

Analogical Reasoning with and about Databases and Knowledge Bases

John F. Sowa

VivoMind Intelligence, Inc.

Fundamental Problem

- * Deduction is the method of reasoning at the foundation of current database and knowledge-based systems.
- * Deduction is precise, predictable, and brittle.
- * If everything is perfect, deduction is perfect.
- * If anything in the system is imperfect, deduction can magnify and propagate the imperfection to the point of a total collapse.
- * When multiple systems interoperate, the likelihood and the danger of imperfection escalates.

Prospects for a Universal Ontology

Attempts to create a universal classification of all concepts:

- * 4th century BC: Aristotle's categories and syllogisms.**
- * 17th century: Universal language schemes by Descartes, Mersenne, Pascal, Leibniz, Newton, Wilkins, and others.**
- * 18th century: More schemes, the grand academy of Lagado, Kant's categories.**
- * 19th century: Roget's Thesaurus, Oxford English Dictionary.**
- * Early 20th century: Many terminologies in many different fields.**
- * 1960s: Computerized versions of the terminologies.**
- * 1970s: ANSI/SPARC Conceptual Schema.**
- * 1980s: Cyc, WordNet.**
- * 1990s: SRKB, ISO Conceptual Schema, Semantic Web, many workshops.**
- * 2000s: Many proposals, no consensus.**

Informal terminologies and dictionaries have been extremely successful.

Formal systems are still research projects.

How do People Avoid the Problems?

They don't always avoid them, but there are some safeguards:

- * People seldom carry out long chains of deductions.**
- * When the conclusion seems odd, they check the facts, ask for advice, and perform a “sanity check.”**
- * They communicate by message passing.**
- * They don't expect every message to be completely understood.**
- * They ask questions, give explanations, negotiate, and compromise.**

Some of the greatest disasters in history were caused by people who demanded and expected computer-like perfection.

Challenges for Interoperable Systems

Greatest mistake: Demand and expect perfection.

- * No two computer programs, even versions of “the same” program, are based on identical axioms (AKA “specifications”).**
- * Every major computer system in use today has specifications that are incomplete, incorrect, unavailable, or nonexistent.**
- * A large system requires an investment of multimillions of dollars, and its lifetime is usually 20 years longer than originally planned (except, of course, when it is so bad that it's never deployed).**
- * The original specifications for any large system become hopelessly obsolete, they are rarely updated, and they are not specified in any formal language.**
- * Supporting programs always adapt to the “500-pound gorilla”.**
- * A graceful migration path from the old to the new is essential.**

Suggested Solution

Treat computers as if they're no better than humans:

- * Don't assume that all systems have the same or similar ontologies or even any explicit ontology of any kind.**
- * Don't require agreement on anything outside the specific task.**
- * Insist on sanity checks and the option to back out gracefully.**
- * Communicate by message passing and dialog, not direct orders.**
- * Develop methods of negotiation, compromise, and voting.**

This approach requires reasoning methods that are more flexible and robust than strict deduction.

Methods of Reasoning

Three methods of formal logic:

1. Deduction. Apply a general principle to infer some fact.

Given: Every bird flies. Tweety is a bird.

Infer: Tweety flies.

2. Induction. Assume a general principle that subsumes many facts.

Given: Tweety, Polly, and Hooty are birds. Fred is bat.

Tweety, Polly, and Hooty fly. Fred flies.

Assume: Every bird flies.

3. Abduction. Guess a new hypothesis that explains some fact.

Given: Every bird flies. Tweety flies.

Guess: Tweety is a bird.

According to Peirce (1902), “Besides these three types of reasoning there is a fourth, analogy, which combines the characters of the three, yet cannot be adequately represented as composite.”

Four Views of Analogical Reasoning

1. By logicians:

Deduction is reasoning from “first principles.”

Analogy is an unsound, but interesting heuristic.

2. By psychologists:

Analogy is a fundamental cognitive mechanism.

Language and reasoning depend heavily on analogy.

3. Theoretical:

Analogy is a very general pattern-matching process.

Deduction, induction, and abduction depend on disciplined uses of analogy.

4. Computational:

A powerful and flexible technique for reasoning, learning, and language processing.

But practicality depends on finding analogies efficiently.

Example of Analogy: How is a cat like a car?

Analogy of Cat to Car	
Cat	Car
head	hood
eye	headlight
cornea	glass plate
mouth	fuel cap
stomach	fuel tank
bowel	combustion chamber
anus	exhaust pipe
skeleton	chassis
heart	engine
paw	wheel
fur	paint

Data from WordNet and other sources are translated to conceptual graphs.

The VivoMind Analogy Engine (VAE) starts at the nodes for Cat and Car and tries to find the longest matching paths through the graphs.

All analogies found are ranked by a semantic-distance measure. The above analogy received the highest score for that pair of words

Operations Performed During Pattern Matching

Following paths from each starting node:

[Cat]→(HasPart)→[Head]→(HasPart)→[Eye]→(HasPart)→[Cornea]

[Car]→(HasPart)→[Hood]→(HasPart)→[Headlight]→(HasPart)→[GlassPlate]

Matching concept nodes with similar types, properties, and relations:

- Head and hood are in the front.
- Eyes and headlights are related to light.
- Cornea and glass plate are transparent.
- Paws and wheels support the body, and there are four of each.

Approximate matching (ignoring the nodes in red):

[Cat]→(HasPart)→[Mouth]→(Flow)→[Esophagus]→(Flow)→[Stomach]→
→(Flow)→[Bowel]→(Flow)→[Anus]

[Car]→(HasPart)→[FuelCap]→(Flow)→[FuelTank]→
→(Flow)→[CombustionChamber]→(Flow)→[Muffler]→(Flow)→[ExhaustPipe]

Note: The source data did not contain information about the pipe from the fuel cap to the fuel tank. That missing part would have matched the esophagus of the cat.

Three Methods Used by VAE

- 1. Matching concept types, in order of increasing semantic distance:**
 - Identical types.
 - Subtype - supertype.
 - Siblings of same supertype.
 - More distant cousins.
- 2. Matching subgraphs:**
 - Match isomorphic subgraphs independent of type labels.
 - Merge adjacent nodes to make them isomorphic.
- 3. Finding metalevel mappings that can relate subgraphs,
— even though they are not isomorphic.**

Methods #1 and #2 are used for the Cat-Car example.

Method #3 (combined with #1 and #2) is used for aligning ontologies.

Need to Align Ontologies

Any topic can be described at many levels of detail with different choices of labels for the concept and relation types.

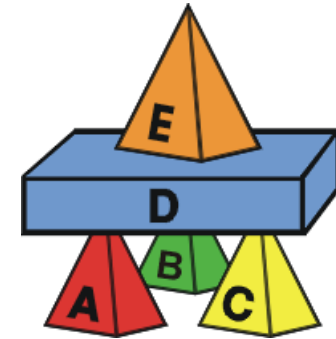
There is no standard upper-level ontology, but there are many important ontologies that must be accommodated:

- Some widely used ontologies, such as Cyc, OpenCyc, SUMO, Dolce, BFO, etc.**
- Many ontologies required by governments and large organizations.**
- Example: The Amazon.com ontology.**
- An enormous number of legacy systems with no explicit ontologies.**

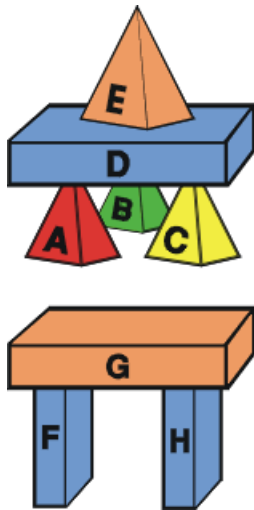
Critical requirement: Enable heterogeneous systems to interoperate despite differences in their underlying ontologies, whether implicit or explicit.

Example of Different Ontologies

- The structure on the right may be described in different ways:
- English sentence: “A red pyramid A, a green pyramid B, and a yellow pyramid C support a blue block D, which supports an orange pyramid E.”
- A relational database would use tables.
- But there are many different options for organizing the tables and choosing labels for the tables and columns.
- These choices lead to different ontologies, which may be structurally very different from one another.



Representation in a Relational DB



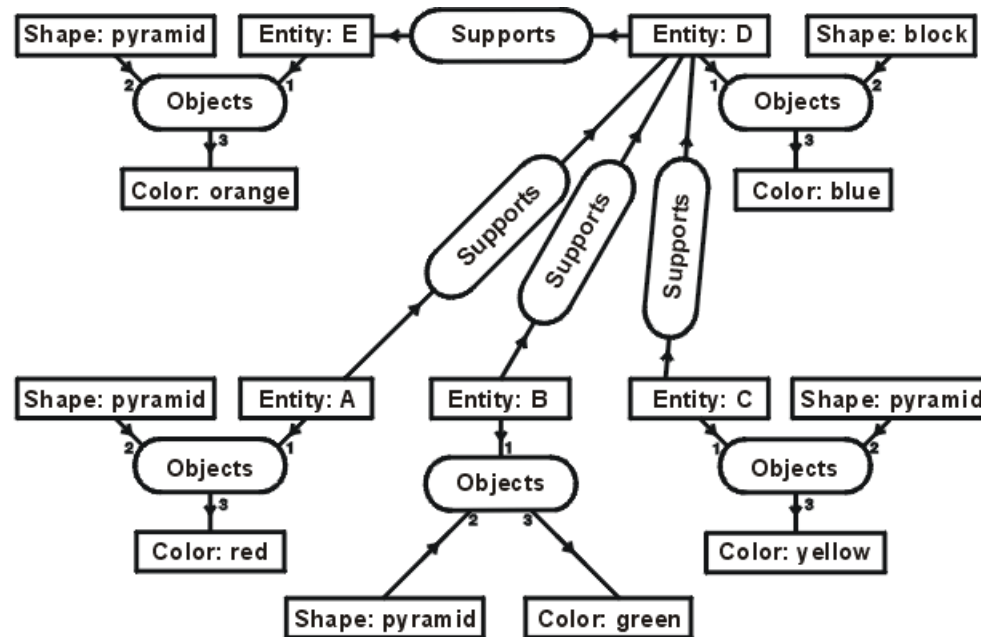
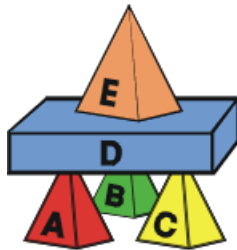
Objects

ID	Shape	Color
A	pyramid	red
B	pyramid	green
C	pyramid	yellow
D	block	blue
E	pyramid	orange
F	block	blue
G	block	orange
H	block	blue

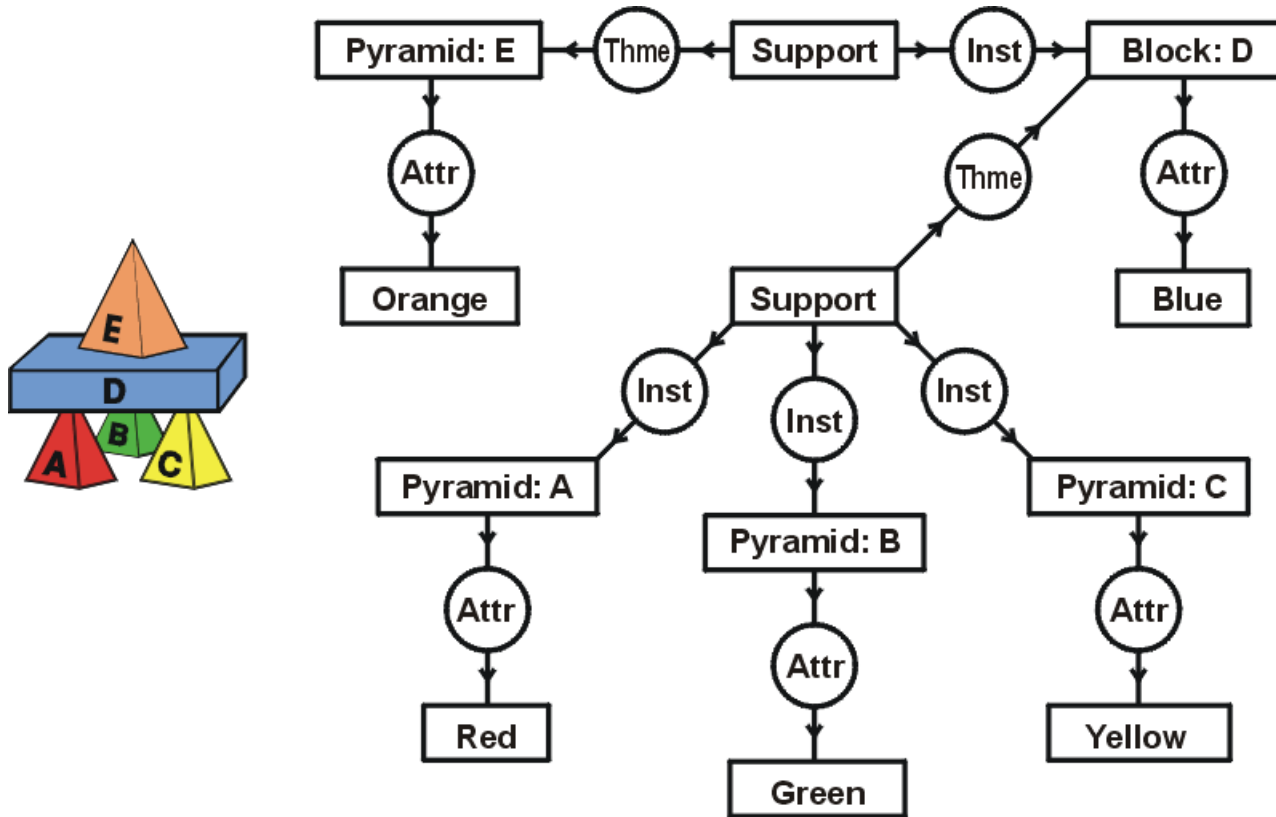
Supports

Supporter	Supportee
A	D
B	D
C	D
D	E
F	G
H	G

Conceptual Graph from Relational DB



Conceptual Graph from English Sentence

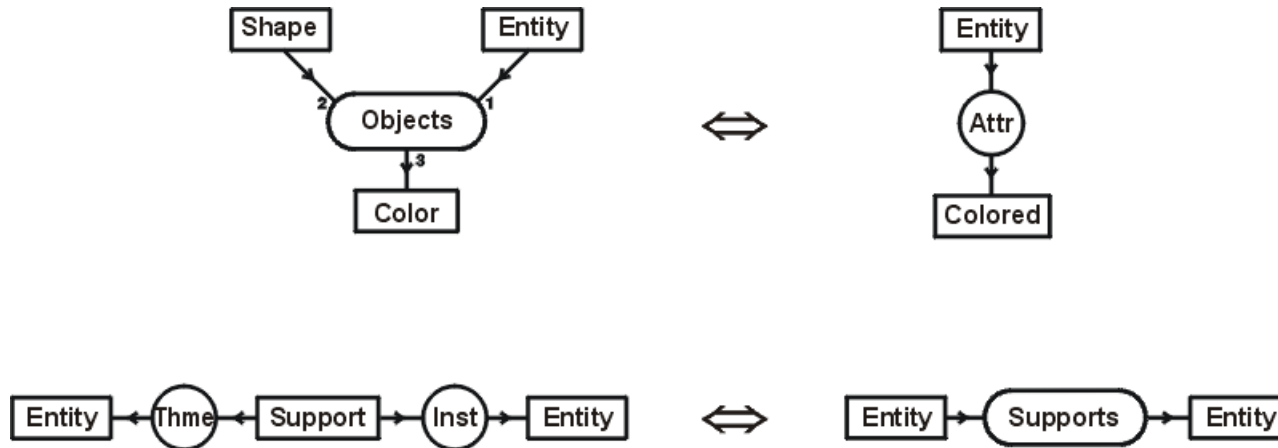


“A red pyramid A, a green pyramid B, and a yellow pyramid C support a blue block D, which supports an orange pyramid E.”

The Two CGs Look Very Different

- * **CG from RDB has 15 concept nodes and 8 relation nodes.**
- * **CG from English has 12 concept nodes and 11 relation nodes.**
- * **All the type labels of concept and relation nodes are different.**
- * **But there are some structural similarities.**
- * **VAE uses method #3 to find them.**

Transformations Found by VAE



Top transformation was applied to 5 subgraphs.

Bottom one was applied to 4 subgraphs.

One application of any given transformation could be due to chance, but 4 or 5 applications are strong evidence for its significance.

The fact that these two transformations completely map each graph to the other is convincing.

Computational Complexity

Research by Falkenhainer, Forbus, & Gentner:

Problem: In a knowledge base of N graphs, the time to find the best analogy to a graph G takes time proportional to N^3 .

If N is 10, N^3 is a thousand.

But if N is a billion, N^3 is 1,000,000,000,000,000,000,000,000,000.

MAC/FAC Solution:

Use a search method to narrow down the likely candidates from some large number N to a much, much, much smaller number n .

With a good search method:

- * The execution time for the search should be proportional to $\log(N)$.
- * It should find some small number of graphs, n .
- * And the most relevant graphs should be among those n .

Question: What kinds of search algorithms have these properties?

Knowledge Signature

An encoding of a conceptual graph with the following properties:

1. Encodes the ontology of the concept and relation types.
2. Encodes the topology (structural connectivity) of the graph.
3. Independent of later additions and extensions to the ontology.
4. Determines a measure of the semantic distance between graphs.
5. Suitable for indexing the graphs in time proportional to
 - * $N \times \log(N)$ for building the index.
 - * $\log(N)$ for using the index to find all graphs within a small semantic distance ϵ from any graph G .

Such encodings are possible, but there are many variations, which are being explored and evaluated.

Criticisms of Logical Deduction

Deduction in mathematics can be precise, but deduction about any empirical subject must depend on prior induction, which is almost always incomplete.

Criticism by the physician, Sextus Empiricus (2nd century AD):

Every human is an animal.

Socrates is human.

Therefore, Socrates is an animal.

If the major premise was derived by checking every human, then Socrates was considered, and the argument is circular.

Otherwise, the induction was incomplete, and the conclusion is uncertain.

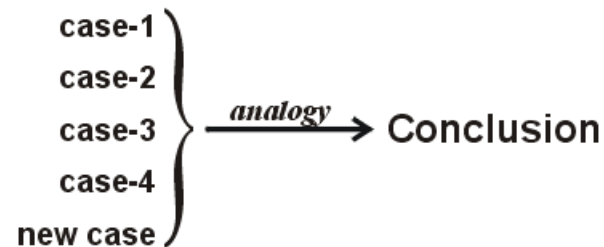
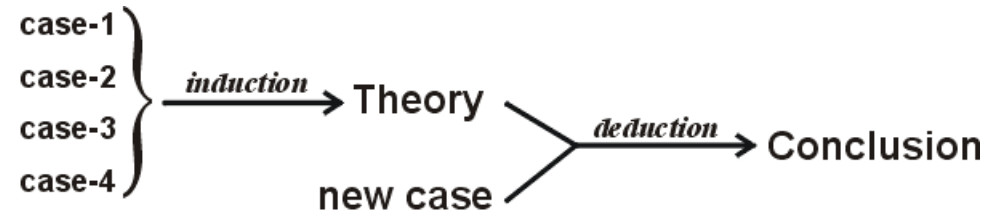
Criticism by the legal scholar, Ibn Taymiyya (14th century AD):

Every empirical theory is derived by induction from cases.

Any deduction from such a theory can be obtained by applying analogical reasoning to exactly the same cases.

Case-Based Reasoning

Ibn Taymiyya's comparison of logical and analogical reasoning:



If the same cases are used, analogy can give the same answer in one step that requires two steps by induction followed by deduction.

What Ibn Taymiyya did not recognize: If the same theory could be reused for many different applications, deduction would be more efficient.

Logic as a Disciplined Use of Analogies

The structure-mapping operations of analogy are used in every kind of logical reasoning:

- * Deduction. Every step requires a unification, which is special case of the structure mappings used in analogies.**
- * Induction. Analogies are used to find common generalizations of multiple instances.**
- * Abduction. The operation of guessing or forming an initial hypothesis, called abduction, requires analogies to find likely causes or explanations.**

In both human reasoning and computer implementations, the same underlying operations can support both logical and analogical reasoning.

An Application of Case-Based Reasoning

A textbook publisher needs to evaluate student answers to math questions.

- * Free-form answers in English sentences.
- * Much harder to evaluate than multiple choice.

Typical question:

The following numbers are 1 more than a square: 10, 37, 65, 82.

*If you are given an integer N that is less than 200,
how would you determine whether N is 1 more than a square?*

Explain your method in three or four sentences.

How could a computer system evaluate student answers?

Determine whether they are correct, incorrect, or partially correct?

And make helpful suggestions about the incorrect answers?

Publisher's Current Procedure

To evaluate new exam questions, the publisher normally gives the exam to a large number of students.

For each question, they get about 50 different answers:

- * Some are completely correct
— but stated in different ways.**
- * Some are partially correct
— and the teacher says what is missing.**
- * Others are wrong
— in many different ways.**

Result: 50 pairs of student answer and teacher response.

Each pair of (answer,response) is a case for case-based reasoning.

Intellitex Parser

VivoMind has developed a parser named Intellitex:

- * Uses a rather simple grammar.**
- * Depends on analogies for interpreting sentences.**
- * Generates conceptual graphs as output.**
- * Robust: always generates some CG as its best guess.**

These properties are important for handling typical student answers, which frequently have poor grammar and incomplete sentences.

Minor errors are not necessarily bad — provided that Intellitex makes the same errors consistently in all cases.

Using VAE to Evaluate Student Answers

Use VAE to compare each new answer to the 50 cases:

1. For all 50 cases, translate student answer to conceptual graphs.
2. Translate each new answer to a new CG.
3. Compare the new CG to the 50 CGs for previous answers.
4. Use the measure of semantic distance to determine the best match.
5. If there is a good match, print out the corresponding response.
6. Otherwise, send the new answer to a teacher to evaluate.

Result:

VAE found a good match for most of the student answers.

For each good match, the previous teacher's response was appropriate.

When VAE failed to find a good match, the new case could be added to the list of cases in order to improve its coverage.

There was no need for the teachers to write rules or programs.

Legacy Re-engineering

Version of Intellitex applied to three languages — English, COBOL, and JCL:

- * 1.5 million lines of COBOL.**
- * Several hundred JCL scripts.**
- * 100 megabytes of English documentation — reports, manuals, e-mails, Lotus Notes, HTML, and transcriptions of oral communications.**

Goal:

- * Analyze the COBOL and JCL to determine:**
 - Data dictionary, data flow diagrams, process architecture diagrams, system context diagrams.**
- * Analyze the English documentation to determine:**
 - Discrepancies between the documentation and the implementation.**
 - English glossary of all terms and their changes over the years.**
 - English captions on the diagrams derived from COBOL and JCL.**

Sample Documentation

Mixture of English words, computer jargon, and names of programs, files, and COBOL variables:

The input file that is used to create this piece of the Billing Interface for the General Ledger is an extract from the 61 byte file that is created by the COBOL program BILLCRUA in the Billing History production run. This file is used instead of the history file for time efficiency. This file contains the billing transaction codes (types of records) that are to be interfaced to General Ledger for the given month. For this process the following transaction codes are used: 32 - loss on unbilled, 72 - gain on uncollected, and 85 - loss on uncollected. Any of these records that are actually taxes are bypassed. Only client types 01 - Mar, 05 - Internal Non/Billable, 06 - Internal Billable, and 08 - BAS are selected. This is determined by a GETBDATA call to the client file. The unit that the gain or loss is assigned to is supplied at the time of its creation in EBT.

Note nonstandard syntax: 32 - loss on unbilled, 06 - Internal Billable.

An Important Simplification

Extremely difficult problem:

- * Trying to understand English specifications.**
- * Mapping the results to an executable COBOL program.**

Much, much easier problem:

- * Mapping COBOL to conceptual graphs.**
- * Using CGs derived from COBOL to analyze English documentation.**
- * Using VAE to find similarities and differences in the CGs.**

Results

Job finished in 8 weeks by two programmers, Arun Majumdar and André LeClerc.

*** Four weeks for customization:**

Design and logistics.

Additional programming for I/O formats.

*** Three weeks to run Intellitex + VAE + extensions:**

24 hours a day on a 750 MHz Pentium III.

VAE handled matches with strong evidence (close semantic distance).

Matches with weak evidence were confirmed or corrected by Majumdar and LeClerc.

*** One week to produce a CD-ROM with integrated views of the results:**

Glossary, data dictionary, data flow diagrams, process architecture, system context diagrams.

Contradiction Found by VAE

From analyzing English documentation:

- * Every employee is a human being.**
- * No human being is a computer.**

From analyzing COBOL programs:

- * Some employees are computers.**

What is the reason for this contradiction?

Quick Patch in 1979

A COBOL programmer made a quick patch:

- * Two computers were used to assist human consultants.**
- * But there was no provision to bill for computer time.**
- * Therefore, the programmer named the computers Bob and Sally, and assigned them employee ids.**

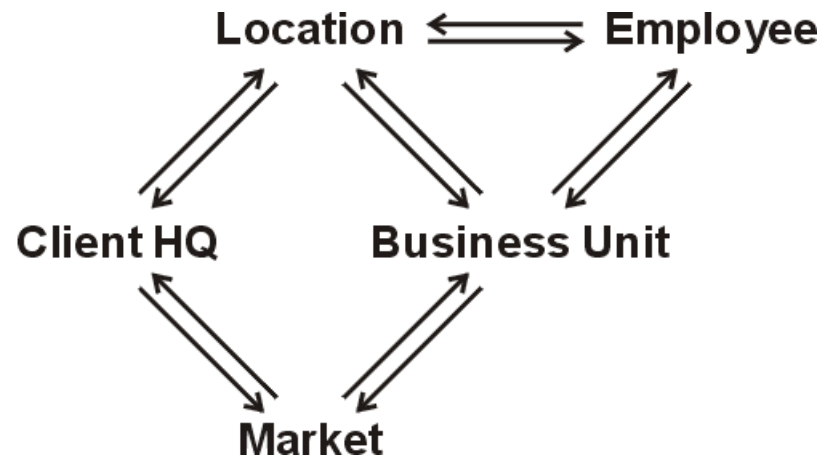
For more than 20 years:

- * Bob and Sally were issued payroll checks.**
- * But they never cashed them.**

VAE discovered the two computer “employees”.

Mismatch Found by VAE

A diagram of relationships among data types in the database:



Question: Which location determines the market?

According to documentation: Business unit.

According to COBOL programs: Client HQ.

Management had been making decisions based on incorrect assumptions.

Conclusions

Deductive methods are good when there are widely applicable theories, as in physics, engineering, and established accounting procedures.

When there are no reliable theories, analogical reasoning is necessary.

Even when good theories are available, analogical reasoning can be a valuable supplement for handling exceptions.

Analogical reasoning can also be used at the metalevel to find mappings between different theories and ontologies.

The semantic distance measures used in analogy can also be used to derive an ontology from either structured or unstructured sources.

For further reading and references, see

<http://www.jfsowa.com/pubs/analog.htm>