

# Tokenization

Vũ Việt Trung

# Tokenization

Hiệp hội Vận tải hàng hóa VN chính thức có văn bản báo cáo và kiến nghị Văn phòng Quốc hội



Hiệp\_1 hội\_3 Vận\_1 tải\_3 hàng\_1 hóa\_3 VN\_0 chính\_1 thức\_3  
có\_0 văn\_1 bản\_3 báo\_1 cáo\_3 và\_0 kiến\_1 nghị\_3 Văn\_1  
phòng\_3 Quốc\_1 hội\_3

- 0: outside (O)
- 1: start (S)
- 2: inside (I)
- 3: end (E)

# Language model

- Xác suất xuất hiện một cụm từ

$$\begin{aligned} P(w_1 \dots w_n) &= P(w_1)P(w_2 \dots w_n \mid w_1) \quad \leftarrow \text{Chain rule } P(x, y) = P(x)P(y \mid x) \\ &= P(w_1)P(w_2 \mid w_1)P(w_3 \dots w_n \mid w_1 w_2) \\ &= P(w_1)P(w_2 \mid w_1)P(w_3 \mid w_1 w_2)P(w_4 \dots w_n \mid w_1 w_2 w_3) \\ &= \dots \\ &= \prod_{i=1}^n P(w_i \mid w_1 \dots w_{i-1}) \end{aligned}$$

# N-grams

- Markov assumption
  - history
- Maximum likelihood estimation (MLE)
  - Count  $f$

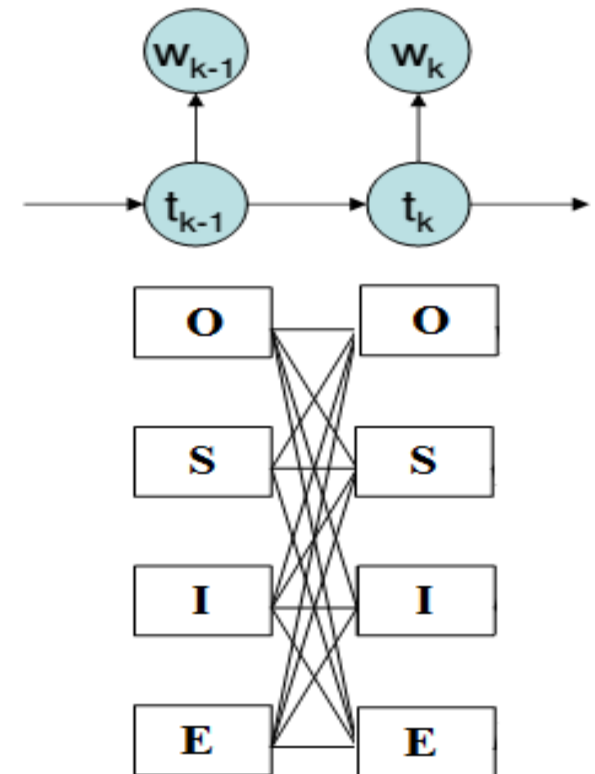
- Unigram: 
$$P(w_1 \dots w_n) = \prod_{i=1}^n P(w_i \mid w_1 \dots w_{i-1})$$
$$\approx \prod_{i=1}^n P(w_i)$$

- Bigram:

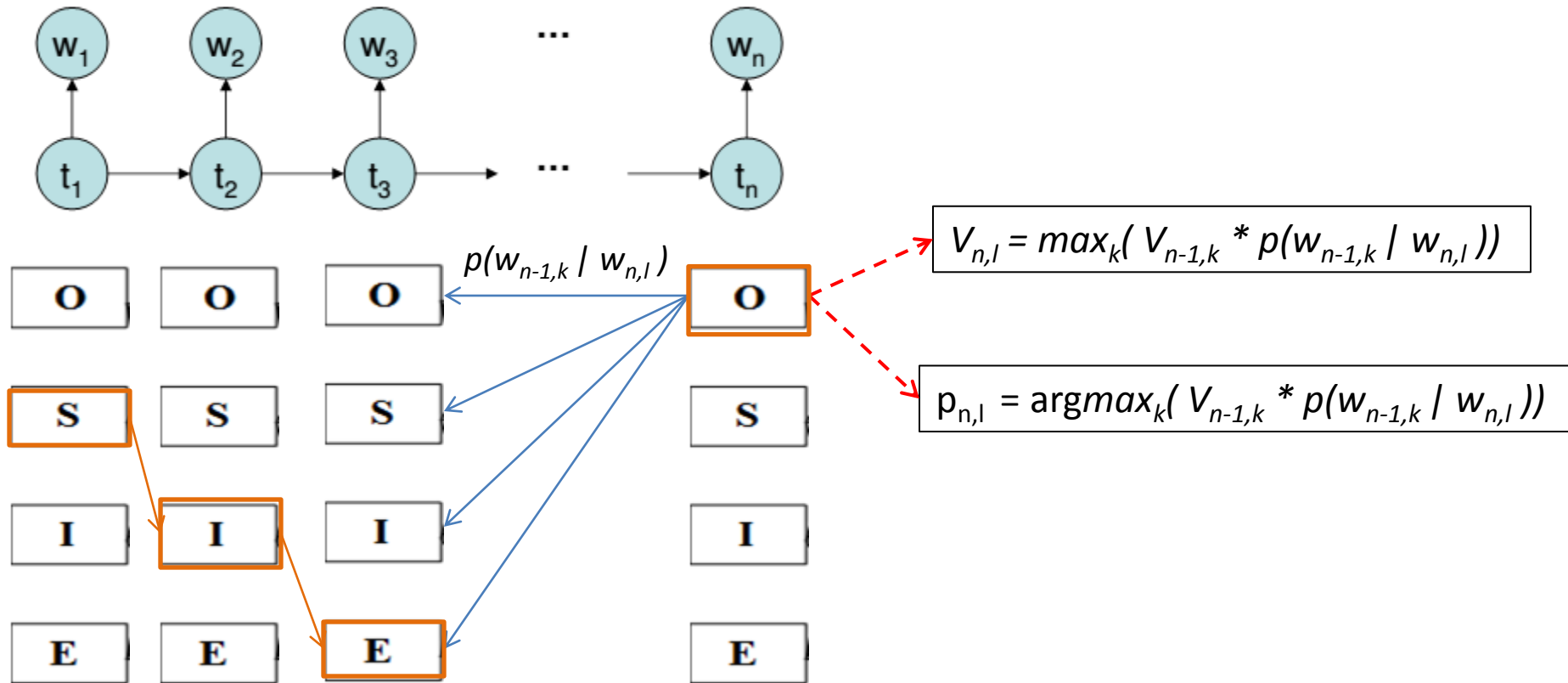
$$P(w_1 \dots w_n) = \prod_{i=1}^n P(w_i \mid w_1 \dots w_{i-1})$$
$$\approx \prod_{i=1}^n P(w_i \mid w_{i-1})$$

# Labeled-bigram

- Dựa trên ý tưởng N-grams:
  - Mỗi từ được gán thêm nhãn (0-3):  
*kiến\_1 nghị\_3*
  - Bài toán tách từ  $\rightarrow$  Bài toán gán nhãn từ sao cho câu có xác suất lớn nhất
  - Tương tự giải thuật trong POS tagging (Viterbi)



# Dynamic Programming



Tìm đường đi có xác suất lớn nhất ( $v_n$ )



Dùng log  $\rightarrow$  nhân thành cộng, số không quá bé  $\rightarrow$  tăng hiệu năng

# Tính trọng số $p(w_{n-1,k} / w_{n,l})$

- Unsmooth

$$P(w_k) = \frac{C(w_k)}{\sum_w C(w_k)}$$

$$P(w_k | w_{k-1}) = \frac{C(w_{k-1}w_k)}{C(w_{k-1})}$$

- Smooth

- Nhiều phương pháp:

- Add one, Good Turing, Interpolation, Katz (backoff), Absolute Discounting,...

- Kneser-Ney – hiệu quả nhất

# Modified Kneser-Ney Smoothing

$$p_{\text{KN}}(w_i | w_{i-n+1}^{i-1}) = \frac{c(w_{i-n+1}^i) - D(c(w_{i-n+1}^i))}{\sum_{w_i} c(w_{i-n+1}^i)} + \gamma(w_{i-n+1}^{i-1}) p_{\text{KN}}(w_i | w_{i-n+2}^{i-1})$$

$$D(c) = \begin{cases} 0 & \text{if } c = 0 \\ D_1 & \text{if } c = 1 \\ D_2 & \text{if } c = 2 \\ D_{3+} & \text{if } c \geq 3 \end{cases}$$

$$\begin{aligned} Y &= \frac{n_1}{n_1 + 2n_2} \\ D_1 &= 1 - 2Y \frac{n_2}{n_1} \\ D_2 &= 2 - 3Y \frac{n_3}{n_2} \\ D_{3+} &= 3 - 4Y \frac{n_4}{n_3} \end{aligned}$$

$$n_1 = N_1(\bullet\bullet)$$

$$p_{\text{KN}}(w_i) = \frac{N_{1+}(\bullet w_i)}{N_{1+}(\bullet\bullet)}$$

$$N_{1+}(\bullet w_i) = |\{w_{i-1} : c(w_{i-1} w_i) > 0\}|$$

$$N_{1+}(\bullet\bullet) = |\{(w_{i-1}, w_i) : c(w_{i-1} w_i) > 0\}|$$

$$\gamma(w_{i-n+1}^{i-1}) = \frac{D_1 N_1(w_{i-n+1}^{i-1} \bullet) + D_2 N_2(w_{i-n+1}^{i-1} \bullet) + D_{3+} N_{3+}(w_{i-n+1}^{i-1} \bullet)}{\sum_{w_i} c(w_{i-n+1}^i)}$$



# Kết hợp vào bài toán tokenize

- Dictionary
  - Theo từng topic,  $p_{\text{dict}}=1$  nếu có xuất hiện
  - $p_{\text{mix}} = 0.5 p_{\text{KN}} + 0.5 p_{\text{dict}}$
- Mix các corpus
  - $p_{\text{KN}} = w_1 * p_{\text{KN1}} + w_2 * p_{\text{KN2}} + \dots + w_n * p_{\text{KNn}}$ 
    - $w_1 + w_2 + \dots + w_n = 1$
  - Trọng số  $w_i$  tỷ lệ với số quy mô của corpus  $i$

# Training & Test

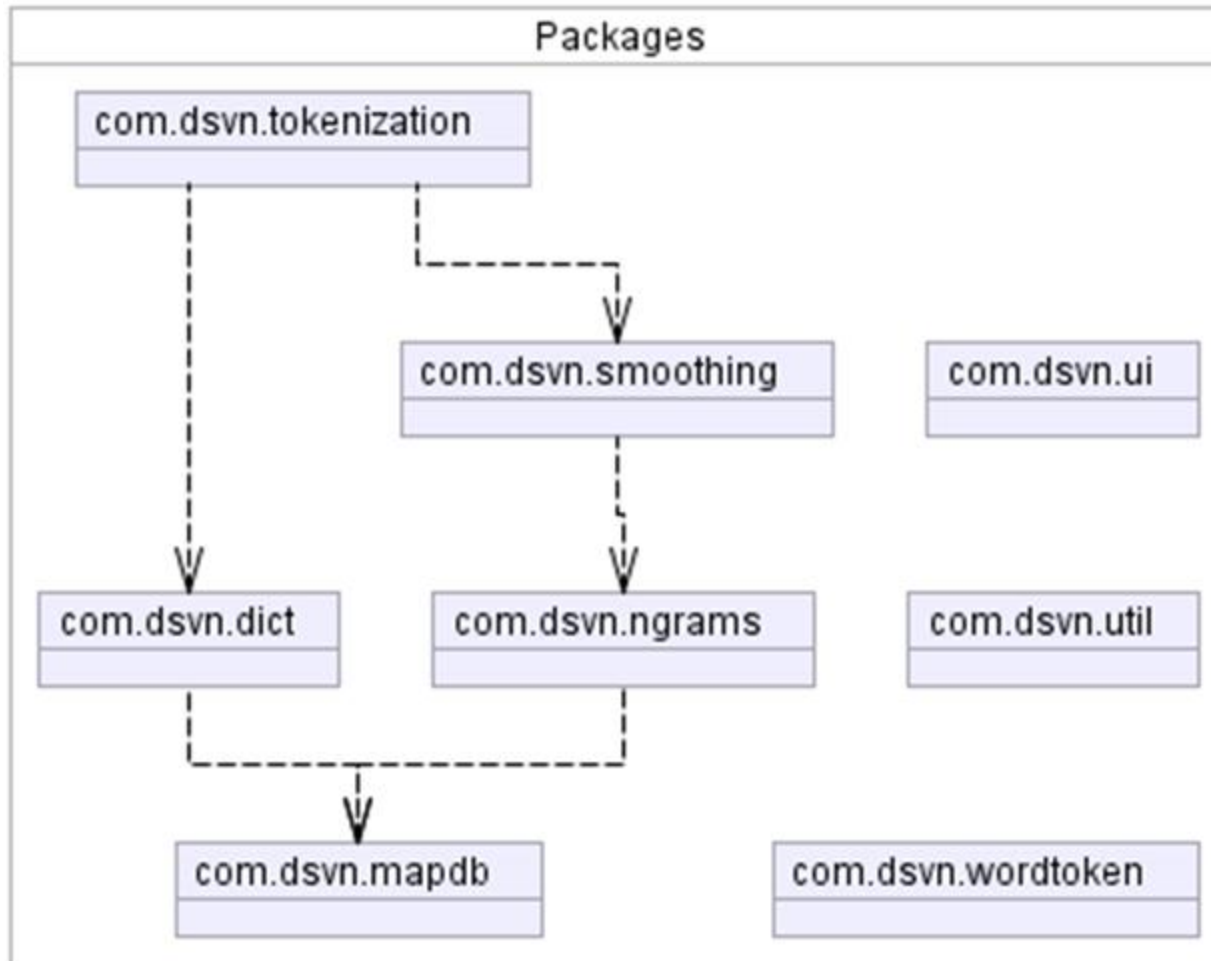
- Chia dữ liệu thành:
  - 80% train (học mô hình  $p_{KN}$ )
  - 10% held-out (điều chỉnh tham số)
  - 10% test
    - Thử gán nhãn từ
      - sau đó có thể sửa lại = heuristic để làm tập training mới
- Đánh giá: **perplexity**
  - pp càng nhỏ càng tốt

$$\begin{aligned}\text{Perplexity}(W) &= 2^{H(W)} \\ &= P(w_1 w_2 \dots w_N)^{-\frac{1}{N}} \\ &= \sqrt[N]{\frac{1}{P(w_1 w_2 \dots w_N)}} \\ &= \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i | w_1 \dots w_{i-1})}}\end{aligned}$$

# Giả thiết đầu vào – đầu ra

- Mỗi dòng là 1 câu
- Thêm 2 từ khóa START và BEGIN để đảm bảo phân bố xác suất
  - <BEGIN> Hôm nay trời đẹp <END>
- Thêm từ khóa <UNK> (xxx ?) cho các từ chưa có trong dictionary

# Implementation



com.dsvn.wordtoken

**MyNode**

**WordNode**

**WordLabel**

**SpecialNode**

com.dsvn.ui

**MapDBManagement**

**BigramManagement**

**DictionaryManagement**

**UnigramManagement**

**TokenizationDemo**

**SmoothedManagement**

