# An overview of language representation

Shiwei Tong

# Outline

- Brief Introduction

- Quality Evaluation

- General Methods
  - Token/ID based methods
  - Morphology based methods

- Summary

# Language Representation

- Basic task of Natural Language Processing (NLP)
- To represent human language (e.g., character/word/sentence/document) in computer language/coding

# Naïve Language Representation

- **Character Encoding**
  - ACSII
  - UTF-8
  - …

- **One-hot Encoding**

- **Shortcoming**
  - Cannot express semantics
    - Thesaurus: person & people
    - Synonyms: dad & father
    - Polysemy: bank & bank
  - Large space

# Language Representation with Semantics

- Token/ID based methods
  - Context2Center: CBOW
  - Center2Context: Skipgram
  - FeatureModel: ELMO, GPT,  BERT
- Morphology based methods
  - Word Morphology: prefix, suffix, root
  - Sentence Morphology: grammar, syntax
  - Document Morphology: Paragraph

# Language Representation Evaluation

- General Semantics Task
  - Word Similarity:
    - (消费者, 顾客)
    - (man, woman)
  - Word Analogy:
    - 国王 – 男人 = 女王 –女人
    - king – man = queen - woman
- Task Specific Evaluation
  - Classification: Precision, Recall, F1, …
  - Translation: BLEU
  - …

# General Semantics Task

- Word Similarity
  - A list of pairs (word1, word2, score)
  - High similarity, high rank: $rank\_score = \cos(word1, word2)$
  - Spearman's correlation
- Word Analogy
  - A list of pairs (head1, tail1, head2, tail2)
  - Accuracy : $\text{head2} = \underset{t \in W}{\text{argmax}}\big(\cos(\text{head1} - \text{tail1} + \text{tail2}, t)\big)$
    - The above is called 3CosAdd
    - The other one expression, which is called 3CosMul, is $\text{head2} = \underset{t \in W}{\text{argmax}} \frac{\cos(t, tail1) \cdot \cos(t, head2)}{\cos(t, head1) + \epsilon}$
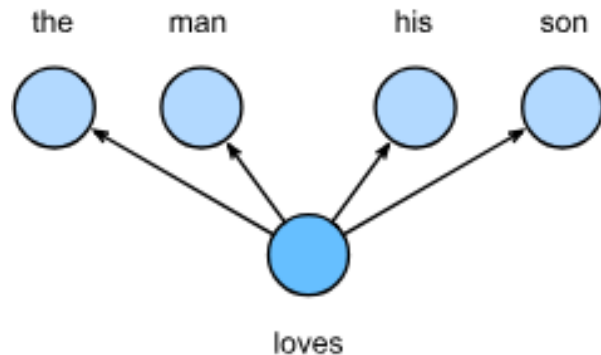
# Word Embedding: a demo case

- Semantics: similarity and analogy

- Naïve way: one-hot
  - Hard to infer the semantics

- Distributed —— Word2vec
  - A sentence: the man loves his son
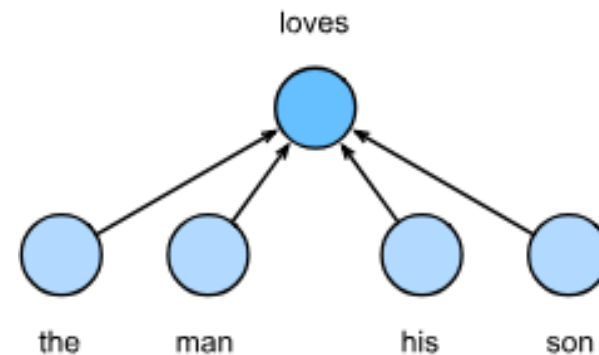  - Assumption: similar words appear in similar contexts

# Token/ID based methods

## Center2Context: Skipgram

- $P(\text{"the"}|\text{"loves"}) \cdot P(\text{"man"}|\text{"loves"}) \cdot \cdots$
- $P(w_o|w_c) = \dfrac{\exp(u_o^T v_c)}{\sum_{i \in \mathcal{V}} \exp(u_i^T v_c)}$
- $\prod_{t=1}^{T} \prod_{-m \leq j \leq m, j \neq 0} P\left(w^{t+j}|w^t\right)$
- $-\sum_{t=1}^{T} \sum_{-m \leq j \leq m, j \neq 0} \log P(w^{t+j}|w^t)$



the    man    his    son

loves

## Context2Center: CBOW

- $P(\text{loves}|\text{"the"}, \text{man"}, \text{his"}, \text{son})$
- $P\left(w_c|w_{o_1,\ldots o_{2m}}\right) = \dfrac{\exp\left(u_c^T \cdot \frac{1}{2m} \Sigma_i^{2m} v_{o_i}\right)}{\sum_{i \in \mathcal{V}} \exp(u_i^T \cdot \frac{1}{2m} \Sigma_j^{2m} v_{o_j})}$
- $P(w_c|W_o) = \dfrac{\exp(u_c^T \bar{v}_o)}{\sum_{i \in \mathcal{V}} \exp(u_i^T \cdot \bar{v}_o)}$
- $\prod_{t=1}^{T} P(w^t|w^{t-m}, \ldots, w^{t+m})$



loves

the    man    his    son

# Problem

- Bias from Frequency
  - "the" vs "microprocessor"
  - Solution
    - Subsample based on frequency
      - Dropout Probability: $P(w_i) = \max(1 - \sqrt{\frac{c}{f(w_i)}}, 0)$
      - $f(w_i)$ is the word frequency, and $c$ is a constant, usually $10^{-3}$ or $10^{-4}$
- High Time Complexity
  - O($|\mathcal{V}|$)
    - $P(w_o|w_c) = \dfrac{\exp(u_o^T v_c)}{\sum_{i \in \mathcal{V}} \exp(u_i^T v_c)}$
    - $P(w_c|W_o) = \dfrac{\exp(u_c^T \bar{v}_o)}{\sum_{i \in \mathcal{V}} \exp(u_i^T \cdot \bar{v}_o)}$

# High Time Complexity

- $O(|\mathcal{V}|)$
  - $P(w_o|w_c) = \dfrac{\exp(u_o^T v_c)}{\sum_{i \in \mathcal{V}} \exp(u_i^T v_c)}$
- Solution
  - Negative sampling
    - $P(w_o|w_c) \rightarrow \max$
    - $P(D = 1|w_o, w_c) = \sigma(u_o^T v_c) \rightarrow max.$ $\sigma$ is the sigmoid function.
    - $P(w_o|w_c) \rightarrow P(D = 1|w_o, w_c) \prod_{k=1, w_k \sim U(w)}^{K} P(D = 0|w_k, w_c)$
    - $U(w) = \dfrac{f^\alpha(w)}{\sum_i f^\alpha(w_i)}$, is the powered unigram distribution. $\alpha$ is a constant (e.g., $\frac{3}{4}$) and $K$ usually is set as 5.
  - Hieratical softmax
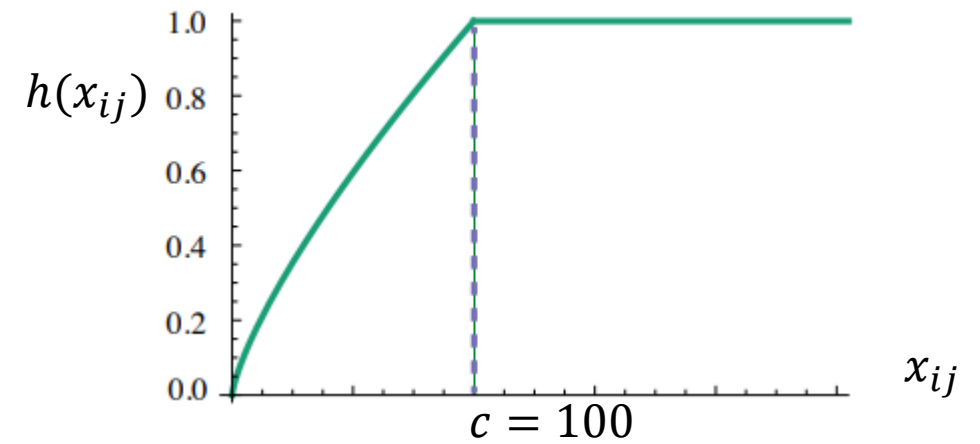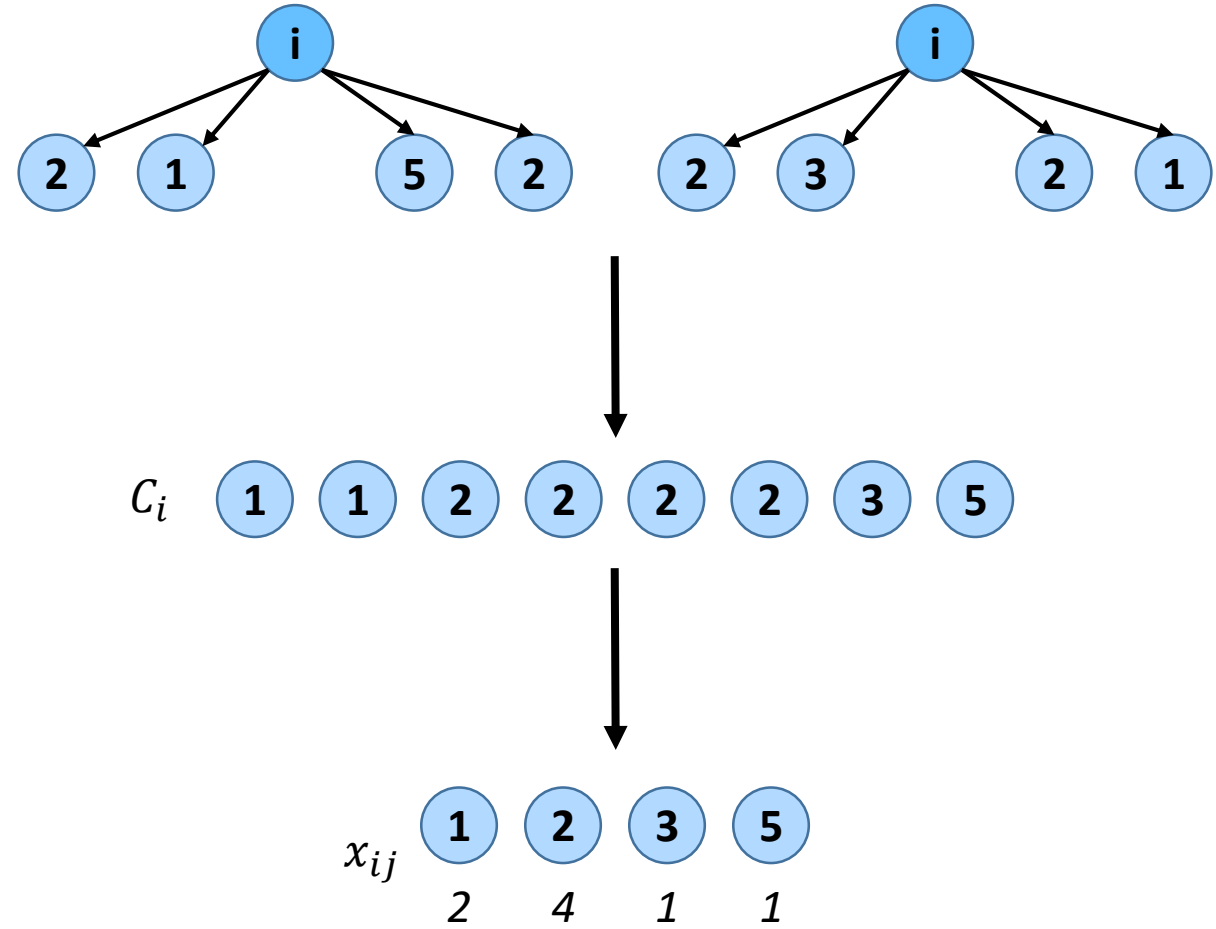    - Space for time

# Global Vectors



- Skipgram
  - Standard
    - $P(w_o|w_c) = \dfrac{\exp(u_o^T v_c)}{\sum_{i \in \mathcal{V}} \exp(u_i^T v_c)}$
    - $-\sum_{t=1}^{T} \sum_{-m \leq j \leq m, j \neq 0} \log P(w^{t+j}|w^t)$
  - Another View
    - $p_{ij} = P(w_j|w_i) = \dfrac{\exp(u_j^T v_i)}{\sum_{k \in \mathcal{V}} \exp(u_j^T v_i)}$
    - $-\sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{V}} x_{ij} \log q_{ij}$
    - $x_i = |C_i|, p_{ij} = \dfrac{x_{ij}}{x_i}$
- GloVe
  - $p'_{ij} = x_{ij}, q'_{ij} = \exp(u_j^T v_i)$
  - $\left(\log q'_{ij} - \log p'_{ij}\right)^2 = \left(u_j^T v_i + b_i + c_j - \log x_{ij}\right)^2$
  - $\sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{V}} h(x_{ij}) \left(u_j^T v_i + b_i + c_j - \log x_{ij}\right)^2$
  - $h(x_{ij}) = \begin{cases} (x_{ij}/c)^\alpha & x_{ij} < c \\ 1 & otherwise \end{cases}$

# Problem

- Polysemy
  - Bank & Bank
  - He walks on the bank
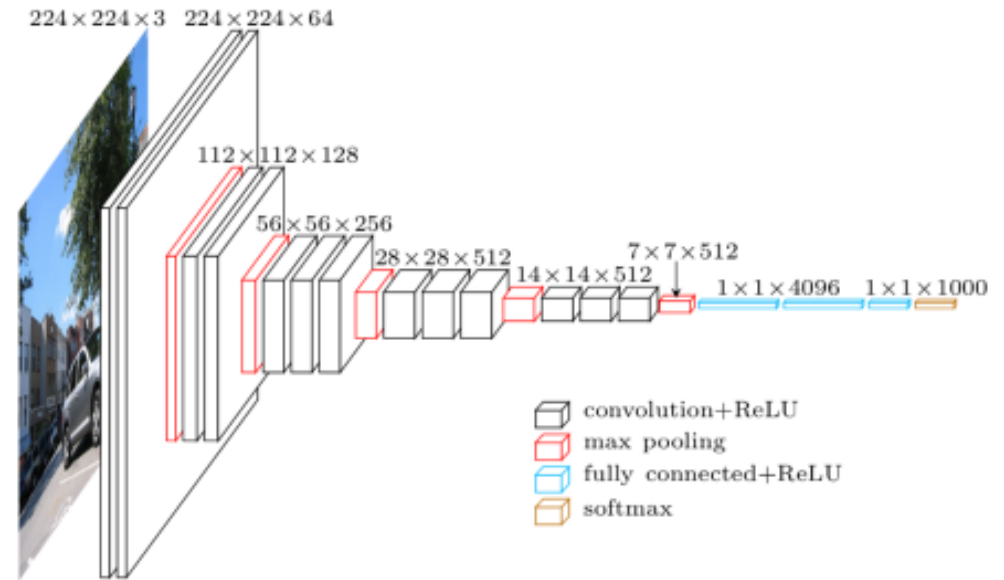  - The bank is robbed
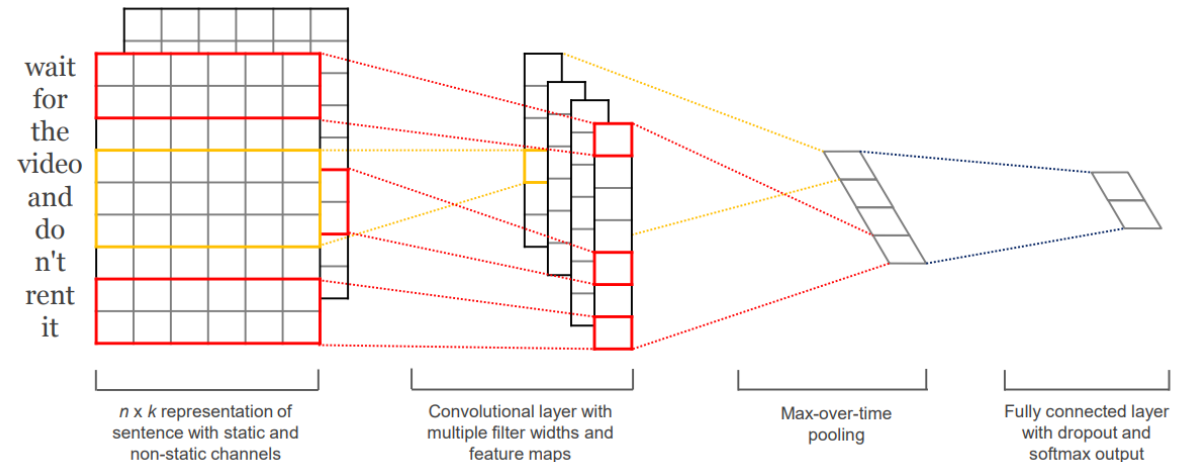
# FeatureModel

- Image
  - VGG16 & VGG19
  - ResNet
- NLP
  - TextCNN
  - ELMO
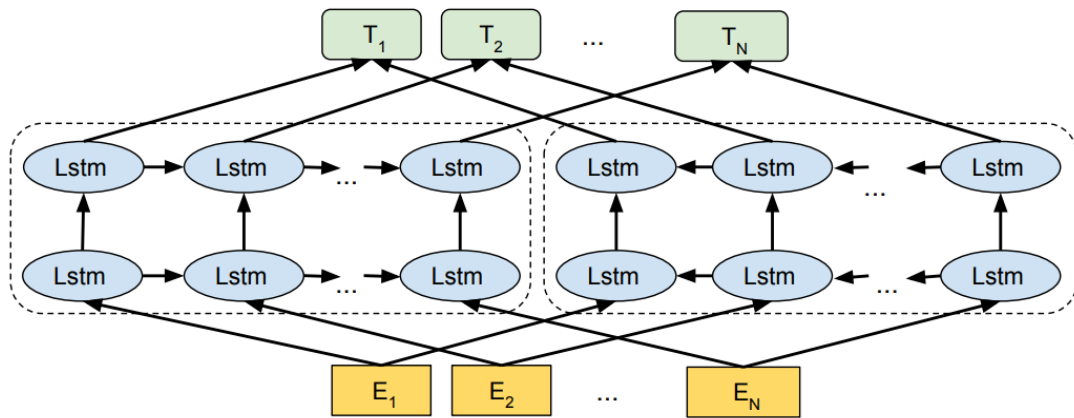  - GPT
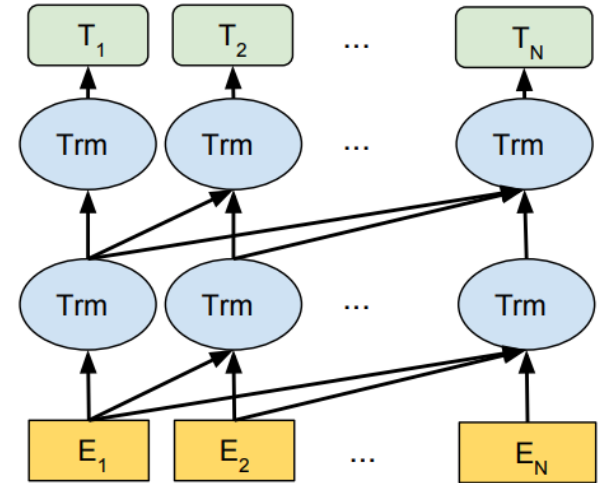  - BERT

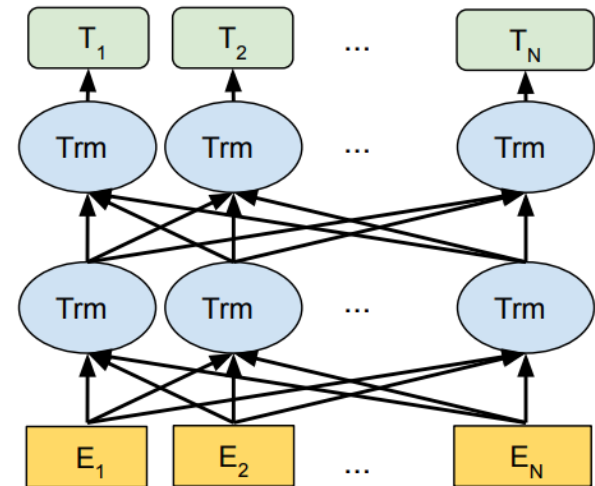

VGG16



TextCNN

# FeatureModel

- ELMO

- GPT

- BERT

# Drawback

- Ignore the abundant information in morphology
  - Word Morphology: prefix, suffix, root
  - Sentence Morphology: grammar, syntax
  - Document Morphology: Paragraph
- OOV Word

# Morphology based methods

- Alphabetic
  - Subword: prefix, suffix, root

- Logogram
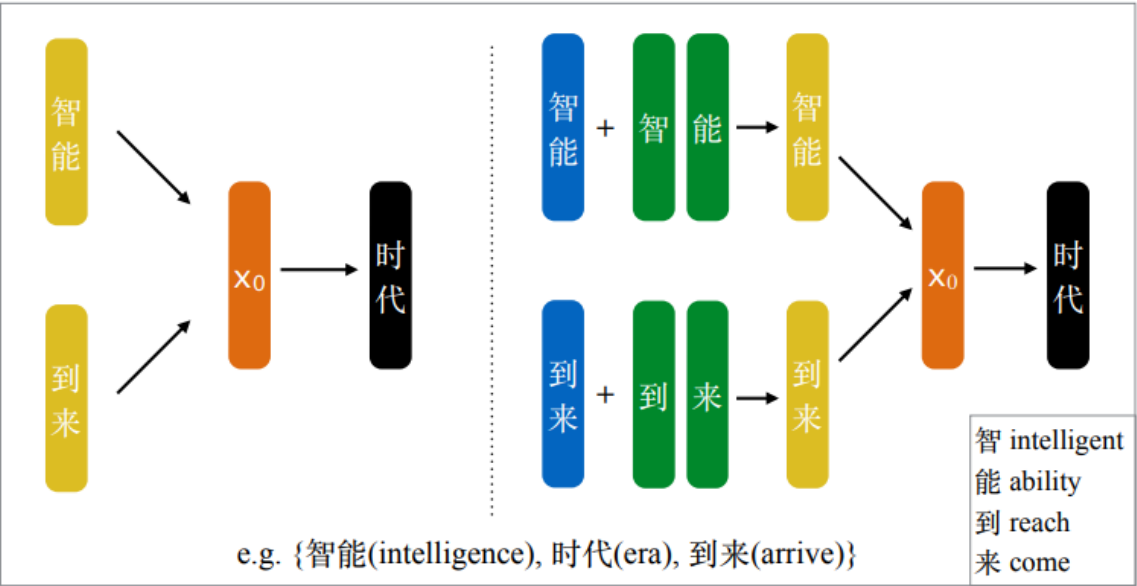  - Character
  - Stroke
  - Glyph

# Morphology Word Embedding for Alphabetic

- sisg
  - fasttext
  - "declare", "clarify" → root word "clar"
  - ("dog", "dogs"), ("interest", "interesting")
  - Subword n-gram
    - where → <where>
    - 3-gram → <wh, whe, her, ere, re>
    - $3 \leq n \leq 6$ and attach a special subword "where" or "<where>"
    - $v_w = \sum_{g \in \mathcal{G}_w} z_g$
    - $P(w_o | w_c) = \frac{\exp(u_o^T v_c)}{\sum_{i \in \mathcal{V}} \exp(u_i^T v_c)}$

# Morphology Word Embedding for Logogram
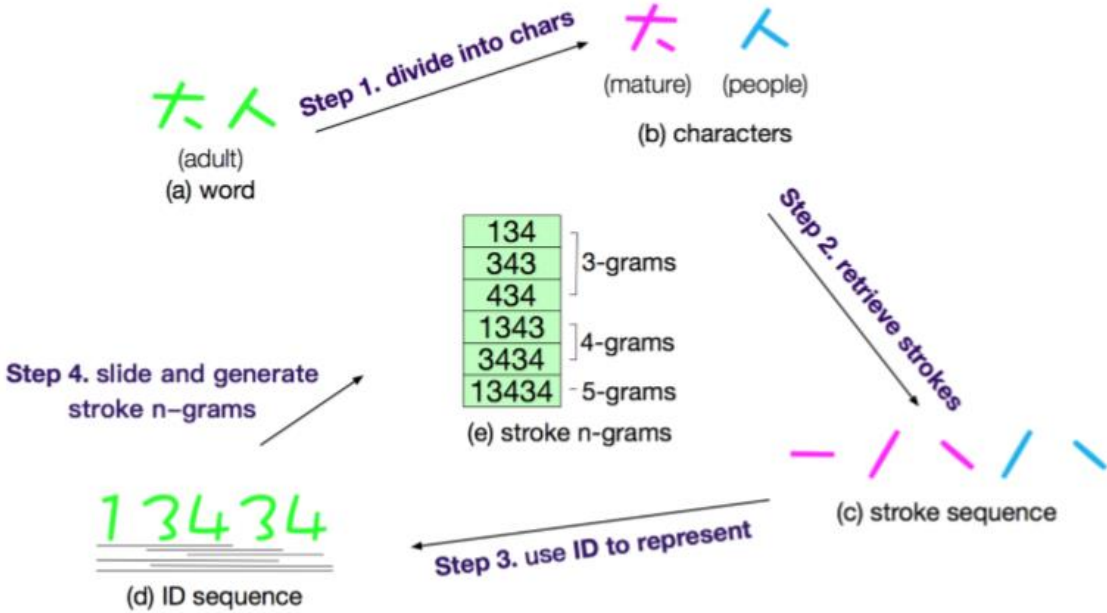
- Character

  - $u_j = w_j \oplus \frac{1}{N_j}\sum_{k=1}^{N_j} c_k$

- Stroke



智 intelligent
能 ability
到 reach
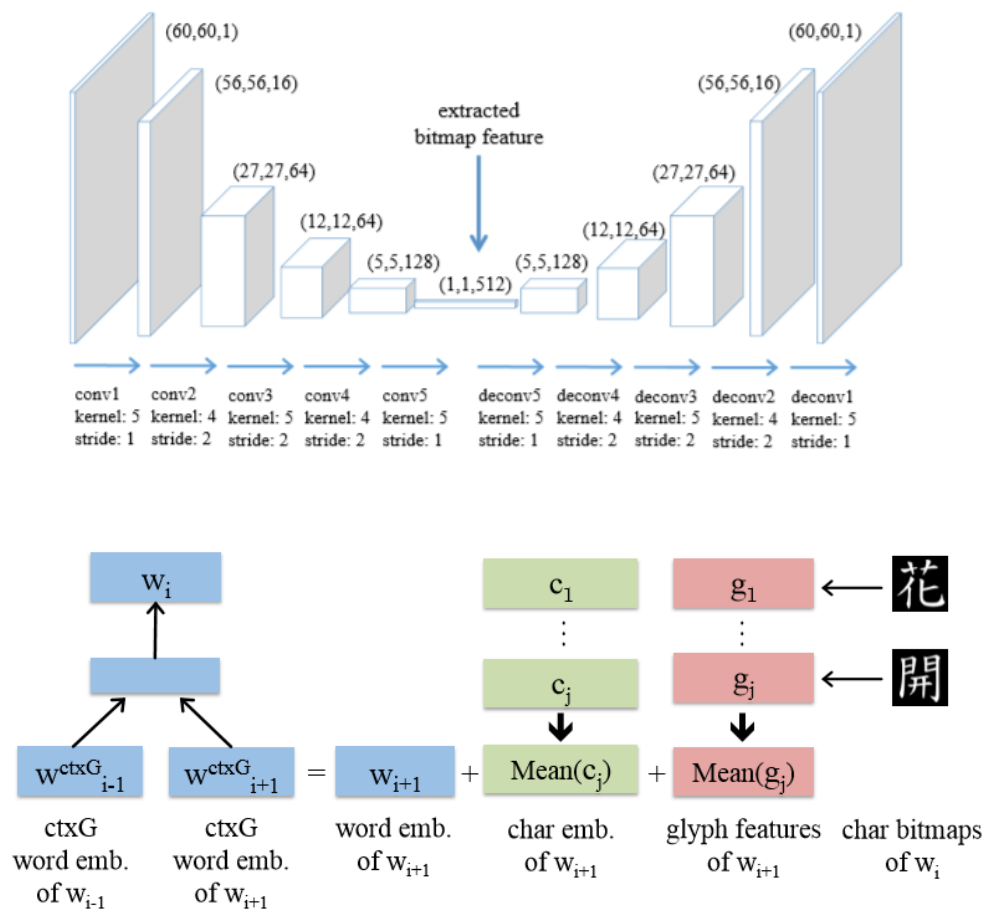来 come

e.g. {智能(intelligence), 时代(era), 到来(arrive)}
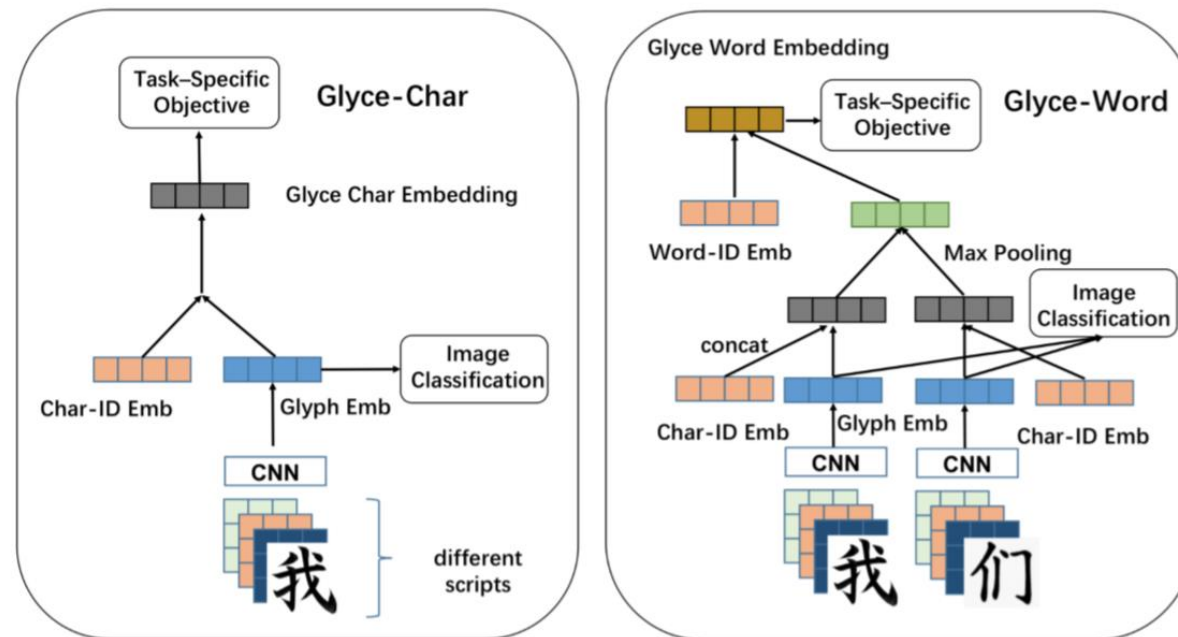
(A) CBOW  (B)Character-enhanced Word Embedding

# Glyph

- GWE

- Glyce

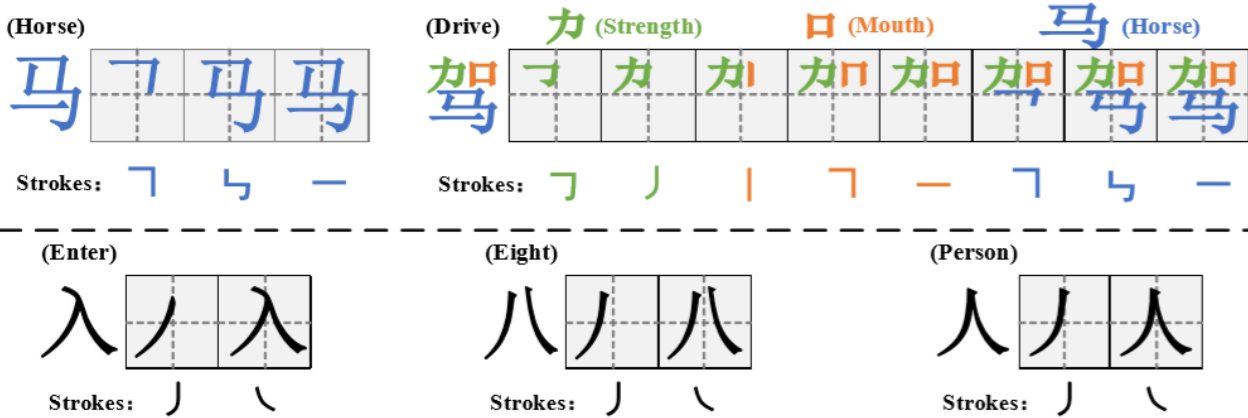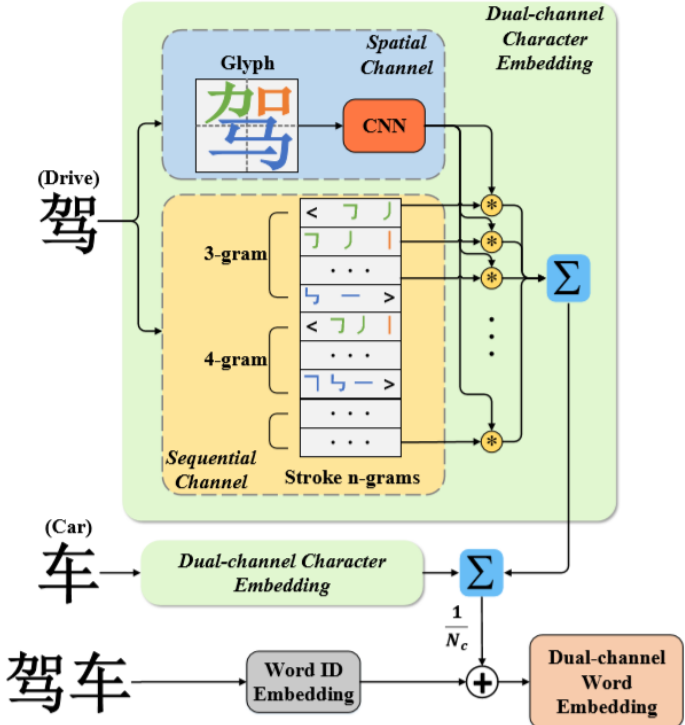# Dual channel view for Morphology Logogram

**Sequential**

- Character

- Stroke

**Spatial**

- Glyph

# Problem

- High Time Complexity

- High Computing Resource

- Weak Interpretability
  - Only can infer similarity and analogy

# Future Work

- Try to compress corpus into knowledge graph
  - Reduce time complexity
  - Reduce computing Resource
  - Strong interpretability
  - High interactivity

# Summary

- Challenge
  - Semantics
    - Distributed Embedding
  - Polysemy
    - Feature Model
  - OOV
- Quality Evaluation
  - General semantics task: word similarity, word analogy
  - Task specific evaluation: classification, translation
- Practical
  - Time Complexity
    - Negative Sampling
  - Interpretability
    - ?

# Reference

- Mikolov T, Sutskever I, Chen K, et al. Distributed representations of words and phrases and their compositionality[C]//Advances in neural information processing systems. 2013: 3111-3119.

- Mikolov T, Chen K, Corrado G, et al. Efficient estimation of word representations in vector space[J]. arXiv preprint arXiv:1301.3781, 2013.

- Pennington J, Socher R, Manning C. Glove: Global vectors for word representation[C]//Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP). 2014: 1532-1543.

- Kim Y. Convolutional neural networks for sentence classification[J]. arXiv preprint arXiv:1408.5882, 2014.

- Peters M E, Neumann M, Iyyer M, et al. Deep contextualized word representations[J]. arXiv preprint arXiv:1802.05365, 2018.

- Radford A, Narasimhan K, Salimans T, et al. Improving language understanding by generative pre-training[J]. URL https://s3-us-west-2. amazonaws. com/openai-assets/research-covers/languageunsupervised/language understanding paper. pdf, 2018.

- Devlin J, Chang M W, Lee K, et al. Bert: Pre-training of deep bidirectional transformers for language understanding[J]. arXiv preprint arXiv:1810.04805, 2018.

# Reference

- Bojanowski P, Grave E, Joulin A, et al. Enriching word vectors with subword information[J]. Transactions of the Association for Computational Linguistics, 2017, 5: 135-146.

- Chen X, Xu L, Liu Z, et al. Joint learning of character and word embeddings[C]//Twenty-Fourth International Joint Conference on Artificial Intelligence. 2015.

- Cao S, Lu W, Zhou J, et al. cw2vec: Learning chinese word embeddings with stroke n-gram information[C]//Thirty-Second AAAI Conference on Artificial Intelligence. 2018.

- Su T R, Lee H Y. Learning chinese word representations from glyphs of characters[J]. arXiv preprint arXiv:1708.04755, 2017.

- Wu W, Meng Y, Han Q, et al. Glyce: Glyph-vectors for Chinese Character Representations[J]. arXiv preprint arXiv:1901.10125, 2019.

# Q&A