# 1 CL

- **Defn**: The study of computer systems for understanding & generating NL.
- **Aims**: get computers to perform human tasks.
- **appli**: sentiment A, Dialogues, ASR.
- **Perspective** [Theory] Describing contents & format of phenomena description, including formating abstract models/frameworks to understand underlying principle governing lang structures and behaviors. [Methods] Developing tools, procedures, formalism to facilitate NLP by creating algorithms, models.

[Tasks] Applying methods to solves specific language tasks.

- **Methodologies**: [rule-based approaches] Explicitly model language by defining set of linguistic rules & structures crafted by human experts. Rely on predefined set of linguistic rules. [Statistical app] Implicitly modelbased on patterns and probabilities derived from large datasets, leveraging statistical models and ML algorithm to learn. [sum] lang is rule based yet show stat regularities. Rule based capture explicit linguistic feature, stat reveals implicit patterns. Interplay shows multidimensional nature of CL. — **RELATION to GL**: ① GL is the study of human lang as a universal and recognisable part of human behaviors and cognitive abilities ② Common foundation in exploration of lang. ③ Approaches & outlook. GL looks at all langs' universals, and principles, CL deals more specifically with technical parts of lang processing / GL covers entirety of human langs, CL focus on practical applies and the dev. GL - theory foundation. CL - real appli.

- **Main subtask**: [NLU] mapping given NL input into useful repres'tns, analysing different aspect of the input: lexicon (pos tagging), morphology (lemmatization), syntax (dependency parsing), discourse (anaphora resolution), pragmatics (sentiment detection). Appli virtual assistant, sentiment analysis [NLG] producing meaningful NL text from some abstract repres'tn: subtle to get it right, very domain-specific.

**2. Encoding**: Writing systems: a system to make permanent marks to repre spoken words so it can be rendered - limit exactly without involving the speaker. as in En. **Alphabetic**: char/letter - sound phoneme, no semantic meaning. abjads: char - all sound, alphabetic: char-consonants as in Hebrew. **Syllabic**: char/symbol - syllable, phonetic sound patterns, x semantic, abjida: families show common consonants but no vowels (Burmese), syllabary: unique symbol for each syllable without systematic organised (片仮名). **Logographic**: char/symbol - entire word/morpheme + meaning. Semantic > phonetics (汉字). **Hybrid**: combine. Alphabetic chars form syllabic chars tgt. 韩文 alphabetic+syllabic + logographic. Emoji is new also do not represent sound-meaning pairings. **2. How lang encode on computers?** Info stored in bits, char repre by bit sequences. With 8 bits (1 byte) per char, encoding systems including Asc... assign unique code to char to repre text in computers.

3. Binary based 2,6 symbol ( & 0. Big Endian leftmost - most significant. $2^4 2^3 2^2 2^1 2^0$

**5. Limitations of corpora** ① Not fully represented / balanced ② absence / infrequency of a construction 不足 indicate grammatical incorrectness, might due to underrepresentation. ③ No data in result ≠ doesn't exist, might be but not well-repre. **Application of corpora**: ① Extracting freq. of pos for words ② Estimating n-gram freq ③ Extracting errors in learner language ④ Extract sentiment marker. 6. **Steps collect data for corpora**: Define purpose → select source → ethical consider → Data collection → preprocessing → annotation → organization → Documentation.

**5. Text classification**. 1. How computers learn. The processing of quantifying specific aspects of data, combining them into feature vectors, and using learning algorithms to analyse and understand patterns into data. **Feature vectors**: numerical repre of features of data, a combination of individual features in a dataset. Train, Dev/Validation, Test set. 2. **Supervised ML**: computer models learn from labelled data. (Through repe) annotation workflow: split data → train model → test and cross-validate → evaluate. Eg. linear/logistic reg, decision tree, classification. **Unsupervised**: learning with data without predefined labels. Workflow: defined feature extraction → apply algorithm on dataset → inspect resulting structure. e.g. K-means Clustering, Hierarchy Clustering.

3. **Evaluate model**:

| | Pred Classif | Plain | |
|---|---|---|---|
| Plain | TP | FN | Recall |
| Ham | FP | TN | True Negative Rate (TNR) |
| | Precision | True Omission | |

$Precision: \frac{TP}{TP+FN}$   $Recall: \frac{TP}{TP+FN}$

$Accuracy: \frac{TN+TP}{TP+FP+FN+TN}$   $TNR: \frac{TN}{FP+TN}$

**4. Tokenization**: Breaking down a text into smallest units (不是总是 words). **Token**: an instance of a sequence of chars that're grouped tgt as a useful semantic unit for processing. **Challenges**: ① contracted/enclitic form "don't" "he's" ② Hyphenated forms "desk - indpt" ③ forms with adjacent to periods "Str." "U.S.A" ④ Slashes "helpful/fun" "http://" ⑤ special characters "!", ",", "?", "!" etc. ⑥ multi-word expression (incl. compound N). "hot dogs" "in spite of". ⑦ Named entities "New York" ⑧ Abbr. "Inc" ⑨ integrated morphological analysis to split & using a lexicon. plus. "U.S.A" "hot dog". **Sentence**: sequence of words that form a complete grammatical unit and convey distinct idea. **Sentence segmentation**: decide sentence boundaries (start/end), punctuation chars help to decide. **Challenge**: ① Different meaning of chars in different language and WS. ② lack of punctuation &/r spaces ③ low-resource language.

---

Decimal-to-binary: $9 \Rightarrow 9/2 = 4$ R0 V | $4/2 = 2$ R0 01. $2/2 = 1$ R0 od. $1/2 = 0$ R1 100|

**Hexadecimal**: 16 symbols

| | | | |
|---|---|---|---|
| 0 | 0 | 0000 | 16进制十六进制数字0 |
| 1 | 1 | 0001 | 表示用十六进制 |
| 2 | 2 | 0010 | |
| 3 | 3 | 0011 | 4. ASCII (American Standard Code for Info Interchange), 7 bits $2^7$ chars |
| 4 | 4 | 0100 | For Latin Alphabet (esp EN). Read: "h" is decimal, 62 hexa. ASCII 0x68 |
| 5 | 5 | 0101 | **Limitations**: no umlants, no accent/other diacritics, no symbols |
| 6 | 6 | 0110 | for currencies other than $, resulting in different modif across countries, |
| 7 | 7 | 0111 | risk of misidentification. 5 **Unicode**: has a single repre for every |
| 8 | 8 | 1000 | char. 8.16.32 bits over 10 billion chars. One univsal encoding for all |
| 9 | 9 | 1001 | **limitations**: take up lots of space & the ver. |
| A | 10 | 1010 | * we need encoding systems to ensure |
| B | 11 | 1011 | compatibility across systems. |
| C | 12 | 1100 | |
| D | 13 | 1101 | |
| E | 14 | 1110 | |
| F | 15 | 1111 | rest 10 字母 |

**5.3 Writer's Aids**. 1. spelling variation across regional differences, different genres. Then explain for writers.

2. **Spelling Error Types**. Non-word: string of char that doesn't exist as a word. Reasons: ① Typographical error (时间) ② spelling confusion ③ Keyboard layout. **Realword**: exist, identified by context. Types: long distance syntax ~, local syntactic ~, semantic ~, repetition ~. 3. **Opaque WS**: no direct correspondence between phonological & char repre. Types: silent chars~ (knight). Similar-sounding (slight), long/short vowels (recieve), double consonant (application). **L1-transfer**: a person's L1 influence on L2. Error Source)

4. **Basic edit operations** (Insertion, Deletion, Substitution, Transposition).
5. **3 step creates spell checker**. ① detects errors (source: word list) → generate candidate corrections (resource: rules, list of similar words). rule-based approach to suggest correction, on the list find candidates. → ranks: how well fitting the context (resource: statistical model). 7. **Minimal Edit Distance**: minimal number of operations required to transfer from non-word to the word. Use DAG (Directed Acyclic Graph) to calculate MED. Why useful? ① Reflect the similarity of contexts, ② Being robust against slight misalignment ③ can be efficiently translated into computer programs, modelling after human/typist operations. **Nodes**: current # of words. **Arcs**: Action with cost =1. $LI \to D$. $\forall S$.

4. **Text as Data**. 1. **Corpus**: structured collection of texts, collected with a specific aim in mind. Usually contains linguistic annotation and metadata. 2. **Metadata**: description of nonprimary data, e.g. creation date, info & the authorship, tags. 3. **Linguistic Annotation** is the process of adding linguistic info to text data, incl. postags, sentence/word boundaries, parse trees… 4. **TTR (Type Token Ratio)**: assess the diversity of vocab in a text = $\frac{\#type}{\#tokens}$ = $\frac{\text{number of unique words}}{\text{number total words}}$. high → high lexical variation (ie. diversity).

---

**6. ASR**. 1. **speech**: acoustic signal produced by human speakers. **Phone**: basic sound that makes up words, repre with special written symbols in IPA. **Phoneme**: smallest units of sound that can change meaning of words, the abstract class repre categories of meaningful sounds in a language. * **prosody** (feature): syllable stress, intonation, voice quality that go beyond individual sounds. Vocal track: physical. **Even**: phones are messy real-world realizations in speech, phonemes are abstract labels we used to categorize sounds better.

2. **How speech repred with computers?** Speech is repred with computers through a process of digitalisation. Speech as sound wave is captured as an analogue signal through a microphone, then through sampling and quantization, these analogue signals are converted into digital data that computers can understand. Computer plays back speech by processing the sequence and using electric current to move these speaker appropriately. **Sampling**: measuring electrical current from microphone is measured thousands of times per sec and stored as sequences of bits and bytes. **Quantization**: convert real-valued numbers into integers for computers to repret. 3. speech processing applications (what task performs, what CL components required): ① **Spoken Dialogue System**: HCI through spoken language, allow users to communicate with computers via voice. [ASR from speech to text | NLU | NLP. ② **Speaker Recognition**: identify identity from voices | MFCC. ③ **Emotion Detection**: | MFCC | Sentiment Analysis | Token joining POS.

4. **ASR**: computer transcribe speech to text. **Input**: speech signal. **output**: string of words user said. **Challenges/difficulties**: ① Ambiguity ② variety in speech (accent, gender, age &) ③ Type of speech (context, purpose) read speech, emotion, spontaneous~ ④ Environmental factors. e.g. background noise, microphone quality. ⑤ Single/Multi speaker overlapping turns. **Noisy channel Model**: analyse audio and find the most likely considering potential noises and errors (assume noisy feature)

**ASR system architecture**: Audio → Audio Preproc → Decoder → output, transform to feature vectors, lexicon, acoustic M, LM.

**Acoustic Model**: likely phonemes given audio. **Lexicon**: mapping written words to corresponding phoneme sequence. **LM**: predict prob. rank candidate. **Lattice**: a huge DAG used to represent the search space. **Decoder**: navigate through lattice, select most likely based on above. LM's purpose: candidate ranking, models likelihood of word sequence. Steps in training speech recogniser: data requirement (ideally time aligned) → feature extraction (MFCC

Mel-Freq-Cepstral-Coefficient) → Training (through lattice, build a statistical models that Decoder will use to explore and rank different path though) → testing (rely on Decoder, Beam search) → evaluation (WER: word Error) → Error analysis (confusion matrix: statistics about which phone ones/words were commonly confused).

## 7. Text Search
1. why TS challenging: Due to ambiguity & lack of specificity in queries, resulting in irrelevant e.g. Kate Smith: lack of additional context, shared name. 2. Search Task: Searching: retrieving info that's wanted (aka info retrieval) ① Question answering: find answer to Q. ② browsing: for music, films, friends, web posts etc didn't even know wanted (recommender system)
3. Information need: the info that the searcher is searching for. (A type of) intent: sth that a user wants (to do). query: a request/Q posted by users to SE to retrieve specific info. 4. General vs specialist users: info that queries: advanced features including specialized syntax could be powerful to put but opaque to general users, such as Roger. 5. Defn of Roger: Strings form a formal language describing patterns of char sequences. Literals: char which are identical to what they match (特殊字符排除) e.g. /car/ ⇒ "car" & "carnival". Wildcards 通配 "匹配任意任何字符" e.g. /c.t/ ⇒ cat, cht, c_t. Escaping 转义 /\$/ = 50$. \d = digits. Modes (/gat/i) 全/任意匹配/g 任选其中之一. Char sets and range 方括号 ] 代表 或, 数字 or 字母 "多个字符逐个匹配". e.g. /[aeiou]/匹配其中之一 [a-z] 代表范围. Quantified repetition: 指定特定 "*" 零次以上 0~n次. "+" 1~n次. "?" 0/1次. 范围数 {min,max} e.g. /a{3,7}/ 指定范围. Groups using parenthese ( )/(abc)+ 分组可嵌套可重复.
Anchors /^start// end$/. 6. Document Index: create index for docs available for searching, allows for efficient and fast retrieval of info, easier to locate specific words. esp with large volumes of texts. Different approaches: ① term by document matrix: shows which words appear in which doc. terms→rows. doc column 如 matrix→ term appear in. 0 → nowhere. ② Inverted Indices List: to avoid inefficient 0s as term-by-doc Matrix would be too sparse. Each term is associated with a list of unique doc IDs. Only store IDs that DS(contain multiple docs), faster processing & more memory efficient. 8. Evaluating search quality: User Exp (survey: happy or not), Objective user interaction logs. ① Precision: % of docs returned that're relevant (e.g. return 200, 100 rele, Precision=100/200). Recall: % of total rele docs that're returned. (e.g. rele 200, return 200). Recall = 200/400. F-measure: combine both. $F_1 = 2 \times \frac{Pre \times Re}{Pre + Re}$. Their limits:
intent: task the user is supposed to achieve with support from assistant. intent slots/templates that need to be filled. intent recognition: to classify user utterance as input. e.g. Siri 分割识别. Task-Indep DS: Design for entertaining/impressing users &/serve as digital companion.
△不同复杂度的 implements) ① Brute force: hand-written QA. e.g. (ELIZA).
② rule-based: fix rules to derive response given preceding turns.
③ corpus-trained chatbox: based on huge dataset to derive most suitable answer. ④ language generation: based on corpus data, built ML models, take previous turns as input too. output highest prob of reply (even not existed in corpus) (e.g. ChatGPT). △Evaluation: Turing test: intelligency indistinguishable from that of human. Winograd Schemas: ambiguous pron, need world knowledge. Microsoft Tay: lack of moral filter. 5. POS Defn: labels assigned to indicate functions. Different tagset & granularity across corpora.

## 10. LLM
1. NN: interconnected nodes organised in layers. Transformer: a type of NN relies on self-attention mechanisms to process data. What's new about LLM? ① Unified model for various tasks, efficient↑.
② parallelization (break down task into small & indep & simult executed) techniques allow for faster training on larger dataset. ③ attention mechanisms enhance ability to focus on relevant parts of input. Large in data size, parameter size, model size.
2. Applications: Visual Assistance (Siri, Alexa). Language Translation (Google Translate). Text generation: (GPT) Chat Bots. Code Gene (GitHub Copilot), content creation. 3. Ethical bias: not diverse data, social biases, whiteness norms, gender biases. How social biases inscribed in data? Data collection, human annotators. Why problematic? Bias output and decisions when applied LLM in real world applications.
Fine-tuning. Due to complexity of language ethical, difficult to ensure only ethical. Try to solve by Reinforcement learning from human feedback. 4. Environment Impact: has negative impact with traing consuming significant power and emitting CO2. Training LLM like GPT 如 car. Carbon intensity

Pre-doesn't account for relevant docs that were not retrieved by SE. Re-hard to know how many rele docs were missing by SE threat life as we often lack complete datasets or ground truth annotations for evaluation. F1: high & most effective. Overload, only 1st page result.
* why ranking is important: allows SE to prioritize rele. On top result page, Pre > Recall.

## 8. CALL (Computer-Assisted Language Learning)
1. L1 characteristics: ① Divergence between linguistic levels: learner mistakes in different parts of language (sound, word form, sentence structure ...) ② categories for native language not applicable: L2 learning mistakes that native won't make. ③ Difficult to determine target hypotheses: learners making multiple mistakes and show inconsistencies. ④ Often more than 1 error.
NLP in language learning: 1/analysing language for learner: NLP search relevant & appropriate examples/texts. 2/analysing learner production.
2. ITS (Intelligent Tutoring System): △Defn: a computer program help learner's learning (automatic immediate feedback, many users 同时). △Goals: 1/close gap between ITS search, Foreign language Teaching (FLT) insights, real-life classroom/address real format education needs using NLP term. 3. NLP appli in context of ITS: ① NLP for well-formed language: Need of users: accurate explanation of grammar rules, vocab, structures. NLP use: analyse interactive exercises and quizzes. ② NLP for mal-formed Needs: correction feedback, practice. NLP used: error detection, generate feedback. interactive exercise by NLP enhanced ITS. 4. Evaluation: same as text search in PRF.

## 9. Dialogue System
1. speech acts: actions performed through language. common grounds: info mutually shared by participants. Turn taking: roles of speakers & hearers. adjacency pairs: sequence of two turns with expected structure. APD in DS: facilitate effective communication (recognize intentions, coherence in conversation/negotiate conversion flow, maintain conversational structure). 2. Grice Maxims: norms guiding conversations ① Quality: say true things. ② Quantity: as much as info necessary, no more, no less. ③ Relevance: say only rele. ④ manner: easily understood. + FLT: norms guiding conversation, ensuring coherence and effectiveness, satisfying user experience.
3. Task-Specific DS: designed to help accomplish a specific concrete task.
varies based on energy sources: coal-powered data centers huge emissions. High computational demands amplify power consumption leading to environment strain.