# Deep-Parsing Natural Languages

## 深度语法分析是自然语言应用的核武器

ArchSummit 全球架构师峰会 2017 (北京)

**李维（京东硅谷研究院）**
12/08/2017 (liweinlp.com)
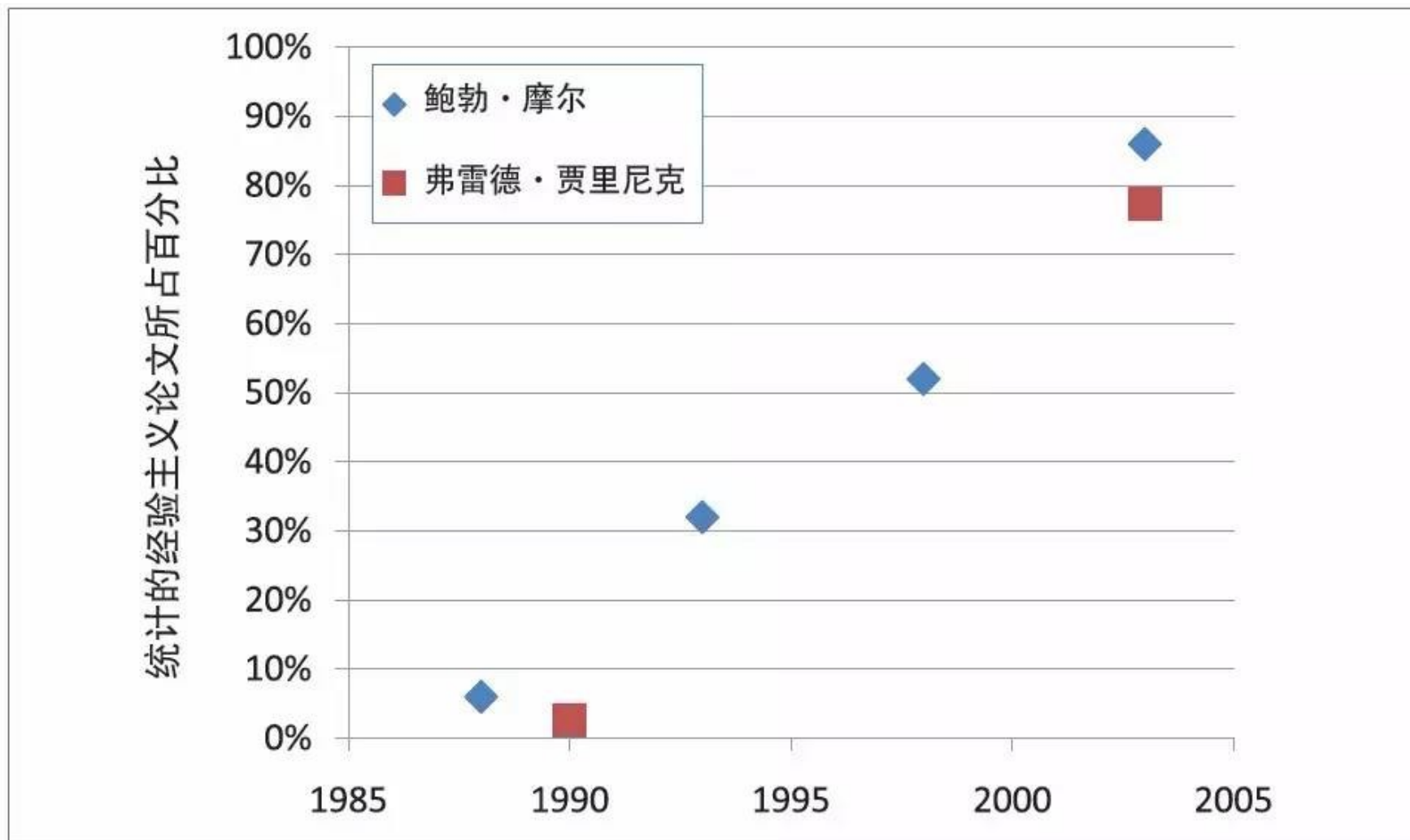
京东

# Outline

- 人工智能的历史和现状简介：从感知到认知
  - 此消彼长的经验主义和理性主义钟摆

- 深度解析（Deep Parsing）是什么?

- NLP 架构纵览

- 核武器应用举例

# Outline: AI History

- 人工智能的历史和现状简介：从感知到认知
  - 此消彼长的经验主义和理性主义钟摆

# NLP Mainstream Since 1990s



Courtesy of Prof. Church: **"A Pendulum Swung Too Far"**
http://blog.sciencenet.cn/blog-362400-988692.html

# Two Basic Approaches to NLP

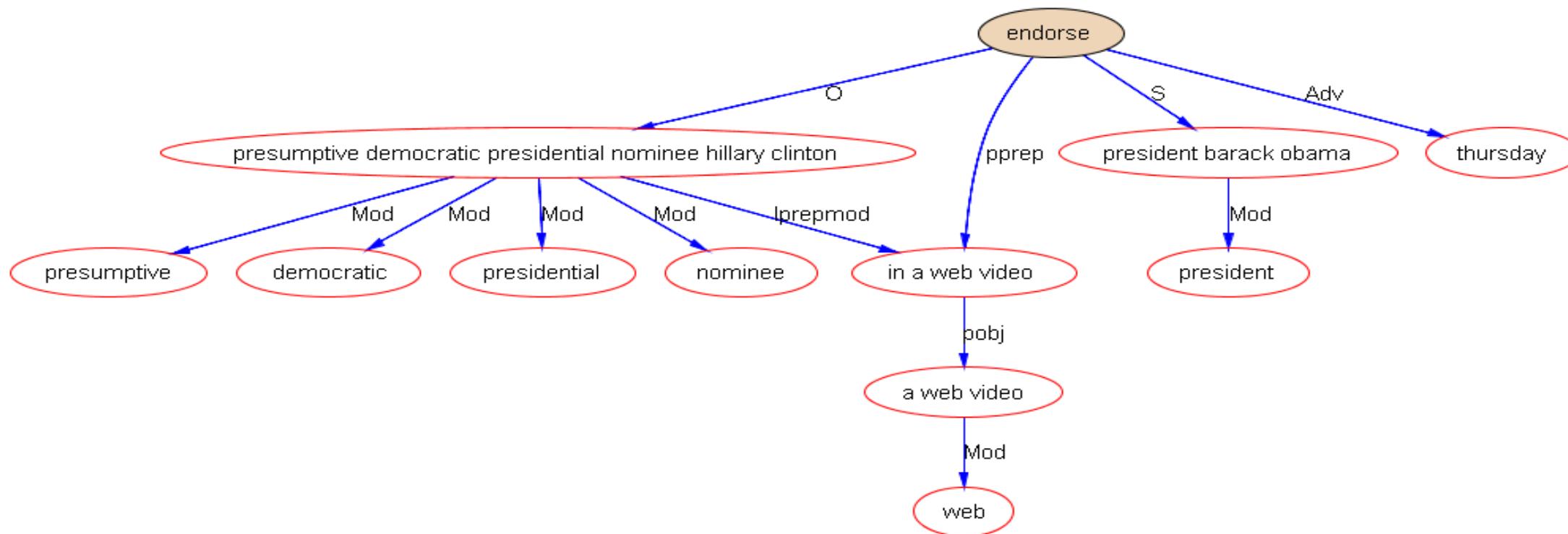| Approach | Pros | Cons |
|---|---|---|
| **Statistical Learning** (based on keywords) | • Good for document-level<br>• High recall<br>• Robust<br>• Easy to scale<br>• Fast development<br>  (if data available) | • Requires large annotation<br>• Coarse-grained<br>• Difficult to debug<br>• Fail in short messages<br>• Only shallow NLP<br>• No understanding |
| **Grammar Engineering** (based on sentence structure) | • Good for sentence level<br>• Handles short messages well<br>• High precision<br>• Fine-grained insights<br>• Easy to debug<br>• Parsing and understanding | • Requires deep skills<br>• Requires scale up skills<br>• Requires robustness skills<br>• Moderate recall (coverage)<br>• Parser development slow |

- Complementary rather than competing
- Hybrid: Best of both worlds
- Balance and configurability between precision and recall
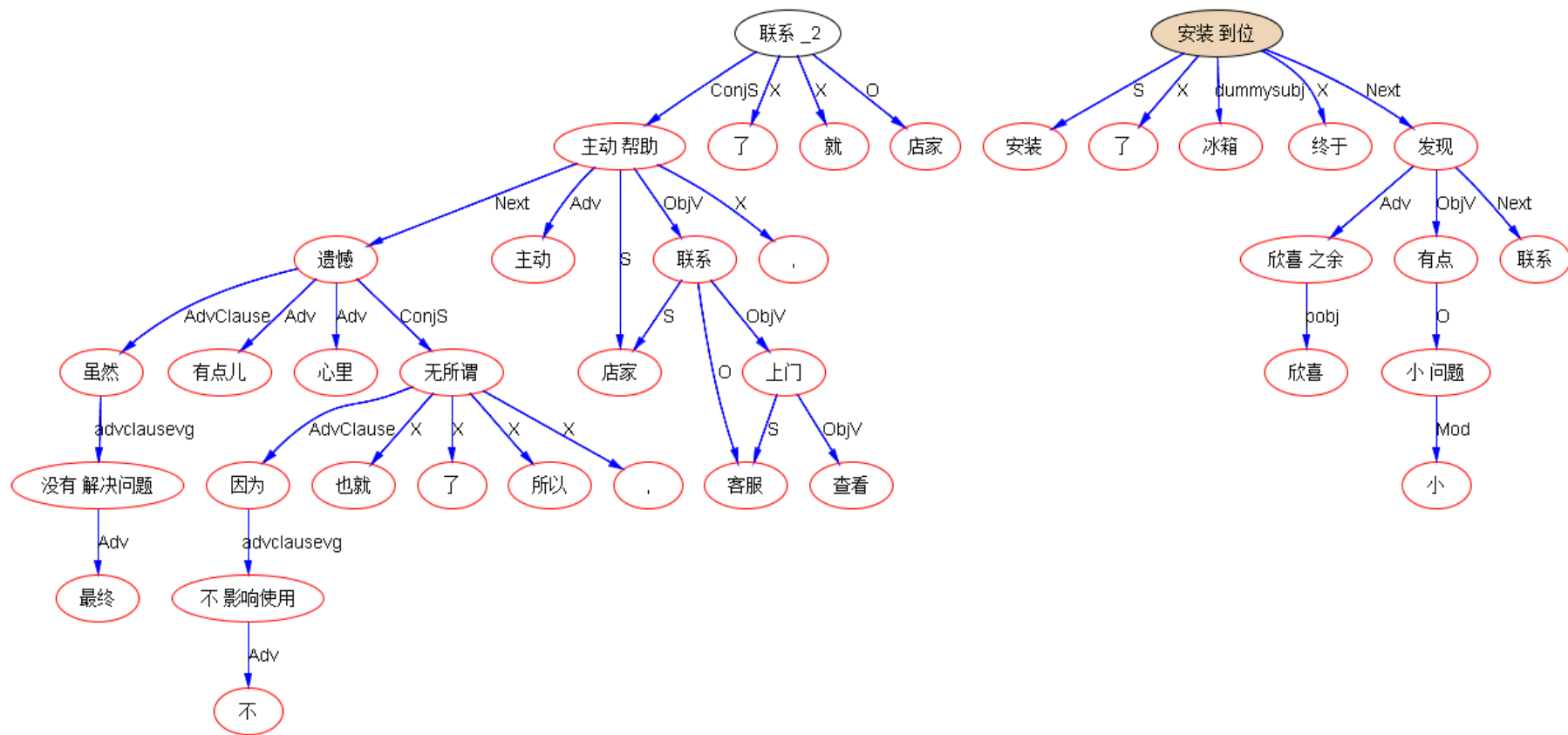
NLP频道
liweinlp.com

- 人工智能的历史和现状简介：从感知到认知
    此消彼长的经验主义和理性主义钟摆

- 深度解析（Deep Parsing）是什么？

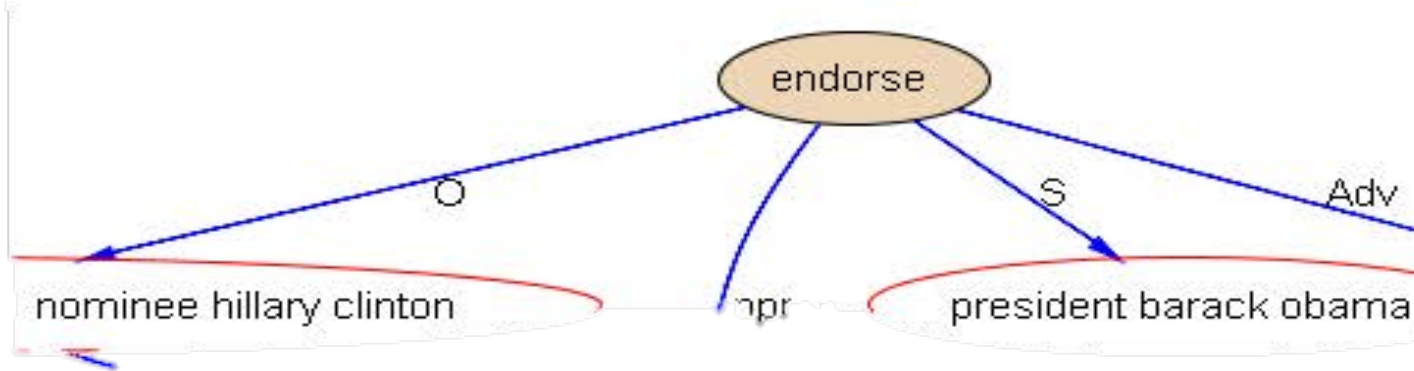# Deep Parsing: Unstructured to Structures



1. Input: President Barack Obama endorsed presumptive Democratic presidential nominee Hillary Clinton in a web video Thursday .

# Why parsing?  Limited Patterns



终于冰箱安装到位了，欣喜之余发现有点儿小问题，就联系了店家，店家主动帮助联系客服上门查看，虽然最终没有解决问题，心里有点儿遗憾，但是因为不影响使用，所以也就无所谓了.

NLP频道
liweinlp.com

# Subtree Pattern: Data to Intelligence 京东

endorse

nominee hillary clinton    npr    president barack obama
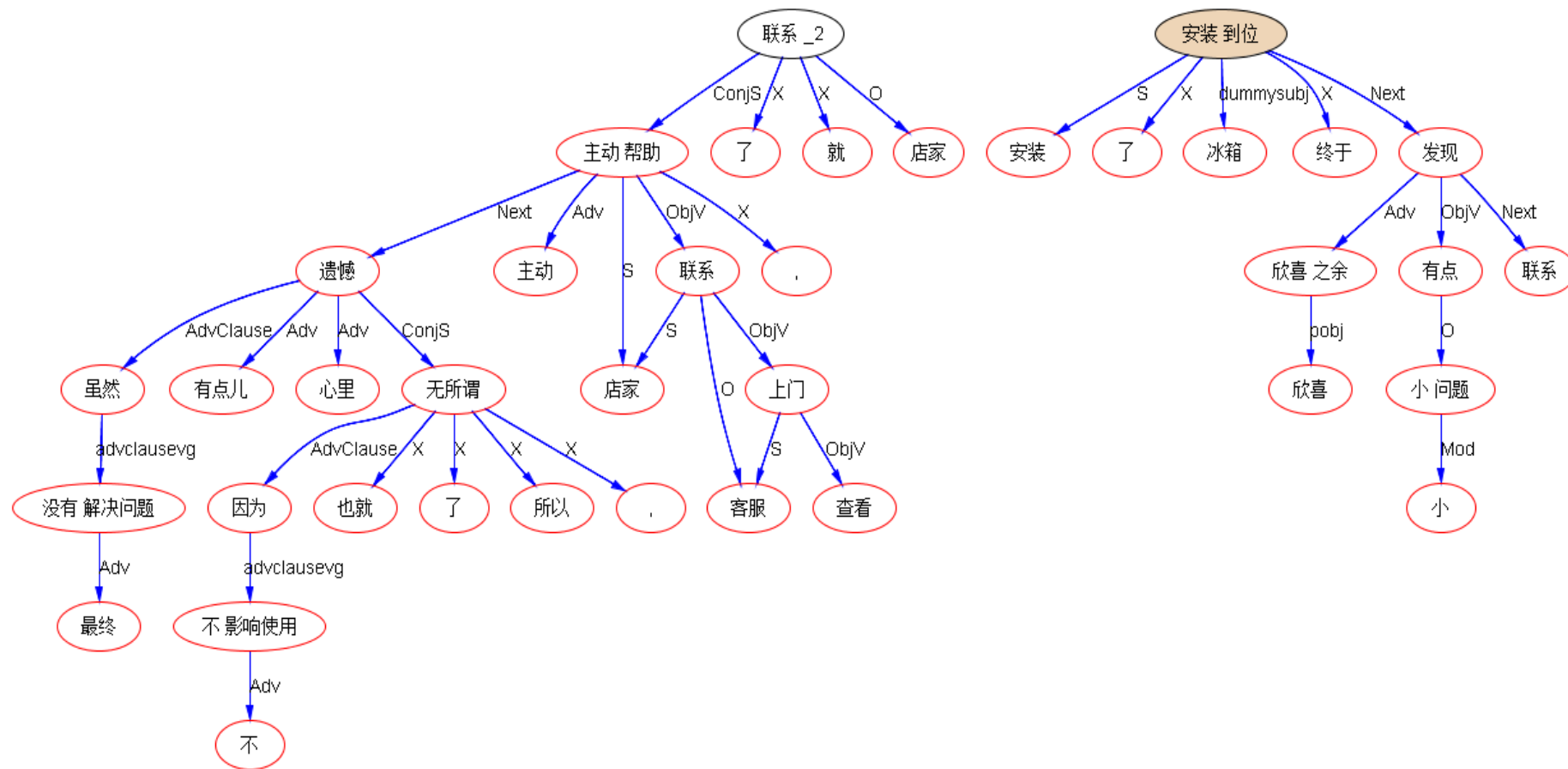
O          S          Adv

SVO Pattern:
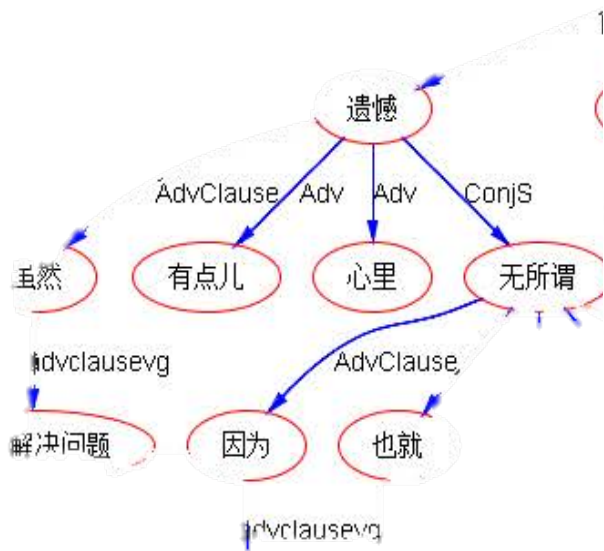
Barack Obama (S)
Endorse (V)
Hillary Clinton (O)

Knowledge Graph

# Deep Parsing: Unstructured to Structures 京东



终于冰箱安装到位了,欣喜之余发现有点儿小问题,就联系了店家,店家主动帮助联系客服上门查看,虽然最终没有解决问题,心里有点儿遗憾,但是因为不影响使用,所以也就无所谓了.

# Subtree Pattern: Data to Intelligence

京东



Inter-Clause Pattern:
虽然 ... 遗憾...无所谓...

mild sentiment

Linear: Infinite number of sentences

Structure: Limited patterns

Data → Intelligence

# Outline: NLP Architectures

- 人工智能的历史和现状简介：从感知到认知
    此消彼长的经验主义和理性主义钟摆

- 深度解析（Deep Parsing）是什么？

- NLP 架构纵览

## Cascaded FSAs break through Chomsky's hierarchy walls
### Robust, linear, F-measure: scale up to big data



NLP频道
liweinlp.com

# Sample Deep Parse Tree (dependency)

京东



终于冰箱安装到位了，欣喜之余发现有点儿小问题，就联系了店家，店家主动帮助联系客服上门查看，虽然最终没有解决问题，心里有点儿遗憾，但是因为不影响使用，所以也就无所谓了.

# Sample Deep Parse Tree (PS flavor)



这个可怜的年轻女孩经过非常困难的历程终于功成名就,成为了职业经理人

这个可怜的年轻女孩经过非常困难的历程终于功成名就

,

成为了职业经理人

这个可怜的年轻女孩
S

经过非常困难的历程终于功成名就
H

成为了
H

职业经理人
O

这个
X

可怜的
M

年轻
M

女孩
H

经过非常困难的历程
mannerR

终于功成名就
H

成为
H

了
X

职业
M

经理人
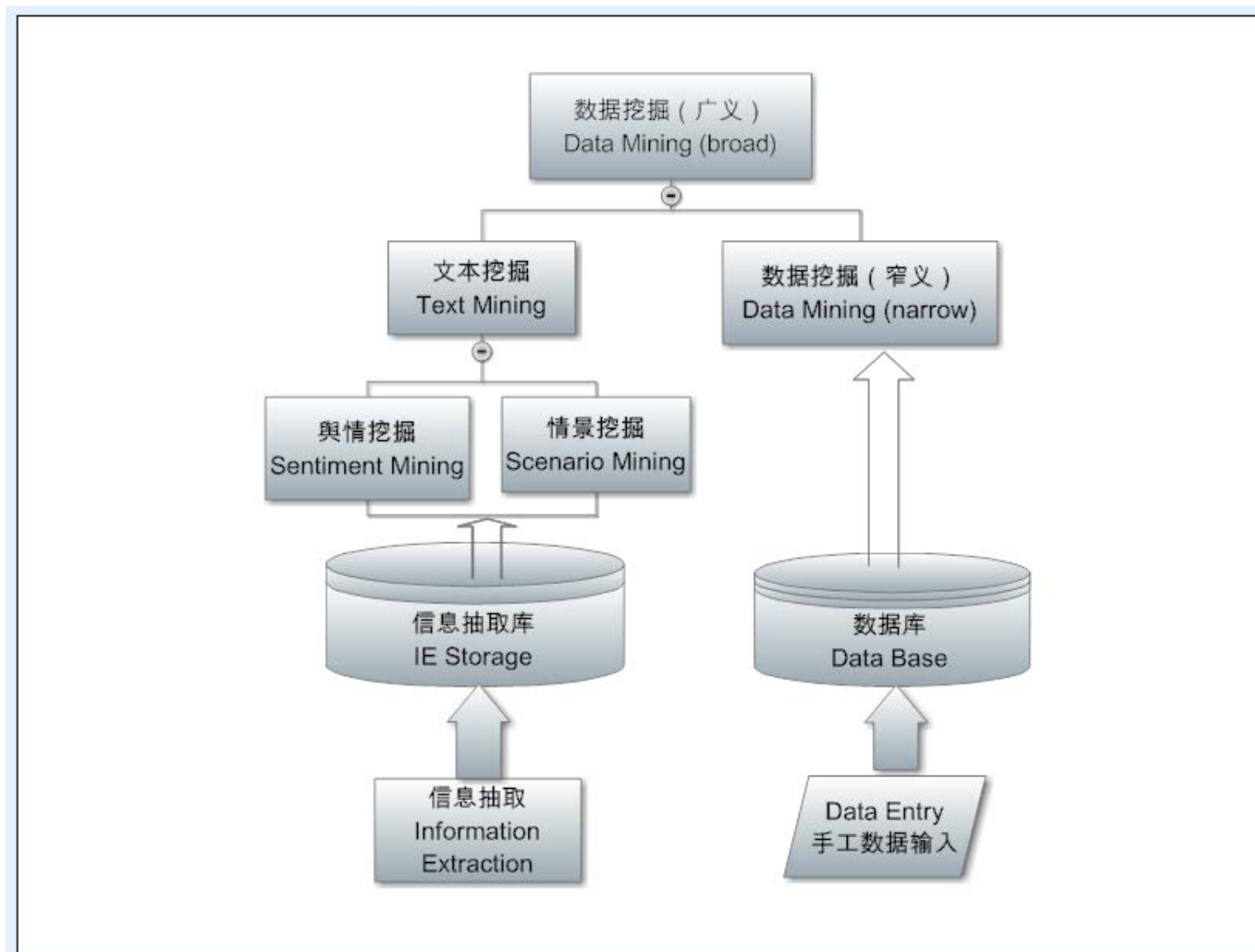H

经过
H

非常困难的历程
O

终于
R

功成名就
H

非常困难
M

的
X

历程
H

非常
R

困难
H

NLP频道
liweinlp.com

Including sentiment analysis （on subjective language）



NLP频道
liweinlp.com

# NLP Architecture 3: Text Mining

# Sample Deep Parse Tree

# Sample Deep Parse Tree

京东



传言中所说的事儿纯属空穴来风。

传言中所说的事儿纯属空穴来风

。

传言中所说的事儿
S

纯属空穴来风
H

传言中
M

所说的事儿
H

纯属
X

空穴来风
H

传言
H

中
X

所说
M

的
X

事儿
H

事
H

儿
X

NLP频道
liweinlp.com

# Outline: NLP Applications

- 人工智能的历史和现状简介：从感知到认知
    此消彼长的经验主义和理性主义钟摆

- 深度解析（Deep Parsing）是什么?

- NLP 架构纵览

- 核武器应用举例
    社媒舆情分析，大数据挖掘，智能搜索，对话系统 ………..

# Sentiment Analysis

**Why deep parsing, not deep learning?**

**Learning without parsing does not work for social media sentiment analysis**

- *Social media is dominated by short messages*
- *Statistical learning breaks in short messages: no sufficient data points*
- *Deep parsing enables linguistic analysis for best precision*
- *Deep parsing enables insights mining 2 magnitudes more efficient*
  - *parsing-supported rule has power of about 100 ngram rules*
- ***Deep learning is a great algorithm but still delinked from parsing***
  - *Parsers trained by deep learning are all research systems*
    - *difficult to adapt to real life text of social media (or other genres)*
    - *knowledge bottleneck: domains where labeled data are insufficient*
  - *Real life deep learning systems are mostly end-to-end, still no structures*

NLP频道
liweinlp.com

# Sentiment Analysis: Bag of Words vs. Parsing 京东

**KEYWORD CHALLENGE**

The iPhone has never been good.

The iPhone has never been <u>this</u> good.

**ASSOCIATION CHALLENGE**

Another reason to switch from Visa to MasterCard

I prefer MasterCard over Visa.

MasterCard is way better than Visa.

**CLASSIFICATION CHALLENGE**

*I had a wonderful day today. Even my instant coffee tastes great. However my Dell laptop doesn't boot again. Maybe I should check out the MacBook. It [MacBook] seems so easy to use.*

Coarse-grained Classification thumbs-up and down: overall tone positive (3 vs 1)
Fine-grained Analysis uncovers "why" behind sentiments:
(1) Instant coffee / tastes great
(2) Dell Laptop / does not boot
(3) Macbook / easy to use

NLP频道
liweinlp.com

# Deep Parsing Supports Deep Sentiments 京东

## Sentiment analysis has different layers

1. sentiment classification: thumbs-up and down (or neutral)

2. sentiment association: to associate a sentiment with a topic or brand as its object

3. deep sentiment insights:

   (i) who has the sentiment?

   (ii) how intense?

   (iii) why?

   (iv) Evaluations, comparisons and contrasts;

   (v) needs and wish-list;

   (vi) positive/negative actions (e.g. adopt / abandon);

   (vii) purchase intent;

   (viii) pros and cons

Most learning systems stop at 1 and sometimes at 2. All 3 can be done via deep parsing.

NLP频道
liweinlp.com

# Illustration: Real-time Polls



Challenges observed:

economy topic at 6:55pm;

China topic at 7:30pm

# Illustration: Stock Market Trends



Topic: HTC

Data 1: Stock Market Performance

Data 2: Chinese social media
(Weibo, Tianya, Facebook, Twitter…)

Time range: 2013/08 – 2014/08

Strong correlation observed

NLP频道
liweinlp.com

# Big Data Mining: Who benefits?

## For businesses: social listening

Consumer in:sights: sentiments and why

Brand image: trends

Competitive research: where do we stand

## For consumers

Purchase decision

Personalized service

## For government

Election campaign

Public opinions on policies and social topics

## Others?

Hot topics or anywhere public opinions are involved

Stock market trends correlation

# For consumers: Purchase Decision

Brand Passion Index for Washers in US Market

Like — Love

Sentiment Range

Kenmore

Whirlpool

GE

LG

Samsung

Maytag

You Wish

Dislike — Hate

Passion Intensity

Brand Passion Index: Measures buzz (size of bubble), passion, and sentiment about topics as expressed in social media.

# For consumers: Purchase Decision



Brand Passion Index for Two Types of Washers

# For consumers: Purchase Decision



Brand Passion Index

Like · Love

Sentiment Range

Maytag front loading · LG front loading

Dislike · Hate

Passion Intensity

# Intelligent Search and Chatbots

Three types of Chatbots:

    1. Domain knowledge QA:  e.g. customer service;
    2. Open domain knowledge QA:  e.g. Who won Nobel Prize in 2015?
    3. Interactive chatting: e.g. just for fun (killing time);
        in time, for comfort (senior people); for mental health counselling

Q:    questions are a subset of language, tractable
      for decoding intent, asking point, and hidden slots

A:    1 and 2 can be based on Knowledge Graph enabled by deep parser;
      3 can be enabled by learning from human chats plus parsing

A mixture/convergence of 3 is possible

# Apply NLP to Verticals: Medicine Domain 京东



**Some Big Data Verticals:**

1. News
2. Social Media

3. Medicine
4. Legal
5. Education
6. Financing

7. Multilingual

NLP频道
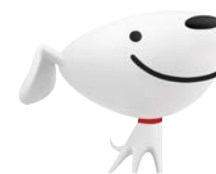liweinlp.com

And we are hiring!

At Beijing & Silicon Valley

NLP频道:
liweinlp.com
www.linkedin.com/in/
liwei4nlp