# CS 481

## *Artificial Intelligence Language Understanding*

**January 17, 2023**
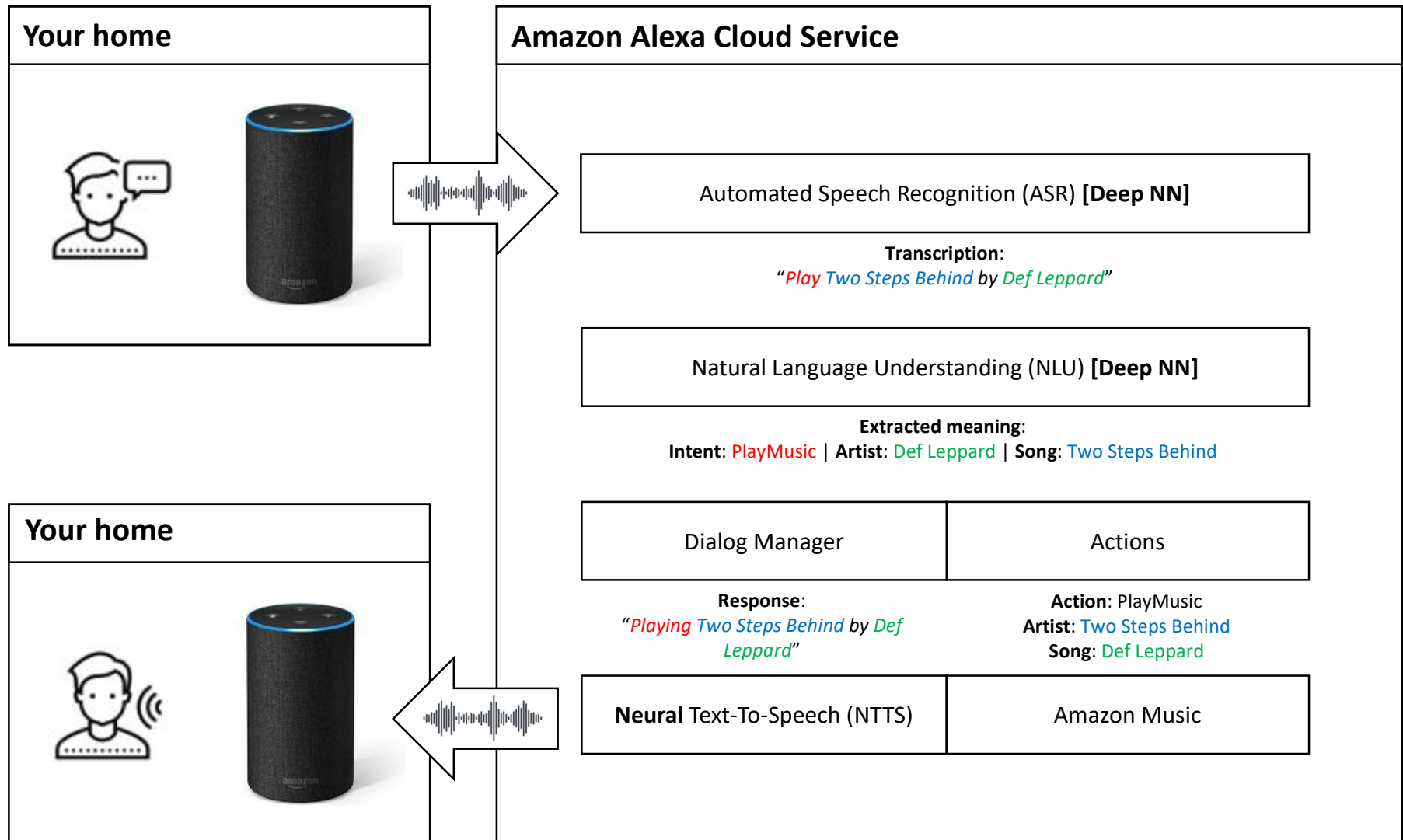
# Announcements / Reminders

- **Please follow the Week 01 To Do List instructions**

- **Quiz #01 due on Sunday 01/22/23 at 11:59 PM CST**

- **Exam dates:**
  - **Midterm:**      **03/02/2023 during Thursday lecture time**
  - **Final:**      **04/27/2023 during Thursday lecture time**

# Plan for Today

- **Introduction to NLP - continued**

- **Language basics - continued**

- **Text pre-processing**

- **Regular Expressions (RegEx)**

  - **Introduction**

  - **RegEx for basic text pre-processing**

- **Python libraries / packages for NLP**

- **Text corpora**

# Voice Assistant: Alexa

**Your home**

**Amazon Alexa Cloud Service**

Automated Speech Recognition (ASR) **[Deep NN]**

**Transcription**:
"*Play Two Steps Behind by Def Leppard*"

Natural Language Understanding (NLU) **[Deep NN]**

**Extracted meaning**:
**Intent**: PlayMusic | **Artist**: Def Leppard | **Song**: Two Steps Behind

**Your home**

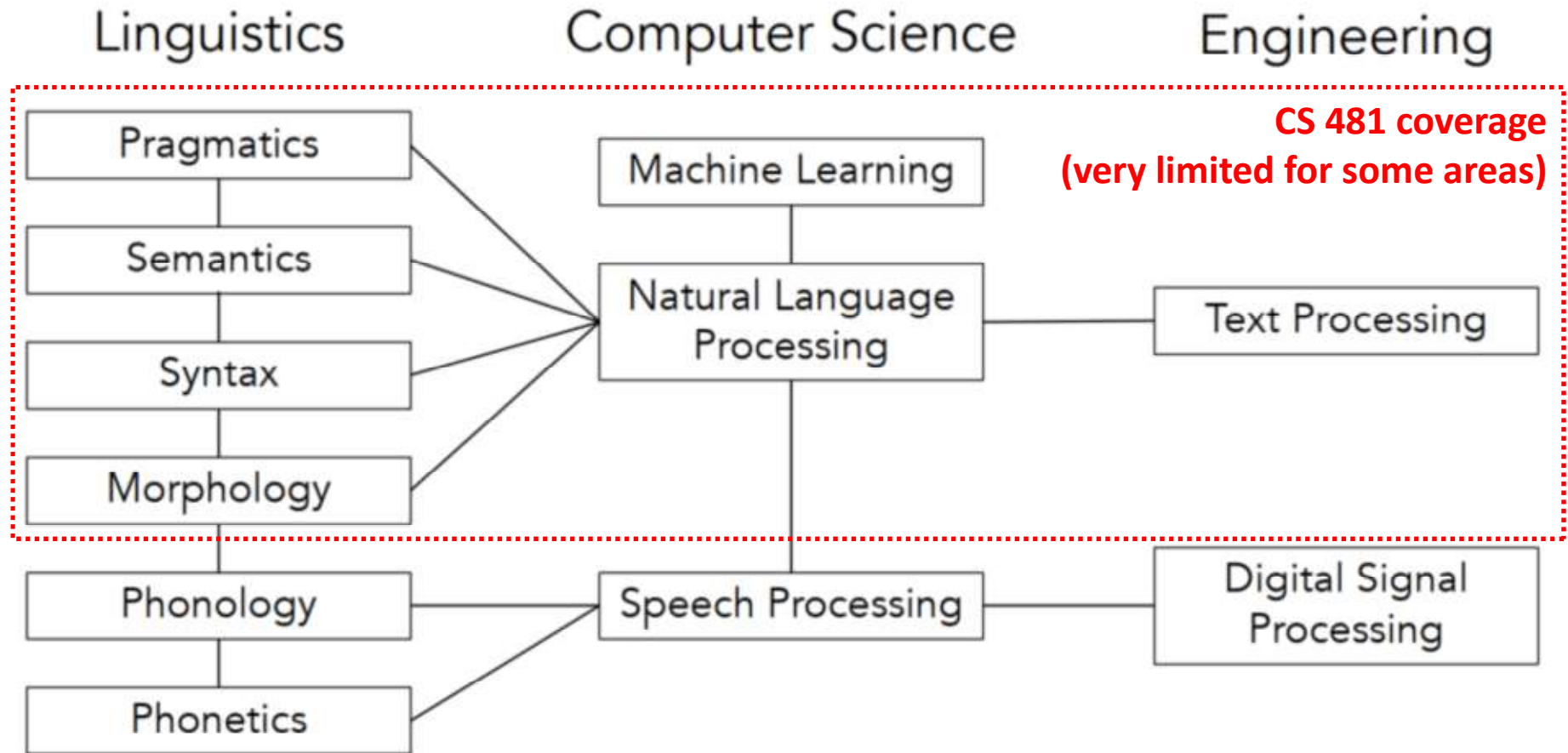| Dialog Manager | Actions |
|---|---|
| **Response**:<br>"*Playing Two Steps Behind by Def Leppard*" | **Action**: PlayMusic<br>**Artist**: Two Steps Behind<br>**Song**: Def Leppard |
| **Neural** Text-To-Speech (NTTS) | Amazon Music |

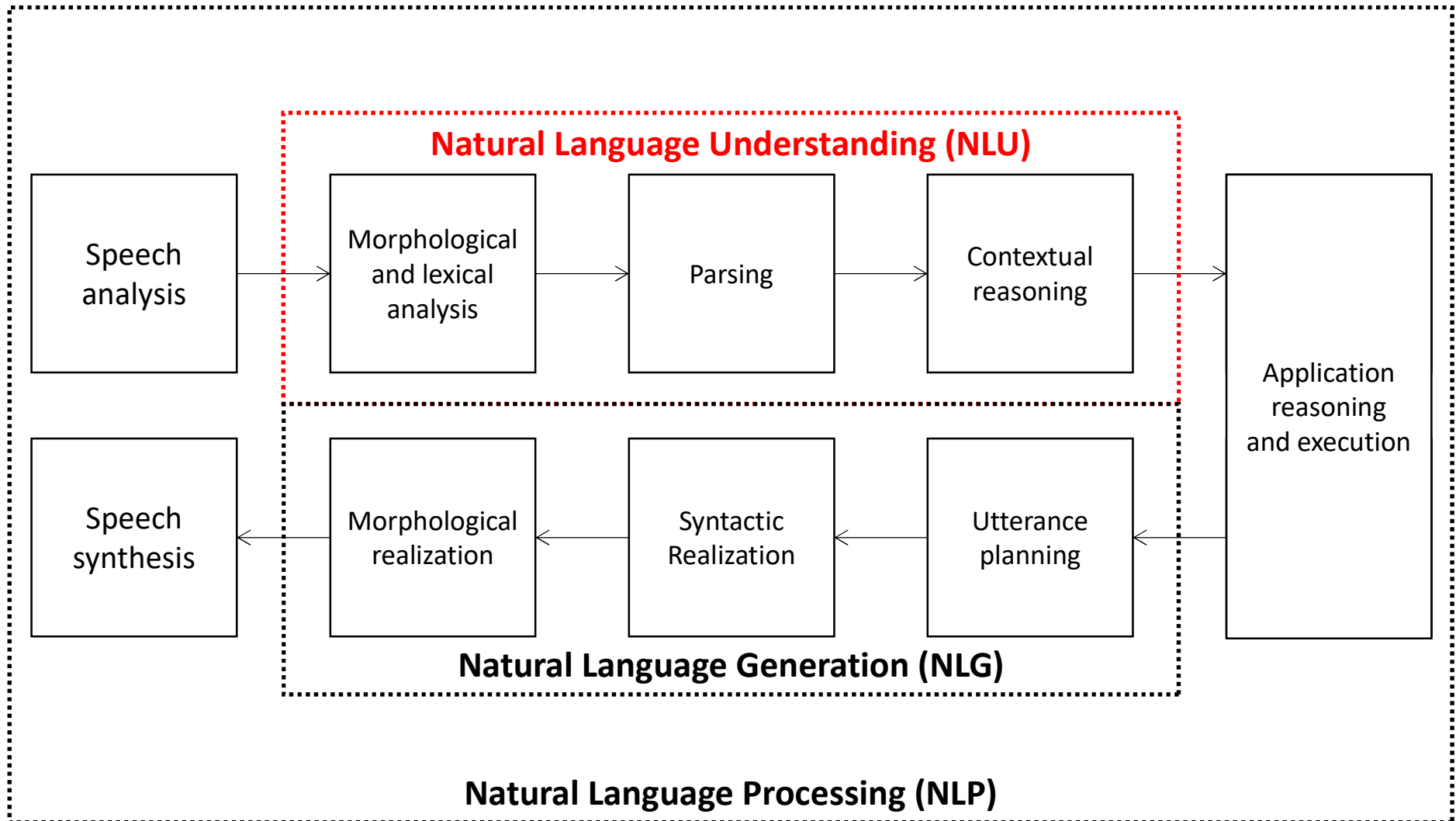# Generative Pre-trained Transformer 3

**What is it?**

Generative Pre-trained Transformer 3 (GPT-3) is an autoregressive **language model that uses deep learning to produce human-like text**. It is the third-generation language prediction model in the GPT-n series (and the successor to GPT-2) created by OpenAI, a San Francisco-based artificial intelligence research laboratory. GPT-3's full version has a capacity of **175 billion machine learning parameters**.
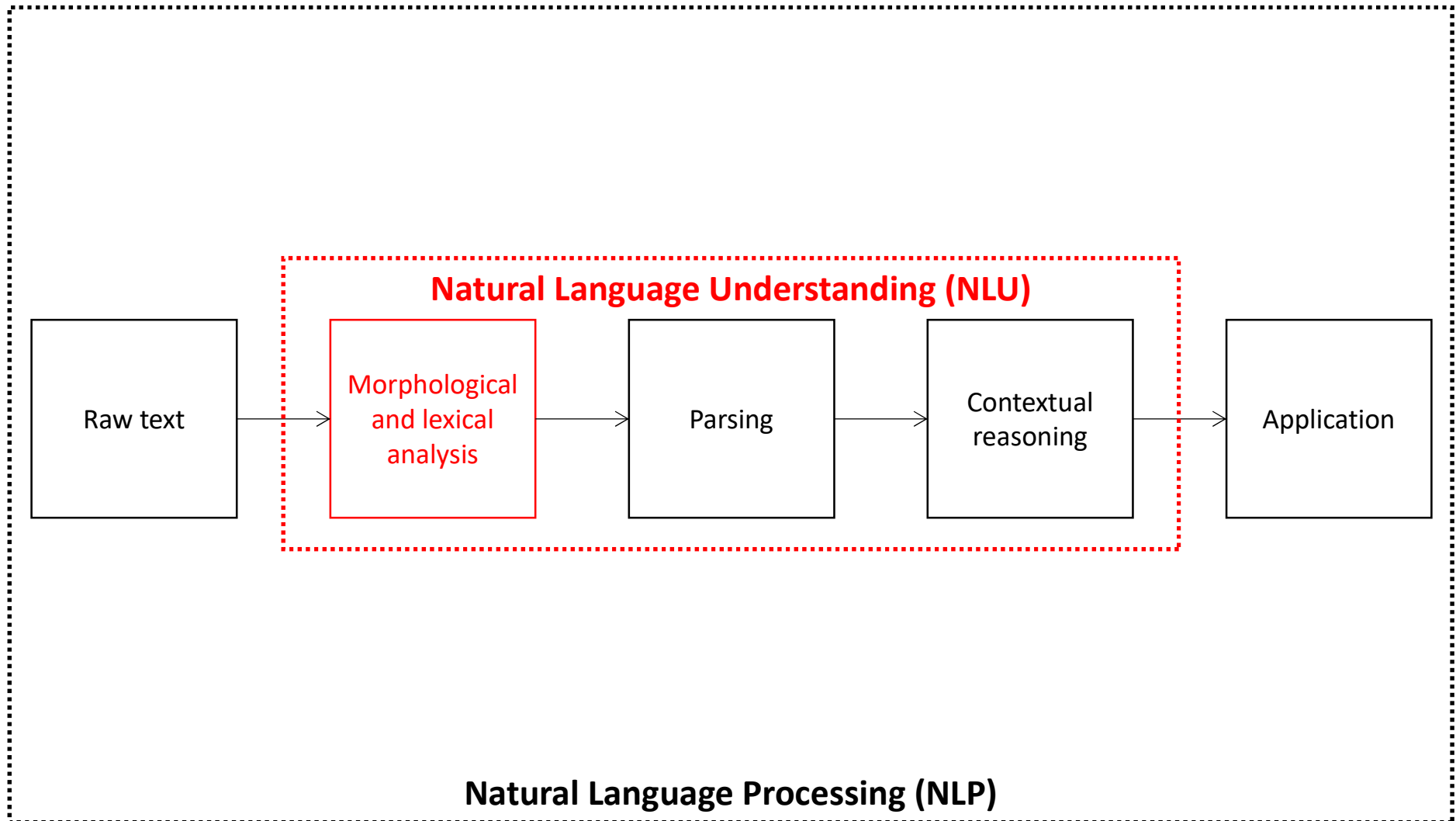
# NLP vs. Adjacent Fields

# Basic NLP Spoken Language Pipeline

**Natural Language Understanding (NLU)**

Speech analysis → Morphological and lexical analysis → Parsing → Contextual reasoning → Application reasoning and execution

Speech synthesis ← Morphological realization ← Syntactic Realization ← Utterance planning ← Application reasoning and execution

**Natural Language Generation (NLG)**

**Natural Language Processing (NLP)**

# Basic NLP Text Processing Pipeline

**Natural Language Understanding (NLU)**

Raw text → Morphological and lexical analysis → Parsing → Contextual reasoning → Application

**Natural Language Processing (NLP)**

# Common Lexical Categories

| Lexical category | Definition* | Example |
|---|---|---|
| Adjective | A word or phrase naming an attribute, added to or grammatically related to a noun to modify or describe it | The *quick red* fox jumped over the *lazy brown* dogs. |
| Adverb | A word or phrase that modifies or qualifies an adjective, verb, or other adverb, or a word group, expressing a relation of place, time, circumstance, manner, cause, degree, etc. | The dogs *lazily* ran down the field after the fox. |
| Conjunction | A word that joins two words, phrases, or clauses | The quick red fox *and* the silver coyote jumped over the lazy brown dogs. |
| Determiner | A modifying word that determines the kind of reference a noun or noun group has, for example *a*, *the*, *very* | *The* quick red fox jumped over *the* lazy brown dogs. |
| Noun | A word used to identify any of the class of people, places, or things, or to name a particular one of these. | The quick red *fox* jumped over the lazy brown *dogs*. |
| Preposition | A word governing, and usually preceding, a noun or pronoun and expressing a relation to another word or element in the clause | The quick red fox jumped *over* the lazy brown dogs. |
| Verb | A word used to describe an action, state, or occurence, and forming the main part of the predicate of a sentence, such as *hear*, *become*, and *happen* | The quick red fox *jumped* over the lazy brown dogs. |

* all definitions are taken from the New Oxford American Dictionary, 2nd Edition

# Lexical Categories: Subcategories

**Nouns:**

- **common nouns represent classes of entities:**
  - *town*, *ocean*, *person*

- **proper nouns represent unique entities:**
  - *London*, *John*, *Eiffel Tower*

- **pronouns are nouns representing other entities (usually mentioned previously):**
  - *he*, *she*, *it*

# Morphology

- **Morphology is a study of the internal structure of words**

- **Words consist of:**
  - **lexeme (root form)**
  - **affixes (suffix, prefix)**

- **Morphology has two categories:**
  - **inflectional - does not create new lexemes (happier)**
  - **derivational - creates new lexemes (unhappy)**

- **Inlectional morphemes carry grammatical meaning (plural -s), but they do not change the meaning of the word**

| Suffix | Example | Verb |
|--------|---------|------|
| -ation | nomination | nominate |
| -ee | appointee | appoint |
| -ure | closure | close |
| -al | refusal | refuse |
| -er | runner | run |

| Suffix | Example | Adjective |
|--------|---------|-----------|
| -dom | freedom | free |
| -hood | likelihood | likely |
| -ist | realist | real |
| -th | warmth | warm |
| -ness | happiness | happy |

| Suffix | Example | Marked form |
|--------|---------|-------------|
| N/A | look | base form |
| -ing | looking | gerund form |
| -s | looks | third person singular |
| -ed | loooked | past tense form |
| -en | taken | past participle |

# Phrases

- **Phrases consist of multiple words**

- **Phrases are rooted by at least one word of a particular type, but can also consist of words and phrases of other types**

- **Phrases can be combined to form clauses that are the minimal units to construct a sentence**

# Phrases and Clauses
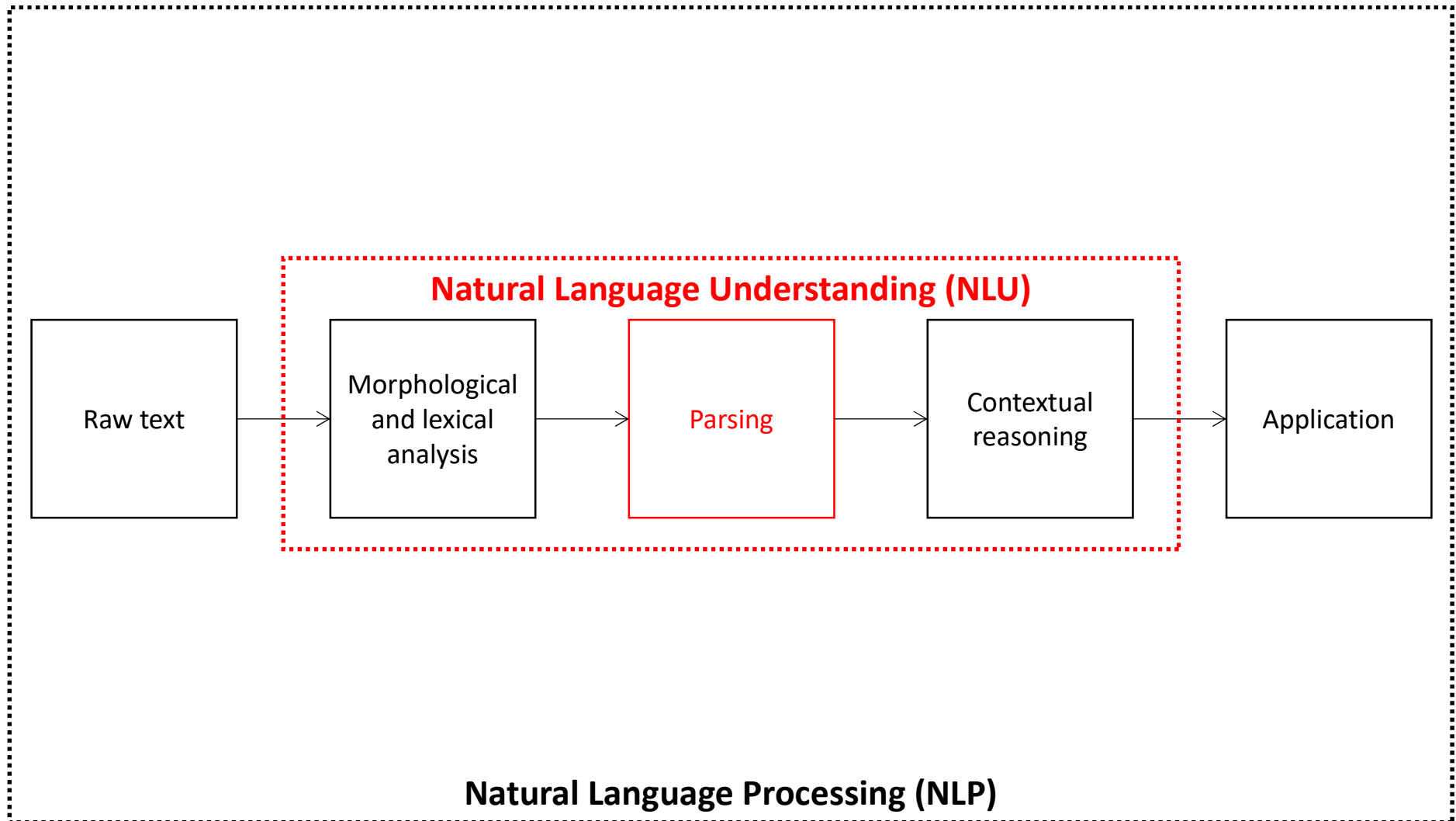
**Every sentence is constructed from phrases and/or clauses.**

- **A phrase is a group of words, but it doesn't contain a subject and a verb.**

  - **Example:** *The big clock*

- **A clause is a group of words that contains a subject and a verb.**

  - **Example:** *The big clock chimed*

**Phrase → Clause → Sentence**

# Common Phrasal Categories

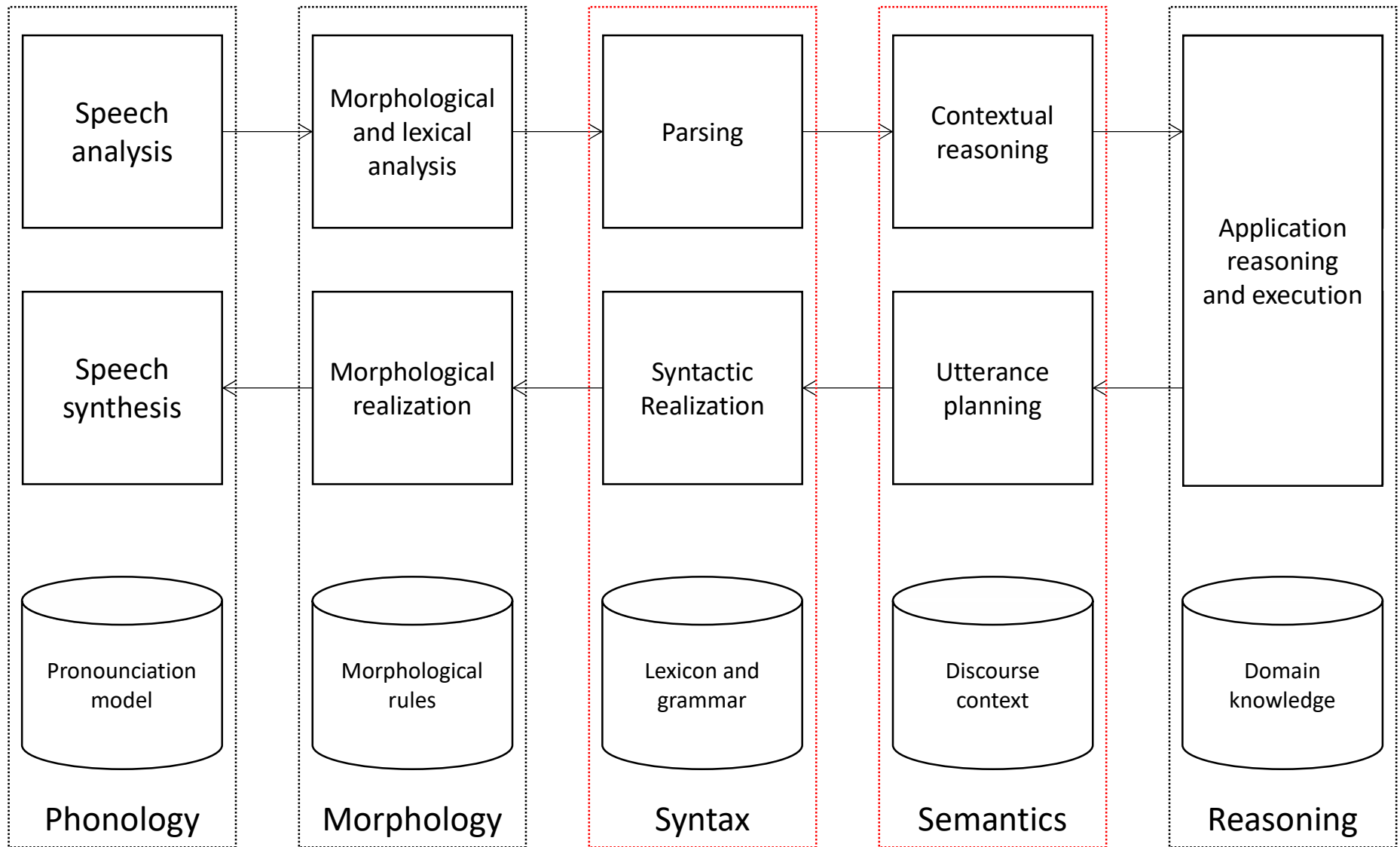| Type | Example | Comments |
|------|---------|----------|
| Adjective | The *unusually red* fox jumped over the *exceptionally lazy* dogs. | The adverbs *unusually* and *exceptionally* modify the adjectives *red* and *lazy*, respectively, to create adjectival phrases. |
| Adverb | The dogs *almost always* ran down the field after the fox. | The adverb *almost* modifies the adverb *always* to create adverbial phrase. |
| Conjunction | The quick red fox *as well as* the silver coyote jumped over the lazy brown dogs. | Though this is somewhat of an exceptional case, you can see that the phrase *as well as* performs the same function as a conjunction such as *and*. |
| Noun | *The quick red fox jumped* over *the lazy brown dogs*. | The noun *fox* and its modifiers *the*, *quick*, and *red* create a noun phrase, as does the noun *dogs* and its modifiers *the*, *lazy*, and *brown*. |
| Preposition | The quick red fox jumped *over the lazy brown dogs*. | The preposition *over* and the noun phrase *the lazy brown dogs* form a prepositional phrase that modifies the verb *jumped*. |
| Verb | The quick red fox *jumped over the lazy brown dogs*. | The verb *jumped* and its modifier the prepositional phrase *over the lazy brown dogs* form a verb phrase. |

# Basic NLP Text Processing Pipeline

**Natural Language Understanding (NLU)**

Raw text → Morphological and lexical analysis → Parsing → Contextual reasoning → Application

**Natural Language Processing (NLP)**

# Parsing

The task of determining the parts of speech, phrases, clauses, and their relationship to one another is called <span style="color:red">parsing</span>.

# Basic NLP Spoken Language Pipeline



| Phonology | Morphology | Syntax | Semantics | Reasoning |
|-----------|-----------|--------|-----------|-----------|
| Speech analysis | Morphological and lexical analysis | Parsing | Contextual reasoning | Application reasoning and execution |
| Speech synthesis | Morphological realization | Syntactic Realization | Utterance planning | |
| Pronounciation model | Morphological rules | Lexicon and grammar | Discourse context | Domain knowledge |

# Knowledge Levels / Forms for NLP

| Level | Description |
|---|---|
| **Phonetic and phonological knowledge** | Concerned with how the words are related to sounds that realize them. Such knowledge is crucial for speech-based systems. |
| **Morphological knowledge** | Concerned with how words are constructed from the basic meaning units called **morphemes**. |
| **Syntactic knowledge** | Concerned with how words can be put together to form correct sentences and determines what structural role each word plays in the sentence and what phrases are subparts of what other phrases. |
| **Semantic knowledge** | Concerned with what the words mean and how these meanings combine in sentences to form sentence meanings. This is the study of context-independent meaning - the meaning a sentence has regardless of the context in which it is used. |
| **Pragmatic knowledge** | Concerned with how sentences are used in different situations and how use affects the interpretation of the sentence. |
| **Discourse knowledge** | Concerned with how the immediately preceding sentences affect the interpretation of the next sentence. This information is especially important for interpreting pronouns and for interpreting the temporal aspects of the information. |
| **World knowledge** | Includes the general knowledge about the structure of the world that language users must have in order to, for example, maintain a conversation. It includes what each language user must know about the other user's beliefs and goals. |

# (English) Syntax

The structure of words and phrases within a sentence:

- Different formalisms, coming from the American (phrase structure) and European (dependency grammar) structuralist traditions

Applications:

- Part-of-speech tagging
- Entity extraction
- Syntactic parsing (Context-Free Grammar)
- Syntactic parsing (dependencies)

Examples:

# Semantics

The representation of meaning in language:

- at different levels: lexical, sentential, textual

- logical formalisms: reference and truth conditions

Applications:

- Word embedding / encoding
- Lexical resources
- Semantic role labeling

Example:

$$\forall x \, [\text{bird}(x) \Rightarrow \text{fly}(x)]$$

# Pragmatics

How language is used to achieve specific intentions:

- conversational implicatures: how I interpret what you say because of what I assume you are trying to do
- speech acts

Applications:

- Speech act labeling
- Discourse structure parsing
- Dialogue systems

Examples:

"I ate **most** of your cookies"

☐

I did not eat **all** of your cookies

"**Where** does your brother live?"

☐

**I do not know where** your brother lives

# Structure / Rank Levels for NLP



"So, what do you think?"

"I disagree..."

Discourse / text level

S
NP          VP
NP          PP
V          NP          AJP
DT   N1   VB   VVG   DT   DPS   N2   PRP   AV   AJ

The lecturer's teaching all his courses with great class

lectur   er   's   teach   ing   course   s

Clause / sentence level

Group / phrase level

Word level

Morpheme level

# Syntax - Semantics - Pragmatics

- **Syntax: what is its "formal" relation structure?**

- **Semantics: what does it "mean"?**

- **Pragmatics: how is it "used"?**

**Consider a hypothetical "first" CS 481 sentence:**

| Sentence | Syntax | Semantics | Pragmatics |
|---|---|---|---|
| Language is one of the fundamental aspects of human behavior and is a crucial component of our lives. | 👍 | 👍 | 👍 |
| Green frogs have large noses. | 👍 | 👍 | 👎 |
| Green ideas have large noses. | 👍 | 👎 | 👎 |
| Large have green ideas nose. | 👎 | 👎 | 👎 |

# Knowledge - Rank Mapping: NLP Tasks

| Rank / Domain | Syntax | Semantics | Pragmatics |
|---|---|---|---|
| **Word** | Parts-of-speech Morphology | Word senses Word similarity | Sentiment analysis |
| **Group / Phrase** | Shallow parsing | Named entity recognition Semantic role labeling | Deixis Coreference |
| **Clause / Sentence** | Parsing | Information extraction Entailment | Speech act interpretation Sentiment analysis |
| **Discourse / Text** | Rhetorical discourse structure | Text categorization Story understanding | Coherence Sentiment analysis |

# Representations and Understanding

Text **understanding** involves **computing** a **representation of the meaning** of sentences and texts.

- The sentence itself <u>is not enough</u> to represent the meaning:
  - word *cook* has a verb and a noun sense,
  - word *catch* can mean a baseball move, a fish, etc.

# Representation Language Properties

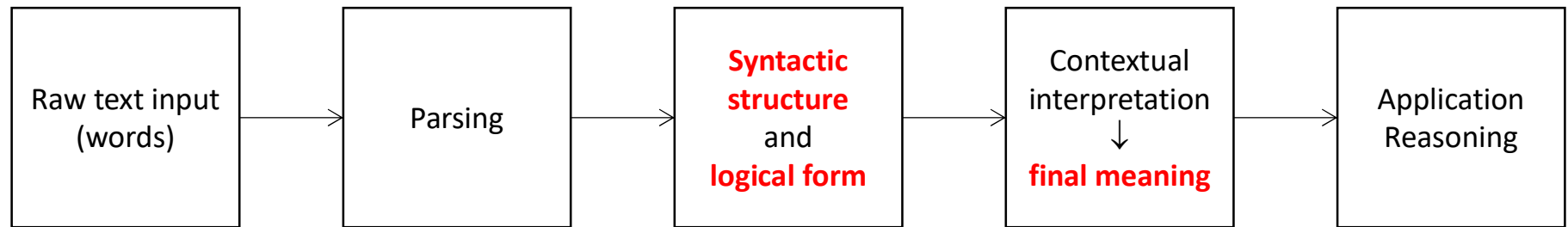Useful representation language have two properties:

- The representation **must be precise and unambiguous.** You should be able to **express every distinct reading of a sentence as a distinct formula** in the representation.

- The representation should **capture the intuitive structure of the natural language sentences** that it represents. For example:

    - sentences that are structurally similar should have similar structural representations, and

    - the meanings of two sentences that are paraphreses of each other should be closely related to each other.

# Syntax - Logical Form - Final Meaning
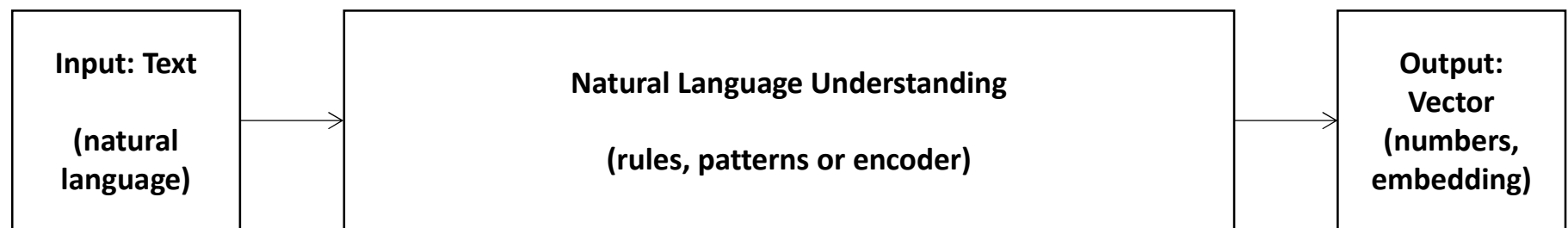
**Representation can be realized using:**

- **sentence syntactic structure: it indicates the way that individual words in a sentence**
  - **are related to each other**
  - **grouped together into phrases**
  - **modify other words**
  - **are of central importance**
- **the logical form: context-independent meaning of a sentence**
- **the final meaning: general knowledge representation meaning | context-dependent meaning**

# NLU: Flow of Information

| Raw text input (words) | → | Parsing | → | **Syntactic structure** and **logical form** | → | Contextual interpretation ↓ **final meaning** | → | Application Reasoning |

# Automated Text Processing

The task of **automatic processing of text** is to **extract a numerical representation of the meaning of that text**. This is the natural language understanding (NLU) part of NLP. The **numerical representation of the meaning of natural language** usually takes the form of a **vector called an embedding**.

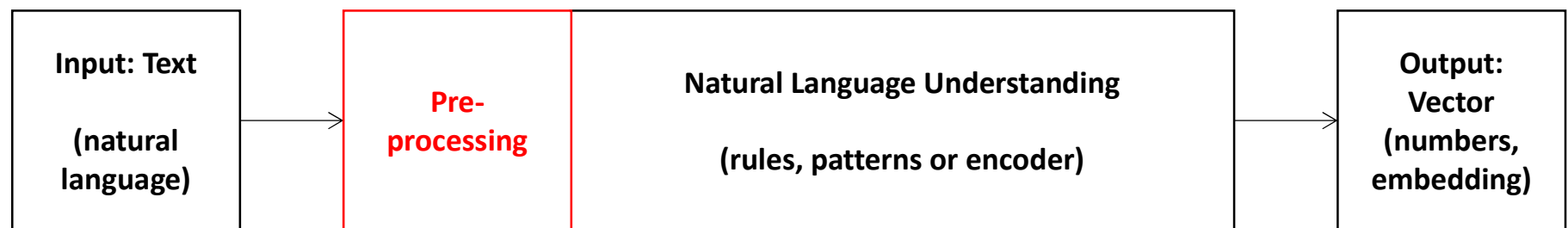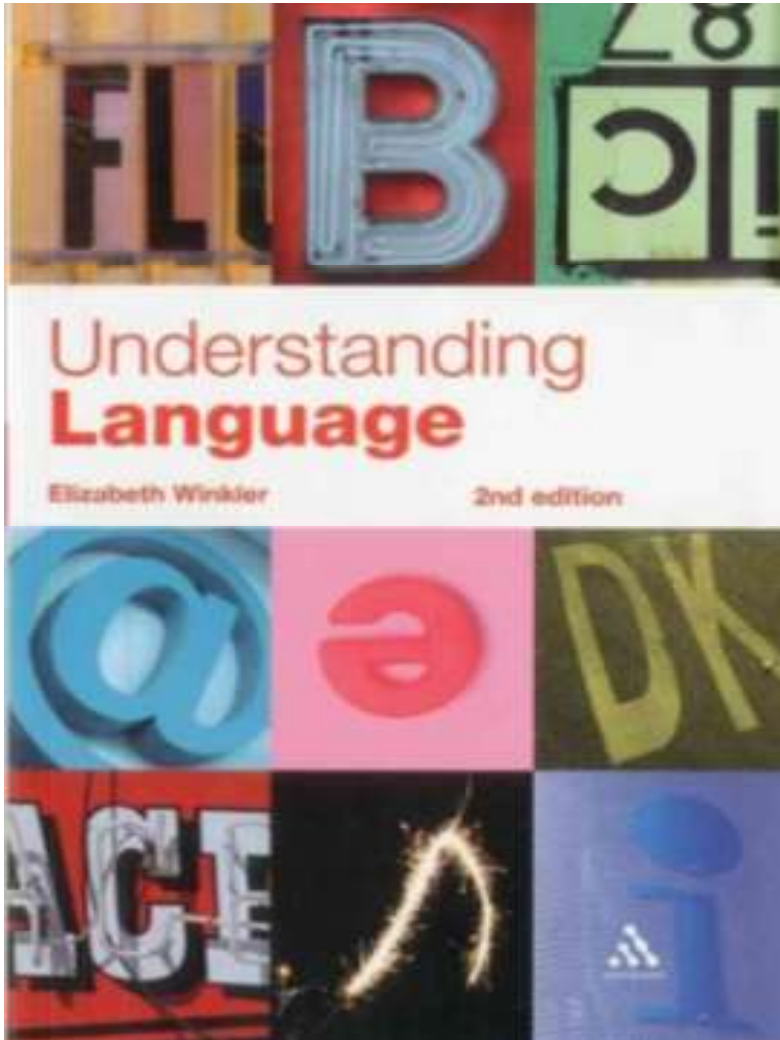| Input: Text<br><br>(natural language) | → | Natural Language Understanding<br><br>(rules, patterns or encoder) | → | Output: Vector (numbers, embedding) |
|---|---|---|---|---|

# Automated Text Processing

The task of **automatic processing of text** is to **extract a numerical representation of the meaning of that text**. This is the natural language understanding (NLU) part of NLP. The **numerical representation of the meaning of natural language** usually takes the form of a **vector called an embedding**.

| Input: Text<br><br>(natural language) | Pre-processing | Natural Language Understanding<br><br>(rules, patterns or encoder) | Output: Vector (numbers, embedding) |
|---|---|---|---|

# Optional Reading

Title:

Understanding Language (2nd edition)

Authors:

Elizabeth Winkler

ISBN:

9781441138965

Publisher:

Continuum Books

Published:

March 29, 2012

# Character Encoding Standards

Encoding standards specify how to interpret individual text characters in documents.

Well-known standards for English:

- ASCII (American Standard Code for Information Interchange)

- ISO 8859-1 Latin 1

- Unicode UTF-8 (also UTF-16 and others):

    - first 128 characters: ASCII

    - first 256 characters: ISO 8859-1 Latin 1

# Character Encoding: ASCII

| Hex | Value | Hex | Value | Hex | Value | Hex | Value | Hex | Value | Hex | Value | Hex | Value | Hex | Value |
|-----|-------|-----|-------|-----|-------|-----|-------|-----|-------|-----|-------|-----|-------|-----|-------|
| 00 | NUL | 10 | DLE | 20 | SP | 30 | 0 | 40 | @ | 50 | P | 60 | ` | 70 | p |
| 01 | SOH | 11 | DC1 | 21 | ! | 31 | 1 | 41 | A | 51 | Q | 61 | a | 71 | q |
| 02 | STX | 12 | DC2 | 22 | " | 32 | 2 | 42 | B | 52 | R | 62 | b | 72 | r |
| 03 | ETX | 13 | DC3 | 23 | # | 33 | 3 | 43 | C | 53 | S | 63 | c | 73 | s |
| 04 | EOT | 14 | DC4 | 24 | $ | 34 | 4 | 44 | D | 54 | T | 64 | d | 74 | t |
| 05 | ENQ | 15 | NAK | 25 | % | 35 | 5 | 45 | E | 55 | U | 65 | e | 75 | u |
| 06 | ACK | 16 | SYN | 26 | & | 36 | 6 | 46 | F | 56 | V | 66 | f | 76 | v |
| 07 | BEL | 17 | ETB | 27 | ' | 37 | 7 | 47 | G | 57 | W | 67 | g | 77 | w |
| 08 | BS | 18 | CAN | 28 | ( | 38 | 8 | 48 | H | 58 | X | 68 | h | 78 | x |
| 09 | HT | 19 | EM | 29 | ) | 39 | 9 | 49 | I | 59 | Y | 69 | i | 79 | y |
| 0A | LF | 1A | SUB | 2A | * | 3A | : | 4A | J | 5A | Z | 6A | j | 7A | z |
| 0B | VT | 1B | ESC | 2B | + | 3B | ; | 4B | K | 5B | [ | 6B | k | 7B | { |
| 0C | FF | 1C | FS | 2C | , | 3C | < | 4C | L | 5C | \ | 6C | l | 7C | \| |
| 0D | CR | 1D | GS | 2D | - | 3D | = | 4D | M | 5D | ] | 6D | m | 7D | } |
| 0E | SO | 1E | RS | 2E | . | 3E | > | 4E | N | 5E | ^ | 6E | n | 7E | ~ |
| 0F | SI | 1F | US | 2F | / | 3F | ? | 4F | O | 5F | _ | 6F | o | 7F | DEL |

*Source: https://www.sciencebuddies.org/science-fair-projects/references/ascii-table*

# Character Encoding: ISO 8859-1 Latin 1



*Source: https://visual-integrity.com/iso-8859/*

# Character Encoding: Unicode (Sample)



*Source: https://www.vertex42.com/ExcelTips/unicode-symbols.html*

# Textual Data Acquisition



**PDF document**

**Web page**

PDF Parser

Web scraper

Text

**Social media feed**

Feed scraper

OCR

**Image with text**

BTW:
PDF Parser:         https://py-pdf-parser.readthedocs.io/
Web scraping:       https://oxylabs.io/blog/python-web-scraping
Twitter scraper:    https://pypi.org/project/twitter-scraper/
OCR:                https://pypi.org/project/pytesseract/

some related features / tools present in NLP libraries

# Pre-processing: Cleaning Text

Some documents (email threads, comment sections, etc.) may require additional clean-up:

- **special formatting and code:**
  - special characters, HTML tags, hashtags
- **salutations, signatures, addresses, etc.**
  - example: polite phrases that are irrelevant to analysis
- **replies**
  - example: duplicate questions, etc.
- **other**

### We will work with "clean" text only.

# Basic Pre-Processing: Normalization

**Document(s) / text level:**

```
┌──────────────┐     ┌──────────────────────────┐     ┌──────────────┐
│     Text     │ ──> │        Sentence          │ ──> │   Sentences  │
│              │     │ tokenization / segmentation│     │              │
└──────────────┘     └──────────────────────────┘     └──────────────┘
```

**Tokenized sentence level:**

```
┌──────────┐        ┌──────────────────────────┐
│          │ ─────> │  Lowercasing / case folding │
│          │        └──────────────────────────┘
│          │        ┌──────────────────────────┐
│Tokenized │ ─────> │  Removal of punctuation  │
│ sentence │        │     and stop words       │
│          │        └──────────────────────────┘
│          │        ┌──────────────────────────┐
│          │ ─────> │        Stemming          │
│          │        └──────────────────────────┘
│          │        ┌──────────────────────────┐
│          │ ─────> │      Lemmatization       │
└──────────┘        └──────────────────────────┘
```

**Note**: depending on the nature of data, additional pre-processing steps may be required / important.

# Pre-processing: Normalization

**Every NLP task needs to do <span style="color:red">text normalization</span>:**

- **<span style="color:red">Segmenting</span> / tokenizing sentences in a document**

- **Segmenting / <span style="color:red">tokenizing</span> words in sentences**

- **Normalizing word formats**
  - **lowercasing / case folding**
  - **stemming**
  - **lemmatization**
  - **etc.**

# Segmentation and Tokenization

- **Text segmentation (text into sentence): breaking up text into sentences at the appearance of full stops, exclamation and question marks.**
  - **abbreviations, forms of addresses (*Mrs.*), ellipses (*...*), numbers (*.02%*) are problematic**
    - **potential solution: build a "end of sentence" classifier**
  - **tools: `.split()` method, RegEx, NLP library specific**
- **Sentence tokenization (sentence to words): breaking up sentences into tokens based on the presence of whitespaces and punctuation marks (others possible).**
  - **tools: `.split()` method, RegEx, NLP library specific**

# Segmentation with Decision Trees



Lots of blank lines after me?
- YES → E-O-S
- NO → Final punctuation is *?*, *!*, or *:* ?
  - YES → E-O-S
  - NO → Final punctuation is a period?
    - YES → Am I an abbreviation (ex: *etc.*)?
      - YES → not E-O-S
      - NO → E-O-S
    - NO → not E-O-S

E-O-S: End-Of-Sentence

# Tokenizaton: Type vs. Token

- **Type: an element of the vocabulary.**

- **Token: an <span style="color:red">instance of that type</span> in runninga type in text.**

**Vocabulary: a set of types**

**Example:** *"A good course is a course that you like"*

- **9 tokens**

- **7 types (*a* and *course* are repeated)**

- **Type/token ratio (TTR): 7/9**

# Tokenization: Problematic Cases

**NLP applications should handle problematic cases during tokenization:**

- **tokens containing periods:** *Dr., xyz.com*
- **hyphens:** *rule-based*
- **clitics (connected word abbreviations):** *couldn't, we've*
- **numerical expressions and dates:** *(123) 555-5555, August 7th, 2019*
- **emojis, hashtags, email and web addresses (URLs)**

**NLP libraries differ in this regard.**

# Tokenization: Problematic Cases

| | | |
|---|---|---|
| *Finland's capital* | → | *Finland* **or** *Finlands* **or** *Finland's*  **??** |
| *what're, I'm, isn't* | → | *What are***,** *I am***,** *is not* |
| *Hewlett-Packard* | → | *Hewlett Packard* **??** |
| *state-of-the-art* | → | *state of the art* **??** |
| *Lowercase* | → | *lower-case* **or** *lowercase* **or** *lower case* **??** |
| *San Francisco* | → | **one token or two ??** |
| *m.p.h.***,** *PhD.* | → | **??** |

**Other languages (German example):**

- *Lebensversicherungsgesellschaftsangestellter* **is a noun compound for** *"life insurance company employee"*

# Pre-processing: Lowercasing

Some applications (eg. Information Retrieval, search) reduce all letters to lower case:

- users tend to use lower case

- possible exception: upper case in mid-sentence?
    - General Motors
    - Fed vs. fed

For sentiment analysis, topic modeling:

- preserving case is important (US vs. us)

# Pre-processing: Stemming

**Stemming** refers to the process of **removing suffixes and reducing the word to some base** form such that all **different variants of that word can be represented by the same base form** (*car* **and** *cars* **are reduced to** *car*).

- **use a set of rules to accomplish stemming**
  - **if the word ends in "*-es*", remove "*-es*"**
- **final base form may NOT be linguistically correct**
  - *airliner* → *airlin*
- **commonly used by search engines and in text classification**

# Stemming: Before and After

**Before:** | **After:**

*For example compressed and compression are both accepted as equivalent to compress.*

*For exampl compress and compress ar both accept as equival to compress.*

# Pre-processing: Lemmatization

**Lemmatization is a process of mapping all the different forms of a word to its base word, or lemma. In other words: reduce inflections or variant forms to base form**

- **has to find correct dictionary headword form**

# Stemming vs. Lemmatization

| Stemming | Lemmatization |
|---|---|
| | |
| adjust**able** → adjust | was → (to) be |
| | better → good |
| meet**ing** → meet | meeting → meeting |
| stud**ies** → studi | studies → study |
| study**ing** → study | studying → study |

# Tokenization / Lemmatization Example

**Sentence input:**

**Chaplin wrote, directed, and composed music for most of his films.**

**Tokenization:**

| Chaplin | wrote | , | directed | , | and | composed | music | for | most | of | his | films | . |

**Chaplin wrote, directed, and composed music for most of his films.**

**Lemmatization:**

| Chaplin | write | , | direct | , | and | compose | music | for | most | of | he | film | . |

**Chaplin wrote, directed, and composed music for most of his films.**

# Pre-processing: Stop Words

**Very common words** (articles, propositions, pronouns, conjunctions,etc.) that **do not add much information** (but take up space) are called **stop words** and are frequently filtered out.

- **Examples in English:** *an*, *the*, *a*, *for*, *is*
- **Filtering based on the stop (word) list**
  - generated based on **collection frequency**
- **Tools: RegEx + stop list, NLP libraries have their own stop lists**
- **Careful: sometimes it may lead to removing important information**

# Additional Pre-processing Steps

- **Additional normalization**
  - **in addition to stemming, lemmatization:**
    - standardizing abbreviations (eg. expanding), hyphenations, digits to text (9 to nine) conversions, etc.

- **Language detection**

- **Code mixing**
  - embedding of linguistic units such as phrases, words, and morphemes of one language into an utterance of another language

- **Transliteration**
  - converting between different writing systems

# Regular Expressions (RegEx)

A **regular expression** (RegEx, regex or regexp) is a sequence of characters that **specifies a search pattern in text**.

- patterns are used by string-searching algorithms for "find" or "find and replace" operations on strings, or for input validation.

- a technique developed in theoretical computer science and formal language theory. Related to the concept of a regular language (a formal language that can be defined by a regular expression).

- very efficient for pre-processing tasks

- there are two key RegEx libraries:
  - re (built-in)
  - and regex (external: https://pypi.org/project/regex/)

# Regular Expressions: Disjunction

- **Disjunction**: characters within square brackets []

| Pattern | Matches |
|---|---|
| [wW]oodchuck | Woodchuck OR woodchuck |
| [1234567890] | Any digit (1 or 2 or 3 or 4 or 5 or 6 or 7 or 8 or 9 or 0) |

- **Disjunction**: a **range** of characters, square brackets [] and dash -

| Pattern | Matches | Example |
|---|---|---|
| [A-Z] | An upper case letter | Drenched Blossoms |
| [a-z] | A lower case letter | my beans were impatient |
| [0-9] | A single digit | Chapter 1: Down the Rabbit Hole |

# Regular Expressions: Disjunction

- **Disjunction**: a pipe | also means OR

| Pattern | Matches |
|---|---|
| woodchuck|groundhog | woodchuck **OR** groundhog |
| yours|mine | yours **OR** mine |
| a|b|c | same as [abc] or [a-c] |

# Regular Expressions: Negation

- **Negation**: a caret ^ means negation (has to be first within [])

| Pattern | Matches | Example |
|---------|---------|---------|
| [^A-Z] | NOT an upper case letter | O<u>y</u>fn pripetchik |
| [^Ss] | Neither 'S' nor 's' | <u>I</u> have no reason for it |

# Regular Expressions: Optionality

- **Optionality**: the question mark ? marks optionality of previous

| Pattern | Matches | Example |
|---|---|---|
| `woodchucks?` | woodchuck OR woodchucks | nice <u>woodchuck</u>! |
| `colou?r` | color OR colour | beautiful <u>colour</u> |

# Regular Expressions: . wildcard

- **. wildcard**: period . represents ANY character

| Pattern | Matches | Example |
|---------|---------|---------|
| beg.n | any character between *beg* and *n* | begin<br>beg'n<br>begun |

# Regular Expressions: Other Operators

- **Some additional and useful operators:**

| Operator | Expansion | Match | Examples |
|---|---|---|---|
| **\d** | `[0-9]` | any digit | Party of <u>5</u> |
| **\D** | `[^0-9]` | any non-digit | <u>B</u>lue moon |
| **\w** | `[a-zA-Z0-9_]` | any alphanumeric / underscore | <u>D</u>aiyu |
| **\W** | `[^\w]` | a non-alphanumeric | <u>!</u>!!!!! |
| **\s** | `[ \r\t\n\f]` | whitespace (space, tab) | |
| **\S** | `[^\s]` | non-whitespace | <u>i</u>n Chicago |

# Regular Expressions: Backslash

- **Some characters need to be backlashed (operators in RegEx)**

| Pattern | Matches | Comment |
|---|---|---|
| \* | an asterisk * | K*A*P*L*A*N |
| \. | a period . | Dr. Livingston, I presume |
| \? | a question mark | What is the time? |
| \n | a newline character | |
| \t | a tab | |

# Regular Expressions: Anchors

- **Anchors**: anchor regular explations to specific places in a string

| Pattern | Matches | Comment |
|---------|---------|---------|
| `^[A-Z]` | <u>P</u>alo Alto | Start of string anchor (^)<br>**First** character has to be uppercase letter |
| `^[^A-Za-z]` | <u>1</u> "Hello" | Start of string anchor (^)<br>**First** character cannot be a letter |
| `\.$` | The end<u>.</u> | End of string anchor ($)<br>Note the \ before .<br>**Last** character has to be . |
| `.$` | The end<u>?</u><br>The end<u>!</u> | End of string anchor ($)<br>**Last** character can be anything |
| `\b` | | word boundary |
| `\B` | | word non-boundary |

# Regular Expressions: Example

Find an instance of the word 'the' within input string.

RegEx patterns:

- `the` : **will miss capitalized 'The'**

- `[tT]he` : **will match substrings 'the' and 'The' within other words (o*the*r, *the*m)**

- `[^a-zA-Z][tT]he[^a-zA-Z]` : **this will do it**

Fixed two type of errors (to increase precision and recall):

- **Type I: matching strings we shouldn't have (false positive)**

- **Type II: not matching strings we should have matched (false negative)**

# Python re Module / Library

**Python's re module / library is built-in. Documentation:**

```
https://docs.python.org/3/library/re.html
```

**Key functions / methods:**

- `match()`: **checks for matching string at the beginning**
- `search()`: **find first location of a matching string**
- `findall()`: **returns all non-overlapping matches**
- `split()`: **splits string by occurences of a pattern**
- `sub()`: **"replace"**

# Python NLP Libraries / Packages

- **Natural Language Toolkit (NLTK) [more academic]**
- **TextBlob**
- **CoreNLP**
- **Gensim**
- **spaCy [industry / production]**
- **Polyglot**
- **scikit-learn (machine learning)**
- **pyTorch (machine learning)**
- **Pattern**
- **PyNLPl**

# Natural Language Toolkit (NLTK)

"NLTK is a leading platform for building Python programs to work with human language data. It provides easy-to-use interfaces to over 50 corpora and lexical resources such as WordNet, along with a suite of text processing libraries for classification, tokenization, stemming, tagging, parsing, and semantic reasoning, wrappers for industrial-strength NLP libraries, and an active discussion forum."

Link:           https://www.nltk.org/

Anaconda:       https://anaconda.org/anaconda/nltk

Install:        https://www.nltk.org/install.html

# TextBlob

"**TextBlob is a Python (2 and 3) library for processing textual data. It provides a simple API for diving into common natural language processing (NLP) tasks such as part-of-speech tagging, noun phrase extraction, sentiment analysis, classification, translation, and more.**"

Link:          https://textblob.readthedocs.io/en/dev/

Anaconda:      https://anaconda.org/conda-forge/textblob

Install:       https://textblob.readthedocs.io/en/dev/install.html

# CoreNLP

"CoreNLP is your one stop shop for natural language processing in Java! CoreNLP enables users to derive linguistic annotations for text, including token and sentence boundaries, parts of speech, named entities, numeric and time values, dependency and constituency parses, coreference, sentiment, quote attributions, and relations. CoreNLP currently supports 8 languages: Arabic, Chinese, English, French, German, Hungarian, Italian, and Spanish."

Link:          https://stanfordnlp.github.io/CoreNLP/

Anaconda:      https://anaconda.org/auto/corenlp

# Gensim

"Gensim is a Python library for topic modelling, document indexing and similarity retrieval with large corpora. Target audience is the natural language processing (NLP) and information retrieval (IR) community."

Link:            https://github.com/RaRe-Technologies/gensim

Anaconda:        https://anaconda.org/anaconda/gensim

Install:         https://github.com/RaRe-Technologies/gensim

# spaCy

"spaCy is a free, open-source library for advanced Natural Language Processing (NLP) in Python."

Link:           https://spacy.io/

Anaconda:       https://anaconda.org/conda-forge/spacy

Install:        https://spacy.io/usage

# Polyglot

**"Polyglot is a natural language pipeline that supports massive multilingual applications."**

Link: https://polyglot.readthedocs.io/en/latest/index.html

Anaconda: https://anaconda.org/syllabs_admin/polyglot

Install: https://polyglot.readthedocs.io/en/latest/Installation.html

# scikit-learn

"Scikit-learn is a free software machine learning library for the Python programming language. It features various classification, regression and clustering algorithms including support-vector machines, random forests, gradient boosting, k-means and DBSCAN, and is designed to interoperate with the Python numerical and scientific libraries NumPy and SciPy."

Link:          https://scikit-learn.org/stable/index.html

Anaconda:      https://anaconda.org/anaconda/scikit-learn

Install:       https://scikit-learn.org/stable/install.html

# Pattern

**"Web mining module for Python, with tools for scraping, natural language processing, machine learning, network analysis and visualization."**

Link:        https://github.com/clips/pattern

Anaconda:    https://anaconda.org/conda-forge/pattern

Install:     https://github.com/clips/pattern

# PyNLPl

"PyNLPl (**Py**thon **N**atural **L**anguage **P**rocessing **l**ibrary), pronounced as 'pineapple', is a Python library for Natural Language Processing. It contains various modules useful for common, and less common, NLP tasks."

Link:        https://github.com/proycon/pynlpl

Anaconda:    N/A?

Install:     https://github.com/proycon/pynlpl

# Text Corpora

In linguistics, a **corpus** (plural **corpora**) or **text corpus** is a language resource consisting of a large and structured set of texts (nowadays usually electronically stored and processed).