

MASTER
CLASS

POTENT
POTABLES

THE
“BOLD”
WAR

WATSON

TENNIS

AROUND
THE
WORLD

THE ARCHITECTURE OF WATSON HOW DOES IT WORK?

Grady Booch
IBM Fellow

(with deep appreciation to the entire Watson
team for access to their inner sanctum)

What Is Watson Not?

Watson is not

the beginning of Skynet

nor our new computer overlord

nor an advanced search engine

nor a fancy database retrieval system.

What Is Watson?

**Watson is a reasoning system
with a question and answer front end
that processes natural language
across structured and unstructured data
using deep analytic algorithms
that the system learns to combine
in optimal ways.**

How Does Watson Work?

Watson operates by

analyzing a question

generating hypotheses (forward chaining)

collecting evidence (backward chaining)

then presenting its results with scored

levels of confidence.

What Is Unique About Watson?

Watson is unique in that it

attends to heterogeneous sources

postulates multiple possible answers

considers evidence across multiple

dimensions and learns.

Who Is Watson?

**Watson does not
understand
nor does it think.**

What's So Hard About This Problem?

**Real language is filled with nuance, slang, and metaphor*.
Reasoning about open-ended problems requires inferring
context
and meaning
and relevance
from evidence that is often incomplete,
ambiguous, and sometimes contradictory.**

* David Ferrucci, Watson Principle Investigator

Opening the Curtain on Watson

Watson was

built by humans

and it's still just ones and zeros

down at the bottom.

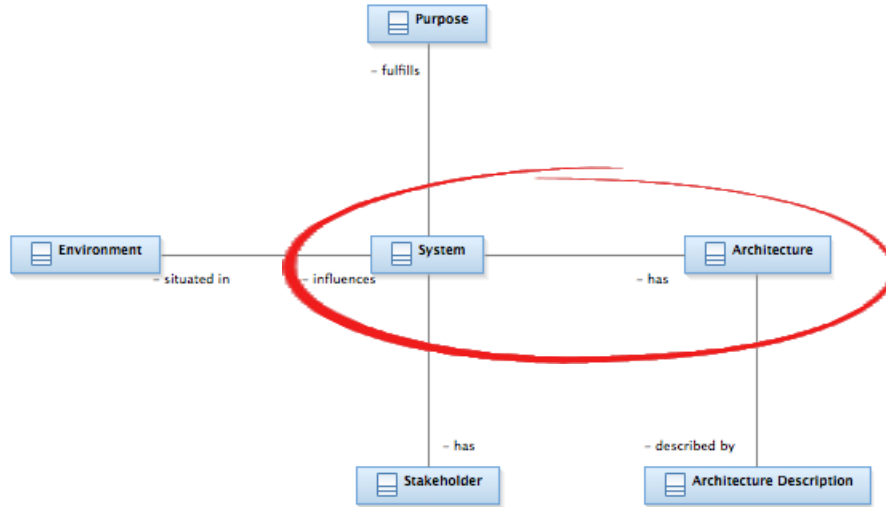
What Is Architecture?

- Architecture as **essence**.

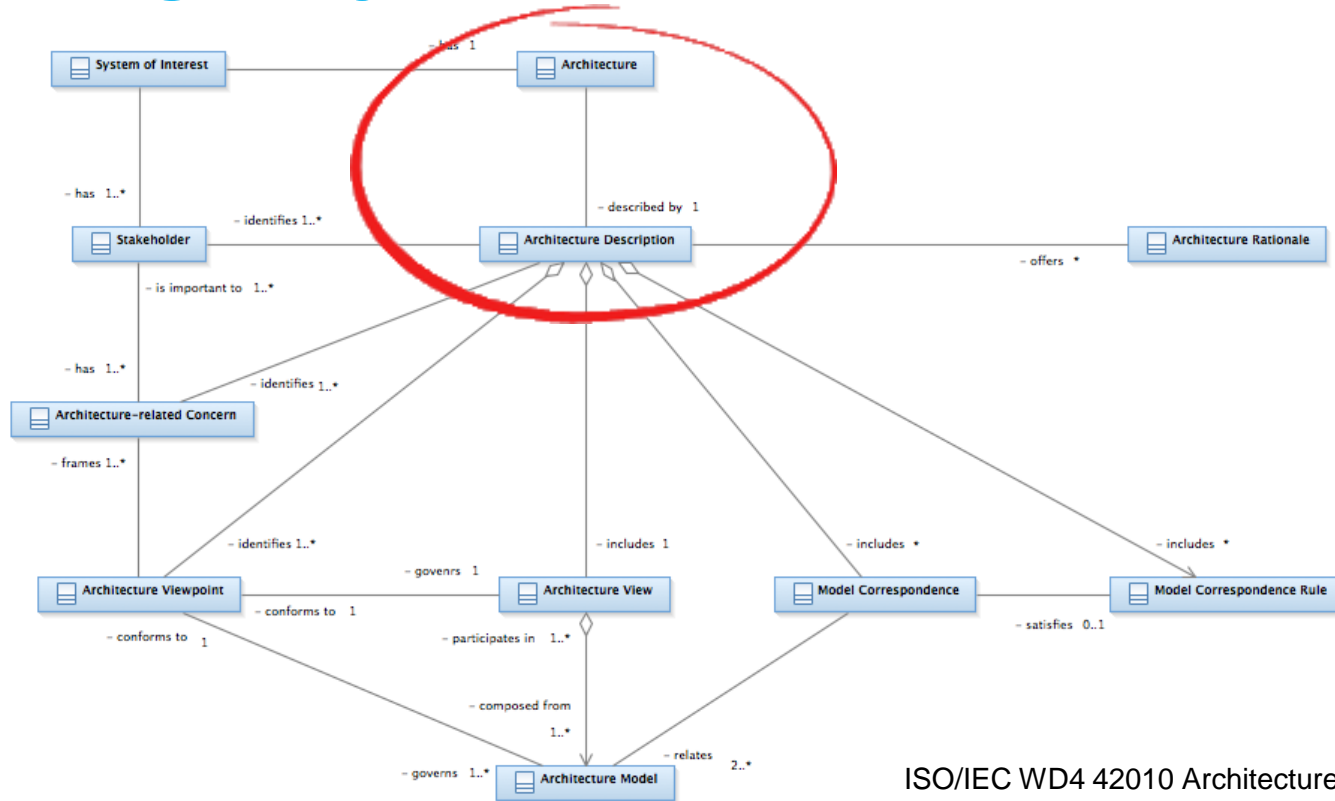
Architecture is the fundamental conception of a system in its environment embodied in elements, their relationships to each other and to the environment, and principles guiding system design and evolution.

- Architecture as **blueprint**.
- Architecture as **literature**.
- Architecture as **language**.
- Architecture as **decision**.

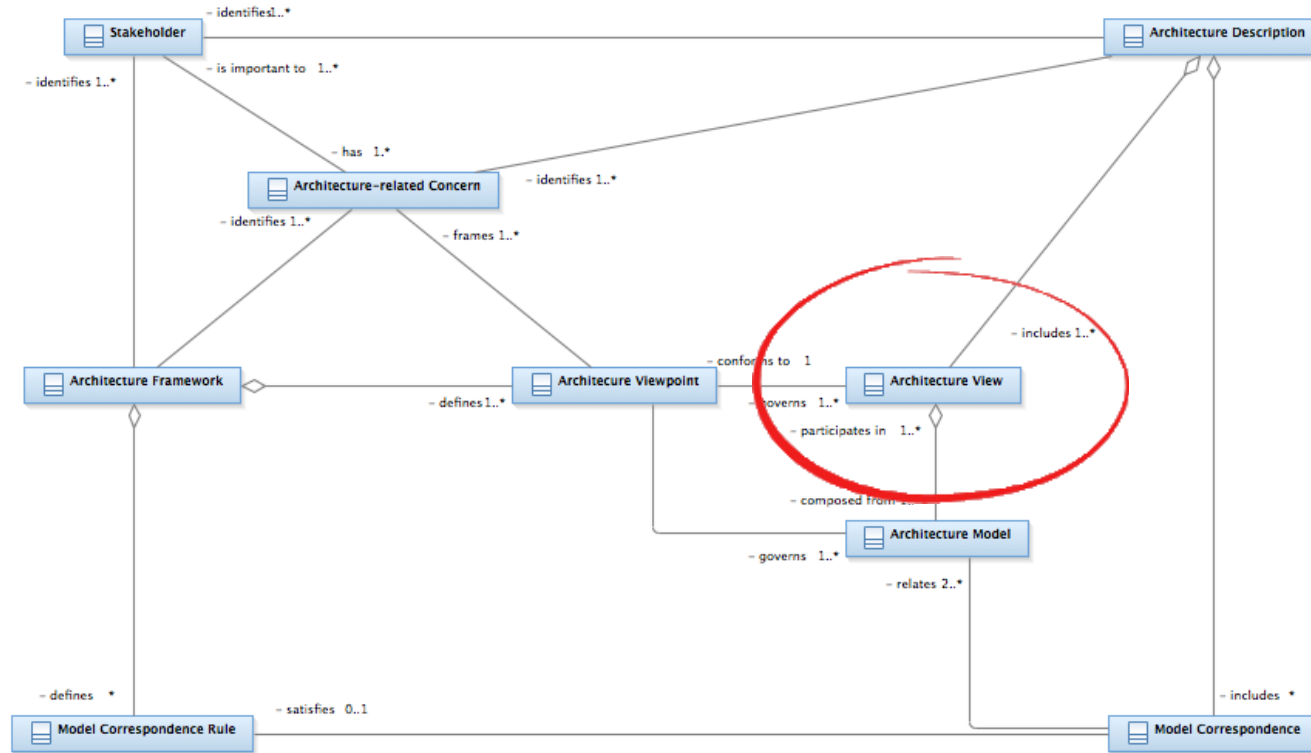
Describing A System's Architecture



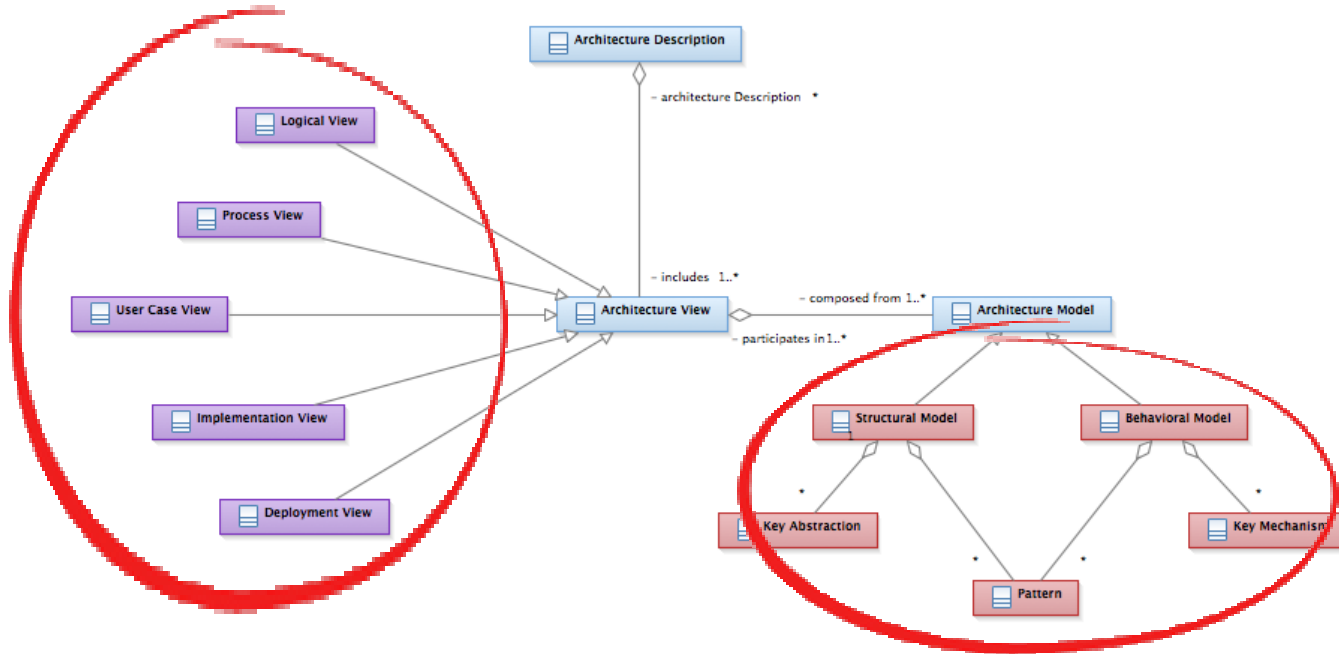
Describing A System's Architecture



Describing A System's Architecture



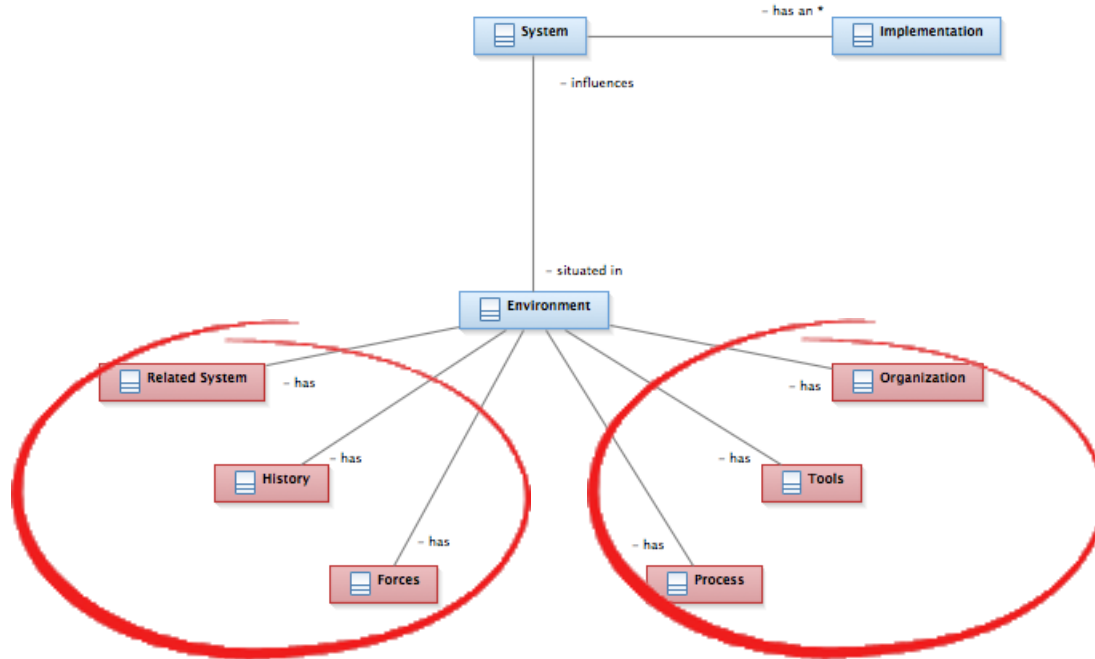
Describing A System's Architecture



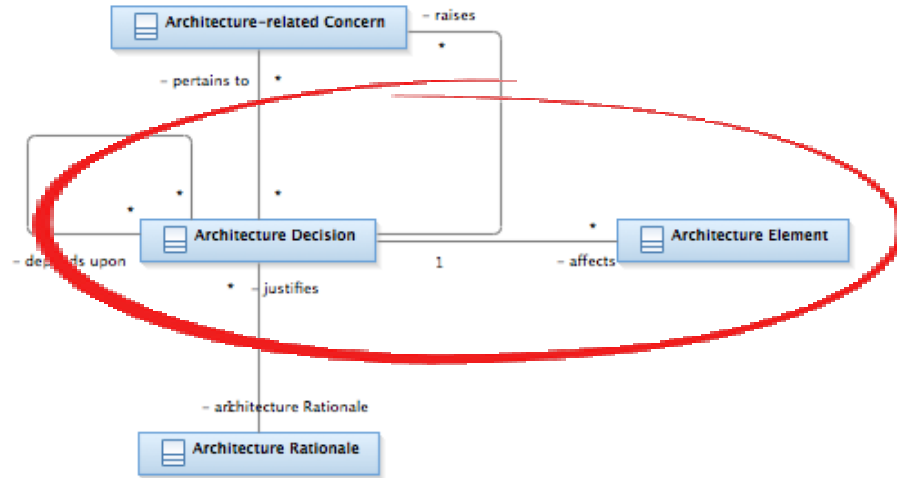
Kruchten, "The 4+1 View Model of Software Architecture"

Booch, *The Handbook of Software Architecture*

Describing A System's Architecture



Describing A System's Architecture



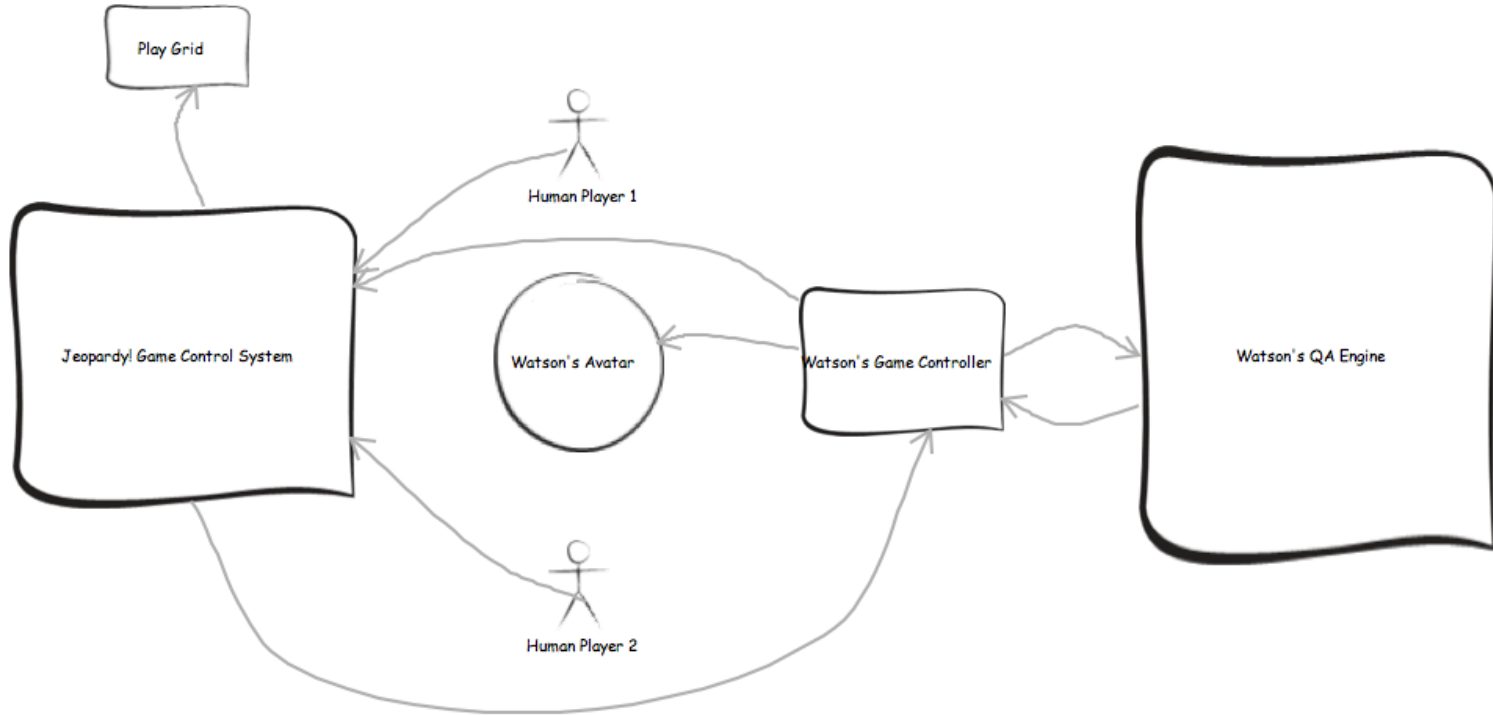
Conducting An Archeological Dig

- Become immersed in the domain.
- Absorb all relevant development documents.
- Interview the development team.
- Study the source code.
- Interpret the architecture.
- Use the architecture description in anger.
- Repeat.

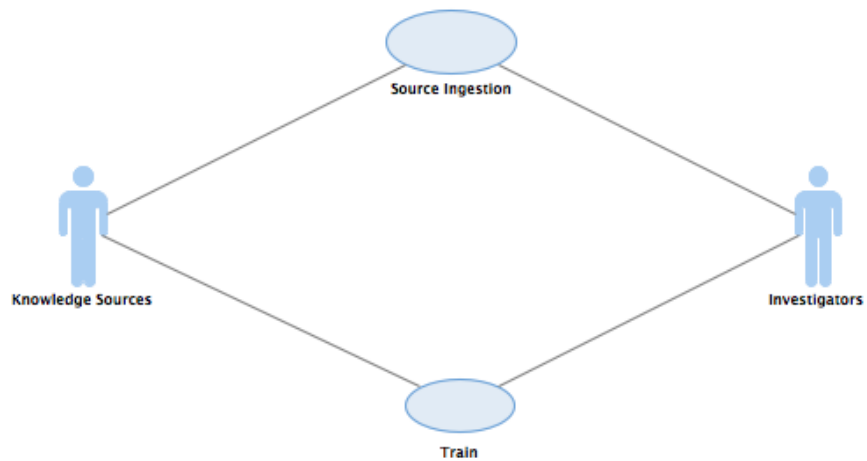
Statistics

Development Team	25 people
Project Duration	4 years
Software	1,000,000+ SLOC 700K Java, 300K C++, plus other bits ~ 130 components
Hardware	90 IBM Power 750 servers 2880 Power7 cores @ 80+ TFLOPS 20 TB memory 10 Gbps network

Architecture: Context



Use Cases



Key Design Decisions: Technical

- Use a pipe and filter architectural style.
- Acquire and apply heterogeneous data sources.
- Consider many possible candidate answers.
- Retrieve and evaluate multiple pieces of evidence in support of each candidate answer.
- Evaluate evidence along multiple dimensions.
- Combine evidence using machine learning.
- Build on UIMA-AS.

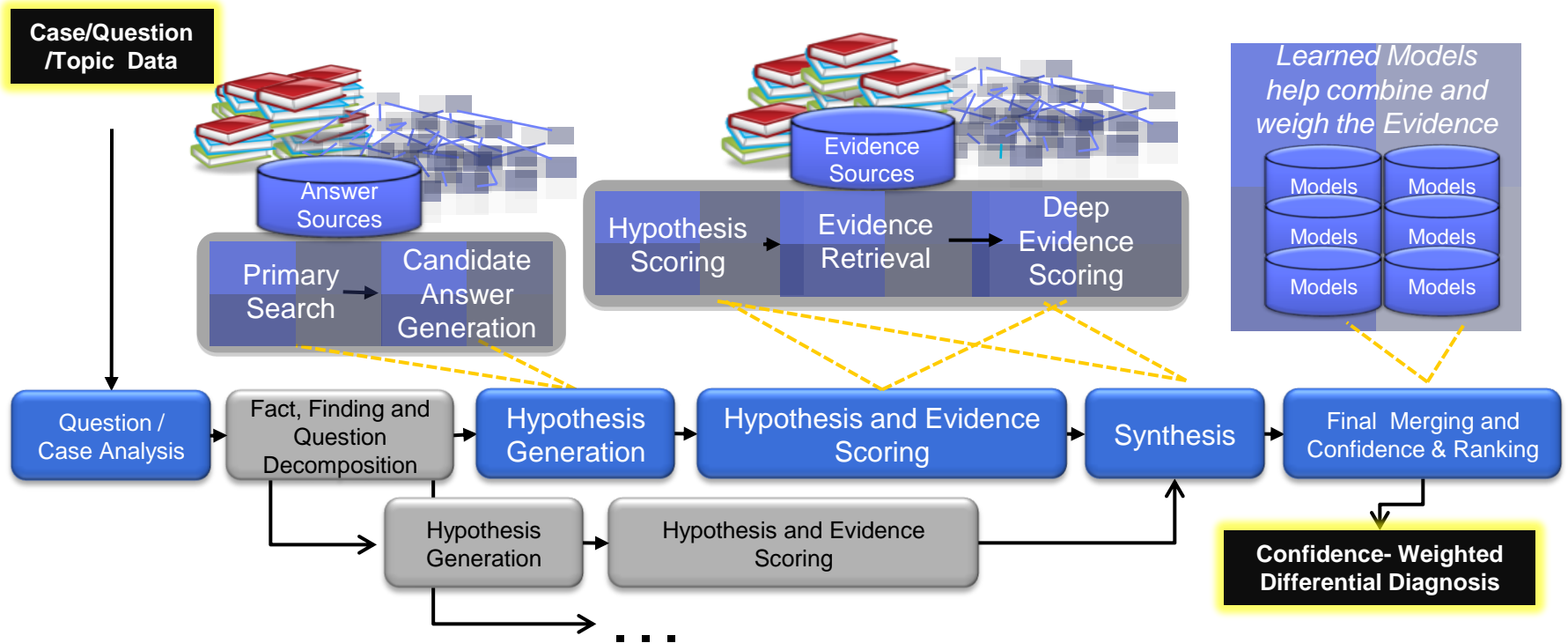
Key Design Decisions: Operational

- Permit the location of data to be configurable according to the needs of a specific deployment scenario.
- Directly map massive logical parallelism to massive deployment parallelism (but with mechanisms for easy reconfiguration).

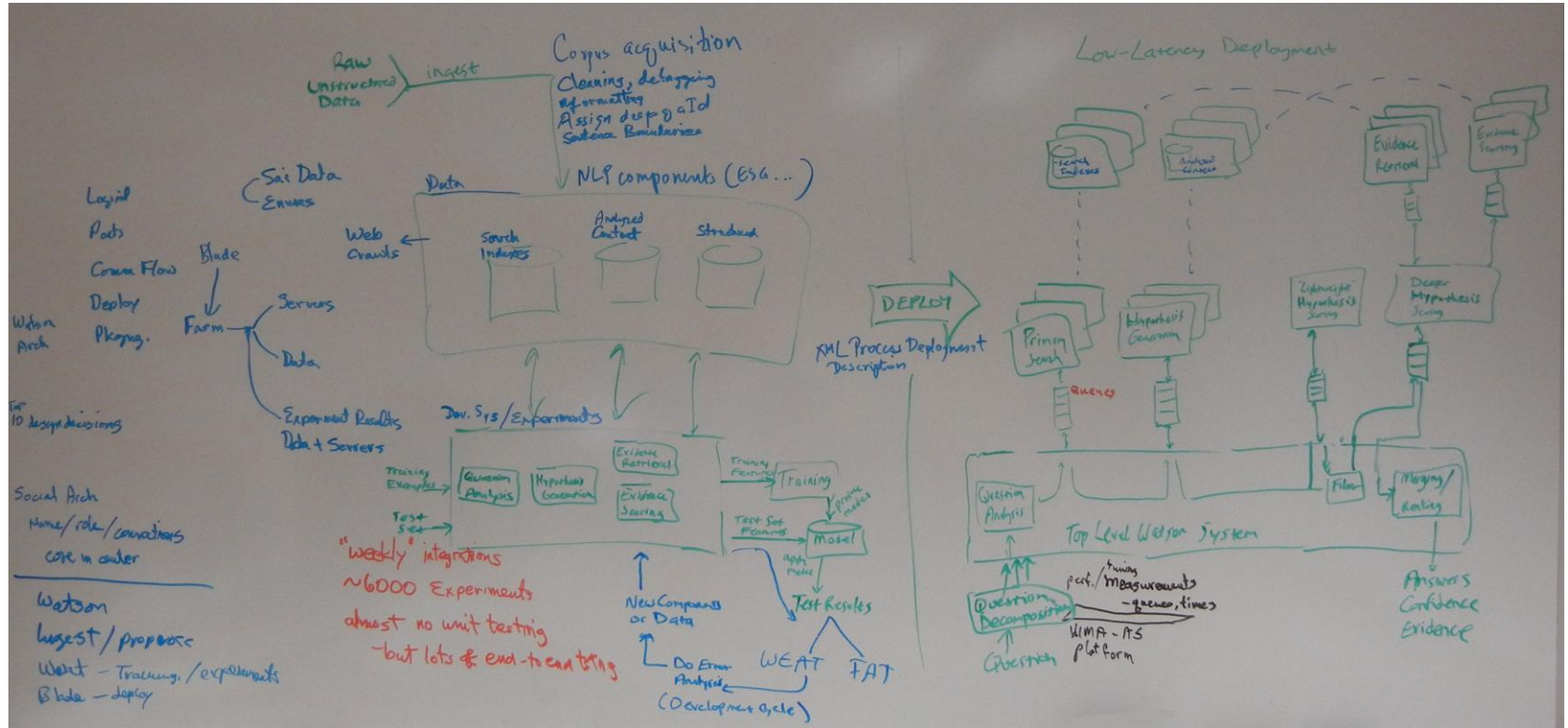
Key Design Decisions: Methodological

- Establish strong end-to-end metrics.
- Preserve considerable meta information for data, algorithms, and processes.
- Invest in tools to analyze Watson's operation.

Architecture: In Their Own Words



Architecture: In Their Own Words



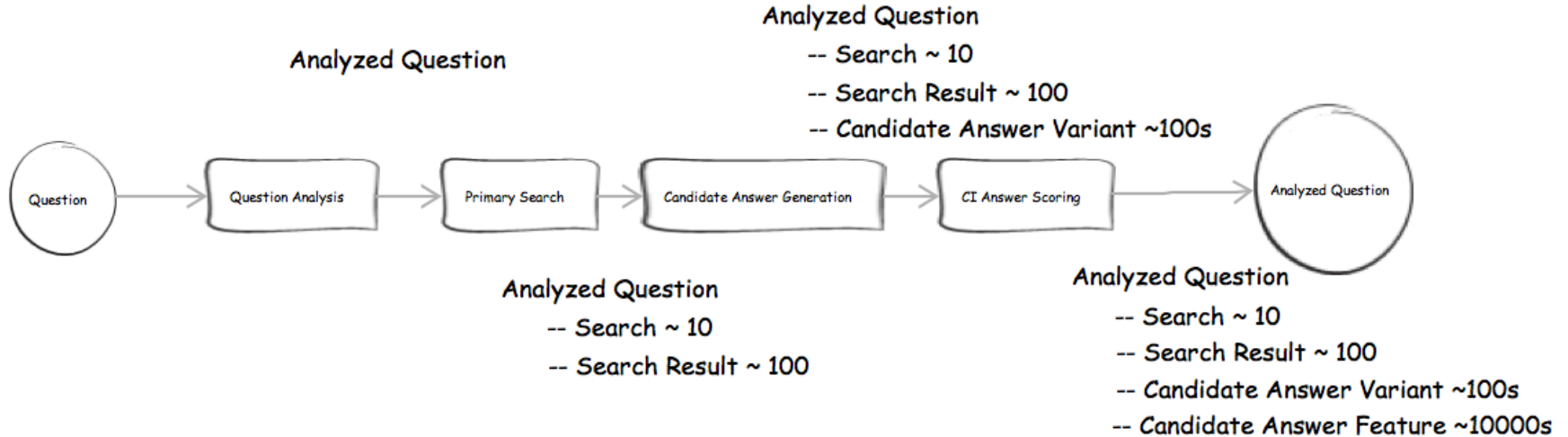
Sources

- Wikipedia/Wikiquote/Wiktionary/Wikibooks (The Free Encyclopedia) @ <http://wikipedia.org>
- YAGO2 (A Spatially and Temporally Enhanced Knowledge Base from Wikipedia) @ <http://www.mpi-inf.mpg.de/yago-naga>
- Dbpedia (Extracting Structured Information from Wikipedia) @ <http://dbpedia.org>
- WordNet (A Lexical Database for English) @ <http://wordnet.princeton.edu>
- Web expansion of many primary sources.
- Various licensed encyclopedias, dictionaries, books of quotations, and wire news.

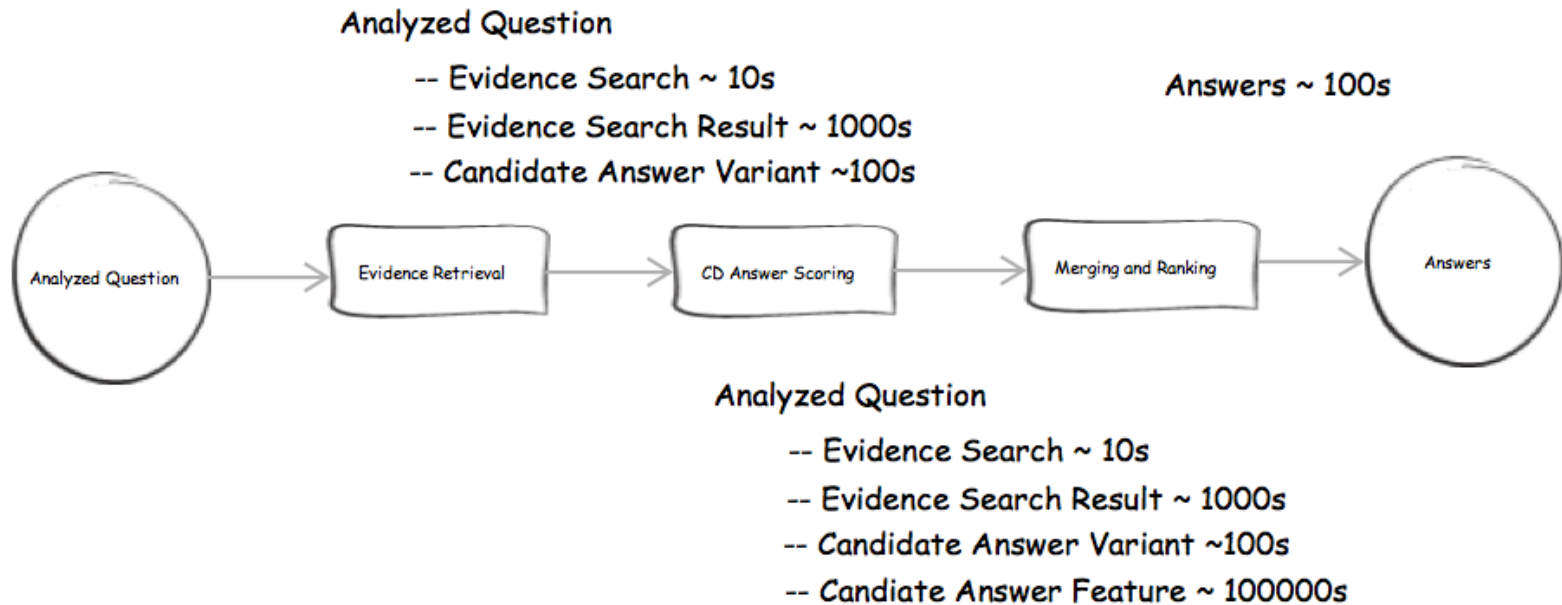
Test And Training Data

- J! Archive (A Fan-created Archive of Jeopardy!) @ <http://www.j-archive.com>

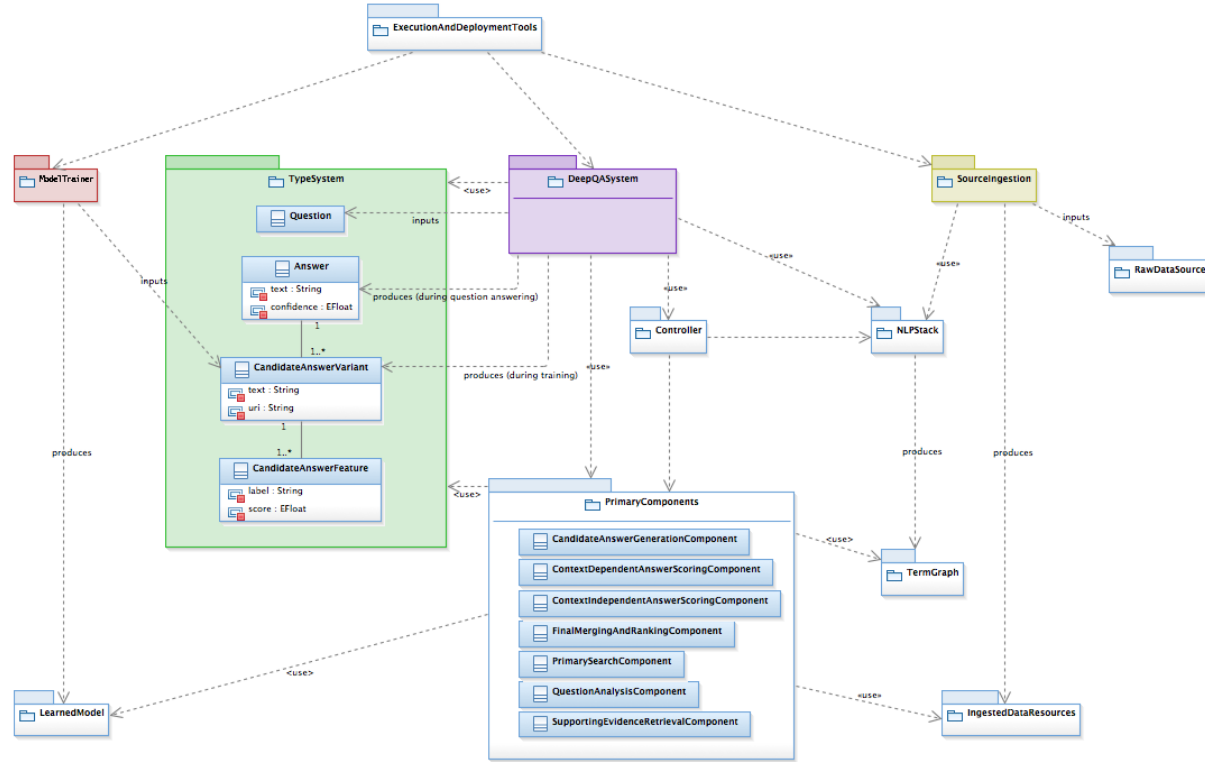
Logical View: Play Part 1



Logical View: Play Part 2



Logical View: Subsystems



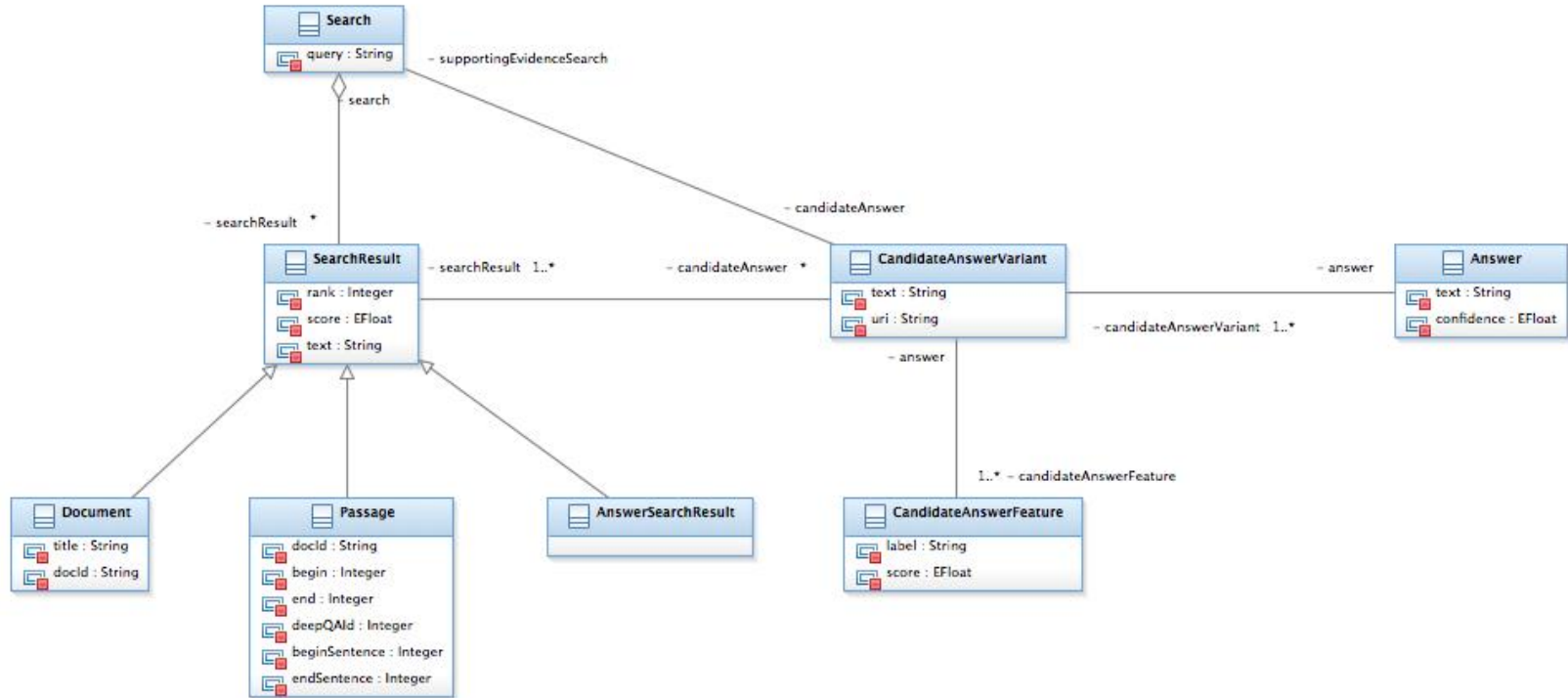
Logical View: Key Abstractions

- UIMA Common Analysis Structure (CAS)
- CoreTypeSystem (aka Data Model).
- QuestionAnalysisTypeSystem.
- Terms.
- IngestedDataResources.
- PrimaryComponents.
- ModelTrainer.

Logical View: CAS

- The common data structure shared by all UIMA analytics to represent unstructured information being analyzed (artifact) as well as metadata produced by the analysis workflow (artifact metadata), encompassing:
 - The artifact (the object being analyzed).
 - The subject of the analysis (one or more views of the artifact).
 - A type system description (including types, subtypes, and features).
 - Metadata (describing the artifact or a region of the artifact).
 - Index repository (supporting efficient access to and iteration over the results of analysis).

Logical View: Core Type System/Data Model

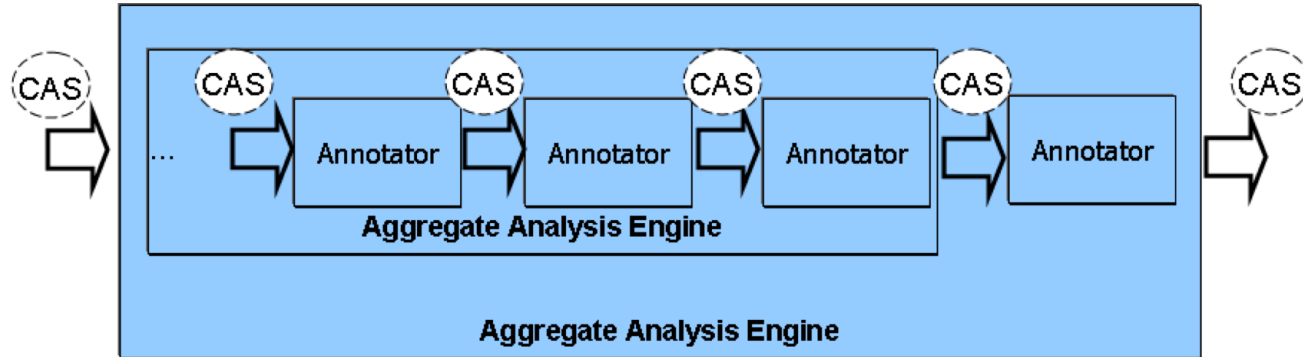


Logical View: Key Mechanisms

- UIMA.
- Question analysis.
- Primary search.
- Candidate answer generation.
- Shallow and deep scoring.
- Merging and ranking.

Logical View: UIMA

- UIMA is a software architecture which specifies component interfaces, data representation, design patterns, and development roles for creating, describing, discovering, composing, and deploying multi-modal analysis capabilities. The principal objective of the UIMA specification is to support interoperability among analytics.



Logical View: UIMA

- Watson uses over 100 annotators:
 - Basic (parsers, co-reference, dictionary lookup, named entity detectors, true casers).
 - LAT (2500 different lexical answer types).
 - Decomposition
 - Special question form identification.
 - Language translation.
 - Temporal logic.
 - Geospatial.

Logical View: UIMA

- Watson's annotators reveal different kinds of features:
 - Temporal (events and people happen during particular times and have likely life extents).
 - Location (events happen in particular places; places are located in or near other places).
 - Passage support (passages relate key entities to a candidate answer).
 - Shallow evidence (passage superficially aligns with question text).
 - Deep evidence (candidate answer is understood to be in the right logical relationship with key entities).

Logical View: UIMA

- Watson's annotators reveal different kinds of features:
 - Classification (answers should be the right type or class).
 - Popularity (answer is popularly associated with key elements in question)
 - Source reliability (sources supporting answer are learned to be reliable).
 - Predicate role (candidate answer plays the right role in key predicate).
 - Document support (document discusses fact in the context of the answer).
 - Hidden link (candidate answer and question entities share a common thread that is logically aligned).
 - Pun (candidate answer and question entities are associated with each other in one or more ways, such as sounds like, part of, synonym of, etc).

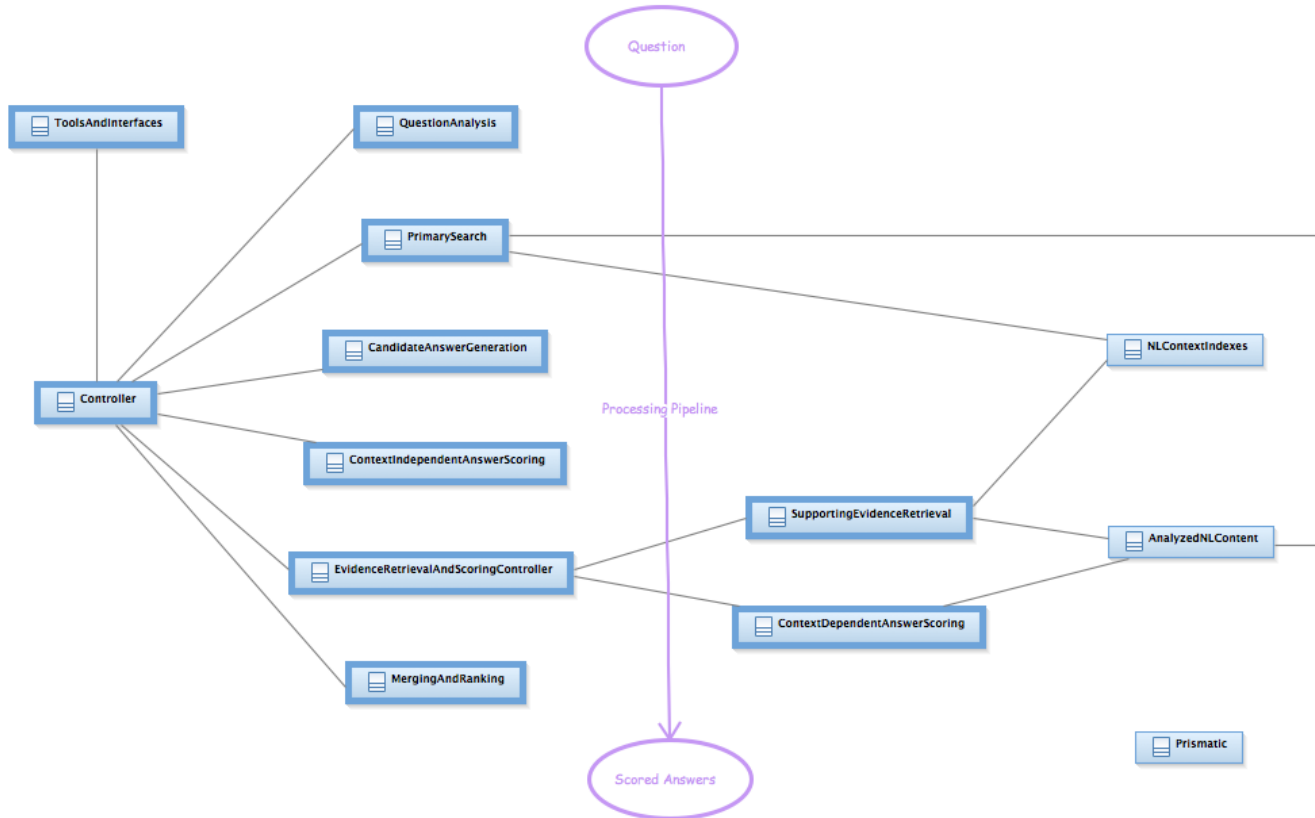
Logical View: Candidate Answer Generation

- Produces candidate answers without limitation on the type of answers produced.
 - Type coercion components attempt to relate the candidate answers to the lexical answer type.

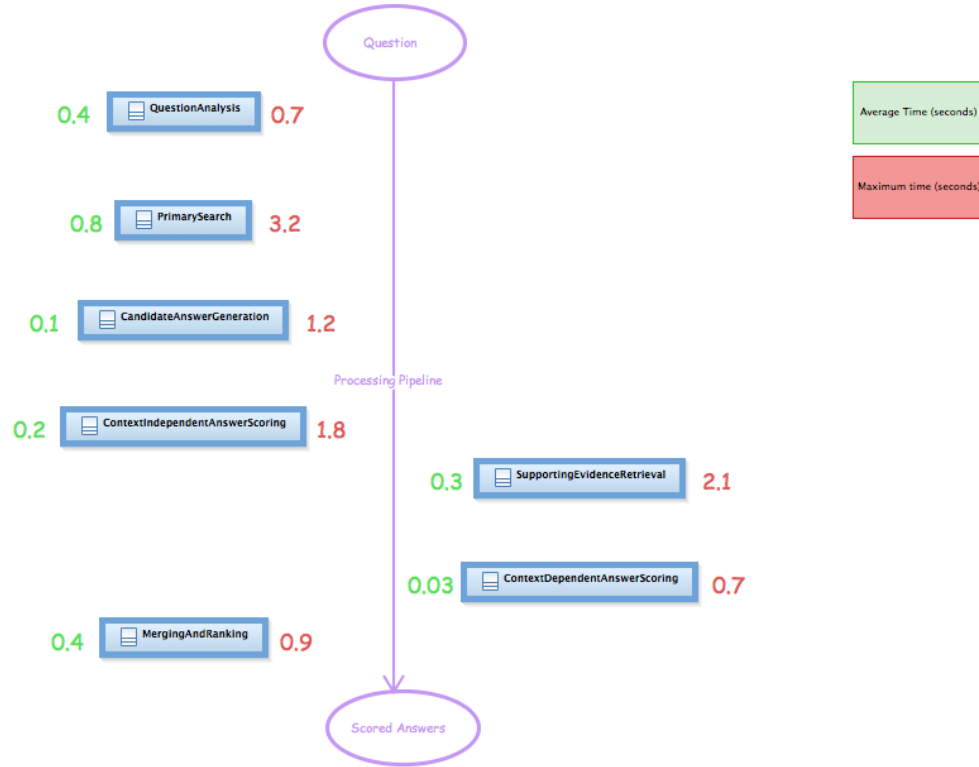
Logical View: Shallow And Deep Scoring

- Conducts both shallow (content independent and deep (context dependent) scoring.
 - Candidate answers that do not pass shallow scoring are passed directly to merging and ranking.
 - Candidate answers that pass shallow scoring continue to learned deep evidence retrieval and scoring, then passed on to merging and ranking.

Process View: Processing Pipeline



Process View: Timing



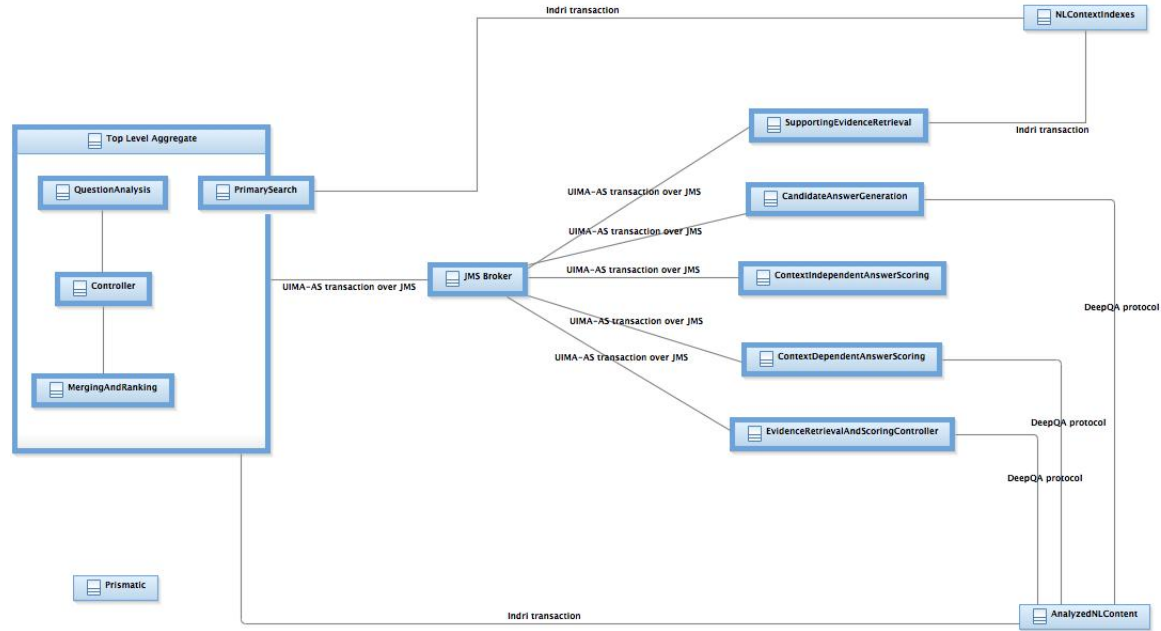
Process View: Key Abstractions

- Concurrent-aware POJOs.
 - May be mixed and matched as to location and grouping.
 - ~70 Mb data flow per question within UIMA-AS transactions.
 - ~140 Mb data flow per question within Indri transactions.

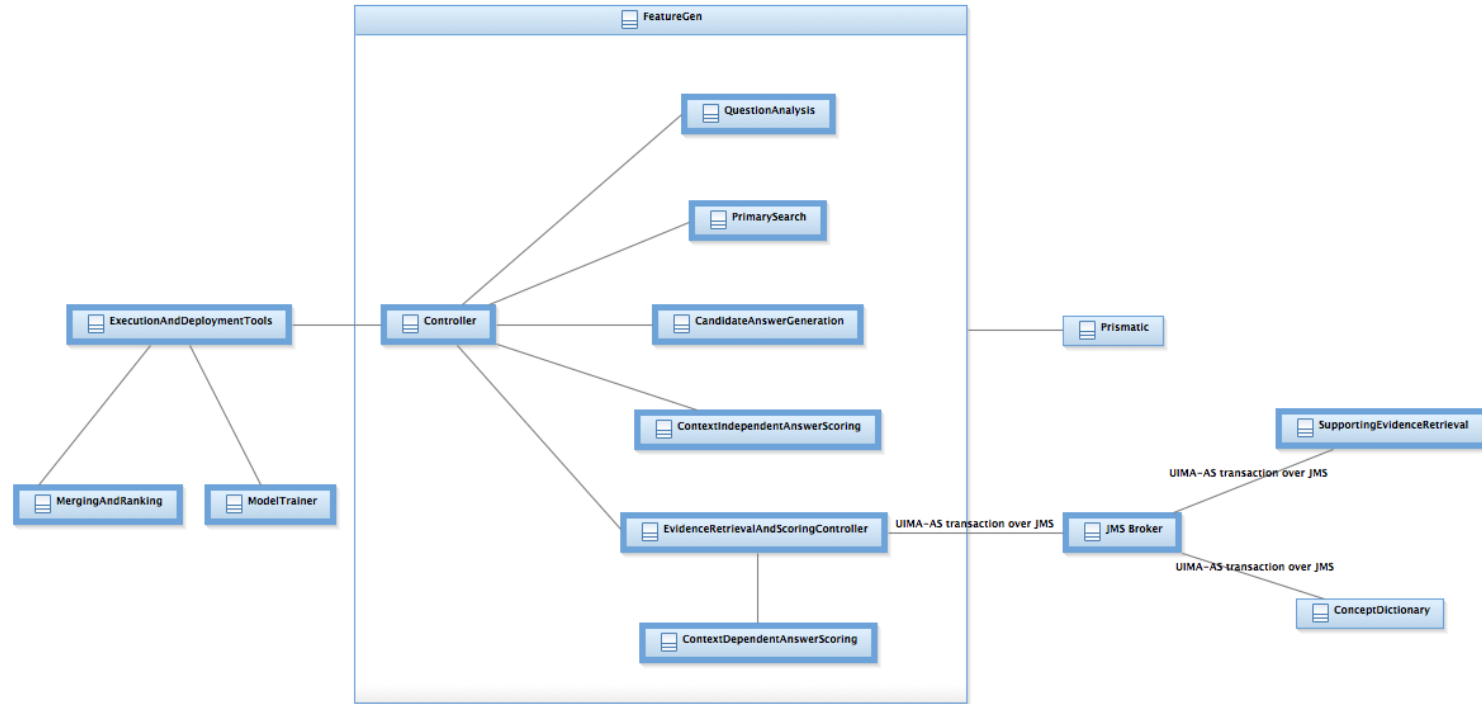
Process View: Key Mechanisms

- UIMA-AS.
 - A set of general capabilities for achieving scale out, built upon UIMA.
- UIMA CAS Multiplier and CAS pools.
 - Expand/consolidate CAS envelopes efficiently across multiple configurable flows.
- Three communication protocols:
 - UIMA-AS transactions across JMS.
 - DeepQA protocol for accessing large in-memory datasets.
 - Indri distributed search protocol.

Process View: Low-latency Production



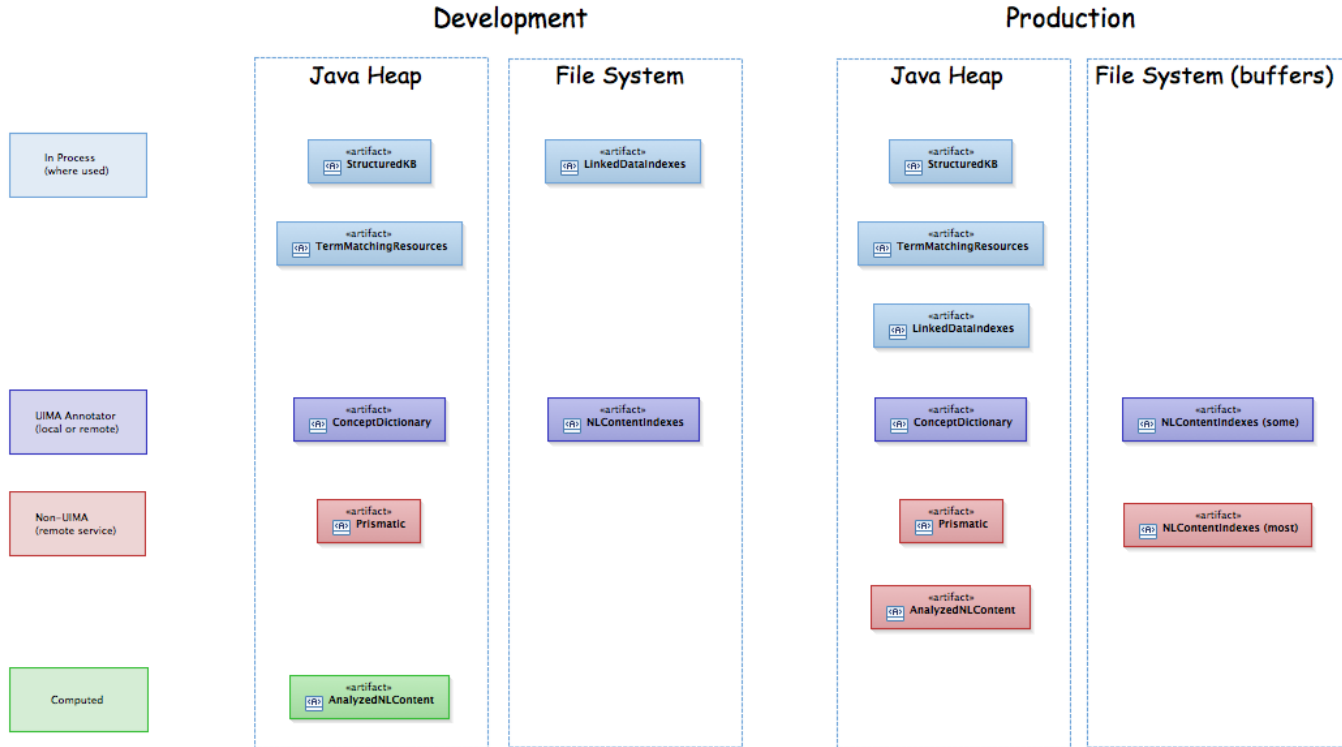
Process View: High-throughput Development



Component View: Code Layering

Execution Scripts/Deployment Tools (bluej.system.ant/bluej.system.as)				
DeepQASystem (bluej.system)				ModelTrainer (system.trainer)
Primary Components (bluej.question_analysis/bluej.subsystem.question_analysis/bluej.retrieval/ bluej.candidate_answer_generation/bluej.answer_scoring)			Development Tools (parts of bluej.tools)	
InternalComponents/TermSource/TermMatcher (sai.matcher/bluej.constrainer/sai.disambiguation/etc)				
Source Ingestion (bluej.tools.corpus_processing/ bluej.corpus_processing)	Evidence Interfaces/Resource Interfaces (bluej.ksp/bluej.rdf/bluej.prismatic/bluej.content_server/bluej.spatial)			Controller (bluej.core)
DeepQA Utilities (bluej.util)	Type System (bluej.model)	NLPStack (sai.text_analysis/bluej.text_analysis/watson.xsg/ bluej.question_analysis.relations)	TermGraph (sai.logical_form.kr)	Base Tools (bluej.tools.corpus_processing/bluej.corpus_processing)
General Utilities (sai.utilities/3rd party JARS)			UIMA (uima)	

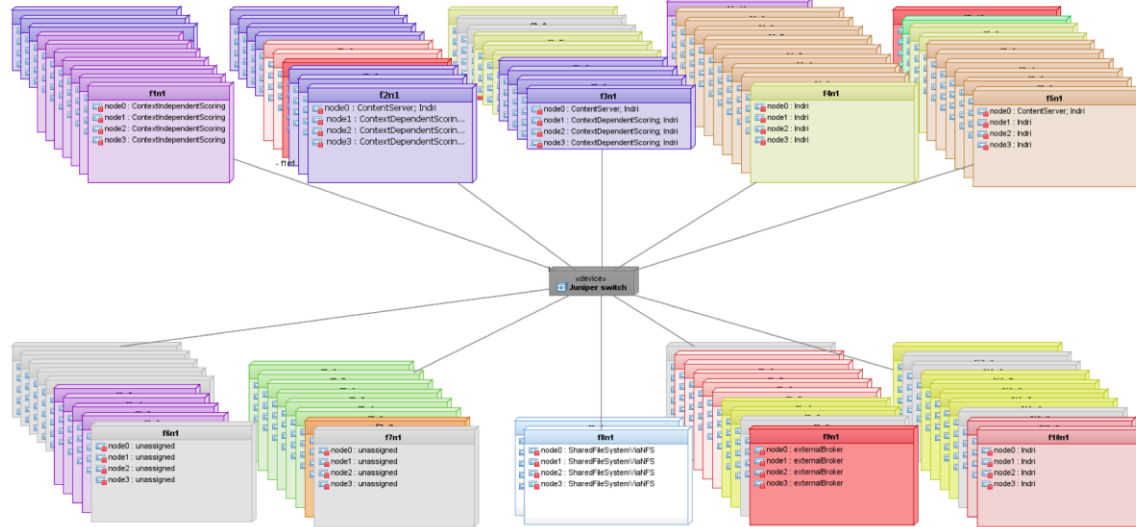
Component View: Data Storage



Deployment View: Watson

- 90 IBM Power 750 servers.
 - 4 Power7 processors/server.
 - 8 cores/processor.
 - 10 TB memory/server.
 - Linux.
- SONAS storage @ 20 TB.
- Juniper switch @ 10 Gbps.
- 2 20-air conditioning units.

Deployment View: Watson



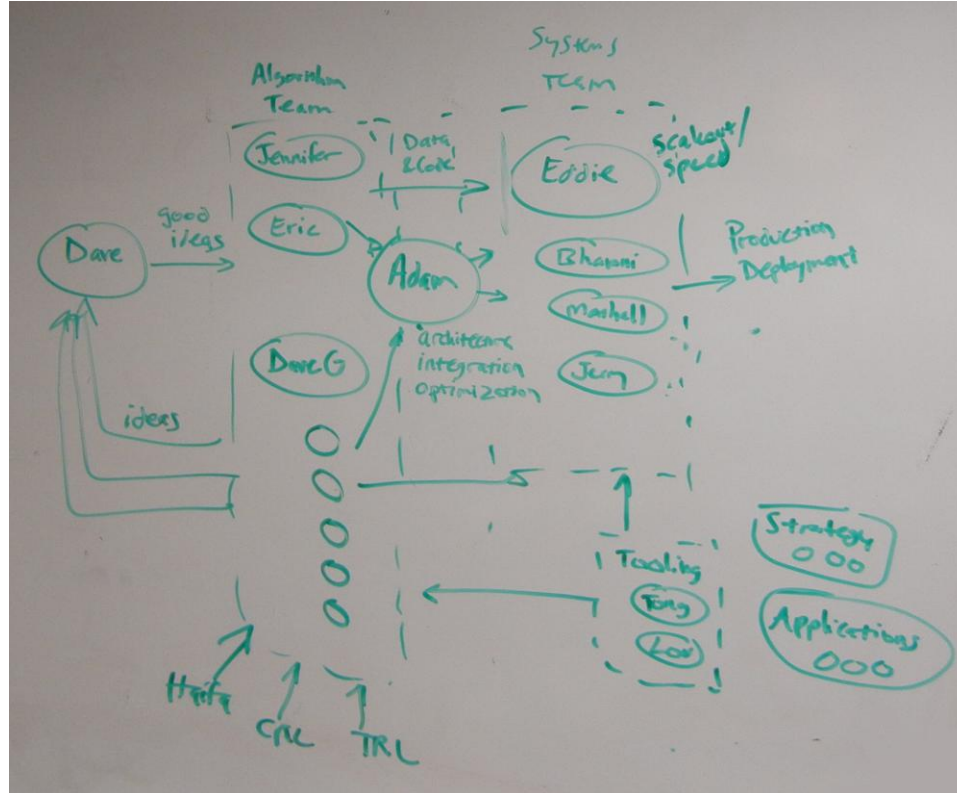
CandidateGeneration	ContentServer*	ContextDependentScoring*	ContextIndependentScoring	Supporting Evidence Retrieval / Evidence Retrieval And Scoring Controller	FileSystem
Indici (NICContentIndexes)	Lucene (PrimarySearch + NICContentIndexes)	Prismatic	TopLevelAggregate	Other	Unassigned

* denotes Indici processes also deployed on same node

Deployment View: Watson



Organization: In Their Own Words

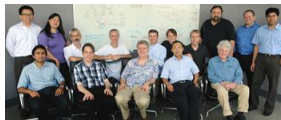


Organization

- David Ferrucci



- Algorithms Team



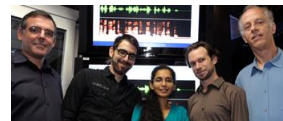
- Strategy Team



- Systems Team



- Speech Team



- Annotations Team

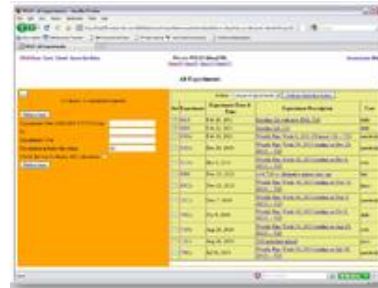


- Project Management



Tools

- Eclipse.
- Subversion -> RTC.
- Watson Error Analysis Tool (WEAT).
- Feature Analysis Tool (FAT).
- BlueJ Automatic Distributed Execution Environment tools (BAIDE)
- Data repository tools.



Question#508175
STRATEGY CHARACTERISTICS: The columns include Charity, Bargain, Used, Eye, Weekly, & Interest. Select a category from the dropdown menu.

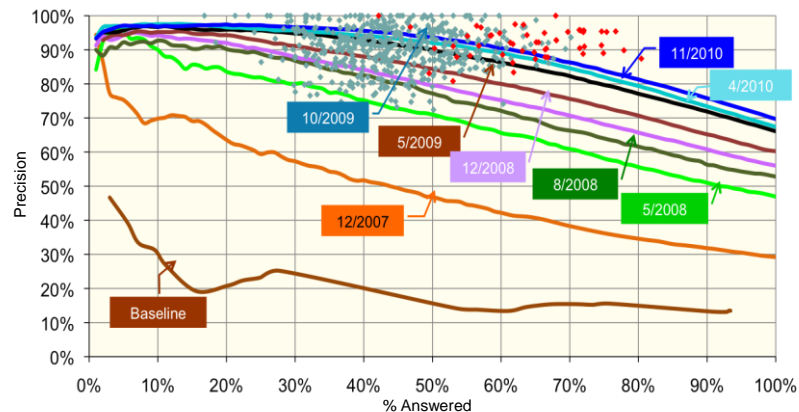
< Compare 'Used' (Selected) and 'Weekly' (Filter)

Feature Description	SP Value	SP Value	Used	Weekly	Interest
Equipment Used			2,444	2,444	
Technical Support			2,444	2,444	2,444
Real Estate			-1,111	-1,111	-1,111
Original Feature Set					
ACCOUNT_SUPPORT			-6,111	-6,111	-6,111
GENERAL_SPECIFIC			8,111	8,111	8,111
PERSONAL_SUPPORT			6,111	6,111	6,111
POPULARITY			-6,111	-6,111	-6,111
SECURITY_SECURITY			-6,111	-6,111	-6,111
TYPE_INDEX			6,111	6,111	6,111
REMOVAL_ADDITION			-6,111	-6,111	-6,111



Process

- Agile development.
- War room setting with continuous collaboration.
- Weekly integration.
- Results driven with end to end regression testing.
- ~ 6,000 experiments
- 10 gigabits of test data/week.



Watson's Future: New Development

- Hygienic.
 - Refactoring; elevating certain features to first-class architecture elements; performance and platform improvements; common configuration management.
- Research.
 - Greater introspection; dialoging; video & speech input; real time data input.
- Business.
 - Product line architecture.



Software. Everywhere.



www.ibm.com/software/rational

© Copyright IBM Corporation 2011. All rights reserved. The information contained in these materials is provided for informational purposes only, and is provided AS IS without warranty of any kind, express or implied. IBM shall not be responsible for any damages arising out of the use of, or otherwise related to, these materials. Nothing contained in these materials is intended to, nor shall have the effect of, creating any warranties or representations from IBM or its suppliers or licensors, or altering the terms and conditions of the applicable license agreement governing the use of IBM software. References in these materials to IBM products, programs, or services do not imply that they will be available in all countries in which IBM operates. Product release dates and/or capabilities referenced in these materials may change at any time at IBM's sole discretion based on market opportunities or other factors, and are not intended to be a commitment to future product or feature availability in any way. IBM, the IBM logo, Rational, the Rational logo, Telelogic, the Telelogic logo, and other IBM products and services are trademarks of the International Business Machines Corporation, in the United States, other countries or both. Other company, product, or service names may be trademarks or service marks of others.