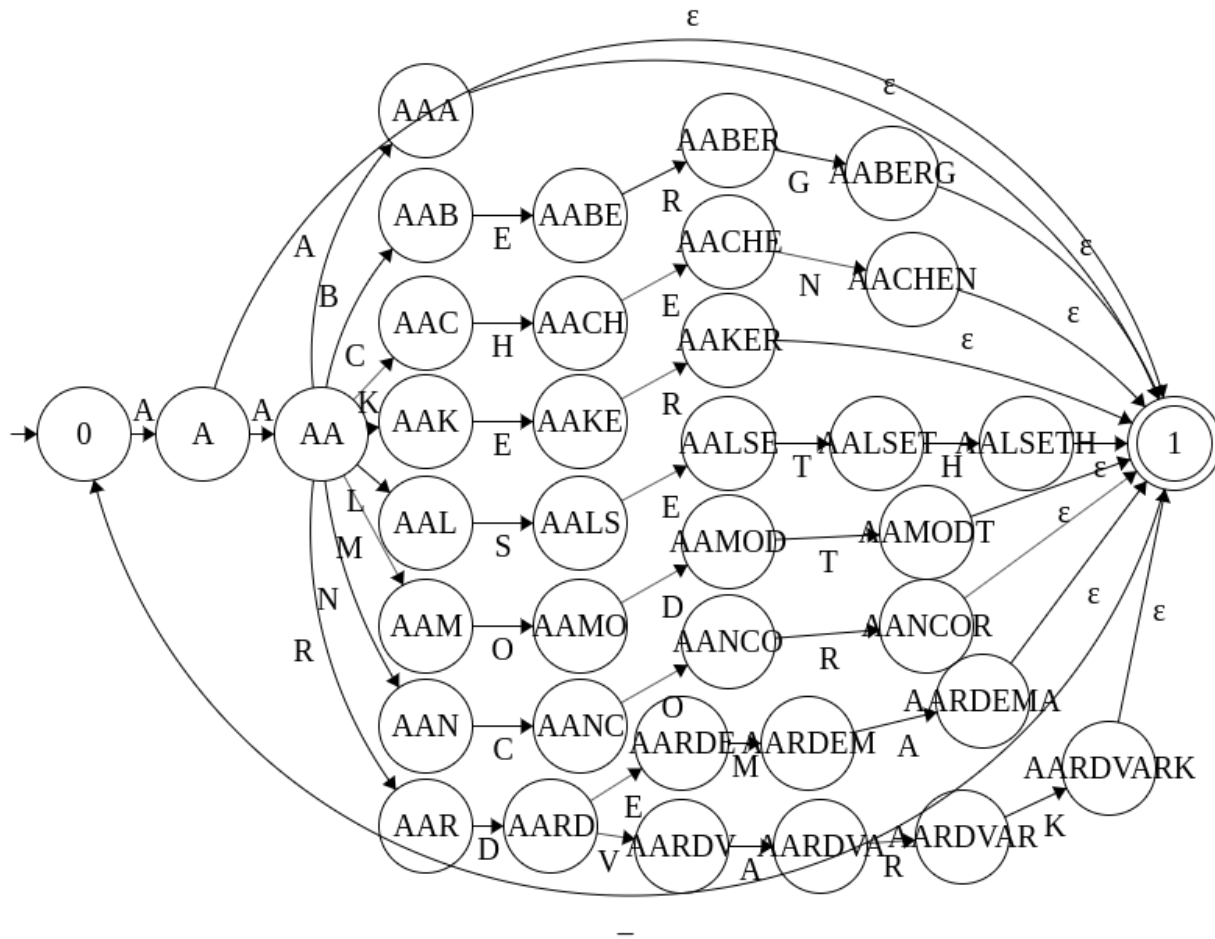


NLP HW1

1. Finite-state acceptors

1-1.



1-2.

(a)

Number of states in result: 256331

Number of arcs in result: 361732

Number of paths in result (valid for acyclic only; a cycle means infinitely many):

105402.000000005

There are 361732 transitions, and 256331 states.

(b)

Strings:

Input line 1: LIST_THE_FLIGHTS_FROM_BALTIMORE_TO_S
EATTLE_THAT_STOP_IN_MINNEAPOLIS

(102 states / 101 arcs reduce-> 79/78)

LIST_THE_FLIGHTS_FROM_BALTIMORE_TO_SEATTLE
_THAT_STOP_IN_MINNEAPOLIS

Input line 2: DOES_THIS_FLIGHT_SERVE_DINNER

(48 states / 47 arcs reduce-> 35/34)

DOES_THIS_FLIGHT_SERVE_DINNER

Input line 3: I_NEED_A_FLIGHT_TO_SEATTLE_LEAVING_F
ROM_BALTIMORE_MAKING_A_STOP_IN_MINNEAPOLI
S

(119 states / 118 arcs reduce-> 93/92)

I_NEED_A_FLIGHT_TO_SEATTLE_LEAVING_FROM_BA
LTIMORE_MAKING_A_STOP_IN_MINNEAPOLIS

Input line 4: I_NEED_TO_HAVE_DINNER_SERVED

(48 states / 47 arcs reduce-> 35/34)

I_NEED_TO_HAVE_DINNER_SERVED

Input line 5: I_HAVE_TWO_FRIENDS_THAT_WOULD_LIKE_
TO_VISIT_ME_ON_WEDNESDAY_HERE_IN_WASHINGT
ON_D_C

(128 states / 127 arcs reduce-> 100/99)

I_HAVE_TWO_FRIENDS_THAT_WOULD_LIKE_TO_VISI
T_ME_ON_WEDNESDAY_HERE_IN_WASHINGTON_D_C

strings.bad:

Input line 1: I_WUNT_TO_LEEVE_MONDAY_MORNING

(0 states / 0 arcs)

Empty or invalid result of composition with transducer "english.fsa".

Input line 2: NOW_I_NEAD_A_FLIGHT_ON_TOOSDAY_FROM
_PHEENIX_TO_DETROIT

(0 states / 0 arcs)

Empty or invalid result of composition with transducer "english.fsa".

Input line 3: WHICH_ONES_LEEVE_IN_THE_MORNING

(0 states / 0 arcs)

Empty or invalid result of composition with transducer "english.fsa".

Input line 4: WHICH_ONES_ARRIVE_ERLY_IN_THE_DAY

(0 states / 0 arcs)

Empty or invalid result of composition with transducer "english.fsa".

Input line 5: I _ N E A D _ A _ F L I G H T _ F R O M _ P H E E N I X _ T O _ D E T R
O I T _ L E E V I N G _ M O N D A Y _ E E V E N I N G

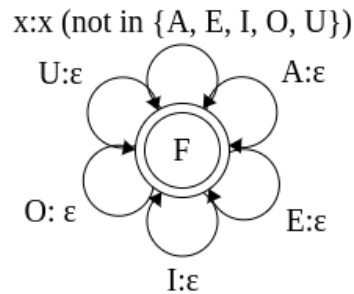
(0 states / 0 arcs)

Empty or invalid result of composition with transducer "english.fsa".

2. Finite-state transducers

2-3.

(a)



(b)

L S T _ T H _ F L G H T S _ F R M _ B L T M R _ T _ S T T L _ T H T _ S T P _ N _ M
N N P L S
D S _ T H S _ F L G H T _ S R V _ D N N R
_ N D _ _ F L G H T _ T _ S T T L _ L V N G _ F R M _ B L T M R _ M K N G _ _ S T
P _ N _ M N N P L S
_ N D _ T _ H V _ D N N R _ S R V D
_ H V _ T W _ F R N D S _ T H T _ W L D _ L K _ T _ V S T _ M _ N _ W D N S D Y
_ H R _ N _ W S H N G T N _ D _ C

(c)

B L D N G
B L D N G U
B L D N U G
B L D N G U U
B L D N U G U
B L D N G U U U
B L D N U G U U
B L D N G U U U U
B L D N U G U U U

BLDNGUUUUU

2-4.

(a)

LST_TH_FLGHTS_FRM_BLTMR_T_STTL_THT_STP_N_M
NNPLS
DS_THS_FLGHT_SRV_DNNR
_ND__FLGHT_T_STTL_LVNG_FRM_BLTMR_MKNG__ST
P_N_MNNPLS
_ND_T_HV_DNNR_SRV
_HV_TW_FRNDS_THT_WLD_LK_T_VST_M_N_WDNSDY
_HR_N_WSHNGTN_D_C

(b)

The accuracy is 1.2865497076%

(c)

Because automaton could generate an infinite number of words by randomly inserting each vowel, and the majority of words are meaningless which looks randomly organized. And automaton will select the first word which is output itself (without vowel letters). Therefore, it is nearly impossible to get the correct word for each non-vowel word (except words like “S H Y”,etc.). So the accuracy of vowel restorer is very low.

3. Combining FSAs and FSTs

3-5.

(a)

First of all, we generated a FSA for file *vocab*. Then, based on that FSA, we did some modifications in python file. If the input is one of the vowels, the output will be modified to **e**, instead of the input itself. Finally, we combine the FSA with FST.

We could also achieve combination by using command:

```
cat strings | carmel -slibOEWk 1 english.fsa remove-vowels.fst
```

(b)

The accuracy is 33.0994152047%

make.py is our python file for implementing our combined vowel restorer.

vowels_remove.fst is our combined FST. *strings.restored.mod1* is the restored samples of the implementation.

(c)

The accuracy is still low and the results are not satisfactory. People are based on grammar, context meaning and the correctness of the word to restore the vowels for

English words. However, in this case, our combined FSA and FST automaton only based on word correctness. It will lead to random guess to restore the vowels without thinking about the grammatical problem and sentence meaning, when several words are same after running *vowels.remove*.

3-6.

(a)

First of all, we could generate *english.fsa* by using test file (*strings*) instead of the *vocab* file which includes much more words. Our restorer will be more likely to recover correctly since the fewer number of words could be derived in smaller *english.fsa*.

For the further improve, we add weights in each path, and automaton will know the most possible vowel letter which will follow by the state (identified string). Hence automaton is able to find the most likely word based on the frequency of words in *strings* which will improve the accuracy a lot.

(b)

Question6.py is our python file for implementing our improved vowel restorer.

Vowels_remove.wfst is our improved FST with adding probabilities in it.

strings.restored.mod2 is the restored samples of the implementation.

(c)

The accuracy, after implementing two ideas above, improves to 91.35%.