

CS 188: Artificial Intelligence

Language

Pieter Abbeel – UC Berkeley

Slides from Dan Klein

What is NLP?



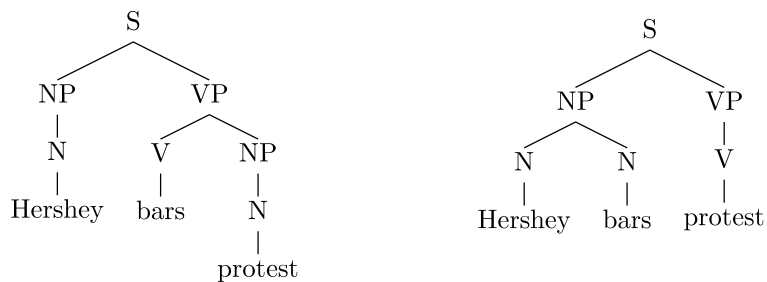
- **Fundamental goal: analyze and process human language, broadly, robustly, accurately...**
- **End systems that we want to build:**
 - Ambitious: speech recognition, machine translation, information extraction, dialog interfaces, question answering...
 - Modest: spelling correction, text categorization...

23

Problem: Ambiguities

- **Headlines:**
 - Enraged Cow Injures Farmer With Ax
 - Hospitals Are Sued by 7 Foot Doctors
 - Ban on Nude Dancing on Governor's Desk
 - Iraqi Head Seeks Arms
 - Local HS Dropouts Cut in Half
 - Juvenile Court to Try Shooting Defendant
 - Stolen Painting Found by Tree
 - Kids Make Nutritious Snacks
- **Why are these funny?**

Parsing as Search

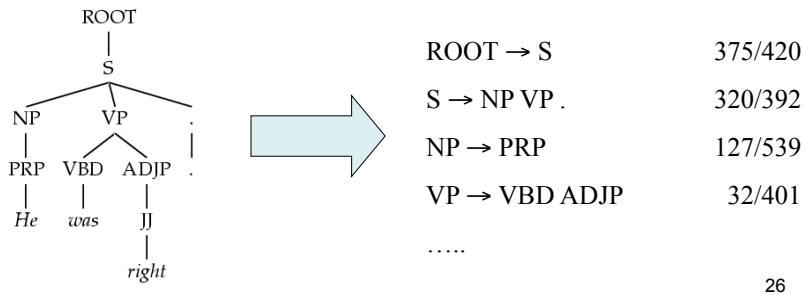


Hershey bars protest

25

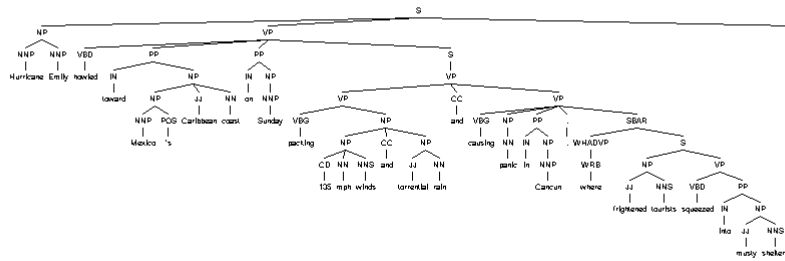
Grammar: PCFGs

- Natural language grammars are very ambiguous!
- PCFGs are a formal probabilistic model of trees
 - Each “rule” has a conditional probability (like an HMM)
 - Tree’s probability is the product of all rules used
- Parsing: Given a sentence, find the best tree – search!



26

Syntactic Analysis



Hurricane Emily howled toward Mexico 's Caribbean coast on Sunday packing 135 mph winds and torrential rain and causing panic in Cancun, where frightened tourists squeezed into musty shelters .

27

Machine Translation

"Il est impossible aux journalistes de rentrer dans les régions tibétaines"

Bruno Philip, correspondant du "Monde" en Chine, estime que les journalistes de l'AFP qui ont été expulsés de la province tibétaine du Qinghai "n'étaient pas dans l'illégalité".

Les faits Le dalaï-lama dénonce l'"enfer" imposé au Tibet depuis sa fuite, en 1959

Vidéo Anniversaire de la rébellion



"It is impossible for journalists to enter Tibetan areas"

Philip Bruno, correspondent for "World" in China, said that journalists of the AFP who have been deported from the Tibetan province of Qinghai "were not illegal."

Facts The Dalai Lama denounces the "hell" imposed since he fled Tibet in 1959

Video Anniversary of the Tibetan rebellion: China on guard



- Translate text from one language to another
- Recombines fragments of example translations
- Challenges:
 - What fragments? [learning to translate]
 - How to make efficient? [fast translation search]



The Problem with Dictionary Look-ups

| | |
|----|--|
| 顶部 | /top/roof/ |
| 顶端 | /summit/peak/ top /apex/ |
| 顶头 | /coming directly towards one/ top /end/ |
| 盖 | /lid/ top /cover/canopy/build/Gai/ |
| 盖帽 | /surpass/ top / |
| 极 | /extremely/pole/utmost/ top /collect/receive/ |
| 尖峰 | /peak/ top / |
| 面 | /fade/side/surface/aspect/ top /face/flour/ |
| 摘心 | / top /topping/ |

Example from Douglas Hofstadter

A Brief and Biased History



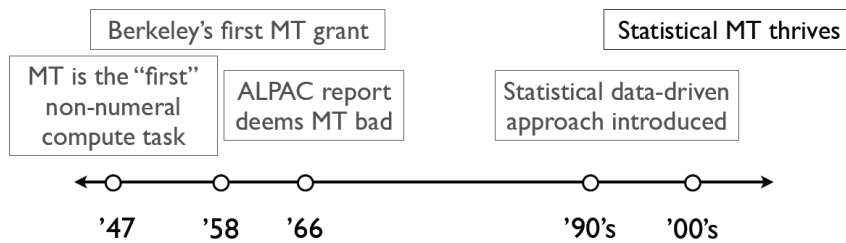
Warren Weaver

When I look at an article in Russian, I say: "This is really written in English, but it has been coded in some strange symbols. I will now proceed to decode."



John Pierce

"Machine Translation" presumably means going by algorithm from machine-readable source text to useful target text... In this context, there has been no machine translation...



Data-Driven Machine Translation

Target language corpus:

I will get to it soon

See you later

He will do it

Sentence-aligned parallel corpus:

Yo lo haré mañana
I will do it tomorrow

Hasta pronto
See you soon

Hasta pronto
See you around

Machine translation system:

Yo lo haré pronto

NOVEL SENTENCE

Model of translation

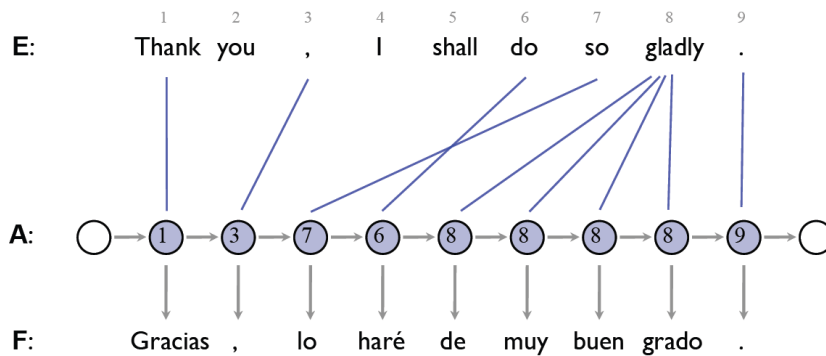
I will do it soon

Learning to Translate

| | | CLASSIC SOUPS | | Sm. | Lg. | |
|----|---|---------------|-------------------------------------|-------------------------------------|---|---|
| 清 | 燉 | 雞 | 湯 | 57. | House Chicken Soup (Chicken, Celery, Potato, Onion, Carrot) 1.50 2.75 | |
| 雞 | 飯 | 湯 | 58. | Chicken Rice Soup 1.85 3.25 | | |
| 雞 | 麵 | 湯 | 59. | Chicken Noodle Soup 1.85 3.25 | | |
| 廣 | 東 | 雲吞 | 湯 | 60. | Cantonese Wonton Soup 1.50 2.75 | |
| 蕃 | 茄 | 香 | 湯 | 61. | Tomato Clear Egg Drop Soup 1.65 2.95 | |
| 雲吞 | 湯 | 62. | Regular Wonton Soup 1.10 2.10 | | | |
| 酸 | 辣 | 湯 | 63. | Hot & Sour Soup 1.10 2.10 | | |
| 蛋 | 花 | 湯 | 64. | Egg Drop Soup 1.10 2.10 | | |
| 雲吞 | 湯 | 65. | Egg Drop Wonton Mix 1.10 2.10 | | | |
| 豆 | 腐 | 菜 | 湯 | 66. | Tofu Vegetable Soup NA 3.50 | |
| 雞 | 玉 | 米 | 湯 | 67. | Chicken Corn Cream Soup NA 3.50 | |
| 蟹 | 肉 | 玉 | 米 | 湯 | 68. | Crab Meat Corn Cream Soup NA 3.50 |
| 海 | 鮮 | 湯 | 69. | Seafood Soup NA 3.50 | | |

Example from Adam Lopez

The HMM Model

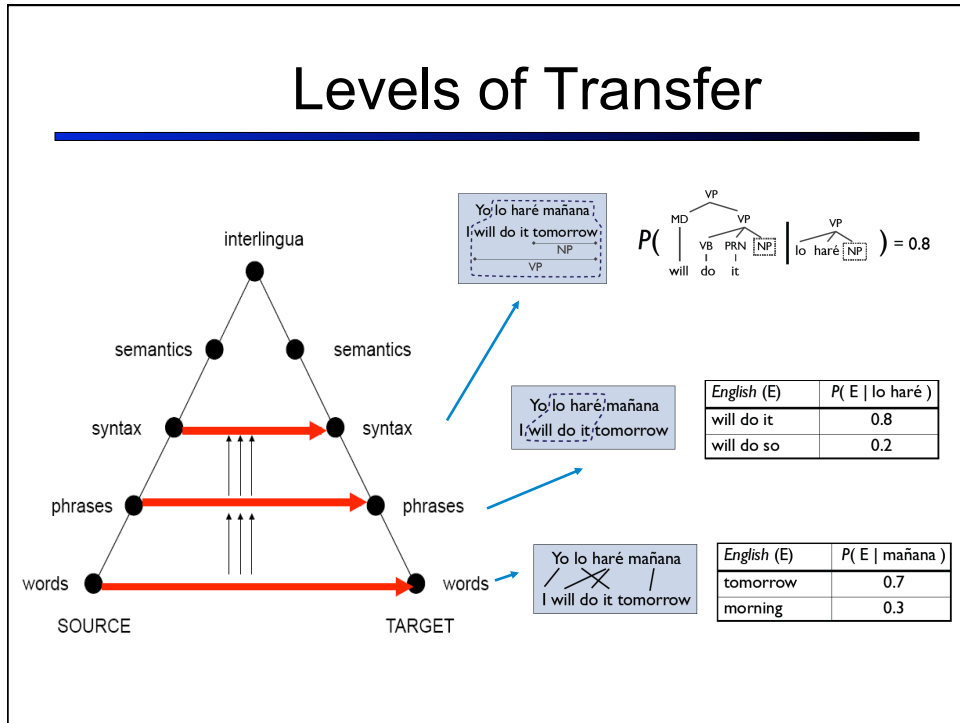


Model Parameters

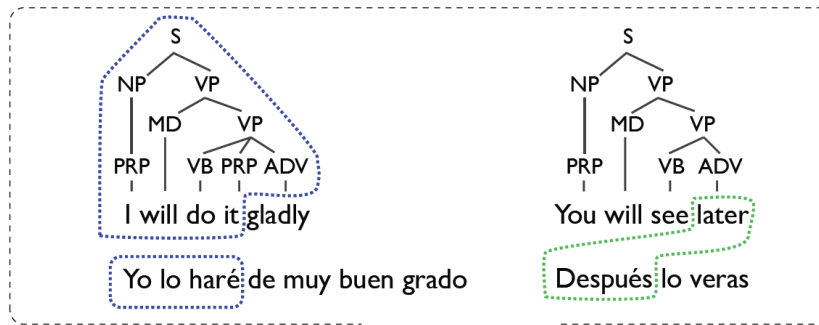
Emissions: $P(F_1 = \text{Gracias} \mid E_{A_1} = \text{Thank})$

Transitions: $P(A_2 = 3 \mid A_1 = 1)$

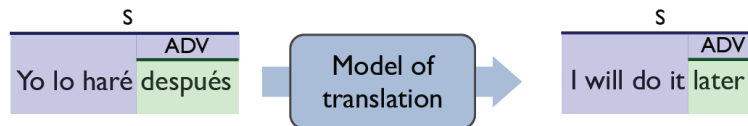
Levels of Transfer



Machine Translation

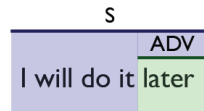
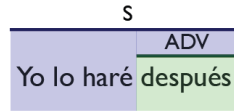


Machine translation system:



A Statistical Translation Model

Synchronous Derivation



Synchronous Grammar Rules

S → ⟨ Yo lo haré ADV ; I will do it ADV ⟩

ADV → ⟨ después ; later ⟩

A Statistical Model

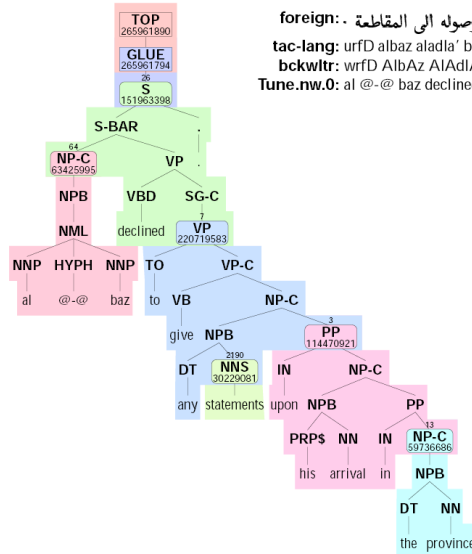
Translation model components
factor over applied rules

How well are these rules
supported by the data?

Language model factors over n-grams

How well is this output sentence
supported by the data?

Example Syntax-Based Translation



foreign: - ورفض الباز الادلاء باى تصريحات فور وصوله الى المقاطعة .
 tac-lang: urfD alBaz aladla' baá tSryHat fur uSulh alá almqaT'e .
 bckwlr: wrfD AlBaz AlAdlA' bAY tSryHAT fur wSwlh AlY AlmqATEp .
 Tune.nw.0: al @-@ baz declined to make any statements upon his arrival in the province .

[demo: MT]