

To run in Colab (skip otherwise)

```
!curl -sS https://apertium.projectjj.com/apt/install-release.sh | sudo
bash
!apt install apertium-all-dev lexd
```

Make sure that `hfst-guess bezhta.guesser.hfst` is in the same directory

Coverage

Unrecognized tokens

```
def guess_word(word):
    output = os.popen(f"echo {word} | hfst-guess bezhta.guesser.hfst -
n 100").read().rstrip('\n\n')
    parses = []
    for el in output.split('\n'):
        parses.append(':'.join(el.split('\t')))

    return parses

def check_guesser_coverage(path):
    with open(path, 'r', encoding='utf-8') as file:
        file = file.read()
        words = re.findall(r"(?P<num>\d+) \^(?P<word>[а-яōāēīŷ|"]
+)\./.*\$", file)
        wd = {}
        for word in words:
            guess = guess_word(word[1])
            if guess == ['']:
                guess = []
            wd[word[1]] = {'number': int(word[0]), 'guess': guess,
'len_guess': len(guess)}

        n_recog = 0
        n_unrecog = 0
        for word in wd.keys():
            if wd[word]['len_guess']:
                n_recog += wd[word]['number']
            else:
                n_unrecog += wd[word]['number']

        print('recog: ', n_recog)
        print('unrecog: ', n_unrecog)
        print('coverage: ', n_recog/(n_recog+n_unrecog))

    return wd
```

```

gospel = check_guesser_coverage('unrecog-gospel.txt')
recog: 4019
unrecog: 995
coverage: 0.8015556441962505

gospel = check_guesser_coverage('unrecog-gospel.txt')
recog: 4019
unrecog: 995
coverage: 0.8015556441962505

gospel = check_guesser_coverage('unrecog-gospel.txt')
recog: 4019
unrecog: 995
coverage: 0.8015556441962505

gospel = check_guesser_coverage('unrecog-gospel.txt')
recog: 4019
unrecog: 995
coverage: 0.8015556441962505

```

Full corpora

```

def check_guesser_coverage_full(path):
    with open(path, 'r', encoding='utf-8') as file:
        un = []
        total = 0.0
        tokens_recognised = 0.0
        unique = 0.0
        lines = file.readlines()
        for l in lines:
            parses = [guess_word(i.strip()) for i in l.split()]
            for w in parses:
                if w[0] != '':
                    tokens_recognised += 1
                    total += 1
        print('tokens: ', total)
        print('tokens recognised: ', tokens_recognised)
        print('tokens unrecognised: ', total - tokens_recognised)
        print('token coverage: ', tokens_recognised / total)

luke_full = check_guesser_coverage_full('text-Luke.txt')

tokens: 19739.0
tokens recognised: 16011.0
tokens unrecognised: 3728.0
token coverage: 0.8111353158721313

```

```

prov_full = check_guesser_coverage_full('text-Prov.txt')

tokens: 11644.0
tokens recognised: 9192.0
tokens unrecognised: 2452.0
token coverage: 0.7894194434902095

text_1_full = check_guesser_coverage_full('text-turkey.txt')

tokens: 86.0
tokens recognised: 71.0
tokens unrecognised: 15.0
token coverage: 0.8255813953488372

text_2_full = check_guesser_coverage_full('text-life.txt')

tokens: 116.0
tokens recognised: 90.0
tokens unrecognised: 26.0
token coverage: 0.7758620689655172

```

Accuracy

All tokens

```

def guess_word_acc(word):
    output = os.popen(f"echo {word} | hfst-guess bezhta.guesser.hfst -
n 100").read().rstrip('\n\n')
    parses = []
    for el in output.split('\n'):
        if output.split('\n')[0] != '':
            parse = el.split('\t')[1]
            parses.append(re.sub('\[GUESS_CATEGORY=\w+\]', '', parse))
    if parses == ['']:
        parses = []
    return parses

def check_guesser_acc(path):
    with open(path, 'r', encoding='utf-8') as file:
        file = file.read()
        fully_guessed = 0.0
        recog = 0.0
        total = 0.0
        tags_guessed = 0.0
        seen = []
        words = re.findall(r"^(?P<word>[а-яōāēñ̄|ʰ]+)\/(?P<parse>.*)\
$", file)
        for w in words:

```

```

        if w not in seen:
            guesses = guess_word_acc(w[0])
            if len(guesses) != 0:
                if w[1] in guesses:
                    fully_guessed += 1
                #else:
                #    #print('FAIL')
                #    #print('standard:', w[0], w[1])
                #    #print(guesses)
            guessed_tags = [re.findall(r"<.*>", l)[0] for l in
guesses]

            true_tags = re.findall(r"<.*>", w[1])[0]
            if true_tags in guessed_tags:
                tags_guessed += 1
            recog += 1
            total +=1
            seen.append(w)
            print('total:', total)
            print('recognised:', recog)
            print('types recognised:', recog/total)
            print('fully_guessed:', fully_guessed)
            print('fully_guessed_all:', fully_guessed / total)
            print('fully_guessed_recog:', fully_guessed / recog)
            print('tags_guessed:', tags_guessed / total)

check_guesser_acc('text-1-gold.txt')

total: 76.0
recognised: 64.0
types recognised: 0.8421052631578947
fully_guessed: 26.0
fully_guessed_all: 0.34210526315789475
fully_guessed_recog: 0.40625
tags_guessed: 0.34210526315789475

check_guesser_acc('text-2-gold.txt')

total: 79.0
recognised: 68.0
types recognised: 0.8607594936708861
fully_guessed: 23.0
fully_guessed_all: 0.2911392405063291
fully_guessed_recog: 0.3382352941176471
tags_guessed: 0.2911392405063291

```

Perfect tokens

```

def check_guesser_acc_perfect(path):
    with open(path, 'r', encoding='utf-8') as file:
        file = file.read()

```

```

    fully_guessed = 0.0
    recog = 0.0
    total = 0.0
    tags_guessed = 0.0
    words = re.findall(r"^(?P<word>[a-яōāēñŷ!"]+)\/(.*<(n|v|num|
dem)>.*)\$", file)
    seen = []
    for w in words:
        if w not in seen:
            guesses = guess_word_acc(w[0])
            if len(guesses) != 0:
                if w[1] in guesses:
                    fully_guessed += 1
                #else:
                #print('FAIL')
                #print('standard:', w[0], w[1])
                #print(guesses)
            guessed_tags = [re.findall(r"<.*>", l)[0] for l in
guesses]
            true_tags = re.findall(r"<.*>", w[1])[0]
            if true_tags in guessed_tags:
                tags_guessed += 1
            recog += 1
            total += 1
            seen.append(w)

    print('total:', total)
    print('recognised:', recog)
    print('types recognised:', recog/total)
    print('fully_guessed:', fully_guessed)
    print('fully_guessed_all:', fully_guessed / total)
    print('fully_guessed_recog:', fully_guessed / recog)
    print('tags_guessed:', tags_guessed / total)

```

check_guesser_acc_perfect('text-1-gold.txt')

```

total: 52.0
recognised: 46.0
types recognised: 0.8846153846153846
fully_guessed: 26.0
fully_guessed_all: 0.5
fully_guessed_recog: 0.5652173913043478
tags_guessed: 0.5

```

check_guesser_acc_perfect('text-2-gold.txt')

```

total: 51.0
recognised: 46.0
types recognised: 0.9019607843137255
fully_guessed: 23.0

```

```

fully_guessed_all: 0.45098039215686275
fully_guessed_recog: 0.5
tags_guessed: 0.45098039215686275

check_guesser_acc_perfect('text-2-gold.txt')

total: 51.0
recognised: 46.0
types recognised: 0.9019607843137255
fully_guessed: 23.0
fully_guessed_all: 0.45098039215686275
fully_guessed_recog: 0.5
tags_guessed: 0.45098039215686275

```

Verbs & nouns separately

```

def check_guesser_acc_nouns(path):
    with open(path, 'r', encoding='utf-8') as file:
        file = file.read()
        fully_guessed = 0.0
        recog = 0.0
        total = 0.0
        tags_guessed = 0.0
        words = re.findall(r"^(?P<word>[а-яōāēīŷ|ʱ]+)\/(.*<(n)>.*)\$",
file)
        seen = []
        for w in words:
            if w not in seen:
                print(w)
                guesses = guess_word_acc(w[0])
                if len(guesses) != 0:
                    if w[1] in guesses:
                        fully_guessed += 1
                    #else:
                    #    print('FAIL')
                    #    print('standard:', w[0], w[1])
                    #    print(guesses)
                guessed_tags = [re.findall(r"<.*>", l)[0] for l in
guesses]
                true_tags = re.findall(r"<.*>", w[1])[0]
                if true_tags in guessed_tags:
                    tags_guessed += 1
                recog += 1
                total += 1
                seen.append(w)

        print('total:', total)
        print('recognised:', recog)
        print('types recognised:', recog/total)
        print('fully_guessed:', fully_guessed)

```

```

print('fully_guessed_all:', fully_guessed / total)
print('fully_guessed_recog:', fully_guessed / recog)
print('tags_guessed:', tags_guessed / total)

check_guesser_acc_nouns('text-1-gold.txt')

('базайгаъой', 'базай<n><obl><cum><ess>', 'n')
('бикълабашейоълыи', 'бикълабашейоълыи<n>', 'n')
('ли', 'ли<n><abs>', 'n')
('аьл|аъаъш', 'аьл|<n><obl><in><ess><abl>', 'n')
('миц', 'миц<n><abs>', 'n')
('бикълабашейолъил|она', 'бикълабашейолъи<n><abs><quot><add>', 'n')
('ъаъ"гъаъ', 'ъаъ"гъаъ<n><abs>', 'n')
('бикълабашейолъина', 'бикълабашейолъи<n><abs><add>', 'n')
('гъекъар', 'гъекъар<n><abs>', 'n')
('оълоъхъаън', 'оълоъхъаън<n><abs>', 'n')
('клималид', 'клима<n><obl><ins>', 'n')
('клима', 'клима<n><abs>', 'n')
('цлитлад', 'цлитл<n><obl><ins>', 'n')
('роъыл', 'роъыл<n><abs>', 'n')
('инкар', 'инкар<n><abs>', 'n')
('къимат', 'къимат<n><abs>', 'n')
('гемо', 'гемо<n><abs>', 'n')
('гъикматна', 'гъикмат<n><abs><add>', 'n')
total: 18.0
recognised: 12.0
types recognised: 0.6666666666666666
fully_guessed: 3.0
fully_guessed_all: 0.16666666666666666
fully_guessed_recog: 0.25
tags_guessed: 0.16666666666666666

check_guesser_acc_nouns('text-2-gold.txt')

('дунналнакъодā', 'дуннал<n><abs><add><irr>', 'n')
('аллагъ', 'аллагъ<n><abs>', 'n')
('интернетна', 'интернет<n><abs><conj>', 'n')
('клетлослъина', 'клетлослъи<n><abs><add>', 'n')
('телевизорлана', 'телевизор<n><obl><erg><add>', 'n')
('ийо', 'ийо<n><abs>', 'n')
('аболкъодā', 'або<n><obl><dat><irr>', 'n')
('руслан', 'руслан<n><abs>', 'n')
('оъмроъ', 'оъмроъ<n><abs>', 'n')
('эркенлъина', 'эркенлъи<n><abs><add>', 'n')
('або', 'або<n><abs>', 'n')
('загъматлъи', 'загъмат<n><abs>', 'n')
('хунзахъ', 'хунзахъ<n><abs>', 'n')
('жо', 'жо<n><abs>', 'n')
('оъмроълаъкъаъ', 'оъмроъ<n><obl><sup><ess>', 'n')
('оъмроънаъ', 'оъмроъ<n><abs><add>', 'n')

```

```

('абона', 'абон<n><abs><add>', 'n')
('аьгаьрлѝна', 'аьгаьрлѝ<n><abs><add>', 'n')
('заманна', 'заман<n><abs><add>', 'n')
('телевизорлиѝ', 'телевизор<n><obl><in><ess>', 'n')
('интернетлиѝ', 'интернет<n><obl><in><ess>', 'n')
('аьгаьрлѝли', 'аьгаьрлѝли<n><abs>', 'n')
total: 22.0
recognised: 19.0
types recognised: 0.8636363636363636
fully_guessed: 3.0
fully_guessed_all: 0.13636363636363635
fully_guessed_recog: 0.15789473684210525
tags_guessed: 0.13636363636363635

def check_guesser_acc_verbs(path):
    with open(path, 'r', encoding='utf-8') as file:
        file = file.read()
        fully_guessed = 0.0
        recog = 0.0
        total = 0.0
        tags_guessed = 0.0
        words = re.findall(r"^(?P<word>[а-яōāēīȳ!"]+)\/(.*<(v)>.*)\$",
file)
        seen = []
        for w in words:
            if w not in seen:
                guesses = guess_word_acc(w[0])
                if len(guesses) != 0:
                    if w[1] in guesses:
                        fully_guessed += 1
                    else:
                        print('FAIL')
                        print('standard:', w[0], w[1])
                        print(guesses)
                guessed_tags = [re.findall(r"<.*>", l)[0] for l in
guesses]
                true_tags = re.findall(r"<.*>", w[1])[0]
                if true_tags in guessed_tags:
                    tags_guessed += 1
                recog += 1
                total += 1
                seen.append(w)

        print('total:', total)
        print('recognised:', recog)
        print('types recognised:', recog/total)
        print('fully_guessed:', fully_guessed)
        print('fully_guessed_all:', fully_guessed / total)
        print('fully_guessed_recog:', fully_guessed / recog)
        print('tags_guessed:', tags_guessed / total)

```



```
check_guesser_acc_verbs('text-1-gold.txt')
```

FAIL

```
standard: йуьч|йаьгъеч|е <IV>уьч|<IV>аьгъ<v><neg><pfv.cvb>
['йуьч|<antip><pl>гъо<v><neg><pfv.cvb>',
 'йуьч|<antip><pl>гъ<v><neg><pfv.cvb>', 'йуьч|йаьгъо<v><neg><pfv.cvb>',
 'йуьч|йаьгъе<v><neg><pfv.cvb>', 'йуьч|йаьгъи<v><neg><pfv.cvb>',
 'йуьч|йаьгъоь<v><neg><pfv.cvb>', 'йуьч|йаьгъа<v><neg><pfv.cvb>']
```

FAIL

```
standard: богъльол <III>ов<v><ant.cvb>
['б<m>гъ<v><ant.cvb>', 'бог<v><ant.cvb>', 'богъ<n><obl><dat>',
 'богъльо<n><obl><dat>', 'богъл<n><in><ess><lat>',
 'богъл<n><obl><in><ess><lat>', 'богъль<n><obl><dat>',
 'богълхъ<n><obl><dat>', 'богъль<n><dat>', 'богъль<obl><dat>',
 'богъльоа<n><obl><dat>', 'богъльое<n><obl><dat>',
 'богъльонзил<n><obl><dat>', 'богъльоам<n><obl><dat>',
 'богъльоо<n><obl><dat>', 'богъльохаь<n><obl><dat>',
 'богъльо<sup>H</sup><n><obl><dat>']
```

FAIL

```
standard: йуьч|на <II>уьч|<v><pfv.cvb>
['йуьч|на<n><pl><abs>', 'йуьч|<v><pfv.cvb>', 'йуьч|ни<n><pl><abs>',
 'йуьч<num><dstr><abs>', 'йуьча<num><dstr><abs>', 'йуьч|<num><part>',
 'йуьч|н<n><pl><abs>', 'йуьч|наь<n><pl><abs>', 'йуьч|на<v><imp>',
 'йуьч|не<v><imp>', 'йуьч|ни<v><imp>', 'йуьч|но<v><imp>',
 'йуьч|не<n><pl><abs>', 'йуьч|но<n><pl><abs>', 'йуьч|ноь<n><pl><abs>']
```

FAIL

```
standard: бойч|е <III>ов<v><neg><pfv.cvb>
['бо<v><neg><pfv.cvb>']
```

FAIL

```
standard: йиHылна <IV>иHыл<v><pfv.cvb>
['йиHыл<num><part>', 'йиHыл<v><pfv.cvb>', 'йиHылна<n><pl><abs>',
 'йиHылни<n><pl><abs>', 'йиHылн<n><pl><abs>',
 'йиHылнаь<n><pl><abs>', 'йиHыл<dstr><abs>',
 'йиHыл<num><dstr><abs>', 'йиHылна<v><imp>', 'йиHылне<v><imp>',
 'йиHылни<v><imp>', 'йиHылно<v><imp>', 'йиHылне<n><pl><abs>',
 'йиHылно<n><pl><abs>', 'йиHылноь<n><pl><abs>']
```

FAIL

```
standard: йиHилйуьголь <IV>иHыл<IV>уьго<v><simul.cvb>
['йиHыл<II>уьго<v><simul.cvb>', 'йиHыл<IV>уьго<v><simul.cvb>',
 'йиHыл<nhpl>уьго<v><simul.cvb>', 'йиHилй<I>уьго<v><simul.cvb>',
 'йиHилйуьго<n><obl><cont><ess>', 'йиHилйуь<obl><cont><ess>',
 'йиHилйуьг<n><obl><cont><ess>', 'йиHилйуьгхъ<n><obl><cont><ess>',
 'йиHилйуьг<n><cont><ess>', 'йиHилйуьг<obl><cont><ess>',
 'йиHилйуьгоа<n><obl><cont><ess>', 'йиHилйуьгое<n><obl><cont><ess>',
 'йиHилйуьгонзил<n><obl><cont><ess>',
 'йиHилйуьгоам<n><obl><cont><ess>', 'йиHилйуьгоо<n><obl><cont><ess>',
 'йиHилйуьгохаь<n><obl><cont><ess>',
 'йиHилйуьгоH<n><obl><cont><ess>', 'йиHилйуьго<n><in><ess>',
 'йиHилйуьгола<n><obl><in><ess>', 'йиHилйуьголе<n><obl><in><ess>',
 'йиHилйуьголнзил<n><obl><in><ess>',
```

```
'йи"ъилийугъолам<n><obl><in><ess>', 'йи"ъилийугъоло<n><obl><in><ess>',  
'йи"ъилийугъолхӕ<n><obl><in><ess>', 'йи"ъилийугъол" <n><obl><in><ess>']  
FAIL
```

```
standard: бӕхъна <hpl>a<pl>хъ<v><pfv.cvb>  
['бӕхъ<v><pfv.cvb>', 'бӕхъна<n><pl><abs>', 'бӕх<n><in><ess><add>',  
'бӕха<n><obl><in><ess><add>', 'бӕхе<n><obl><in><ess><add>',  
'бӕхнзил<n><obl><in><ess><add>', 'бӕхам<n><obl><in><ess><add>',  
'бӕхо<n><obl><in><ess><add>', 'бӕххӕ<n><obl><in><ess><add>',  
'бӕх" <n><obl><in><ess><add>', 'бӕхъни<n><pl><abs>',  
'бӕхънаъ<n><pl><abs>', 'бӕха<num><dstr><abs>', 'бӕхъ<num><part>',  
'бӕхън<n><pl><abs>', 'бӕхъ<dstr><abs>', 'бӕхъна<v><imp>',  
'бӕхъне<v><imp>', 'бӕхъни<v><imp>', 'бӕхъно<v><imp>',  
'бӕхъне<n><pl><abs>', 'бӕхъно<n><pl><abs>', 'бӕхъноъ<n><pl><abs>']
```

total: 29.0

recognised: 29.0

types recognised: 1.0

fully_guessed: 22.0

fully_guessed_all: 0.7586206896551724

fully_guessed_recog: 0.7586206896551724

tags_guessed: 0.7586206896551724

check_guesser_acc_verbs('text-2-gold.txt')

FAIL

```
standard: йовал <IV>ов<v><inf>  
['йо<v><antip><inf>', 'йо<v><pl><inf>', 'йо<v><pst.ptsp><obl><dat>',  
'йо<n><pl><obl><dat>', 'йо<n><obl><dat>', 'й<m>в<v><inf>',  
'йова<n><pl><obl><dat>', 'йова<n><obl><dat>', 'йо<antip><v><inf>',  
'йо<pl><v><inf>', 'йо" <n><pl><obl><dat>', 'йов<n><pl><obl><dat>',  
'йови<n><pl><obl><dat>', 'йоваъ<n><pl><obl><dat>',  
'йово<n><obl><dat>', 'йова<v><inf>', 'йове<v><inf>', 'йови<v><inf>',  
'йово<v><inf>', 'йове<n><pl><obl><dat>', 'йово<n><pl><obl><dat>',  
'йовоъ<n><pl><obl><dat>', 'йов<dat>', 'йов<n><obl><dat>',  
'йовче<n><obl><dat>', 'йованзил<n><obl><dat>', 'йовйо<n><obl><dat>',  
'йова<n><dat>', 'йоваа<n><obl><dat>', 'йоваам<n><obl><dat>',  
'йовахӕ<n><obl><dat>', 'йова" <n><obl><dat>']
```

FAIL

```
standard: йегӕйо <IV>егӕ<v><pst>
```

```
['йегӕ<v><pst>']
```

FAIL

```
standard: бегӕкъала <III>ера<v><neg.opt>  
['<III>ера<v><proh><conj>', '<hpl>ера<v><proh><conj>',  
'бере<v><proh><conj>', 'беги<v><proh><conj>', 'бего<v><proh><conj>',  
'бегоъ<v><proh><conj>', 'бегӕ"къе<v><imp><conj>',  
'бегӕ"къо<v><imp><conj>', 'бегӕкъа<n><pl><abs><conj>',  
'бегӕкъо<n><pl><abs><conj>', 'бегӕкъ<n><pl><abs><conj>',  
'бегӕкъаъ<n><pl><abs><conj>', 'бегӕкъа<n><pl><obl><genII>',  
'бегӕкъо<n><pl><obl><genII>', 'бегӕкъ<n><pl><obl><genII>',  
'бегӕкъаъ<n><pl><obl><genII>', 'бегӕкъо<n><obl><genII>',  
'бегӕкъ<n><obl><genII>', 'бегӕ"къе<v><imp><prmI>',
```

```
'бегā"кѡ<v><imp><prmI>', 'бегāкѡ<n><pl><abs>',  
'бегāкѡала<n><pl><abs>', 'бегāкѡали<n><pl><abs>',  
'бегāкѡ<v><imp><conj>', 'бегāкѡи<n><pl><abs><conj>',  
'бегāкѡи<n><pl><obl><genII>', 'бегāкѡ<n><obl><genII>',  
'бегāк<n><obl><genII>', 'бегāкѡ<v><imp><prmI>',  
'бегāкѡ"ѡа<antip><v><imp>', 'бегāкѡ"ѡал<v><imp>',  
'бегāкѡ<v><imp><conj>', 'бегāкѡи<v><imp><conj>',  
'бегāкѡе<n><pl><abs><conj>', 'бегāкѡѡ<n><pl><abs><conj>',  
'бегāкѡе<n><pl><obl><genII>', 'бегāкѡѡ<n><pl><obl><genII>',  
'бегāкѡ<genII>', 'бегāкѡче<n><obl><genII>',  
'бегāкѡанзил<n><obl><genII>', 'бегāкѡѡ<n><obl><genII>',  
'бегāкѡ<v><imp><prmI>', 'бегāкѡи<v><imp><prmI>',  
'бегāкѡе<antip>ле<v><imp>', 'бегāкѡе<pl>ле<v><imp>',  
'бегāкѡи<antip>л<v><imp>', 'бегāкѡи<pl>л<v><imp>',  
'бегāкѡало<n><pl><abs>', 'бегāкѡ<n><genII>',  
'бегāкѡаа<n><obl><genII>', 'бегāкѡаам<n><obl><genII>',  
'бегāкѡаѡ<n><obl><genII>', 'бегāкѡ"а<n><obl><genII>',  
'бегāкѡал<n><pl><abs>', 'бегāкѡалаѡ<n><pl><abs>', 'бегāкѡала<v><imp>',  
'бегāкѡали<v><imp>', 'бегāкѡало<v><imp>', 'бегāкѡале<n><pl><abs>',  
'бегāкѡалоѡ<n><pl><abs>', 'бегāкѡал<dst><abs>']  
total: 20.0  
recognised: 20.0  
types recognised: 1.0  
fully_guessed: 17.0  
fully_guessed_all: 0.85  
fully_guessed_recog: 0.85  
tags_guessed: 0.85
```