# Gotong Royong in NLP Research
# A Mobile Tool for Collaborative Text Annotation in Indonesia

Lisa Madlberger          Ade Romadhony          Ayu Purwarianti

# Introduction

- The Indonesian language Bahasa Indonesia counts as a "low-resource" language

- Machine learning technology advanced the development of NLP tools in Indonesia

  **BUT:**

  Machine learning based NLP methods
  depend on the availability of annotated training data

# Annotated Training Data

Example – Named Entity Recognition

**MAY DAY: Buruh KSPI Ancam Mogok Kerja Jika Tuntuan Tak Digubris.**

**Time**    **Per**    **Org**

*May Day: KSPI workers threaten to strike if their demands are ignored*

# The Problem

| Token | ACTOR | TRIGGER | TARGET | LOCATION | TIME |
|---|---|---|---|---|---|
| MAY | | | | | x |
| DAY | | | | | x |
| : | | | | | |
| Buruh | x | | | | |
| KSPI | x | | | | |
| Ancam | | x | | | |
| Mogok | | | | | |
| Kerja | | | | | |
| Jika | | | | | |
| Tuntuan | | | | | |
| Tak | | | | | |
| Digubris | | | | | |
| Angkot | | | | | |
| di | | | | | |
| bogor | | | | | |
| pada | | | | | |
| mogok | | | | | |
| kerja | | | | | |
| , | | | | | |
| jalanan | | | | | |
| tuh | | | | | |

Manual annotation of data is
- tedious and
- time-consuming

# Solution Approach



Gotong Royong



Mobile First Culture

**The first Mobile Collaborative Annotation Tool**
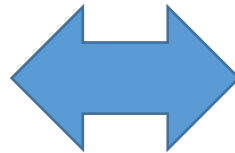
# Gotong Royong

# Mobile First Culture



In today's Indonesia, 93% of online users access the Internet via their smartphone (Andrews et al., 2015)

# Solution

**katakita.click**

**The First Mobile Collaborative Annotation Tool**

**Colleagues**

**Students**

**Family**

**Friends**

# Existing Systems

Shortcomings for their application in Indonesia

1) Not Mobile Friendly

# Existing Systems

BRAT



BRAT, P. Stenetorp, 2012

# Existing Systems

GATE

# Existing Systems

Shortcomings for their application in Indonesia

1) Not Mobile Friendly

2) Interface does not support Bahasa Indonesia

# The Solution

- We propose a tool

- that makes data annotation more efficient

- allows data to be annotated
  by several users at the same time

- and can be used anywhere, anytime
  – using a mobile phone

Klik tombol label dan kemudian klik token kata

#NewsGibol | Otamendi | Siapkan | Aksi | Mogok | | | http://t.co/z1zFrzvU1G

Orang/Organisasi | Lokasi | Waktu | Lainnya

Tidak relevan | Hapus | >

# Example Binary Classification



Is this Tweet related to **labour strikes or protests?**

Dear mahasiswa yg lg 'aksi' demo, bagus sih merjuangin hak rakyat. TAPI YA GAK NUTUPIN JALAN JUGA. HUH. Pengguna jalan raya juga RAKYAT btw.

Ya    Tidak    Lanjut
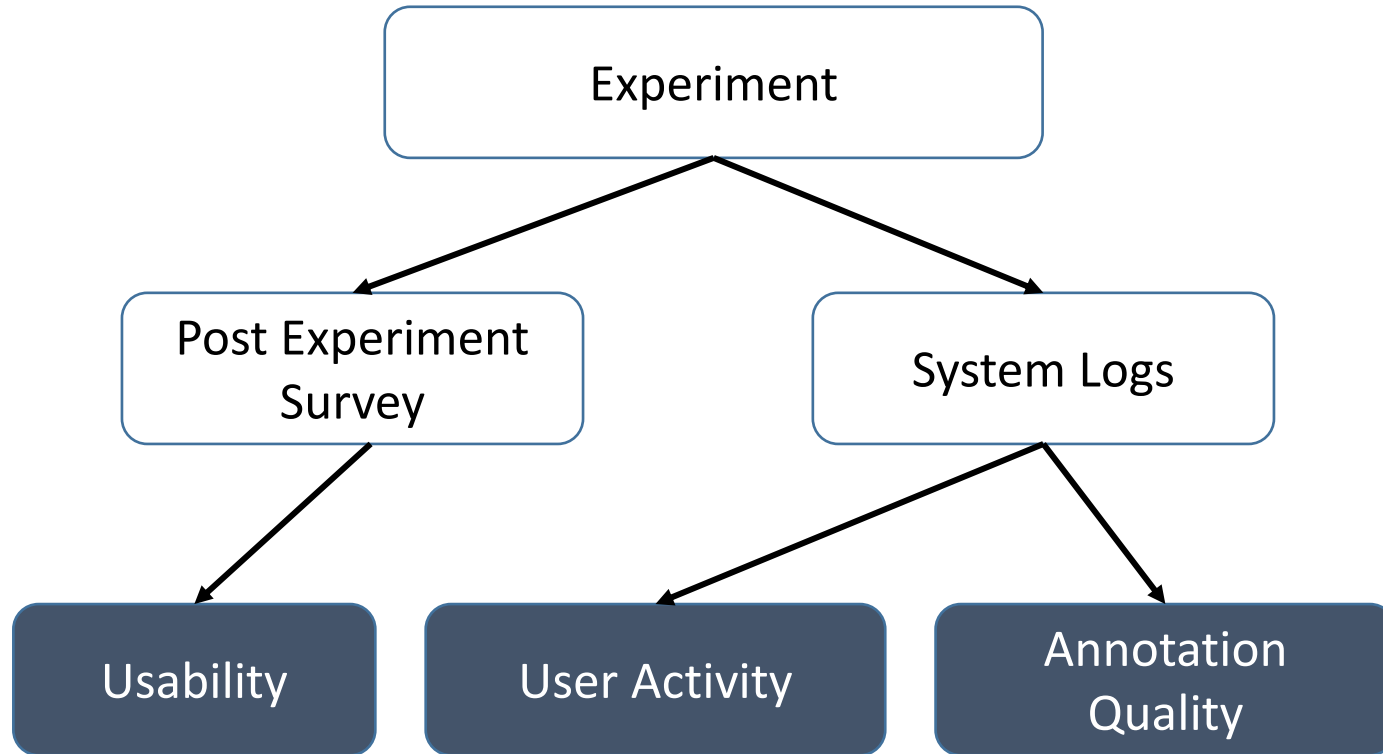
# Experimental Evaluation

- 15 Indonesian students/alumni from 5 Universities

- Labelled 100 Tweets each

- Using one of three NLP data annotation tasks:

  - Named Entity Recognition
  - Semantic Role Labeling
  - Binary Classification

- In one week, using KataKita on their mobile phones

# Evaluation Criteria

# Usability

Strongly disagree                          Strongly agree

I could use KataKita from mobile phone so I can annotate anytime and anywhere.

When I use KataKita, I need to wait couple of minutes until all the tokens were loaded on the screen.

KataKita annotation guideline is easy to understand.

I think KataKita is too complicated.

I think KataKIta is easy to use.

I think I need technical support to use KataKita

I imagine that most of KataKita users could learn to use KataKita quickly
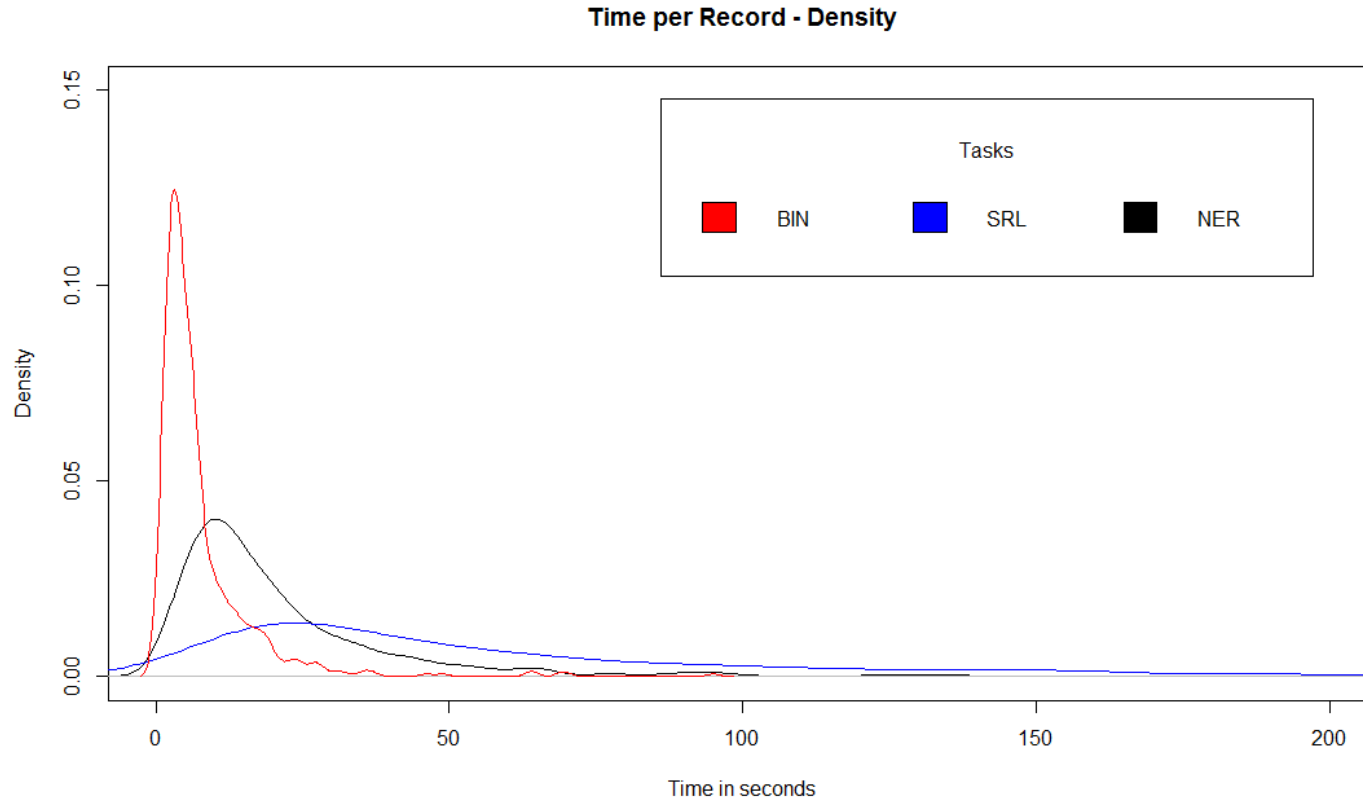
I think KataKita is impractical to use

I feel very confident when doing the annotation using KataKita.

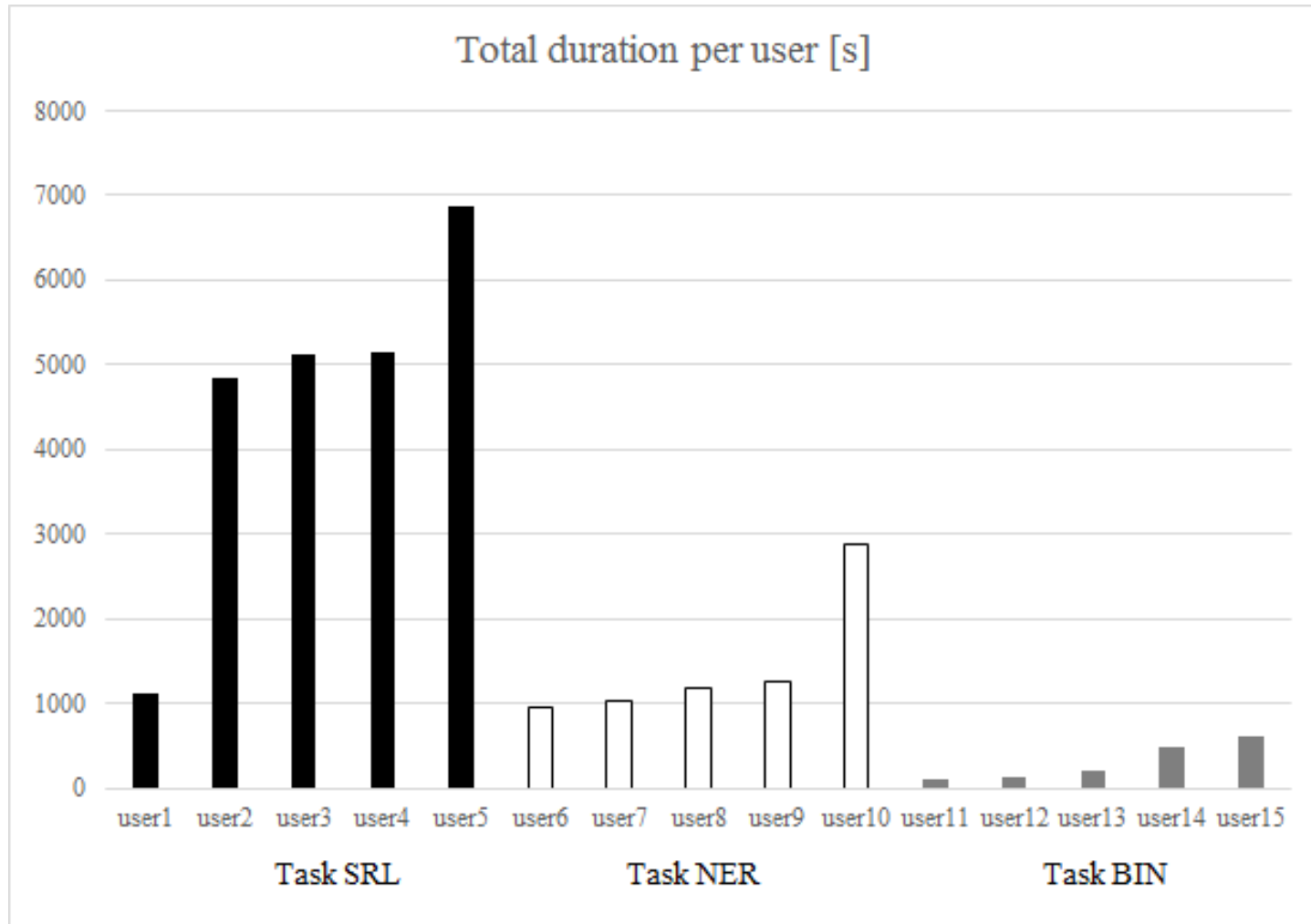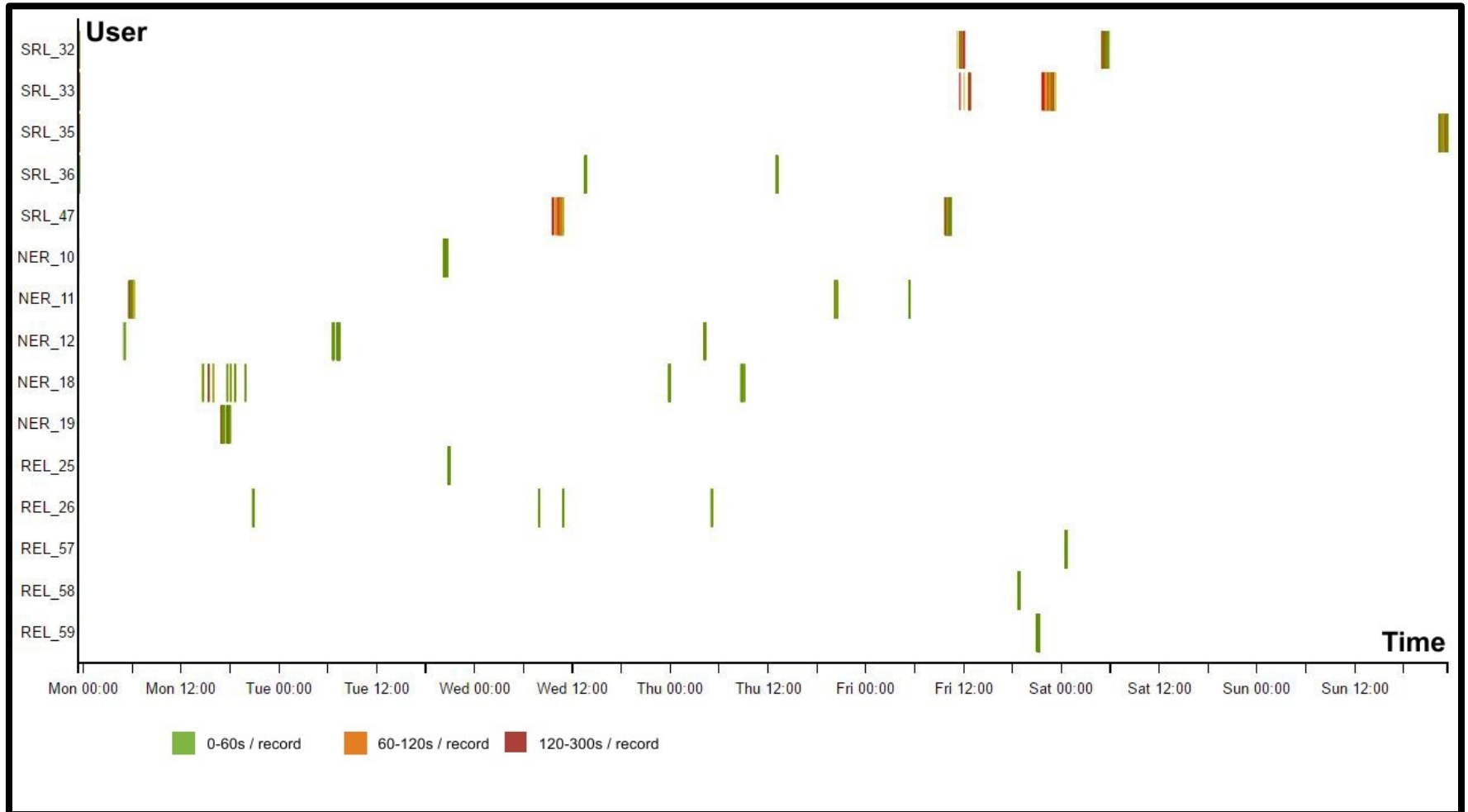I must learn a lot of things before using KataKita.

Time per Record - Density

**Median Time**

| Binary Classification | 5s |
|---|---|
| Named Entity Recognition | 17s |
| Semantic Role Labelling | 41s |

# Experimental Evaluation

# Annotation Quality

| Task | Fleiss' Kappa | Interpretation |
|---|---|---|
| Binary Classification | 0.45 | Moderate Agreement |
| Named Entity Recognition | 0.22 | Low Agreement |
| Semantic Role Labelling | 0.41 | Moderate Agreement |

*0 = no agreement ,1 = perfect agreement*

- How to improve annotation quality? What are the factors and user attributes influencing the quality?

- How to present guidelines and provide training on the phone?

**Questions?**
Please contact us!
ade_romadhony@students.itb.ac.id
lisa.madlberger@tuwien.ac.at

https://github.com/strikesensor/

Lisa Madlberger          Ade Romadhony          Ayu Purwarianti