# Automatic Identification of Figurative Language

Presented by: Saša Milić
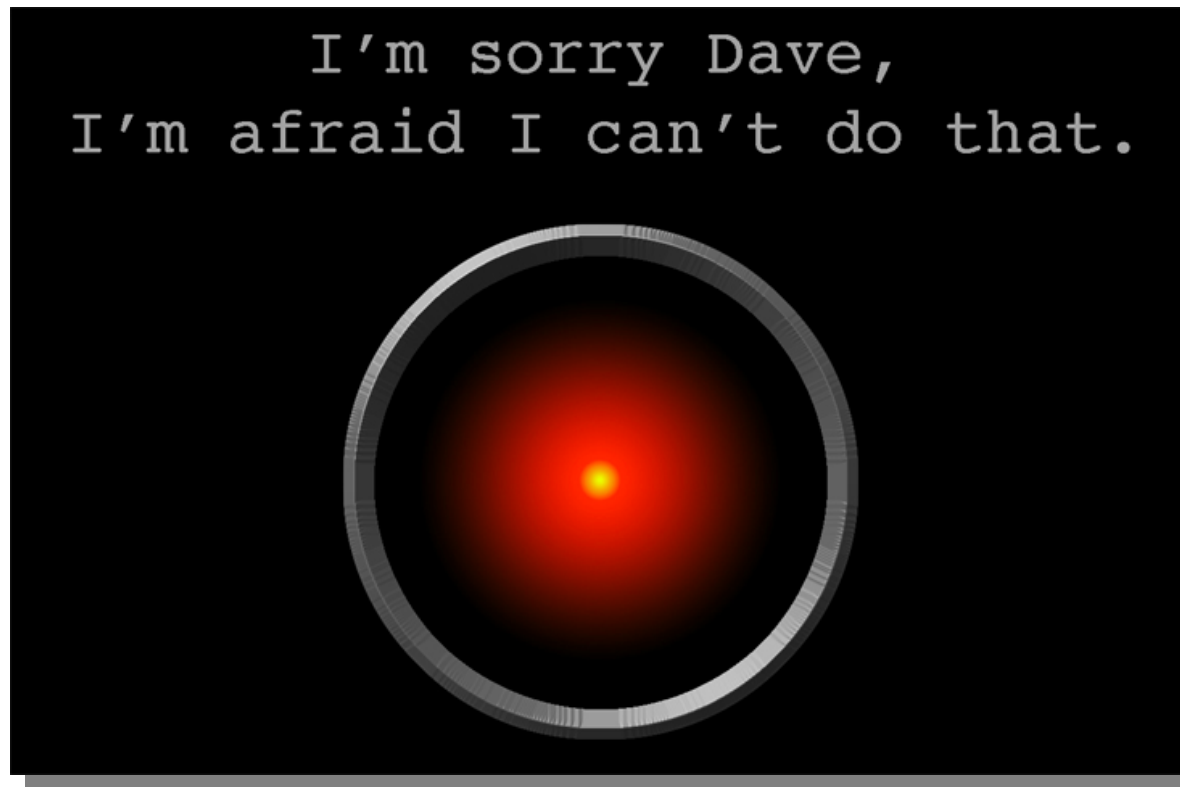
*Special thanks to:*
*My supervisor Suzanne Stevenson*
*and my mentor Afsaneh Fazly*
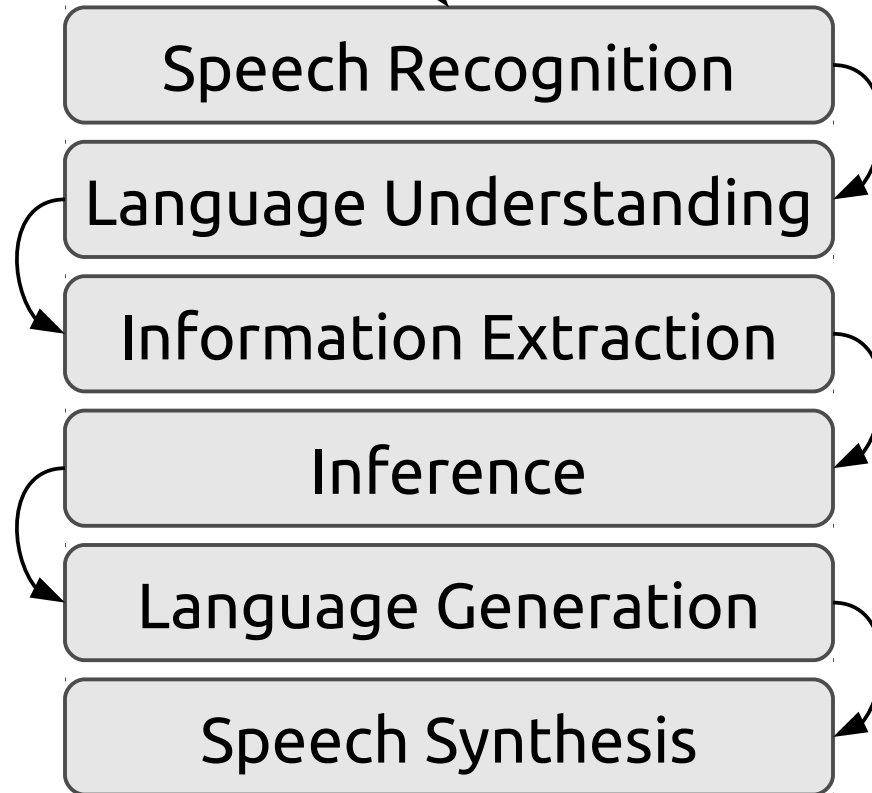
*June 19, 2013*

# A Conversant Computer?

UNIVERSITY OF TORONTO
FACULTY OF ARTS & SCIENCE

# The Minimal Requirements

# Computational Linguistics (CL)

**Understand**

- acquisition
- comprehension
- production

of *human language* from a *computational* perspective

**Apply**

focus on *practical outcomes* of modeling human language

UNIVERSITY OF TORONTO
FACULTY OF ARTS & SCIENCE

# Applications of CL

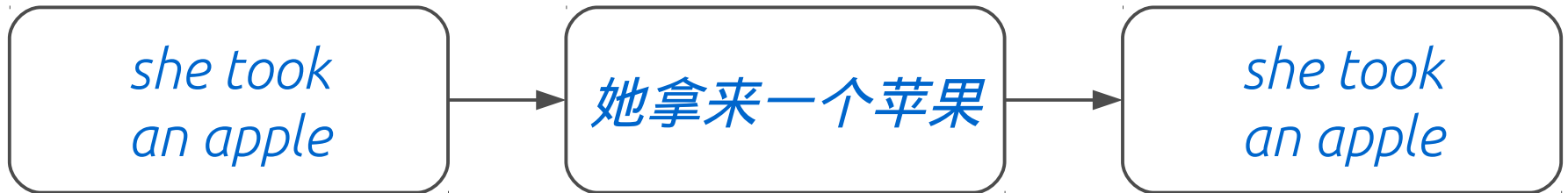- Grammar and style checking

- Apple's *Siri*

- Search Engine

- Machine translation

# Google Translate : An Informal Experiment

- Translating a <u>literal phrase</u>:

| *she took an apple* | → | 她拿来一个苹果 | → | *she took an apple* |

- Translating a <u>multiword expression</u>:

| *she took a walk* | → | 她散步了 | → | *she walks up* |

# Difficulty with Multiword Expressions

- Multiword expression:
  - two or more words that together form a *single unit* of meaning
    - *"frying pan"*
    - *"keep an eye out for"*
    - *"shoot the breeze"*

- overall meaning ≠ sum of the meaning of the components

# Light Verb Construction (LVC)

- A multiword expression (in our case, `verb` + `noun`) where the `noun` determines the primary meaning of the whole

| LVC | "give a sigh" | "make a decision" | "take a walk" |
|---------|-----------------|---------------------|-------------------|
| Literal | "give a present" | "make a cake" | "take an apple" |

- *Again:*
  - overall meaning ≠ sum of the meaning of the components

- *However:*
  - the component meanings still contribute something to the overall meaning

UNIVERSITY OF TORONTO
FACULTY OF ARTS & SCIENCE

# Identifying LVCs

- Which of the following is a light verb construction?

  - *He gave a donation.*

  - *It took place over there.*

  - *He gave her an advantage.*

- Motivates the question: can we do better than a simple binary classification?

UNIVERSITY OF TORONTO
FACULTY OF ARTS & SCIENCE

# A More Appropriate Measure

- Binary decision-making vs graded decision-making
  - *"Is this an LVC?"* vs *"How acceptable is this as an LVC?"*

- More formally:
  - What is the probability that some verb + noun combination forms an LVC?

- New measure: **Acceptability**

# Measuring Acceptability

- Linguistic studies suggest that a measure of LVC acceptability should incorporate both **frequency** and **semantic similarity**.

- **Hypothesis**:
  - a *novel* LV + noun is considered more acceptable if the noun is **similar** to a noun in a **high-frequency** LVC

- Example:
  - How acceptable is "*take a saunter*"?

UNIVERSITY OF TORONTO
FACULTY OF ARTS & SCIENCE

# "take a saunter"

take + {
- stroll, hike, walk ...
- shower, bath, wash, ...
- apple, banana, durian, ...
- ...
}

C( take ) = { ⬤, ⬤, ⬤, ⬤, ... }

C( v ): set of semantic classes of nouns that can occur with verb v

# *"take a saunter"*

$$P(\text{ saunter belongs to } \boxed{\begin{array}{c} \textit{stroll,} \\ \textit{hike, walk,} \\ \textit{...} \end{array}} ) = ?$$

$P(\text{saunter} \mid \bigcirc) = \text{high}$

$P(\text{n} \mid \text{c})$: probability that noun **n** belongs to class **c**

# *"take a saunter"*

$$P\left(\text{take} + \begin{array}{c} stroll, \\ hike, walk, \\ ... \end{array} = LVC \right) = ?$$

$$P_{LVC}\left( \bigcirc \mid \text{take} \right) = \text{high}$$

$P_{LVC}(c \mid v):$  probability that class $c$ forms

acceptable LVCs with $v$

# Measuring Acceptability

- Acceptability:
  - A *probabilistic* measure

- Components
  - $C(v)$
  - $P(n|c)$
  - $P_{LVC}(c|v)$

# Estimating Probabilities

▫ We can't know the true probabilities.  So we estimate.

▫ In order to estimate $P_{LVC}(c|v)$ we need to know:

– $P_{LVC}(n|v)$
  • for all $n$ in class $c$

– Estimate **directly**
  • *Why can't we do this for novel LVCs?*
– Estimate **indirectly**

UNIVERSITY OF TORONTO
FACULTY OF ARTS & SCIENCE

# Estimating Probabilities

- We use a machine learning algorithm to estimate this *directly* for frequent combinations :
  - $P_{LVC}(n \mid v)$

- Using ~25 features drawing on linguistic properties of LVCs
  - Examples:
    - frequencies
    - association
    - syntactic behavior

UNIVERSITY OF TORONTO
FACULTY OF ARTS & SCIENCE

# Some Features of LVCs

- We expect the noun and the verb in an LVC to have strong associativity

- We expect LVCs to have a preference for indefinite determiners (*"a"*, *"an"*, …)
  - consider:
    - *"make <u>a</u> speech"* **vs** *"make <u>the</u> speech"*
  - Which one occurs more often?
    - ~16 million vs ~2 million Google hits

UNIVERSITY OF TORONTO
FACULTY OF ARTS & SCIENCE

# Evaluation

- Obtain human ratings (on some scale) of LVC acceptability

- **Goals**:
  - to introduce a more appropriate (*linguistically-motivated*) measure for identifying LVCs
  - to be able to predict LVC acceptability of novel expressions