

# Building *Watson*

## *An overview of the DeepQA Project*

David Ferrucci

*Watson* Principal Investigator and IBM Fellow  
DeepQA Team @ IBM Research



# A Grand Challenge Opportunity



- **Drive Important Scientific Advances**

- Envision new ways for computers to impact society & science

- **Be Relevant to IBM Customers**

- Enable better, faster decision making over unstructured and structured content
- *Business Intelligence, Knowledge Discovery and Management, Government, Compliance, Publishing, Legal, Healthcare, Product Support, etc.*

- **Capture the Broader Imagination**

- The Next *Deep Blue*

# Want to Play Chess or Just Chat?

## ■ Chess

- A finite, mathematically well-defined search space
- Limited number of moves and states
- All the symbols are completely grounded in the mathematical rules of the game



## ■ Human Language

- Words by themselves have no meaning
- Only grounded in **human cognition**
- Words navigate, align and communicate an infinite space of intended meaning
- Computers can **not** ground words to human experiences to derive meaning



# Easy Questions?

$$=(\text{LN}(12,546,798 * \pi))^3 / 34,600.47 = 0.155$$

Select *Payment* where *Owner*="David Jones" and *Type(Product)*="Laptop",

Owner	Serial Number
David Jones	45322190-AK

Invoice #	Vendor	Payment
INV10895	MyBuy	\$104.56

Serial Number	Type	Invoice #
45322190-AK	LapTop	INV10895

David Jones  
 ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓  
 David Jones

=

Dave Jones ≠  
 ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓  
 David Jones

# Hard Questions?

Computer programs are natively **explicit**, **fast** and **exacting** in their calculation over numbers and symbols....But **Natural Language** is implicit, highly contextual, ambiguous and often imprecise.

Person	Birth Place
A. Einstein	ULM

Structured

Unstructured

## ■ Where was X born?

*One day, from among his city views of Ulm, Otto chose a water color to send to Albert Einstein as a remembrance of Einstein's birthplace.*

Person	Organization
J. Welch	GE

## ■ X ran this?

*If leadership is an art then surely Jack Welch has proved himself a master painter during his tenure at GE.*

# The Jeopardy! Challenge: *A compelling and notable way to drive and measure the technology of automatic Question Answering along 5 Key Dimensions*

**Broad/Open  
Domain**

**Complex  
Language**

**High  
Precision**

**Accurate  
Confidence**

**High  
Speed**

**\$200**

If you're standing, it's the direction you should look to check out the wainscoting.

**\$1000**

The first person mentioned by name in 'The Man in the Iron Mask' is this hero of a previous book by the same author.

**\$600**

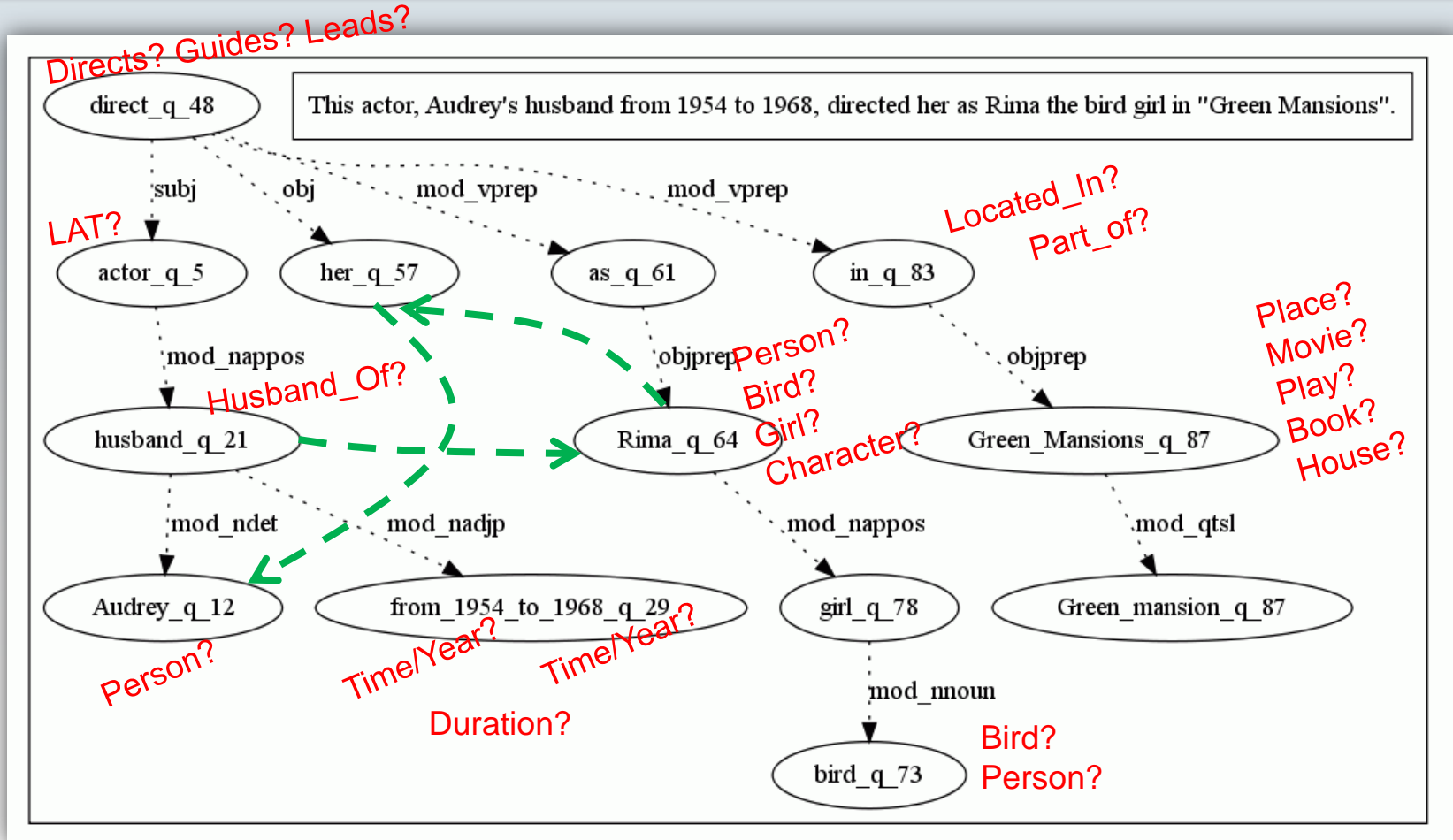
In cell division, mitosis splits the nucleus & cytokinesis splits this liquid *cushioning* the nucleus

**\$2000**

Of the 4 countries in the world that the U.S. does not have diplomatic relations with, the one that's farthest north

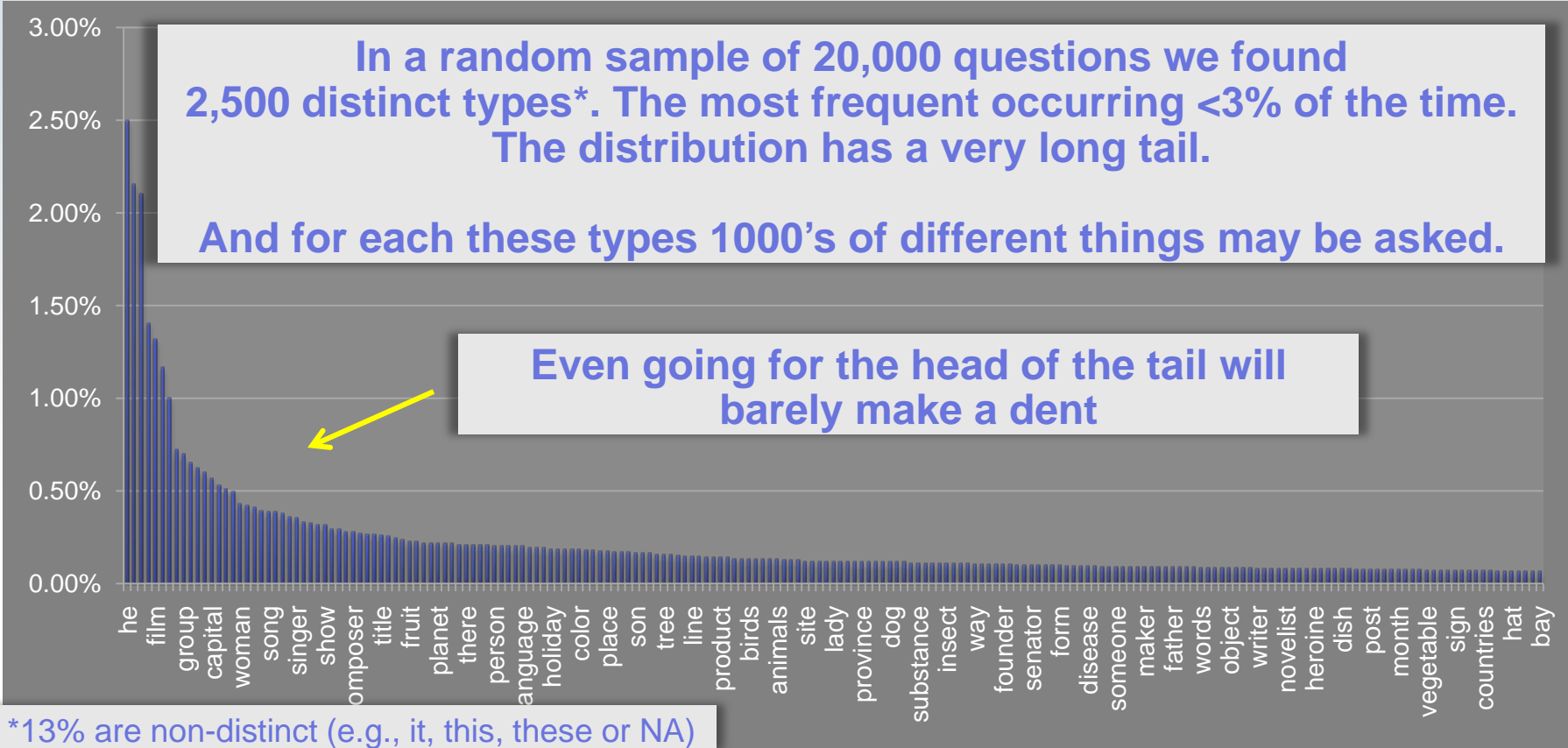
# The Possibilities Multiply

**This actor, Audrey's husband from 1954 to 1968, directed her as Rima the bird girl in "Green Mansions"**



# Broad Domain

**We do NOT attempt to anticipate all questions and build specialized databases.**



**Our Focus is on reusable NLP technology for analyzing volumes of *as-is* text. Structured sources (DBs and KBs) are used to help interpret the text.**



# Inducing Frames

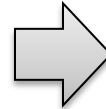
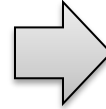
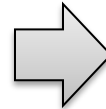
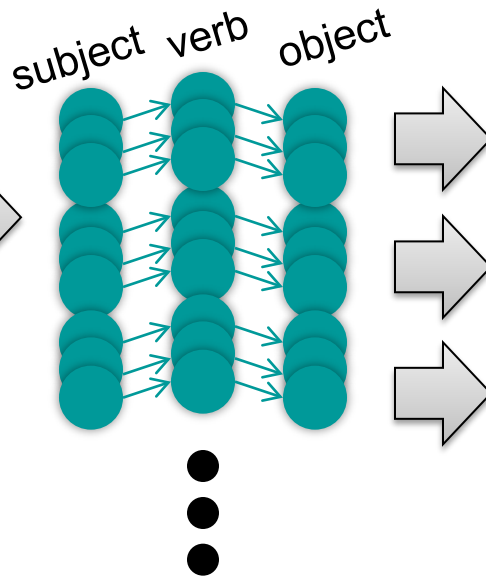
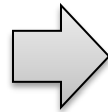
Sentence  
Parsing

Generalization &  
Statistical Aggregation

**Volumes of Text**

***Syntactic Frames***

***Semantic Frames***



- Inventors patent inventions (.8)
- Officials Submit Resignations (.7)
- People earn degrees at schools (0.9)
- Fluid is a liquid (.6)
- Liquid is a fluid (.5)
- Vessels Sink (0.7)
- People sink 8-balls (0.3)
- (pool game (0.8))

## Evaluating Possibilities and Their Evidence

In cell division, mitosis splits the nucleus & cytokinesis splits this **liquid** *cushioning* the nucleus.

- *Organelle*
- *Vacuole*
- *Cytoplasm*
- *Plasma*
- *Mitochondria*
- *Blood ...*

- Many candidate answers (CAs) are generated from many different searches
- Each possibility is evaluated according to **different dimensions of evidence**.
- **Just One** piece of evidence is if the CA is of the **right type**. In this case a "liquid".

Is("Cytoplasm", "liquid") = 0.2↑

Is("organelle", "liquid") = 0.1

Is("vacuole", "liquid") = 0.2

Is("plasma", "liquid") = 0.7

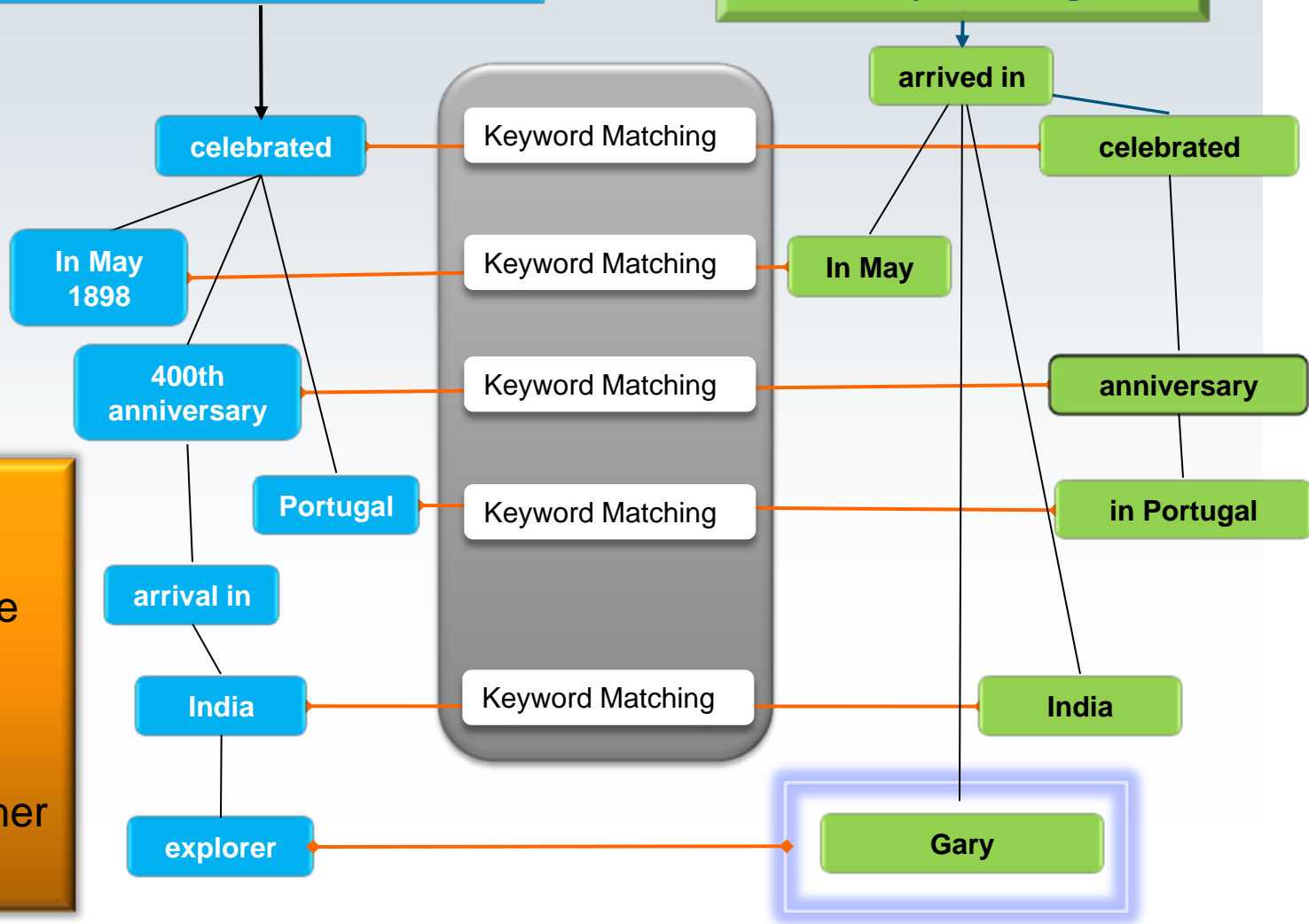
"Cytoplasm is a **fluid** surrounding the nucleus..."

Wordnet → Is\_a(Fluid, Liquid) → ?

Learned → Is\_a(Fluid, Liquid) → yes.

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

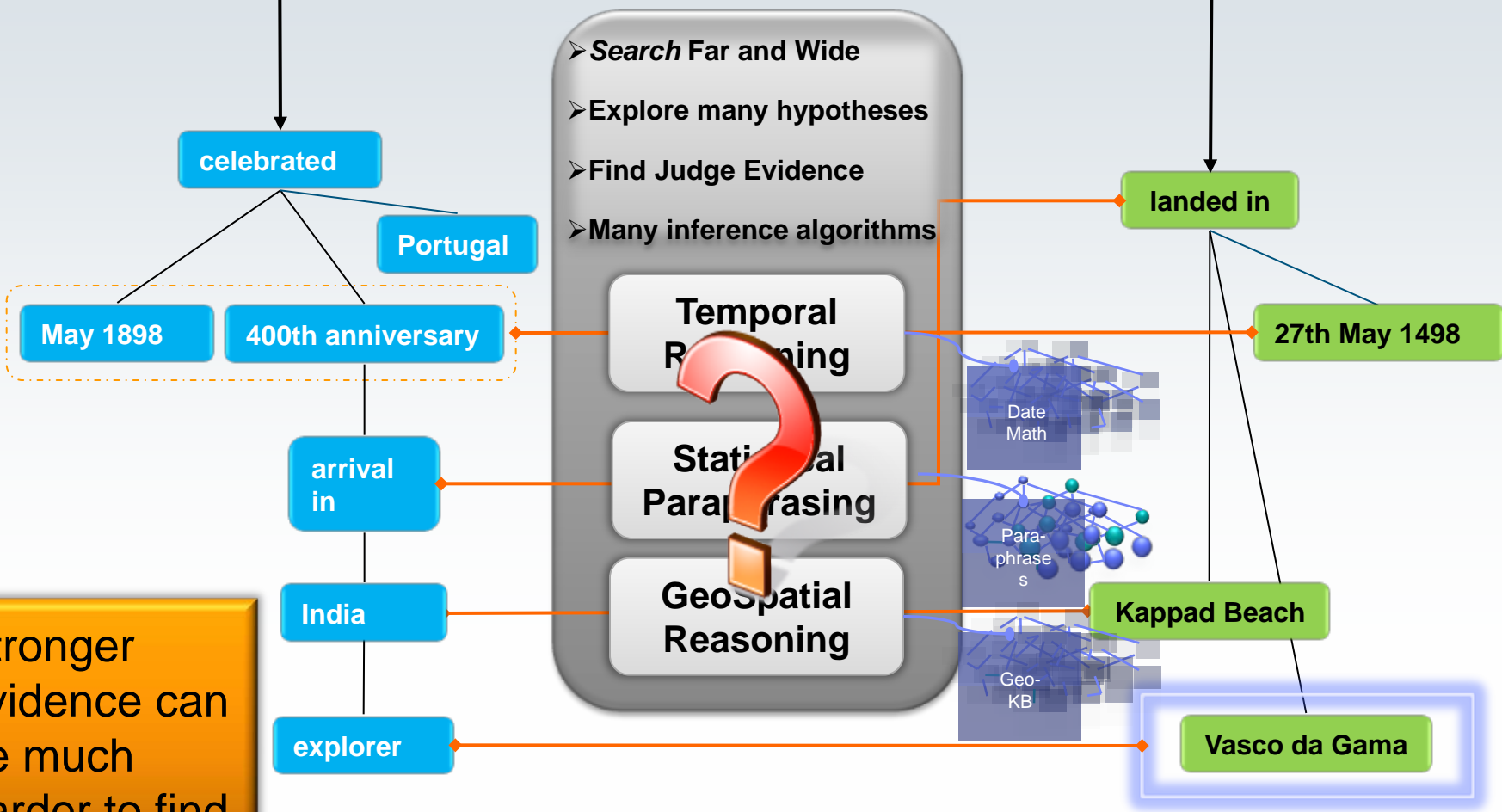
In May, Gary arrived in India after he celebrated his anniversary in Portugal.



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

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

On the 27<sup>th</sup> of May 1498, Vasco da Gama landed in Kappad Beach

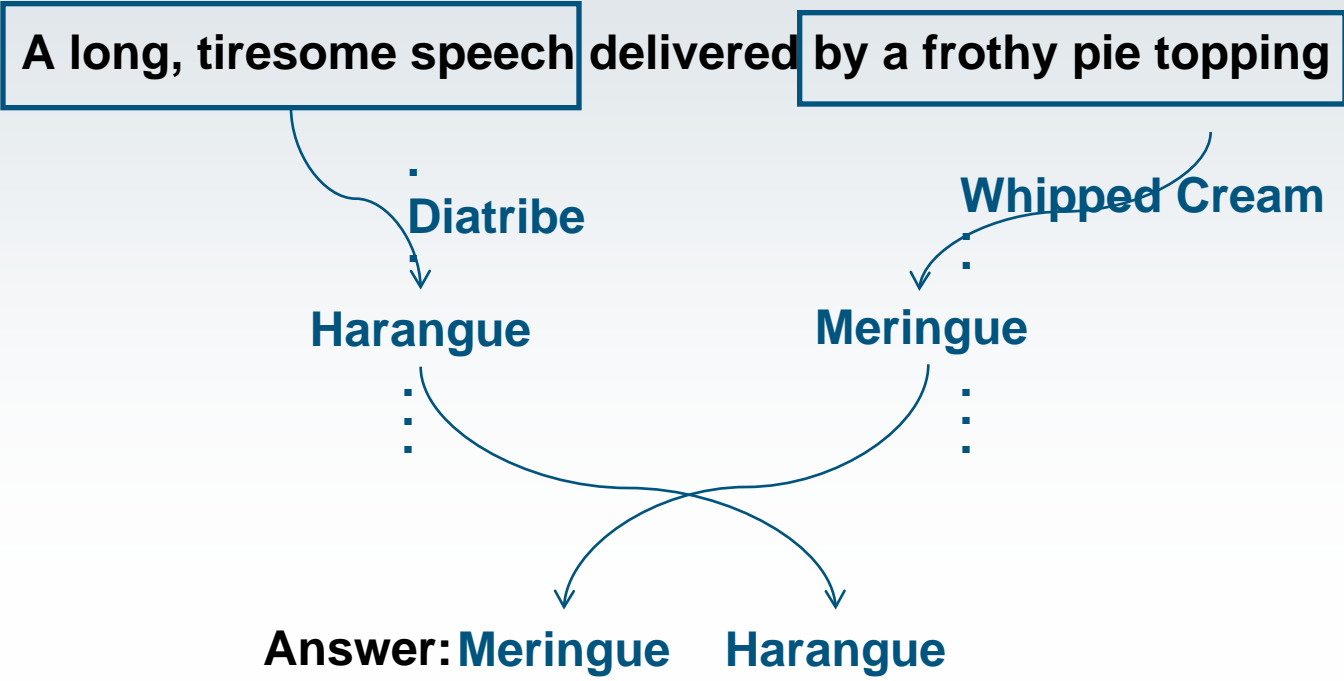


Stronger evidence can be much harder to find and score.

The evidence is still not 100% certain.

# Not *Just* for Fun

Some Questions require Decomposition and Synthesis



Category: Edible Rhyme Time

## Divide and Conquer (Typical in Final Jeopardy!)

**Must identify and solve sub-questions from different sources to answer the top level question**

**Lyndon B Johnson**

**In 1968**

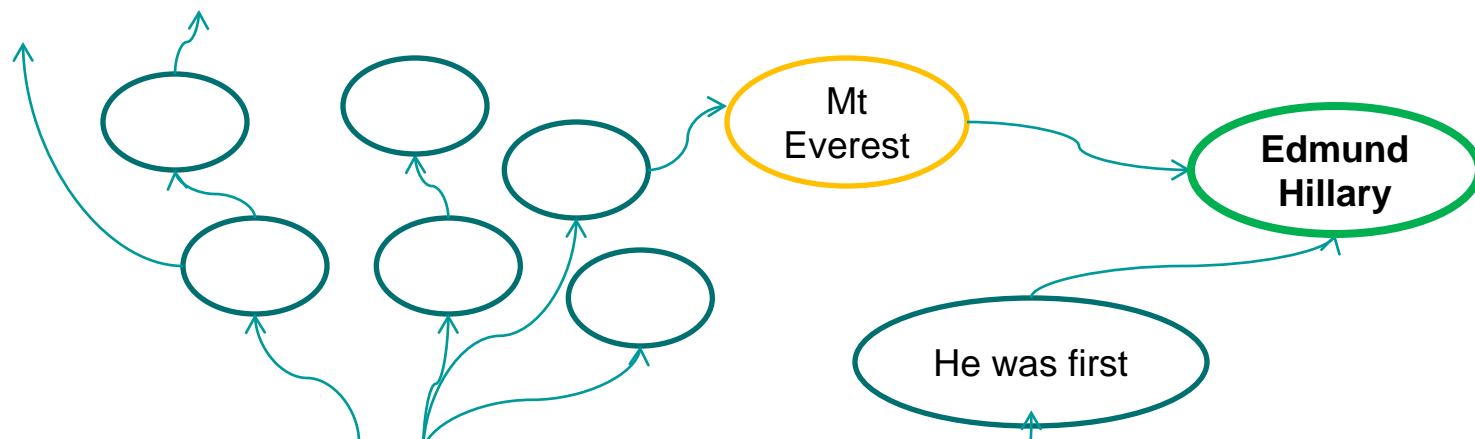
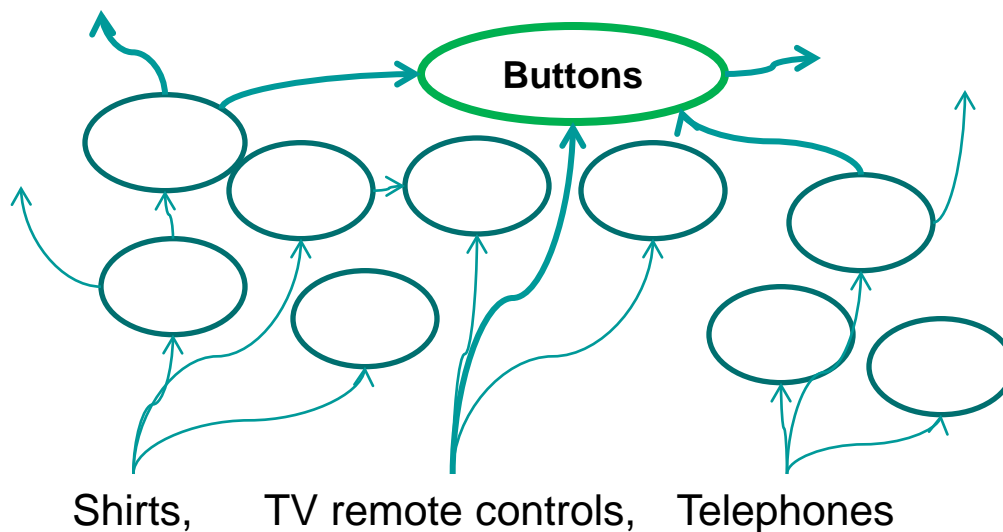
this man was U.S. president.

When "60 Minutes" premiered | this man was U.S. president.

?

**The DeepQA architecture attempts different *decompositions* and recursively applies the QA algorithms**

# The Missing Link

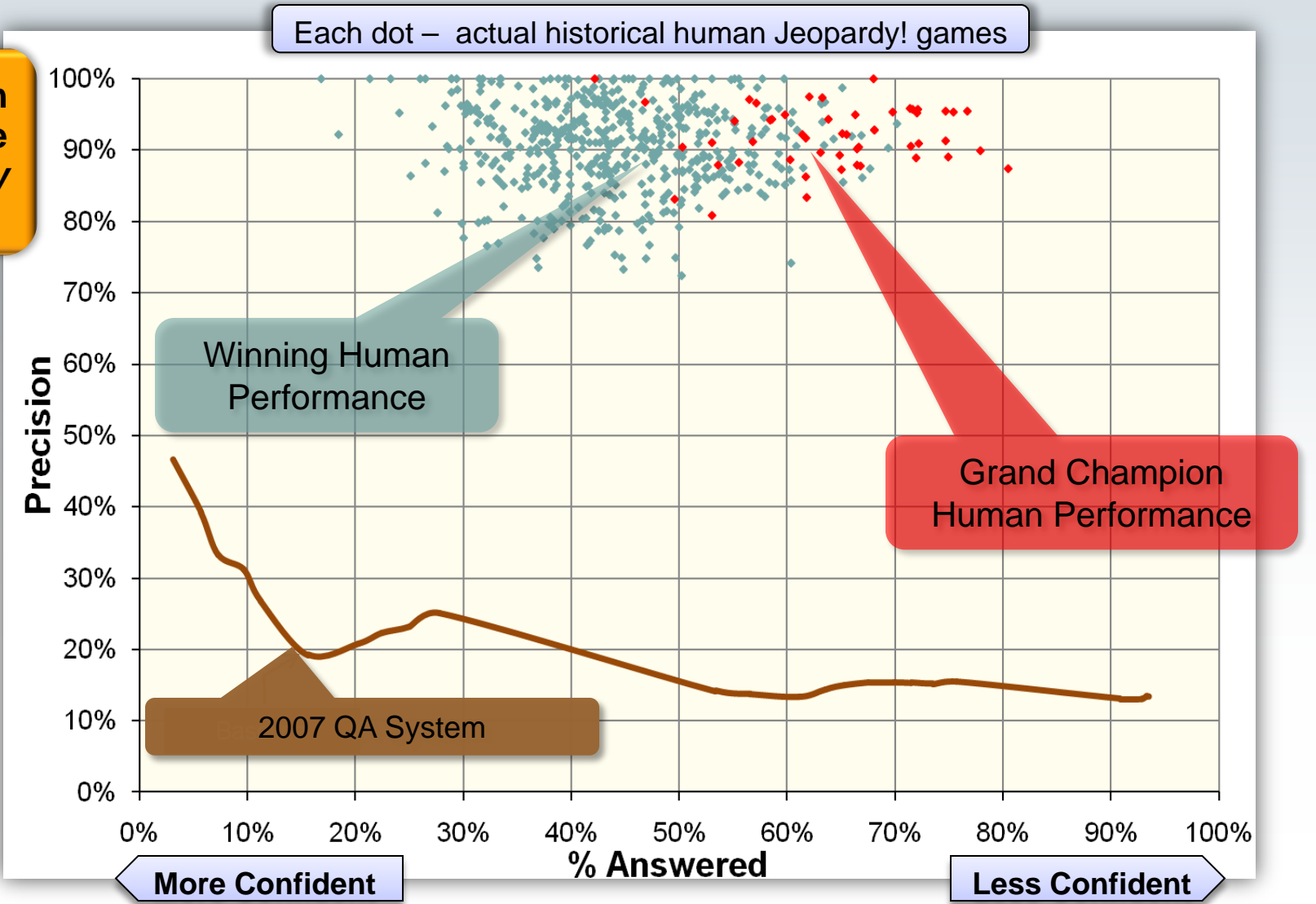


On hearing of the discovery of George Mallory's body, he told reporters he still thinks he was first.

# What It Takes to compete against Top Human Jeopardy! Players

*Our Analysis Reveals the Winner's Cloud*

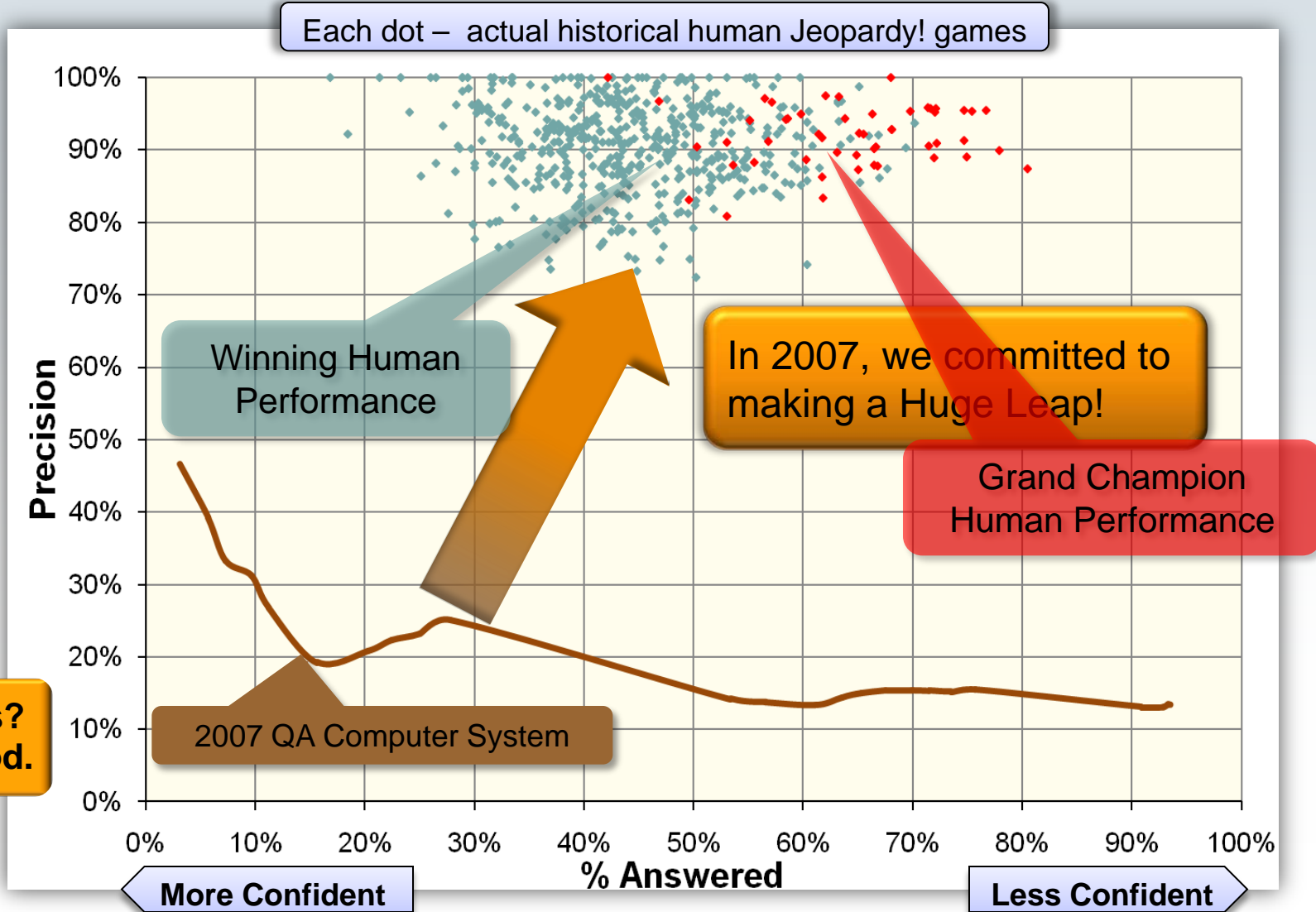
**Top human players are remarkably good.**



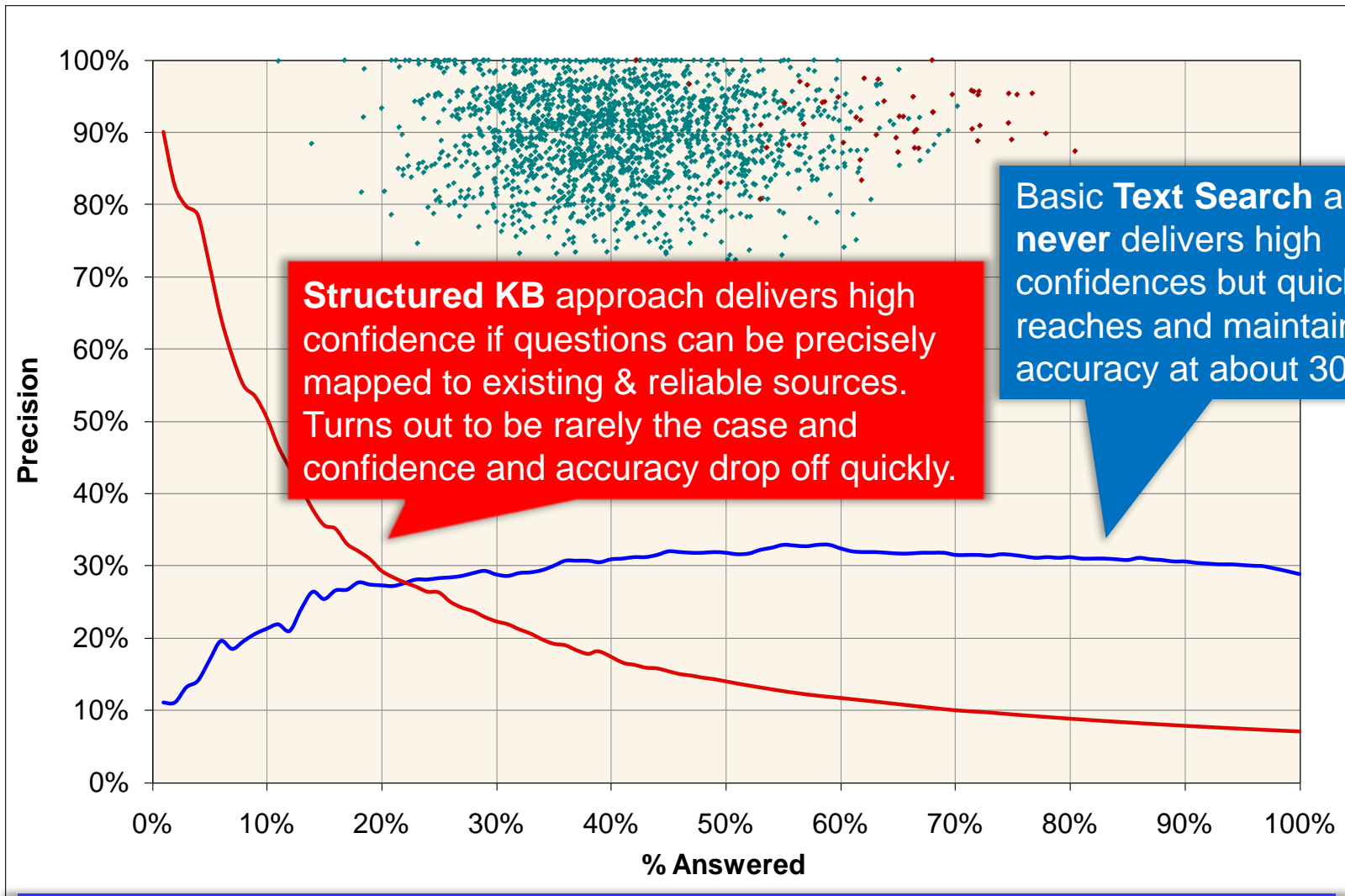


# What It Takes to compete against Top Human Jeopardy! Players

*Our Analysis Reveals the Winner's Cloud*



## Early Experiments on Focused Content



Must combine deep and shallow semantic analysis over structured & unstructured content to drive up precision, recall and confidence.

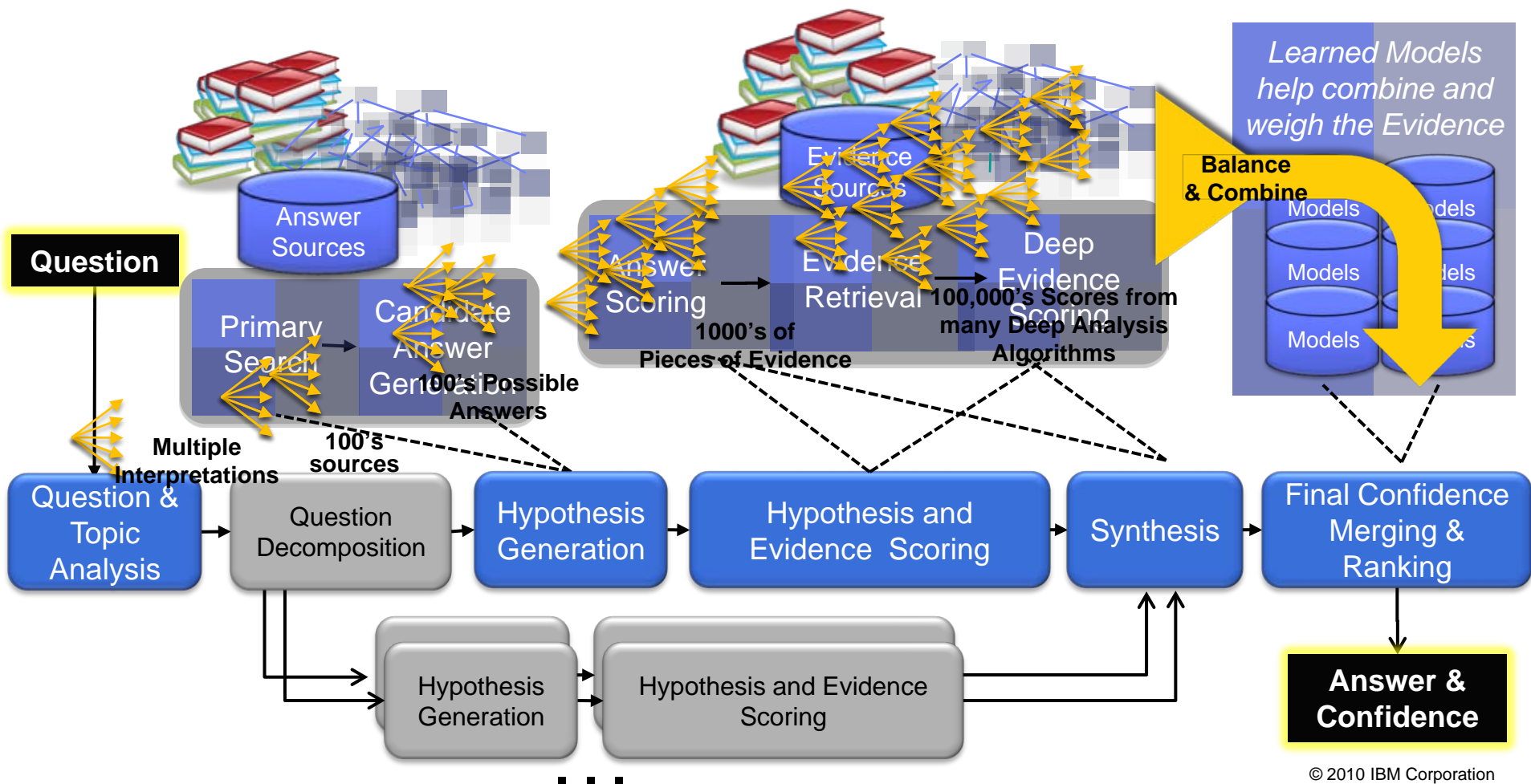
## A Few Guiding Principles

- **Specific Large Hand-Crafted Models Won't Cut It**
  - *Too Slow, Too Narrow, Too brittle, Too Bias*
  - *Need to acquire and analyze information from **As-Is Knowledge sources***
- **Intelligence from many diverse methods**
  - ***Many diverse algorithms must be combined.** No single one expected to solve the whole problem. Each addressing different weaknesses.*
  - *Relative impact of many overlapping methods must be **learned***
- **Massive Parallelism is a Key Enabler**
  - *Pursue many competing independent hypotheses over large data*
  - *Efficiency will demand simultaneous threads of evidence evaluation*

## Massively Parallel Probabilistic Evidence-Based Architecture

Generates and scores many hypotheses using a combination of 100's of **Natural Language Processing, Information Retrieval, Machine Learning and Reasoning Algorithms**.

These gather, evaluate, weigh and balance different types of **evidence** to deliver the answer with the best support it can find.



# A METHODOLOGY FOR RAPID ADVANCEMENT

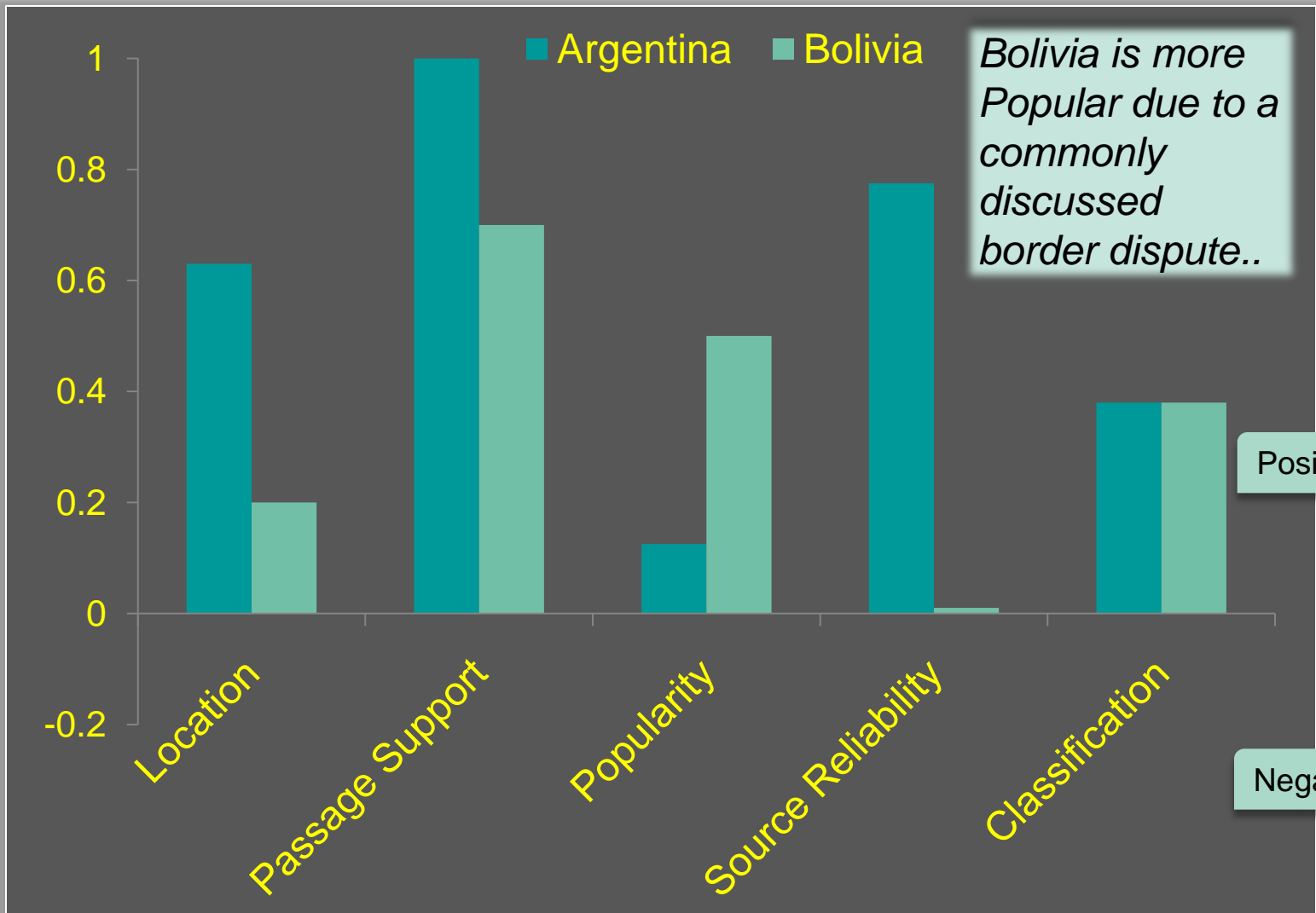
## Methodology

- **Goal-Oriented System-Level Metrics and Longer-Term Incentives**
  - Headroom Analysis: Estimate potential impact on System-Level metrics before “go”.
  - End-to-end evaluation vs. Isolated Component Metrics
  - 3-5 Year Big Team Result First vs. Incremental, Isolated Publications
  - Agile – Shift resources to project getting results
- **Extreme Collaboration**
  - Implemented “One Room” to optimize team work and communication
  - Immediate access to the right “expert”, spontaneous discussions
  - Real-time Documentation of ideas/status on team wiki
- **Diverse Skills ALL IN One Room**
  - Natural Language Processing, Information Retrieval, Knowledge Representation and Reasoning, Machine Learning, Computational Linguistics, Intelligent Systems Architecture, Software Engineering.
- **Disciplined Engineering and Evaluation**
  - Bi-weekly End-to-End Integration Runs & Evaluations (Large Compute Resources)
  - >10 GBs of error analysis output made accessible via Web-Based Tool
  - Feature Analysis Tools For Deeper Error Analysis

*>8000 Documented experiments performed in 4 years*

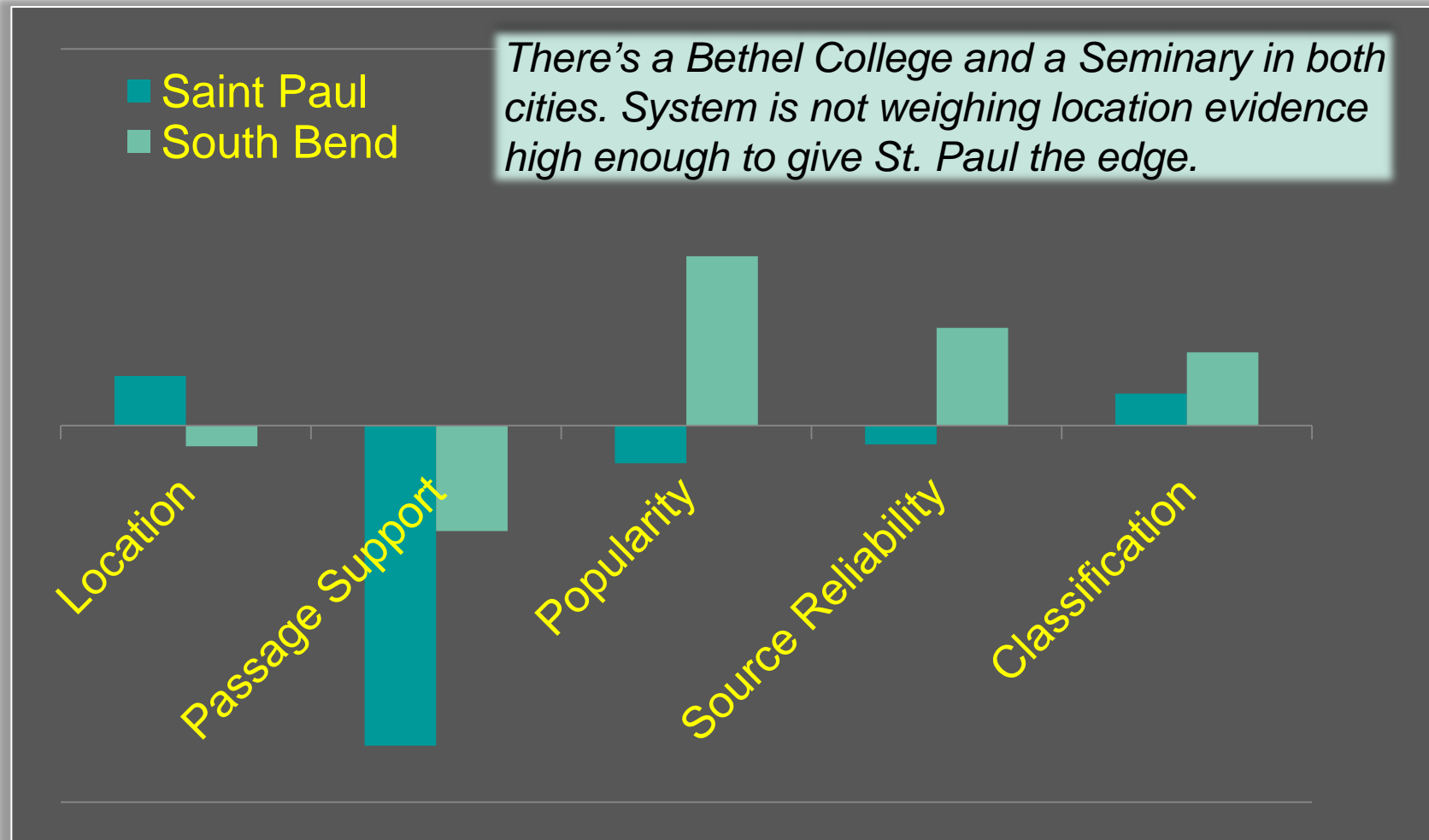
# Grouping Features to produce *Evidence Profiles*

**Clue:** Chile shares its longest land border with this country.



## Evidence: Time, Popularity, Source, Classification etc.

**Clue:** You'll find Bethel College and a Seminary in this "holy" Minnesota city.



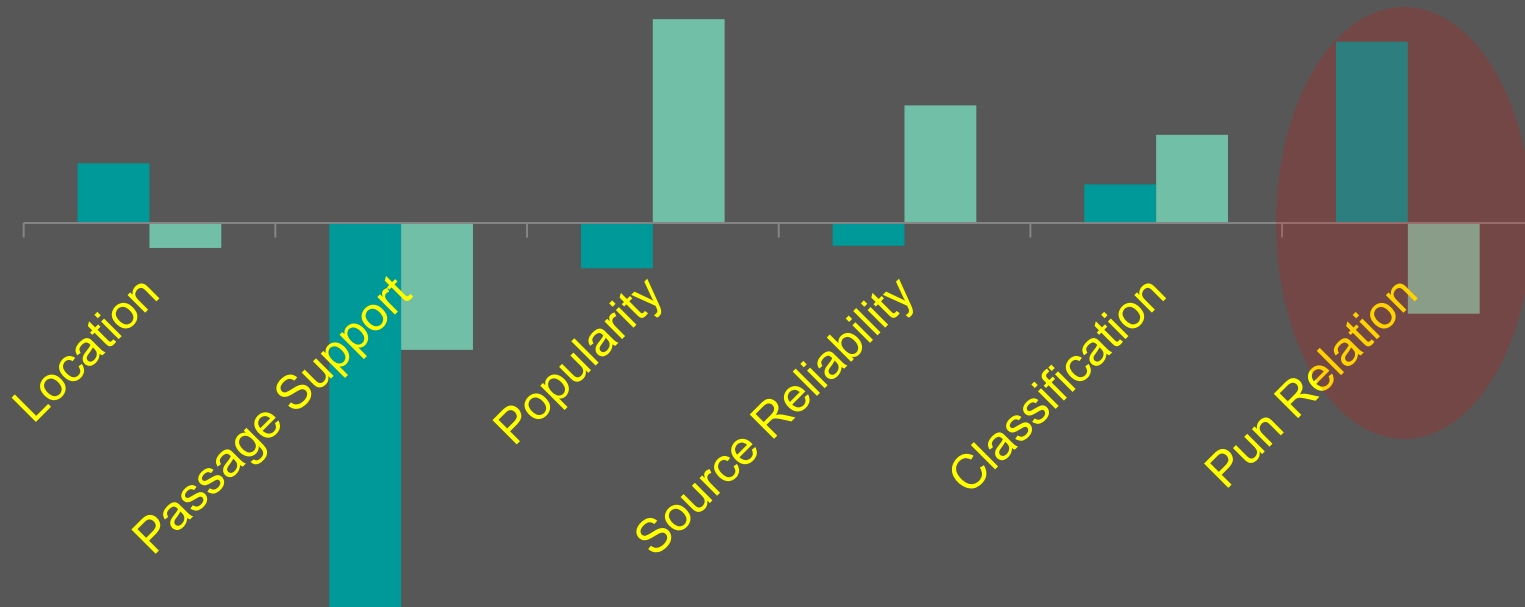


## Evidence: Puns

**Clue:** You'll find Bethel College and a Seminary in this "holy" Minnesota city.

- Saint Paul
- South Bend

*Humans may get this based on the pun since St. Paul since is a "holy" city. We added a Pun Scorer that discovers and scores Pun relationships.*



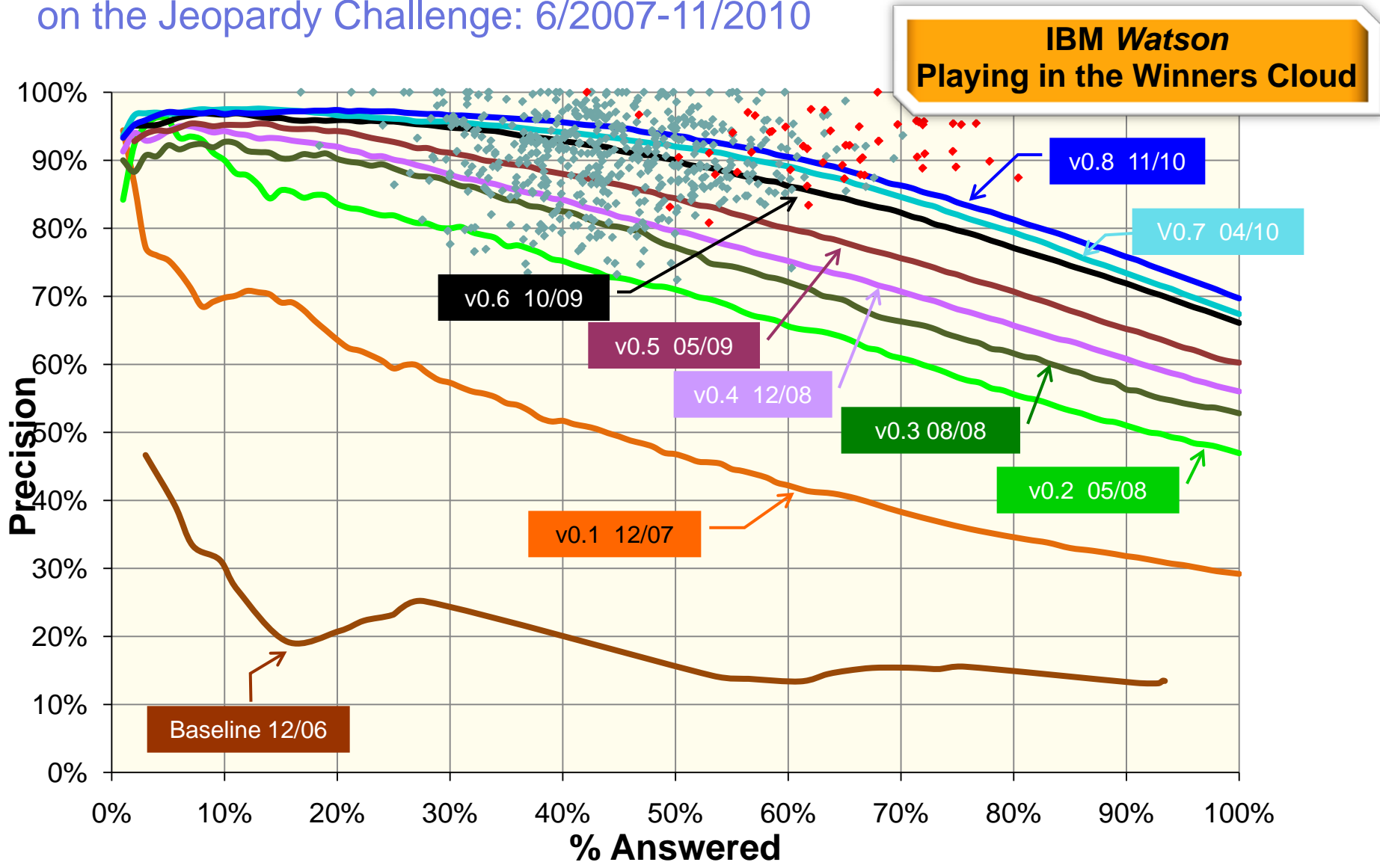
# In-Category Learning

## CELEBRATIONS OF THE MONTH

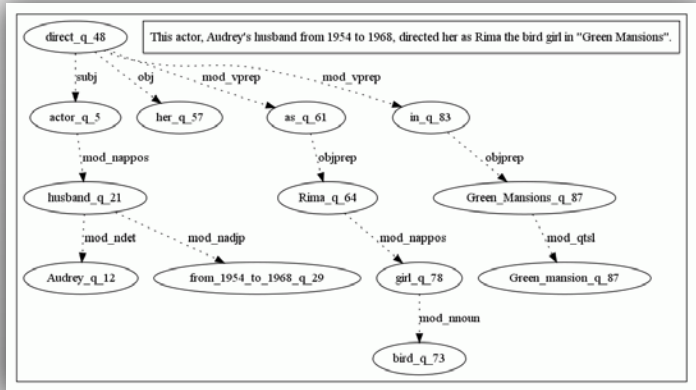
- What the Jeopardy! Clue is asking for is NOT always obvious
- Watson tries infer the type of thing being asked for from the previous answers.
- In this example after seeing correct answers, Watson starts to dynamically learn, using Bayesian inference, the probability that the answer type is a “**month**”.

Clue	Type	Watson's Answer	Correct Answer
D-DAY ANNIVERSARY & MAGNA CARTA DAY	day	Runnymede	<b>June</b>
NATIONAL PHILANTHROPY DAY & ALL SOULS' DAY	day	Day of the Dead	November
NATIONAL TEACHER DAY & KENTUCKY DERBY DAY	Day/month(.2)	Churchill Downs	May
ADMINISTRATIVE PROFESSIONALS DAY & NATIONAL CPAS GOOF-OFF DAY	day / month(.6)	April	April
NATIONAL MAGIC DAY & NEVADA ADMISSION DAY	day / month(.8)	October	October

# DeepQA: Incremental Progress in Answering Precision on the Jeopardy Challenge: 6/2007-11/2010

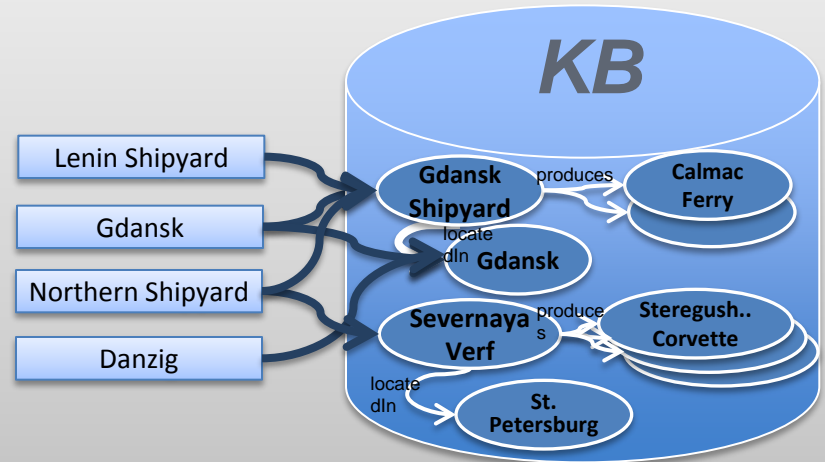


## Question Processing



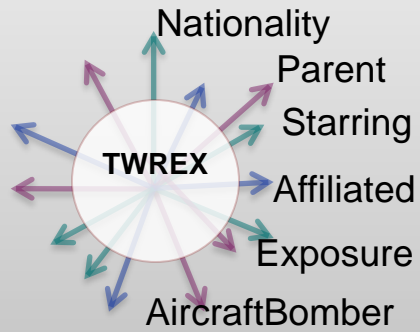
Dependency parse, Focus/LAT detection: 6000 rules, Decomp.  
 Eval on Jeopardy!: Parser: **92.4% acc**, Stat LAT detection: 96.8%

## KAFE: Knowledge From Extracted Content



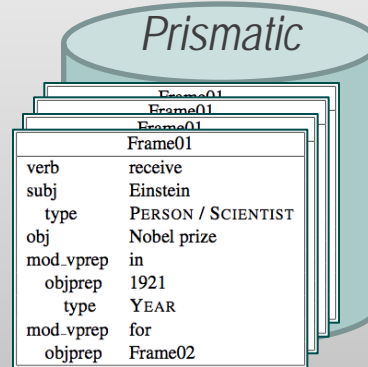
Entity Disambiguation, Entity Typing, Type Disambiguation  
 Eval on **Wikipedia Disambig Task**, F-score 92.5

## Relation Extraction



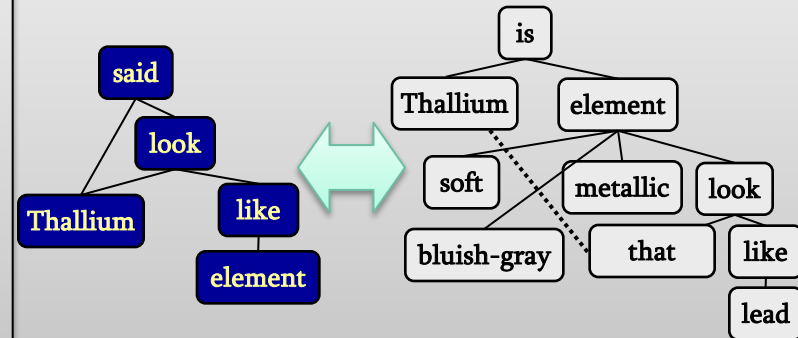
Semantic Relation Repository 7000 rels  
 Eval on ACE: **F-score 73.2** (leading score)

## Linguistic Frame Extraction



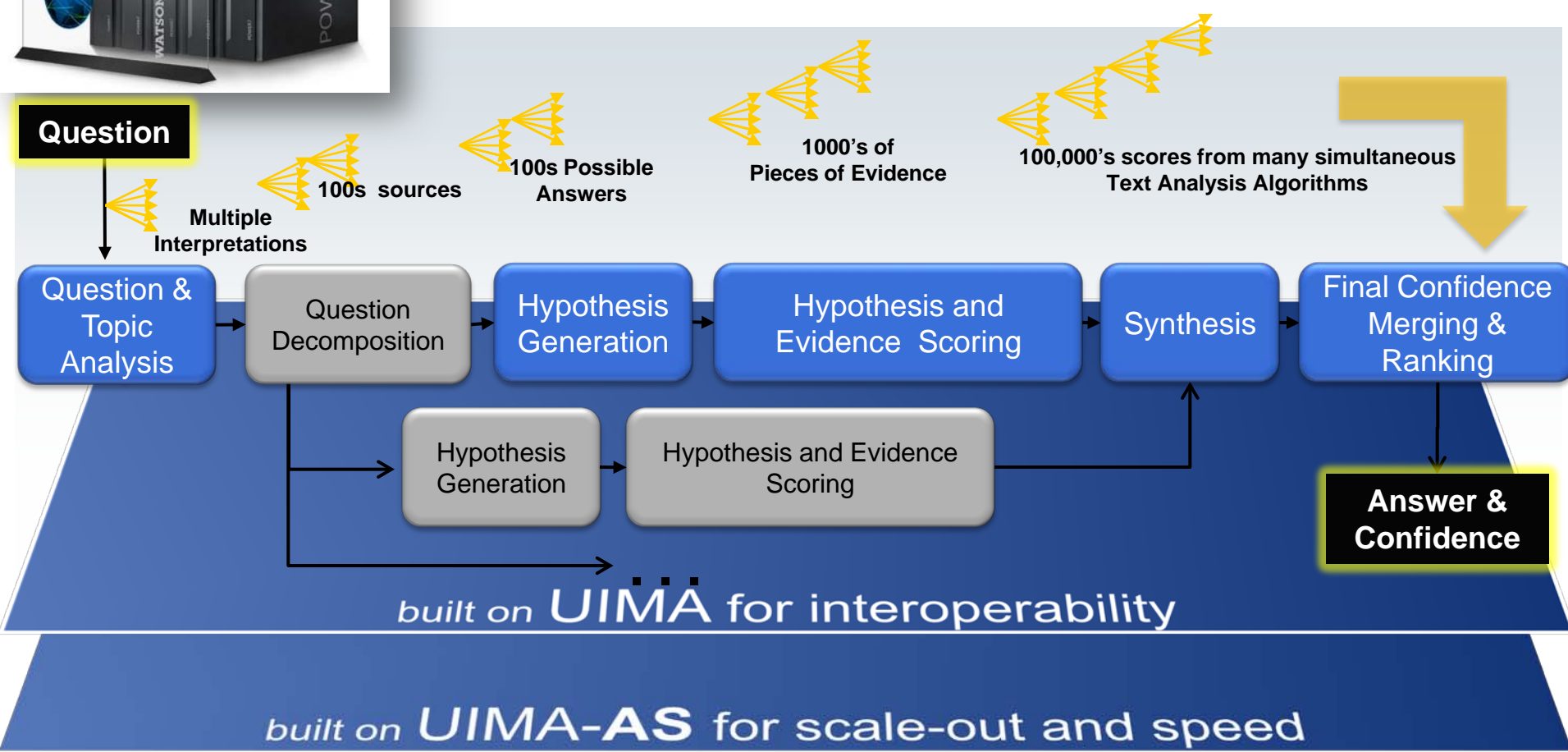
>1 Billion Frames. Mining from ClueWeb;  
 SVO/isa/etc. cuts,  
 Intensional/Extensional representation

## Passage Matching Ensemble



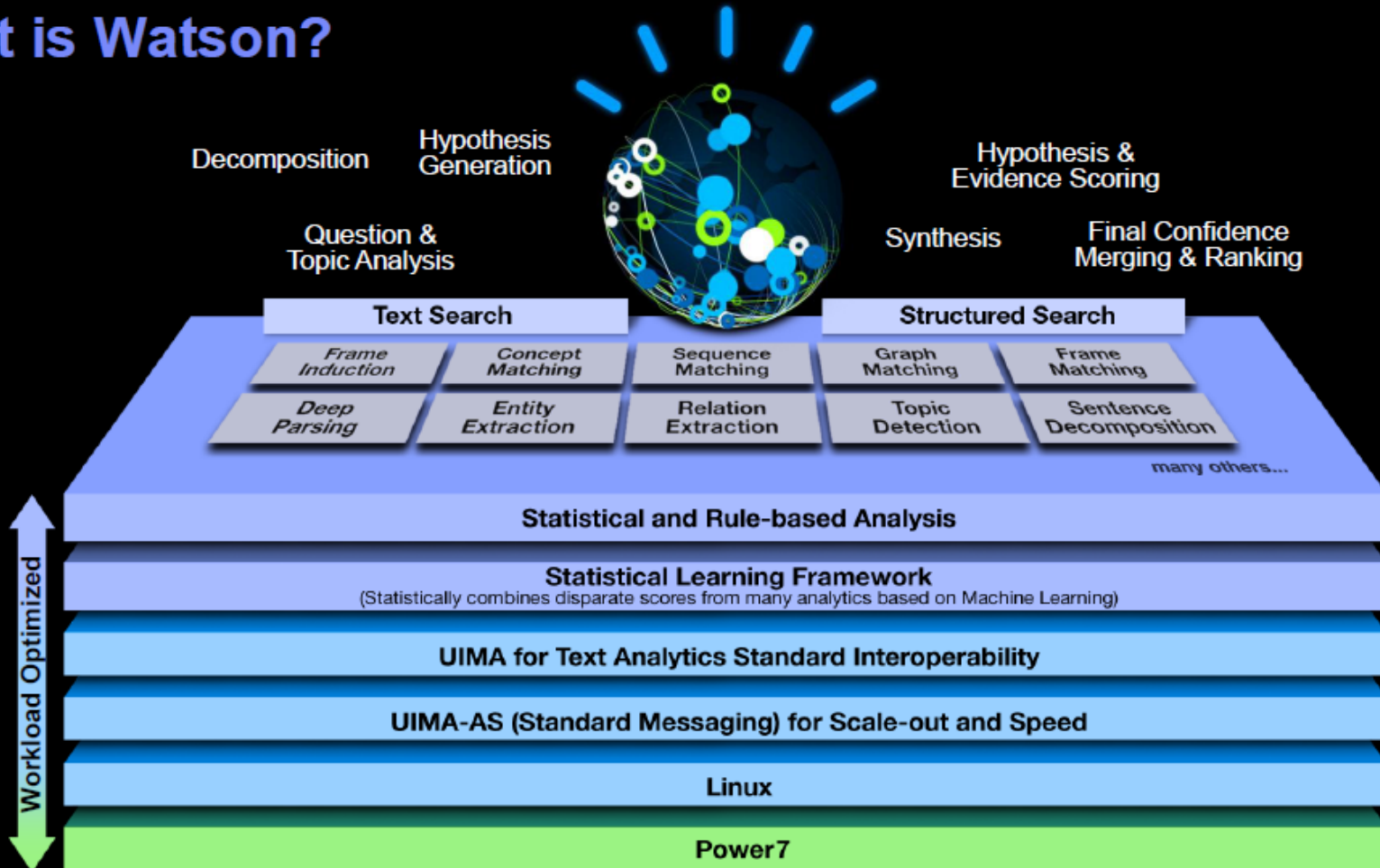
Synonymy, Temporal/Geographic Reasoning, Linguistic Axioms  
 Eval on RTE 2010 **Text Entailment**: F-score 48.8 (leading score)

One Jeopardy! question can take 2 hours on a single 2.6Ghz Core  
Optimized & Scaled out on 2880-Core IBM Power750's using UIMA-AS,  
*Watson* is answering in 2-6 seconds.



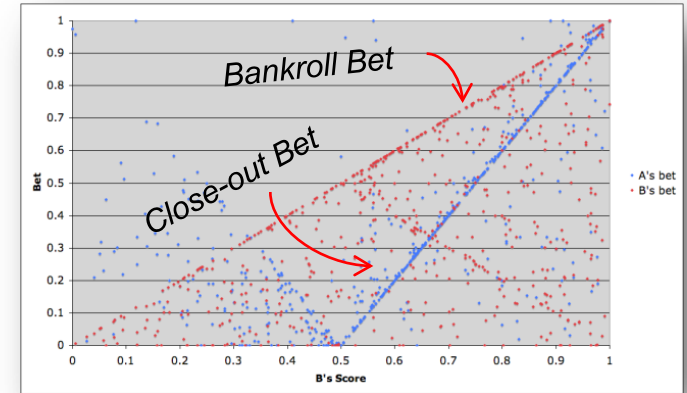
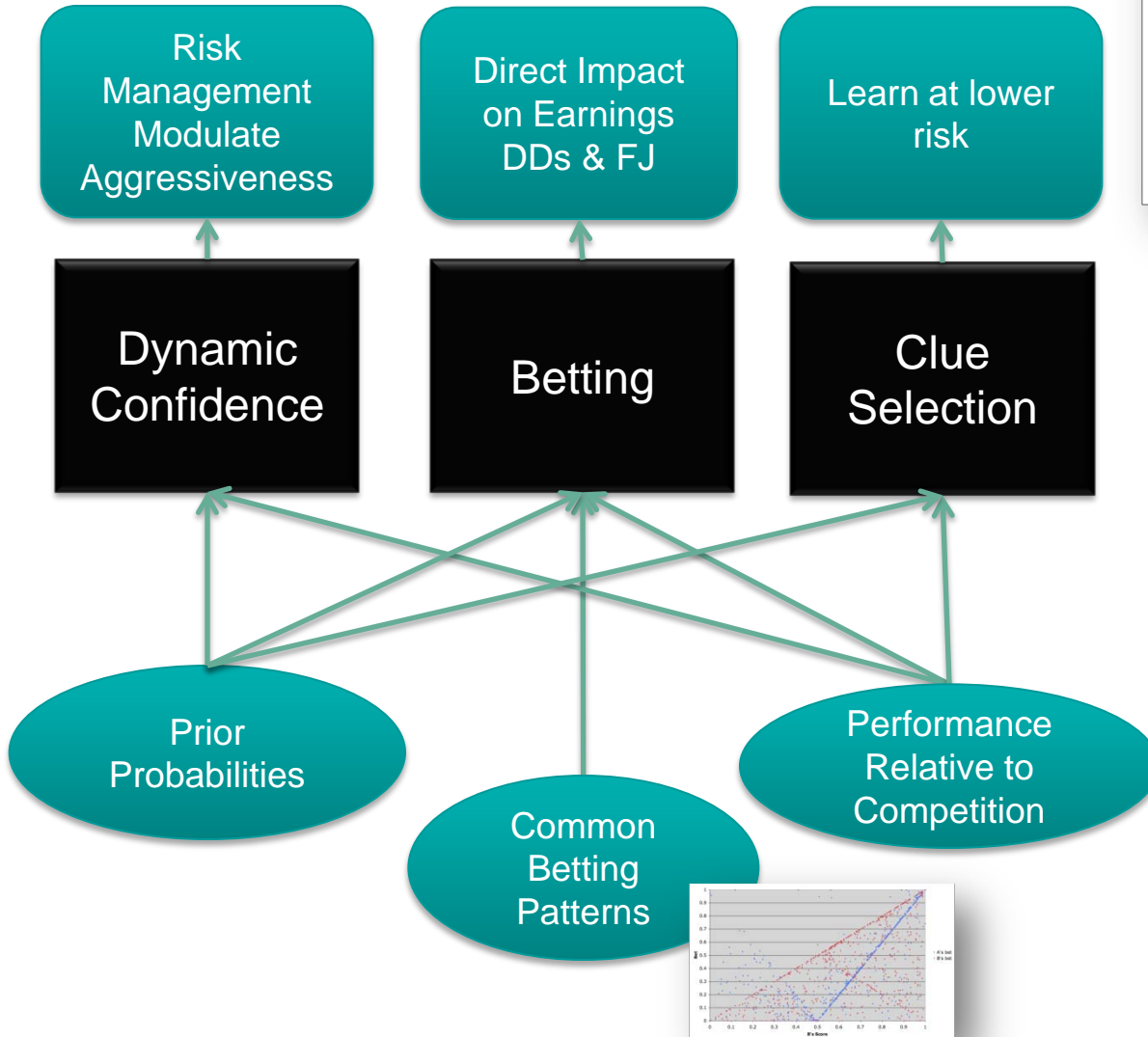
# Run-Time Stack

## What is Watson?



# Game Strategy

## Managing the Luck of the Draw



# With Precision, Accurate Confidence and Speed, the rest was History

The image shows a Jeopardy! game board with two contestants, IBM Watson (left) and a human player (right). The board features three columns of questions and answers. The central column has a neural network graphic above the question. The background is blue with the word 'THINK' and various translations. The contestants' names and scores are listed at the bottom.

Contestant	Score
Emily Dickinson	\$24,000
Walt Whitman	\$77,147
Barnard	\$21,600

Question	Answer
Who is Stoker? (FOR ONE WELCOME OUR NEW COMPUTER OVERLORDS)	\$1,000
Who is Bram Stoker?	\$17,973
WHO IS BRAM STOKER?	\$5600

Contestant	Confidence
Emily Dickinson	99%
Walt Whitman	60%
Barnard	10%



# Potential Business Applications



**Healthcare / Life Sciences:** Diagnostic Assistance, Evidenced-Based, Collaborative Medicine

**Tech Support:** Help-desk, Contact Centers



**Enterprise Knowledge Management and Business Intelligence**

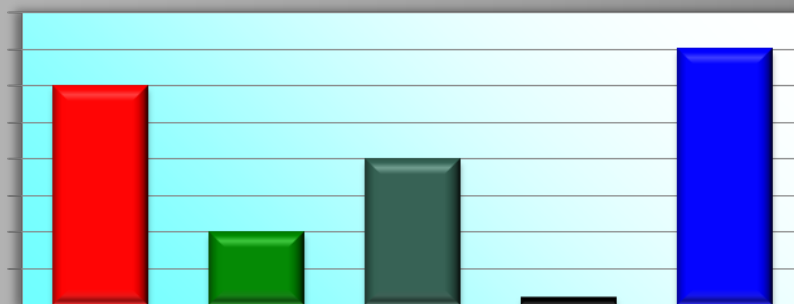
**Government:** Improved Information Sharing and Security



## *Evidence Profiles from disparate data is a powerful idea*

- Each dimension contributes to supporting or refuting hypotheses based on
  - **Strength of evidence**
  - **Importance of dimension for diagnosis** (learned from training data)
- Evidence dimensions are combined to produce an overall confidence

### Evidence Profile for UTI Diagnosis



Positive Evidence

Negative Evidence

### Overall Confidence

### Confidence

0 0.5 1

# DeepQA in Continuous Evidence-Based Diagnostic Analysis

IBM Research



## Analysis

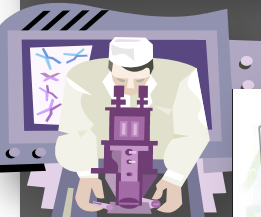
Considers and synthesizes a broad range of evidence improving quality, reducing cost



Symptoms



Family History  
Patient History  
Medications  
Tests/Findings



Notes/Hypotheses



Diagnosis Models	Symb	Fam	Hist	Meds	Find	Confidence
Renal failure	Low	Low	Low	Low	Low	Low
UTI	High	Low	Low	Low	High	High
Diabetes	Low	Low	Low	High	Low	High
Influenza	High	Low	Low	Low	Low	High
hypokalemia	Low	Low	Low	Low	Low	Low
esophagitis	Low	Low	Low	High	Low	Low

**Most Confident Diagnosis: Diabetes**

Huge Volumes of Texts, Journals, References, DBs etc.

# The Core Technical Team\*



Researchers and Engineers in **NLP, ML, IR, KR&R and CL** at  
IBM Labs and a growing number of universities

PI: David Ferrucci

## Systems & Speed

Eric Brown
Jerry Cwiklik
Pablo Duboue
Eddie Epstein
Tong Fin
Dan Gruhl
Bhavani Iyer
Adam Lally
Burn Lewis
Marshall Schor

## Core Algorithms

Eric Brown	Radu Florian	Dafna Sheinwald
Sugato Bagchi	David Gondek	Siddarth Patwardhan
Bran Boguraev	Aditya Kalyanpur	Kohichi Takeda
David Carmel	Hiroshi Kanayama	Yue Pan
Art Ciccolo	Adam Lally	John Prager
Jennifer Chu-Carroll	Tony Levas	Chris Welty
Bonaventura Coppola	Michael McCord	Wlodek Zadrozny
James Fan	Bill Murdock	Lei Zhang
David Ferrucci	Yuan Ni	Chang Wang
Achille Fokoue	Zhao Ming Qiu	

## Strategy

David Gondek
Jon Lenchner
Gerry Tesauro
James Fan
John Prager

## Speech

Andy Aaron
Raul Fernandez
Miroslav Novak
Andrew Rosenberg
Roberto Sicconi

## University Collaborations & Students

<b>Eric Nyberg (CMU)</b>	James Allen (UMASS)	Andy Schlaikjer (CMU)
<b>Nico Schlaefer (CMU)</b>	Ed Hovy (USC)	Saurav Sahay (GT)
<b>Manas Pathak (CMU)</b>	Bruce Porter (UT)	Rutu Mulkar-Mehta (USC)
<b>Hideki Shima (CMU)</b>	Pallika Kanani (UMASS)	Doo Soon Kim (UT)
Chang Wang (UMASS)	Boris Katz (MIT)	
Barbara Cutler (RPI)	<b>Alessandro Moschitti (Trento)</b>	

## Data Annotation

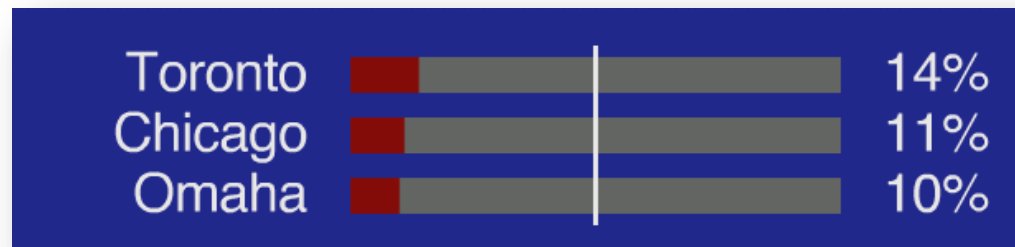
Karen Ingraffea
Matt Mulholland

The broader team that contributes to delivering Watson for the "Stage", to compete in Jeopardy Games is growing and changing rapidly.

Our Apologies if this is not up to date.

\*NOT full-time Equivalents. Names listed if contributed some time to that part of project.

# TORONTO?



## Categories are not as simple as they seem

Watson uses statistical machine learning to discover that **Jeopardy!** categories are only weak indicators of the answer type.

### U.S. CITIES

St. Petersburg is home to Florida's annual tournament in this game popular on shipdecks  
**(Shuffleboard)**

Rochester, New York grew because of its location on this  
**(the Erie Canal)**

### Country Clubs

From India, the shashpar was a multi-bladed version of this spiked club  
**(a mace)**

A French riot policeman may wield this, simply the French word for "stick"  
**(a baton)**

### Authors

Archibald MacLeish? based his verse play "J.B." on this book of the Bible  
**(Job)**

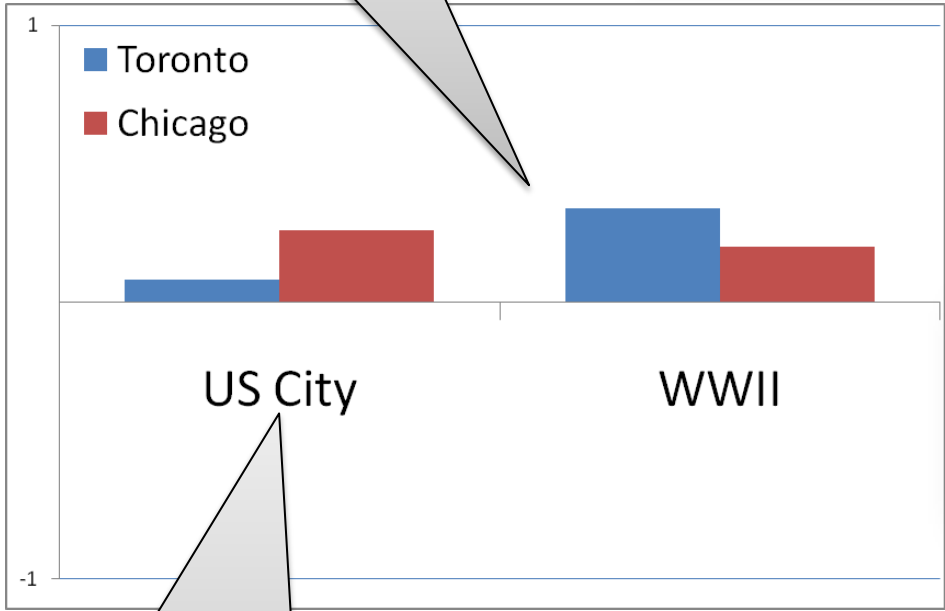
In 1928 Elie Wiesel was born in Sighet, a Transylvanian village in this country  
**(Romania)**

# Toronto vs. Chicago

## US CITIES

Its largest airport is named for a World War II hero; its second largest, for a World War II battle

Low because of weak evidence in content



Overall confidence was below threshold for both answers

Toronto		14%
Chicago		11%
Omaha		10%

Low because being a **US City** is not a strong requirement simply based on Jeopardy! category

When being a **US City** is a more trusted requirement, **Chicago** is in first position but still weak, because Watson is not convinced about the WWII evidence it could find.

## US CITIES

This **US City's** largest airport is named for a World War II hero; its second largest, for a World War II battle



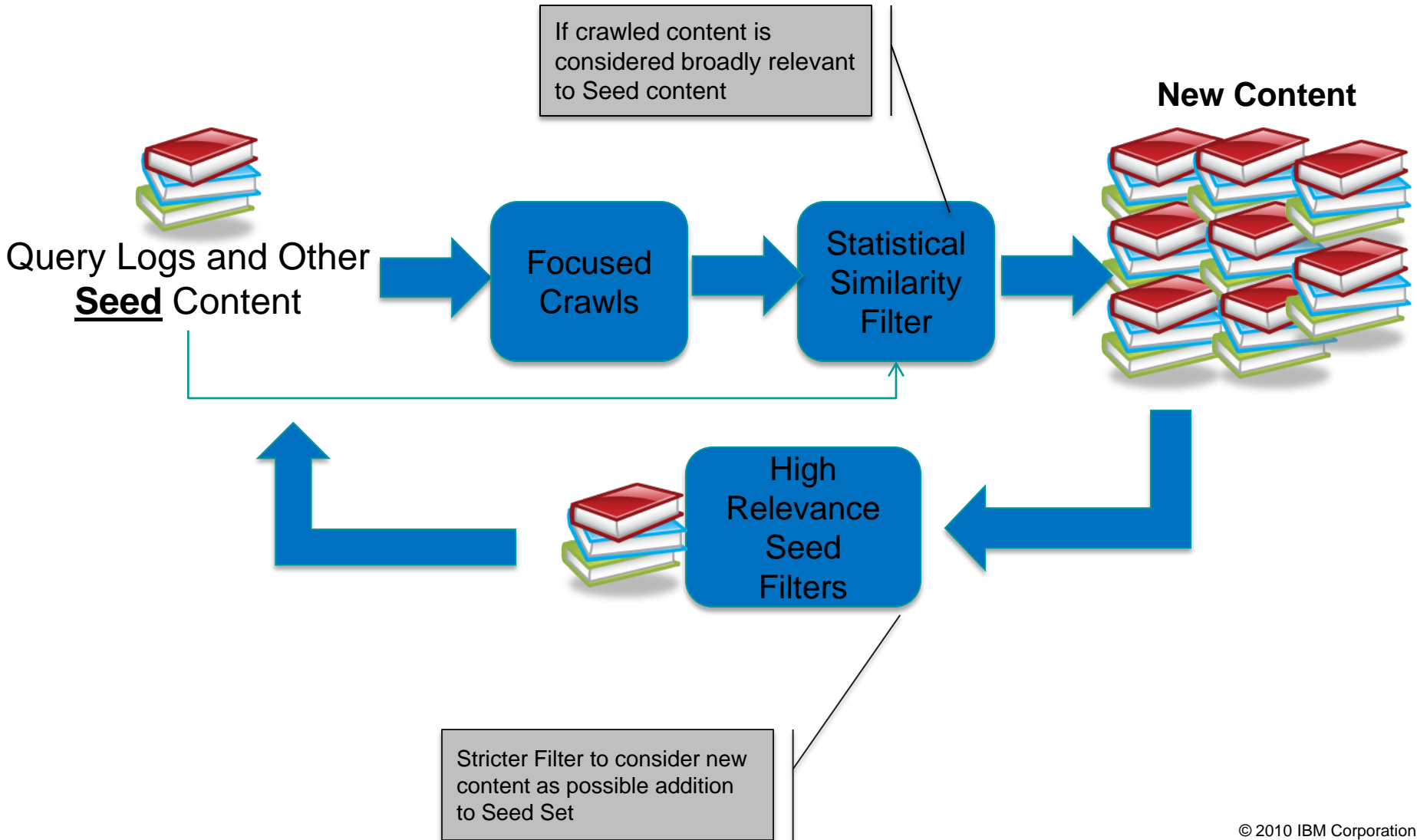


**THANK YOU**

# BACKUPS

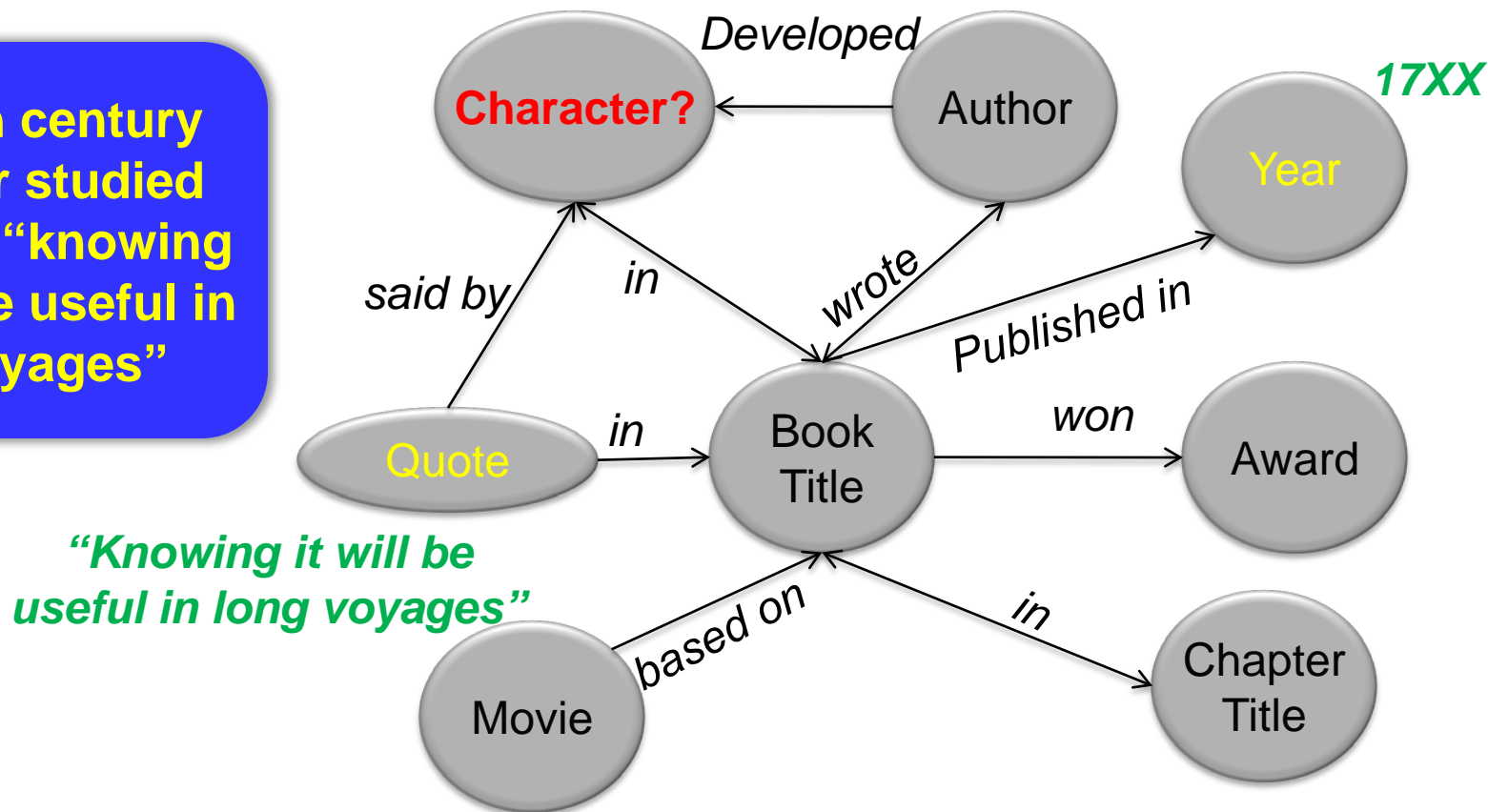
# Semi-Automatic Content Expansion

## Collaborative work with CMU



# Using Frames to Find Alternatives Paths

This 18th century character studied medicine, “knowing it would be useful in long voyages”



- No *said* relation indicating who uttered the **quote** BUT...
- Can discover the quote is contained in a specific **book**, however
- The **frame** tells us how to relate a **character** to a **book** to a **year** ...
- Filling out more of the frame can build confidence in selecting a character as the right answer