

ဒီ Jupyter Notebook က GitHub မှာ ကျွန်တော်တင်ပေးထားတဲ့ Sylbreak Python ပရိုဂရမ် <https://github.com/ye-kyaw-thu/sylbreak/blob/master/python/sylbreak.py> (<https://github.com/ye-kyaw-thu/sylbreak/blob/master/python/sylbreak.py>) ကို Jupyter Notebook, Python 3 Kernel မှာ function တစ်ခုအနေနဲ့ဆောက်ပြီး သုံးတဲ့ပုံစံကို နမူနာအနေနဲ့ ပြသထားတာ ဖြစ်ပါတယ်။

```

1 # Regular Expression Python Library ကို သုံးလို့ရအောင် import လုပ်တာ
2 import re
3
4 # စာလုံးတွေကို အုပ်စုဖွဲ့တာ (သို့) variable declaration လုပ်တာ
5 # တကယ်လို့ syllable break လုပ်တဲ့ အခါမှာ မြန်မာစာလုံးချည်းပဲ သပ်သပ် လုပ်ချင်တာဆိုရင် enCh
6 myConsonant = "က-အ"
7 enChar = "a-zA-Z0-9"
8 otherChar = "ဣတ္ထိဉ္ဇီဧဩဠြေဿဋ္ဌ၍၏ဝ-၉။၊-/:-@[ -`{ -~\s"
9 ssSymbol = '၊'
10 ngaThat = 'င'
11 aThat = 'အ'
12
13 # Regular expression pattern for Myanmar syllable breaking
14 # *** a consonant not after a subscript symbol AND
15 # a consonant is not followed by a-That character or a subscript sym
16 # မြန်မာစာကို syllable segmentation လုပ်ဖို့အတွက်က ဒီ RE pattern တစ်ခုတည်းနဲ့ အဆင်
17 BreakPattern = re.compile(r"((?<!" + ssSymbol + r")["+ myConsonant +
18
19 # sylbreak function ဆောက်တဲ့ အပိုင်း
20 def sylbreak(line):
21     line = re.sub(r"\s+", "", line)
22     line = BreakPattern.sub(r" " + r"\1", line)
23     return line

```

```
In [2]: 1 # sylbreak function ကိုခေါ်သုံးကြည့်ရအောင်
        2
        3 sylbreak("မြန်မာစာသည် တိုစာ။ တိုစာကို သုတေသန လုပ်ပါ။")
```

```
Out[2]: ' မြန်မာ စာ သည် တို့ စာ ။ တို့ စာ ကို သု တေ သ န လုပ် ပါ ။ '
```

စိတ်ထဲမှာ ပေါ်လာတာကို ကောက်ရေး(ပြီးတော့ syllable segmentation လုပ်ခိုင်းလိုက်တာပါ။ :)

နောက်ထပ် ဥပမာအနေနဲ့ [Wikipedia Myanmar](https://my.wikipedia.org/wiki/%E1%80%A1%E1%80%AC%E1%80%81%E1%80%AE%E1%80%9) မှာရေးထားတဲ့ အာဒိမီးဒီးစ် ရဲ့ အတ္ထုပ္ပတ္တိအကျဉ်း (<https://my.wikipedia.org/wiki/%E1%80%A1%E1%80%AC%E1%80%81%E1%80%AE%E1%80%9> ထဲမှာရေးထားတဲ့ စာကြောင်းတွေကို sylbreak နဲ့ ဖြတ်ကြည့်ရအောင်။

```
In [3]: 1 syllbreak("အာဒီမီးဒီးစ်ကို ဘီစီ ၂၈၇ ခန့်က ရှေးဟောင်း မဂ္ဂနာဂရေစီယာပြည်လက်အောက်ခံ စစ္စလီပြည်
2
3 " ")
```

Out[3]: ' အာ ဒီ မီး ဒီးစ် ကို ဘီ စီ ၂ ၈ ၇ ခ န် က ရှေး ဟောင်း မဂ္ဂ နာ ဂ ရေ စီ ယာ ပြည် လက် အောက် ခံ စစ္စ လီ ပြည် ဆိုင် ရာ ကျူးစ် မြို့ တွင် မွေး ဖွား ခဲ့ သည် ။ ဘိုင် ဇန် တိုင်း ဂ ရီ ခေတ် က သ မိုင်း ပ ညာ ရှင် ဂျန် ဇီ ဇီ ၏ မှတ် တမ်း အ ရ အာ ဒီ မီး ဒီးစ် သည် အ သက် ၇ ၅ နှစ် အ ထိ နေ ထိုင် သွား ရ ကြောင်း သိ ရ သည် ။ အာ ဒီ မီး ဒီးစ် သည် သူ ၏ တီ ထွင် မူ တစ် ခု ဖြစ် သော သဲ နာ ရီ နှင့် ပတ် သက် ၍ ရေး သား ထား သော S a n d R e c k o n e r s အ မည် ရှိ စာ တမ်း များ တွင် သူ ၏ ဖ ခင် အ မည် ကို နက္ခတ္တ ဗေဒ ပ ညာ ရှင် ဖီး ဒီး ယပ်စ် ဟု ဖော် ပြ ထား သည် ။ သ မိုင်း ပ ညာ ရှင် ပ လူး တပ် ရေး သား သော ခေတ် ပြိုင် ပုဂ္ဂိုလ် ထူး ကြီး များ စာ အုပ် တွင် အာ ဒီ မီး ဒီးစ် သည် ဆိုင် ရာ ကျူးစ် ဘု ရင် ဒု တိ ယ မြောက် ဟီ ရိုး နှင့် ဆွေ မျိုး တော် စပ် ကြောင်း ဖော် ပြ ထား သည် ။ သူ ငယ် ရွယ် စဉ် က အီ ဂျစ် ပြည် အ လက် ဇနီး ယား မြို့ တွင် ပ ညာ ဆည်း ပူး ခဲ့ သည် ဟု ယူ ဆ ရ သည် ။ ဘီ စီ ၂ ၁ ၂ တွင် အာ ဒီ မီး ဒီးစ် သေ ဆုံး ခဲ့ သည် ။ ရောမ စစ် ဗိုလ် ချုပ် မား ကပ်စ် က လ ဝေ ဒီး ယပ်စ် မာ ဆဲ လပ်စ် က နှစ် နှစ် ကြာ ဝိုင်း ရံ ပိတ် ဆို့ ပြီး နောက် ဆိုင် ရာ ကျူးစ် မြို့ ကို သိမ်း ပို က် လိုက် သည် ။ ထို အချိန် တွင် အာ ဒီ မီး ဒီးစ် သည် ရော် မက် ထ ရီ ပုစ္ဆာ တစ် ပုဒ် ကို စဉ်း စား အ ဖ ြေ ရှာ နေ ခိုက် ဖြစ် သည် ။ ရောမ စစ် သား က သူ့ အား ဖမ်း ဆီး လိုက် ပြီး ဗိုလ် ချုပ် မာ ဆဲ လပ်စ် နှင့် တွေ့ ဆုံ ရန် ပြော ဆို ရာ သူ က သူ ၏ ပုစ္ဆာ စဉ်း စား နေ ဆဲ ဖြစ် ၍ မ တွေ့ လို ကြောင်း ငြင်း ဆို သည် တွင် ရောမ စစ် သား က ဒေါ သ ထွက် ကာ ဓား ဖြင့် ထိုး သတ် လိုက် သည် ဟု ပ လူး တပ် က ရေး သား ခဲ့ သည် ။ ဗိုလ် ချုပ် မာ ဆဲ လပ်စ် သည် အာ ဒီ မီး ဒီးစ် သေ ဆုံး သွား သ ည့် အ တွက် များ စ ဝှာ နှ မြော တ သ ဖြစ် ရ သည် ။ အာ ဒီ မီး ဒီးစ် အား ပ ညာ ရှင် တစ် ယောက် အ ဖြစ် သိ ရှိ ထား သ ဝေ ကြောင့် မ သတ် ရန် ကြို တင် အ မိ န် ပေး ထား ခဲ့ သည် ။ “ ငါ့ စက် ဝိုင်း တွေ ပေါ် တက် မ နင်း ပါ နဲ့ ” ဟု သော စ ကား ကို အာ ဒီ မီး ဒီးစ် နောက် ဆုံး ပြော ဆို ခဲ့ သည် ဟု အချို့ က ယူ ဆ ကြ သ ဝေ လည်း သ မိုင်း ပ ညာ ရှင် ပ လူး တပ် ရေး သော စာ အုပ် တွင် မူ မ ပါ ရှိ ပေ ။ အာ ဒီ မီး ဒီးစ် ၏ ဂူ ဗိမ္မာန် တွင် ထု လုံး ရှည် မှန် တစ် ခု အ တွင်း စက် လုံး တစ် ခု ကို ထ ည့် သွင်း ထား သ ည့် ရုပ် တု တစ် ခု ကို စိုက် ထူ ထား သည် ။ အာ ဒီ မီး ဒီးစ် သေ ဆုံး ပြီး နှစ် ပေါင်း ၁ ၃ ၇ နှစ် အ ကြာ ဘီ စီ ၇ ၅ တွင် ရောမ ခေတ် နိုင် င် ရေး သု ခ မိန် ဆီ ဇာ ရီ က အာ ဒီ မီး ဒီးစ် အ ကြောင်း ကြား သိ ရ ၍ သူ ၏ အုတ် ဂူ အား ရှာ ဖွေ ခဲ့ သည် ။ ခြုံ နွယ် ပိတ် ပေါင်း များ ဖုံး အုပ် နေ သော အာ ဒီ မီး ဒီးစ် ၏ အုတ် ဂူ ကို ဆိုင် ရာ ကျူးစ် မြို့ အ နီး တွင် ရှာ ဖွေ တွေ့ ရှိ ခဲ့ ပြီး သ န် ရှင်း ရေး ပြု လုပ် ကာ အုတ် ဂူ ပေါ် မှ စာ သား များ ကို ဖတ် ရှု သွား သည် ။ ဆိုင် ရာ ကျူးစ် စစ် ပွဲ အ ပြီး နှစ် ပေါင်း ၇ ၀ အ ကြာ တွင် ပို လီး ဘီး ယပ်စ် ရေး သား သော ဆိုင် ရာ ကျူးစ် စစ် ပွဲ အ ကြောင်း စာ အုပ် တွင် အာ ဒီ မီး ဒီးစ် နှင့် ပ တ် သက် သော အ ကြောင်း များ ပါ ရှိ ၍ သ မိုင်း ပ ညာ ရှင် ပ လူး တပ် က ထပ် မံ ရေး သား နိုင် ခဲ့ ခြ င်း ဖြစ် ပါ သည် ။ ဆိုင် ရာ ကျူးစ် မြို့ ကာ ကွယ် ရေး အ တွက် စစ် ပွဲ ဝင် စက် ကိ ရီ ယာ လက် နက် ဆန်း များ ကို လည်း အာ ဒီ မီး ဒီးစ် က တီ ထွင် ပေး ခဲ့ ကြောင်း အ ဆို ပါ စာ အုပ် တွင် ဖော် ပြ ပါ ရှိ ပါ သည် ။ '

Typing order

မြန်မာစာနဲ့ ပတ်သက်တဲ့ NLP (Natural Language Processing) အလုပ် တစ်ခုခု လုပ်ဖို့အတွက် syllable segmentation လုပ်ကြမယ်ဆိုရင် တကယ်တမ်းက မလုပ်ခင်မှာ၊ မြန်မာစာ စာကြောင်းတွေရဲ့ typing order အပါအဝင် တခြား ဖြစ်တတ်တဲ့ အမှားတွေကိုလည်း cleaning လုပ်ရပါတယ်။ အဲဒီလိုမလုပ်ရင် syllbreak က ကျွန်တော် အကြမ်းမျဉ်း သတ်မှတ်ထားတဲ့ မြန်မာစာ syllable unit တွေအဖြစ် မှန်မှန်ကန်ကန် ဖြတ်ပေးနိုင်မှာ မဟုတ်ပါဘူး။ မြန်မာစာ စာကြောင်း တွေထဲမှာ ရှိတတ်တဲ့အမှား တွေက တကယ့်ကို အများကြီးပါ။ တချို့ အမှားတွေက မျက်လုံးနဲ့ကြည့်ယုံနဲ့ မခွဲခြားနိုင်တာမျိုး တွေလည်း ရှိပါတယ်။ ဒီနေရာမှာတော့ အမှားအမျိုးအစားတွေထဲက တစ်မျိုးဖြစ်တဲ့ typing order အမှား တစ်မျိုး၊ နှစ်မျိုးကို ဥပမာအနေနဲ့ရှင်းပြရင်း၊ အဲဒီလိုအခြေအနေမျိုးမှာ ဖြစ်တတ်တဲ့ syllbreak က ထွက်လာမယ့် အမှား output တွေကိုလည်း လေ့လာကြည့်ကြရအောင်။

အောက်မှာ သုံးပြထားတဲ့ "ခန့်" က "ခ န ့်" (ခခွေး နငယ် အောက်မြစ် အသတ်) ဆိုတဲ့ အစီအစဉ် အမှားနဲ့ ရိုက်ထားတာဖြစ် ပါတယ်။ အဲဒါကြောင့် syllbreak က ထွက်လာတဲ့အခါမှာ "ခခွေး" နဲ့ "နငယ် အသတ် အောက်မြစ်" က ကွဲနေတာဖြစ်ပါတယ်။

```
In [4]: 1 sylbreak("ဘီစီ ၂၈၇ ခန့်")
```

```
Out[4]: ' ဘီ စီ ၂ ၈ ၇ ခ န် '
```

တကယ်တမ်း မှန်ကန်တဲ့ "ခန့်" ရဲ့ typing order က "ခ န ျ ့" (ခခွေး နငယ် အသတ် အောက်မြစ်) ပါ။
အမြင်အားဖြင့်ကတော့ မခွဲနိုင်ပေမဲ့၊ မှန်ကန်တဲ့ typing order နဲ့ ရိုက်ထားရင်တော့ "ခန့်" ဆိုပြီး syllable တစ်ခုအနေနဲ့ ရိုက်ထုတ်ပြပေးပါလိမ့်မယ်။

```
In [5]: 1 sylbreak("ဘီစီ ၂၈၇ ခန့်")
```

```
Out[5]: ' ဘီ စီ ၂ ၈ ၇ ခန့် '
```

နောက်ထပ် typing order အမှားတစ်ခုကို ကြည့်ကြရအောင်။

```
In [6]: 1 sylbreak("ထည့်သွင်းထားသည့်ရုပ်တု")
```

```
Out[6]: ' ထ ည့် သွင်း ထား သ ည့် ရုပ် တု '
```

"ညကြီး အောက်မြစ် အသတ်" ဆိုတဲ့ မှားနေတဲ့ အစီအစဉ်ကို "ညကြီး အသတ် အောက်မြစ်" ဆိုပြီး ပြောင်းရိုက်ပြီးတော့ sylbreak လုပ်ကြည့်ရင်တော့ အောက်ပါအတိုင်း "ထ" နဲ့ "ည့်", "သ" နဲ့ "ည့်" တွေက ကွဲမနေတော့ပဲ မှန်မှန်ကန်ကန်ဖြတ်ပေး ပါလိမ့်မယ်။

```
In [7]: 1 sylbreak("ထည့်သွင်းထားသည့်ရုပ်တု")
```

```
Out[7]: ' ထည့် သွင်း ထား သည့် ရုပ် တု '
```

တချို့အမှားတွေကတော့ ဂရုစိုက်ရင် မျက်စိနဲ့ မြင်နိုင်ပါတယ်။
ဥပမာ "ဥ" (အက္ခရာ ဥ) နဲ့ "ဥ" (ညကလေး) ကိုမှားရိုက်တဲ့ကိစ္စပါ။

သို့သော် ကျွန်တော်မြန်မာစာကြောင်းတွေအများကြီးကို ကိုင်တွယ်အလုပ်လုပ်တဲ့အခါတိုင်းမှာ ဒီလိုအမှားက အမြဲတမ်းကို ပါ တတ်ပါတယ်။

ဖောင့် (font) မှာလည်း မှန်မှန်ကန်ကန်ခွဲထားမယ်ဆိုရင်၊ အမှန်က ညကလေးဆိုရင် အမြီးက ရှည်ပါတယ်။ စာရိုက်သူအများစု က သတိမပြုမိတဲ့ အကြောင်းအရင်း တစ်ခုကလည်း တချို့ text editor တွေမှာ "အက္ခရာ ဥ" နှင့် ညကလေး "ဥ" ကို ကွဲပြား အောင် မပြသပေးနိုင်လို့ပါ။

```
In [8]: 1 sylbreak("ကာရီသည်ဒီနှစ်၏ပါရမီရှင်တစ်ဦးနှင့်ထိုက်တန်သောအမျိုးသမီးအဆိုရှင်ဖြစ်သည်။")
```

```
Out[8]: ' ကာ ရီ သည် ဒီ နှစ် ၏ ပါ ရ မီ ရှင် တစ် ဦး နှ င့် ထိုက် တန် သော အ မျိုး သ မီး အ ဆို ရှင် ဖြစ် သ ည့် ။ '
```

ဝီကီပီးဒီးယားက မှားနေတဲ့ "ညကလေး" ကို "အက္ခရာ ဥ" နဲ့ပြန်ပြင်ရိုက်ထားတဲ့ စာကြောင်းနဲ့ နောက်တစ်ခေါက် syllable ဖြတ်ထားတာက အောက်ပါအတိုင်းဖြစ်ပါတယ်။ "ညကလေး" နဲ့ "အက္ခရာ ဥ" အမှားကိစ္စမှာတော့ syllable segmentation ဖြတ်တဲ့အပိုင်းမှာတော့ ထူးထူးခြားခြား အပြောင်းအလဲ မရှိပါဘူး။

```
In [9]: 1 sylbreak("ကာရီသည်ဒီနှစ်၏ပါရမီရှင်တစ်ဦးနှင့်ထိုက်တန်သောအမျိုးသမီးအဆိုရှင်ဖြစ်သည်။")
```

```
Out[9]: ' ကာ ရီ သည် ဒီ နှစ် ၏ ပါ ရ မီ ရှင် တစ် ဦး နှ င့် ထိုက် တန် သော အ မျိုး သ မီး အ ဆို ရှင် ဖြစ် သ ည့် ။ '
```

Note

- syllbreak မှာ သုံးထားတဲ့ မြန်မာစာ syllable unit (အဖြတ်အတောက်တွေ) က ကျွန်တော်လုပ်ခဲ့တဲ့ NLP (Natural Language Processing) သုတေသနအလုပ်တွေဖြစ်တဲ့ Machine Translation, Automatic Speech Recognition, Text to Speech, POS tagging စတဲ့ အလုပ်တွေအတွက် လုပ်ရကိုင်ရ အဆင်အပြေဆုံး ပုံစံအတိုင်း တကယ့်ကို simple unit အနေနဲ့ဖြတ်ထားတာ ဖြစ်ပါတယ်။ အဲဒါကြောင့် "ဘီစီ ၂၈၇" ကို "ဘီ စီ ၂ ၈ ၇"၊ "နက္ခတ္တဗေဒ" ကို "နက္ခတ္တ ဗေ ဒ"၊ နောက်ပြီးတော့ မြန်မာစာတွေနဲ့ အတူတူရောပါနေတဲ့ "Sand Reckoners" ကို "S a n d R e c k o n e r s" ဆိုပြီး ဖြတ်ထားပါတယ်။ ဆိုလိုတာက ပါဠိဆင့်တွေကို ဖြေတဲ့ အလုပ် (ဥပမာ နက္ခတ္တကို နက် ခတ် တ)၊ ဂဏန်းစာလုံးတွေကို တွဲတဲ့အလုပ် (ဥပမာ ၂၈၇)၊ ရောပါနေတဲ့ အင်္ဂလိပ်စာလုံးတွေကို ဖယ်ပစ်တာ၊ နဂိုအတိုင်းပဲ တွဲထားတာ (ဥပမာ Sand Reckoners) မျိုးတွေကို တမင်တကာ လုပ်မထားတာပါ။ အကြောင်းအရင်းကတော့ အမျိုးမျိုးရှိပါတယ်။ ဥပမာ ပါဠိဆင့်တွေကို ဖြေပစ်လိုက်ရင် လိုအပ်တဲ့အခါမှာ နဂိုပုံစံအတိုင်းပြန်ရအောင် ပြန်ဆင့်ရပါတယ်။ အဲဒီအလုပ်က လွယ်မလိုလိုနဲ့ လက်တွေ့မှာတော့ ပါဠိဆင့် ပြန်ဆင့်ပေးရတဲ့အလုပ်အတွက် processing time နဲ့ အဲဒီကနေထွက်လာမဲ့ error တွေကို ရှာဖွေရတဲ့အလုပ်၊ ပြန်ပြင်ပေးရတဲ့အလုပ်တွေကို ရှောင်ချင်လို့ ဖြစ်ပါတယ်။ [syllbreak \(https://github.com/ye-kyaw-thu/syllbreak\)](https://github.com/ye-kyaw-thu/syllbreak) က ကျွန်တော်တို့ မြန်မာစာကို unicode နဲ့သာ မှန်မှန်ကန်ကန်ရေးထားရင်၊ "Regular Expression တစ်ကြောင်းထဲနဲ့ လွယ်လွယ်ကူကူ ဖြတ်လိုရကြောင်း" နောက်ပြီးတော့ အဲဒါက "တကယ်လည်း မြန်မာစာ NLP သုတေသနအလုပ်တွေအတွက် အသုံးဝင်ကြောင်း"၊ ဒါ့အပြင် "အခြေခံကြတဲ့ မြန်မာဝဏ္ဏ (syllable unit) အနေနဲ့လည်း ရပါလိမ့်မယ်" ဆိုတဲ့ message ကိုပေးထားတာပဲ ဖြစ်ပါတယ်။ ကျွန်တော်ရဲ့ syllable unit တွေက မြန်မာစာအနေနဲ့ကြည့်ရင် ပြည့်စုံမှန်ကန်တယ်လို့ မဆိုလိုပါဘူး။ ကိုယ်လုပ်မဲ့ အလုပ်၊ develop လုပ်နေတဲ့ application ပေါ်ကို မူတည်ပြီးတော့ လက်ရှိ ကျွန်တော်ပြင်ဆင်ပေးထားတဲ့ Regular Expression ကို ကြိုက်သလို ဖြည့်စွက်တာ၊ ပြင်သုံးတာကိုလုပ်နိုင်ပါတယ်။
- Python 3.4 ကနေစပြီးတော့ "ur" (Unicode + Raw text) ဆိုပြီးတွဲရေးတာကို support မလုပ်ပါဘူး။ သို့သော် "u" တစ်လုံးတည်း "r" တစ်လုံးတည်း သုံးတာကိုတော့ ခွင့်ပြုပါတယ်။

<https://stackoverflow.com/questions/26063899/python-version-3-4-does-not-support-a-ur-prefix>
(<https://stackoverflow.com/questions/26063899/python-version-3-4-does-not-support-a-ur-prefix>)

- Jupyter Notebook နဲ့ ပတ်သက်တဲ့ installation လုပ်ပုံလုပ်နည်း၊ အသုံးပြုပုံနဲ့ ပတ်သက်ပြီး မြန်မာလိုလေ့လာချင်တဲ့ သူများအတွက် ကျွန်တော့်ရဲ့ [Tutorial \(https://github.com/ye-kyaw-thu/Tutorials\)](https://github.com/ye-kyaw-thu/Tutorials) မှာ လေ့လာနိုင်ပါတယ်။