

SE 953482 Natural Language Processing for SE 66/2 NLP Overview

Asst. Prof. Pree Thiengburanathum

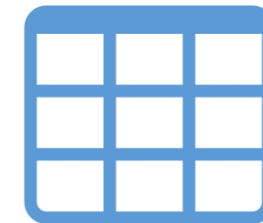
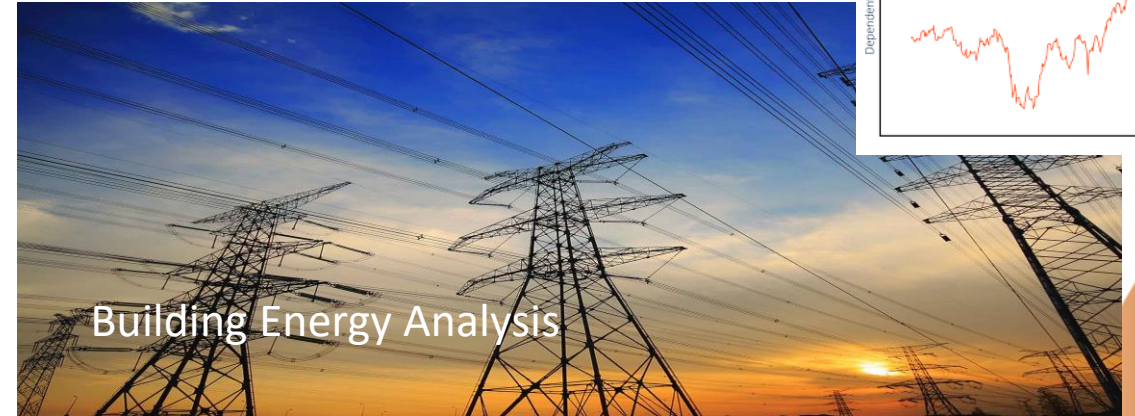
Agenda

- Course outline
- Intro to unstructured data (NLP)

Education

- **Ph.D. in (Computing and Informatic) Bournemouth University, 2013**
– *Mar 2018, Faculty of Science and Technology, Department of Computing and Informatics, Bournemouth University, UK.*
- **ERASMUS MUNDUS Research fellow, Feb 1, 2010- Dec 1, 2010,**
Universite De Lyon 2, France
- **Master of Science in Computer Science, Fall 2008, University of Colorado at Denver, Denver, Colorado, USA.**
- **Bachelor of Science in Computer Science, 2005, Colorado State University, Fort Collins, Colorado, USA.**

Area of interests



Master students Thesis/IS

- (In progress) 2023 – Co-supervisor Trpnakorn, Zeroshot learning in Thai Sentiment Analysis
- (In progress) 2023- Nuttawut, Cryptocurrency and Stock Price Prediction Using Sentiment Analysis
- (In progress 2023- SIRIYAPORN Development of Predictive Model for Cyprinid Herpes Virus 2 Infection of Goldfish
- (*In progress*) 2022 – Sukanya, An Analysis of Personal Factors Influencing the Electricity Consumption in Chiang Mai University's Dormitories during the Post-COVID-19 Recovery.
- (*In progress*) 2022 – Jukkrit, Thai Spell Correction Using Pre-train BERT for End-to-End Word Extraction based on Tesseract
- (*In progress*) 2022 – Parinya, Sentiment Analysis and visualization from online restaurant review using NLP
- (*Graduated*) 2022 – Patcharapol Y. Development of Electricity Consumption Forecasting Model for Campus-Scaled Building Using Machine Learning
- (*Graduated*) 2021 – Decheng Yang, A Comparative Study of Open-Source Crawlers Based on Robustness and Scalability Testing on E-Commerce Websites
- (*Graduated*) 2021 - Li Ye, A Real-Time Bus Arrival Time Prediction System Based on Spark Framework and Machine Learning Approaches: a case study in Chiang Mai

Course outline

Course learning outcomes (CLOs): Students are able to

- 1 Explain the process in NLP and techniques.
2. Use the appropriate models and metrics tools for the right problem.

Course Description:

NLP overview, word tokenization and text preprocessing, text extraction methods, machine-learning models in NLP, Deep-learning models in NLP, Transformer, model evaluation and explain-ability, evaluation metrics, NLP-based systems, and case studies

Course outline (cont.)

- Pre-requisite SE 233 (953233), SE 411 (953411)

1.	NLP Overview	3
2.	Word Tokenization, Text preprocessing	6
3.	Text extraction methods	6
4.	Machine-learning models in NLP	6
5.	Deep-learning models in NLP	6
6.	Transformers	3
7.	Evaluation metrics and explain-ability	3
8.	NLP-based Systems	3
9.	Case studies and Project	9

Course outline (cont.)

- **Grading system:**
 - Midterm Examination 30%
 - Final Examination 30%
 - Workshops 15% (6 works)
 - Programming Assignments 25% (4-5 works)
 - **I have right to adjust grading system based on the student performance.*
 - CMU-based (i.e., a grade A cut at 80%)

Grade policy:

- Any late assignment submissions will either be penalized (at least 50% reduction) or NOT be accepted.
- If a student is late more than 15 minutes in either lab or lecture, you will be regarded as absence.
- If a student needs to be absent with legitimate causes, please notify the lecturer or TA before the date of absence.
- The student who has come to class less than 80% will NOT allow to take the FINAL EXAM.
- The student who does not take the final exam gets “F” for this course. o The work that does not strictly follow the instruction is not accepted.

Course outline (cont.)

- Course communication
 - MS Team Channel
 - pree.t@cmu.ac.th, or direct message me at Team
 - Room 415-1 (T, Th)
- Programming language and tools
 - Python3
 - Anaconda framework (optional)

Textbooks

(main) Practical Natural Language Processing: A Comprehensive Guide to Building Real-World NLP Systems 1st Edition

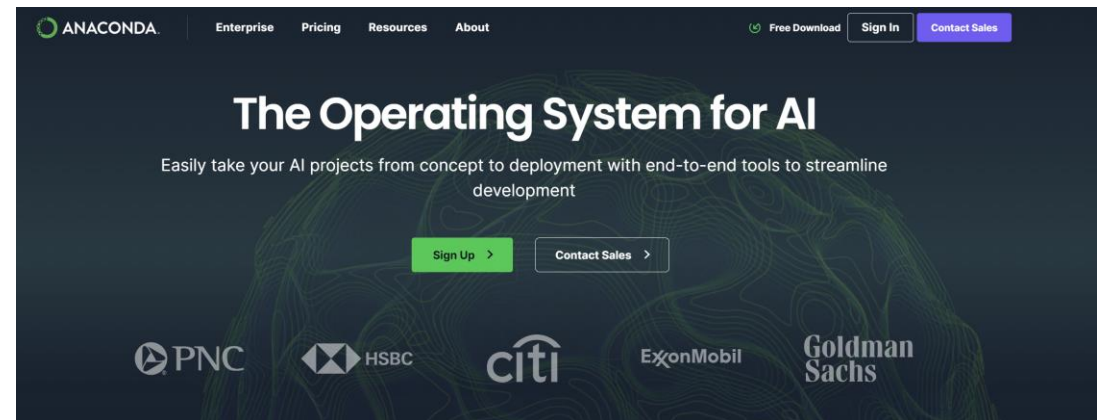
by Sowmya Vajjala (Author), Bodhisattwa Majumder (Author), Anuj Gupta (Author), Harshit Surana (Author)

(optional) Getting Started with Natural Language Processing

- by Ekaterina Kochmar (Author)

Anaconda/Miniconda

- Pre-install libraries
- Good package management
- Easy to install, maintain, and export
- Cross platforms
- <https://www.anaconda.com/>



Dev tools for the course

The Kaggle logo is presented as a stylized button. It consists of a dark blue rounded rectangle in the background, with a lighter blue rounded rectangle centered on top. The word "Kaggle" is written in a black, sans-serif font in the center of the light blue area.

Kaggle

The Colab logo is presented as a stylized button. It consists of a dark blue rounded rectangle in the background, with a lighter blue rounded rectangle centered on top. The word "Colab" is written in a black, sans-serif font in the center of the light blue area.

Colab

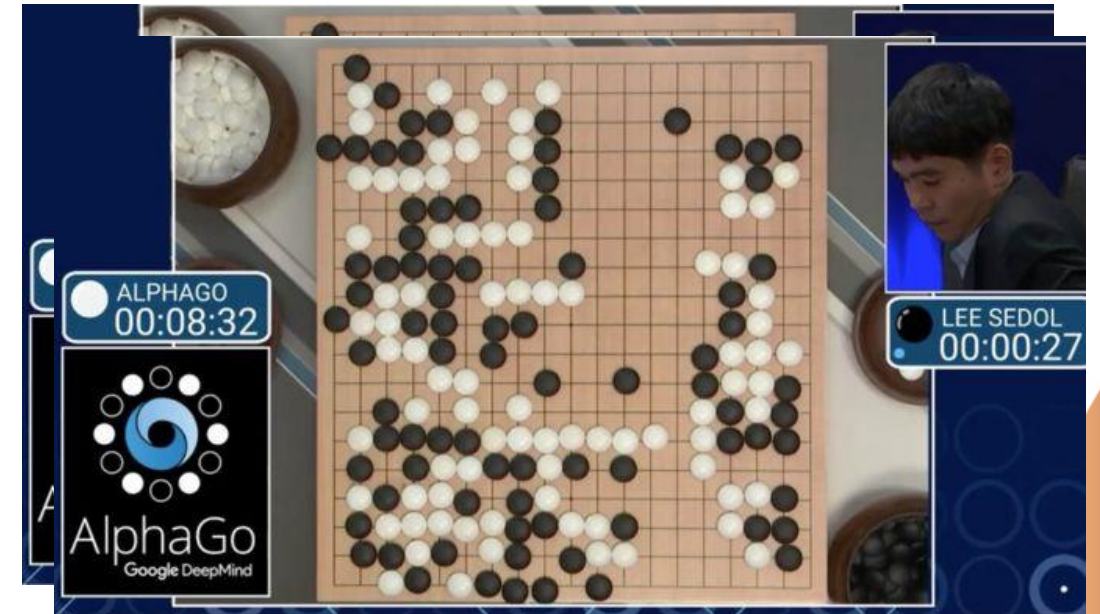
The Github logo is presented as a stylized button. It consists of a dark blue rounded rectangle in the background, with a lighter blue rounded rectangle centered on top. The word "Github" is written in a black, sans-serif font in the center of the light blue area.

Github

Introduction (cont.)

- Data continues to grow exponentially
 - Estimated to be 2.5 MTB a day
 - Grow to 40 BTB by 2020 (50 * of 2010)
- **Approx. 80%** of data is estimated to be unstructured/text-rich data
 - >4.5 billion web pages
 - >40 million articles (5 million in English)
 - >500 million tweets a day, 200 billion a year
 - >1.5 trillion queries on Google a year

Machine Versus Men



Machine Versus Men (cont.)

- Can machine beat the best of man in what man is supposed to be the best at?
- <https://www.youtube.com/watch?v=YgYSv2KSyWg>
- Watson, which is called DeepQA
- Watson had access to 200 million pages of structured and unstructured content consuming four terabytes of disk storage.

DeepQA overall architecture

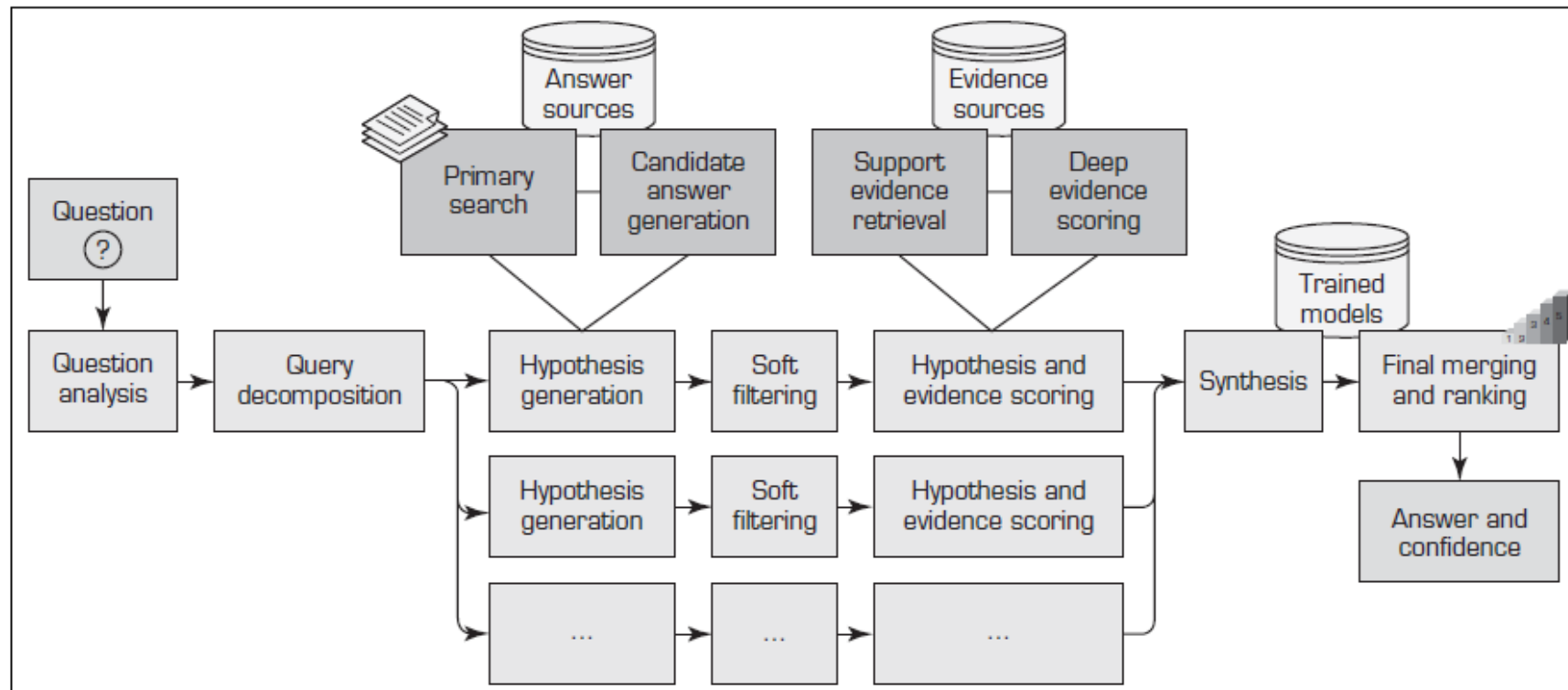


FIGURE 7.1 A High-Level Depiction of DeepQA Architecture.

Meta NLP 2022

- LaMDA – conversation robot, BERT+GPT-3
- FLORES-101 – Largest language transition dataset

~5x

Flores supports translation evaluation in 10 100 directions while the Talks data set only supports 2 162 directions

TALKS DATA SET Ye et al, 2018 46

Amharic	Esperanto	Japanese	Portuguese
Arabic	Estonian	Kannada	Russian
Asturian	Filipino	Kazakh	Serbian
Basque	Finnish	Lithuanian	Serbian
Belarusian	French	Macedonian	Silesian
Burmese	Galician	Malagasy	Sinhala
Catalan	Georgian	Malay	Slovak
Cebuano	Greek	Marathi	Telugu
Chinese	Gujarati	Mauritian	Thai
Czech	Haitian Creole	Creole	Turkish
Dutch	Hindi	Occitan	Vietnamese
English	Igbo	Pashto	

FLORES

Afrikaans	Greek	Macedonian	Serbian
Amharic	Gujarati	Malay	Sorani Kurdish
Arabic	Hausa	Malayalam	Spanish
Armenian	Hebrew	Maltese	Swahili
Assamese	Hindi	Māori	Swedish
Asturian	Hungarian	Marathi	Tajik
Azerbaijani	Icelandic	Mongolian	Tamil
Belarusian	Igbo	Nepali	Telugu
Bengali	Indonesian	Northern Sotho	Thai
Bosnian	Irish	Norwegian	Turkish
Bulgarian	Italian	Nyanja	Ukrainian
Burmese	Japanese	Occitan	Umbundu
Catalan	Javanese	Oriya	Urdu
Cebuano	Kabuverdianu	Oromo	Uzbek
Chinese Sim.	Kamba	Pashto	Vietnamese
Chinese Trad.	Kannada	Persian	Welsh
Croatian	Kazakh	Polish	Wolof
Czech	Khmer	Portuguese	Xhosa
Danish	Korean	Punjabi	Yoruba
Dutch	Kyrgyz	Romanian	Zulu
Estonian	Lao	Russian	
Filipino	Latvian	Serbian	
Finnish	Lingala	Shona	
French			
Fula			
Galician			



Meta NLP 2023



- GPT-3 4,096 and 2,049 tokens
- GPT-4 8,192 and 32,768 tokens
- Introduced “Multimodal” (e.g., can also understand image)
- Better and solve math problems
- Even more languages (with low-resources)
- Able to include reference and source of the text generated

Structured vs Unstructured data

1	Indicator ID	Dimension List	Timeframe	Numeric Value	Missing Value Flag	Confidence Inte
2	214390830	Total (Age-adjusted)	2008	74.6%		73.8%
3	214390833	Aged 18-44 years	2008	59.4%		58.0%
4	214390831	Aged 18-24 years	2008	37.4%		34.6%
5	214390832	Aged 25-44 years	2008	66.9%		65.5%
6	214390836	Aged 45-64 years	2008	88.6%		87.7%
7	214390834	Aged 45-54 years	2008	86.3%		85.1%
8	214390835	Aged 55-64 years	2008	91.5%		90.4%
9	214390840	Aged 65 years and over	2008	94.6%		93.8%
10	214390837	Aged 65-74 years	2008	93.6%		92.4%
11	214390838	Aged 75-84 years	2008	95.6%		94.4%
12	214390839	Aged 85 years and over	2008	96.0%		94.0%
13	214390841	Male (Age-adjusted)	2008	72.2%		71.1%
14	214390842	Female (Age-adjusted)	2008	76.8%		75.9%
15	214390843	White only (Age-adjusted)	2008	73.8%		72.9%
16	214390844	Black or African American only (Age-adjusted)	2008	77.0%		75.0%
17	214390845	American Indian or Alaska Native only (Age-adjusted)	2008	66.5%		57.1%
18	214390846	Asian only (Age-adjusted)	2008	80.5%		77.7%
19	214390847	Native Hawaiian or Other Pacific Islander only (Age-adjusted)	2008	DSU		
20	214390848	2 or more races (Age-adjusted)	2008	75.6%		69.6%

Figure 1.1 An Excel table is an example of structured data.

← << >> Delete Move Spam ↑ ↓ ×

New team of UI engineers

CDA@engineer.com

To xyz@program.com

Today 10:21 ★

An investment banking client of mine has had the go ahead to build a new team of UI engineers to work on various areas of a cutting-edge single-dealer trading platform.

They will be recruiting at all levels and paying between 40k & 85k (+ all the usual benefits of the banking world). I understand you may not be looking. I also understand you may be a contractor. Of the last 3 hires they brought into the team, two were contractors of 10 years who I honestly thought would never turn to what they considered "the dark side."

This is a genuine opportunity to work in an environment that's built up for best in industry and allows you to gain commercial experience with all the latest tools, tech, and processes.

There is more information below. I appreciate the spec is rather loose – They are not looking for specialists in Angular / Node / Backbone or any of the other buzz words in particular, rather an "engineer" who can wear many hats and is in touch with current tech & tinkers in their own time.

For more information and a confidential chat, please drop me a reply email. Appreciate you may not have an updated CV, but if you do that would be handy to have a look through if you don't mind sending.

← Reply << Reply to All → Forward

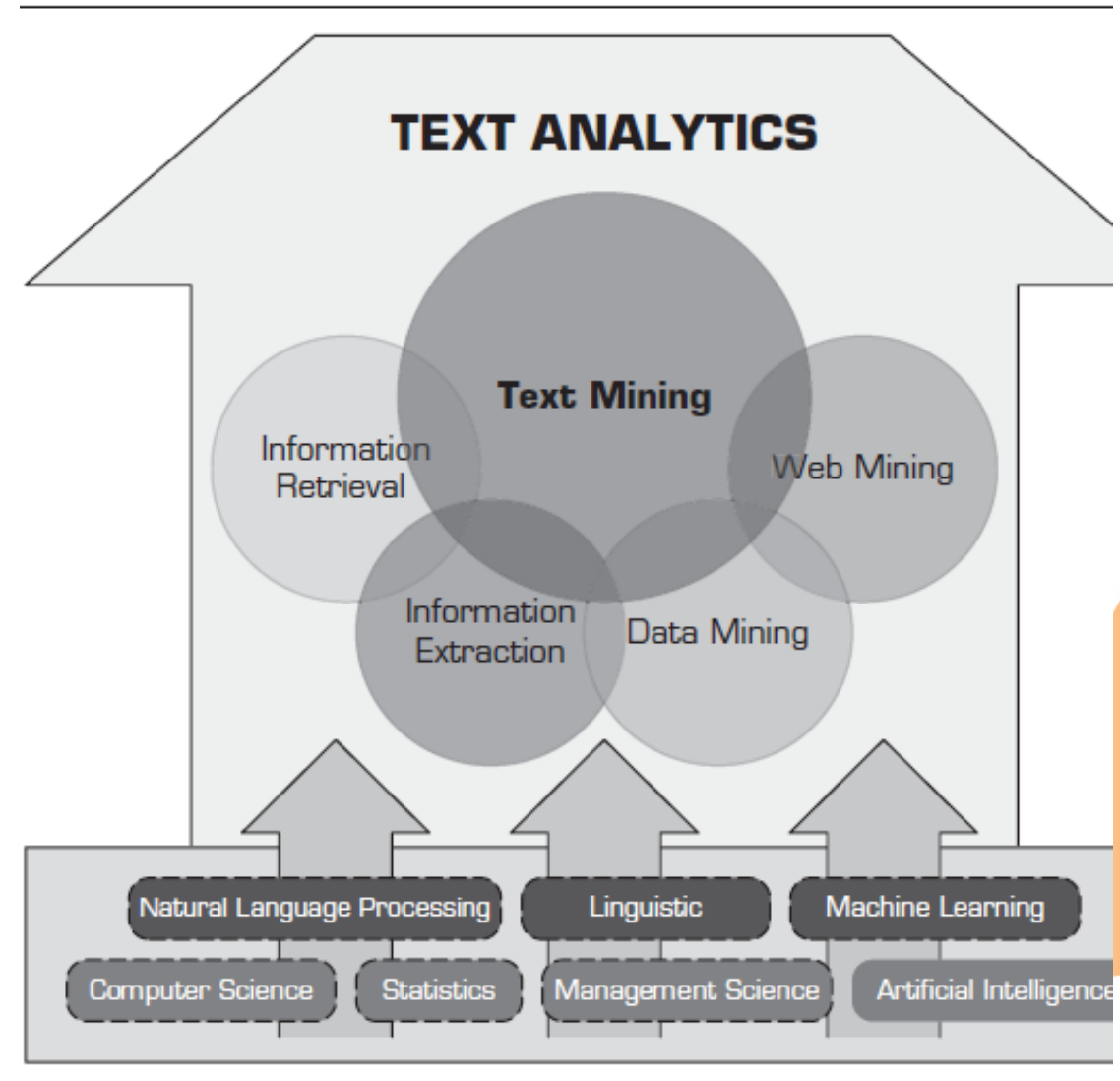
Figure 1.2 Email is simultaneously an example of unstructured data and natural language data.

Text Analytics

- The vast majority of business data is stored in text documents that are virtually unstructured.
- Text analytics is a broader concept that includes **information retrieval** (e.g., searching and identifying relevant documents for a given set of key terms) as well as **information extraction**, **data mining**, and **Web mining**,
- Whereas text mining is primarily focused on discovering new and useful knowledge from the textual data sources.

Text Analytics (cont.)

- Text Mining is a derivative of Data Mining.
- Sentimental Analysis is a derivative of Text Mining.



Text Mining

- AKA. Text data mining/knowledge discovery in textual database
- Large amount of unstructured data
- Word, PDF, XML, etc.
- Benefit domains:
 - In the areas where very large amount of textual data are being generated.
 - Could you give some example?

Text Mining (cont.)

- Law – court orders
- Academic research – research article
- Finance – quarterly report
- Medicine – discharge summaries
- Biology- molecular interactions
- Marketing – customer comments
- Social Science – Web board, Twitter , etc.
- Technology – e-mail platforms? Is Gmail the smartest?

Text Mining applications

- **Information extraction** – identify the key phrases and relationship within text
- **Topic tracking** – predict/recommend other document of interest to user.
- **Summarization** – summarizing a document to save time of the reader.
- **Categorization** – identify the main themes of a document and put to the right themes
- **Clustering** – group similar documents without having a predefined set of categories
- **Concept linking** – connects related documents by identify their shared concepts.
- **Question answering** – find the best answer to a given question

Text Mining applications (cont.)

- **Intention mining/recognition/detection** – discover user intention based on comments, reviews, tweets, blogs
- **Concept mining** – extract idea and concept from large static social media
- **Sentiment Analysis** – categorize text to sentiment polarity (pos, neg, neu)
- Topic modeling – uncover the topical structure of a large collection of docs.

NLP applies in Software Engineering

- **Code Generation and Understanding:** NLP techniques can be used to convert natural language commands into code and to help developers understand complex codebases.
- **Automated Documentation:** NLP can be used to generate and maintain technical documentation based on code changes.
- **Bug Tracking and Analysis:** It can assist in categorizing and prioritizing bugs by analyzing bug reports.
- **Customer Support:** NLP can power chatbots and support systems that interact with users to solve technical problems.
- **Code Reviews:** NLP can automate some aspects of code reviews by summarizing changes and identifying potential issues.

Level of difficulty

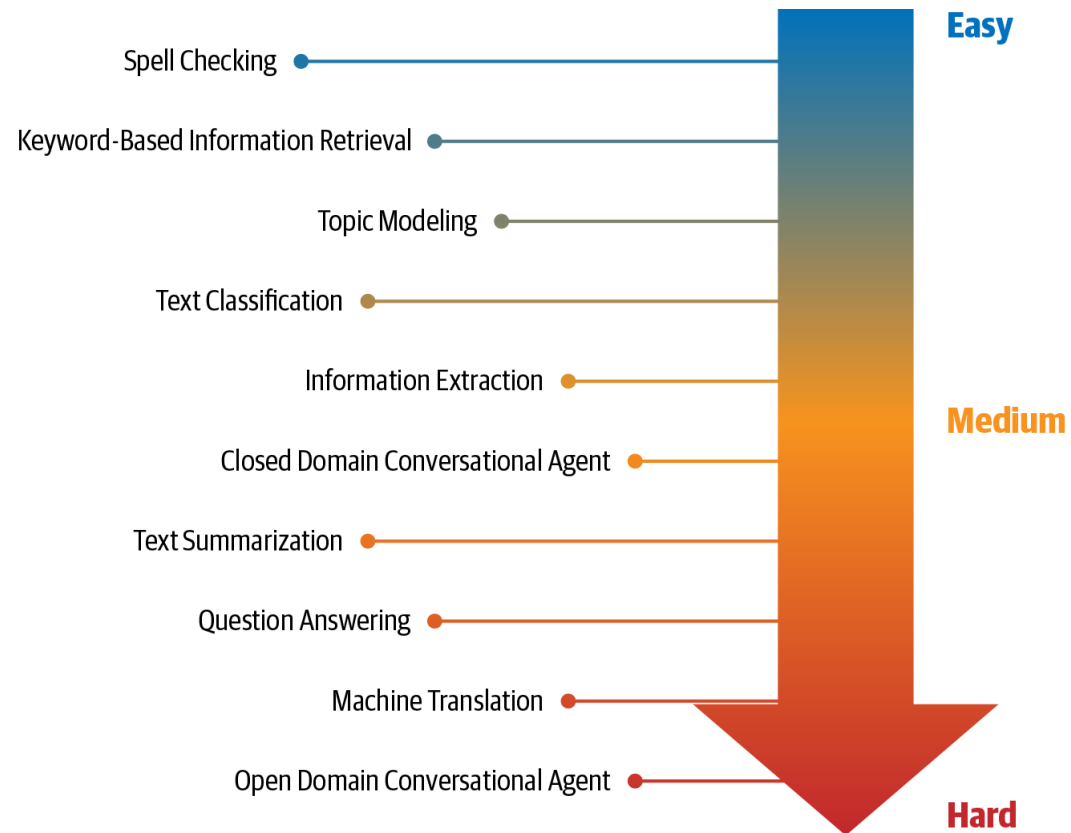


Figure 1-2. NLP tasks organized according to their relative difficulty

Intro Natural Language Processing (NLP)

- Natural language vs Programming language
 - NL - human share information with human
 - PL – human tells machines what to do
- NLP – Machine can now process natural language (i.e., interpreter)

Intro NLP (cont.)

- Subfield of AI and Computation Intelligent/linguistics.
- Try to understanding the natural human language.
- Moving forward to syntax-drive (word counting) to true understanding and processing of NLP
- Considering grammatical, semantic constraint and context.

NLP practical applications

- **Editing** – spelling, grammar, style
- **Dialog** – Chatbot, assistant, scheduling
- **Email** – spam filter, classification, prioritization
- **Text mining** – Summarization, knowledge extraction
- **News** – event detection, fact checking, fake news detection
- **Attribution** – plagiarism detection, literacy forensics,
- **Creative writing** – Movie scripts, poetry, song lyrics.
- **Search**- web, documents, autocomplete
- **Chatbot** – Ambiguous commands, Q/A, scheduling

Intro Natural Language Processing (NLP)

- Bag-of-words (classical method)
 - Text, sentences, paragraph, or document -> words
 - Classification model <spam/legitimate>
 - One bag is filled with words found in spam messages (Viagra, stock, buy)
 - Another bag is filled with words related to user's friend or workplace.
- Human do not use words without some order or structure
 - Semantic and syntactic structure
- Text mining need to look for ways beyond the bag-of-words.

Challenge in NLP (non-practical)

- **Part of speech tagging (POS-tagging)**- identify Adverb verb, noun in the sentence.
- **Text segmentation** - Chinese/Thai/Other languages.
- **Word sense disambiguation** – a word may has more than one meaning.
- **Syntactic ambiguity** – grammar is ambiguous
- **Imperfect or irregular input** – typos , grammar errors

Text mining terminologies

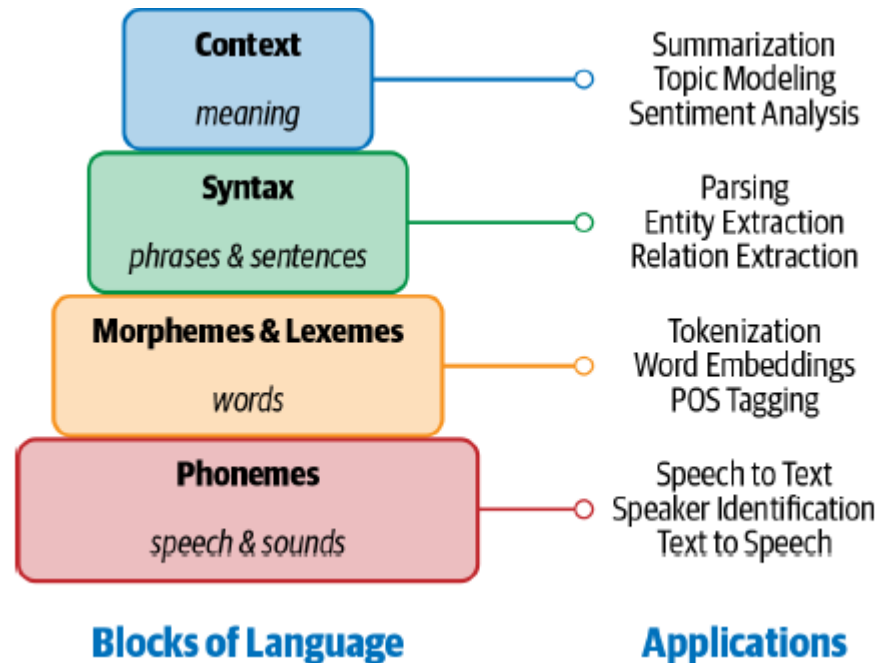
- ***Unstructured data (versus structured data).*** – human readable
- ***Corpus*** = dataset
- ***Terms*** – single word or multiword phrase from corpus
- ***Concepts*** – features generated from documents.
- ***Stemming*** – process of reducing inflected word to their root
- ***Stop words***- words that are filtered out after processed.
- ***Synonyms and polysemes*** – syntactically different/identical words
- ***Tokenizing*** – block of text
- ***Word frequency*** - #time that word occur in the document

Text mining terminologies (cont.)

- ***Morphology*** – *form and formation of words in a language*
- ***Word frequency*** – *the number of times a word is found*

Language

- “Language is a structured system of communication that involves complex combinations of its constituent components, such as characters, words, sentences, etc.”



Language (cont.)

- Phonemes – smallest unit of sound in language
 - English has 44 phonemes (single letter or combo)
 - Useful in apps like speech reg, speech-to-text/text-to-speech
- Morphemes – smallest unit of language that has meaning
 - Combination of phonemes
 - Cats = Cat + s
 - Unbreakable = Un + break + able

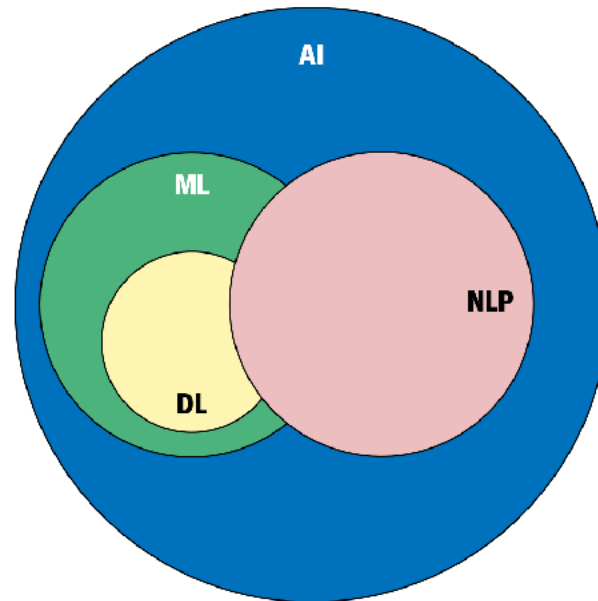
Language (cont.)

- Syntax – a set of rules to construct grammatically correct sentences
 - *Runs she fast.*
 - *She runs fast.*

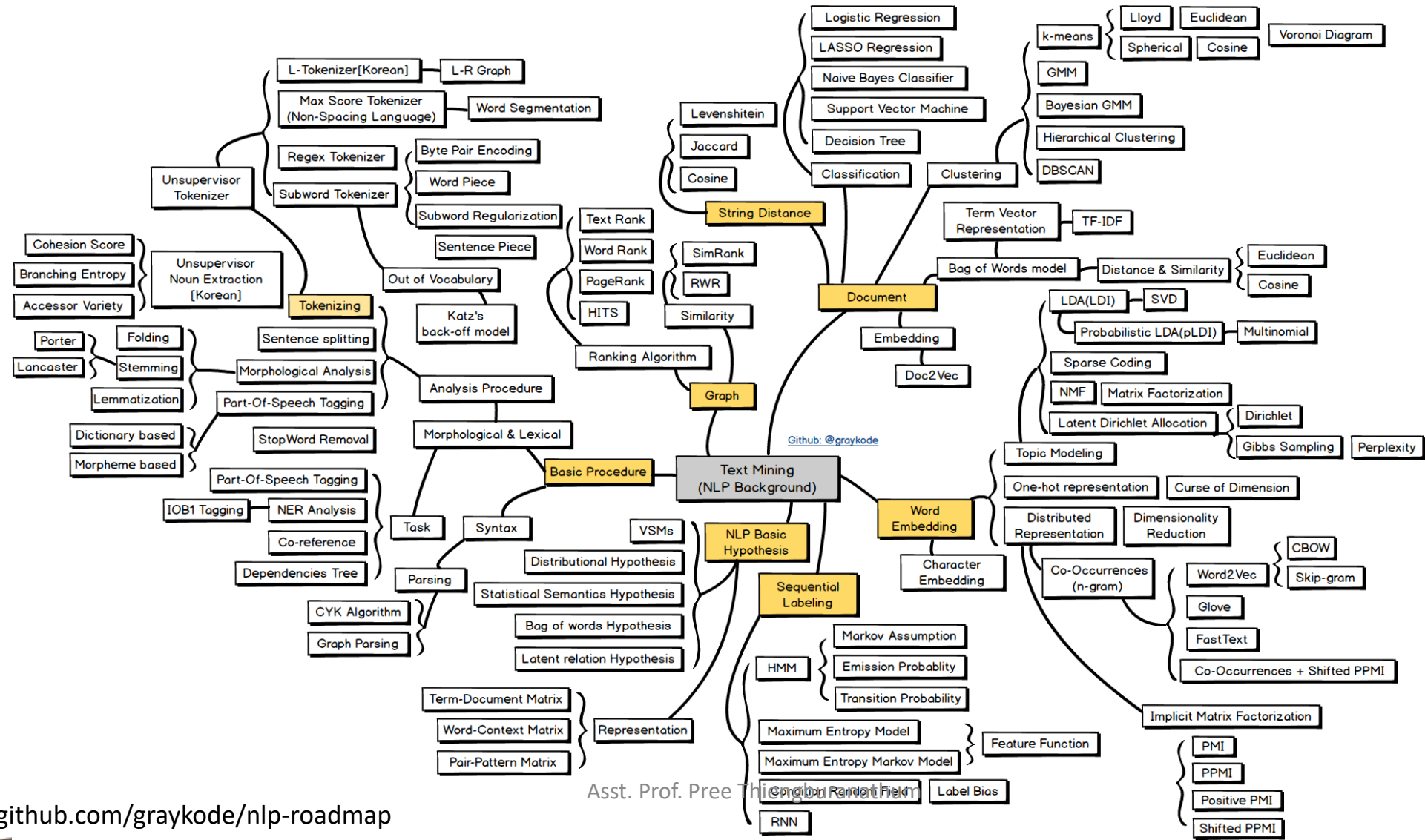
Why NLP is challenges?

- The **ambiguity** and **creativity** of human language
- Ambiguity - uncertainty of meaning. Most human languages are inherently ambiguous
 - “I made her duck.”
 - “Call me a taxi.”
 - “The teacher said the test would be difficult tomorrow.”
- Creativity – language is not a rule-based driven.
 - Various styles, dialects
 - Poem is a great example.

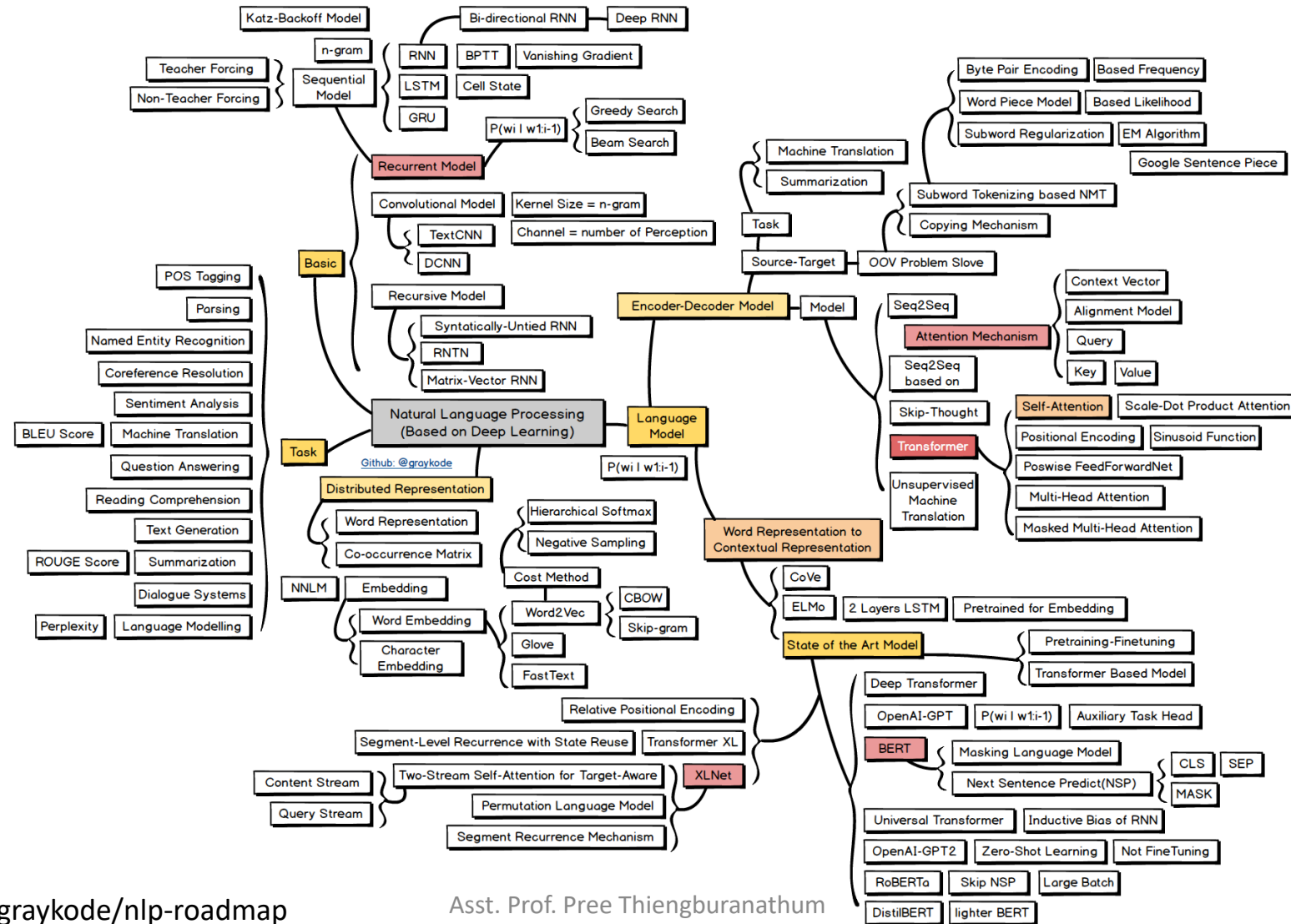
NLP related fields of studies.



NLP road map



NLP road map (cont.)



Case study 2021 / NLP application in Action

- AI chat bot with Stress detection
- [Read more at:](#)
 - <https://arxiv.org/abs/1911.00133>

Conversational Agents

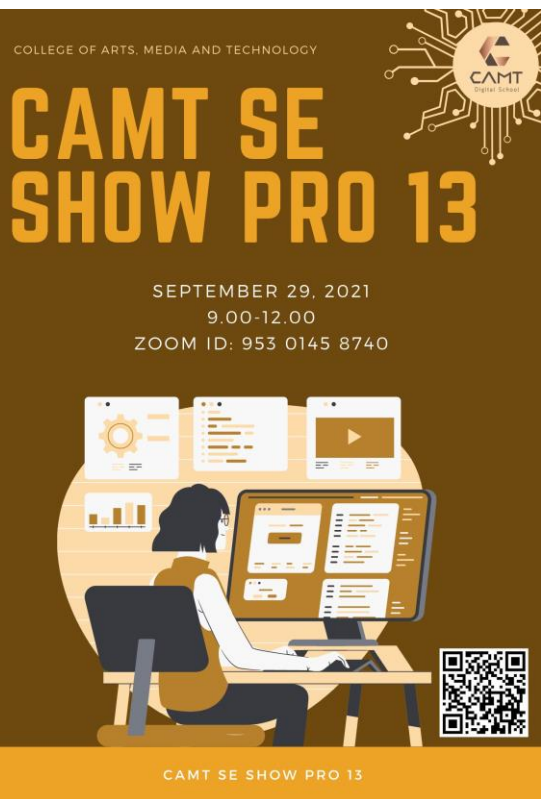
- AKA. Dialogue System, Dialogue Agents, Chatbots
- Personal Assistants on phones or other platforms
 - Alexa, SIRI, Google Assistant, Cortana
- Playing music, setting times and clock
- Chatting for fun
- Booking, scheduling reservation
- Clinical uses for mental health (like my project)

Conversational Agents (cont.)

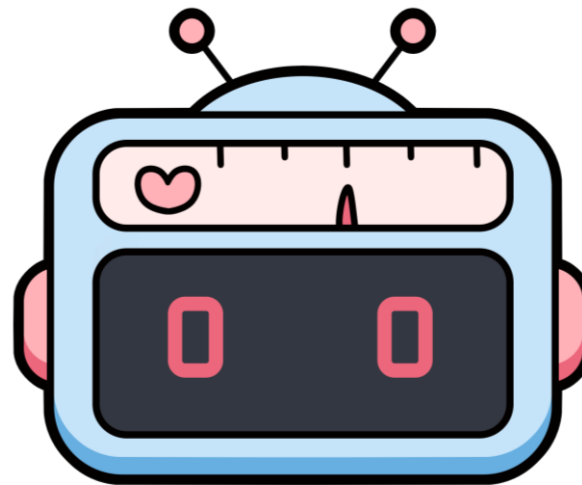
- Chatbots
 - Mimic informal human chatting
 - For fun, even for therapy
- Task-based
 - Personal assistant
 - Book seat in restaurant, movie theater, flights

Chatbot Arch.

- Rule-based
 - Pattern-action rules (ELIZA)
- Corpus-based(data-driven)
 - Information Retrieval
 - Model-based (My chatbot)



AI Chatbot with Stress Detection



Group name: *FTW_SD*

Presenter: *Pakin Kampeera (612115005)*

Aoxue Gui (612115501)

Project Advisor: *Dr. Pree Thiengburanathum*

Asst. Prof. Pree Thiengburanathum

Background

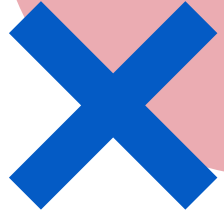
- **Stress hurts human health.** (Hathaway, 2012)

- **Most people have stress problems.**
(Mental Health Foundation, 2018; APA, 2020; United Nations, 2019; KFF, 2020)

- **Factors that affect people's active medical treatment:**
 - **Cost**
 - **Inconvenience of making an appointment**
 - **Obstacles to confiding in strangers**

- **The medical expenses invested to alleviate people's psychological problems are far from enough worldwide.** (United Nations, 2019)
- **The pandemic has caused a surge in the number of people calling the mental health hotline.** (CDC, 2020)





Motivation

we hope to develop a "self-help" tool to let limited psychological counselling services help those who need it more

Aim

Develop a chatbot to detect the user stress during communication and a dashboard to visually shows the statistical data.



Feature Overview

05.

Data cleaning, Data preprocessing, Data analysis

- Analyze sentences.

06.

Report generate

- View own stress at last 6 months

07. information

Notification

- View notification message.



Feature Overview

05.

Data cleaning, Data preprocessing, Data analysis

- Analyze sentences.

06.

Report generate

- Generate intuitive report to see stress

07. change notification during the last six notification message.

Asst. Prof. Pree Thiengburanathum



Feature Overview

01.

Authentication

- Register
- Login
- Logout
- Reset password

03.

User management

- View account information
- Change user

02.

Dashboard

- View statistic data
- View sentences table

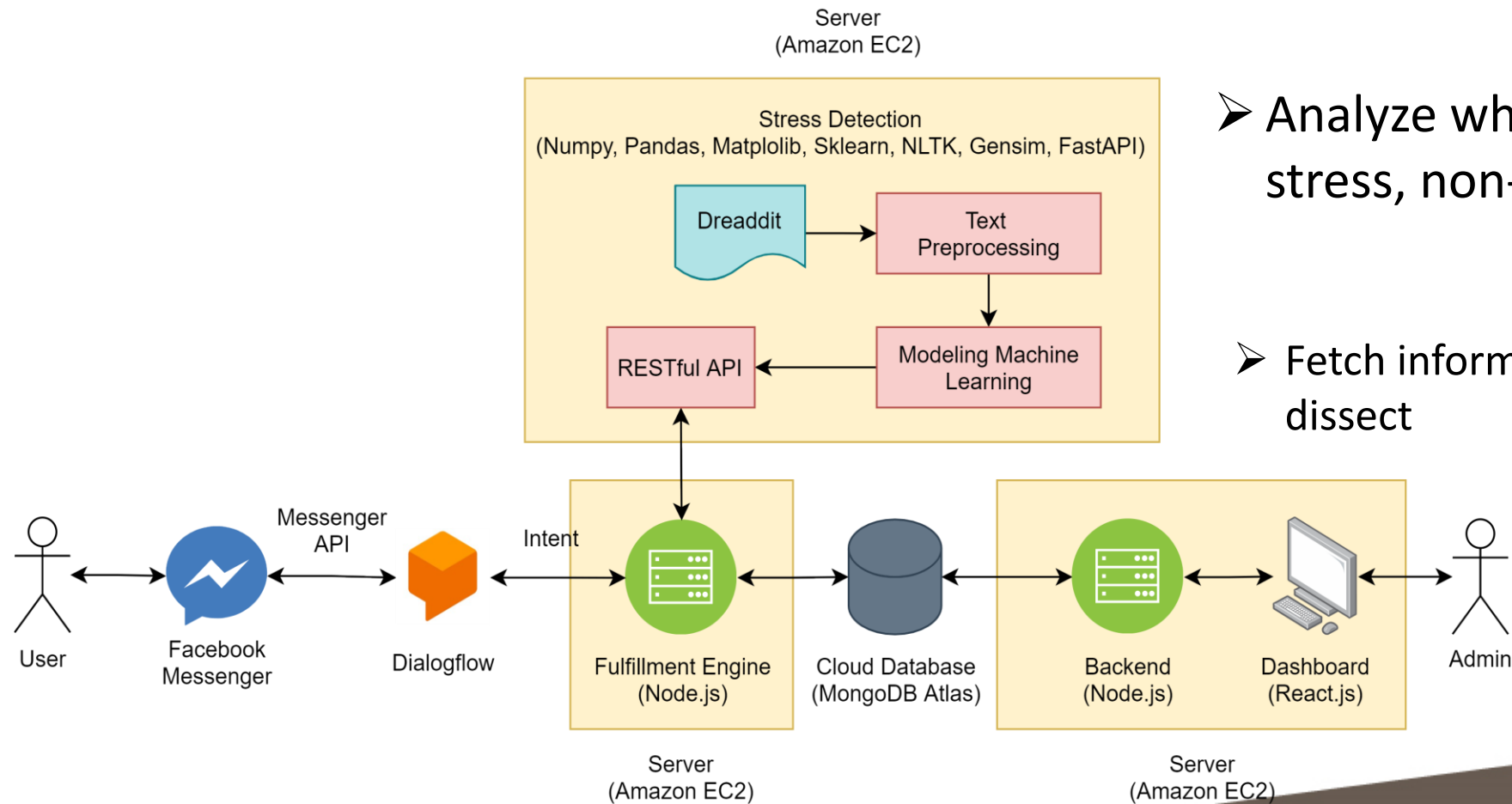
04.

Export table to Facebook messenger Chatbot

- Search username or message
- Communicate with Chatbot



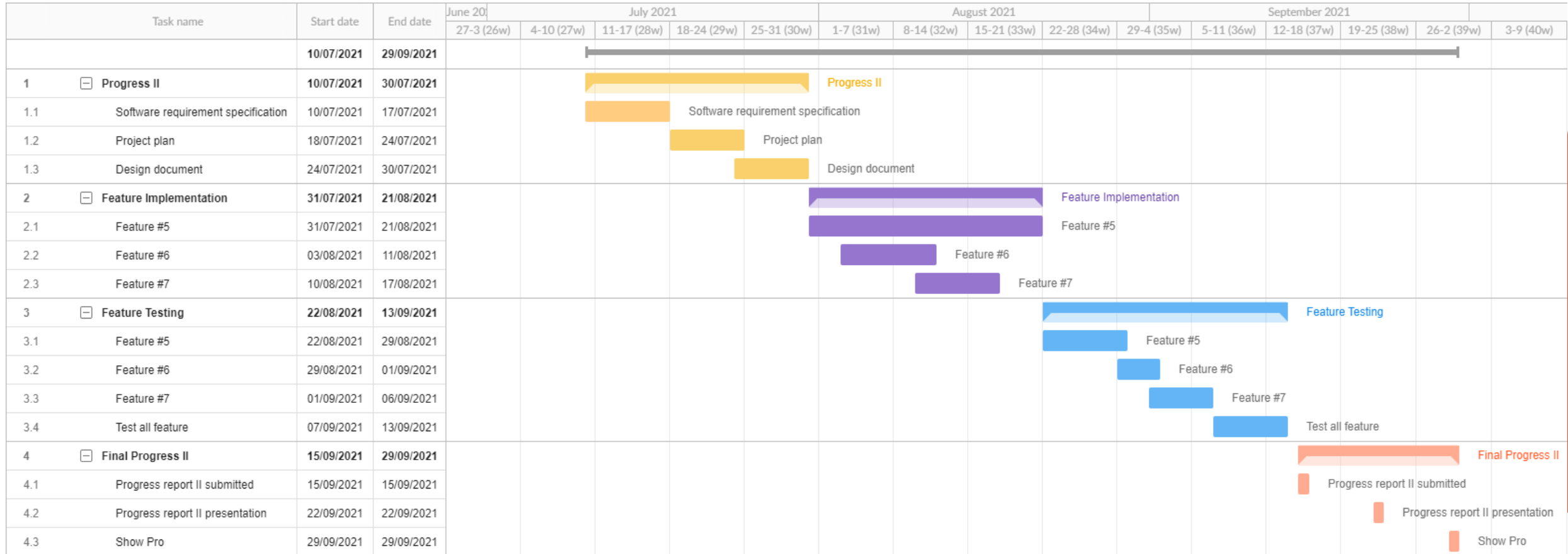
Product Perspective



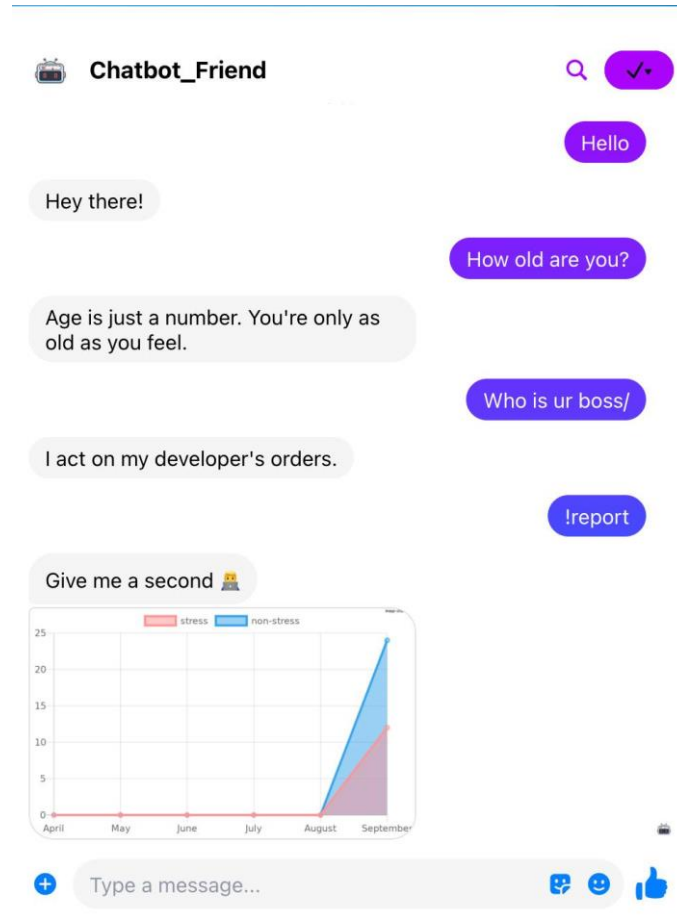
➤ Analyze whether the sentences are stress, non-stress and cannot tell

➤ Fetch information from database to dissect

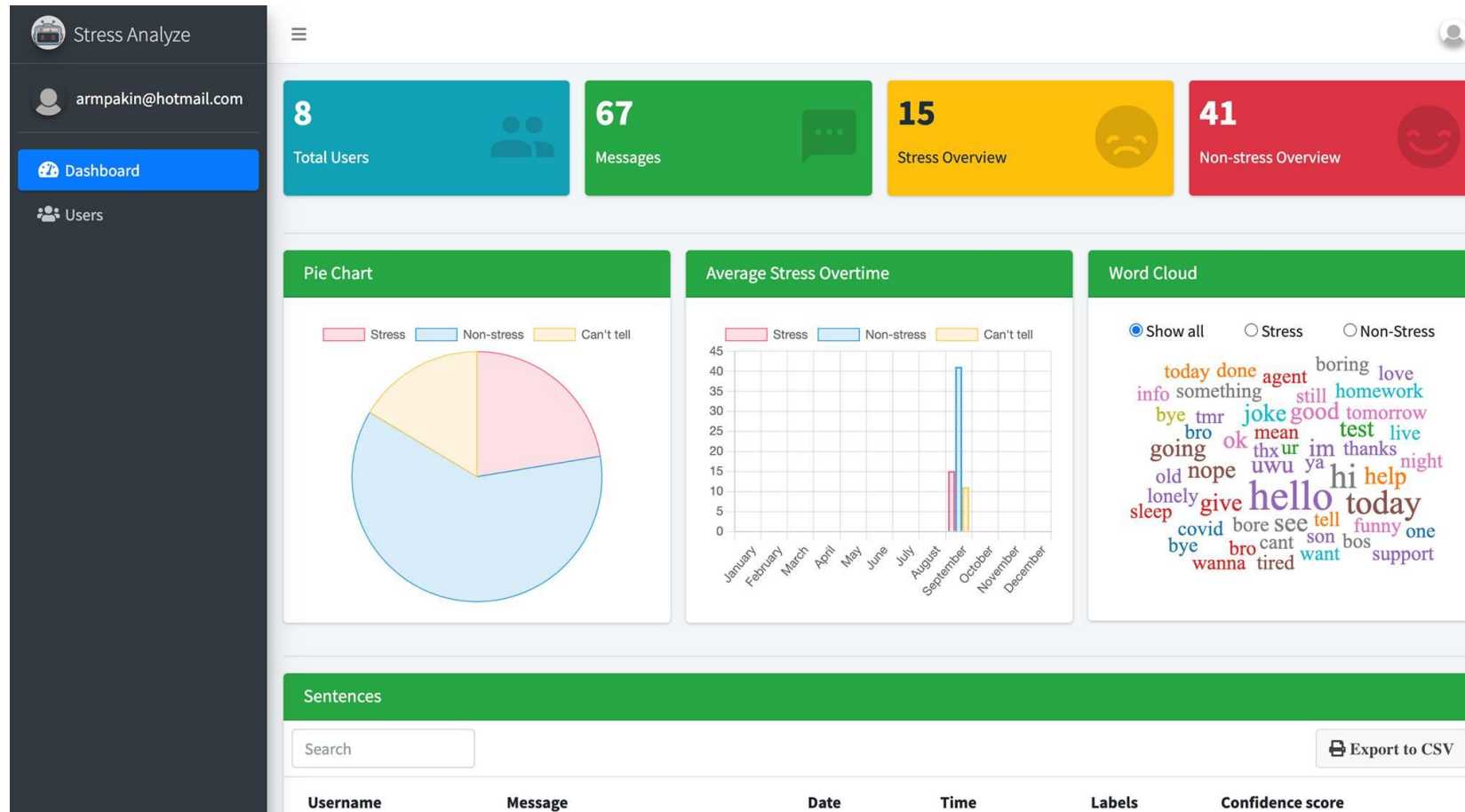
Milestone



User interface



User interface



User interface

Sentences

Search

Export to CSV

Username	Message	Date	Time	Labels	Confidence score
Pathomsakul Supamanee	UwU mean uwu	9/26/2021	11:37:08 PM	stress	0.5428668336031935
Pathomsakul Supamanee	UwU	9/26/2021	11:36:57 PM	can't tell	-
Pathomsakul Supamanee	r u winning?	9/26/2021	11:36:49 PM	can't tell	-
Pathomsakul Supamanee	hello my son	9/26/2021	11:36:42 PM	non-stress	0.13118900402470252
Pakin Kampeera	see you tmr bye	9/26/2021	9:53:01 PM	stress	0.5542619158800078
Pakin Kampeera	no bro, its still covid 19	9/26/2021	9:52:52 PM	stress	0.9356098723333347
Pakin Kampeera	today is so boring	9/26/2021	9:52:34 PM	non-stress	0.42549680705190424
Pakin Kampeera	How is it going today?	9/26/2021	9:52:22 PM	non-stress	0.42549680705190424
Pakin Kampeera	Hello	9/26/2021	9:52:08 PM	non-stress	0.060582713902607
Pakin Kampeera	see you tomorrow	9/26/2021	9:47:22 PM	stress	0.6573186070424237

10

<<

<

3

4

5

6

7

>

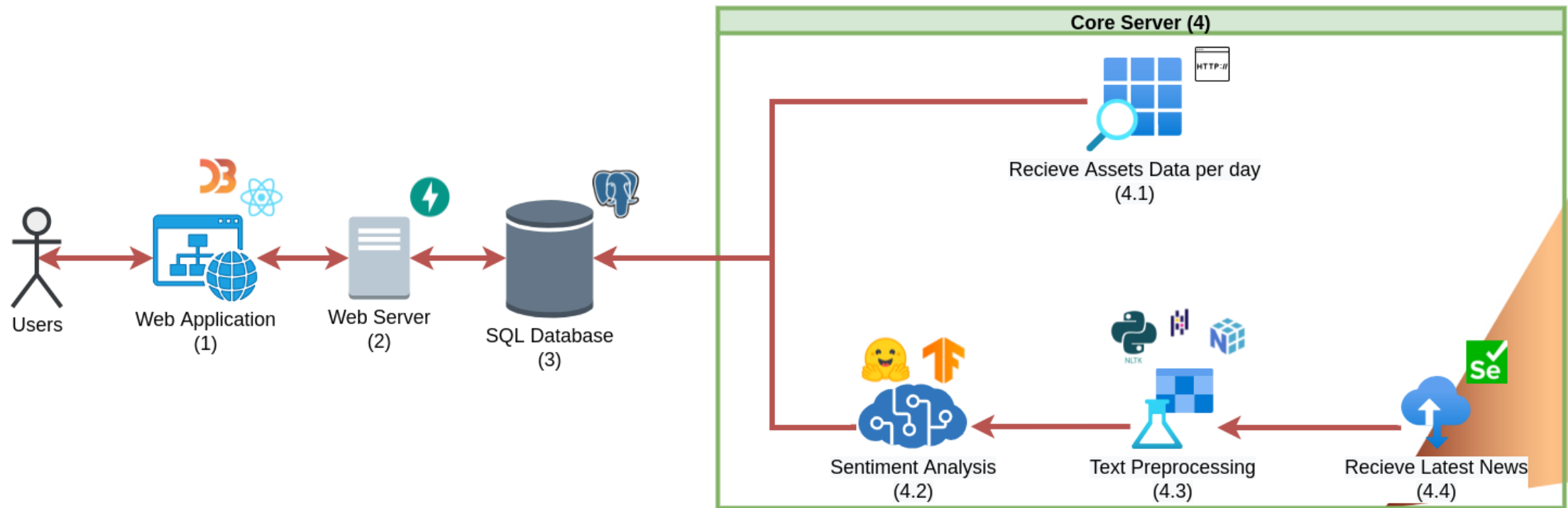
Copyright © 2021 All rights reserved

Version 0.1.0


Cryptocurrency News Sentiment Analysis

Thanatorn Kanthala 622115014

Pongpanoth Panya 622115024




UI

 Centiment
 Search API Logout pongpanoth213@gmail.com

The cryptocurrency market has come under a lot of pressure from regulators around the world. At the beginning of the final quarter, America's Commodity Futures Trading Commission (CFTC) and federal prosecutors charged one of the largest Bitcoin trading platform BIMEK with facilitating unregistered trading among other

Positive Impact
17.65% [READ MORE](#)

Crypto.com Coin, IOTA, VeChain Price Analysis:
06 October
October 7, 2020 2:30 AM




Bitcoin stood at \$10,700 and Ethereum struggled to climb past the \$355 level. The market as a whole had a bearish outlook, although many assets showed short-term bullishness. Crypto.com Coin was bearish in both the short and medium-term, while IOTA formed a bullish continuation pattern that could take some more time to

Negative Impact
23.81% [READ MORE](#)

Prior to DMALINK, Micheal spent more than three years at ADSS, having originally joined the broker as Relationship Manager.

Positive Impact
28.57% [READ MORE](#)


US Congress Considering Bill That Would Significantly Boost Blockchain's Legitimacy in Court of Law
October 7, 2020 2:15 AM



Congressional Blockchain Caucus co-chair Rep. David Schweikert has introduced a bill that will recognize digital signatures created through blockchain as valid and enforceable under federal law. Since the Electronic Signatures in Global and National Commerce Act, also known as "the e-signature bill," was signed into law in

Positive Impact
25.90% [READ MORE](#)

< previous
 1
 2
 3
 4
 5
 ...
 12
 13
 14
 next >

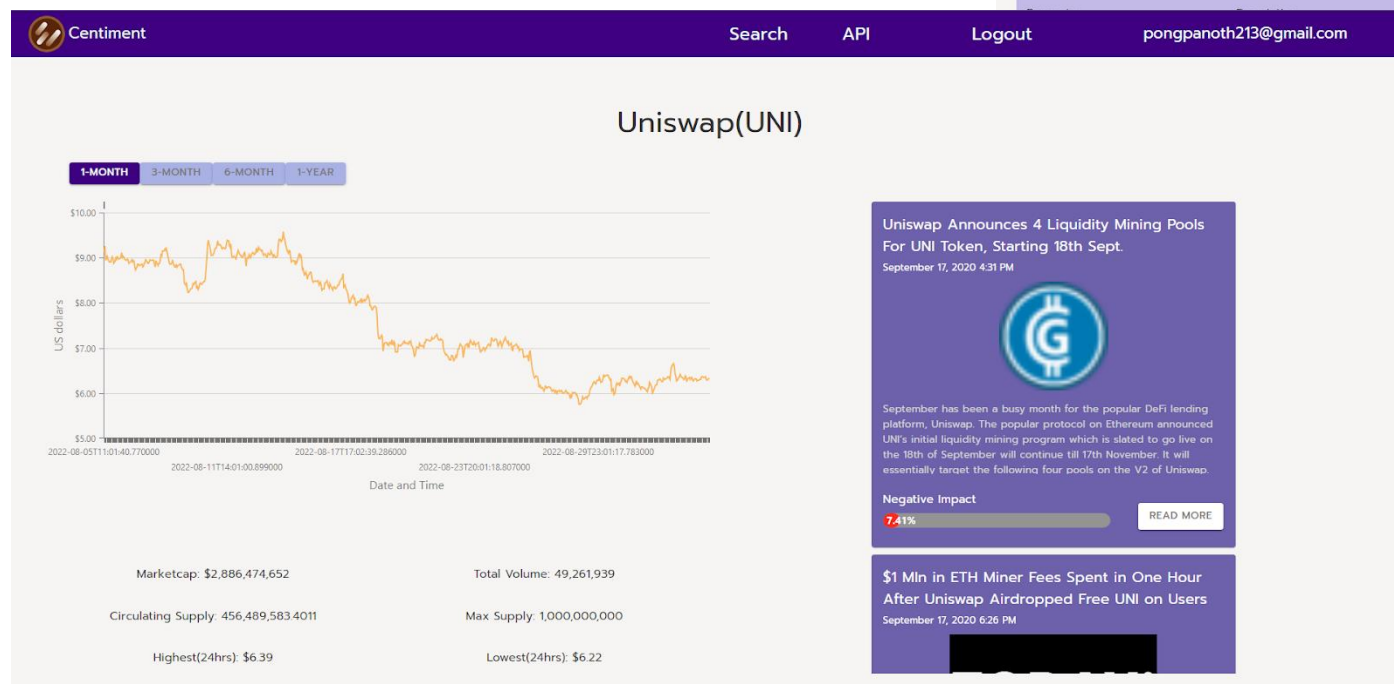
 Centiment
 Search API Logout pongpanoth213@gmail.com

Cryptocurrency List

ID	Coin Name	Coin Symbol
1	Bitcoin Cash	BCH
2	Uniswap	UNI
3	Defi Token	DEFI
4	Function X	FX
5	Binance USD	BUSD
6	Tether	USDT
7	Bitcoin	BTC
8	Ethereum	ETH
9	VeChain	VET
10	Utrust	UTK

< previous
 1
 2
 3
 4
 5
 ...
 8
 9
 10
 next >

UI



Centiment Search API Logout pongpanoth213@gmail.com

API Documentation

Get News

`/news?limit={limit}&offset={offsetNo}`

Parameter	Description
limit [integer]	The maximum number of entries to return. If the value exceeds the maximum, then the maximum value will be used
offsetNo [integer]	The (one-based) offset of the first item in the collection to return.

Search Coin Name or Symbol

`/tags?search={searchText}&limit={limit}`

GPT-Baker

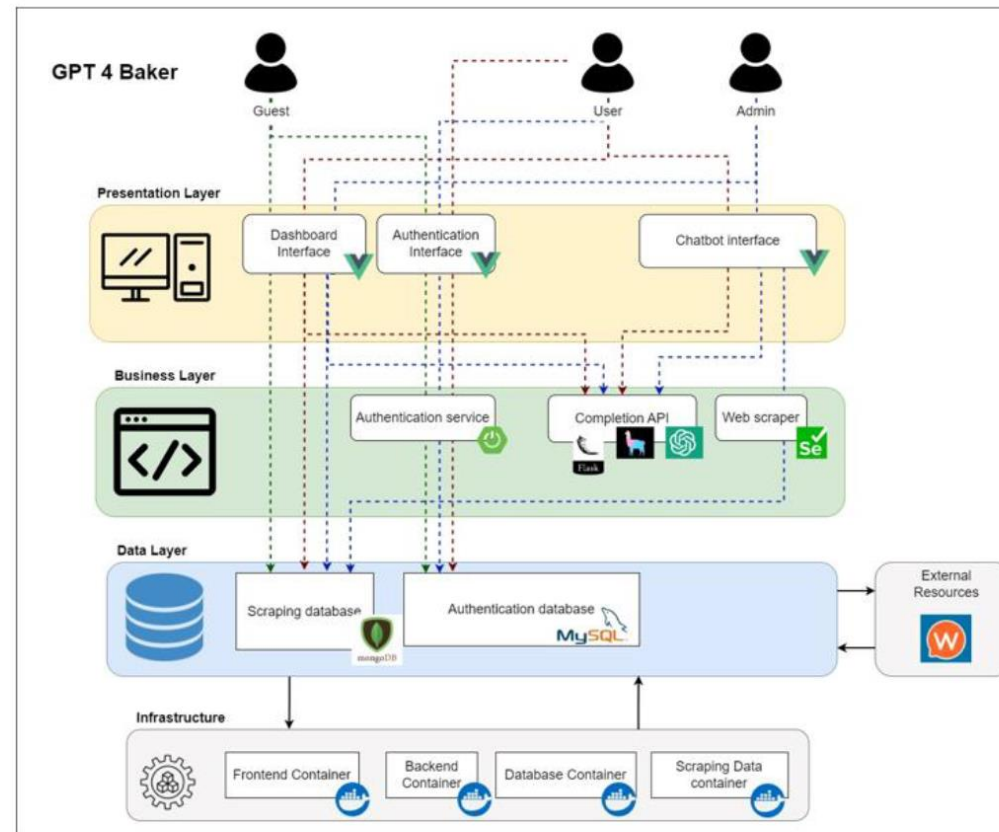


Figure 1: GPT 4 Baker System Architecture

Class activity (if we have time)

- <https://www.kaggle.com/datasets/uciml/sms-spam-collection-dataset>
- Answer the following question:
- How many rows/sample?
- What is the longest sample?
- How many word?
- What is the average word length?

References

- Lane, H., Hapke, H., & Howard, C. (2019). *Natural Language Processing in Action: Understanding, analyzing, and generating text with Python* (1st edition). Manning.
- Turban, E., Delen, D., & Sharda, R. (n.d.). *Business Intelligence, Analytics, and Data Science: A Managerial Perspective*.
- Turcan, E., & McKeown, K. (2019). Dreddit: A Reddit Dataset for Stress Analysis in Social Media. *ArXiv:1911.00133 [Cs]*. <http://arxiv.org/abs/1911.00133>
- **Practical Natural Language Processing: A Comprehensive Guide to Building Real-World NLP Systems 1st Edition**