



ĐẠI HỌC BÁCH KHOA HÀ NỘI
VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

Dịch máy

Viện CNTT & TT – Trường ĐHBKHN

Ví dụ

- Au sortir de la saison 97/98 et surtout au debut de cette saison 98/99...
- With leaving season 97/98 and especially at the beginning of this season 98/99...

Các vấn đề

1. Xử lý sự giống và khác nhau giữa các ngôn ngữ
 - Hình vị: # số âm tiết/từ:
 - *Ngôn ngữ đơn âm tiết (tiếng Việt, Trung Quốc) – 1 tiếng/từ*
 - *Ngôn ngữ đa âm tiết (Siberian Yupik), 1 từ = cả 1 câu*
 - Mức độ phân chia âm tiết

Các vấn đề

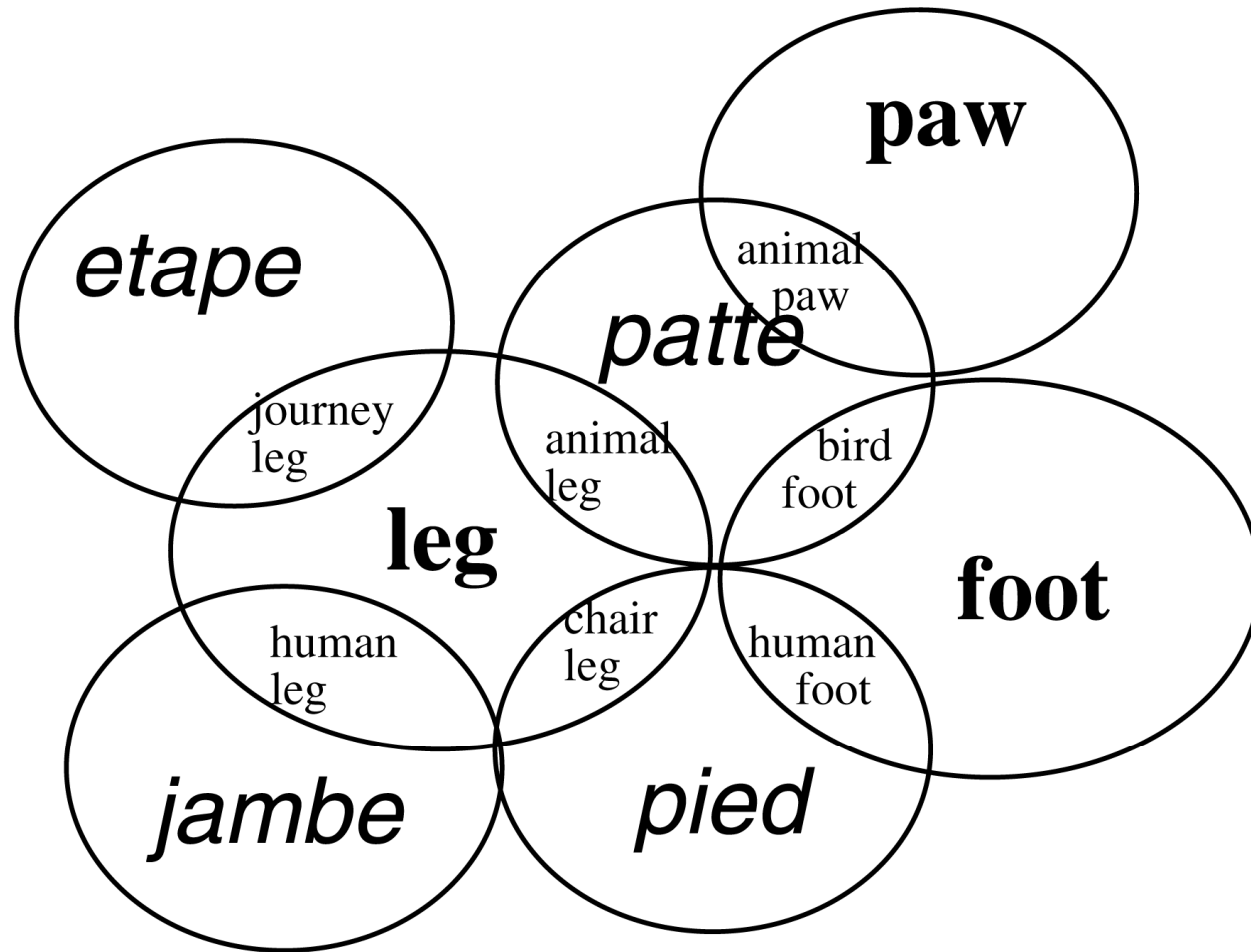
2. Cú pháp: trật tự từ trong câu

- *To Yukio; Yukio ne*
- Tiếng Anh – tiếng Việt:
 - *The* (affix1) *red* (affix2) *flag* (head)
 - *Lá cờ* (head) *đỏ* (affix2) *ấy* (affix1)

3. Các nét riêng biệt

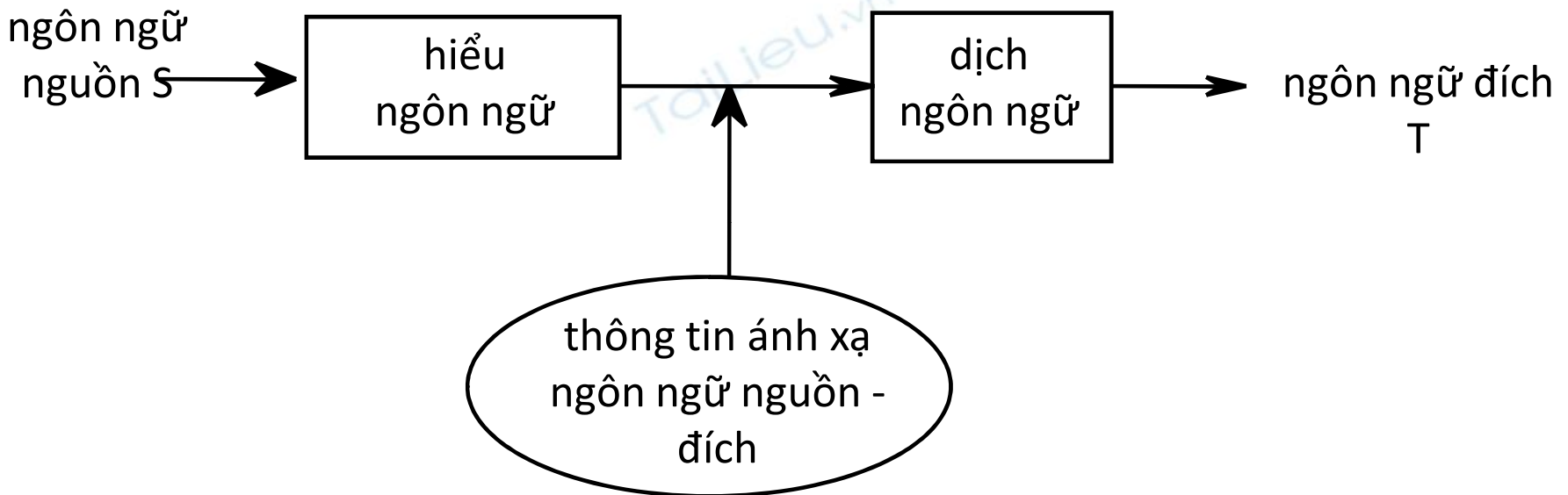
English	brother	Vietnamese	anh em
English	wall	German	wand (inside) mauer(outside)
German	berg	English	hill mountain

Không gian khái niệm

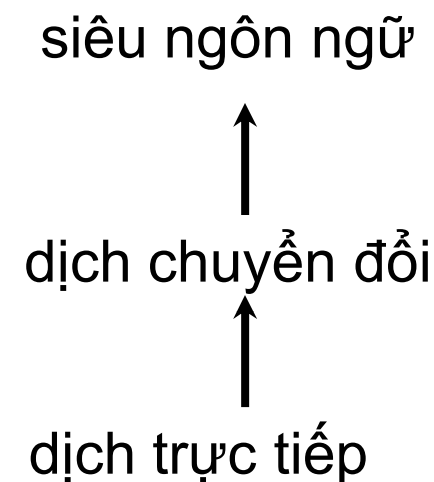
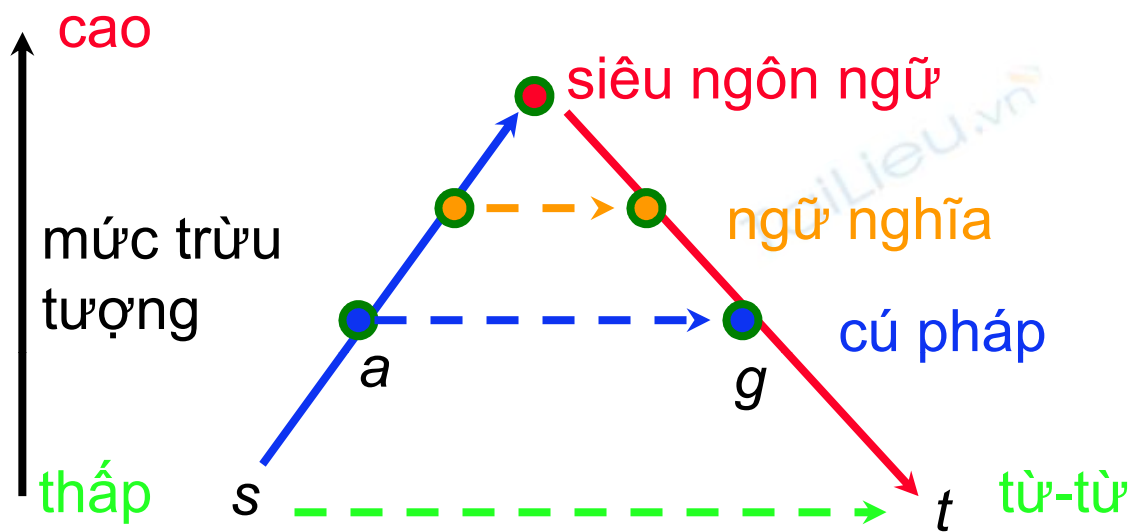


Khoảng trống từ vựng: tiếng Nhật không có từ nào nghĩa *privacy*;
tiếng Anh không có từ ứng với *yakoko* (lòng hiếu thảo)

Ba khối chính trong dịch máy



Các phương pháp dịch máy

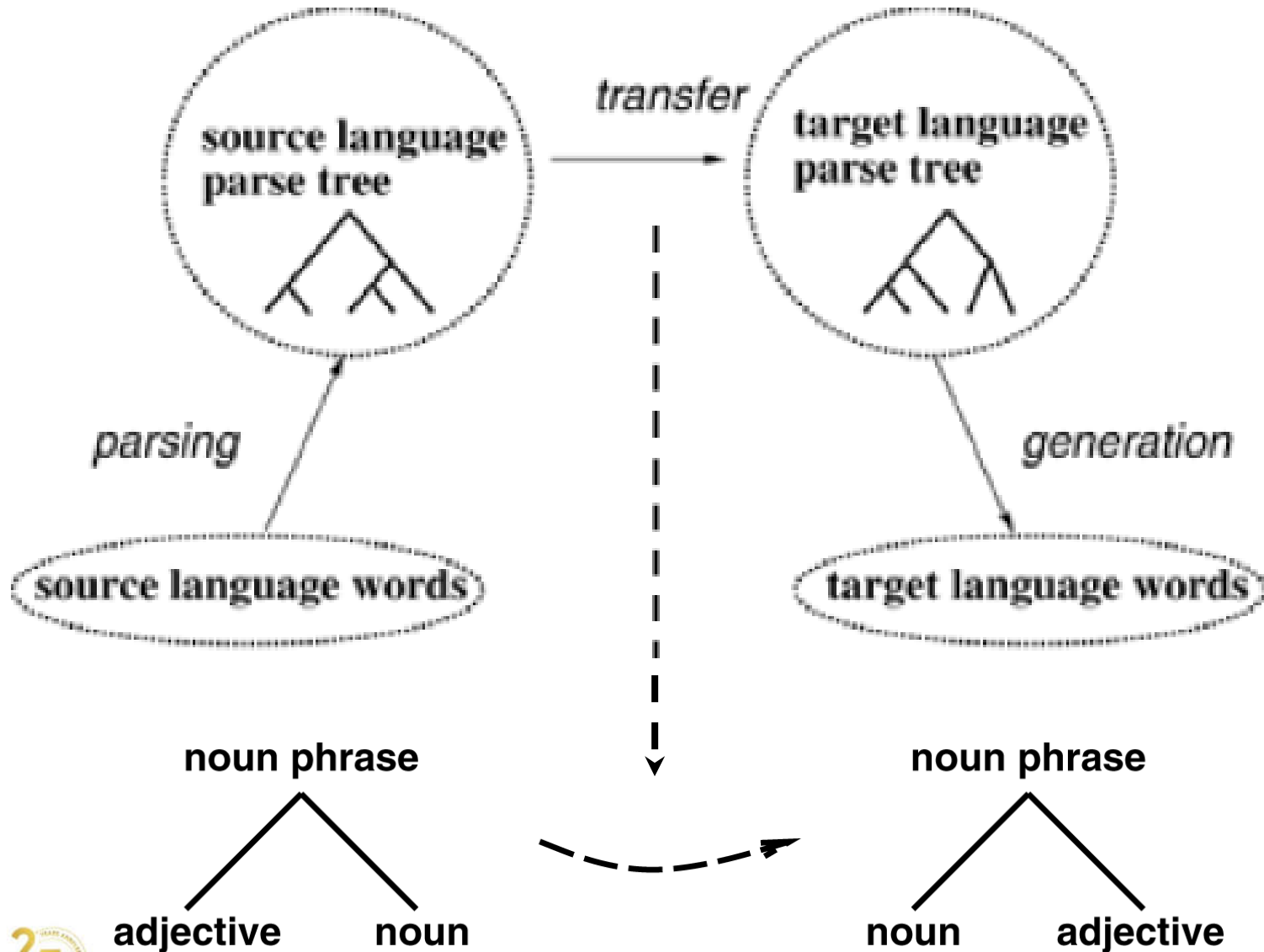


$$a = a(s)$$

$$g = f(a(s)); f - \text{hàm chuyển đổi}$$

$$t = g(f(a(s)))$$

Sơ đồ chuyển đổi



Luật chuyển đổi

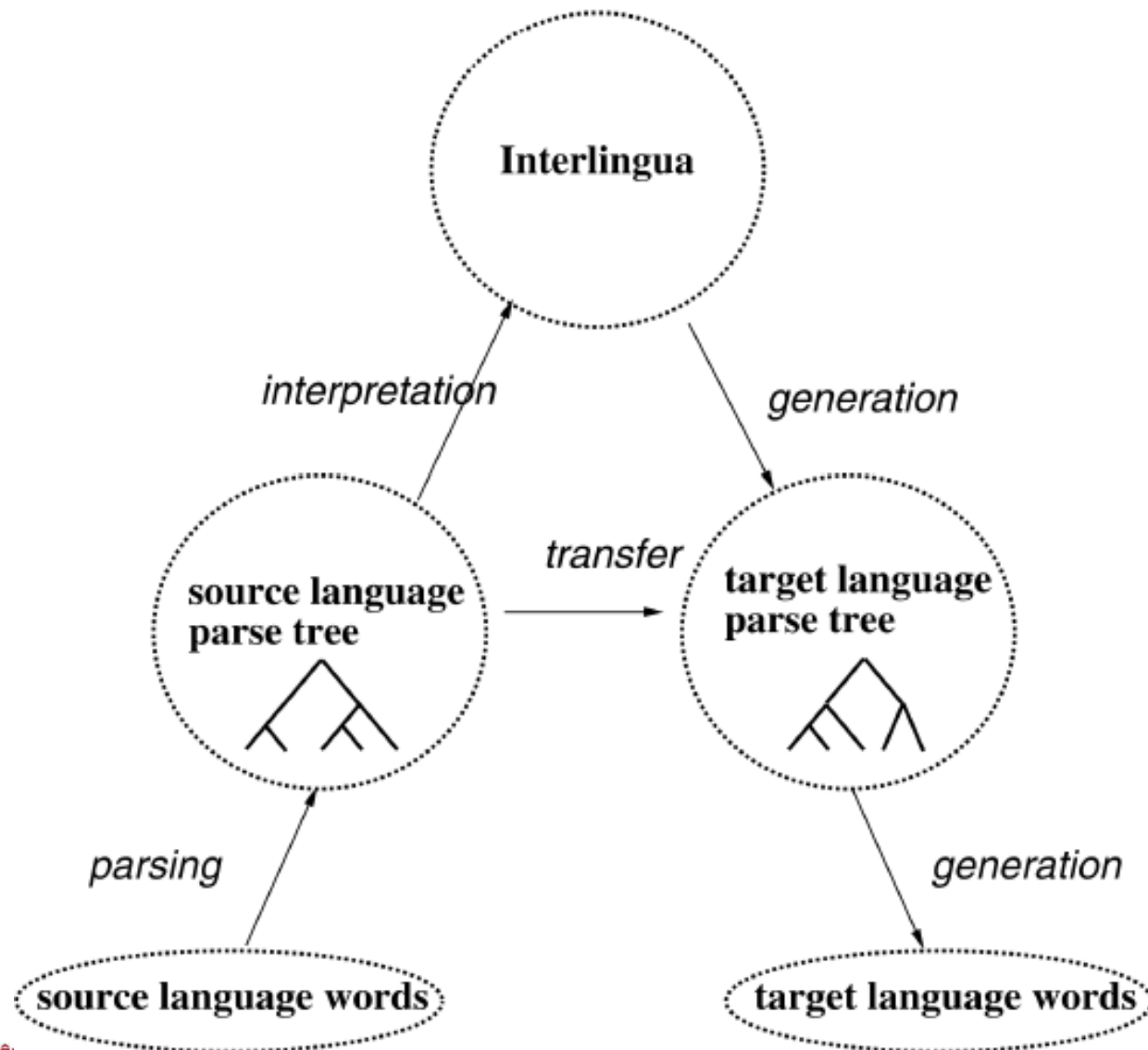
English to French:

1. NP \rightarrow Adjective₁ Noun₂
 \Rightarrow
NP \rightarrow Noun₂ Adjective₁

Japanese to English:

2. Existential-There-Sentence \rightarrow There₁ Verb₂ NP₃ Postnominal₄
 \Rightarrow
Sentence \rightarrow (NP \rightarrow NP₃ Relative-Clause₄) Verb₂
3. NP \rightarrow NP₁ Relative Clause₂
 \Rightarrow
NP \rightarrow Relative-Clause₂ NP₁

Sơ đồ chuyển đổi



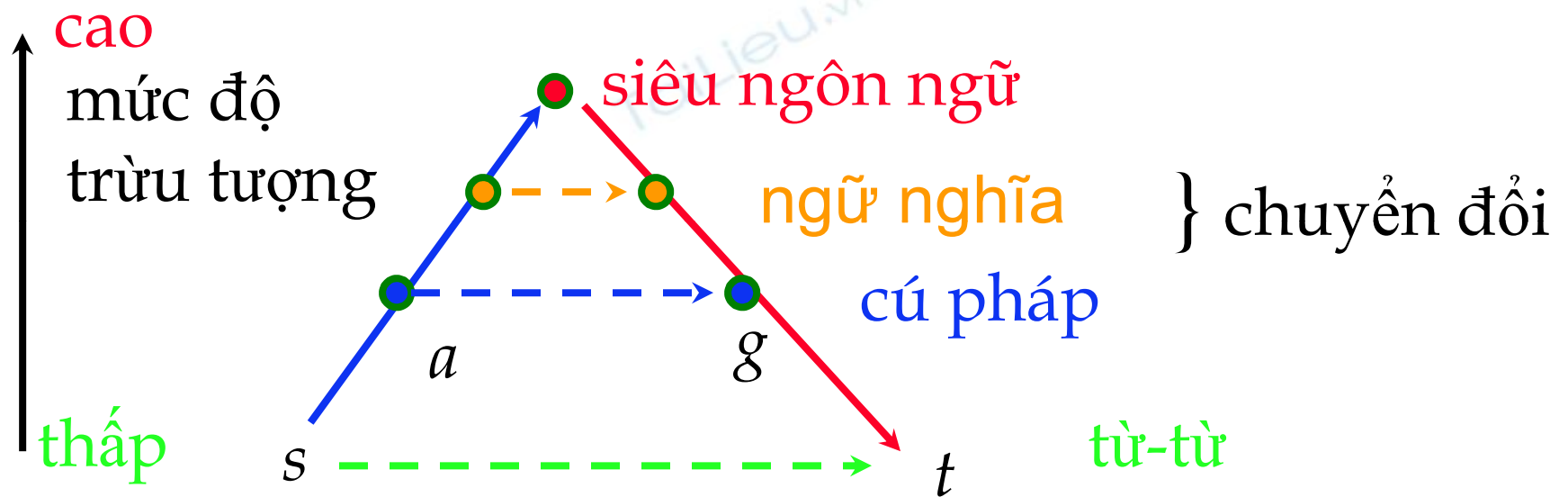
Cách tiếp cận siêu ngôn ngữ: sử dụng nghĩa

- Chuyển đổi: các luật chuyển đổi từ ngôn ngữ này sang ngôn ngữ khác
- Đối tượng/sự kiện (ontology)

event	gardening						
agent	<table><tr><td>man</td><td></td></tr><tr><td>number</td><td>sg</td></tr><tr><td>definiteness</td><td>indef</td></tr></table>	man		number	sg	definiteness	indef
man							
number	sg						
definiteness	indef						
aspect	progressive						
tense	past						

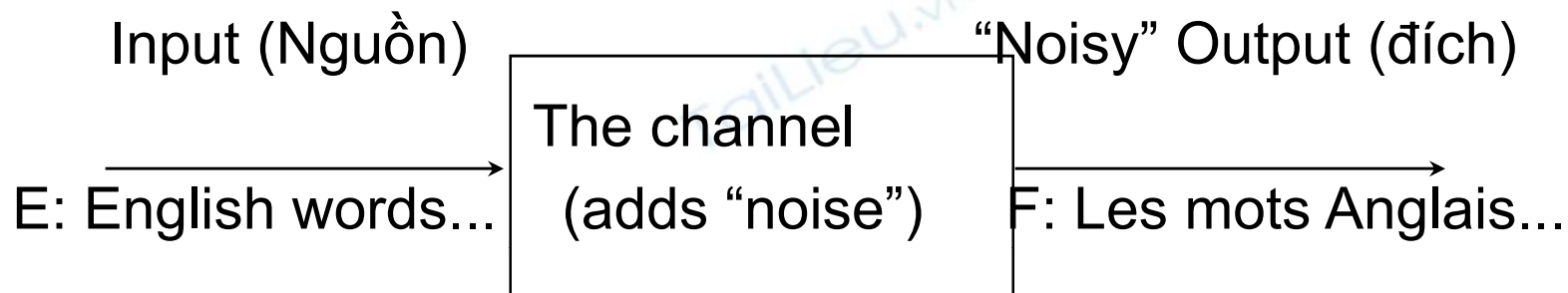
Dịch máy thống kê

Các kiểu dịch máy



ý tưởng

- Coi việc dịch như bài toán kênh có nhiễu



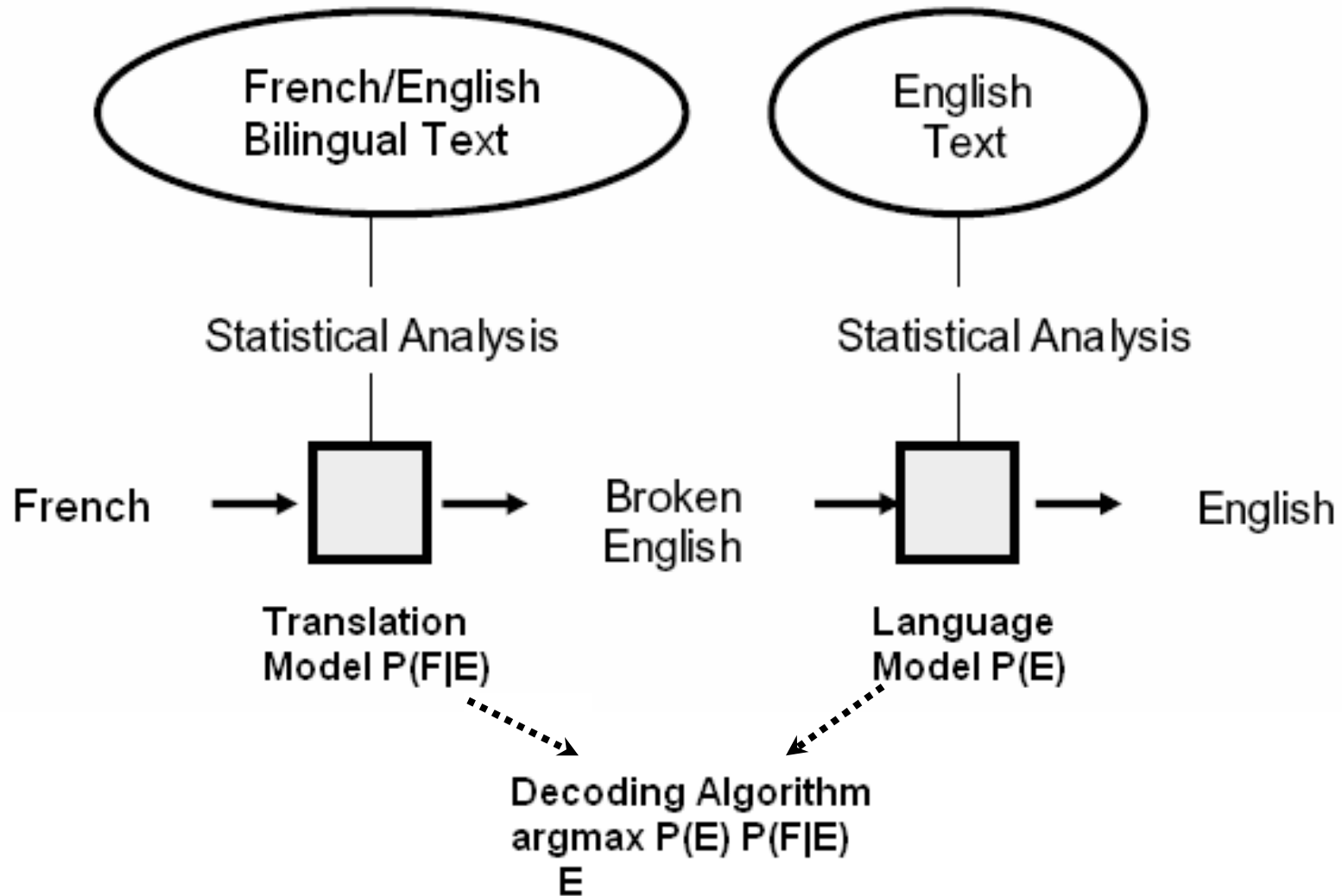
- Mô hình dịch: $P(E|F) = P(F|E) P(E) / P(F)$

- Khôi phục lại \underline{E} khi biết \underline{F} :

Sau khi đơn giản hóa ($P(F)$ không đổi):

$$\operatorname{argmax}_E P(E|F) = \operatorname{argmax}_E P(F|E) P(E)$$

Dịch máy thống kê



Các yếu tố

- **Mô hình ngôn ngữ - Language Model (LM)**: xác suất thấy 1 câu tiếng Anh (E) (xác suất tiên nghiệm):
 $P(E)$
- **Mô hình dịch - Translation Model (TM)**: câu đích trong tiếng Pháp (F) khi có câu tiếng Anh:
 $P(F|E)$
- Thủ tục tìm kiếm:
 - Cho F, tìm E tốt nhất sử dụng mô hình ngôn ngữ LM và mô hình dịch TM.
- Vấn đề: thiếu dữ liệu!
 - Ta không thể tạo từ điển câu $E \leftrightarrow F$
 - Thậm chí bình thường ta không thấy 1 câu lặp lại 2 lần

Ý tưởng giống hàng

- Mô hình dịch TM không quan tâm đến chuỗi đúng các từ tiếng Anh
- Sử dụng cách tiếp cận gán nhãn:
 - 1 từ tiếng Anh (“tag”) ~ 1 từ tiếng Pháp (“word”)
 - không thực tế: thậm chí số từ trong 2 câu không bằng nhau
 - sử dụng “giống hàng”.

Ý tưởng giống hàng

- Các tập ngữ liệu sử dụng giả thiết:
 - Dữ liệu song song (dịch $E \leftrightarrow F$)
- Giống hàng câu
 - Phát hiện câu
 - Giống hàng câu
- Giống hàng từ
 - Tách từ
 - Giống hàng từ (với 1 số ràng buộc)

Giống 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.