

# Морфологічна розмітка

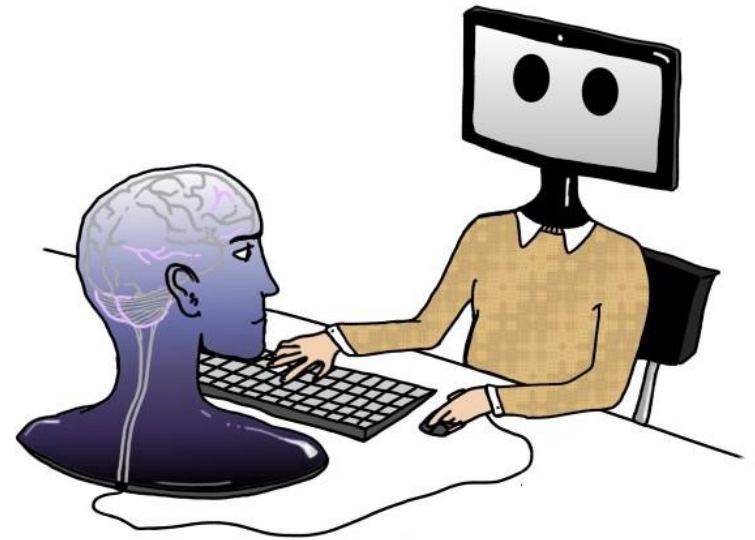
Чернівецький національний університет  
імені Юрія Федьковича

Наукова студентська конференція  
2016 рік

Доповідач: Дубінін Данило  
Науковий керівник: доц. Сопронюк Т.М.

# Обробка природних мов

- Загальний напрямок
  - Інформатики
  - Штучного інтелекту
  - Математичної лінгвістики
- Для взаємодії
  - Людина ↔ Машина ↔ Людина



# Завдання

- Машинний переклад
- Розпізнавання мови
- Видобування даних (аналіз тексту)
- Синтез мовлення
- Інформаційний пошук, витяг
- Реферування
- Розв'язання лексичної багатозначності

# Приклад Wolfram|Alpha

```
WolframAlpha["How many vehicles are there?", ExcludePods -> {"Input"}]|
```

Summary:

Show details

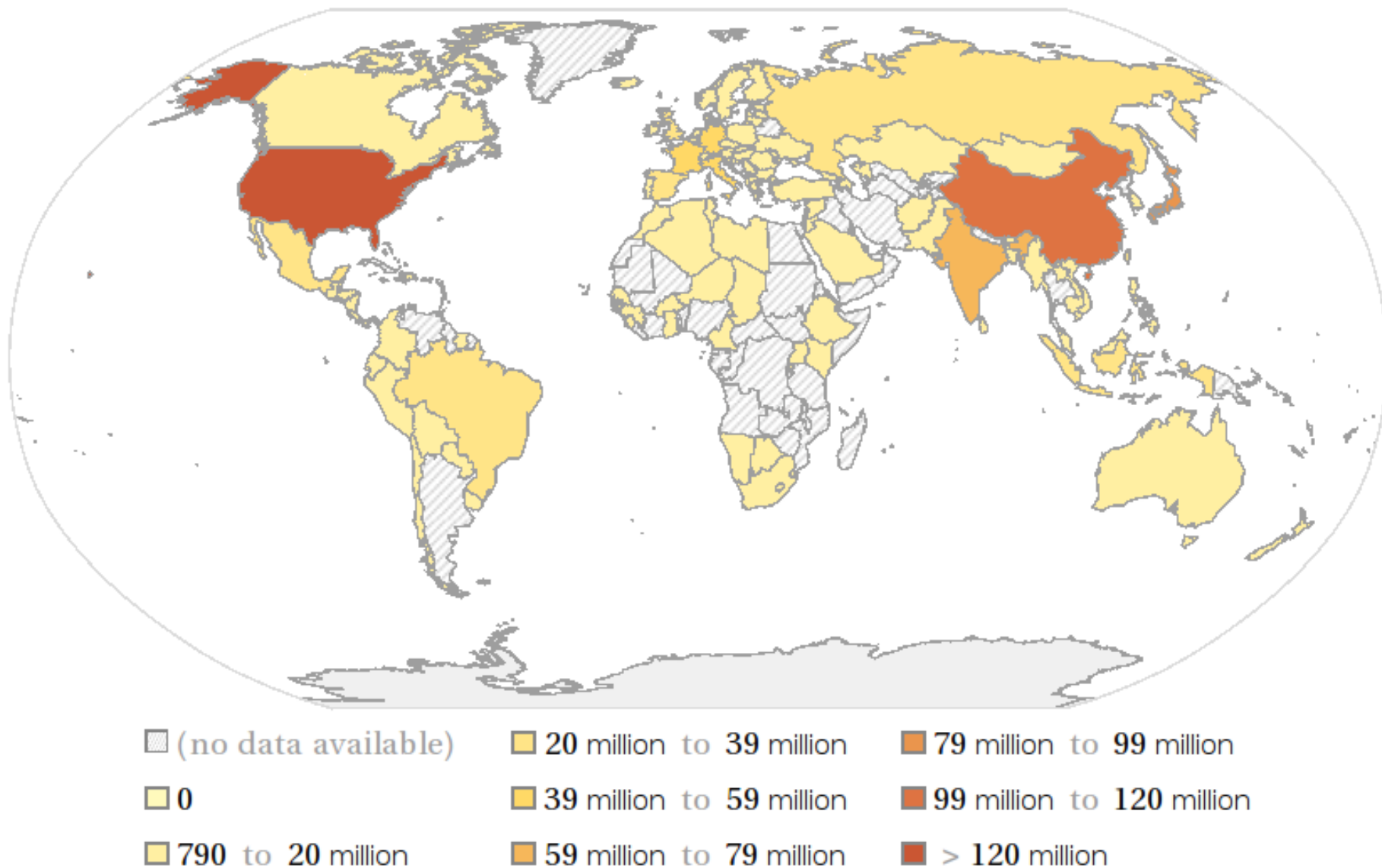


total	1.1 billion vehicles
median	1.1 million vehicles
highest	254.4 million vehicles (United States)
lowest	790 vehicles (Comoros)

(2001, 2002, 2003, 2004, 2005, 2006, 2007, and 2009 estimates)

(based on 123 values; 117 unavailable)

# Приклад Wolfram|Alpha



# Приклад Wolfram|Alpha

- **Make oxygen from CO<sub>2</sub>?**

CO<sub>2</sub> (carbon dioxide) + Na<sub>2</sub>O<sub>2</sub> (sodium peroxide) → O<sub>2</sub> (oxygen) + Na<sub>2</sub>CO<sub>3</sub> (soda ash)

---

CO<sub>2</sub> (carbon dioxide) + KO<sub>2</sub> (potassium superoxide) → O<sub>2</sub> (oxygen) + K<sub>2</sub>CO<sub>3</sub> (pearl ash)

---

CO<sub>2</sub> (carbon dioxide) + H<sub>2</sub>O (water) + KO<sub>2</sub> (potassium superoxide)  
→ O<sub>2</sub> (oxygen) + KHCO<sub>3</sub> (potassium bicarbonate)

---

CO<sub>2</sub> (carbon dioxide) + H<sub>2</sub>O<sub>2</sub> (hydrogen peroxide)  
→ O<sub>2</sub> (oxygen) + H<sub>2</sub>O (water) + CO (carbon monoxide)

---

CO<sub>2</sub> (carbon dioxide) + HOOCCH<sub>2</sub>CH<sub>2</sub>COOH (succinic acid)  
+ C<sub>14</sub>H<sub>19</sub>N<sub>3</sub>O<sub>7</sub>S (Deacetylcephalosporin C) →  
O<sub>2</sub> (oxygen) + C<sub>5</sub>H<sub>6</sub>O<sub>5</sub> (2-oxoglutaric acid) + C<sub>14</sub>H<sub>19</sub>N<sub>3</sub>O<sub>6</sub>S (deacetoxcephalosporin-c)

# Машинний переклад

- Методи
  - Правил
  - Статистичні
  - Гібридні
- Найскладніші проблеми
  - Лексична неоднозначність
  - Власні назви

# Лексична багатозначність



Гриф





# Видобування даних

- Виявлення закономірностей у великих необроблених масивах даних



# Історія

- I (1940-1960)
  - Georgetown-IBM experiment (1954)
  - Russian → English
- II (1960-1970)
  - SHRDLU
- III (1970-1980)
  - Штучний інтелект
  - Логічне програмування
- IV (1990-)
  - Статистичні методи

# I. Georgetown-IBM experiment

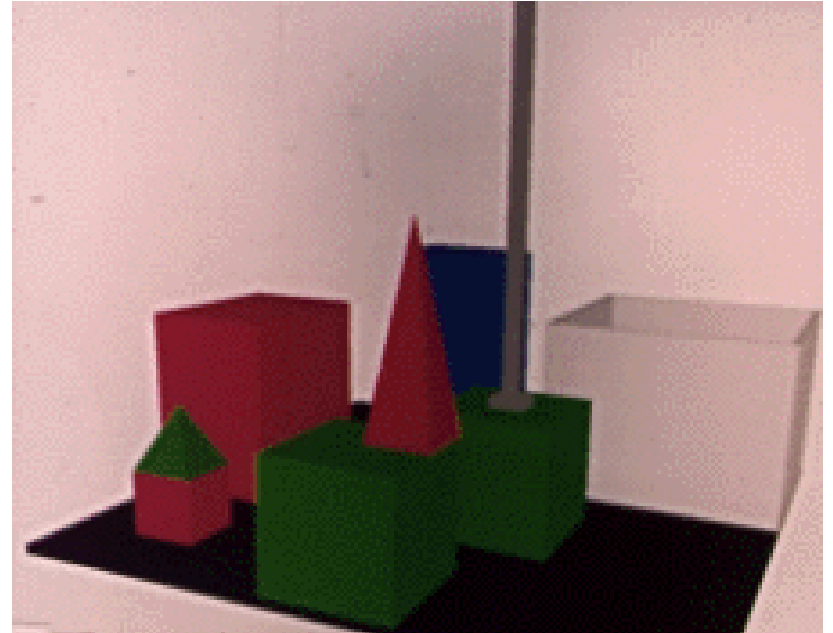
- Перший машинний переклад оснований на правилах
  - Перестановка слів
  - Вставка/вилучення слів
  - І т.п.

`Mi pyeryedayem mislyi posryedstvom ryechyi.`

`We transmit thoughts by means of speech.`

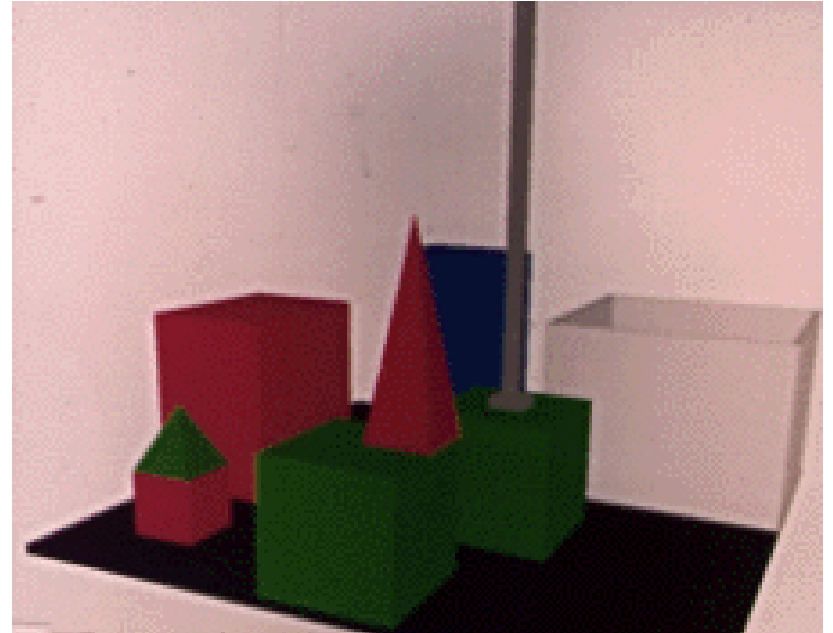
## II. SHRDLU

- P: Pick up a big red block.
- C: OK.
- P: Grasp the pyramid.
- C: I DON'T UNDERSTAND WHICH PYRAMID YOU MEAN.



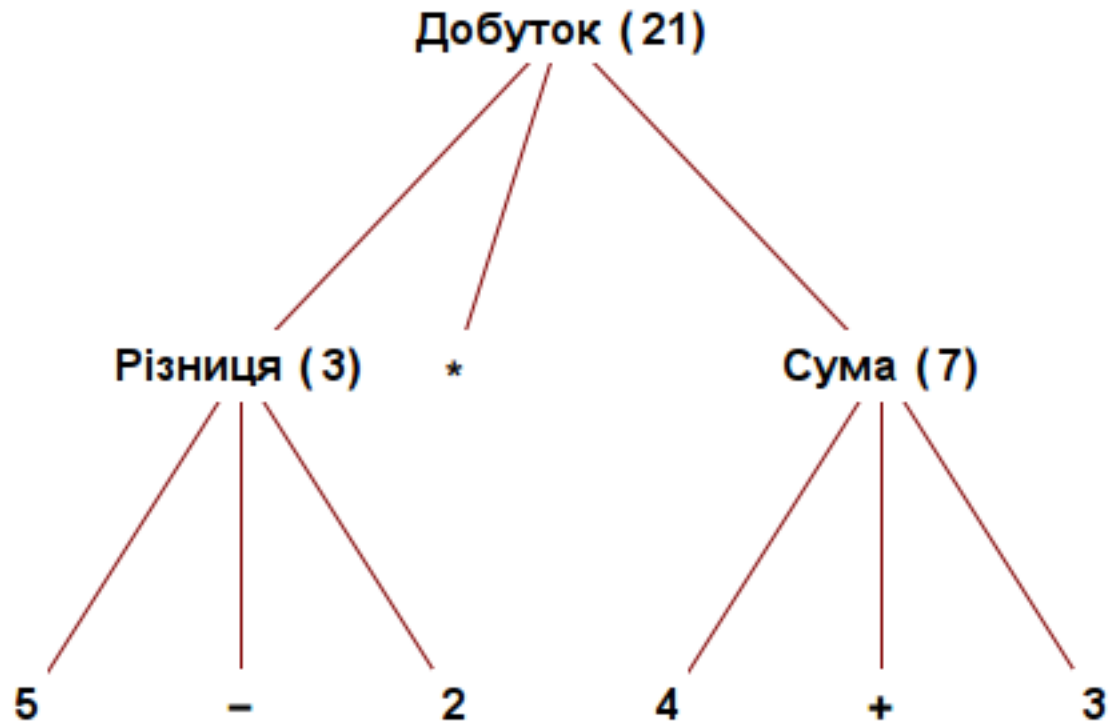
## II. SHRDLU

- P: How many blocks are not in the box?
- C: FOUR OF THEM.
- P: Is at least one of them narrower than the one which i told you to pick up?
- C: YES, THE RED CUBE.



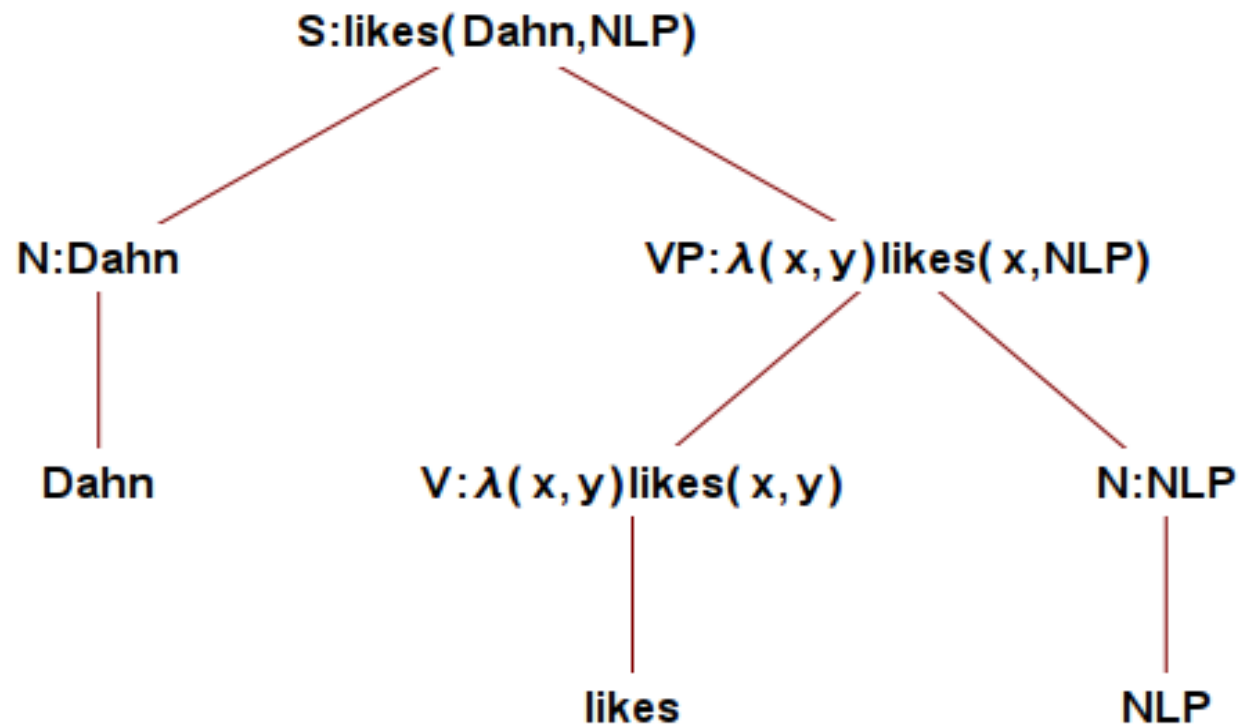
# III. Синтаксичний розбір

- Матиматичні вирази? – Просто!
  - $(5-2)*(4+3)$



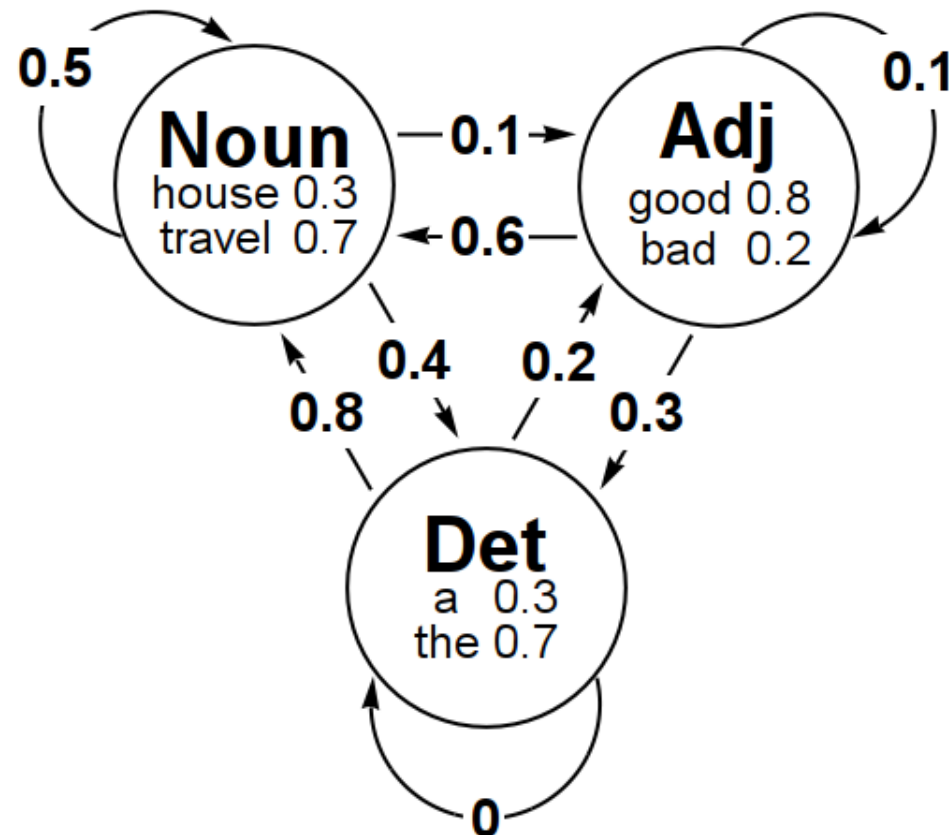
# III. Синтаксичний розбір

- Речення природної мови?
- Предикативне числення!



# IV. Статистичні методи

- Стохастичні контекстно-залежні граматики
- Прихована марковська модель (ППМ)





# ШІ-повні задачі (AI-complete)

- Комп'ютерний зір
- Розуміння природної мови
- Розв'язування реальних проблем за непередбачуваних обставин

# ШІ-повні задачі (AI-complete)

- Комп'ютерний зір
- Розуміння природної мови
- Розв'язування реальних проблем за непередбачуваних обставин (навігація, планування, експертні системи, і т.п.)

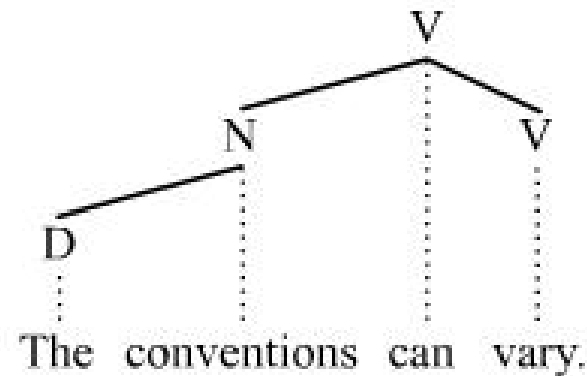


# Морфологічна розмітка

- Дано речення: послідовність слів ***x***
- Визначення частини мови слів ***y***
  - ADJ, ADP, ADV, CONJ, DET, NOUN, NUM, PRON, PRT, VERB
  - . – пунктуація
  - <S>, </S> – початок/кінець речення
  - X – невідома частина мови

# Навіщо

- Розв'язання лексичної неоднозначності
  - Покращення результату використовуючи оточуючі частини мови як характеристики
- Переклад
  - Додаткове знання частин мови не заважає
- Синтаксичний розбір



# Схема

- Обчислення параметрів ППМ
  - Машинне навчання
- Алгоритм Вітербі
  - Динамічне програмування

# Навчання

- *Лінгвістичний корпус: Brown corpus*
  - складений у 1960-х
  - 500 текстів
  - 1 млн. слів
  - людьми визначенні частини мови всіх слів
  - Once/ADP the/DET principle/NOUN was/VERB established/VERB ,/. the/DET increase/NOUN in/ADP state-owned/ADJ vehicles/NOUN came/VERB rapidly/ADV ./.



# Прихована марковська модель

- ПММ – марковський процес із *неспостережуваними станами*.
- *Спостереження*
  - Слова речення
  - **The quick brown fox jumps over the lazy dog.**
- *Стани*
  - Частина мови
  - **DET ADJ ADJ NOUN VERB PREP DET ADJ NOUN**

# ПММ морфологічної розмітки

- Дано ланцюжки пар:  $\langle \text{слово}, \text{частина мови} \rangle$ 
  - $(x)_i, (y)_i$
- $f : X \rightarrow Y$  ?
- Вивчаємо розподіл  $\mathbb{P}(y|x)$
- Знаходимо розв'язок у вигляді

$$f(x) = \operatorname{argmax}_y \mathbb{P}(y|x)$$



# ПММ морфологічної розмітки

$$y = \operatorname{argmax}_y$$

$$\underbrace{\prod_i \mathbb{P}(x_i \mid y_i)}_{\text{спостереження}}$$

Слова

$$\times \underbrace{\prod_i \mathbb{P}(y_i \mid y_{i-1})}_{\text{стани}}$$

Частини мови

# Алгоритм Вітербі

- Таблиця динамічного програмування

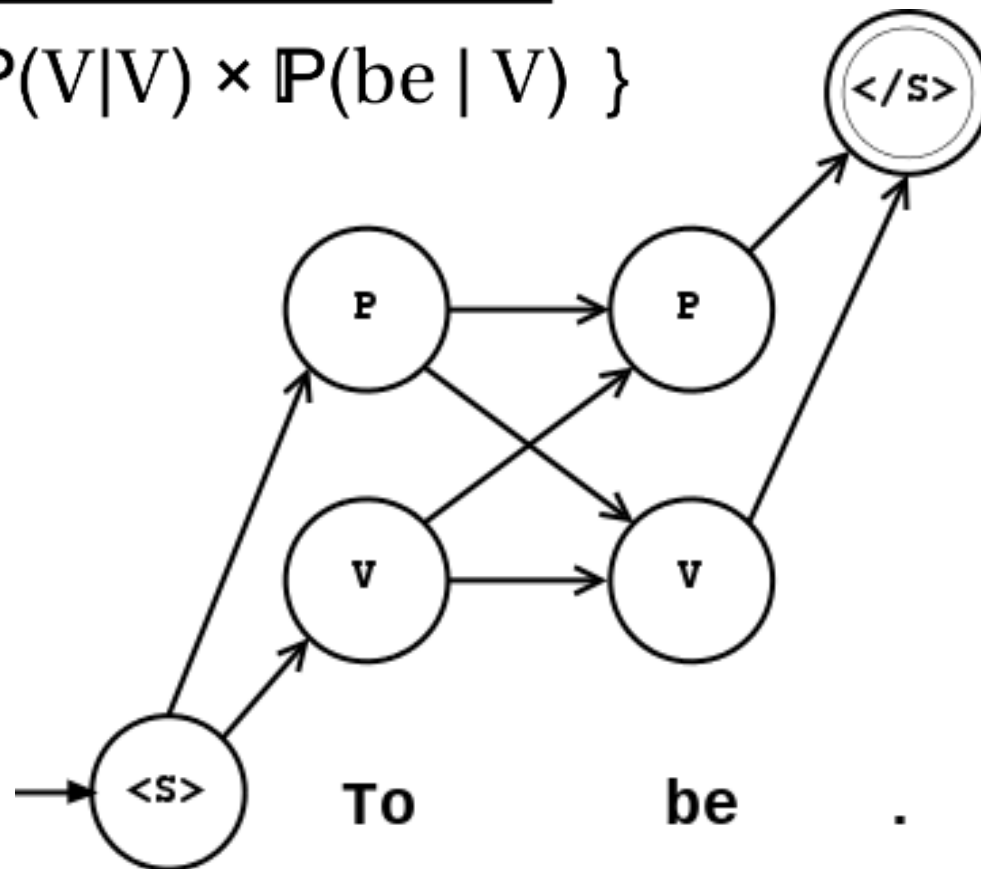
$$\pi_k(A,B) = \max_{\substack{\langle y_1 \dots y_k \rangle : \\ y_{k-1}=A, y_k=B}} \prod_{i=1}^k \mathbb{P}(x_i | y_i) \mathbb{P}(y_i | y_{i-1})$$

- Рекурсивне обчислення

$$\pi_k(Y,Z) = \max_X ( \pi_{k-1}(X,Y) \times \mathbb{P}(Z | X,Y) \times \mathbb{P}(x_k | Z) ) \quad \pi_0 = 1$$

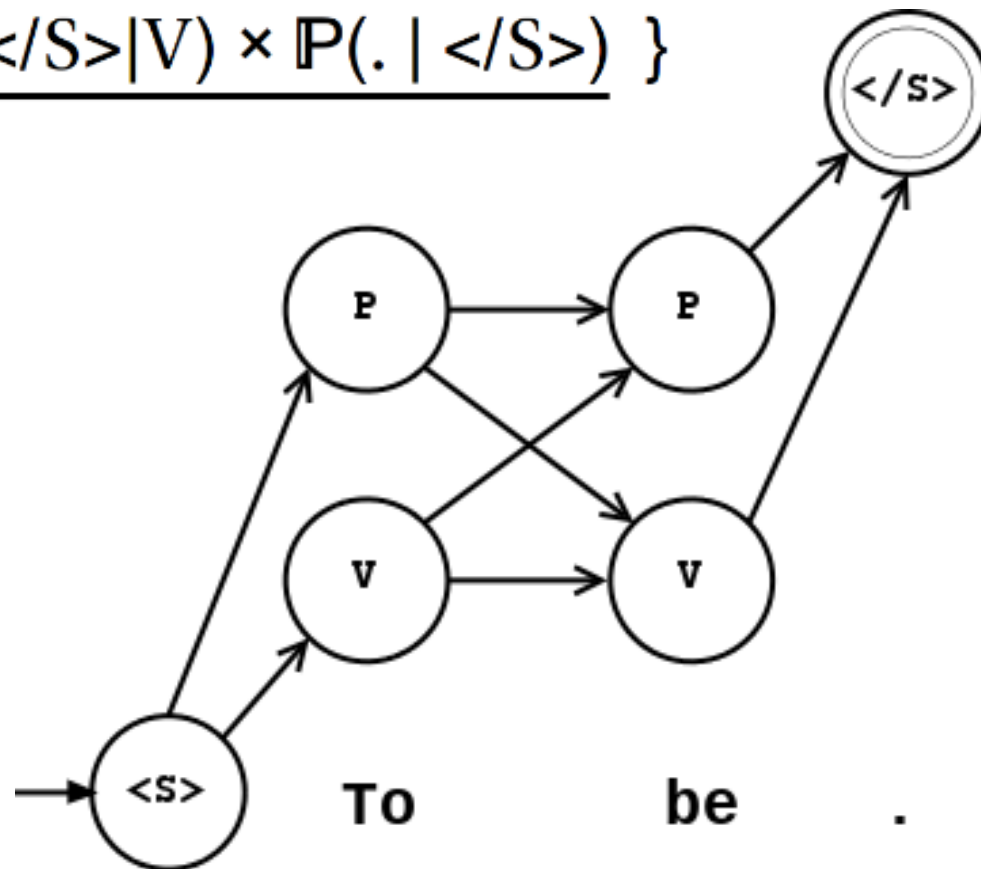
# Приклад

$$\pi_2(V) = \max \left\{ \frac{\pi_1(P) \times \mathbb{P}(V|P) \times \mathbb{P}(\text{be} | V),}{\pi_1(V) \times \mathbb{P}(V|V) \times \mathbb{P}(\text{be} | V)} \right\}$$



# Приклад

$$\pi_3(</S>) = \max \{$$
$$\pi_2(P) \times \mathbb{P}(</S>|P) \times \mathbb{P}(. \mid </S>),$$
$$\underline{\pi_2(V) \times \mathbb{P}(</S>|V) \times \mathbb{P}(. \mid </S>) } \}$$



# Точність

- Власна реалізація
  - ~ 93%
- Межа професійних реалізацій
  - ~ 97% (2014 рік)
- Лінгвісти
  - ~ 98%



MAX-PLANCK-GESELLSCHAFT



У дїї

Sentence > I've given my great presentation for scientific student conference in CHNU.

Tagged : I/PRON 've/VERB given/VERB my/DET great/ADJ presentation/NOUN for/ADP scientific/ADJ student/NOUN conference/NOUN in/ADP CHNU/NOUN ./.

У дїї

Sentence > I love those exciting  
travels we had last year!

Tagged : I/PRON love/VERB those/DET  
exciting/ADJ travels/NOUN we/PRON  
had/VERB last/ADJ year/NOUN !/.

У дії

Sentence > My best friend travels a lot.

Tagged : My/DET best/ADJ  
friend/NOUN travels/VERB a/DET  
lot/NOUN ./.



# Але статистика не завжди права

чорний кіт тікає від білої собаки

Translate

[Ukrainian](#) > [English](#)

black cat runs away from a white dog

чорний кіт тікає від білого собаки

Translate

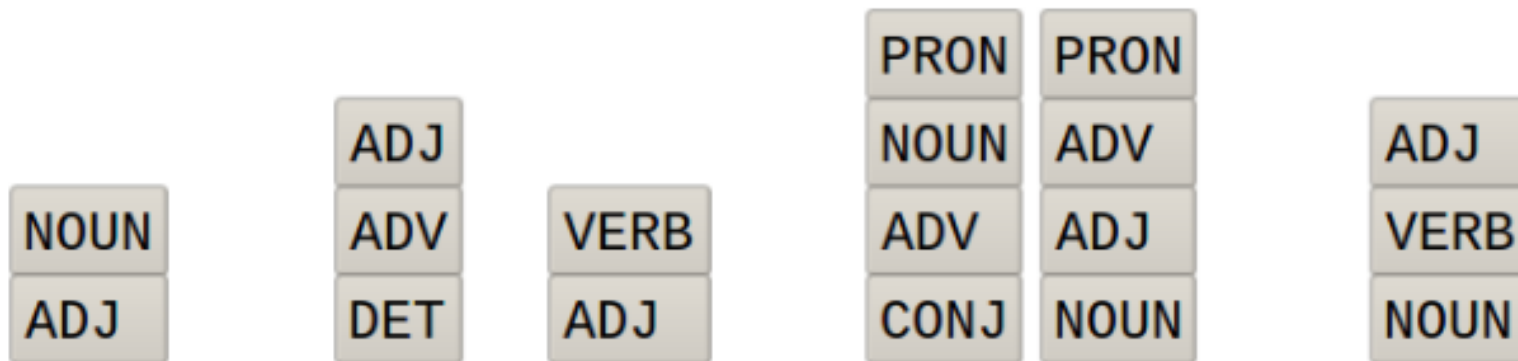
[Ukrainian](#) > [English](#)

black cat runs away from white dogs

# Висновок

- Математичний підхід до лінгвістики!
- Прихована марковська модель ефективна
  - Для легких підзадач
- Для складніших задач питання відкриті
  - Поєднання математики та лінгвістики?
  - Використання більших навчальних баз та обчислюваних можливостей?

# Тестування?



*Perfection* is *finally attained* not *when there* is no *longer* anything to add,  
but when there is no longer anything to take away.

Дякую :)



Perfection/NOUN is/VERB finally/ADV  
attained/VERB not/ADV when/ADV there/PRT  
is/VERB no/DET longer/ADJ anything/NOUN  
to/PRT add/VERB ,/. but/CONJ when/ADV  
there/PRT is/VERB no/DET longer/ADJ  
anything/NOUN to/PRT take/VERB away/ADV ./.