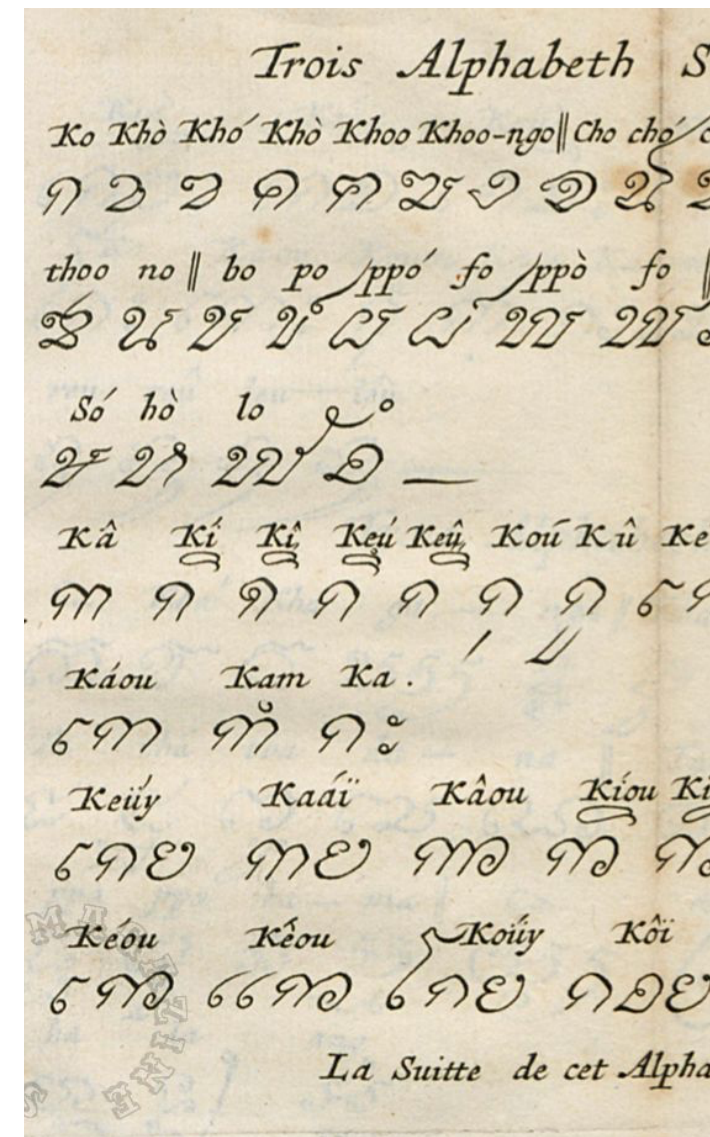# A Brief Specification of LST20 Corpus

**Prachya Boonkwan, Vorapon Luantangsrisuk,
Sitthaa Phaholphinyo, Kanyanat Kriengket,
Dhanon Leenoi, Charun Phrombut, Monthika Boriboon,
Krit Kosawat, and Thepchai Supnithi**

Language and Semantic Technology Lab (LST)
National Electronics and Computer Technology Center (NECTEC)
Bangkok, Thailand

# Outline

- Executive summary

- Agreement of usage & download

- Annotation schemes

- Statistics

# Executive Summary

# Executive Summary

- ## LST20 Corpus
  - ### Dataset for training fundamental Thai language processing tasks
  - ### Featured linguistic information
    - **Word boundaries** for word segmentation
    - **Named entities** for named entity recognition
    - **Clause boundaries** for clause segmentation
    - **Sentence boundaries** for sentence segmentation
    - **News genres** for document classification
  - ### CoNLL-2003 format: tab-separated columns

| | |
|---|---|
| Words | 3,163,034 |
| Named entities | 288,020 |
| Clauses | 248,181 |
| Sentences | 74,180 |
| Distinct words | 46,692 |
| Genres | 15 |
| News articles | 3,745 |

Available at
https://aiforthai.in.th

# Agreement of Usage & Download
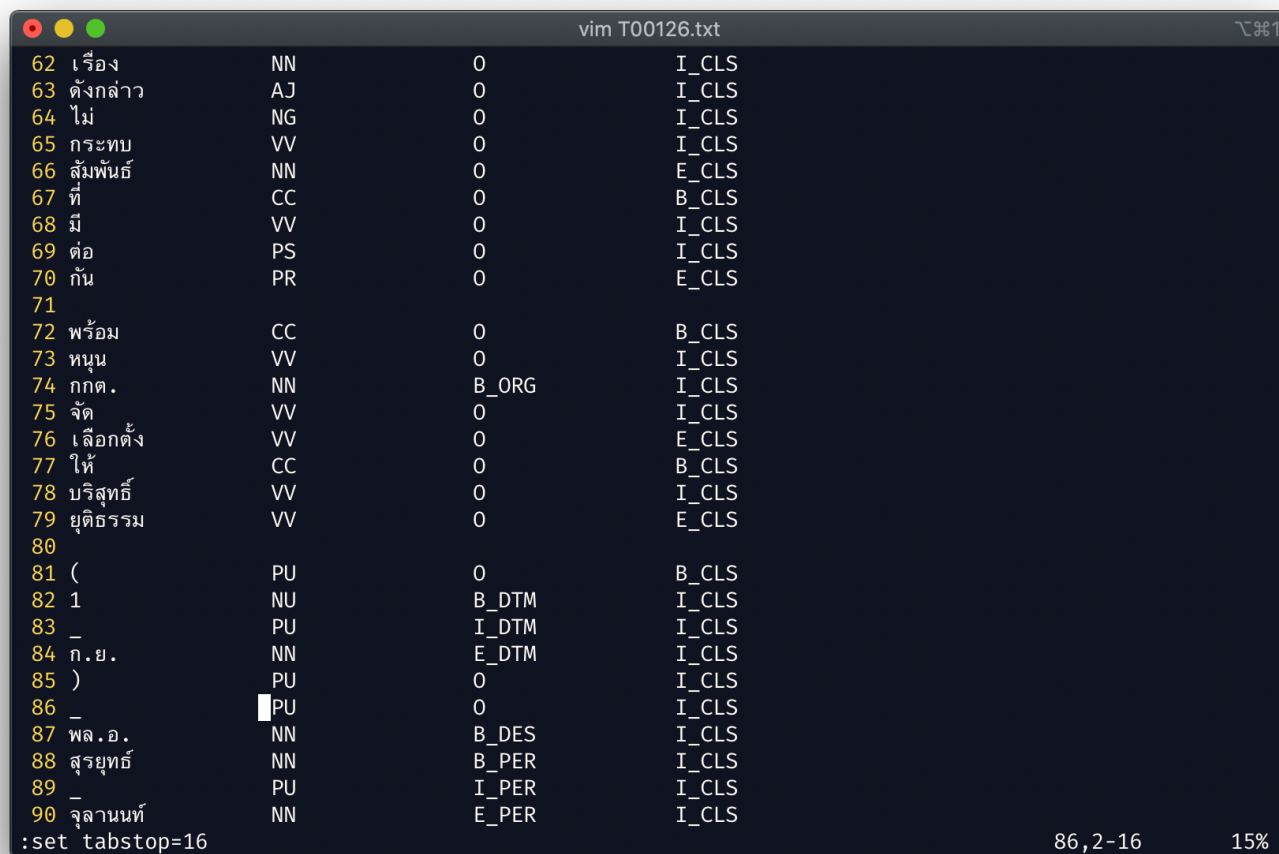
# Agreement of Usage

- Data consortium
  - Non-commercial use, research, open source
    - Free of charge
    - Citing our technical paper is **always** required
  - Commercial use
    - **Option 1 (in-kind):** Contributing a dataset of 150K-300K words completely annotated with our annotation schemes within 1 year
    - **Option 2 (in-cash):** Purchasing the dataset

- **Download**
  - https://aiforthai.in.th
  - Registration is required before downloading
  - Much more web APIs and demos are also available therein

# Annotation Schemes

# Format: CoNLL-2003 Style

```
vim T00126.txt
62 เรื่อง          NN        O           I_CLS
63 ดังกล่าว        AJ        O           I_CLS
64 ไม่             NG        O           I_CLS
65 กระทบ          VV        O           I_CLS
66 สัมพันธ์        NN        O           E_CLS
67 ที่             CC        O           B_CLS
68 มี              VV        O           I_CLS
69 ต่อ             PS        O           I_CLS
70 กัน             PR        O           E_CLS
71
72 พร้อม           CC        O           B_CLS
73 หนุน            VV        O           I_CLS
74 กกต.            NN        B_ORG       I_CLS
75 จัด             VV        O           I_CLS
76 เลือกตั้ง        VV        O           E_CLS
77 ให้             CC        O           B_CLS
78 บริสุทธิ์         VV        O           I_CLS
79 ยุติธรรม         VV        O           E_CLS
80
81 (               PU        O           B_CLS
82 1               NU        B_DTM       I_CLS
83 _               PU        I_DTM       I_CLS
84 ก.ย.            NN        E_DTM       I_CLS
85 )               PU        O           I_CLS
86 _               PU        O           I_CLS
87 พล.อ.           NN        B_DES       I_CLS
88 สุรยุทธ์         NN        B_PER       I_CLS
89 _               PU        I_PER       I_CLS
90 จุลานนท์        NN        E_PER       I_CLS
:set tabstop=16                    86,2-16        15%
```

- Four columns
  - Word
  - POS tag
  - Named entity
  - Clause boundary
- Notes
  - Each column is separated by a tab
  - Empty line marks sentence boundary

# POS Tagset

| Tags | Names | Brief Descriptions |
|------|-------|--------------------|
| AJ | Adjective | Attribute, modifier, or description of a noun |
| AV | Adverb | Word that modifies or qualifies an adjective, verb, or another adverb |
| AX | Auxiliary | Tense, aspect, mood, and voice |
| CC | Connector | Conjunction and relative pronoun |
| CL | Classifier | Class or measurement unit to which a noun or an action belongs |
| FX | Prefix | Inflectional (nominalizer, adjectivizer, adverbializer, and courteous verbalizer), and derivational |
| IJ | Interjection | Exclamation word |
| NG | Negator | Word of negation |

| Tags | Names | Brief Descriptions |
|------|-------|--------------------|
| NN | Noun | Person, place, thing, abstract concept, and proper name |
| NU | Number | Quantity for counting and calculation |
| PA | Particle | Politeness, intention, belief, question |
| PR | Pronoun | Word used to refer to an element in the discourse |
| PS | Preposition | Location, comparison, instrument, exemplification |
| PU | Punctuation | Punctuation mark |
| VV | Verb | Action, state, occurrence, and word that forms the predicate part |
| XX | Others | Unknown category |

* Green texts = content word | Black texts = function word | Red texts = undesirable (yet they still exist)

# Named Entities

| Tags | Names | Descriptions |
|------|-------|--------------|
| TTL | Title | Family relationship, social relationship, and permanent title |
| DES | Designation | Position and professional title |
| PER | Person | Name of a person or family |
| ORG | Organization | Name of organization, office, or company |
| LOC | Location | Name of land according to geo-political borders |
| BRN | Brand | Name of brand, product, and trademark |
| DTM | Date and time | Time or a specific period of time |
| MEA | Measurement | Measurement unit and quantity of things |
| NUM | Number | The number of a measurement unit |
| TRM | Terminology | Domain-specific word |

- BIEO tagging convention
  - B_$tag$ = beginning
  - I_$tag$ = intermediate/inside
  - E_$tag$ = ending
  - O_$tag$ = outside
- Notes
  - BIEO convention can be converted to BIO by replacing all E_$tag$'s with I_$tag$

**Green texts = personal entity** | **Red texts = collective entity** | **Black texts = referential entity**

# Clause Boundary

- **Definition:** clause = a part of sentence that contains at least one verb

- Syntactic clues for clause boundary
  - Subordinate connector e.g. ซึ่ง, ถ้า, ว่า
  - Cohesive marker e.g. อย่างไรก็ตาม, นอกจากนี้
  - List marker e.g. เช่น, ได้แก่, ตามลำดับ
  - Particle e.g. ครับ, นะ
  - Question adverbs e.g. อย่างไร, ไหม

# Sentence Boundary

- **Definition:** sentence = (1) a group of at least one clause, or (2) a phrase acting as a topic

- Clues for sentence segmentation

  - Topic shift denoted with a cohesive marker: **break**

  - Subject shift between two adjacent clauses: **break**

  - Direct speech "………": **concatenate** inside quotes

  - Indirect speech: **break** when the subject changes

  - Item list: **concatenate** all items

  - Particle e.g. ครับ, นะ: **break**

# Statistics

# News Genres

| | |
|---|---|
| Politics | Royal |
| Crime and Accident | Disaster |
| Economics | Urban development |
| Entertainment | Environment |
| Sports | Culture |
| International | Weather forecast |
| Science, Technology, and Education | General |
| Health | |

- Specification
  - **Size:** 3,745 articles
  - **Sources:** Thairath, Dailynews, Manager, Matichon, Nation, Prachachat Business
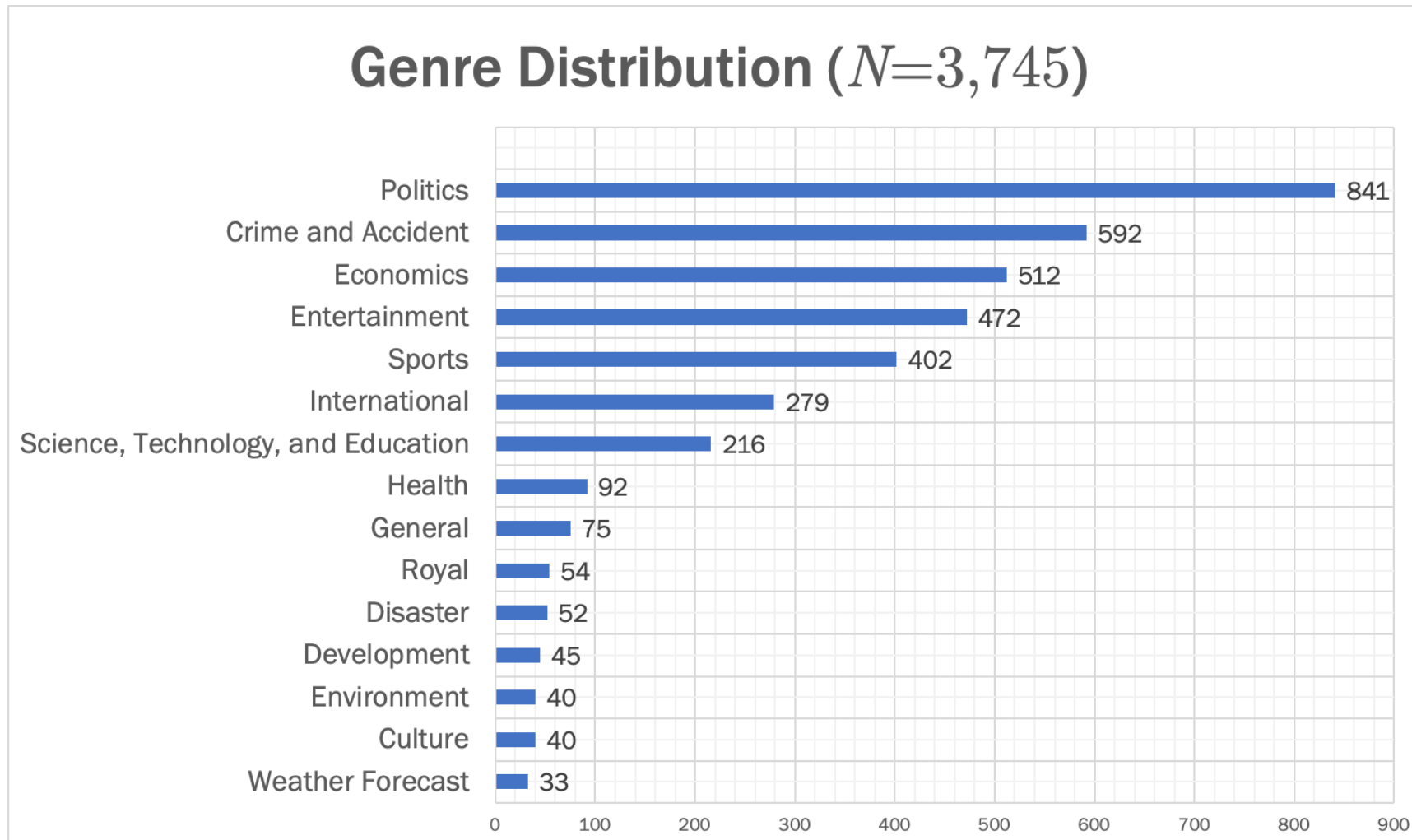  - **Duration:** JAN 2003 to DEC 2009

# Corpus Size

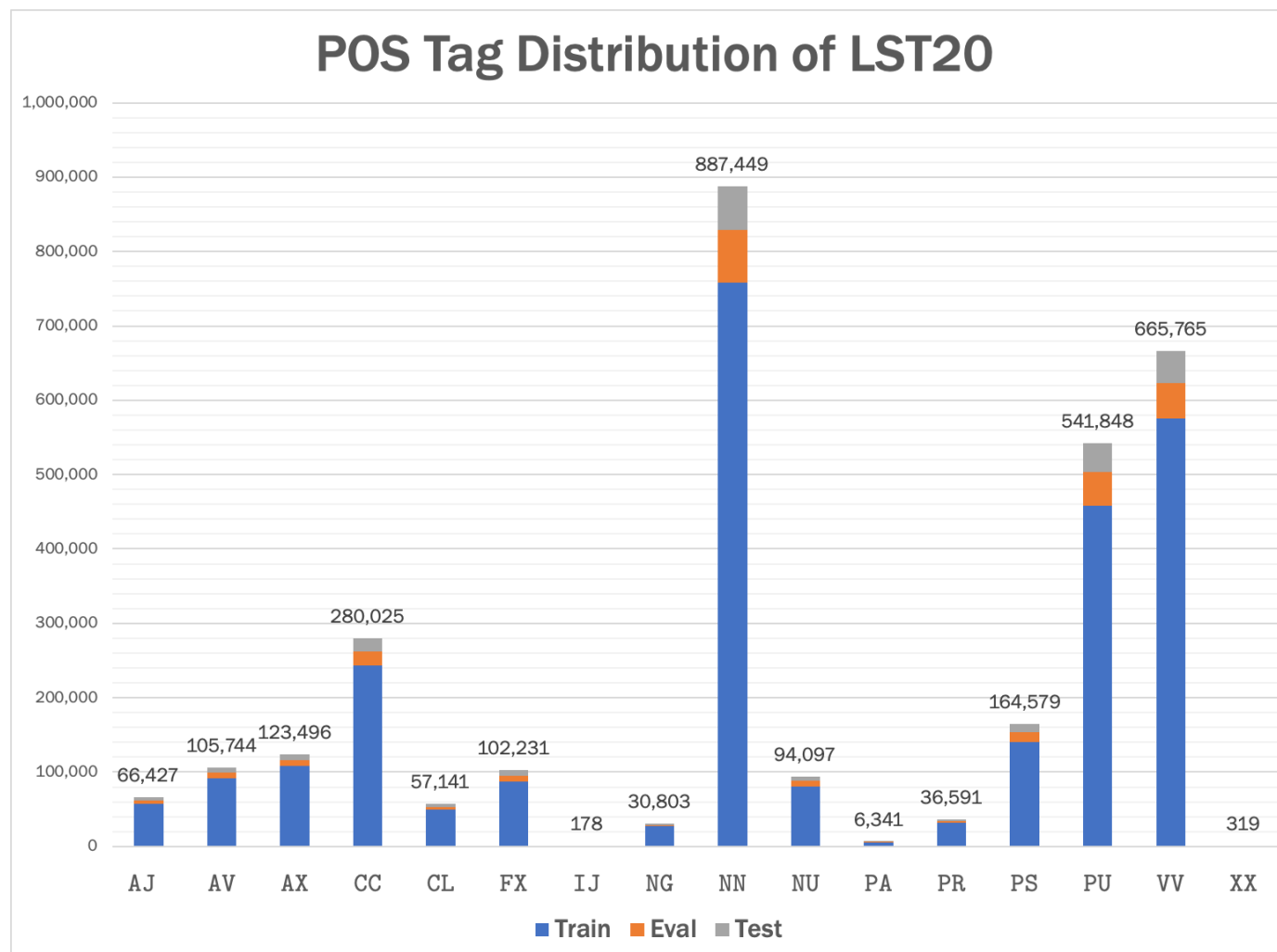| | Training | Evaluation | Testing | Total |
|---|---|---|---|---|
| Words | 2,714,848 | 240,891 | 207,295 | 3,163,034 |
| Named entities | 246,529 | 23,176 | 18,315 | 288,020 |
| Clauses | 214,645 | 17,486 | 16,050 | 248,181 |
| Sentences | 63,310 | 5,620 | 5,250 | 74,180 |
| Distinct words | 42,091 | (oov) 2,595 | (oov) 2,006 | 46,692 |
| Breaking spaces* | 63,310 | 5,620 | 5,250 | 74,180 |
| Non-breaking spaces** | 402,380 | 39,920 | 32,204 | 475,504 |

* Breaking space = space that is used as a sentence boundary marker
** Non-breaking space = space that is not used as a sentence boundary marker

# Statistics: New Genres



**Genre Distribution ($N=3{,}745$)**

| Genre | Count |
|---|---|
| Politics | 841 |
| Crime and Accident | 592 |
| Economics | 512 |
| Entertainment | 472 |
| Sports | 402 |
| International | 279 |
| Science, Technology, and Education | 216 |
| Health | 92 |
| General | 75 |
| Royal | 54 |
| Disaster | 52 |
| Development | 45 |
| Environment | 40 |
| Culture | 40 |
| Weather Forecast | 33 |

# Statistics: POS Distribution



POS Tag Distribution of LST20

# Thank You