

Named Entity Recognition on Indonesian Microblog Messages



Natanael Taufik, Alfian F. Wicaksono, Mirna Adriani

Information Retrieval Lab.

Faculty of Computer Science, Universitas Indonesia

Outline

- Named Entity Recognition for Entity Names
- Background
- Dataset
- Feature Representation
- Features
- Result
- Conclusion

Named Entity Recognition for Entity Names

Task to identify named entity from text.

Scope:

- Person
- Organization
- Location

Named Entity Recognition for Entity Names

BBC Sign in News Sport Weather Shop Earth

NEWS

Home Video World Asia UK Business Tech Science Magazine

US Election 2016 Results States A-Z



US ELECTION 2016

Trump presidency: Your questions answered

14 November 2016 | US Election 2016

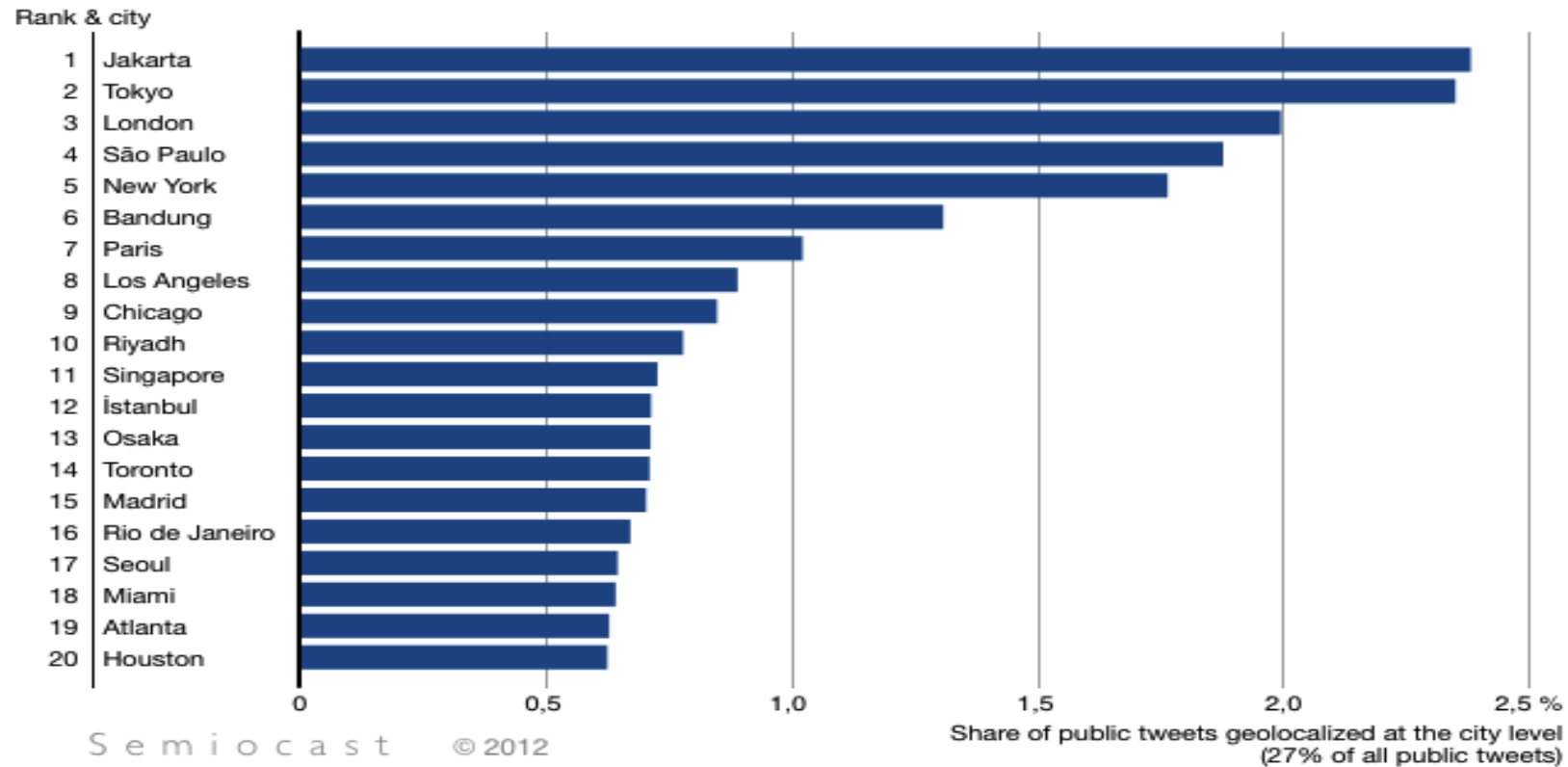


Background



Background

Top 20 cities by number of posted tweets
(among 10.6B public tweets posted in June 2012)



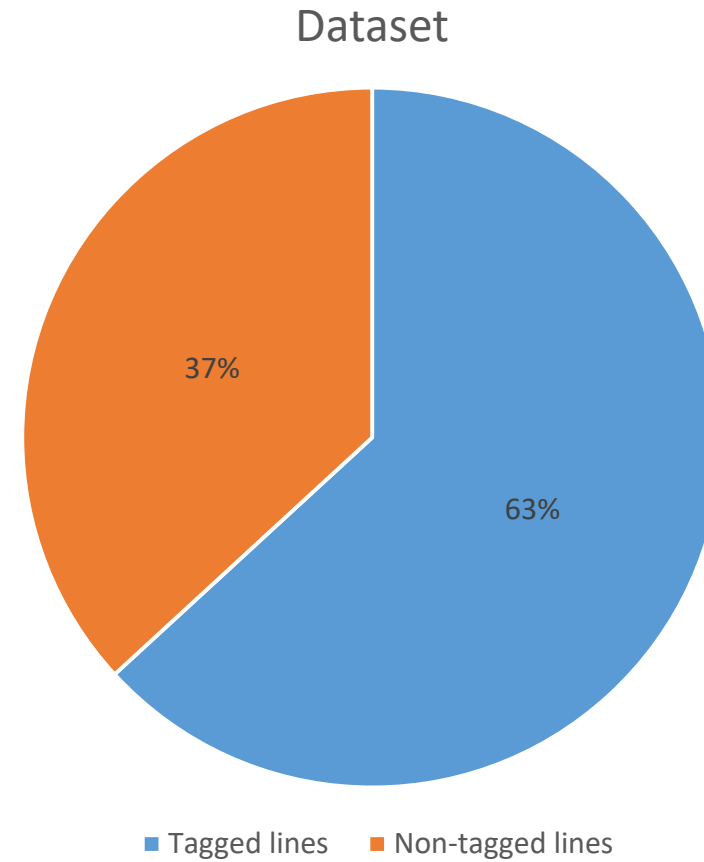
NER for Indonesian Microblog

404 - Not Found

Dataset

	# count
Tagged lines	379
Non-Tagged lines	221
Total	600

Type	# count
Person	225
Location	257
Organization	204



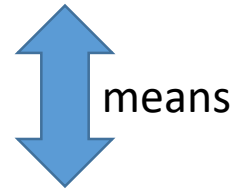
Feature Representation



Feature 1	Feature 2	Feature 3	Label	
Feature 1	Feature 4	Label		
Feature 1	Feature 5	Label		
Feature 1	Label			
Empty line				
Feature 1	Feature 3	Label		
Feature 1	Feature 2	Label		
...		

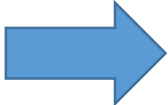
Feature Representation(example)

Jokowi Diminta Perhatikan Angkutan Umum



Jokowi Asked To Pay Attention To Public Transportation

Feature Representation(example)

Source		Features	Label
Jokowi		capitalized	Person
Asked		capitalized	Other
To		capitalized	Other
Pay		capitalized	Other
Attention		capitalized	Other
To		capitalized	Other
Public		capitalized	Other
Transportation		capitalized	Other

Feature Representation(example)

<B_ENAMEX TYPE="PERSON">Jokowi<E_ENAMEX> Asked To Pay Attention
To Public Transportation

Features

- Word

Jokowi, written as is: **“jokowi”**

- Last 3 letters

Example of family names: *Setiwan, Hendrawan, Himwan*

- Word length

Written as follow: **“wordLength:6”**

Features

- Pattern function_[1]:

[A-Z] => "A"

[a-z] => "a"

[0-9] => "0"

other character => "-"

Jokowi => Aaaaaaa

=> Aa

Result: ["Aaaaaaa", "Aa"]

Features

- Inside bracket

Example:

You're like a moon, can be seen, but can't be owned(Ari, Dewi)

Feature: **“insideBracket”**

Features

- Part-of-speech

Train Stanford Log-linear Part-Of-Speech Tagger using tagged tweets made by Canggihabrata and Bressan

Example:

“Susilo Bambang Yudhoyono”

$f(\text{“Yudhoyono”}) = [\text{“1stLeftPOS-NNP”}, \text{“2ndLeftPOS-NNP”}]$

- Surrounding words

$f(\text{“Yudhoyono”}) = [\text{“1stLeftWord-Bambang”}, \text{“2ndLeftWord-Susilo”}]$

Features

- Lookup list – common location and stopwords^[1]

$f(\text{“Jawa”}) = [\text{“isRegion”}]$

$f(\text{“dan”}) = [\text{“isStopword”}]$

- Non-standard word list^[1]

$f(\text{“nggak”}) = [\text{“tidak”}]$

Result

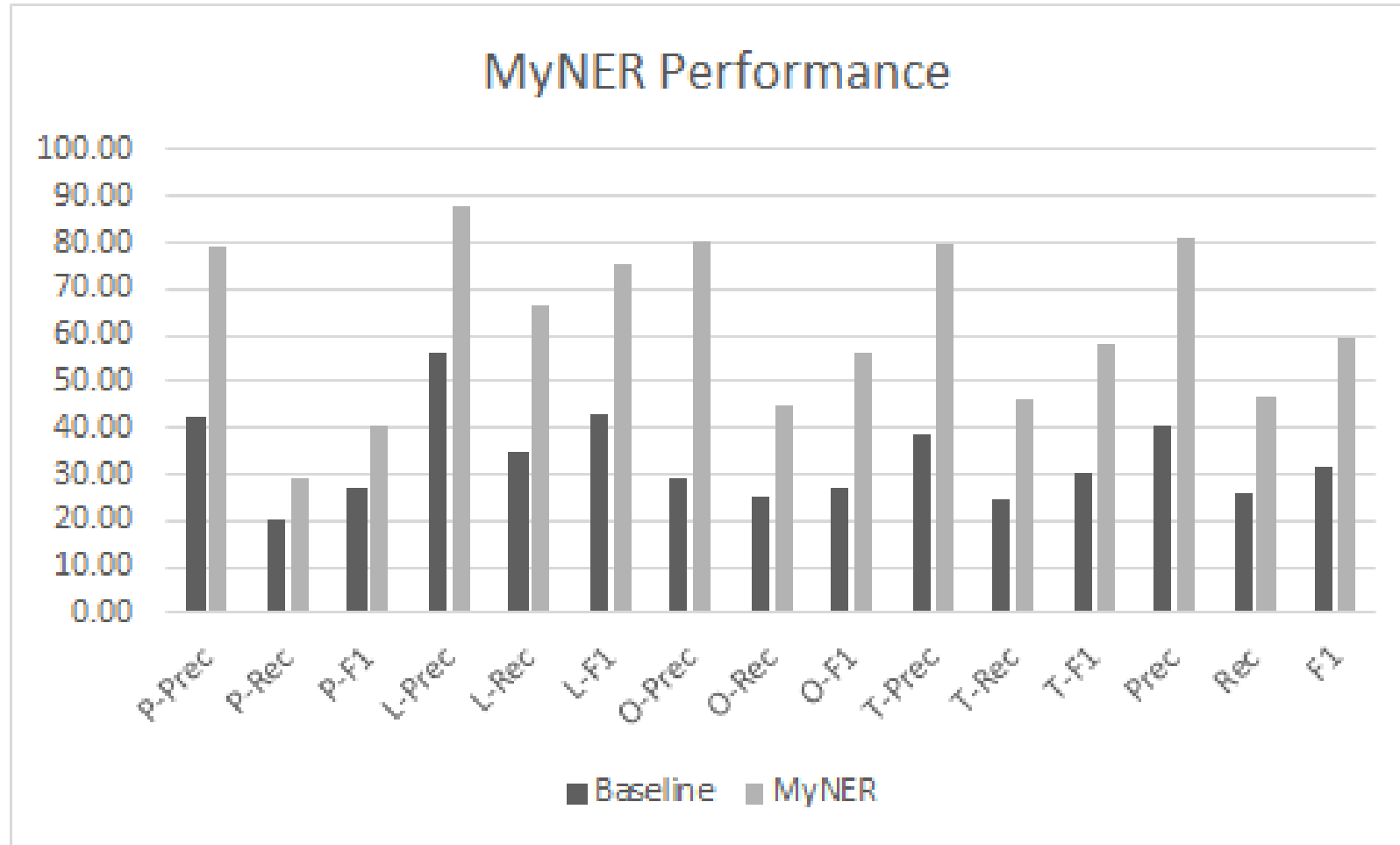
Baseline

Research by G. Wahyudi, using rule based approach based on contextual and morphological information as well as part-of-speech of the words.

Result

	Type	Precision(%)	Recall(%)	F1(%)
Baseline	PER	42.45	20.00	27.19
	LOC	56.33	34.63	42.89
	ORG	29.38	25.49	27.30
MyNER (10-fold cross validation)				
MyNER (10-fold cross validation)	PER	79.02	29.16	40.75
	LOC	88.04	66.16	75.13
	ORG	80.27	44.63	56.35

Result



Conclusion

- NER on Indonesian microblog is challenging because it's sort and written in a non-standard way
- Specific model built for Indonesian microblog messages outperforms model built for formal Indonesian language
- Now we have a NER for Indonesian microblog messages
- Further improvements are needed since the best results do not seem sufficient for higher level applications
- Our model is based on sequence labeling task that employs Conditional Random Fields as our machine learning algorithm

Terima kasih 😊

Thank you 😊

謝謝 😊

ありがとうございました 😊

Any question?