

JACK AND JANET IN SEARCH OF A THEORY OF KNOWLEDGE

Eugene Charniak
Artificial Intelligence Laboratory
Massachusetts Institute of Technology

Abstract

In order to answer questions about children's stories one needs a great deal of "common sense" knowledge. A model is presented which gives a rough organization to this knowledge along with specifications as to how the information will be accessed. This rough model is then used as a basis for tight arguments about narrow issues (primarily using examples concerning piggy banks.) The paper is intended as an illustration of how one might go about constructing a theory of knowledge.

Acknowledgements

This paper is based on portions of an MIT Ph.D. thesis submitted to the department of Electrical Engineering. The thesis is reproduced as AI Technical Report 266. As in my thesis I would like to thank all the people at the MIT Artificial Intelligence Laboratory who listened to and argued with me on many occasions.

The work reported herein was conducted at the Artificial Intelligence Laboratory, a Massachusetts Institute of Technology research program supported in part by the Advanced Research Projects Agency of the Department of Defense and monitored by the Office of Naval Research under Contract Number H00014-70-A-0362-0003.

1 Introduction

Let us consider the problem of constructing an abstract model of story comprehension. To determine what the model, or program, has "understood" about what it has read, we will ask it questions. So a typical story might start:

- (1) Janet needed some money. She got her piggybank (PB) and started to shake it. Finally some money came out.

Some typical questions would be:

- (2) Why did Janet get the PB?
- (3) Did Janet get the money?
- (4) Why was the PB shaken?

Questions (2) - (4) are not answered explicitly in the text. That is, the story did not say "Janet got her PB because she ..." The story does not even contain a full implicit answer; one cannot logically deduce an answer from the statements in the story without using general knowledge about the

world such as:

- (5) One can often get money from PBs.
- (6) The hard part of getting money from a PB is getting it out. Once that is done one can be said to have the money.
- (7) Shaking helps get money out of a PB.

So in order to understand a children's story we need a theory of every day knowledge. This theory would have to answer questions like "What is the knowledge we have?" and "How is it organized so we can get at the necessary information when it is needed?" Note that this latter question assumes that we have some specific task or tasks in mind, in our case answering questions about children's stories.

The rest of this paper divides into two parts. In the first part a rough description of a model of children's story comprehension will be presented. In the second section we will assume the model presented in the first and look at some narrow questions concerning the organization and content of our knowledge about piggy banks.

2 A Model of Children's Story Comprehension

The model presented here is solely concerned with deduction and does not consider problems of natural language per se. In particular it does not deal with syntax or those problems on the boundary between syntax and deduction like disambiguation of word senses and determination of noun phrase referents. (However, my Ph.D. thesis considers the noun phrase problem in some detail.)

So we will assume that as the story comes into the program it is immediately translated into an internal representation which is convenient for doing deduction. The internal representation will be a group of "assertions" each assertion being a predicate on an arbitrary number of arguments. Putting an assertion into the data base is to "assert" it. The model will try to "fill in the blanks" of the story on a line by line basis. That is, as it goes along, it will try to make connections between events in the story (usually causal connections) and fill in missing facts which seem important such as Janet's now having the money in (1).

2.1 Demons and Base Routines

Consider a fact like:

- (8) If "it is (or will be) raining" and
if "person P is outside"
then "P will get wet"

We have an intuitive belief that (8) is a fact about "rain", rather than, say, a fact about "outside." Many things happen outside and getting wet is only one of them. On the other hand only a limited number of things happen when it rains.

We will embody this belief in our system by associating (8) with "rain" so that only when "rain" comes up in the story will we even consider using rule (8). We will say that rain is the "topic concept" of (8). To put this another way, when a concept is brought up in a story, the facts associated with it are "made available" for use in making deductions. (We will also say that the facts are "put in" or "asserted." So, if "circus", say, has never come up, the program will not be able to make deductions using those facts associated only with "circus."

Note however that we are not saying that "rain" has to be mentioned explicitly in the story before we can use (8). It is only necessary that there be a "rain" assertion put into the data base. Other parts of the story may provide facts which cause the program to assert that it is raining. For example:

- (9) One afternoon Jack was outside playing ball with Bill. Bill looked up and noticed that the sky was getting dark. "I think we should stop" said Bill. "We will get wet if we keep playing."

Here, the sky's getting dark in the afternoon suggests that it is going to rain. If this is put into the data base it will be sufficient to bring in facts associated with "rain."

Also note that a topic concept need not be a single "key word." A fact may not become available to the system until a complex set of relations appears in the data base. A fact may be arbitrarily complex, and in particular may activate other facts depending on the presence or absence of certain relations in the story.

Looking Forward, Looking Back. When a fact is made available we might not have all the information needed to make use of the fact. Since we are making deductions as we go, if the necessary information comes in after the rule has been asserted we want to make the deduction when the information comes in. So we might have:

- (10) Jack was outside. It was raining.
(11) It was raining. Jack was outside.

In (10) there is no problem. When we introduce "rain" we have sufficient information to use (8) and deduce that Jack is going to get wet. But in (11), we only learn that Jack is outside after we have mentioned rain. If we want to use (8) we will need some way to have our fact "look forward" in the story. To do this we will break a fact into two parts, a pattern and a

body (an arbitrary program). We will execute the body of the fact only when an assertion is in the data base which matches the pattern. (We will also say that the assertion "excites" the fact.) In (8) the pattern would be "someone outside." Then in (11) when we introduce (8) no assertion matches the pattern. But the next line creates a matching assertion, so the fact will be excited. We will say that a fact is "looking forward" when its topic concept appears before the assertion which matches the pattern. When the assertion which matches the pattern comes first we will say that the fact is "looking backward" (as in 10).

We can see how important looking forward is with a few examples.

- (12) "Janet was thinking of getting Jack a ball for his birthday. When she told Penny, Penny said, 'Don't do that. Jack has a ball.'" Here we interpreted the line "Jack has a ball" as meaning that he did not want another. The common sense knowledge is the fact that in many cases having an X means that one will not want another X. This piece of information would probably be filed under "things to consider when about to get something for somebody else." Naturally it was an earlier line which mentioned that Janet was thinking of getting Jack a ball.
- (13) "Bill offered to trade his pocket knife for Jack's dog Tip. Jack said 'I will ask Janet. Tip is her dog too.'" The last line is interpreted as the reason Jack will ask Janet. This requires information about the relation between trading and ownership.
- (14) "Janet wanted to get some money. She found her piggy bank and started to shake it. She didn't hear anything." The last line means that there was nothing in the piggy bank on the basis of facts about piggy banks.

In each of these cases it is an earlier line which contains the information which is used to assign the interpretation. So in (12) there is nothing inherent in the line "Jack has a ball" which means "don't get him another." If there were, something in the line would also have to key a check for the following situations:

- (15) Bill and Dick wanted to play baseball. When Jack came by Bill said "There is Jack. He has a ball."
- (16) Tom asked his father if he would buy him a ball. "Jack has a ball," said Tom.

- (17) Bill's ball of string was stuck in the tree. He asked Jane how he could get it out. Jane said "You should hit it with something. Here comes Jack. He has a ball."

Those familiar with Planner might notice that our "facts" look quite similar to Planner antecedent theorems, with the exception that our facts can "look back" as well as "look forward." Antecedent theorems are only designed to look forward. I initially formulated facts as antecedent theorems because I was so impressed with the need to "look forward." However, rather than call the facts antecedent theorems, I call them "demons" since it is a shorter name.

Specification and Removal of Demons.

It should be emphasized that the model does not "learn" the information contained in the demons. This information is put in by the model maker. On the other hand, the demons are not specific to the story in the sense that they mention Jack, or "the red ball." Rather, they talk about "a person X" who at one point in the story could be Jack, at another, Bill. We will assume a mechanism for binding some of the variables of the demon ("specifying" the demon) at the time the demon is asserted.

We want demons to be active only while they are relevant to the story. A story may start by talking about getting a present for Jack, but ultimately revolve around the games played at his party. We will need some way to remove the "present getting" demons when they have outlived their usefulness. (An irrelevant but active demon not only wastes time and space, but can cause us to misinterpret a new line.) As a first approximation we will assume that a demon is declared irrelevant after a given number of lines have gone by.

Base Routines. So far we have said that demons are asserted when the proper concept has been mentioned. But this implies that there is something attached to the concept name telling us what demons should be put in.

If we look at a particular example, say (13), it is Bill's offer to trade which sets up the context for the rest of the fragment. I will assume that the information to do so is in the form of a program. Such routines, which are available to set up demons, will be called "base routines."

These base routines will be responsible for more than setting up demons. Suppose we are told that Jack had a ball, and Bill a top. Then Jack traded his ball to Bill for the top. One question we might ask is "who now has the top?" Naturally since questions of "who has what" are important in understanding stories we will want to keep tabs on such information. In this particular case, it must again be the "trade" statement which tells us to switch possession of the objects. Every time a trade occurs we will want to exchange objects, so whenever we see "trade" we execute the "trade" base routine. Of course, the program can't be too simple-minded, since it must also handle "I will trade..." and perhaps even "Will you trade

...?"

A good test as to whether a given fact should be part of a base routine or a demon is whether we need several lines to set it up or whether we can illustrate the fact by presenting a single line. (Naturally several lines could be made into one by putting "and's" between them, but this is dodging the point. I am only suggesting an intuitive test.) So we saw that "Jack has a ball" was not enough by itself to tell us that Jack does not want another ball. Hence this relation is embodied by a demon, not a base routine. On the other hand, often a single line can tell us quite a bit as in "Jack was on second base." This indicates that the base routine for "second base" can often tell us that we are talking about a baseball game.

2.2 Bookkeeping and Fact Finders

Updating and Bookkeeping. Up to this point we have introduced two parts of the model, demons and base routines. In this section we will introduce the remaining two parts.

Again let us consider the situation when Jack had a ball, Bill a top, and they traded. When we say that Bill now has the ball, it implies that Jack no longer does. That is to say, we must somehow remove the fact that Jack has the ball from the data base. Actually we don't want to remove it, since we may be asked "Who had the ball before Bill did." Instead, we want to mark the assertion in some way to indicate that it has been updated. We will assume that there is a separate section, pretty much independent of the rest of the model, which is responsible for doing such updating. We will call this section "bookkeeping."

Fact Finders. But even deciding that one statement updates another requires special knowledge. Suppose we have:

- (18) Jack was in the house. Sometime later he was at the store.

If we ask "Is Jack in the house?" we want to answer "no, he is at the store." But how is bookkeeping going to figure this out? There is a simple rule which says that (<state> A B) updates (<state> A C) where C is not the same as B. So (AT JACK FARM!) would update (AT JACK NEW-YORK). But in (18) we can't simply look for Jack AT <someplace which is not the store>, since he is !!! the house. To make things even worse, we could have:

- (19) Jack was in the house. Sometime later he was in the kitchen.

To solve this problem we will need:

- (20) To establish that PERSON is not at location LOC.

Find out where PERSON is, call it X.
 If X = LOC, then theorem is false so return "No."
 If X is part of LOC then return "No."
 If LOC is part of X, then try to find a different X.
 Else return "Yes."

In (18) the bookkeeper would try to prove that Jack is not at the store, and it would succeed by using (20) and the statement that Jack is in the house. Bookkeeper would then mark the earlier statement as updated.

Theorems like (20) are called "fact finders." Like demons, fact finders have a pattern and a body. A particular fact finder is called when either a demon, base routine or bookkeeping wants to establish a goal which matches the fact finder's pattern. This is different from demons which are called when we encounter a given fact.

The basic idea behind fact finders is that they are used to establish facts which are comparatively unimportant, so that we do not want to assert them and hence have them in the data base. So in (18) we do not want to assert "Jack is not in the house" as well as "Jack is at the store." In the same way we will have a fact finder which is able to derive "<person> knows <fact>" by asking such questions as "was the <person> there when <fact> was mentioned or took place?" Again, this information is easily derivable, and not all that important, so there would seem to be no reason to include it explicitly in the data base.

3 Some Narrow Questions

In section 1 we stated that our theory of knowledge should answer questions like "how do we access the information" and "what is the information." In this portion of the paper we will look at two problems, one of each kind. We start with a question of information access.

3.1 Demon-Demon Interaction

In the description of the model it was stated that demons are excited when an assertion enters the data base which matches the demon's pattern. In this section we will present evidence that, given the model of part 1, we must also allow demons to excite other demons. I call this "demon-demon interaction."

A Demon About PBs. Before we can talk about demon-demon interaction we need to establish the need for two particular demons. Suppose we were given:

- (21) Janet needed money. She got her piggy bank.
 (22) Janet got her PB. "I want a nickel" she said.

Were we asked what Janet is going to do with the PB we would say, "get money from it." The obvious way to handle this is with a demon which is declared relevant when we see a person getting a PB. Naturally this would be done by the PB base routine. This demon would be looking for "<getter> need <money>" (i.e., that would be its pattern). When excited the demon would assert that there is a causal relation between "need money" and "get PB." We will call this demon PB-NEED-MONEY.

Now we might claim that "want money" should put in a demon to look for "get PB" rather than vice versa. However, this seems to be less reasonable, since there are many ways of getting money, but only a very limited number of reasons a person gets a PB.

A Demon About Buying. Buying things often requires money. So if we saw

- (23) Janet was going to buy some candy. She needed some money.
 (24) Janet needed some money. She was going to buy some candy.

we would assert that the reason she needs money is to buy candy. Since the "need money" statement can occur on either side of the "buy" statement, it is clear that we want "buy" to put in a demon which says "If the 'buyer' needs money, it is because of the 'buying'." We will call this the BUY-NEED-MONEY demon.

We do not want "need money" to put in a demon looking for "buy." There are many things one can do with money, bribe a juror, pay rent, take a taxi, tip a bellboy, etc. It would seem more obvious to have the various events state that they (usually) require money, rather than have "need money" look for all of them. Nor does it seem reasonable to claim that all of these events ("bribing", etc.) are really versions of "buying." To express "take a taxi" as "buy some of the taxi driver's time plus the temporary use of the automobile in order to convey oneself from one location to another" seems somewhat forced. Most of this "definition" comes from knowledge of economics rather than taxis.

Evidence of the Phenomenon. But now consider the following fragment:

- (25) Janet was going to buy some candy. She went to get her PB.

Here we want to assert that Janet gets her PB because she wants money. But this time there is no previous assertion which says that she needs money. Of course, what is at work here is the fact that buying requires money. But, we have represented this fact as a demon, and so far we have no way for two demons to interact with one another. Now both PB-NEED-MONEY and BUY-NEED-MONEY will have the same pattern. So we can account for (25) by allowing a demon to excite other demons with the same pattern. Naturally this is what I mean by demon-demon interaction.

A Restriction on Demon-Demon Interaction. Demon-demon interaction is

probably more complex than we have indicated so far. Consider:

- (26) Janet was going to buy some candy. She was also going to buy some fruit.

In (26) both occurrences of "buy" will activate BUY-NEED-MONEY demons. (Though we did not comment on this earlier, the idea of specifying demons as mentioned in 2.1 obviously requires separate copies of a demon to be able to exist simultaneously.) However, (26) does not imply that Janet really needs money. For all we know she has as much as she needs in her pocket. If demon-demon interaction were as simple as we have made it out to be, the two instances of BUY-NEED-MONEY would join up to produce a "need money" assertion. So it is not sufficient for two demons to be looking for the same pattern.

Looking at example (25) we note that one of the demons gave a reason why Janet might need money, and the second suggested that needing money was the cause of a certain action. So we have:

Will buy --> Need money --> Will get PB

To put this in everyday terms, in (25) we have both a motive for needing money (buying), and a result of needing the money (go and get PB). In (26) we have two motives. The natural suggestion is that demon-demon interaction be restricted to cases where we have both motive and result.

How do we recognize when we have both motive and result? As it stands now one demon looks pretty much like any other. We might just try to label all demons as "motive" or "result" with respect to their pattern. On the other hand it might be possible to derive "motive" and "result" from more basic considerations. At any rate, it seems premature to formalize such concepts at this point. We simply don't know enough.

Capturing Generalizations. Before moving on I should point out that the kind of argument used in this section (and the next also) is a "capture the generalization" type argument commonly found in linguistics. We could have created a new demon to explain (25). It would have said, "If a person gets his PB look for him planning to buy something." However, this would be missing the generalization that "motives" and "results" always act this way. So far I have only given one example to support the "demon-demon interaction generalization", but in in the next section we will see another.

3.2 Putting Money into a Piggy Bank

In this section we will look at a possible demon associated with piggy banks and argue that the deduction it would account for can be better handled by demon-demon interaction between two other demons. In effect we will be trying to determine, on an extremely small scale, what people know.

A Piggy Bank Problem. One fact we know about PB's is that they are good places to keep money. This fact seems to come into play in:

- (27) Penny said to Janet, "Don't take your money with you to the park. (You will lose it.) Go and get your PB!"
- (28) After Janet helped Ms. Jones with her groceries Ms. Jones gave her a dime. Jack came along and said, "Come with me to the park, Janet." "OK," said Janet. "But first I am going home to find my PB. I do not want to take the money to the park."
- (29) Janet put some money on the sink. Mother said, "If you leave the money there it may fall in the drain. Let's find your PB."

In each case the natural question is, "Why should Janet get her PB?" Now we might try to construct a "piggy bank" demon which responds to some common element in (27) - (29) and then make the necessary assertions. A close look at the examples even gives a start at what such a common element might be, say "a particular location for the money is negatively evaluated." We will call this demon PB-BAD-PLACE. The trouble with such a solution would be that it would not account for:

- (30) Janet said, "I am going to put my money away. I will get my PB."
- (31) Janet helped Ms. Jones with her groceries. Ms. Jones gave Janet a dime. Jack came along and said, "Janet, let's go to the park." "OK," said Janet. "But I want to put my money in a safe place. I am going to get my PB."

Now there is nothing saying that our demon needs to account for (30) and (31). However, it seems quite obvious that we are using the same information in all the examples above. The only difference is that in (27) - (29) we are expressing the need for a "safe place" by making negative comments about another location. If this is a single fact we would like a single demon to express it. The trouble is finding what (27) - (31) have in common.

A Non-Piggy Bank Problem. In the course of looking at examples like (27) - (31) I noted examples like:

- (32) Penny said to Janet, "Don't take your money with you to the park. Put it on the shelf."
- (33) After Janet helped Ms. Jones with her groceries Ms. Jones gave her a dime. Jack came along and said "Come with me to the park, Janet." "OK," said Janet. "But first I am going to put my money in the house. I do not want to take the money to the park."

- (34) Janet put some money on the sink. Mother said, "If you leave the money there it may fall in the drain." Janet put the money in a drawer.
- (35) Janet said "I am going to put my money away. I will put it in my toy box."
- (36) Janet helped Ms. Jones with her groceries. Ms. Jones gave Janet a dime. Jack came along and said "Janet, let's go to the park." "OK," said Janet. "But I want to put my money in a safe place." Then Janet went into the house and put the money in her room.

These examples exactly mirror (27) - (31), except that (32) - (36) don't mention PB's. Naturally, in these examples the question to ask is "Why did Janet put the money in the drawer?", "in the house?", etc.

Such examples tend to indicate that the problem facing us is wider than just PB's. We will name this wider problem the "put away" problem. However it is not the case that our problem with PB's can be completely reduced to the "put away" problem. So while in the non-piggy bank examples we mention that Janet has or actually intends to "put" the money some place, in the PB examples all we needed to say was that Janet was going to get the PB. To put this another way, our knowledge of PB's allowed us to interpret "get PB" as meaning that Janet was going to put money into it. However our knowledge of houses or shelves does not allow us to make similar deductions in (32) - (36).

The Put-Away Demon. Ignoring piggy banks for the moment, what would a solution to (32) - (36) look like? We will have some demon, called the PUT-AWAY demon, which is activated by lines like:

- (37) Don't leave the money by the sink.
 (38) I do not want to take my money to the park.
 (39) I will put my money away.

These lines will put in a demon looking for "put away" and the demon will assert that the reason for putting the thing away is (37) - (39). Ultimately we will want a theory of why people put things away (i.e., what lines put in the "put away" demon), and how to determine what constitutes a "put away" location. However, (32) - (36) clearly show that the problem is distinct from the question of what we know about PBs.

The Piggy Bank Demon. What we will now see is that if we assume the PUT-AWAY demon, all the examples in (27) - (31) fall out easily, plus a few others which we haven't even looked at yet. But first we need to consider a new PB demon entitled PB-MONEY-IN. It is parallel to PB-NEED-MONEY, but while the latter was for recognizing that money was going to be taken out of the PB, PB-MONEY-IN is for recognizing that money is going to be put in. It says "if you see that the person wants some money to be in the PB then the

reason he is getting the PB is to put it in." (Actually this theorem is true of a wide class of containers, but that does not affect the argument at hand.) This demon will account for examples like:

- (40) Ms. Jones gave Janet a dime. Janet went to get her PB. "I want the money to be in my PB," she thought.
- (41) Janet got her PB and dropped some money in.
- (42) After Ms Jones gave Janet a dime, Jack came by and asked Janet if she wanted to go to the park. "OK," said Janet. "I will go home first and get my PB." Soon Janet came back and said "My money is in the PB, let's go!"

Demon-Demon Interaction. Now, if we assume demon-demon interaction as discussed in section 3.1, PB-MONEY-IN plus PUT-AWAY will interact to solve all the examples from (27) to (31). Let us see how this will happen.

First note that the restrictions we placed on demon-demon interactions are met here. First both demons have the same pattern, e.g., "money is in PB." (Actually the pattern for PUT-AWAY is "<object> is in <appropriate location>" however <object> will be bound to the money at the time the demon is asserted, and <appropriate location> will match PB when the demon is excited.) Secondly, we need both a motive and a result before we can "combine" demons. In the case at hand, PUT-AWAY is a motive for having the money in the PB, and "get PB" is a result of intending to put money in the PB.

Saving Money. Finally, note that our solution extends to the following case:

- (43) Janet got a dime from Ms. Jones. She said "I am saving my money to buy a bicycle. I am going home to get my PB."

Here we know that Janet is going to put the money in the PB because of the "save" statement. However, we immediately note that we have cases like:

- (44) Janet got a dime from Ms. Jones. Janet told her "I am saving my money to buy a bicycle. I am going home to put the money away. (I am going home to put the money in my drawer.)"

Naturally, (44) indicates that "save" must activate PUT-AWAY. If this is the case, then (43) is accounted for in exactly the same manner as all the initial examples. While the reader may not be surprised at this result, I am, since initially I thought that the relationship of "save" with piggy banks would need a separate PB demon.

4 Conclusion

The two halves of this paper stand in contrast to each other. The presentation of the model (section 2) is general (in theory covering all of children's stories), but vague and full of covert appeals to the reader's intuition. Section 3 on the other hand is narrow, only talking about small portions of our knowledge of PBs, but tightly reasoned (hopefully).

How by themselves the conclusions of section 3 are not that important. Of course, if we could pin down one hundred facts the way we pinned down one in section 3.2 then we would have the beginnings of a theory of knowledge. But I did not write this paper to tell of one fact about PBs. Rather I view the paper as an illustration of how one might go about the task of constructing a theory of knowledge.