

Open-Tamil

Text Processing Library in Python

A Muthiah, T Srinivasan, M Annamalai



13th Tamil Internet Conference – 2014, Puducherry, India





T Shrinivasan

tshrinivasan@gmail.com

GNU/Linux Evangelist

Editor : <http://Kaniyam.com>

Ex-Co-ordinator : <http://ilugc.in>

Publisher : <http://FreeTamilEbooks.com>

Blogger : <http://Goinggnu.wordpress.com>

Coder : <http://github.com/tshrinivasan>

Photographer : <http://commons.wikimedia.org/wiki/Special>ListFiles/Tshrinivasan>

தொகுத்தல் [தொகு]

ஏனைய பாரம்பரியக் கலைக்களஞ்சியங்களைப் போல்லாது, விக்கிப்பீடியா வெளித்தொகுப்புக்களை ஏற்கிறது. எனினும், முக்கியமான அல்லது குழப்பம் விளைவிக்கும் ஆபத்துடைய சில கட்டுரைகள் தொகுக்க முடியாமல் பாதுகாக்கப்பட்டுள்ளன.

தமிழ் எழுத்துகளை

[29] மேலும், கட்டுரையை வாசிக்கும் எந்தவொரு வாசகரும் கணக்கொன்று இல்லாமலேயே கட்டுரைகளை தொகுக்க முடியும். எனினும் வெவ்வேத விழிப் புதிப் புதில் இட்டோன்று வித்தியாசமாகக் கடைப்பிடிக்கப்படுகிறது. தொடர்மைக், ஆங்கிலப் புதிலில் பதிப்புச்செய்து பார்முடுப்புது கட்டுரையொன்றை உருவாக்க முடியும்.

[30] எந்தவொரு சட்டுரையையும் அதனை கருவகியுள்ளாரா, வேறு பயன்ரோ உரிமை கொண்டாட முடியாது என்பதோடு, எந்தவொரு அங்கீகரிக்கப்பட்ட அதிகாரத் தரப்பும் அதனை ஆராய முடியாது. அதற்குப் பதிலாகத் தொகுப்பாளர்கள் தமிழடையோன கருத்தொருமிப்பின் அடிப்படையில் கட்டுரைகளின் உள்ளடக்கங்களையும் அமைப்பையும் ஏற்றுக்கொள்ள வேண்டும்.

[31] வழுமையாக, கட்டுரையொன்றில் மேற்கொள்ளப்படும் தொகுப்பு எனது படி என்று பயிற்சியிடப்படும். எனவே, இது கட்டுரைகளில் துவியமின்மை, கருத்துக் கோட்டுகள் அல்லது காப்புரிமைத் தகவல்கள் இடம்பெறலாக. ஒவ்வொரு மூத்துப் பதில்களுக்கும் வெவ்வேறு நிலைக்கூட்டுப்பாடுகளைக் கொண்டிருப்பதுடன், இக்கொள்கைகளிலும் திருத்தங்களைக் கொண்டுவரலாம். உதாரணமாக, செருமானிய விக்கிப்பீடியாவில் கட்டுரைத் தொகுப்புக்கள் சில மேற்பார்வையிடல்களுக்குப் பின் உறுதிப்படுத்தப்படுகின்றன.

[32] பல்வேறு சோதனை ஒட்டங்கள் மற்றும் கலந்துரையாடல்களுக்குப் பின் டிசம்பர் 2012 அன்று "மாற்றங்களுக்கான காத்திருப்பு" முறைமை ஆங்கில விக்கிப்பீடியாவில் அறிமுகப்படுத்தப்பட்டது.

[33] இம்முறைமையின் கீழ், சர்ச்சைக்குரிய அல்லது குழப்பம் விளைவிக்கக்கூடிய ஆபத்துடைய கட்டுரைகளின் மீதான புதிய பயனர்களின் தொகுப்புக்கள், அங்கீகரிக்கப்பட்ட ஒரு விக்கிப்பீடியா பயனரின் மேற்பார்வையின் பின்னரே வெளியிடப்படும்.

விக்கிப்பீடியாவுக்கு உதவும் மென்பொருட்கள் பங்களிப்பாளர்களுக்கு உதவிகரமாக இருக்கும். ஒவ்வொரு கட்டுரையிலும் காணப்படும் "வரலாற்றைக் காட்டவும்" பக்கம் திருத்தங்களைப் (திருத்தங்கள் அவதுறான தகவல்கள், குற்ற அச்சுறுத்தல் அல்லது காப்புரிமை மீறல் போன்றன மீளமைக்கப்படக் கூடியனவாய் இருப்பினும்) பதிவு செய்யும். இப் பக்கத்தைப் பயன்படுத்தவதன் மூலம் பயனர்கள், விரும்பத்தகாத தொகுப்புக்களை மீளமைக்கவோ அல்லது இழக்கப்பட்ட தகவல்களை மீளப்பெறவோ முடியும். ஒவ்வொரு கட்டுரைக்குமான பேச்சுப்பக்கம் பல்வேறு பயனர்களும் தமிழுளர்களை நடத்தப்பட உதவுகிறது.

[34] முக்கியமாக, தொகுப்பாளர்கள் பேச்சுப்பக்கத்தைப் பயன்படுத்தி கருத்தொருமிப்புப் பெற முடியும். சிலவேளைகளில் இதற்காக வாக்கெடுப்பும் நடத்தப்படும்.

“ பிரபலமான நகைச்சவையொன்று கூறுகிறது, "விக்கிப்பீடியாவிலுள்ள பிரச்சினை என்னவென்றால் அது பயன்பாட்டு ரீதியில் சிறந்தது. ஆனால், கொங்கா தீயில் பயனற்று." ”

—மீக்கா ஸியோக்காக [28]

நிறுவுதல்

- Python package
 - Python Package installer (pip)
 - <https://pypi.python.org/pypi/Open-Tamil/>
- Git-Hub collaboration
 - Open-Tamil core repo
<https://github.com/arcturusannamalai/open-tamil/>
- Social blogs
 - <http://ezhillang.wordpress.com/>

ഉയിരമുത്തുക്കൾ

ଅ	ଆ	ଇ	ା	ୁ	ୟା	ଏ	ୟ	ଐ	ଓ	୦	ୟୋ	ଔଯା
a	ā	i	ī	u	ū	e	ē	ai	o	ō	au	
[a]	[a:]	[i]	[ī]	[u,w]	[ū]	[e]	[ē]	[ay]	[o:]	[ō:]	[au]	

எழுத்துக்களைக் கண்டறிதல்

அ	ஆ	இ	ஈ	உ	ஊ	எ	ஏ	ஐ	ஓ	ஔ	ஓள்
a	ā	i	ī	u	ū	ē	ē	ai	ō	ō	au
[a]	[ā:]	[i]	[ī:]	[u, y]	[ū:]	[ē]	[ē:]	[ai]	[ō]	[ō:]	[au]

Grantha letters

க	k [k, g, x, y, h]	த	t [t, d, ð]	ல	l [l]	வி	v [v]	ஷி	ʃ [ʃ]
---	-------------------	---	-------------	---	-------	----	-------	----	-------

ஙி	ŋ [ŋ]	ஞி	ñ [ñ]	வி	v [v]	ஷி	s [s]
----	-------	----	-------	----	-------	----	-------

ச	c [tʃ, dʒ, ʃ, s]	ஞ	p [p, b, ð]	ழி	z, ð, r [ɹ]	ஸி	s [s]
---	------------------	---	-------------	----	-------------	----	-------

ஞி	ñ [ñ]	மி	m [m]	ளி	! [!]	ஹி	h [h]
----	-------	----	-------	----	-------	----	-------

ஞி	t [t, d, ð]	யி	y [j]	றி	r, R [r, t, ð]	கஷி	kṣ [kṣ]
----	-------------	----	-------	----	----------------	-----	---------

னி	n [n]	ரி	r [r]	னி	ñ, N [ñ]
----	-------	----	-------	----	----------

ஃஃ = äytam - turns p into f and j into z, e.g. பிஃஃ fi [fi:] ஜிஃஃ zi [zi]

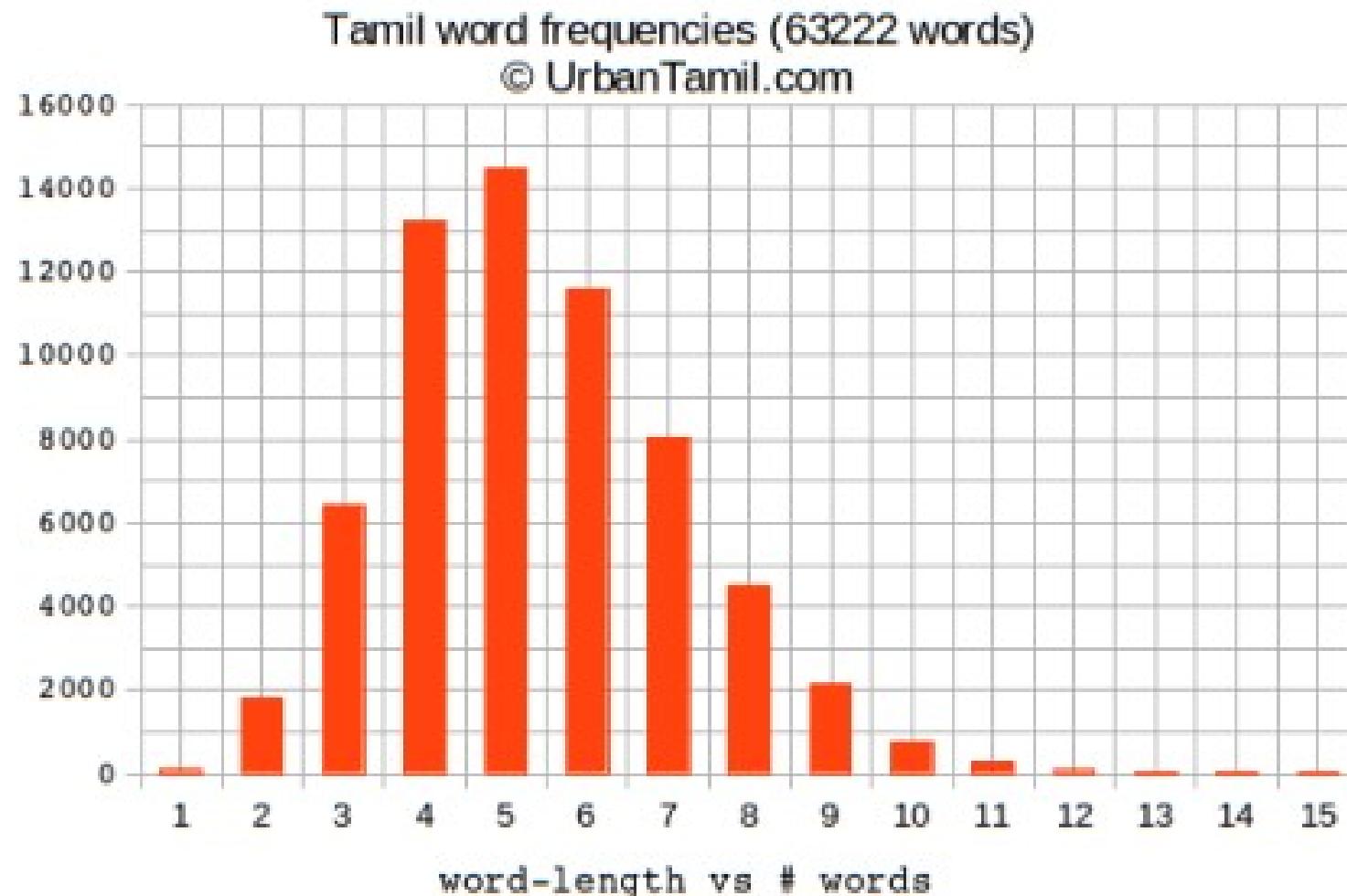
வார்த்தைகளின் நீளம் காணுதல்





வகைப்படுத்துதல்

வார்த்தைகளின் பயன்பாடு



ഉത്തര - IPA

IPA symbols for Tamil Letters

அ	ஆ	இ	ஏ	உ	ஊ
a	ā	i	ē	u	ū

எ	ஏ	ஐ	ஓ	ஃ	ஓா
e	ē	ai	o	ō	au

க	ங	ச	ஞ	ஞ	ண்ட
ka	ṅa	ca	ña	ṭa	ṇa

த	ந	ஞ	ம	ஞ	ர
ta	na	pañ	ma	ña	ra

உரை - IPA

```
shrinivasan@shrinivasan-laptop:~/Dev/open-tamil/examples/txt2ipa$ python demo_txt2ipa.py
input unicode text வணக்கம் தமிழகம்
after ipa பளி்க்கலம் தளவிளக்கலம்
after broad பனிக்கம் தனவிளக்கம்
```

எழுத்துரு மாற்றம்

"அச்சோடு! அவர்கள் எடுப்பது எனக்கு பெரிதல். ஆனால் அவதாரங்களின் பெற்றோரைப் படைப்பதில் நீ ஒத்துழைக்க வேண்டும். இராமாவதாரத்தில் அன்னையாக இனியவள் கோசலையைப் படைத்தாய். கிருஷ்ணவதாரத்தில் என் தாயாக நல்லவள் வகுதேவியையும், என் வளர்ப்புத் தாயாக களங்கமில்லா ஆச்ச படைத்தாய். ஆனால் கவியுகத்தில் நீ படைத்துதெல்லாம் சுதா, ஜெயா, ரமா, ரேகா, உஷா போன்றவர்கள். முதலில் நீ எனக்கு கோசலையைப் போல், வகுதே சோதையைப் போல் ஒரு தாயைப் படைத்து டூலோகத்தில் உலவ விடு. அதன் பிறகு கல்கி அவதாரம் பற்றி பேசு" உணர்ச்சிப் பெருக்கில் பேசி விட்டு உலகளந்து ஒரேயே மொட்டு அமர்ந்தான்.

பிரம்மதேவா! அவதாரம் எடுப்பது எனக்கு பெரிதல். ஆனால் அவதாரங்களின் பெற்றோரைப் படைப்பதில் நீ ஒத்துழைக்க வேண்டும். இராமாவதாரத்தில் அன்னையாக இனியவள் கோசலையைப் படைத்தாய். கிருஷ்ணவதாரத்தில் என் தாயாக நல்லவள் வகுதேவியையும், என் வளர்ப்புத் தாயாக களங்கமில்லா ஆச்ச படைத்தாய். ஆனால் கவியுகத்தில் நீ படைத்துதெல்லாம் சுதா, ஜெயா, ரமா, ரேகா, உஷா போன்றவர்கள். முதலில் நீ எனக்கு கோசலையைப் போல், வகுதே சோதையைப் போல் ஒரு தாயைப் படைத்து டூலோகத்தில் உலவ விடு. அதன் பிறகு கல்கி அவதாரம் பற்றி பேசு" உணர்ச்சிப் பெருக்கில் பேசி விட்டு உலகளந்து ஒரேயே மொட்டு அமர்ந்தான்.

பிரம்மதேவன் வெட்கத்தால் தனது நான்கு தலைகளும் குனிய தேவர்களாகிய எங்களோடு வைகுந்தத்தை விட்டு விலகினான்.

மறுபடி பழையபடி திருமகள் பாதம் வருட, கருநாகம் குடை பிடிக்க, கார்முகில்வண்ணன் கண்களை மூடி முடிவில்லா யோகமாம் நித்திரையில் மூழ்கினான்.

எழுத்துரு மாற்ற வகைகள்

- | | |
|----------------|----------------|
| 1. Anjal | 13. Tam |
| 2. Bhamini | 14. Tscii |
| 3. Boomi | 15. Pallavar |
| 4. Dinakaran | 16. Indoweb |
| 5. Dinamani | 17. Koeln |
| 6. Dinathanthy | 18. Libi |
| 7. Kavipriya | 19. Oldvikatan |
| 8. Murasoli | 20. Webulagam |
| 9. Mylai | 21. Diacritic |
| 10. Nakkeeran | 22. Shreelipi |
| 11. Roman | 23. Softview |
| 12. Tab | 24. Tace |
| | 25. Vanavil |

தீருங்குறி மாற்றம்

Tamil^{[1][2]}

Official Unicode Consortium code chart [PDF](#)

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
U+0B8x			ஃ	ஃ		அ	ஃ	இ	ஃ	ஃ	உள்				எ	ஏ
U+0B9x	ஃ		ஃ	ஃ	ஓள்	க				வ	ச		ஃ	ஞ	ஞ	ஞ
U+0BAx				ண	த				ந	ன	ப			ம	ஷ	
U+0BBx	ர	ற	ல	ள	ழ	வ	ஸ	ஷ	ஸ	ஹ				ஊ	ஒ	
U+0BCx	ஃ	ஃ	ஃ	ஃ			ஃ	ஃ	ஃ	ஃ	ஃ	ஃ	ஃ	ஃ	ஃ	
U+0BDx	ஃ						ஓள்									
U+0BEx						ஓ	க	உ	ங	ச	ஈ	க	ங	ஏ	ஏ	ஏ
U+0BFx	ஃ	ா	த	உ	ஈ	ங	ய	ங	ங	ங	ங					

Notes

1.^ As of Unicode version 7.0

2.^ Grey areas indicate non-assigned code points

ngram

Google books Ngram Viewer

Graph these comma-separated phrases: case-insensitive

between and from the corpus with smoothing of



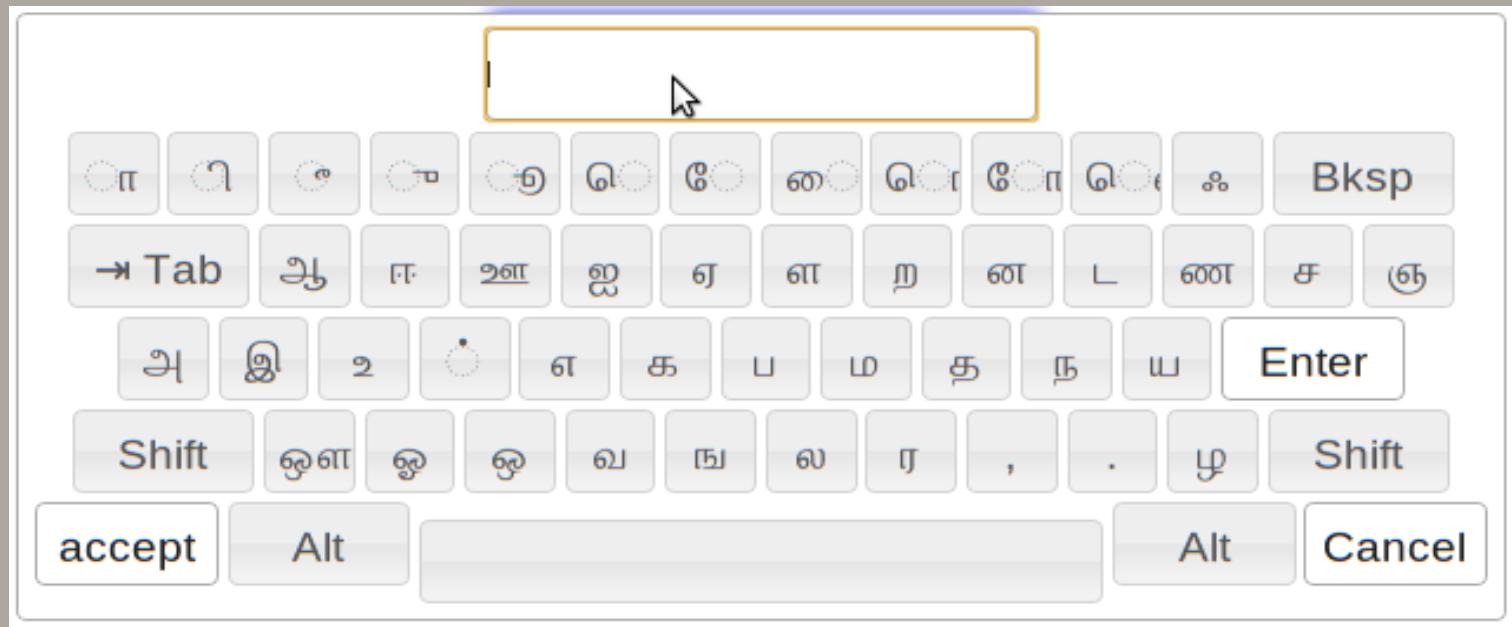
ଓଲିପେଯର୍‌ପ୍ପ

The vowels:											
a	ଅ	aa, A	ଆ	i	ଇ	ee, I	ଉ	u	ଔ	oo, U	ଓ
e	ଏ	ae, E	ଏ	ai	ଇ	o	ଓ	oa, O	ା	au	ାଳେ
Ahh, H	ହ										

The Consonants:										
g, k, kh, c	କ	nG	ଙୁ	ch	ଚ	j	ଝୁ	nY	ଞୁ	ଞୁ
d, t	ଟ	nN	ଣୁ	dh, th	ତୁ	N	ନୁ	n	ଣୁ	ଣୁ
b, bh	ପ	m	ମ	y	ପୁ	r	ରୁ	R	ରୁ	ରୁ
l	ଲ	L	ଲୁ	zh	ଛୁ	v, w	ବୁ			
sh	ଷୁ	s	ସୁ	h	ହୁ	f	ଫୁ	ଫୁ		

வார்த்தைகளை
திருப்புதல்

தமிழ்99 தட்டச்சுப் பலகை



- JQuery
- JQuery UI based
- Free to use on web
- e.g. Www.Urbantamil.com

கட்டற்ற/திறவுற்று மென்பொருள்

- Multi-licensed
 - MIT, and other OSS
- Multi-language
 - C, Python, JavaScript, C
-



Source : <https://github.com/arcturusannamalai/open-tamil>

பயன்பாடு

1. Websites:

- 1. Ezhil Language
- 2. UrbanTamil

2. Installs on Python

- 1. 1000+ downloads on PIP

எழில்
நிரல்
மொழி

தமிழ் மொழி
நிரல்

நவீன், அகராதி
சமுக,

UrbanTamil.com

நகர்ப்புற அகராதி

பங்களிப்போர்



muthuspost@gmail.com &

tshrinivasan@gmail.com



msathia@gmail.com



arulalant@gmail.com

உதாரணங்கள்

எழுத்துக்களை எண்ணுதல்

```
>>> import tamil  
>>> print len(tamil.utf8.get_letters(u'திருவாலவாயுடையார்திருவிலையாடற்புராணம்'))  
20
```

Transliterate

azhagi>> thamizh vaazhga
தமிழ் வாழ்க

வார்த்தைகளை திருப்புதல்

```
>>> print tamil.utf8.reverse_word(u'உலகத் தமிழ் இணைய மாநாடு')  
மாநாடு இணைய முழுமொத்த வினாக்களை எடுத்து
```

எழுத்துகளாகப் பிரித்தல்

```
>>> a = tamil.utf8.get_letters(u'உலகத் தமிழ் இணைய மாநாடு')
>>> for letter in a:
...     print letter
...
உ
ல
க
த
கி
ங்
இ
ணை
ய
மா
நா
டு
```

வார்த்தைகளாகப் பிரித்தல்

```
>>> sentence = tamil.utf8.get_words(u'உலகத் தமிழ் இணைய மாநாடு')
>>> for word in sentence:
...     print word
...
உலகத்
தமிழ்
இணைய
மாநாடு
```

எழுத்து இடத்தைக் கண்டுபிடித்தல்

```
>>> print tamil.utf8.word_intersection( u'தேடுக' , u'தடங்கல்' )  
[(2, 3)]
```

தில்கி - ஓரங்குறி மாற்றம்

```
shrinivasan@shrinivasan-laptop:~/Dev/open-tamil/examples/txt2unicode$ python demo_tscii2utf8.py  
tscii %c0A6UÀ÷  
«ÓÇçÀ %c0iÌÈÙ  
unicode திருவள்ளுவர்  
அருளிய திருக்குறள்
```

ழுரங்குறி - தில்கி மாற்றம்

```
shrinivasan@shrinivasan-laptop:~/Dev/open-tamil/examples/txt2unicode$ python demo_utf8_2_tscii.py  
tscii original input %c0A8ÙÀ÷ «ÓÇ¢À %c0iÌÈû  
from tscii2unicode திருவள்ளுவர் அருளிய திருக்குறள்  
from unicode2tscii %c0A8ÙÀ÷ «ÓÇ¢À %c0iÌÈû
```

உரை - IPA

```
shrinivasan@shrinivasan-laptop:~/Dev/open-tamil/examples/txt2ipa$ python demo_txt2ipa.py
input unicode text வணக்கம் தமிழகம்
after ipa ॥வி॒த்'ங்கலம் ॥வா॒யிர்லா॒ம்
after broad ॥வெ॒யிர்க்கம் ॥வெ॒யிர்க்கே॒ம்
```



நன்றி

கிரியேடிவ் காமன்ஸ் படங்கள் மூலம்

- <http://upload.wikimedia.org/wikipedia/commons/d/d1/Tamil-Encoding-UnicodePUA-TACE16-chart.png>
- <http://pixabay.com/p-2261>
- http://upload.wikimedia.org/wikibooks/ta/c/c8/Tamil_vwl.gif
- http://upload.wikimedia.org/wikibooks/ta/f/f8/Tamil_cons.gif
- http://upload.wikimedia.org/wikipedia/commons/a/ae/Metal_movable_type.jpg
- http://kamalasurabhi.org/basic_lessons/ipa_tamil.gif
- http://en.wikipedia.org/wiki/Tamil_script
- https://books.google.com/ngrams/graph?content=tamil&year_start=1900&year_end=2014&corpus=15&smoothing=0&share=&direct_url=t1%3B%2Ctamil%3B%2Cc0
- http://www.tamildictionary.org/tamil_transliteration.php
- <http://blog.ravidreams.net/wp-content/uploads/2006/12/tamil99.jpg>
- http://commons.wikimedia.org/wiki/File:Nandri_%E0%AE%A8%E0%AE%A9%E0%AF%8D%E0%AE%B1%E0%AE%BF%29.png
- <http://about.me/SathiaNMahadevan>
- <http://pixabay.com/p-96286>

Creative Commons Attribution Share-Alike License

