# Beyond Jeopardy!$^{TM}$
# Adapting Watson$^{TM}$ to new domains using Distributional Semantics

Alfio Massimiliano Gliozzo (gliozzo@us.ibm.com)
IBM T.J. Watson Research Center
Watson Technologies
Yorktown Heights, NY

# Outline

- Watson™ and the JEOPARDY!™ challenge

- Distributional Semantics for Domain Adaptation

# Automatic Open-Domain Question Answering
## *A Long-Standing Challenge in Artificial Intelligence to emulate human expertise*

- Given
  - Rich **Natural Language Questions**
  - Over a **Broad Domain of Knowledge**

- Deliver
  - **Precise Answers:** Determine what is being asked & give precise response
  - **Accurate Confidences:** Determine likelihood answer is correct
  - **Consumable Justifications:** Explain why the answer is right
  - **Fast Response Time:** Precision & Confidence in <3 seconds

# Informed Decision Making: Search vs. Expert Q&A

## Decision Maker

Has Question

Distills to 2-3 Keywords

Reads Documents, Finds Answers

Finds & A...

## Search Engine

Finds Documents containing Keywords

Delivers Documents based...

## Expert

## Decision Maker

Asks NL Question

Considers Answer & Evidence

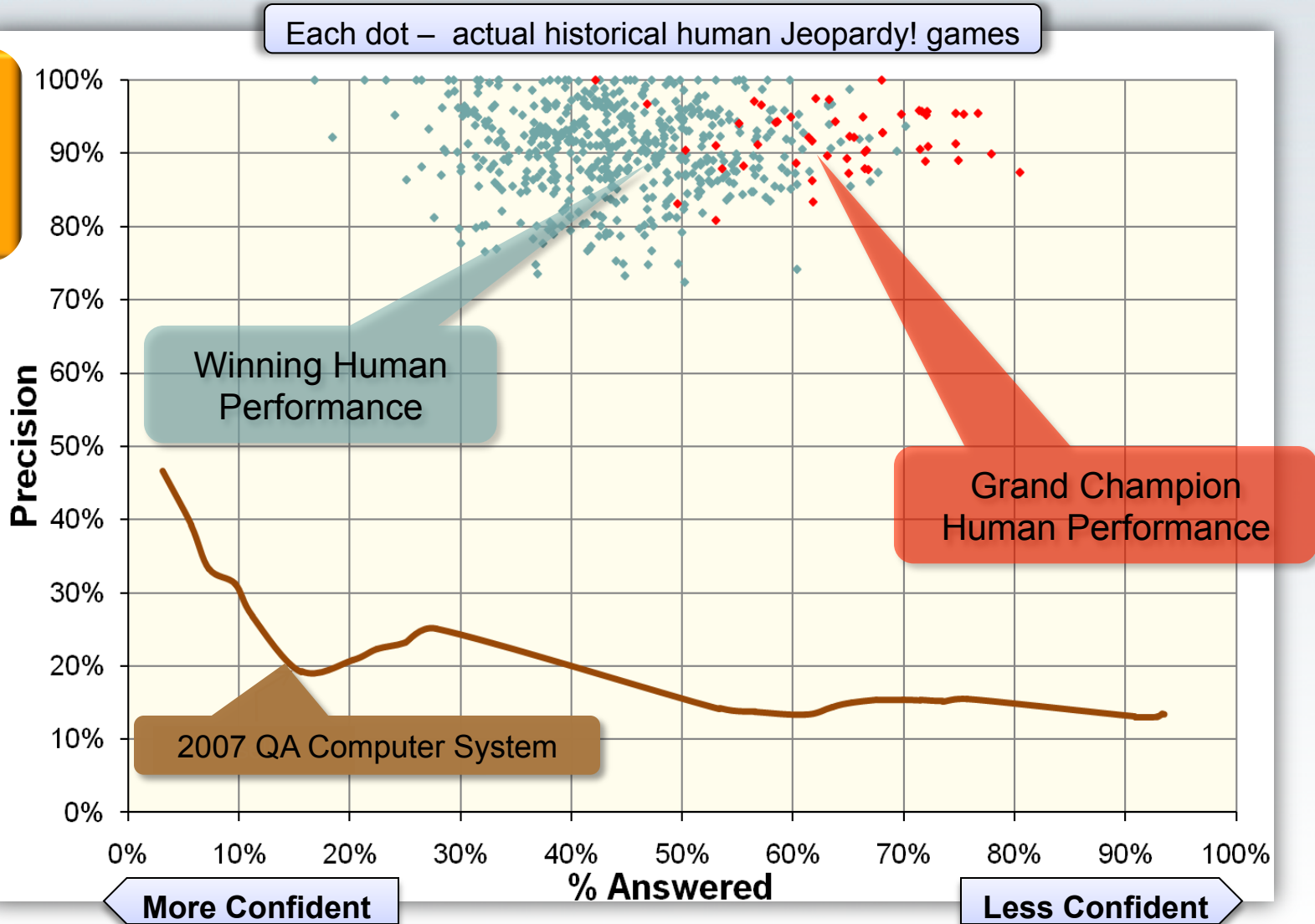Understands Question

Produces Possible Answers & Evidence

Analyzes Evidence, Computes Confidence

Delivers Response, Evidence & Confidence

**IBM WATSON**

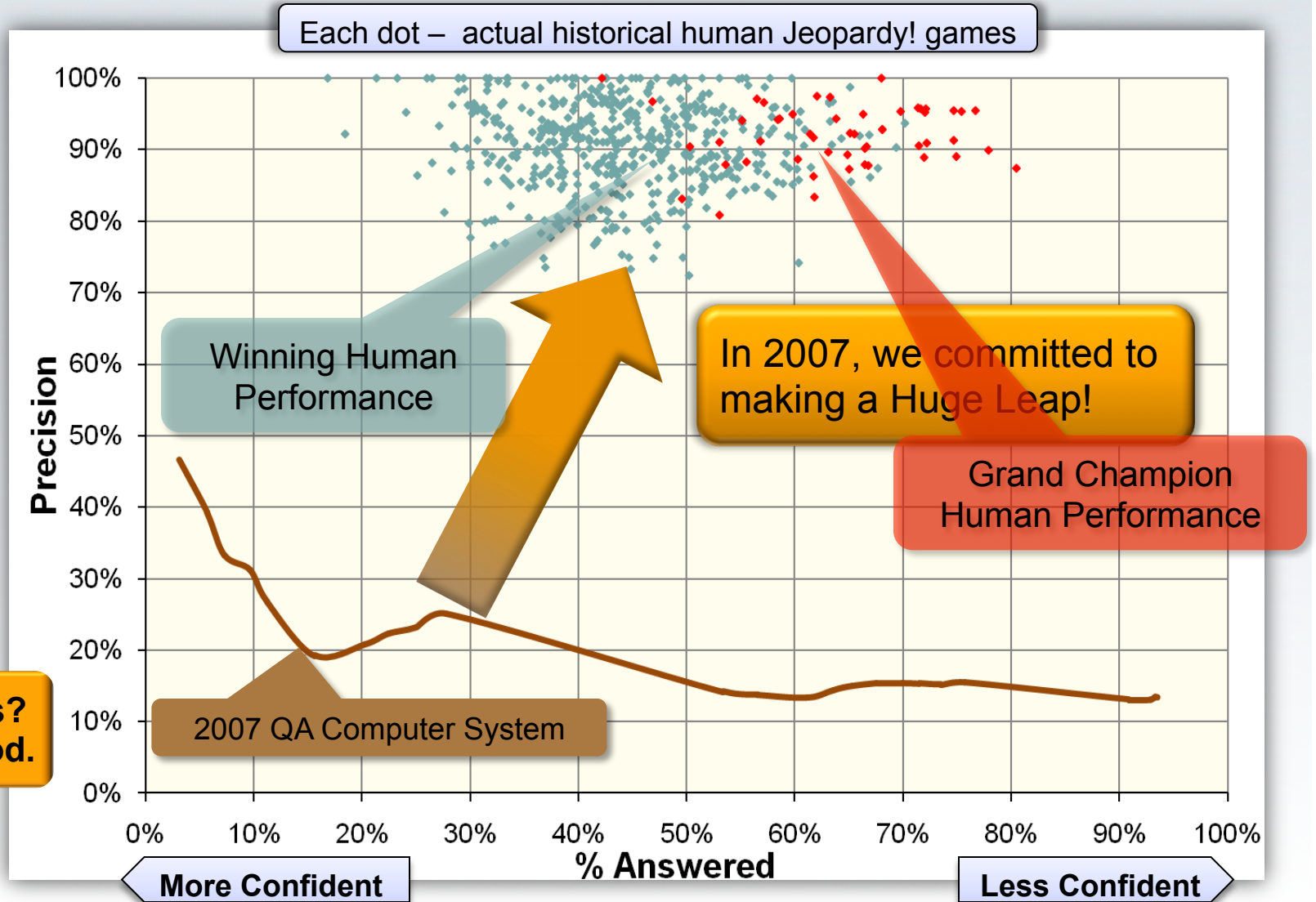# What It Takes to compete against Top Human Jeopardy!™ Players
*Our Analysis Reveals the **Winner's Cloud***

Each dot – actual historical human Jeopardy! games

**Top human players are *remarkably* good.**

**Precision**

Winning Human Performance

Grand Champion Human Performance

2007 QA Computer System

**% Answered**

More Confident

Less Confident

© 2011 IBM Corporation

What It Takes to compete against Top Human Jeopardy!™ Players
*Our Analysis Reveals the **Winner's Cloud***

# Example Question

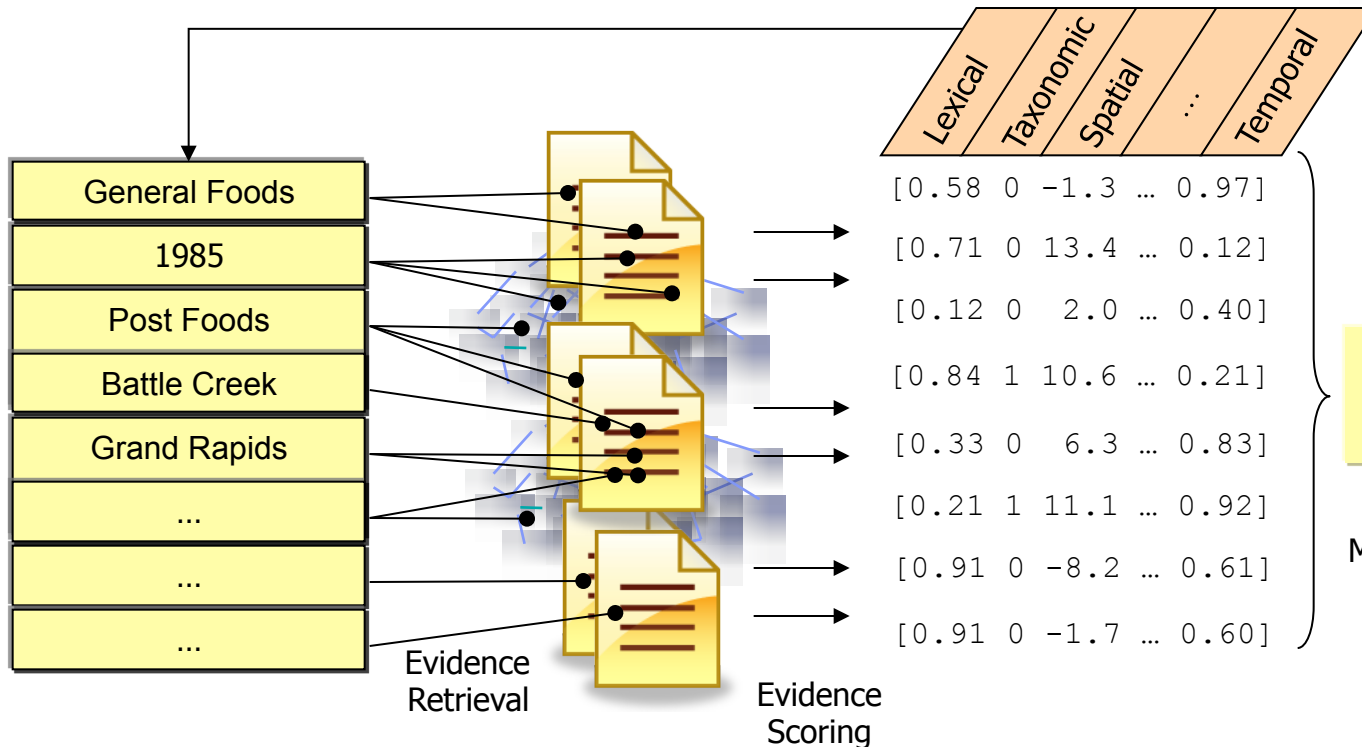**In 1894 C.W. Post created his warm cereal drink Postum in this Michigan city**

Question Analysis →

**Keywords:** 1894, C.W. Post, created …
**Lexical AnswerType:** (Michingan city)
**Date**(1894)
**Relations:**
Create(Post, cereal drink)
…

Primary Search →

Related Content
(Structured & Unstructured)

Candidate Answer Generation

| Lexical | Taxonomic | Spatial | … | Temporal |

| General Foods | [0.58 0 -1.3 … 0.97] |
| 1985 | [0.71 0 13.4 … 0.12] |
| Post Foods | [0.12 0  2.0 … 0.40] |
| Battle Creek | [0.84 1 10.6 … 0.21] |
| Grand Rapids | [0.33 0  6.3 … 0.83] |
| … | [0.21 1 11.1 … 0.92] |
| … | [0.91 0 -8.2 … 0.61] |
| … | [0.91 0 -1.7 … 0.60] |

Evidence Retrieval

Evidence Scoring

Merging & Ranking

1)  Battle Creek (0.85)
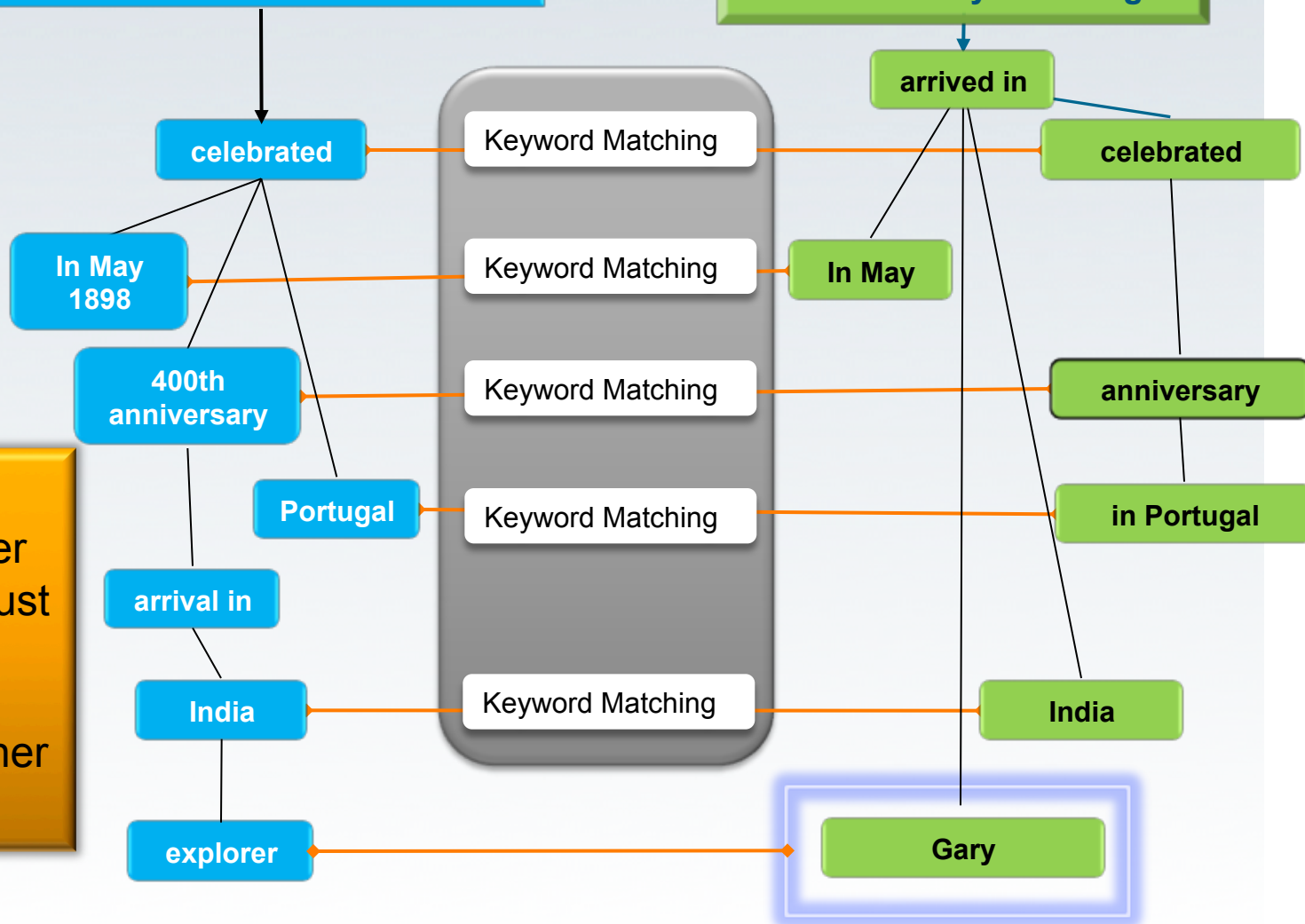2)  Post Foods  ( 0.20)
3)  1985        (0.05)

© 2011 IBM Corporation

# Keyword Evidence



In May 1898 Portugal celebrated the 400th anniversary of this explorer's arrival in India.

In May, Gary arrived in India after he celebrated his anniversary in Portugal.

celebrated — Keyword Matching — celebrated
In May 1898 — Keyword Matching — In May
400th anniversary — Keyword Matching — anniversary
Portugal — Keyword Matching — in Portugal
India — Keyword Matching — India
arrival in
explorer — Gary
arrived in

Evidence suggests "Gary" is the answer BUT the system must learn that keyword matching may be weak relative to other types of evidence

8

© 2011 IBM Corporation

# Why Semantics? Deeper Evidence

**In May 1898 Portugal celebrated the 400th anniversary of this explorer's arrival in India.**

**On the 27th of May 1498, Vasco da Gama landed in Kappad Beach**

- *Search* Far and Wide
- Explore many hypotheses
- Find Judge Evidence
- Many inference algorithms

celebrated

Portugal

May 1898

400th anniversary

landed in

27th May 1498

Date Math

Temporal Reasoning

Statistical Paraphrasing

Para-phrases

GeoSpatial Reasoning

Kappad Beach

Geo-KB

Stronger evidence can be much harder to find and score.

arrival in

India

explorer

Vasco da Gama

The evidence is still not 100% certain.

# Compare Experiments

**Now Playing in the Winners Cloud**

# Outline

- Watson™ and the JEOPARDY!™ challenge

- Distributional Semantics for Domain Adaptation

# Adaptation: What do we have in a new domain?

**Content Adaptation**

**New Text Content**
*Structure and ingest text content*

PubMed

**Training Adaptation**

**New "Questions"**
*Train the system on target scenarios*

58-year-old woman presenting to her primary care physician after several days of dizziness, anorexia, dry mouth, increased ... this ... and for ... ... the ... She had also had a f ... would "get stuck" wh ... reported no pain in h ... no cough, shortness ... Her family history inc ... in her mother, Grave ...

*What inflammation is characterized by nasal mucosal atrophy and foul-smelling crusts in the nasal passages?*
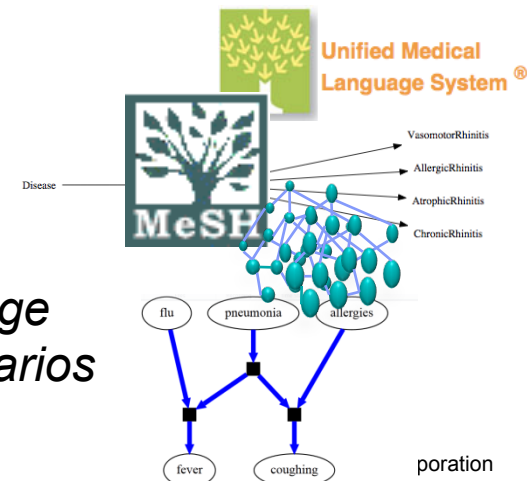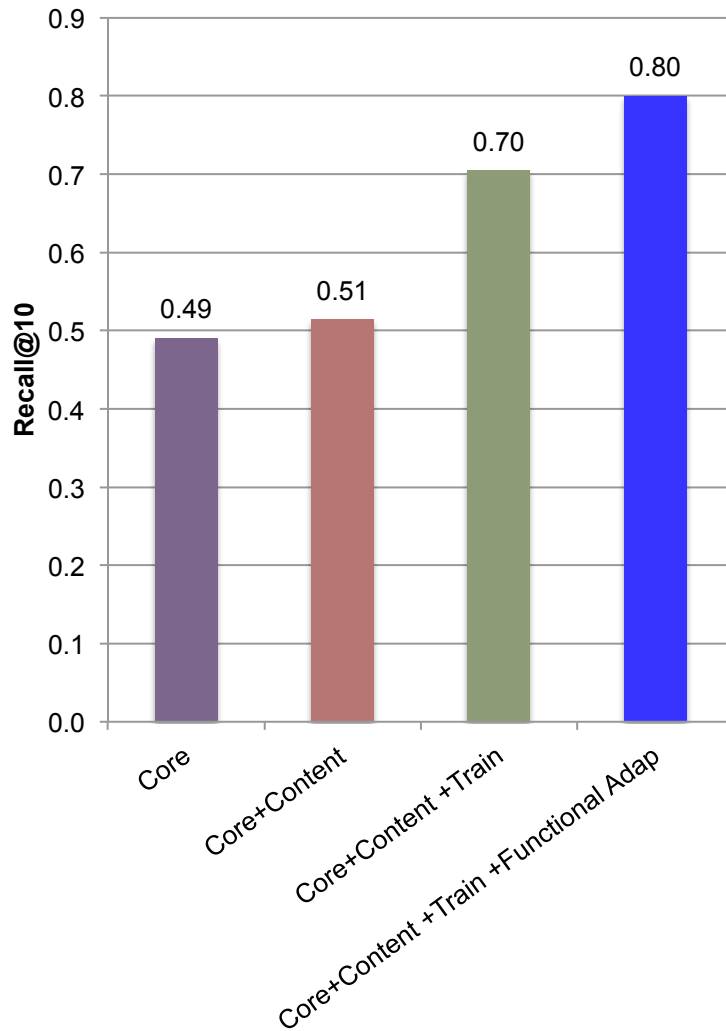
**Functional Adaptation**

**New Concepts / Reasoning / Discourse**
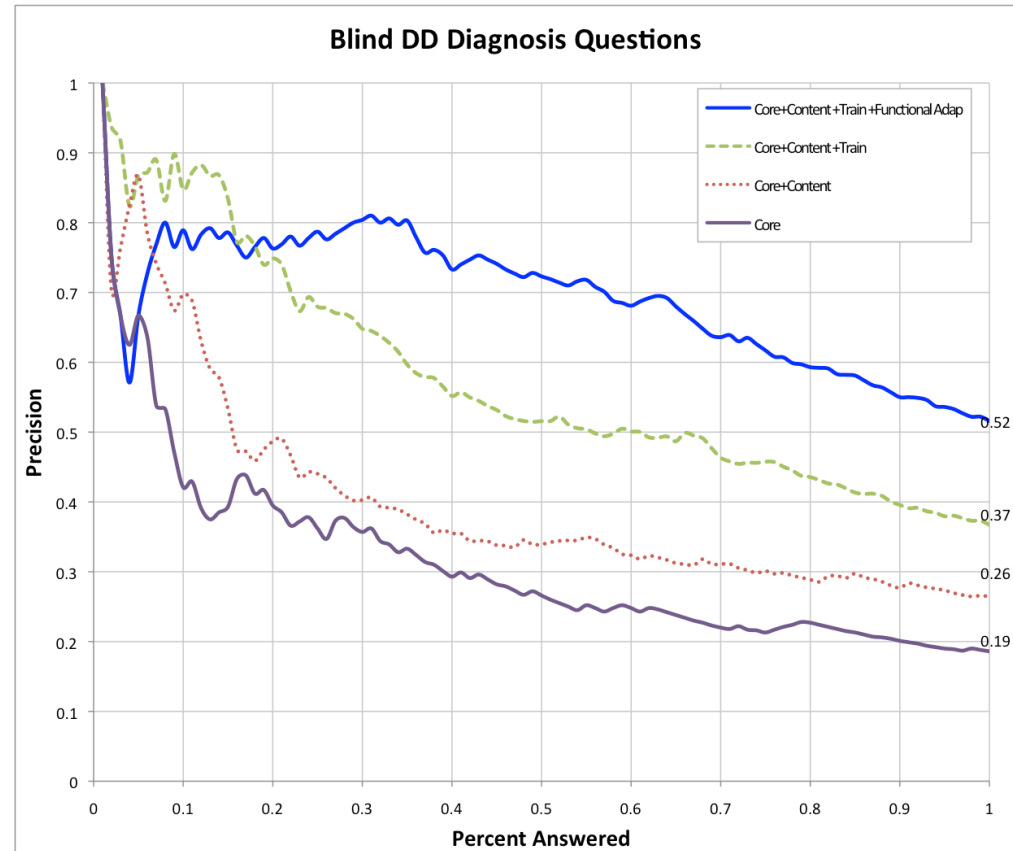*Enhance the functional capabilities with domain-specific*
   *- Concepts: entities, relations from domain modeling*
   *- Reasoning: domain axioms and background knowledge*
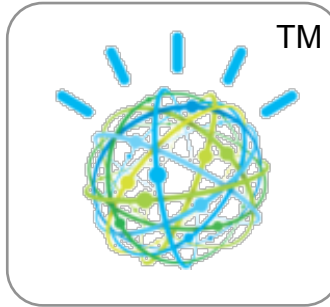   *- Discourse: algorithms for domain text / problem scenarios*

Unified Medical Language System ®

MeSH

Disease → VasomotorRhinitis, AllergicRhinitis, AtrophicRhinitis, ChronicRhinitis

flu   pneumonia   allergies

fever   coughing   poration

# Medical Adaptation - Results

Accuracy: 52%

**IBM WATSON**

TM

*Watson considers…*

What neurological condition contraindicates the use of bupropion?

contraindicate

neurological condition

use

of

Bupropion (C0085208)

NLP Stack

Tokenize /Lemmatize
Named Entity Detection
Dependency Parsing
Coreference Detection
Negation Detection
Relation Detection
Frame Extraction
Topic Detection

```
contraindicated
_drug (X,
buproprion)
```

*Structured Content*
UMLS

Knowledge Acquired from corpora

14

© 2011 IBM Corporation

# IBM WATSON



*Watson considers…*
**Unstructured Content**

What neurological condition contraindicates the use of bupropion?

Wellbutrin - noradrenergic antidepressant. contraindicated in adults with seizure disorders due to possible lowering of seizure threshold

Bupropion is contraindicated in epilepsy, seizure disorder; anorexia/bulimia (eating disorders), patients' use of antidepressant drugs (MAO inhibitors) within 14 days,

contraindicate

neurological condition

use

of

Bupropion (C0085208)

Patients with preexisting *seizure disorder* should not use bupropion due to a higher-than-proportional increase in the possibility of seizure as the dose is increased.

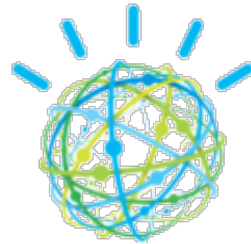contraindicate

Buproprion

in

epilepsy

**Matching Framework**

**NLP Stack**
Tokenize /Lemmatize
Named Entity Detection
Dependency Parsing
Coreference Detection
Negation Detection
Relation Detection
Frame Extraction
Topic Detection

**NLP Stack**
Tokenize /Lemmatize
Named Entity Detection
Dependency Parsing
Coreference Detection
Negation Detection
Relation Detection
Frame Extraction
Topic Detection

**IBM WATSON**

What neurological condition **contraindicates the use of** bupropion?

*Watson considers…*
**Unstructured Content**

Bupropion is contraindicated in epilepsy, seizure disorder; anorexia/bulimia (eating disorders), patients' use of antidepressant drugs (MAO inhibitors) within 14 days,

contraindicate

use

of

neurological condition

Bupropion (C0085208)

contraindicate

buproprion

seizure disorder

in

epilepsy

anorexia

bulimia

**NLP Stack**

Tokenize /Lemmatize
Named Entity Detection
Dependency Parsing
Coreference Detection
Negation Detection
Relation Detection
Frame Extraction
Topic Detection

Need to consider the type ("**neurological condition**") of the answer for possible candidates:
• Epilepsy
• Seizure disorder
• Anorexia
• Bulimia

**NLP Stack**

Tokenize /Lemmatize
Named Entity Detection
Dependency Parsing
Coreference Detection
Negation Detection
Relation Detection
Frame Extraction
Topic Detection

*Structured Content*
UMLS

Knowledge Acquired from corpora

**IBM WATSON**

What neurological condition contraindicates the use of bupropion?

contraindicate

neurological condition

use

of

Bupropion (C0085208)

TM

*Watson considers…*
**Unstructured Content**

Wellbutrin - noradrenergic antidepressant. contraindicated in adults with seizure disorders due to possible lowering of seizure threshold

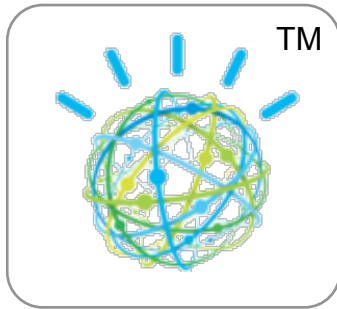Use background medical knowledge (**Wellbutrin** is a brand name of **bupropion**)

**NLP Stack**

Tokenize /Lemmatize
Named Entity Detection
Dependency Parsing
Coreference Detection
Negation Detection
Relation Detection
Frame Extraction
Topic Detection

**NLP Stack**

Tokenize /Lemmatize
Named Entity Detection
Dependency Parsing
Coreference Detection
Negation Detection
Relation Detection
Frame Extraction
Topic Detection

*Structured Content*

UMLS

Knowledge Acquired from corpora

17

**IBM WATSON**

TM

What neurological condition contraindicates the use of bupropion?

contraindicate

neurological condition

use

of

Bupropion (C0085208)

Patients with preexisting *seizure disorder* should not use bupropion due to a higher-than-proportional increase in the possibility of seizure as the dose is increased.

Consider paraphrases in medical language:
(**should not use = contraindicate**)

**NLP Stack**

Tokenize /Lemmatize
Named Entity Detection
Dependency Parsing
Coreference Detection
Negation Detection
Relation Detection
Frame Extraction
Topic Detection

**NLP Stack**

Tokenize /Lemmatize
Named Entity Detection
Dependency Parsing
Coreference Detection
Negation Detection
Relation Detection
Frame Extraction
Topic Detection

*Structured Content*
UMLS

Knowledge Acquired from corpora

- Domain adaptation is difficult!

- Requires:
  - Deeply skilled research team across all the key disciplines (ML, NLP, IR, KR)
  - Domain Experts (Doctors) for annotation/vetting and design reasoning strategies
  - Collaboration between the two groups !
  - Background knowledge for new domains (e.g. UMLS) and analytics exploiting that
  - Rigorous methodological discipline (e.g., blind test!)

- Future Challenge: Scalable and cost effective functional adaptation process
  - Acquiring Domain Knowledge from Text
    - Taxonomy induction
    - Statistical Paraphrasing
    - Sense Induction/ Unsupervised WSD
  - Using the same analytics (e.g. matching, tycor) across domains

- We call it Distributional Semantics!

# The Distributional Semantics Paradigm

- The challenge: Fully Unsupervised Computational Semantics
  - Input: few Gigabytes of raw text in a specific domain
  - Output: Semantic Analyzer having the following capabilities
    - Term/Text Similarity beyond Keyword Matching
    - WSD, Lexical Substitution
    - Matching: terms, relations
    - Linking text to knowledge bases
  - Radical Approach:
    - Mining (clustering) big data
    - No Rules, No labeled data
- Making it scalable (Hadoop)
  - More text = more hardware = same time
  - Fast semantic parsing
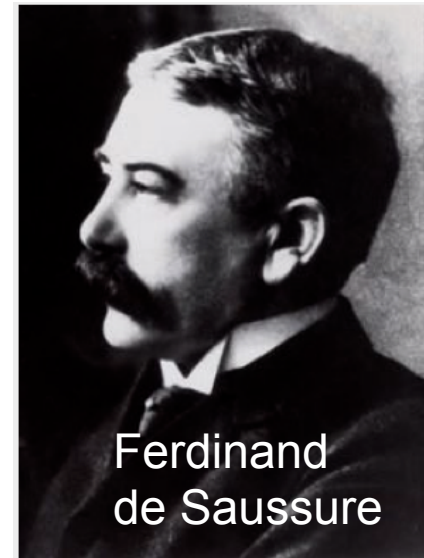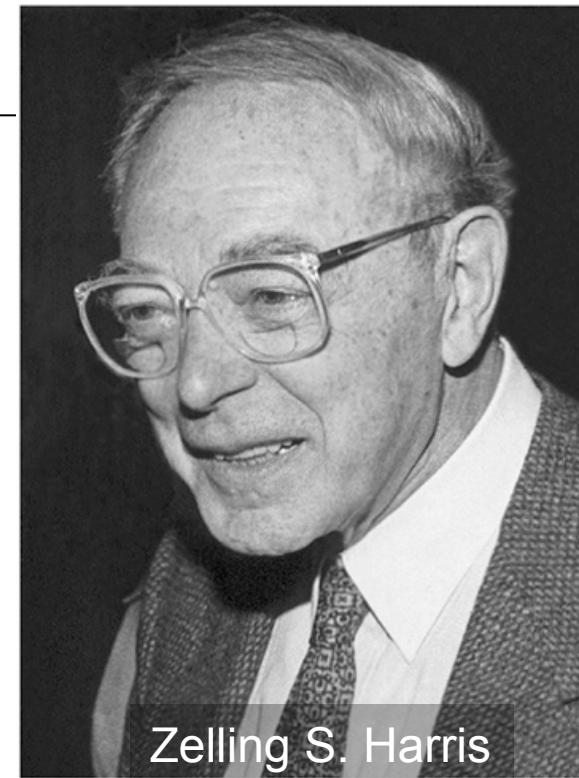  - Web size Distributional Semantics to capture background knowledge

# Distributional Hypothesis and Structuralism

The **Distributional Hypothesis** in linguistics is the theory that words that occur in similar contexts tend to have similar meanings (paradigmatic relations).
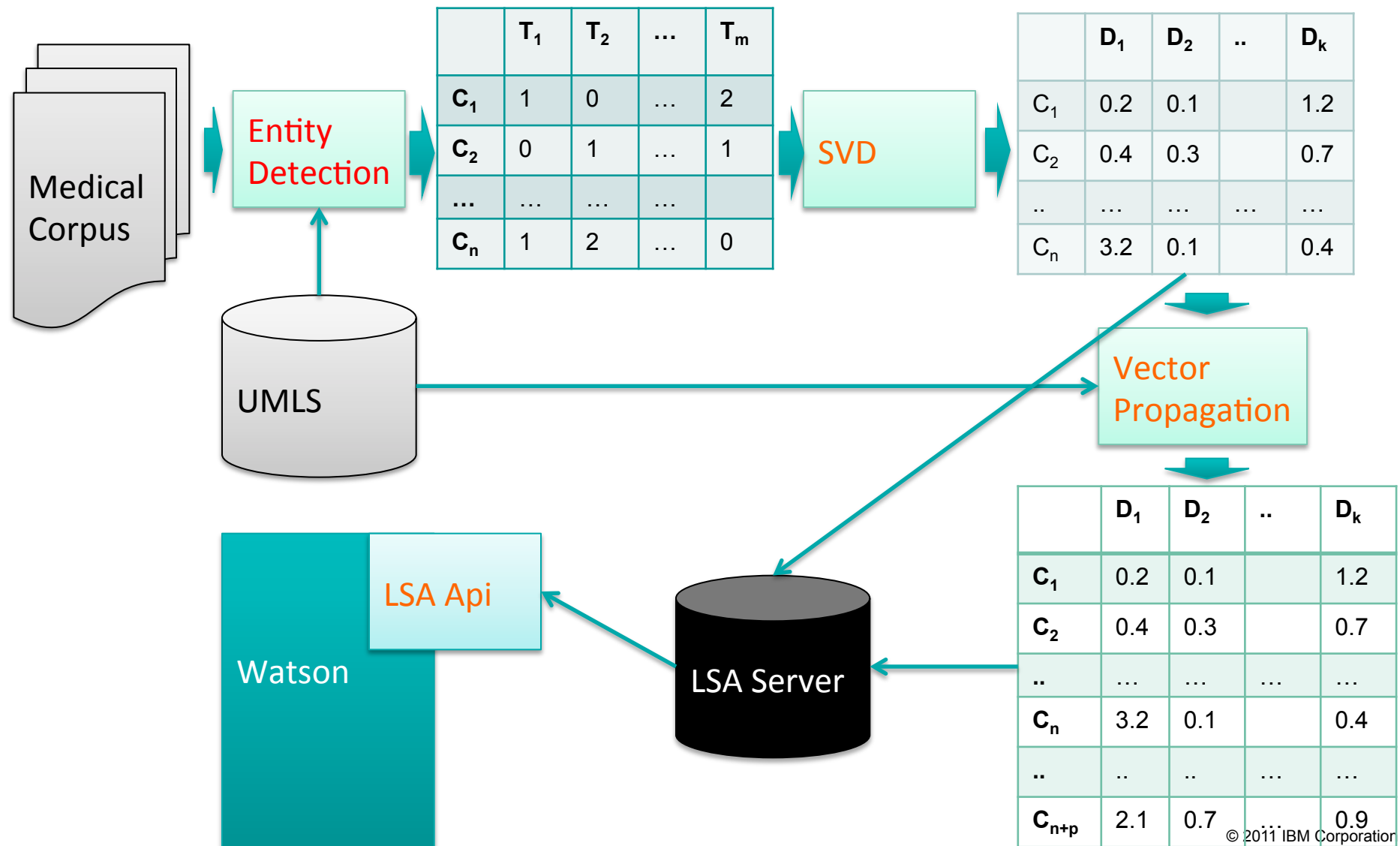
The Distributional Hypothesis is the basis for Distributional Semantics.

It states that the meaning of a word can be defined in terms of its context (properties).

- other words in the same sentence/document (bag of words)
- words in the immediate neighbors
- words along dependency paths
- Predicate Argument Structure
- Frame

➔ any process that builds a structure on sentences can be used as a source for properties

Z. Harris. (1954). Distributional Structure. Word 10 (2/3)

Zelling S. Harris

Ferdinand de Saussure

IBM WATSON

## Latent Semantic Analysis (2.0)



Medical Corpus → Entity Detection

|     | $T_1$ | $T_2$ | ... | $T_m$ |
|-----|-------|-------|-----|-------|
| $C_1$ | 1 | 0 | ... | 2 |
| $C_2$ | 0 | 1 | ... | 1 |
| ... | ... | ... | ... | ... |
| $C_n$ | 1 | 2 | ... | 0 |

→ SVD →

|     | $D_1$ | $D_2$ | .. | $D_k$ |
|-----|-------|-------|----|-------|
| $C_1$ | 0.2 | 0.1 | | 1.2 |
| $C_2$ | 0.4 | 0.3 | | 0.7 |
| .. | ... | ... | ... | ... |
| $C_n$ | 3.2 | 0.1 | | 0.4 |

UMLS

Vector Propagation

|     | $D_1$ | $D_2$ | .. | $D_k$ |
|-----|-------|-------|----|-------|
| $C_1$ | 0.2 | 0.1 | | 1.2 |
| $C_2$ | 0.4 | 0.3 | | 0.7 |
| .. | ... | ... | ... | ... |
| $C_n$ | 3.2 | 0.1 | | 0.4 |
| .. | .. | .. | ... | ... |
| $C_{n+p}$ | 2.1 | 0.7 | | 0.9 |

LSA Api

Watson

LSA Server

# LSA 2.0: Factors related to suicide

| | Accuracy | Precision@70 |
|---|---|---|
| LSA | +0.66% | +0.47% |
| LSA 2.0 | +1.13% (4.46%) | +1.58%(5.229%) |

"Suicide"

| | | | |
|---|---|---|---|
| DANGER OF HARM TO SELF | 0.94843552 | Feeling hopeless | 0.69763276 |
| Depressive Symptoms | 0.85787663 | CYCLOTHYMIC REACTION | 0.6956163 |
| marked mood shift | 0.83171128 | Mental health counselor | 0.6916423 |
| loss of interest in activity | 0.83171128 | Demoralization | 0.68469489 |
| Other mood affective disorders | 0.80852182 | Ability to maintain self-esteem | 0.67854127 |
| Mood Disorders | 0.80852182 | Normal mood | 0.67817024 |
| Bipolar affective disorder, current episode manic | 0.79134531 | Despondency | 0.67736145 |
| Depressive disorder NEC in SNOMEDCT | 0.78274978 | Other and unspecified episodic mood disorder | 0.67540516 |
| change in self-esteem | 0.77332301 | Loss of interest | 0.67413379 |
| (Depression: [episode, unspecified] or [NOS (& | 0.76803559 | Suicidal | 0.67144792 |
| Self Esteem | 0.72473412 | pleasurable emotion | 0.67024476 |
| self-esteem as an AODC | 0.7247341 | Mood (psychological function) | 0.67023923 |
| AODE on self-esteem | 0.7247341 | Mood:-:Point in time:^Patient:- | 0.66983514 |
| | | Suicidal behavior | 0.6680896 |
| | | Adjustment disorder with depressed mood | 0.6555974 |
| | | Depression aggravated | 0.6528632 |
| | | Coping with Chronic Illness Topics | 0.64542571 |
| | | Mental Health and Behavior | 0.6454257 |
| | | Recurrent depression | 0.64434724 |
| | | Other specified episodic mood disorder | 0.64310002 |
| | | Melancholia | 0.64063775 |
| | | Mild recurrent major depression | 0.63696897 |

LSA 2.0

# The @@ operation

**SENTENCE**:

*I suffered from a cold and took aspirin.*

**STANFORD COLLAPSED DEPENDENCIES:**

*http://nlp.stanford.edu:8080/parser/*

nsubj(suffered, I); nsubj(took, I); root(ROOT, suffered); det(cold, a); prep_from(suffered, cold); conj_and(suffered, took); dobj(took, aspirin)

**WORD-PROPERTY PAIRS:**

| Jo | | |
|---|---|---|
| suffered | nsubj(@@, I) | 1 |
| took | nsubj(@@, I) | 1 |
| cold | det(@@, a) | 1 |
| suffered | prep_from(@@, cold) | 1 |
| suffered | conj_and(@@, took) | 1 |
| took | dobj(@@, aspirin) | 1 |

| Bim | | |
|---|---|---|
| I | nsubj(suffered, @@) | 1 |
| I | nsubj(took, @@) | 1 |
| a | det(cold, @@) | 1 |
| cold | prep_from(suffered, @@) | 1 |
| took | conj_and(suffered, @@) | 1 |
| aspirin | dobj(took, @@) | 1 |

# JoBimText
## Linking Language to Knowledge with Distributional Semantics

- www.jobimtext.org
- Open Source Software
  - Apache License
  - SourceForge
- Contributors
  - TU Darmstadt, Germany, FG Language Technology
    - **Chris Biemann** (**Bim**) , Martin Riedl
  - IBM T.J. Watson Research - Watson Technologies
    - **Alfio Gliozzo** (**Jo**) , Michael Glass, Bonaventura Coppola
- What's there
  - Scalable Distributional Similarity (Hadoop)
  - UIMA based text processing implementing @@ operation on different languages/NLP
  - Fast and Scalable Knowledge Management
  - Sense Clustering, WSD, lexical substitution, Thesauri induction, Paraphrasing, Entity Linking, …
  - Machine Learning: CRF, Chinese Whisper Clustering, …

Unstructured Information Management Architecture
An Apache Project.

hadoop

- Input:
  - Watson Medical Corpus
    - ~ 2 Gigabytes of text
    - UMLS
- Preprocessing:
  - Medical Extended Slot Grammar (ESG) Parser
    - Dependency Parser
    - Medical Adaptation of the Jeopardy Parser
  - TWREX
    - Relation Extraction system adapted to UMLS relations
- @@ system:
  - Terms are represented by
    - syntactic dependencies
    - TWREX relations
- Unsupervised learning on a Small Hadoop Cluster
- Watson Analytics for Answer Scoring, Matching, Passage Scoring
- Demo

# References

- Ferrucci et al., **Building Watson: An Overview of the DeepQA Project,** AI Magazine, 2010

- *Ferrucci et al.,* **Watson: Beyond Jeopardy!,** 2011 RC25270, to appear in Artificial Intelligence Journal.

- Deep QA publications website
  - http://researcher.ibm.com/view_grouppubs.php?grp=2099

- Videos on Watson
  - http://www-03.ibm.com/innovation/us/watson/index.html

**VOLUME 56, NUMBER 3/4, MAY/JUL. 2012**

**IBM Journal of Research and Development**

**Including IBM Systems Journal**

| | |
|---|---|
| Bram Stoker | 62% |
| Dracula | 40% |
| Herman Melville | 9% |

This Is Watson

- http://ieeexplore.ieee.org/xpl/tocresult.jsp?reload=true&isnumber=6177717