

Report 83-36
Stanford -- KSL

Scientific DataLink

What's New? A Semantic Definition of
Novelty.
Russell Greiner, Michael R. Genesereth,
Jun 1983

card 1 of 1

Stanford Heuristic Programming Project
Report No. HPP-83-36

June 1983

What's New? A Semantic Definition of Novelty

Russell Greiner and Michael R. Genesereth

Heuristic Programming Project
Department of Computer Science
Stanford University

This is a slightly expanded version of the paper which appeared in the

Proceedings of IJCAI-83.

1 Introduction

A central process in any learning experience is the incorporation of a *new* fact into an existing theory. Often the goal of that process is more specific, to learn some new fact about some concept. But what does it mean to claim that a sentence is *new*, and even more interesting, what qualifies as a novel fact about some concept? Despite the vast interest in learning and the abundance of related papers (cf. [Dietterich 81a], [Buchanan 78], [Michalski 83], [Dietterich 81b], [Dietterich 82]), no one has rigorously defined what it means to be "new", either in general or with respect to a single concept.

This paper attempts to fill that gap. Our goal is to obtain a *semantic* rather than a *syntactic* understanding of novelty. This preference stems from our belief that a semantic account (one based on the possible interpretations of the theory) provides important insight into the phenomenon of novelty. It also means we may be able to generalize these results to other logics and languages.

The results of this research are relevant (and useful) to many different fields. The primary importance of this work is in providing a first stab at describing the different senses of novelty. In addition to the applications a complete and adequate definition of newness would have as an analytic tool, there are possible applications in knowledge acquisition, representation, and discourse analysis. Many of these stem from the intimate connection between novelty and the intuitive notion of "aboutness". (Section 4 elaborates on each of these.)

This report discusses two kinds of novelty. Section 2 describes newness of a sentence with respect to a theory. Section 3 uses this result to address the more difficult task of determining when a sentence conveys something new about a particular concept. While the first kind of novelty is fairly easy to capture, the second requires a consideration of the interconnections among facts within a theory. Section 4 justifies why this undertaking is relevant and describes how these results may eventually be used. The concluding Section 5 lists several outstanding research issues.

Each of the appendices discusses some relevant digression. In particular, Appendix A motivates and illustrates our intended meaning of this second kind of novelty, by presenting a typical example.

2 New with respect to a Theory

This section addresses the issue of what it means for a sentence σ to be new with respect to a

theory¹ Th ; this is the relation $N(Th, \sigma)$. Intuitively, we want σ to be new if it (somehow) further specifies something about the world. Alternatively we can think of a new sentence as providing some additional constraints, which remove some possible worlds ([Moore 80]) from consideration.

We first consider a semantic definition of newness: σ is new with respect to a theory if it eliminates some possible interpretation of that theory.² That is, given any theory Th , in the language L , there is a set of models $I^{Th} = \{I_j\}$, in which each I_j maps the symbols of L into objects or sets of tuples of objects in the "real world" in the standard way. That is, each constant is mapped onto an object, and each n -ary relation onto its characteristic set -- those n -tuples that satisfy it. Notice that this means that the universe is fixed beforehand and that these ranges can overlap. (Appendix A elaborates this point.)

Adding additional sentences to a theory can only restrict the set of possible interpretations: $Th \subseteq Th'$ means that $I^{Th'} \subseteq I^{Th}$. This leads to the proposal that

$$\text{Defn 1: } N_{\text{Sem}}(Th, \sigma) \Leftrightarrow I^{Th+\sigma} \subsetneq I^{Th}.$$

This same definition can be expressed syntactically (in terms of logical deducibility rather than semantic validity) using

$$\text{Defn 2: } N_{\text{Syn}}(Th, \sigma) \Leftrightarrow (\sigma \notin Th) \wedge (\neg\sigma \notin Th).$$

As these are equivalent (whenever the language of the theory remains fixed; see [Enderton 72].) we will simply use N . (Appendix B sketches a proof for this claim.)

3 New with respect to a Concept (and a Theory)

In many situations it is not enough to realize that an assertion is new; rather, one often wants to claim that it is a new fact *about some concept*. With this in mind, we define the ternary relation $\text{New}(Th, s, \sigma)$ to mean that the assertion σ expresses a new fact about the concept s with respect to the theory Th . Notice we will be dealing with a syntactic symbol, as opposed to its referent. The "learning step" involves adding this sentence σ (along with all of its deductive consequences) to the theory.

What should go into a definition of $\text{New}(Th, s, \sigma)$? Clearly, a necessary condition is that σ be a new fact with respect to the entire theory; that is, $N(Th, \sigma)$. But beyond that, we want to capture the

¹We are taking a slightly unorthodox *syntactic* view of theory: viz., a theory is a consistent and deductively closed set of axioms. We will also assume that the deductive system is complete.

²We are only concerned with "true interpretations", which map symbols into referents within a *legal* model. We chose "interpretation" rather than "model" to emphasize that we are dealing with a mapping rather than its range.

sense in which σ further specifies the concept s , or enables the derivation of additional relevant conclusions about s .

This section will present a definition of **New** by proposing a series of "increasingly more nearly correct" descriptions. For simplicity, the examples are taken from propositional logic. (We will later discuss the ways this is, and is not, a limitation.) Each proposal will begin with the intuition involved, followed by a formal description of this conjecture, and finally (in all but the final case) an example of the failure of this conjecture and an analysis of this shortcoming.

Conjecture 1: Syntactic Method. Most statements which relay information about some concept will contain the symbol that refers to that concept. This leads to the proposed syntactic solution: The sentence σ conveys new information about the symbol A if the token "A" is lexically included in the string of tokens which form σ , denoted with the assertion $\text{LexInclude}("A", \sigma)$ — e.g., $\text{LexInclude}("A", "A \wedge B")$. For the reasons mentioned above, we will further insist that $N(\text{Th}, \sigma)$. Formally,

Defn 3: $\text{New}_{\text{Syn}}(\text{Th}, s, \sigma) \Leftrightarrow N(\text{Th}, \sigma) \wedge \text{LexInclude}(s, \sigma)$.

Unfortunately, this syntactic condition is neither necessary nor sufficient. To see that it is not necessary, realize that we want $\text{New}(\{A \Leftrightarrow B\}, "A", "B")$ to be true, since asserting B in this situation means that A must now be true, which had not been the case before that assertion.

To show insufficiency is a little trickier. Should $\text{New}(\{A \vee B\}, "A", "A \Rightarrow B")$ be true? We argue the answer is no: In this context, asserting $A \Rightarrow B$ is the same as asserting B , which we know says nothing new about A . We clearly need a more powerful method for specifying novelty.

Conjecture 2: Fewer Interpretations.

(Subtitle: Model Theory, to the rescue!)

Using the notion of possible interpretations discussed in Section 2, we can define the "interpretation range" of a particular symbol. Let the term $I_j(s)$ designate the "real world" referent of the symbol s , given by the interpretation I_j — here, either \underline{I} or \underline{E} .³ We use this to define the interpretation range of the symbol s , $I^{\text{Th}}(s)$, by

Defn 4: $I^{\text{Th}}(s) = \{I_j(s) \mid I_j \in I^{\text{Th}}\}$.

As additional facts can only further restrict the range of possible interpretations for any symbol, we have $I^{\text{Th}'}(s) \subseteq I^{\text{Th}}(s)$ whenever $\text{Th} \subseteq \text{Th}'$. This leads to our second conjecture,

³The underbar notation denotes the referent of the corresponding linguistic symbol.

Defn 5: $\text{New}_{\text{FI}}(\text{Th}, s, \sigma) \Leftrightarrow I^{\text{Th} + \sigma}(s) \subsetneq I^{\text{Th}}(s)$.

This New_{FI} definition seems, at first, adequate. In addition to paralleling the **N** situation, it also resonates nicely with the ideas of Shannon's Information Theory, in which information is tied to the reduction of uncertainty in the distribution of possible values of a signal. (See [Gallager 78].) It also handles the two cases used above to discredit New_{Syn} .

Unfortunately, this New_{FI} requirement does not include all desired cases. There are some sentences that do convey new information in the informal sense outlined above but that do not satisfy this constraint: Start with an empty theory, $\text{ThA} \leftarrow \{\}$,⁴ in the language $L = \{A, B\}$. The four possible interpretations are shown in Figure 1.

	A	B
I_0	F	F
I_1	F	T
I_2	T	F
I_3	T	T

Figure 1: Interpretations of A and B in ThA.

By inspection, $I^{\text{ThA}}("A") = \{T, F\}$. Now, form $\text{ThB} \leftarrow \text{ThA} + "A \Leftrightarrow B"$. While this leaves only two of the four original interpretations, I_0 and I_3 , $I^{\text{ThB}}("A")$ remains $\{T, F\}$, indicating that $A \Leftrightarrow B$ said nothing New_{FI} about A.

Although New_{FI} rejects this $A \Leftrightarrow B$, we still believe it should be considered a new fact about A in this situation: If we later learn $\neg B$, we will be able to infer that $\neg A$, a conclusion that would not have followed without that sentence. That is, any sentence that makes A's range of interpretations dependent on some other symbol (in the sense that $A \Leftrightarrow B$ made A dependent on B) also "feels" new.

So there are at least two ways a statement can be new:

- It directly limits the interpretation range of A, or
- It establishes (or increases) a dependency of A on some other symbols (as that may, in turn, lead to a reduction of the above type).

While New_{FI} covers the first case exactly, it fails in the second.

⁴The notation $T \leftarrow \psi$ means the theory T is assigned the deductive closure of the set, ψ ; and $\text{Th} + \sigma$ refers to the deductive closure of $\text{Th} \cup \{\sigma\}$.

Conjecture 3: Partial Interpretations. To define dependency requires an understanding of what it means for one symbol to depend on some other symbols. The $A \Rightarrow B$ case above is clearly one instance of this. In addition to such singular dependencies (of A on one other symbol,) A may depend on a combination of symbols. (Consider the assertion $A \Rightarrow (B \Rightarrow C)$. Fixing any assignment to B alone, A will still be "arbitrary"; that is, it could be either \underline{I} or \underline{E} , depending on C . However, if both B and C are fixed, then A is fully determined.)

We saw that σ is a new fact about A if it increases A 's dependency on some n -tuple of symbols $\langle s_1, \dots, s_n \rangle$, that is, if asserting σ restricts the set of assignments available to A , given some assignment $\langle \underline{s}_1, \dots, \underline{s}_n \rangle$ to the symbols $\langle s_1, \dots, s_n \rangle$. For example, we noted that A was more dependent on B in $\text{Th}B$ than in $\text{Th}A$. We see this by considering the assignment of B to \underline{E} . In $\text{Th}A$, A 's value could be either \underline{E} or \underline{I} , independent of this assignment to B . However, A can no longer be assigned \underline{I} in $\text{Th}B$, given this assignment to B ; its value is now restricted to \underline{E} .

As this " $\langle s_1, \dots, s_n \rangle$ assignment to $\langle \underline{s}_1, \dots, \underline{s}_n \rangle$ " reflects an assignment of only a subset of the symbols, $\{s_i\} \subseteq L$, we will call it a *partial interpretation*. We can associate with each partial interpretation the equivalent class of full interpretations that agree on the assignment of a set of symbols. Formally, take any function that maps some of the symbols of the language into the universe, U — that is, any $\varphi: \xi \rightarrow U$, where $\xi \subseteq L$. We can use this to define the equivalence class $[[\varphi^{\text{Th}}]]$ as the set of interpretations that consistently extend φ — that is, it includes each interpretation that agrees with φ 's assignment of each symbol in φ 's domain and assigns every other element of L in some consistent manner.

Defn 6: $[[\varphi^{\text{Th}}]] = \{I \in I^{\text{Th}} \mid \forall x \in \text{Domain}[\varphi]. \varphi(x) = I(x)\}$.

The assignments of s that are consistent with the partial interpretation $[[\varphi^{\text{Th}}]]$ are just

Defn 7: $[[\varphi^{\text{Th}}]](s) = \{I(s) \mid I \in [[\varphi^{\text{Th}}]]\}$.

With this notation, we can state that A 's dependency on the symbols $\xi = \text{Domain}[\varphi]$ increases if the set of possible values of A consistent with the partial interpretation, $[[\varphi^{\text{Th}}]]("A")$, decreases but remains non-empty. (Seeing it vanish means that there are no values of A that are consistent with this assignment, which means that no models can be derived by extending this partial interpretation.)

To test if σ is new, therefore, consider all of the assignments available to A for each partial interpretation, before and after adding this purportedly new sentence σ to the theory. If the number of possible referents of A decreases for any partial interpretation (and remains non-empty), we will declare that σ is new.

Defn 8: $\text{New}_{PI}(\text{Th}, s, \sigma) \Leftrightarrow \exists \varphi. \llbracket \varphi^{\text{Th}+\sigma} \rrbracket(s) \subsetneq \llbracket \varphi^{\text{Th}} \rrbracket(s) \wedge \llbracket \varphi^{\text{Th}+\sigma} \rrbracket(s) \neq \{\}$

A few notes:

1. We will say that the particular function φ whose partial interpretation $\llbracket \varphi^{\text{Th}} \rrbracket$ decreased in the above equation is a "witness" to σ 's novelty (with respect to A).
2. This definition subsumes the $\text{New}_{FI}(\text{Th}, A, \sigma)$ condition. This follows from the fact that $\llbracket \{\}^{\text{Th}} \rrbracket(A)$ is equal to $I^{\text{Th}}[A]$. (Note this $\{\}$ mapping is the "null mapping", whose domain is empty.)
3. Realize that if $A \in \text{Domain}[\varphi]$, then $\llbracket \varphi^{\text{Th}} \rrbracket(A)$ would contain a single member. As there are no non-empty proper subsets of such singleton sets, it is sufficient to use $\text{Domain}[\varphi] \subseteq L - \{A\}$, rather than $\text{Domain}[\varphi] \subseteq L$.
4. Consider the set of "almost complete interpretations" $\llbracket \varphi^{\text{Th}} \rrbracket$, whose domain includes every symbol of L except A ; that is, $\text{Domain}[\varphi] = L - \{A\}$. While it may appear that these partial interpretations are sufficient — that one of these would witness any new σ — the counterexample below shows that is not the case.

A tableau helps to visualize this definition. The left tableau in Figure 2 corresponds to the theory $\text{ThC} \leftarrow \{A \Leftrightarrow B\}$ in the language $L = \{A, B, C\}$, and the one on the right to $\text{ThD} \leftarrow \text{ThC} + "A \Leftrightarrow C"$. Each row is indexed by a mapping φ and each column by an assignment to A . A tableau position is tagged with a "1" if this assignment of $\text{Domain}[\varphi] \cup \{A\}$ is consistent — that is, if there is any full interpretation associated with this position — and a "0" otherwise.

Finding a witness to a sentence's novelty reduces to finding a row, r , in which a "1" is flipped to "0" but which does not vanish — that is, r must retain a "1" in some position. The fifth and sixth rows below (labeled with the " $\langle C, F \rangle$ " and " $\langle C, I \rangle$ " mappings) each satisfy this property, showing that $\text{New}_{PI}(\text{ThC}, "A", "A \Leftrightarrow C")$. Notice that none of the top four rows, which correspond to those "almost complete interpretations" has that property, demonstrating the point of Item 4 above. (These rows do form an adequate spanning set, though, as all the other rows can be derived by ORing together appropriate sets of these.)

4 Uses

Even in its current unmechanized form, this definition can be used effectively as an analytic tool with which to understand many existing learning programs. Eventually, we hope to develop a "NewP" predicate or possibly a pair of operational (multivalued) functions: "New σ ", which generates New_{PI} sentences from a given theory and symbol, and "NewSYM", which returns the symbols for which a given sentence is New_{PI} . This section lists several ways this definition (and its operationalizations) can be used.

A			A	
E	I		E	I
1	0	→ { <B,E>, <C,E> }	1	0
1	0	→ { <B,E>, <C,I> }	0	0
0	1	→ { <B,I>, <C,E> }	0	0
0	1	→ { <B,I>, <C,I> }	0	1
1	1	→ { <C,E> }	1	0
1	1	→ { <C,I> }	0	1
1	0	→ { <B,E> }	1	0
0	1	→ { <B,I> }	0	1
1	1	→ { }	1	1

Figure 2: Partial interpretations for ThC and ThD.

- *Analytic Tool.* An adequate definition of newness would help us identify the sources (and recipients) of novelty within learning programs. For example, the teacher provides the ARCH program ([Winston 75]) with the new facts that enable it to learn. LEX's problem solver and critic are the sources of novelty for the rest of the system ([Mitchell 81]). AM ([Lenat 82]) has no clear source of novelty. This definition may also help us understand the distinction between compositional new terms — such as AM's definition of prime numbers — and other new terms, such as Bacon's use of intrinsic properties ([Langley 79]). Finally, it may lead to a definition of learning not based exclusively on performance.
- *Learning and Knowledge Acquisition.* An adequate (i.e., computable) definition of novelty might suggest ways of learning a topic more effectively. For example, it could focus the learner's efforts on those aspects of the domain where he has the greatest potential for acquiring something new. This information would help a knowledge-base builder decide which concepts need to be better understood, helping him to direct the dialogue. An analysis of a symbol's dependencies (defined above) might then be used to generate appropriate "probe" sentences to help understand this still vague concept.
- *Representation.* How should a given proposition be indexed? In general each concept should point to all the relevant facts that are *about* that concept. The most obvious approach, based strictly on lexical inclusion, is inadequate. For example, one would want to index " $x + 1 = 0$ " by "x" and not by "+", whereas " $x + y = y + x$ " should be associated with "+" and not with "x".

So how does one determine those concepts that a given fact is really about? We claim that "aboutness" is intimately tied to "newness" in the sense that σ is *about* a concept c whenever this σ expresses something *new* about c with respect to the appropriate diminished theory (which excludes $g[s]$ and all of its consequences).

- *Linguistics.* The basic purpose of communication is for the speaker, S, to transmit a set of *new* facts, usually about some specific topic. To understand this process, we have to know what it means for a fact to be *new* to H and then how S (and H) can use this meta-fact when constructing (or understanding) the message.⁵

⁵This research stemmed from the authors' efforts to understand a particular type of communication, analogy, in terms of this model.

5 Conclusion

While we have yet to solve all the issues associated with this model of novelty, this paper would be incomplete if it did not address the following topics:

- *Applicability*. The New_{PI} relation described above is applicable to any symbol in predicate calculus as well as propositional logic. In particular, the same formalism we saw work for constant symbols works adequately for relation symbols, albeit with an even larger tableau.
- *"Assertional Novelty"*. The novelty we discussed above, New_{PI} , is "definitional", in that its goal is to specify more precisely the referent of a given symbol. Another source of novelty comes from specifying some attribute of the concept; we label such facts "assertionally novel". (See [Woods 75].)

These two categories are distinct: Imagine the symbol RDG had been totally determined, in the sense that the set $I^{\text{Th}}(\text{"RDG"})$ had but a single member. As such, there is nothing New_{PI} we can say about RDG. Despite this certainty, you still might not know what his hair color is. That is, $\text{HairColor}(\text{RDG Brunette})$ might be true in one interpretation, whereas others might hold that $\text{HairColor}(\text{RDG Blond})$. Clearly $\text{HairColor}(\text{RDG Blond})$ is a New_{PI} fact about HairColor ; however, most people would also want this to be a new fact about RDG as well — that is, $\text{New}_{\text{Assert}}(\text{Th, "RDG", "HairColor(RDG Blond)"})$.

- *Intensional, not Extensional*. This paper has dealt exclusively with extensional phenomena, where novelty was determined with respect to the extensions of the symbols. Another approach is intensional — based on descriptions.
- *Deductively Closed*. Probably the most serious criticism of this work is its dependency on a complete deductive system and the requirement that each theory be deductively closed. New-sounding statements can also be used to focus the hearer's attention on some facts he already knew, rather than expose him to new facts. It should be possible to extend this formalism to handle such resource-limited deducibility. Then we could address topics like monotonic novelty and information obsolence.

Each of the issues mentioned above suggests a research task — that of plugging each limitation.

The two issues we find most pressing are:

- Finding an equivalent but syntactical formulation of the semantical New_{PI} relation, in the same manner that N_{Syn} matched N_{Sem} . We hope this will lead to one or more operational versions, of the types mentioned in the beginning of Section 4.
- Expanding this New_{PI} definition to work with deductive systems that are incomplete. (This reiterates the last issue shown above.)

Our basic thesis is that σ is a new fact about A , with respect to the theory Th , if, *under some set of circumstances*, σ limits the number of interpretations of A . New_{PI} achieves this by examining every partial interpretation, testing each to see if A loses a possible interpretation in that situation. This "partial interpretation" definition of context is clearly as general as possible. Furthermore, by counterexample, we have shown that this extreme generality is necessary.

These issues all seem fertile ground for further investigation. Such explorations may very well lead to interesting and usable new results.

Acknowledgments

We would especially like to thank Tom Dietterich for his significant contributions to this work. The following people also contributed: Professor Bruce Buchanan, Dr. Lew Creary, Jim Davidson, Dr. Johan deKleer, Dianne Kanerva, Dr. Jussi Ketonen, Jock Mackinlay, Dr. Robert Moore, Yoram Moses, and Ben Moszkowski. Many useful comments were provided by the reviewers. This basic research was funded by ARPA Contract #MDA903-80-C-0107.

Bibliography

- [Buchanan 78] Buchanan, B. G., Mitchell, T. M., Smith, R. G. and Johnson, C. R. Jr.
Models of Learning Systems.
In *Encyclopedia of Computer Science and Technology*, . Dekker, 1978.
- [Dietterich 81a] Dietterich, T. G. and Buchanan, B. G.
The Role of the Critic in Learning Systems.
Technical Report STAN-CS-81-891, Computer Science Department, Stanford University, December, 1981.
- [Dietterich 81b] Dietterich, T. G. and Michalski, R. S.
Inductive Learning of Structural Descriptions.
Artificial Intelligence 16, 1981.
- [Dietterich 82] Dietterich, T. G., London, R., Clarkson, K., and Dromey, G.
Learning and Inductive Inference.
In Cohen, P. and Feigenbaum, E.A. (editors), *The Handbook of Artificial Intelligence*, . William Kaufman, Inc., Los Altos, CA, 1982.
- [Enderton 72] Enderton, Herbert B.
A Mathematical Introduction to Logic.
Academic Press, Inc., New York, 1972.
- [Gallager 78] Gallager, Robert G.
Information Theory and Reliable Communication.
John Wiley and Sons, Inc., New York, 1978.
- [Langley 79] Langley, Pat.
Rediscovering Physics With BACON.3.
In *6-IJCAI*, pages 505-507. Tokyo, August, 1979.
- [Lenat 82] Lenat, Douglas B.
AM: Discovery in Mathematics as Heuristic Search.
In Davis, Randall and Lenat, Douglas B. (editors), *Knowledge-Based Systems in Artificial Intelligence*, . McGraw-Hill International Book Company, San Francisco, 1982.
- [Michalski 83] Michalski, Ryszard S., Carbonell, Jaime G., and Mitchell, Tom M. (editors).
Machine Learning: An Artificial Intelligence Approach.
Tioga Publishing Company, Palo Alto, CA, 1983.
- [Mitchell 81] Mitchell, Thomas M., Utgoff, Paul E., Nudel, Bernard and Banerji, Ranan.
Learning Problem-Solving Heuristics through Practice.
In *IJCAI-7*, pages 127-134. UBC, 1981.
- [Moore 80] Moore, Robert C.
Reasoning about Knowledge and Action.
Technical Note 191, SRI International, October, 1980.
- [Winston 75] Winston, P. H.
Learning Structural Descriptions from Examples.
McGraw-Hill Book Company, New York, 1975, chapter 5.
- [Woods 75] Woods, W. A.
What's in a Link: Foundations for Semantic Networks.
In D. G. Bobrow & A. M. Collins (editors), *Representation and Understanding*. Academic Press, 1975.

A. A Typical Novelty Situation

This appendix presents a series of "novelty situations", designed to motivate and illustrate our intended meaning of the ternary form of novelty, $\text{New}(\text{Th}, \text{c}, \sigma)$. In each of these, we argue why a certain fact feels "new" with respect to a particular concept, in the context of a particular theory. These examples, expressed in predicate calculus, parallel the propositional logic ones given in the text. They also indicate the assumptions which underlie our intended use of semantic interpretations.

Begin with an initial theory which contains the following three facts:

$$\text{Th}_0 \leftarrow \left\{ \begin{array}{l} (\text{Person RDG}) \\ (\text{Article WN}) \\ (\forall \$P, \$A. [(\text{Author } \$A \$P) \wedge (\text{Laconic } \$A)] \\ \Rightarrow (\text{ShortPaper } \$P) \end{array} \right\}$$

(It may include other sentences as well; but nothing which relates to either RDG or WN.)

Now augment this theory by adding the sentence (In WN IJCAI-83), to form

$$\text{Th}_1 \leftarrow \text{Th}_0 + \text{"(In WN IJCAI-83)"}$$

This fact seems new with respect to WN, but not with respect to RDG. We can express this by

$$\begin{array}{l} \neg \text{New}(\text{Th}_0, \text{"RDG"}, \text{"(In WN IJCAI-83)"}) \\ \text{New}(\text{Th}_0, \text{"WN"}, \text{"(In WN IJCAI-83)"}) \end{array}$$

Later we learn that (Author WN RDG), building

$$\text{Th}_2 \leftarrow \text{Th}_1 + \text{"(Author WN RDG)"}$$

This "authorship" statement seems like a new fact about both WN and RDG — hence,

$$\begin{array}{l} \text{New}(\text{Th}_1, \text{"WN"}, \text{"(Author WN RDG)"}) \\ \text{New}(\text{Th}_1, \text{"RDG"}, \text{"(Author WN RDG)"}) \end{array}$$

We are then told that

$$\text{Th}_3 \leftarrow \text{Th}_2 + \text{"(Laconic RDG)"}$$

Not only does this express something about RDG, but, using the "laconic people write short papers" rule above, we can deduce that (ShortPaper WN). For this reason we want both

$$\begin{array}{l} \text{New}(\text{Th}_2, \text{"RDG"}, \text{"(Laconic RDG)"}) \\ \text{New}(\text{Th}_2, \text{"WN"}, \text{"(Laconic RDG)"}) \end{array}$$

Finally we form

$$\text{Th}_4 \leftarrow \text{Th}_3 + \text{"(Topic WN (Topic WN1))"}$$

We will discuss how this sentence is novel, later.

Some comments are in order.

- A given sentence may be new with respect to one concept and not new with respect to others

(in the context of the same theory):

$\neg \text{New}(Th_0, "RDG", "(In WN IJCAI-83)")$ and
 $\text{New}(Th_0, "WN", "(In WN IJCAI-83)")$.

- The same sentence can be new with respect to several distinct concepts:

$\text{New}(Th_1, "WN", "(Author WN RDG)")$ and
 $\text{New}(Th_1, "RDG", "(Author WN RDG)")$.

- A given sentence may be new with respect to a concept in the context of one theory but not new, with respect to the same concept, in the context of another theory:

$\text{New}(Th_0, "WN", "(In WN IJCAI-83)")$ but
 $\neg \text{New}(Th_1, "WN", "(In WN IJCAI-83)")$.

- A sentence may be new with respect to a symbol not included in that sentence:

$\text{New}(Th_2, "WN", "(Laconic RDG)")$.

(The conclusion of Section 3's first conjecture reiterates this point.)

The definition of New_{FI} , (which is discussed in the second conjecture in Section 3,) is based on the number of interpretations of the symbol in question. For this to make any sense we have to assume (1) that the models involved "overlap", in the sense that each interpretation maps symbols into the same range of possible elements, and (2) that some symbols of the language are "fixed" — that is, their respective extensions are understood and unambiguous. Realize we have implicitly assumed this all along: in this example we "knew" the desired referent of the constant symbol *IJCAI-83* and the predicate symbol *Person*. *IJCAI-83* obviously refers to a particular book; or, stated another way, it designates the same "thing" in every model. Similarly (*Person RDG*) means that *RDG*'s denotation is a person.⁶ Given that these symbols, among others, are fixed, we can build the table shown in Figure A.

That is, all Th_0 knew about *RDG* was that (*Person RDG*), which means that *RDG* could refer to any of the 4.5 billion people in the world. Similarly *WN* could be any of the world's 10 million articles. (As mentioned above, this is only meaningful if the concept *Person* (resp. *Article*) is well defined.)

On hearing that (*In WN IJCAI-83*), we know that *WN* must be one of only about 250 articles in the *IJCAI-83* proceedings. This "10 million to 250" reduction implies that (*In WN IJCAI-83*) is a New_{FI} fact about *WN*. As $\|I^{Th1}("RDG")\| = \|I^{Th2}("WN")\| = 4.5E9$, we see that

⁶This "some symbols are fixed" assumption is not as major restriction as it may seem. In fact, this is the typical of most communication situations. Usually both speaker and hearer have a common "consensual domain" of understanding, in which most concepts are unambiguous and are commonly understood — as the predicate *Person* is in day-to-day discourse. The typical goal of these communication acts is to relay a (new) fact about some other, less-well-defined concept, defining it in terms of those well-understood concepts.

<i>th</i>	<i>facts</i>	$ i^{th}("RDG") $	$ i^{th}("WN") $
	(Article WN) (Person RDG) $\forall \$P, \$A. \text{ IF } [(\text{Author } \$P \$A) \wedge (\text{Laconic } \$A)]$ THEN (ShortPaper \$P)		
...			
Th ₀		4.5E9	1E7
	(In WN IJCAI-83)		
Th ₁		4.5E9	300
	(Author WN RDG)		
Th ₂		300	250
	(Laconic RDG)		
Th ₃		15	10
	(Topic WN (Topic WN1))		
Th ₄		15	10

Figure 3: Number of Interpretations of RDG and WN, in various Theories

(In WN IJCAI-83) was not a New_{FI} fact about RDG, as desired.

We similarly see that

$\text{New}(Th_1, "RDG", "(Author WN RDG))$
 $\text{New}(Th_2, "RDG", "(Laconic RDG))$
 $\text{New}(Th_2, "WN", "(Laconic RDG))$

Unfortunately this analysis leads to the conclusion that $\neg \text{New}_{FI}(Th_3, "WN", "(Topic WN (Topic WN1)))$, which does NOT correspond to our intuitions: While this statement does not restrict the number of interpretations of the symbol, (here WN,) it does establish a dependency of this concept on another, (here WN1). We feel that statements like this, which establish a dependency of one concept on another, should be considered novel. Our argument is that later, when we learn that (Topic WN1 Novelty), we will be able to infer (Topic WN Novelty), which in turn means that WN is then totally specified — that is, that it must refer to one specific article, the original version of the one you are reading now.

Hence, we see that New_{FI} is slightly inadequate, as there are some sentences which our intuitions claim should be deemed novel, which it declares are not. This leads to the New_{PI} relation, defined in

Conjecture 3 of Section 3.

B. Equivalence of N_{Syn} and N_{Sem}

This appendix sketches the proof that N_{Syn} and N_{Sem} are equivalent.

$$N_{Syn}(Th, \sigma) \Rightarrow N_{Sem}(Th, \sigma):$$

Define $Th' \leftarrow Th + \sigma$. As $I_{Th'} \subseteq I_{Th}$, it is sufficient to show that $I_{Th'} \neq I_{Th}$. Assume that σ did not eliminate any member of I_{Th} ; i.e. that this σ "just happened" to be true under all of those interpretations. That would mean that σ is true in all models of Th ; i.e. $Th \models \sigma$. This means that $Th \vdash \sigma$, by Godel's Completeness Theorem. But this directly contradicts the initial assumption that $Th \not\vdash \sigma$, from $N_{Syn}(Th, \sigma)$. Hence there must be some interpretation, I_j , in which this σ is false; and this I_j must therefore be thrown out of $I_{Th'}$, proving that $I_{Th'} \neq I_{Th}$, as desired.

$$N_{Sem}(Th, \sigma) \Rightarrow N_{Syn}(Th, \sigma):$$

As $I_{Th'} \neq I_{Th}$, clearly $Th' \neq Th$, which means that $\sigma \notin Th$. Had $\neg\sigma \in Th$, then Th' would be inconsistent. $I_{Th'}$ would either be undefined, or include every possible interpretation (if you insist); and, in either case, not be contained within I_{Th} . ■

C. Preliminary Connection to Syntactic Methods

Before we can begin to mechanize this novelty relation, we must first express it in a syntactic form. This chapter will present a first shot at this proof-based formulation, for the proposition case. (We feel that some parts of this will scale up, *mutatis mutandis*, to the full predicate calculus case.)

Recall that $New_{pI}(Th, s, \sigma)$ is true when some interpretation of the symbol s , which was possible in some context of the theory Th , is ruled out in the enhanced theory $Th + \sigma$, within the corresponding, still-possible context.

That is,

$$New_{pI}(Th, s, \sigma) \Leftrightarrow \exists \rho, \alpha. \left(\begin{array}{l} (Th \not\vdash \rho \Rightarrow \neg\alpha(s)) \wedge \\ (Th, \sigma \vdash \rho \Rightarrow \neg\alpha(s)) \wedge \\ (Th, \sigma \not\vdash \neg\rho) \end{array} \right)$$

where ρ is a sentence and α is a unary predicate. WLOG this ρ can be a conjunct of assignments to various proposition symbols, i.e. each term is either a symbol or its negation. As such we can think of it as referring to some row of the tableau. (If only expressibility was so trivial in the predicate calculus case...)

Returning to the $\text{New}_{\rho_1}(\text{ThC}, "A", "A \Leftrightarrow C")$ example used in Conjecture 3 of Section 3, we see that setting

$$\begin{array}{lll} \text{Th} & \equiv & \{ A \Leftrightarrow B \} \\ \rho & \equiv & \neg C \quad \text{--- describing the fifth row of the tableaux of Figure 2} \\ \alpha(s) & \equiv & s \quad \text{--- meaning that } \alpha(A) = A \end{array}$$

leads to the conclusion that $\sigma = A \Leftrightarrow C$ is indeed a New_{ρ_1} fact about A in this theory.

Proof of equivalency to the semantic form:

The first two clauses, $\text{Th} \not\vdash (\rho \Rightarrow \neg \alpha(s))$ and $\text{Th}, \sigma \vdash (\rho \Rightarrow \neg \alpha(s))$, mean that there is some interpretation of s , expressed by $\alpha(s)$, which is possible given the theory Th , within the context defined by ρ . In this propositional case we can syntactically specify any partial interpretation (read "row of the tableau",) as a simple sentence.

As there is something NOT derivable from $\text{Th} + \sigma$, that larger theory must be consistent; meaning that Th must have been consistent, and σ not contradictory. The final clause means that ρ is still satisfiable — i.e. that that row has not vanished. ■

Unfortunately, the predicate calculus case is not nearly so simple. First, there are other types of novelty we must worry about here. (Recall that the only type of (extensional) novelty in this propositional case is this definitional case.) Furthermore, even for this definitional situation is complicated by the issue of expressibility: unlike propositional logic, where any interpretation can be syntactically expressed, there are "partial interpretations" which cannot be expressed in predicate calculus.

FILMED FROM BEST AVAILABLE COPY

**Copyright © 1985 by KSL and
Comtex Scientific Corporation**