

Manual for PANGKAT: A Dedicated Tokenizer for the Tagalog Language

PANGKAT was designed to be implemented as a Python module, which can be imported into more advanced Tagalog NLP applications as its tokenizer.

```
module > main.py > ...
1  import pangkat
2
3  # Initialize PANGKAT, as it is implemented as a Python Class
4  PANGKAT = pangkat.PANGKAT()
5
6  # To use PANGKAT, call its labelTokens function
7  # Its parameter is the file containing the data to be tokenized
8  # The labelTokens function returns the arrays of the resulting tokens and their labels
9  # Results are stored in independent arrays for both short and longer unit tokenization
10 tokenList, labelList, longerTokenList, longerLabelList = PANGKAT.labelTokens("input.txt")
```

For instance, the results of PANGKAT are as follows for tokenizing the following sentence:

“MANILA, Philippines — Pinababalik ni Hon. Antonio Kho, Jr. and P60 billion na budget ng Philippine Health Insurance Corporation (PhilHealth) matapos makwestiyon ang malaki-laking budget nito na 'di umano'y hindi napapakinabangan ng publiko.”

```
PS C:\Users\Justin\Desktop\PANGKAT\module> python main.py
PANGKAT is loaded

Short Unit Tokenization Results:

[['MANILA', ',', 'Philippines', '-', 'Pinababalik', 'ni', 'Hon', '.', 'Antonio', 'Kho', ',', 'Jr', '.', 'and', 'P60',
 'billion', 'na', 'budget', 'ng', 'Philippine', 'Health', 'Insurance', 'Corporation', '(', 'PhilHealth', ')', 'mata',
 'pos', 'makwestiyon', 'ang', 'malaki', '-', 'laking', 'budget', 'nito', 'na', '"', 'di', 'umano', '"', 'y', 'hindi',
 'napapakinabangan', 'ng', 'publiko', '.']]
[['B-LOC', 'I', 'I', 'O', 'O', 'O', 'B-PER', 'I', 'I', 'I', 'I', 'I', 'I', 'O', 'B-MWE', 'I', 'O', 'O', 'O', 'B-ORG',
 'I', 'I', 'I', 'I', 'I', 'I', 'O', 'O', 'O', 'B-MWE', 'I', 'I', 'O', 'O', 'O', 'B-MWE', 'I', 'I', 'I', 'I', 'O', 'O', 'O', 'O', 'O']]

Longer Unit Tokenization Results:

[['MANILA,Philippines', '-', 'Pinababalik','ni', 'Hon.Antonio Kho,Jr.', 'and', 'P60 billion', 'na', 'budget', 'ng',
 'Philippine Health Insurance Corporation (PhilHealth)', 'matapos', 'makwestiyon', 'ang', 'malaki-laking', 'budget',
 'nito', 'na', "'di umano'y", 'hindi', 'napapakinabangan', 'ng', 'publiko', '.']]
[['NE-LOC', 'W', 'W', 'W', 'NE-PER', 'W', 'MWE', 'W', 'W', 'W', 'NE-ORG', 'W', 'W', 'W', 'MWE', 'W', 'W', 'W', 'MWE',
 'W', 'W', 'W', 'W', 'W']]
```

For better visualization, the following images present the results using PANGKAT's graphical user interface for both short and longer unit tokenization, respectively:

PANG-KAT: A Dedicated Tokenizer for the Tagalog Language

SHORT UNIT	LONGER UNIT
TOKENS	LABELS
MANILA	B-LOC
,	I
Philippines	I
—	O
Pinababalik	O
ni	O
Hon	B-PER
,	I
Antonio	I
Kho	I
,	I
Jr	I
,	I
ang	O
P60	B-MWE
billion	I
na	O
budget	O
ng	O
Philippine	B-ORG
Health	I
Insurance	I
Corporation	I
(I
PhilHealth	I

Download CSV

Download JSON

PANG-KAT: A Dedicated Tokenizer for the Tagalog Language

SHORT UNIT	LONGER UNIT
TOKENS	LABELS
MANILA, Philippines	NE-LOC
—	W
Pinababalik	W
ni	W
Hon. Antonio Kho, Jr.	NE-PER
ang	W
P60 billion	MWE
na	W
budget	W
ng	W
Philippine Health Insurance Corporation (PhilHealth)	NE-ORG
matapos	W
makwestiyon	W
ang	W
malaki-laking	MWE
budget	W
nito	W
na	W
'di umano'y	MWE
hindi	W
napapakinabangan	W
ng	W
publiko	W
.	W

Download CSV

Download JSON