



Abstract

Processing Thai language linguistically and computationally has been very challenging because of the linguistic characteristics and writing system of Thai. While developing a Thai dependency treebank, such as UD Thai-PUD, the number of attempts to resolve such challenges has not been quite successfully done. One of the main reasons is that there are no Thai-specific guidelines of the relevant annotation layers. Through the Thai analysis influenced by some other languages, annotating Thai corpora for a Thai treebank might be inconsistent, and might not present the actual structures of Thai. The purpose of the project therefore is to develop a new Thai UD treebank with the Thai-specific guidelines of the relevant annotation layers: word segmentation, sentence segmentation, and POS tagging to enable consistent annotation, to analyze the syntactic structures of Thai through the UD as well as SUD frameworks and to present the actual structures of Thai without the analysis influenced by some Indo-European grammar.

Thai Language

Sample of Written Thai Text

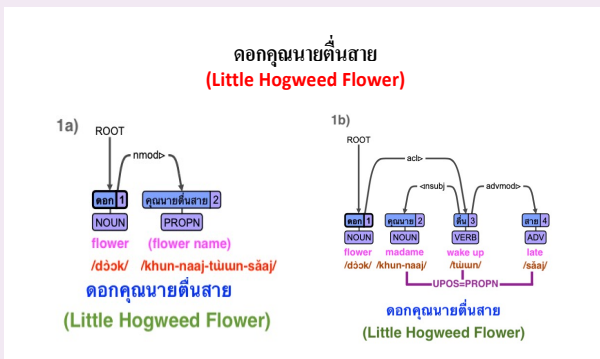
- SVO and isolating language with 4 tone markers but 5 tone sounds
- No grammatical markers
- Postmodifiers mostly used to express additional information
- Cardinal numbers and quantifiers placed before a modified noun
- Thai words – very fluid and going through grammaticalization
- Compounding mostly used to coin new words
- A very long list of classifiers used to classify nouns
- Context – a crucial clue of syntactic and semantic information
- Language use and style – relying on the hierarchical society and cultural aspects

- Most words not delimited
- No capital letters in Thai; a very little of Thai and standard punctuation used
- Whitespace used between clauses and phrases; double whitespace used to end a sentence

ดอกกุหลาบต้นสายริมทางเหลืองแต่ศรีอยุธยาแต่หัวแห่งเสียดิน ถ้าปล่อยทิ้งไว้คงจะแห้งตายไปเอง แต่ที่ร้านข้าวมันไก่ สองตาดายจะลงแปลงปลูกใหม่ทุกครั้ง เราสองคนไม่ได้แวะไปนาน ต่อมาวันหนึ่งเห็นต้นกุหลาบต้นสายริมทางเหลืองแต่ศรีอยุธยาแต่หัวแห่งเสียดิน อยู่ริมทางเกิดคิดถึงระหงษ์นางแทนที่ อย่างกับมาติดที่ ทุกอย่างใหม่หมดจดไม่เหลือเค้าเดิม ไม้รูปไปไหน ย้ายไปไหนทำไมแล้ว หรือจะเจ็บป่วยไม่สบาย หรือจะ... ขอย้ายให้เป็นอย่างนั้น ความกลัวก็พาความคิดเล็ดลอดไป ชายคนนั้นคงไม่ได้ยินมีหรือรับฟังหรือ ข้าวมันไก่ซึ่งงานพร้อมมีรูปไปโรงโกล่งใหญ่อีกแล้ว ห้องฉันหรือใครครากัดเพื่อด้วยความคิดหวัง ความกลัว ก็วางใจไม่อาจรู้ชื่อในน้ำลายค้ายิ้ม มันทนอดเหนียวอยู่เต็มปากเต็มคำคือ หน้ะมีติดตายจนเห็นประหลาดปานกลายเป็นกระเพาะใหญ่ยักษ์ที่กำลังเปื่อยคด ซุกกลิ่น ข้อนเรอออกมาเคี้ยวแล้วกลืนเข้า เคี้ยวกลืน เคี้ยวกลืน เข้าแล้วเข้าแล้วไม่หยุดหย่อน ความพิศพิศพิศอยู่ในกระเพาะจนคนที่กำลังเดินผ่าน คนที่ไม่รู้จัก คนที่ไม่ถูกมองเห็น นมผง นมข้นสำเร็จรูป ขนมปัง ปลากระป๋อง ข้าวสาร กล้วยเป็นหลักฐานของกลางทางแต่ตรงหน้าผู้ต้องหาลักขโมย เจ้าของในหน้านั้นจะเป็นใครก็ได้ โชคร้ายในหน้าข้าวราวยาว ฉันนึกถึงความลำบากแค้นที่ย้ายเคยเล่า สมัยนั้นคนเราลำบากยากจนกันมาก... แต่มันก็ยังไม่ว่าสมใจฉันเลยขยาย ถ้าหวาน ประชาชนอย่างเราจะยอมให้สิ้นรสชาตลิ้นในทันทีสุด แต่ถ้าเข้ม ไม่ว่าจะสักแค่ไหนเราก็จะฉีกลิ้นมันลงไปเหมือนที่ผ่านๆ มา

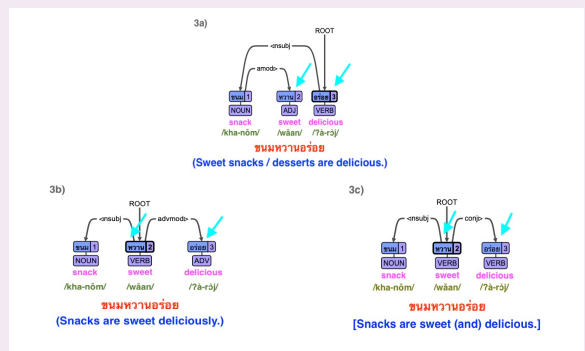
Word Segmentation

- Due to the way of how to write a Thai text, it is difficult to decide whether to segment each word separately or combine them as one single word.



POS Tagging

- Due to Thai words being very fluid, it is difficult to categorize Thai words into each part of speech.



Sentence Segmentation

- Due to the way of how to write a Thai text, it is difficult to find sentence boundaries even with some whitespace between phrases and clauses used in a text.

2a) **ช**างเลี้ยงวัวลูกไปล้างมือในลำห้วยแล้วกลับมาล้างที่เก่า **ล**้างหาขวดเหล้าในยามที่เหลือน้อย **ก**ันขวดออกมาระคดคิม แล้วลงมือแกะห่อข้าว **บ**ั้นข้าวก้อนหนึ่งวางบนใบไม้ที่เตล็ดไว้เมื่อไปล้างมือ **ใ**้ใช้ปลายนิ้วจกน้ำพริกปลาทูกรอบวางลงบนข้าว **น**ำไปวางไว้บนราคะตาคที่ไหลพันดินเพราะน้ำชะเนในฤดูน้ำหลากให้เข้าที่ **เ**กกลับมาล้างอีกครึ่งลงมือจกน้ำพริกอย่างกับข้าว **จ**กเข้าปากแล้วกินผักกูดตาม

Plans for Thai UD Treebank

- Collecting around 5,000 Thai sentences from different genres of written Thai texts to build a Thai corpus for a Thai UD treebank
- Manually and/or automatically annotating the corpus with the UD framework, including the SUD framework
- Also contributing Thai annotations to the treebank of UD Thai-TDT; correcting the Thai corpus of UD Thai-PUD by following the guidelines to be developed

Promising Procedures and Results

- Instead of following the conventional sequence of annotation, starting firstly with sentence segmentation by depending mainly on the context of a text to gather what is discussed in the same topic
- And later looking for linguistic clues, such as topic changing, sentence subjects, proper nouns, personal pronouns, conjunctions, discourse markers, final particles, etc., to find a sentence boundary
- Segmenting the words of each segmented sentence manually and/or automatically with the Thai Tokenization tool, https://huggingface.co/spaces/pythainlp/newmm_online by PyThainlp
- POS tagging with UPOS by analyzing the true structures of Thai with the framework of the syntactic distribution
- Syntactically annotating the tagged words and segmented sentences with the UD framework, including the SUD one
- To develop the Thai-specific guidelines out of these annotation layers to perform better consistent annotations; and hopefully to introduce the Thai-specific guidelines to any Thai annotations and/or studies