What Is

# HUNSPELL

?

# Hunspell Framework

Hunspell is a spell checker and morphological analyzer designed for languages with rich morphology and complex word compounding or character encoding and can use UTF-8 encoded dictionaries.

- Used by many popular applications.
  - Open Office
  - Mozila products - ThunderBird/ FireFox /SeaMonkey
  - Opera 10
  - Google Chrome
  - Apple's Snow Leopard

# Hunspell in Brief

Requires two files to define the language that it is spell checking.

- Dic file - a dictionary containing words for the language
- Aff file - an "affix" file that defines the meaning of special flags in the dictionary.

# Hunspell in Brief Cont...

More on Affix files.

- SET - setting the character encodings of affixes and dictionary files
- TRY - sets the change characters for suggestions
- REP - replacement table for multiple character corrections
- PFX - defines prefixes
- SFX - suffix classes

# Hunspell in Brief Cont…

```
SET UTF-8
TRY esianrtolcdugmphbyfvkwzESIANRTOLCDUGMPHBYFVKWZ'
REP 2
REP f ph REP ph f
PFX        A      Y     1
PFX        A      0      re .
SFX        B      Y      2
SFX        B      0      ed    [^y]
SFX        B      y      ied   y
```

**Aff File**

```
hello
try/B
work/AB
```

**Dic File**

"hello", "try", "tried", "work", "worked", "rework", "reworked".

File    Edit    Help

Generate    Crawl    Analyze

Generate complete list of words by entering dic/aff pair.

Single Words

Base Word     [                                            ]    [ Generate ]

Compound Syntax  [                                    ]
                 [                                    ]
                 [                                    ]

Files

Dic File    /home/buddhika/Desktop/si_LK/test/dic/test.dic    [ Browse ]

                                                              [ Generate ]

අංක
අංකය
අංකයක්
අංකයෙක්
අංකයත්
අංගෙ
අංකයත්ට
අංකල්
අංකතය
අංගෙට
අංකුර

                                          [ Analyze ]    [ Save to File ]

1.2 – Manually processing and error correcting

Category    verbs    ▼    Modify    Delete    |    [            ]    Add New Category

| Option | Strip | Append | Condition |
|--------|-------|--------|-----------|
| SFX | කි | ම | කි |
| SFX | කි | ඩ | කි |
| SFX ▼ | කි | තා | කි |

PFX

SFX

Delete Row    Add New Row    Save Changes

2. Rule building for various categories to expand words

Save Settings    Cancel

File   Edit   Help

Generate | Crawl | Analyze

Use this area to crawl a given source and extract words to generate a wordlist.

Crawl

URL   http://www.dinamina.lk/2010/08/19/   [Start] [Pause]

සජීව
ජාල
තවගමුව
මාපිටිගම
පාලම
ඡ තක
අසිතියට
ලො'කයම
විශ්ව
ගම්මානයක්
කරමට
කුඩා
වත
සමයේ
සංවර්ධනයේ
පැරැ.ම
තිසා
අප
රටේ
අසල්වැසි
ගම්මාන
පවා
දුරස්ථ
එම

[Analyze] [Save to File]

3. Word harvesting

2010/08/19/_art.asp?fn=s1008192

| General | Crawl | Parsing | Categories | Word Lists |
|---------|-------|---------|------------|------------|

**Crawl Settings**

Domain

http://www.dinamina.lk

Max Depth

5

Max Pages

1000

Max Words

100000

Crawler Name

Sinhala Crawler

3. Word harvesting (Cont)

Save Settings     Cancel

## Add Words

Use this section to add new words to the dictionary.

Word to add: ගම්මානයක     [Suggest]

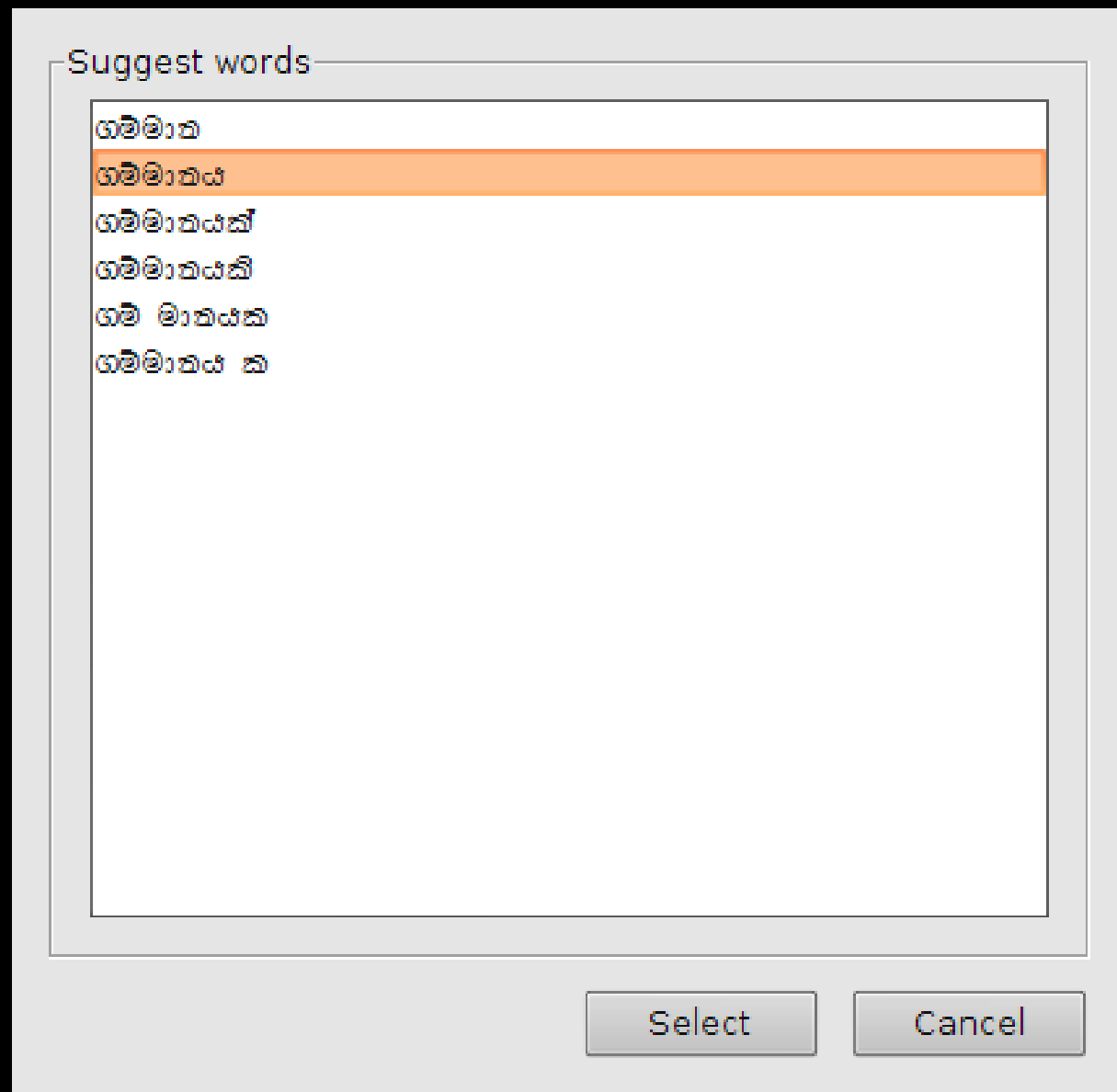Category: None ▼    [Generate]    [Modify Categories]

[Add New Word]

[Modify Selected]

[Remove Selected]

[Add Directly]    [Add From Table]    [Cancel]

4. Find Base words

## Suggest words

- ගම්මාත
- ගම්මාතය
- ගම්මාතයක්
- ගම්මාතයකි
- ගම් මාතයක
- ගම්මාතය ක

[ Select ]  [ Cancel ]

4. Find Base words (Cont)

## Add Words

Use this section to add new words to the dictionary.

Word to add: ගම්මානය          [ Suggest ]

Category: [ noun          ] [ ▼ ]   [ Generate ]   [Modify Categories]

| None |
| verbs |
| **noun** |

[ Add New Word ]

[ Modify Selected ]

[ Remove Selected ]

[ Add Directly ]   [ Add From Table ]   [ Cancel ]

5. Expanding words

## Add Words

Use this section to add new words to the dictionary.

Word to add: ගම්මාතය්

[Suggest]

Category: noun ▼ [Generate] [Modify Categories](#)

ගම්මාතය්
ගම්මාතය්
ගම්මාතය්
ගම්මාතය්ය්

[Add New Word]

[Modify Selected]

[Remove Selected]

[Add Directly] [Add From Table] [Cancel]

5. Expanding words (cont)

6. Ensuring correctness through word harvesting and frequency matching

| General | Crawl | Parsing | Categories | Word Lists |
|---------|-------|---------|------------|------------|

### Manage Lists

Word List        Blocked Word List   ▼

Highlighted Colour

Save Settings    Cancel