







Gióng hàng câu

The old man is happy. He has fished many times. His wife talks to him. The fish are jumping. The sharks await.

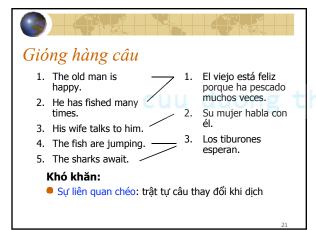
El viejo está feliz porque ha pescado muchos veces. Su mujer habla con él. Los tiburones esperan.

19

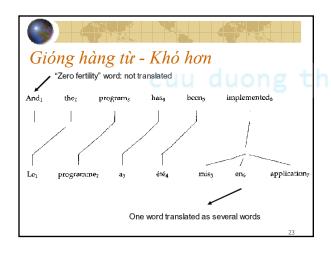


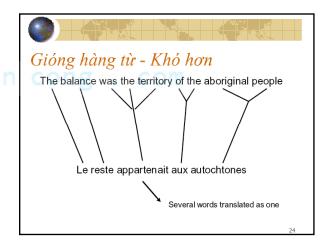
Gióng hàng câu

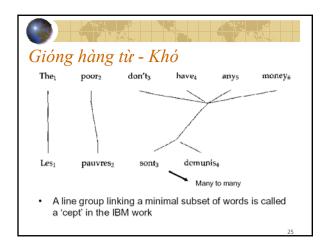
- 1. The old man is happy.
- He has fished many times.
- 3. His wife talks to him.
- 4. The fish are jumping.
- 5. The sharks await.
- 1. El viejo está feliz porque ha pescado muchos veces.
- 2. Su mujer habla con él.
- 3. Los tiburones esperan.

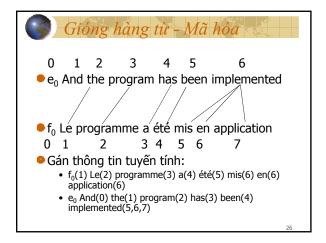


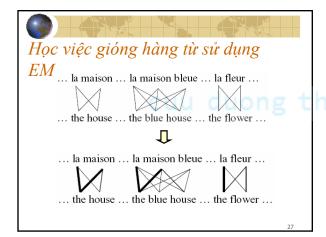


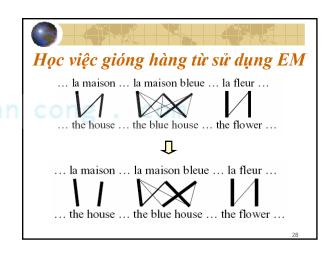


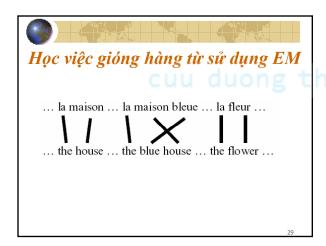


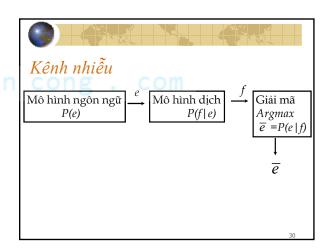














Các thành phần của mô hình dịch

- Giả thiết
 - Individual translations are independance
 - 1 từ tiếng Anh n từ tiếng Pháp
 - 1 từ tiếng Pháp (0-1) từ tiếng Anh

$$P(f | e) = \frac{1}{Z} \sum_{a_1}^{l} \cdots \sum_{a_{m=0}}^{l} \prod_{j=1}^{m} P(f_j | e_{a_j})$$

- f_j từ j trong f;
 a_j vị trí trong e được gióng hàng với f_j
 e_{aj} từ trong e được gióng hàng với f_j
 Z là hằng số chuẩn hóa
 a_j = 0: từ j trong câu tiếng Pháp được gióng hàng với một từ rỗng (không dịch sang)
- m độ dài của f



Ví du

- P(Jean aime Marie/ John loves Mary)
- Gióng hàng (Jean, John), (aime, loves), (Marie, Mary), ta có 3 xác suất P(Jean|John) x P(aime|loves) xP(Marie|Mary)



Giải mã

$$\bar{e} = \arg\max_{e} P(e \mid f)$$

$$= \arg\max_{e} \frac{P(e)P(f \mid e)}{P(f)}$$

$$= \arg\max_{e} P(e)P(f \mid e)$$

Vấn đề: không gian tìm kiếm vô hạn Meo:

- tìm kiếm dùng ngăn xếp: xây dựng dần, lưu trong stack các phần đã dịch
- sử dụng một số độ đo về độ phù hợp, vd., chamber/house, (nhưng có thể đi sai đường nếu 1 từ thường xuất hiện với từ khác, như commune/house, vì có Chambre de Communes (hạ nghị viên)



Thuật toán EM

E-step

- Khởi tạo giá trị P(w_f/w_e) ngẫu nhiên
- Tính số lần tìm thấy w_f trong tiếng Pháp khi có w_e trong

$$z_{w_f, w_e} = \sum_{(e, f) \text{s.t.} w_e = e, w_f = f} P(w_f \mid w_e)$$

Đánh giá lại xác suất dịch prs từ giá trị z trên:

$$P(w_f \mid w_e) = \frac{z_{w_f, w_e}}{\sum_{v} z_{v, w_e}}$$

tổng được tính trên tất cả các từ tiếng Pháp v



Đánh giá

Đánh giá dưa trên tập ngữ liệu Hansard:

- 48% câu tiếng Pháp được dịch đúng
- 2 loai lõi:
 - Dịch sai nghĩa:
 - Permettez que je donne un example à chambre
 - Let me give an example in the House (incorrect decoding)
 - (Let me give the House an example)
 - Dịch sai ngữ pháp:
 - Vous avez besoin de toute l'aide disponsible
 - You need all of the benefits available (ungrammatical
 - (You need all the help you can get)



Lý do

- Hiện tượng méo: từ tiếng Anh ở đầu câu được gióng hàng với từ tiếng Pháp ở cuối câu - hiện tượng này giảm xác suất giống hàng
- Hiện tượng sinh (fertility): sự tương ứng giữa từ tiếng Anh và tiếng Pháp (1-to-1, 1-to-2, 1-to-0,
 - Vd, fertility(farmers) trong tâp ngữ liêu = 2, vì từ này khi dịch sang tiếng Anh thường gồm 2 từ: les argiculteurs
 - To go → aller



Lý do

- Các giả thiết độc lập: các câu ngắn được ưu tiên hơn vì có ít xác suất hơn (khi nhân) ⇒ nhân kết quả với 1 hằng số tỉ lệ thuân với đô dài câu
- Phụ thuộc dữ liệu luyện: 1 thay đổi nhỏ trong dữ liệu luyên gây ra thay đổi lớn trong các giá trị ước lượng
 - Vd, P(le/the) thay đổi từ 0.610 xuống 0.497
- TÍnh hiệu quả. Bỏ các câu > 30 từ, vì làm không gian tìm kiếm tăng theo cấp số mũ
- Thiếu tri thức ngôn ngữ



Thiếu tri thức ngôn ngữ

- Không lưu thông tin về các ngữ: ví du không gióng hàng được "to go" và "aller"
- Không có ràng buộc cục bộ:

Eg, is she a mathematician

- Âm vị. Các từ tạo bởi các âm vị khác nhau được coi là các ký hiệu riêng biệt
- Dữ liệu thưa. Các đánh giá cho các từ ít gặp không chính xác



Các hệ thống gióng hàng khác

- Các tập ngữ liêu sử dung giả thiết:
 - Dữ liệu song song (dịch $E \leftrightarrow F$)
- Gióna hàna câu
 - Phát hiên câu
 - · Gióng hàng câu
- Gióng hàng từ
 - Tách từ
 - Gióng hàng từ (với 1 số ràng buộc)



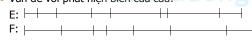
Phát hiện biên của câu

- Sử dụng luật, danh sách liệt kê:
- Dấu kết thúc câu:
 - Dấu ngắt đoạn (nếu được đánh dấu)
 - 1 số ký tự: ?, !, ;
 - Vấn đề: dấu chấm \.'
 - Kết thúc câu (... left yesterday. He was heading to...)
 - Dấu chấm thập phân : 3.6 (three-point-six)
 - Dấu chấm hàng nghìn: 3.200
 - Viết tắt: cf., e.g., Calif., Mt., Mr.
 - Vân vân: ...
 - 1 số ngôn ngữ: 2nd ~ 2.
 - Ký hiệu đầu: A. B. Smith
- Phương pháp thống kê: vd Maximum Entropy



Gióng hàng câu

Vấn đề với phát hiện biên của câu:



- Đầu ra mong đơi: Các phân mảnh với cùng số lượng mánh liên tiếp nhau.
- Gióng hàng:



Kết quả: 2-1, 1-1, 1-1, 2-2, 2-1, 0-1

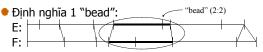
Các phương pháp gióng hàng

- Nhiều phương pháp (xác suất hoặc không)
 - Dựa trên độ dài ký tự
 - Dựa trên độ dài từ
 - "cùng gốc" (sử dụng nghĩa từ)
 - Sử dụng từ điển (F: prendre ~ E: make, take)
 - Sử dụng khoảng cách từ (độ tương tự): tên, số, từ vay mượn, từ gốc Latin
- Kết quả tốt nhất:
 - Thống kê, dựa trên từ hoặc dựa trên ký tự



Gióng hàng dựa trên độ dài

Định nghĩa bài toán như việc tính xác suất:
 argmax_A P(A|E,F) = argmax_A P(A,E,F) (E,F cố định)



Lấy xấp xỉ:

 $P(A,E,F) \cong \Pi_{i=1..n}P(B_i),$

Trong đó B_i là 1 bead; $P(B_i)$ không phụ thuộc vào phần còn lai của E_iF_i

42



Nhiệm vụ gióng hàng

Định nghĩa:

- Cho P(A,E,F) ≅ Π_{i=1..n}P(B_i),
 tìm cách chia (E,F) thành n bead B_{i=1..n}, sao cho tối đa xác suất P(A,E,F) trên tập luyện.
- $B_i = {}_{p:q}\alpha_i$, với $p:q \in \{0:1,1:0,1:1,1:2,2:1,2:2\}$ mô tả phép gióng hàng
- Pref(i,j) xác suất của cách gióng hàng tốt nhất từ điểm đầu cho đến (i,j)

44



Định nghĩa đệ qui

- Khởi tạo: Pref(0,0) = 0.
- Pref(i,j) = max (Pref(i,j-1) $P(0:1\alpha_k)$, Pref(i-1,j) $P(1:0\alpha_k)$, Pref(i-1,j-1) $P(1:1\alpha_k)$, Pref(i-1,j-2) $P(1:2\alpha_k)$, Pref(i-2,j-1) $P(1:1\alpha_k)$, Pref(i-2,j-2) $P(1:2\alpha_k)$
- E: • F: Pref(i,j-1)

,



Xác suất của 1 Bead

- Định nghĩa P(_{p:q}α_k):
 - $\underline{\mathbf{k}}$ đề cập đến "bead" kế tiếp, với các đoạn của câu p và q, độ dài \mathbf{l}_{ke} và \mathbf{l}_{kf} .
- Sử dụng phân bố chuẩn cho các độ dài khác nhau: $P(p_{p:q}\alpha_k) = P(\delta(l_{k,e},l_{k,f},\mu,\sigma^2),p;q) \cong P(\delta(l_{k,e},l_{k,f},\mu,\sigma^2))P(p;q)$ $\delta(l_{k,e},l_{k,f},\mu,\sigma^2) = (l_{k,f}-\mu l_{k,e})/\sqrt{l_{k,e}}\sigma^2$
- Đánh giá P(p:q) từ tập dữ liệu nhỏ, hoặc đoán và đánh gía lai sau khi gióng hàng
- Từ có thể được dùng như dấu hiệu tốt hơn để định nghĩa $P(_{p:a}a_k)$.



Gióng hàng từ

- Nếu chỉ dựa trên độ dài, không thực hiện được:
 - từ có thể bị đảo trật tự, các phép dịch thường có độ dài khác nhau
- Ý tưởng:
 - Đưa ra vài mô hình dịch đơn giản.
 - Tìm các tham số bằng cách xét tất cả các cách gióng hàng.
 - Sau khi có tham số, tìm cách gióng hàng tốt nhất khi có các tham số này.

Thuật toán gióng hàng từ

Khởi tạo với tập ngữ liệu gióng hàng câu. Cho (E,F) là 1 cặp câu (là 1 bead).

- 1. Khởi tạo ngẫu nhiên p(f|e), $f \in F$, $e \in E$.
- 2. Đếm trên tập ngữ liệu:

 $c(f,e) = \Sigma_{(E,F);e \in E,f \in F} p(f|e)$

với ∀ cặp gióng hàng (E,F), kiểm tra xem e có trong E và f có trong F không. Nếu đúng, bổ sung p(f|e).

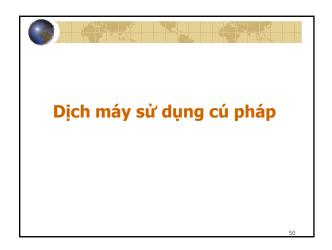
3. Đánh giá lại:

 $p(f|e) = c(f,e) / c(e) [c(e) = \Sigma_f c(f,e)]$

4. Lặp đến khi p(f|e) thay đổi ít.

ii p(i je) tilay doi it.





Tại sao dùng cú pháp
Cần thông tin ngữ pháp
Cần các ràng buộc khi sắp lại câu
Khi chèn các từ chức năng vào câu, cần đặt ở vị trí chính xác
Khi dịch từ cần sử dụng từ có cùng từ loại với nó



Mô hình dựa trên cú pháp

Cây cú pháp
(tiếng Anh)

Mô hình dịch

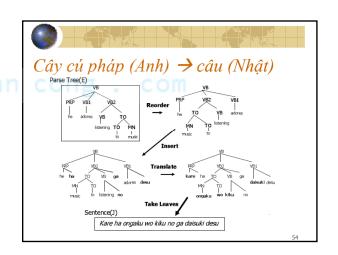
Tiền xử lý câu tiếng Anh bằng bộ PTCP

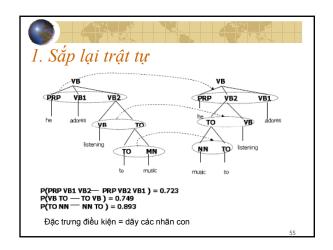
Thực hiện các phép tính xác suất trên cây cú pháp

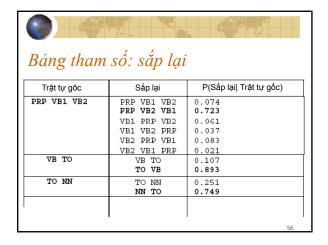
Sắp lại trật tự các nút

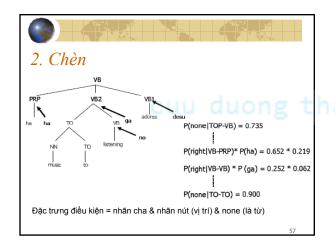
Chèn nút mới vào

Dịch các từ ở lá

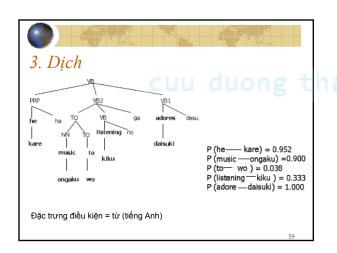


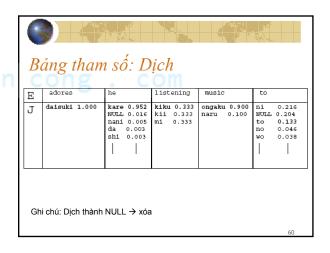




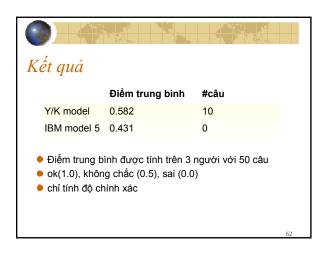


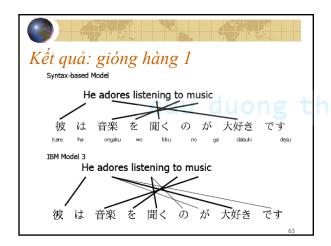


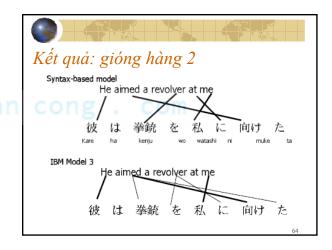






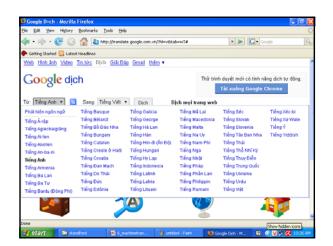


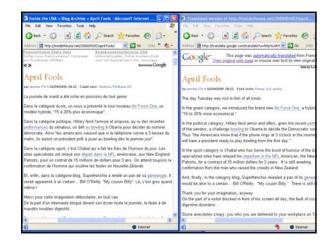
















cuu duong than cong . com