# Developing (and utilizing) an Indonesian Treebank

Arawinda Dinakaramani, Fam Rashel, Andry Luthfi, Bayu Distiawan, and Ruli Manurung

Faculty of Computer Science, Universitas Indonesia

The Second Wordnet Bahasa Workshop

Nanyang Technological University, 15-16 January 2016

1

# Outline

- **Background**
- Annotation process
- Outputs
- Making use of the treebank

# At the previous workshop...



- 10k Indonesian sentences from the PAN Localization parallel corpus (http://www.panl10n.net/indonesia)
- 23 POS tagset
- +/- 250k tokens (incl. MWE from http://kateglo.com)
- Rule-based tagger (utilizes MorphInd: http://septinalarasati.com/work/morphind)
- Released under Creative Commons BY-NC-SA 4.0
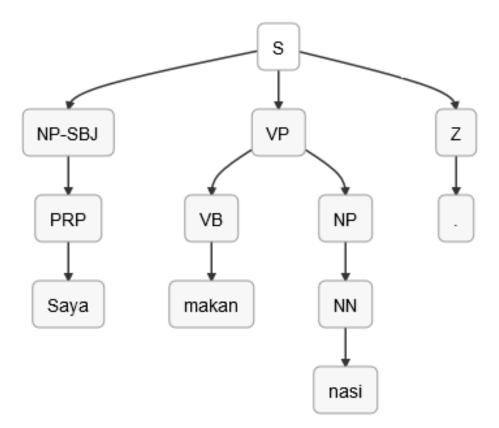


http://bahasa.cs.ui.ac.id/postag
https://github.com/famrashel/idn-tagged-corpus
https://github.com/andryluthfi/indonesian-postag

# Next goal: building a treebank

- A **treebank** is a corpus of sentences complete with annotated syntactic structure.

- Useful as training data for statistical parsers.

- Example:

# Bracketing Guidelines

- Our goal: treebank of the first 1000 sentences of the POS tagged corpus.

- Use POS tags as a starting point.

- Adopt Penn Treebank bracketing guidelines (Bies et al., 1995) where possible.

- Consult authoritative Indonesian grammar references (Alwi et al., 2003; Sneddon et al., 2010).

# Outline

- Background
- **Annotation process**
- Outputs
- Making use of the treebank

# Data preparation

- Convert from POS tagged corpus format to initial bracketing (forest of singleton POS tag trees).
- Example:

*Pembahasan      tadi      masih dalam tahap awal.*
  discussion   previous    still    in    stage  early

```
Pembahasan      NN
tadi            PR
masih           MD
dalam           IN
tahap           NN
awal            NN
.               Z
```

into bracketed file format:

```
(NN (Pembahasan))(PR (tadi))(MD (masih))(IN
(dalam))(NN  (tahap))(NN (awal))(Z (.))
```

# Annotation Process

- 3 annotators parsed the first 100 sentences of our corpus.
- In conjunction with development of bracketing guidelines.
- Sample:

```
(S (PP-TMP Selama
      (NP bertahun-tahun))
   (NP-SBJ monyet)
   (VP mengganggu
      (NP warga Delhi))
   .)
```

- Keep track of all arising issues, resolve among annotators.
  – Consistent phrase structure bracketing
  – Sentence alignment (split & merge)
  – Incorrect POS tags
  – MWE

# Notes of issues

# **Annotation Process: Multi-phase**

- Re-annotate the first 100 sentences.

- Annotate the next 100 sentences.

- Annotate remainder of the 1000 sentences.
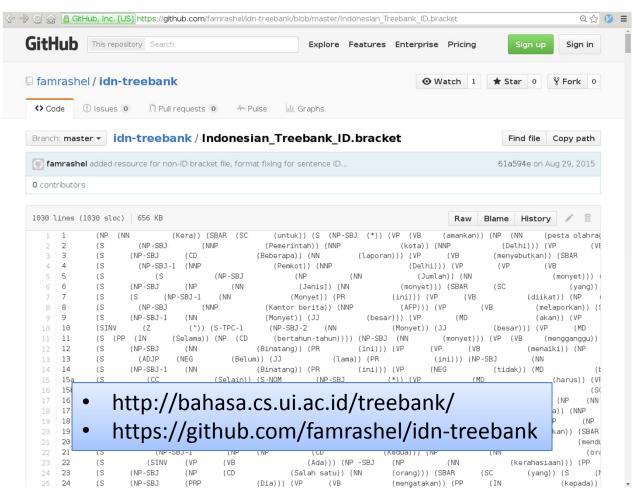
# Outline

- Background

- Annotation process

- **Outputs**

- Making use of the treebank

# The Treebank



- 1000 sentences
- 2 variants: with/without sentence IDs – for mapping to POS tagged corpus
- Creative Commons BY-NC-SA 4.0

- http://bahasa.cs.ui.ac.id/treebank/
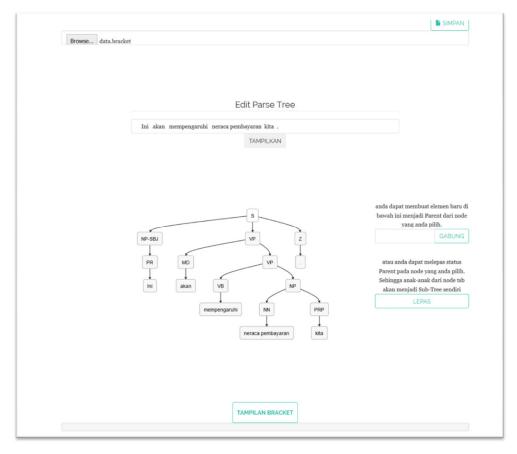- https://github.com/famrashel/idn-treebank

# Indonesian Treebank Bracketing Guidelines

- Guidelines to annotate Indonesian sentence structure in developing Indonesian Treebank.
- Rules for bracketing clauses and sentences:
  - Simple active/passive declarative, imperative, interrogative, inversion, subordinative, coordination, direct/indirect quote, etc.
- Rules for bracketing phrasal structures:
  - Phrasal structures: Adjectival phrases (ADJP), Adverbial phrases (ADVP), Conjunctor phrases (CONJP), Noun phrases (NP), Numeral phrases (QP), Prepositional phrase (PP), Verb phrase (VP), Unlike coordinated phrase (UCP)
- Syntactic category labels and function tags from the Penn Treebank bracketing guidelines.
- POS tags from our Indonesian POS tagset.

# Web-Based Annotation Tool



JavaScript only, runs locally, single user
**https://github.com/andryluthfi/annotation-tools-lightweight**

Client-server using database, multiple concurrent user, agreement checking
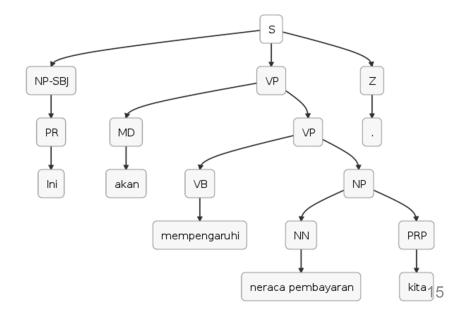**https://github.com/andryluthfi/annotation-tools**

# Web-Based Annotation Tool

- Direct input by user, or load from `.bracket` file
- Resulting annotation saved to `.bracket` file.
- Example:

  *Ini akan mempengaruhi neraca pembayaran kita.*

  this will        impact      balance      payment      us

```
(S (NP-SBJ (PR (Ini)))
(VP (MD (akan)) (VP (VB
(mempengaruhi)) (NP (NN
(neraca pembayaran))(PRP
(kita))))) (Z (.)))
```

# Outline

- Background

- Annotation process

- Outputs: treebank, guidelines, tools
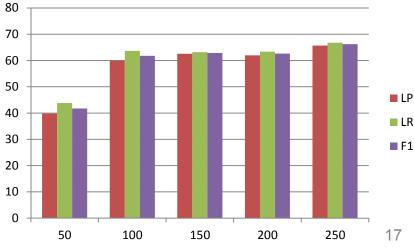
- **Making use of the treebank**

# Teaching tool

- 300 sentence treebank used for undergraduate NLP class assignment

- Each student asked to annotate 10+5 sentences ☺

- Experiment on training Stanford Parser with varying parameters

# *Text Mining Systemic Risk Prioritization* (TM-SRP)

- Detect economic risks stated in financial news articles.

- Domain experts from macroprudential policy dept. of Indonesian central bank constructed model of 31 economic risks and related keywords.

- Baseline approach: matching of keyword occurrence in a single sentence.
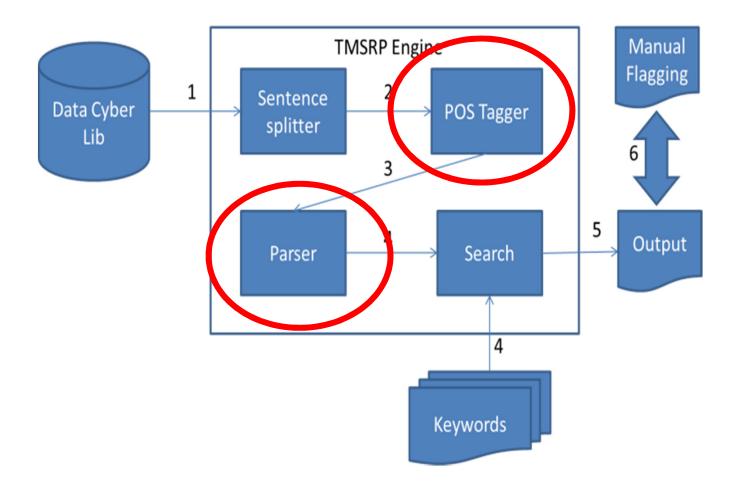
# Problem with Keyword Matching

- Example risk: Global Interest Rate
  - Keyword 1: suku bunga (interest rate)
  - Keyword 2: naik (increasing)

- *Setelah    naik    menjadi  presiden,  Jokowi*
   after    ascend  become president, Jokowi

   *memerintahkan untuk menurunkan suku_bunga   BI*
   instruct            to      lower     interest rate   BI

**Idea:** Utilize syntactic structure from probabilistic parser. Only match keywords in corresponding syntactic relations.

# Proposed Approach

# POS Tagger Domain Adaptation

- Lots of domain-specific terms not found in the training data.
  - "nilai tukar" (exchange rate)
  - "daya beli" (purchasing power)
  - etc.

# Pattern matching

- Focus on each subtree that has root label "S". If a sentence has several clauses, the search will focus on each clause.

- Differentiate 2 types of keywords:

  - Simple Node: Keyword can appear anywhere in a phrase. Mostly for "noun" keywords

  - Head Node: Keyword must appear at the beginning of a phrase. Mostly for "verb" keywords.

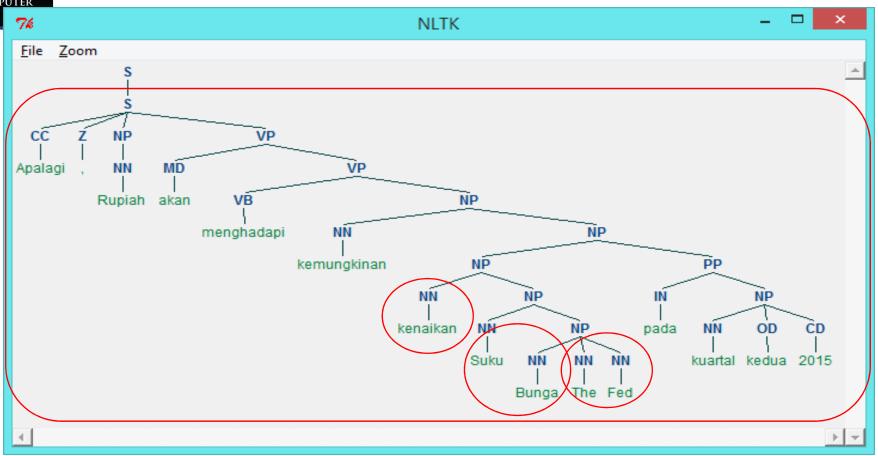- Find a negation label on each sub-tree "S".

# Search Engine

```
"1":{
    "idRisk":"1",
    "jenis":"risiko",
    "kategori":"sumber0",
    "components":[
        {
            "keywords":"Federal Open Market Committee/NP,The Fed/NP",
            "find":"node"
        },
        {
            "keywords":"Suku Bunga/NP,Interest Rate/NP,Fed Fund Rate/NP",
            "find":"node"
        },
        {
            "keywords":"Kenaikan/NP/NN,Naik/VP/VB,Increase/VP/VB,Penaikan/NP/NN",
            "find":"head"
        }
    ]
},
```

# Search Engine



keyword1: **The Fed** ; keyword2: **Suku Bunga**;  keyword3: **Kenaikan**

# Evaluation

- Evaluation judgments provided by domain experts → manually labelled sentences for risk

- Precision: 77.15%

- Recall: 91.76%

# References

- A. Bies, M. Ferguson, K. Katz, and R. MacIntyre. 1995. "Bracketing Guidelines for Treebank II Style Penn Treebank Project". https://catalog.ldc.upenn.edu/docs/LDC99T42/prsguid1.pdf. Last Access: September 2013.

- A. Dinakaramani, F. Rashel, A. Luthfi, and R. Manurung. 2014. "Designing an Indonesian Part of Speech Tagset and Manually Tagged Indonesian Corpus". In Proceedings of the 2014 International Conference on Asian Language Processing.

- H. Alwi, S. Dardjowidjojo, H. Lapoliwa, and A. Moeliono. 2003. Tata Bahasa Baku Bahasa Indonesia. Third Edition. Balai Pustaka, Jakarta.

- J. Sneddon, A. Adelaar, D. Djenar, and M. Ewing. 2010. Indonesian Reference Grammar. Second Edition. Allen & Unwin, Crows Nest.

- M. Marcus, B. Santorini, and M.A. Marcinkiewicz. 1993. Building a large annotated corpus of English: the Penn Treebank. Computational Linguistics, Vol. 19, No. 2, pp. 313-330.