           _documentLambda = documentLambda;
34           _foregroundLambda = (1 − _collectionLambda);
35           env = new indri::api::QueryEnvironment();
36     env->addIndex(indexName);
37
38           assert( _documentLambda >= 0.0 && _documentLambda <= 1.0 );
39           assert( _collectionLambda >= 0.0 && _collectionLambda <= 1.0 );
40           assert( _contextLambda >= 0.0 && _contextLambda <= 1.0 );
41
42           _collectionComponent = _collectionLambda * _collectionFrequency;
43         }
44
45         double scoreOccurrence( double occurrences, int contextSize ) {
46           double seen2 = ( double(occurrences) +1.0 ) / ( double(contextSize) + double(env->termCountUnique ()) );
47           double contextFrequency = contextSize ? occurrences / double(contextSize) : 0.0;
48           return (log( _foregroundLambda * contextFrequency + _collectionComponent )+log(seen2));
49         }
50
51         double scoreOccurrence( double occurrences, int contextSize, double documentOccurrences, int documentLength ) {
52           double seen2 = ( double(occurrences) +1.0 ) / ( double(contextSize) + double(env->termCountUnique ()) );
53           double contextFrequency = contextSize ? occurrences / double(contextSize) : 0.0;
54           double documentFrequency = documentLength ? documentOccurrences / double(documentLength) : 0.0;
55           return (log( _contextLambda * contextFrequency + _documentLambda * documentFrequency + _collectionComponent )+log(seen2));
56         }
57       };
58     }
59   }
60
61   #endif // INDRI_JELINEKMERCERLAPLACETERMSCOREFUNCTION_HPP
```

## 3.2 The shell script to Retreival

For example of laplace we will have :

Listing 20: Run index shell scripts

```bash
#!/bin/bash

cd ~/indri-5.8/runquery

if [ -e IndriRunQuery ]
then
    echo "IndriRunQuery [OK] "
    ./IndriRunQuery ~/workspace/WT2G_query_parameter_file/laplace.xml -trecFormat=true > ~/workspace/WT2G_result/output_laplace.txt
    cd ~/workspace/WT2G_result/
    perl trec_eval.pl [-q] qrels.401-450.txt output_laplace.txt > laplace_evaluation.txt
else
    echo "IndriRunQuery [Error] "
fi
```

# 4 Result

## 4.1 Un-interpolated Mean Average Precision

|             | Okapi TF | Laplace | Jelinek Mercer | Jelinek Mercer add Laplace |
|-------------|----------|---------|----------------|----------------------------|
| stemming    | 0.2704   | 0.5123  | 0.0591         | 0.5123                     |
| no stemming | 0.1932   | 0.3205  | 0.0407         | 0.3205                     |

## 4.2 Precision at rank 10 documents

|             | Okapi TF | Laplace | Jelinek Mercer | Jelinek Mercer add Laplace |
|-------------|----------|---------|----------------|----------------------------|
| stemming    | 0.3000   | 0.9000  | 0.1000         | 0.9000                     |
| no stemming | 0.3000   | 0.8000  | 0.0000         | 0.8000                     |

In conclution , if we do not stemming the performace is not good at all.

## 4.3 Access the project

Listing 21: Access the project

```
1
2    1. The source code with customized class :
3  https://github.com/whuang022nccu/IndriLab
4  2. The shell scripts (.sh) , parameter files (.xml) ,and reslut files :
5  https://github.com/whuang022nccu/IndriLab/tree/master/parameter%20files%20
```

# References

[1] Can't Make - error, `https://sourceforge.net/p/lemur/discussion/2106523/thread/3abe18beed/`, 25 04 2019.

[2] Compilation problems with ZLIB, `https://stackoverflow.com/questions/9700414/compilation-problems-with-zlib/18875275#18875275`, 16 03 2012.

[3] IndriBuildIndex Parameters, David Fisher , `https://sourceforge.net/p/lemur/wiki/IndriBuildIndex%20Parameters/`

[4] Specifying Retrieval Parameters, David Fisher , `https://sourceforge.net/p/lemur/wiki/IndriRunQuery/`

[5] indri::query::TermScoreFunction Class Reference , `https://lemur.sourceforge.io/indri/classindri_1_1query_1_1TermScoreFunction.html`