

---

# **Knowledge Representation and Reasoning**

---

from the book of the same name  
by  
Ronald J. Brachman  
and  
Hector J. Levesque

Morgan Kaufmann Publishers, San Francisco, CA, 2004

---

1.

# Introduction

# What is knowledge?

---

Easier question: how do we talk about it?

We say “John knows that ...” and fill the blank with a proposition

– can be true / false, right / wrong

Contrast: “John fears that ...”

– same content, different attitude

Other forms of knowledge:

- know how, who, what, when, ...
- sensorimotor: typing, riding a bicycle
- affective: deep understanding

Belief: not necessarily true and/or held for appropriate reasons

and weaker yet: “John suspects that ...”

Here: no distinction

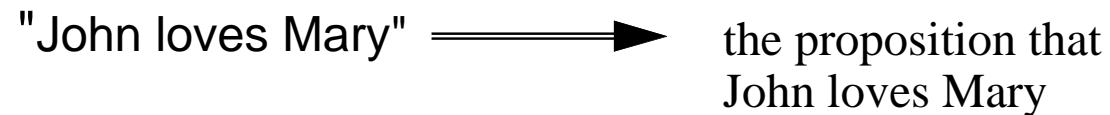
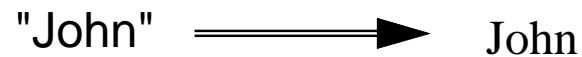
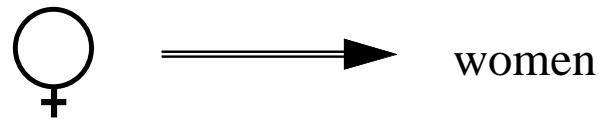
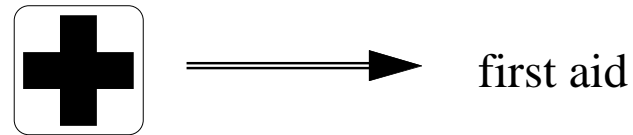
the main idea

taking the world to be one way and not another
--

# What is representation?

---

Symbols standing for things in the world



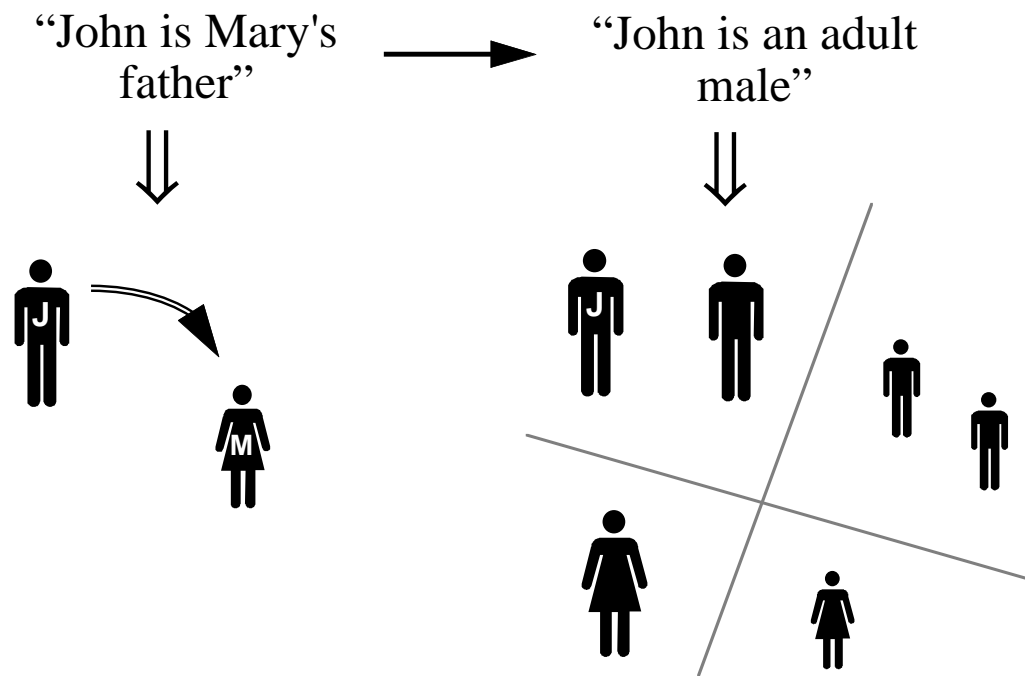
Knowledge representation:

symbolic encoding of propositions believed  
(by some agent)

# What is reasoning?

Manipulation of symbols encoding propositions to produce representations of new propositions

Analogy: arithmetic    "1011" + "10" → "1101"  
                                  ↓          ↓          ↓  
                                  eleven  two      thirteen



# Why knowledge?

---

For sufficiently complex systems, it is sometimes useful to describe systems in terms of beliefs, goals, fears, intentions

e.g. in a game-playing program

“because it believed its queen was in danger, but wanted to still control the center of the board.”

more useful than description about actual techniques used for deciding how to move

“because evaluation procedure P using minimax returned a value of +7 for this position

= taking an intentional stance (Dan Dennett)

Is KR just a convenient way of talking about complex systems?

- sometimes anthropomorphizing is inappropriate  
e.g. thermostats
- can also be very misleading!  
fooling users into thinking a system knows more than it does

# Why representation?

---

Note: intentional stance says nothing about what is or is not represented symbolically

e.g. in game playing, perhaps the board position is represented, but the goal of getting a knight out early is not

KR Hypothesis: (Brian Smith)

“Any mechanically embodied intelligent process will be comprised of structural ingredients that a) we as external observers naturally take to represent a propositional account of the knowledge that the overall process exhibits, and b) independent of such external semantic attribution, play a formal but causal and essential role in engendering the behaviour that manifests that knowledge.”

Two issues: existence of structures that

- we can interpret propositionally
- determine how the system behaves

Knowledge-based system: one designed this way!



# Two examples

---

## Example 1

```
printColour(snow) :- !, write("It's white.").
printColour(grass) :- !, write("It's green.").
printColour(sky) :- !, write("It's yellow.").
printColour(X) :- write("Beats me.").
```

## Example 2

```
printColour(X) :- colour(X,Y), !,
    write("It's "), write(Y), write(".").
printColour(X) :- write("Beats me.").

colour(snow,white).
colour(sky,yellow).
colour(X,Y) :- madeof(X,Z), colour(Z,Y).
madeof(grass,vegetation).
colour(vegetation,green).
```

Both systems can be described intentionally.

Only the 2nd has a separate collection of symbolic structures à la KR Hypothesis

its knowledge base (or KB)

∴ a small knowledge-based system

# KR and AI

---

Much of AI involves building systems that are knowledge-based

ability derives in part from reasoning over explicitly represented knowledge

- language understanding,
- planning,
- diagnosis,
- “expert systems”, etc.

Some, to a certain extent

game-playing, vision, etc.

Some, to a much lesser extent

speech, motor control, etc.

Current research question:

how much of intelligent behaviour is knowledge-based?

Challenges: connectionism, others

# Why bother?

---

Why not “compile out” knowledge into specialized procedures?

- distribute KB to procedures that need it  
(as in Example 1)
- almost always achieves better performance

No need to think. *Just do it!*

- riding a bike
- driving a car
- playing chess?
- doing math?
- staying alive??

Skills (Hubert Dreyfus)

- novices think; experts *react*
- compare to current “expert systems”:  
knowledge-based !

# Advantage

---

Knowledge-based system most suitable for *open-ended* tasks

can structurally isolate *reasons* for particular behaviour

## Good for

- explanation and justification
  - “Because grass is a form of vegetation.”
- informability: debugging the KB
  - “No the sky is not yellow. It's blue.”
- extensibility: new relations
  - “Canaries are yellow.”
- extensibility: new applications
  - returning a list of all the white things
  - painting pictures

# Cognitive penetrability

---

Hallmark of knowledge-based system:

the ability to be *told* facts about the world and adjust our behaviour correspondingly

for example: read a book about canaries or rare coins

Cognitive penetrability (Zenon Pylyshyn)

actions that are conditioned by what is currently believed

an example:

we normally leave the room if we hear a fire alarm

we do not leave the room on hearing a fire alarm

if we believe that the alarm is being tested / tampered

can come to this belief in very many ways

so this action is cognitively penetrable

a non-example:

blinking reflex

# Why reasoning?

---

Want knowledge to affect action

*not* do action  $A$  if sentence  $P$  is in KB

*but* do action  $A$  if world believed in satisfies  $P$

Difference:

$P$  may not be *explicitly* represented

Need to apply what is known in general  
to the particulars of a given situation

Example:

“Patient  $x$  is allergic to medication  $m$ .”

“Anybody allergic to medication  $m$  is also  
allergic to  $m'$ .”

Is it OK to prescribe  $m'$  for  $x$  ?

Usually need more than just DB-style retrieval of facts in the KB

# Entailment

---

Sentences  $P_1, P_2, \dots, P_n$  entail sentence  $P$  iff the truth of  $P$  is implicit in the truth of  $P_1, P_2, \dots, P_n$ .

If the world is such that it satisfies the  $P_i$  then it must also satisfy  $P$ .

Applies to a variety of languages (languages with truth theories)

Inference: the process of calculating entailments

- sound: get only entailments
- complete: get all entailments

Sometimes want unsound / incomplete reasoning

for reasons to be discussed later

Logic: study of entailment relations

- languages
- truth conditions
- rules of inference

# Using logic

---

## No universal language / semantics

- Why not English?
- Different tasks / worlds
- Different ways to carve up the world

## No universal reasoning scheme

- Geared to language
- Sometimes want “extralogical” reasoning

## Start with first-order predicate calculus (FOL)

- invented by philosopher Frege for the formalization of mathematics
- but will consider subsets / supersets and very different looking representation languages



# Knowledge level

---

Allen Newell's analysis:

- Knowledge level: deals with language, entailment
- Symbol level: deals with representation, inference

Picking a logic has issues at each level

- Knowledge level:
  - expressive adequacy,
  - theoretical complexity, ...
- Symbol level:
  - architectures,
  - data structures,
  - algorithmic complexity, ...

Next: we begin with FOL at the knowledge level

---

2.

# The Language of First-order Logic

# Declarative language

---

Before building system

before there can be learning, reasoning, planning,  
explanation ...

need to be able to express knowledge

Want a precise declarative language

- declarative: believe  $P$  = hold  $P$  to be true  
cannot believe  $P$  without some sense of  
what it would mean for the world to satisfy  $P$
- precise: need to know exactly  
what strings of symbols count as sentences  
what it means for a sentence to be true  
(but without having to specify which ones are true)

Here: language of first-order logic

again: not the only choice

# Alphabet

---

## Logical symbols:

- Punctuation: (, ), .
- Connectives:  $\neg$ ,  $\wedge$ ,  $\vee$ ,  $\forall$ ,  $\exists$ , =
- Variables:  $x, x_1, x_2, \dots, x', x'', \dots, y, \dots, z, \dots$   
Fixed meaning and use  
like keywords in a programming language

## Non-logical symbols

- Predicate symbols (like Dog) **Note:** not treating = as a predicate
- Function symbols (like bestFriendOf)  
Domain-dependent meaning and use  
like identifiers in a programming language

Have arity: number of arguments

arity 0 predicates: propositional symbols

arity 0 functions: constant symbols

Assume infinite supply of every arity

# Grammar

---

## Terms

1. Every variable is a term.
2. If  $t_1, t_2, \dots, t_n$  are terms and  $f$  is a function of arity  $n$ , then  $f(t_1, t_2, \dots, t_n)$  is a term.

## Atomic wffs (well-formed formula)

1. If  $t_1, t_2, \dots, t_n$  are terms and  $P$  is a predicate of arity  $n$ , then  $P(t_1, t_2, \dots, t_n)$  is an atomic wff.
2. If  $t_1$  and  $t_2$  are terms, then  $(t_1=t_2)$  is an atomic wff.

## Wffs

1. Every atomic wff is a wff.
2. If  $\alpha$  and  $\beta$  are wffs, and  $v$  is a variable, then  $\neg\alpha$ ,  $(\alpha\wedge\beta)$ ,  $(\alpha\vee\beta)$ ,  $\exists v.\alpha$ ,  $\forall v.\alpha$  are wffs.

The propositional subset: no terms, no quantifiers

Atomic wffs: only predicates of 0-arity:  $(p \wedge \neg(q \vee r))$

# Notation

---

Occasionally add or omit (,), .

Use [,] and {,} also.

Abbreviations:

$(\alpha \supset \beta)$  for  $(\neg\alpha \vee \beta)$

safer to read as disjunction than as “if ... then ...”

$(\alpha \equiv \beta)$  for  $((\alpha \supset \beta) \wedge (\beta \supset \alpha))$

Non-logical symbols:

- Predicates: mixed case capitalized

Person, Happy, OlderThan

- Functions (and constants): mixed case uncapitalized

fatherOf, successor,  
johnSmith

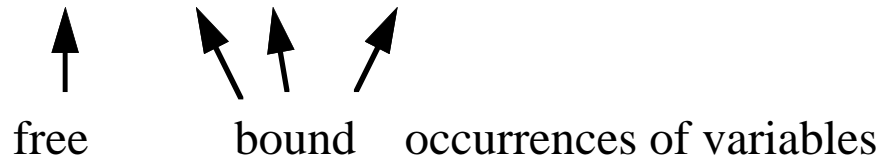
# Variable scope

---

Like variables in programming languages, the variables in FOL have a scope determined by the quantifiers

Lexical scope for variables

$$P(x) \wedge \exists x[P(x) \vee Q(x)]$$



A sentence: wff with no free variables (closed)

Substitution:

$\alpha[v/t]$  means  $\alpha$  with all free occurrences of the  $v$  replaced by term  $t$

Note: written  $\alpha_t^v$  elsewhere (and in book)

Also:  $\alpha[t_1, \dots, t_n]$  means  $\alpha[v_1/t_1, \dots, v_n/t_n]$

# Semantics

---

How to interpret sentences?

- what do sentences claim about the world?
- what does believing one amount to?

Without answers, cannot use sentences to represent knowledge

Problem:

cannot fully specify interpretation of sentences because non-logical symbols reach outside the language

So:

make clear dependence of interpretation on non-logical symbols

Logical interpretation:

specification of how to understand predicate and function symbols

Can be complex!

DemocraticCountry, IsABetterJudgeOfCharacterThan,  
favouriteIceCreamFlavourOf, puddleOfWater27



# The simple case

---

There are objects.

some satisfy predicate  $P$ ; some do not

Each interpretation settles extension of  $P$ .

borderline cases ruled in separate interpretations

Each interpretation assigns to function  $f$  a mapping from objects to objects.

functions always well-defined and single-valued

The FOL assumption:

*this is all you need to know about the non-logical symbols to understand which sentences of FOL are true or false*

In other words, given a specification of

- » what objects there are
- » which of them satisfy  $P$
- » what mapping is denoted by  $f$

it will be possible to say which sentences of FOL are true

# Interpretations

---

Two parts:  $\mathcal{I} = \langle D, I \rangle$

$D$  is the domain of discourse

can be *any* non-empty set

not just formal / mathematical objects

e.g. people, tables, numbers, sentences, unicorns, chunks of peanut butter, situations, the universe

$I$  is an interpretation mapping

If  $P$  is a predicate symbol of arity  $n$ ,

$$I[P] \subseteq D \times D \times \dots \times D$$

an  $n$ -ary relation over  $D$

If  $f$  is a function symbol of arity  $n$ ,

$$I[f] \in [D \times D \times \dots \times D \rightarrow D]$$

an  $n$ -ary function over  $D$

for propositional symbols,

$$I[p] = \{\} \text{ or } I[p] = \{\langle \rangle\}$$

for constants,  $I[c] \in D$

In propositional case, convenient to assume

$$\mathcal{I} = I \in [\text{prop. symbols} \rightarrow \{\text{true, false}\}]$$

# Denotation

---

In terms of interpretation  $\mathcal{I}$ , terms will denote elements of the domain  $D$ .

will write element as  $\|t\|_{\mathcal{I}}$

For terms with variables, the denotation depends on the values of variables

will write as  $\|t\|_{\mathcal{I},\mu}$

where  $\mu \in [Variables \rightarrow D]$ ,  
called a variable assignment

Rules of interpretation:

1.  $\|v\|_{\mathcal{I},\mu} = \mu(v)$ .
2.  $\|f(t_1, t_2, \dots, t_n)\|_{\mathcal{I},\mu} = H(d_1, d_2, \dots, d_n)$   
where  $H = I[f]$   
and  $d_i = \|t_i\|_{\mathcal{I},\mu}$ , recursively

# Satisfaction

---

In terms of an interpretation  $\mathcal{I}$ , sentences of FOL will be either true or false.

Formulas with free variables will be true for some values of the free variables and false for others.

Notation:

will write as  $\mathcal{I}, \mu \models \alpha$     “ $\alpha$  is satisfied by  $\mathcal{I}$  and  $\mu$ ”

where  $\mu \in [Variables \rightarrow D]$ , as before

or  $\mathcal{I} \models \alpha$ , when  $\alpha$  is a sentence

“ $\alpha$  is true under interpretation  $\mathcal{I}$ ”

or  $\mathcal{I} \models S$ , when  $S$  is a set of sentences

“the elements of  $S$  are true under interpretation  $\mathcal{I}$ ”

And now the definition...

# Rules of interpretation

---

1.  $\mathcal{I}, \mu \models P(t_1, t_2, \dots, t_n)$  iff  $\langle d_1, d_2, \dots, d_n \rangle \in R$   
where  $R = I[P]$   
and  $d_i = \llbracket t_i \rrbracket_{\mathcal{I}, \mu}$  as on denotation slide
2.  $\mathcal{I}, \mu \models (t_1 = t_2)$  iff  $\llbracket t_1 \rrbracket_{\mathcal{I}, \mu}$  is the same as  $\llbracket t_2 \rrbracket_{\mathcal{I}, \mu}$
3.  $\mathcal{I}, \mu \models \neg \alpha$  iff  $\mathcal{I}, \mu \not\models \alpha$
4.  $\mathcal{I}, \mu \models (\alpha \wedge \beta)$  iff  $\mathcal{I}, \mu \models \alpha$  and  $\mathcal{I}, \mu \models \beta$
5.  $\mathcal{I}, \mu \models (\alpha \vee \beta)$  iff  $\mathcal{I}, \mu \models \alpha$  or  $\mathcal{I}, \mu \models \beta$
6.  $\mathcal{I}, \mu \models \exists v \alpha$  iff for some  $d \in D$ ,  $\mathcal{I}, \mu\{d;v\} \models \alpha$
7.  $\mathcal{I}, \mu \models \forall v \alpha$  iff for all  $d \in D$ ,  $\mathcal{I}, \mu\{d;v\} \models \alpha$   
where  $\mu\{d;v\}$  is just like  $\mu$ , except that  $\mu(v)=d$ .

For propositional subset:

$$\mathcal{I} \models p \quad \text{iff} \quad I[p] \neq \{\}$$
 and the rest as above

# Entailment defined

---

Semantic rules of interpretation tell us how to understand all wffs in terms of specification for non-logical symbols.

But some connections among sentences are independent of the non-logical symbols involved.

e.g. If  $\alpha$  is true under  $\mathcal{I}$ , then so is  $\neg(\beta \wedge \neg\alpha)$ ,  
no matter what  $\mathcal{I}$  is, why  $\alpha$  is true, what  $\beta$  is, ...

$S \models \alpha$  iff for every  $\mathcal{I}$ , if  $\mathcal{I} \models S$  then  $\mathcal{I} \models \alpha$ .

Say that  $S$  entails  $\alpha$  or  $\alpha$  is a logical consequence of  $S$ :

In other words: for no  $\mathcal{I}$ ,  $\mathcal{I} \models S \cup \{\neg\alpha\}$ .  $S \cup \{\neg\alpha\}$  is unsatisfiable

Special case when  $S$  is empty:  $\models \alpha$  iff for every  $\mathcal{I}$ ,  $\mathcal{I} \models \alpha$ .

Say that  $\alpha$  is valid.

Note:  $\{\alpha_1, \alpha_2, \dots, \alpha_n\} \models \alpha$  iff  $\models (\alpha_1 \wedge \alpha_2 \wedge \dots \wedge \alpha_n) \supset \alpha$   
finite entailment reduces to validity

# Why do we care?

---

We do not have access to user-intended interpretation of non-logical symbols

But, with entailment, we know that if  $S$  is true in the intended interpretation, then so is  $\alpha$ .

If the user's view has the world satisfying  $S$ , then it must also satisfy  $\alpha$ .

There may be other sentences true also; but  $\alpha$  is logically guaranteed.

So what about ordinary reasoning?

Dog(fido)  $\rightsquigarrow$  Mammal(fido) ??

Not entailment!

There are logical interpretations where  $I[\text{Dog}] \not\subseteq I[\text{Mammal}]$

Key idea  
of KR:

include such connections explicitly in  $S$

$\forall x[\text{Dog}(x) \supset \text{Mammal}(x)]$

Get:  $S \cup \{\text{Dog}(\text{fido})\} \models \text{Mammal}(\text{fido})$

the rest is just  
details...

# Knowledge bases

---

KB is set of sentences

explicit statement of sentences believed (including any assumed connections among non-logical symbols)

$KB \models \alpha$      $\alpha$  is a further consequence of what is believed

- explicit knowledge: KB
- implicit knowledge:  $\{ \alpha \mid KB \models \alpha \}$

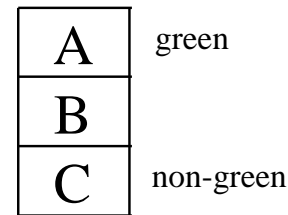
Often non trivial: explicit  $\implies$  implicit

Example:

Three blocks stacked.

Top one is green.

Bottom one is not green.



Is there a green block directly on top of a non-green block?



# A formalization

---

$$S = \{ \text{On}(a,b), \text{On}(b,c), \text{Green}(a), \neg\text{Green}(c) \}$$

all that is required

$$\alpha = \exists x \exists y [ \text{Green}(x) \wedge \neg\text{Green}(y) \wedge \text{On}(x,y) ]$$

Claim:  $S \models \alpha$

Proof:

Let  $\mathcal{I}$  be any interpretation such that  $\mathcal{I} \models S$ .

Case 1:  $\mathcal{I} \models \text{Green}(b)$ .

$\therefore \mathcal{I} \models \text{Green}(b) \wedge \neg\text{Green}(c) \wedge \text{On}(b,c)$ .

$\therefore \mathcal{I} \models \alpha$

Case 2:  $\mathcal{I} \not\models \text{Green}(b)$ .

$\therefore \mathcal{I} \models \neg\text{Green}(b)$

$\therefore \mathcal{I} \models \text{Green}(a) \wedge \neg\text{Green}(b) \wedge \text{On}(a,b)$ .

$\therefore \mathcal{I} \models \alpha$

Either way, for any  $\mathcal{I}$ , if  $\mathcal{I} \models S$  then  $\mathcal{I} \models \alpha$ .

So  $S \models \alpha$ .      QED

# Knowledge-based system

---

Start with (large) KB representing what is explicitly known

e.g. what the system has been told or has learned

Want to influence behaviour based on what is implicit in the KB  
(or as close as possible)

Requires reasoning

deductive inference:

process of calculating entailments of KB

i.e given KB and any  $\alpha$ , determine if  $KB \models \alpha$

Process is sound if whenever it produces  $\alpha$ , then  $KB \models \alpha$

does not allow for plausible assumptions that may be true  
in the intended interpretation

Process is complete if whenever  $KB \models \alpha$ , it produces  $\alpha$

does not allow for process to miss some  $\alpha$  or be unable to  
determine the status of  $\alpha$

---

3.

# Expressing Knowledge

# Knowledge engineering

---

KR is first and foremost about knowledge

meaning and entailment

find individuals and properties, then encode facts sufficient for entailments

Before implementing, need to understand clearly

- what is to be computed?
- why and where inference is necessary?

Example domain: soap-opera world

people, places, companies, marriages, divorces, hanky-panky, deaths, kidnappings, crimes, ...

Task: KB with appropriate entailments

- what vocabulary?
- what facts to represent?

# Vocabulary

---

## Domain-dependent predicates and functions

main question: what are the individuals?

here: people, places, companies, ...

## named individuals

john, sleezyTown, faultyInsuranceCorp, fic, johnQsmith, ...

## basic types

Person, Place, Man, Woman, ...

## attributes

Rich, Beautiful, Unscrupulous, ...

## relationships

LivesAt, MarriedTo, DaughterOf, HadAnAffairWith, Blackmails, ...

## functions

fatherOf, ceoOf, bestFriendOf, ...

# Basic facts

---

Usually atomic sentences and negations

type facts

Man(john),

Woman(jane),

Company(faultyInsuranceCorp)

property facts

Rich(john),

$\neg$ HappilyMarried(jim),

WorksFor(jim, fic)

equality facts

john = ceoOf(fic),

fic = faultyInsuranceCorp,

bestFriendOf(jim) = john

Like a simple database (can store in a table)

# Complex facts

---

## Universal abbreviations

$\forall y[\text{Woman}(y) \wedge y \neq \text{jane} \supset \text{Loves}(y, \text{john})]$

$\forall y[\text{Rich}(y) \wedge \text{Man}(y) \supset \text{Loves}(y, \text{jane})]$

$\forall x \forall y[\text{Loves}(x, y) \supset \neg \text{Blackmails}(x, y)]$

possible to express  
without quantifiers

## Incomplete knowledge

$\text{Loves}(\text{jane}, \text{john}) \vee \text{Loves}(\text{jane}, \text{jim})$

which?

$\exists x[\text{Adult}(x) \wedge \text{Blackmails}(x, \text{john})]$

who?

cannot write down  
a more complete  
version

## Closure axioms

$\forall x[\text{Person}(x) \supset x = \text{jane} \vee x = \text{john} \vee x = \text{jim} \dots]$

$\forall x \forall y[\text{MarriedTo}(x, y) \supset \dots]$

$\forall x[ x = \text{fic} \vee x = \text{jane} \vee x = \text{john} \vee x = \text{jim} \dots]$

limit the domain  
of discourse

also useful to have  $\text{jane} \neq \text{john} \dots$

# Terminological facts

---

General relationships among predicates. For example:

disjoint  $\forall x[\text{Man}(x) \supset \neg \text{Woman}(x)]$

subtype  $\forall x[\text{Senator}(x) \supset \text{Legislator}(x)]$

exhaustive  $\forall x[\text{Adult}(x) \supset \text{Man}(x) \vee \text{Woman}(x)]$

symmetry  $\forall x \forall y [\text{MarriedTo}(x,y) \supset \text{MarriedTo}(y,x)]$

inverse  $\forall x \forall y [\text{ChildOf}(x,y) \supset \text{ParentOf}(y,x)]$

type restriction  $\forall x \forall y [\text{MarriedTo}(x,y) \supset$   
 $\text{Person}(x) \wedge \text{Person}(y) \wedge \text{OppSex}(x,y)]$

sometimes

Usually universally quantified conditionals or biconditionals



# Entailments: 1

---

Is there a company whose CEO loves Jane?

$\exists x [\text{Company}(x) \wedge \text{Loves}(\text{ceoOf}(x), \text{jane})]$  ??

Suppose  $\mathcal{S} \models \text{KB}$ .

Then  $\mathcal{S} \models \text{Rich}(\text{john}), \text{Man}(\text{john})$ ,

and  $\mathcal{S} \models \forall y [\text{Rich}(y) \wedge \text{Man}(y) \supset \text{Loves}(y, \text{jane})]$

so  $\mathcal{S} \models \text{Loves}(\text{john}, \text{jane})$ .

Also  $\mathcal{S} \models \text{john} = \text{ceoOf}(\text{fic})$ ,

so  $\mathcal{S} \models \text{Loves}(\text{ceoOf}(\text{fic}), \text{jane})$ .

Finally  $\mathcal{S} \models \text{Company}(\text{faultyInsuranceCorp})$ ,

and  $\mathcal{S} \models \text{fic} = \text{faultyInsuranceCorp}$ ,

so  $\mathcal{S} \models \text{Company}(\text{fic})$ .

Thus,  $\mathcal{S} \models \text{Company}(\text{fic}) \wedge \text{Loves}(\text{ceoOf}(\text{fic}), \text{jane})$ ,

and so

$\mathcal{S} \models \exists x [\text{Company}(x) \wedge \text{Loves}(\text{ceoOf}(x), \text{jane})]$ .

Can extract identity of company from this proof

## Entailments: 2

---

If no man is blackmailing John, then is he being blackmailed by somebody he loves?

$$\forall x[\text{Man}(x) \supset \neg \text{Blackmails}(x, \text{john})] \supset \\ \exists y[\text{Loves}(\text{john}, y) \wedge \text{Blackmails}(y, \text{john})] \quad ??$$

Note:  $\text{KB} \models (\alpha \supset \beta)$  iff  $\text{KB} \cup \{\alpha\} \models \beta$

Let:  $\mathcal{S} \models \text{KB} \cup \{\forall x[\text{Man}(x) \supset \neg \text{Blackmails}(x, \text{john})]\}$

Show:  $\mathcal{S} \models \exists y[\text{Loves}(\text{john}, y) \wedge \text{Blackmails}(y, \text{john})]$

Have:  $\exists x[\text{Adult}(x) \wedge \text{Blackmails}(x, \text{john})]$  and  $\forall x[\text{Adult}(x) \supset \text{Man}(x) \vee \text{Woman}(x)]$   
so  $\exists x[\text{Woman}(x) \wedge \text{Blackmails}(x, \text{john})]$ .

Then:  $\forall y[\text{Rich}(y) \wedge \text{Man}(y) \supset \text{Loves}(y, \text{jane})]$  and  $\text{Rich}(\text{john}) \wedge \text{Man}(\text{john})$   
so  $\text{Loves}(\text{john}, \text{jane})!$

But:  $\forall y[\text{Woman}(y) \wedge y \neq \text{jane} \supset \text{Loves}(y, \text{john})]$   
and  $\forall x \forall y[\text{Loves}(x, y) \supset \neg \text{Blackmails}(x, y)]$   
so  $\forall y[\text{Woman}(y) \wedge y \neq \text{jane} \supset \neg \text{Blackmails}(y, \text{john})]$  and  $\text{Blackmails}(\text{jane}, \text{john})!!$

Finally:  $\text{Loves}(\text{john}, \text{jane}) \wedge \text{Blackmails}(\text{jane}, \text{john})$   
so:  $\exists y[\text{Loves}(\text{john}, y) \wedge \text{Blackmails}(y, \text{john})]$

# What individuals?

---

Sometimes useful to reduce n-ary predicates to 1-place predicates and 1-place functions

- involves reifying properties: new individuals
- typical of description logics / frame languages (later)

Flexibility in terms of arity:

Purchases(john,sears,bike) or  
Purchases(john,sears,bike,feb14) or  
Purchases(john,sears,bike,feb14,\$100)

Instead: introduce purchase objects

$\text{Purchase}(p) \wedge \text{agent}(p)=\text{john} \wedge \text{obj}(p)=\text{bike} \wedge \text{source}(p)=\text{sears} \wedge \dots$   
allows purchase to be described at various levels of detail

Complex relationships: MarriedTo( $x,y$ ) VS. ReMarriedTo( $x,y$ ) VS. ...

Instead define marital status in terms of existence of marriage and divorce events.

$\text{Marriage}(m) \wedge \text{husband}(m)=x \wedge \text{wife}(m)=y \wedge \text{date}(m)=\dots \wedge \dots$

# Abstract individuals

---

Also need individuals for numbers, dates, times, addresses, etc.

objects about which we ask wh-questions

## Quantities as individuals

$\text{age}(\text{suzy}) = 14$

$\text{age-in-years}(\text{suzy}) = 14$

$\text{age-in-months}(\text{suzy}) = 168$

perhaps better to have an object for “the age of Suzy”, whose value in years is 14

$\text{years}(\text{age}(\text{suzy})) = 14$

$\text{months}(x) = 12 * \text{years}(x)$

$\text{centimeters}(x) = 100 * \text{meters}(x)$

## Similarly with locations and times

instead of

$\text{time}(m) = \text{"Jan 5 2006 4:47:03EST"}$

can use

$\text{time}(m) = t \wedge \text{year}(t) = 2006 \wedge \dots$

# Other sorts of facts

---

## Statistical / probabilistic facts

- Half of the companies are located on the East Side.
- Most of the employees are restless.
- Almost none of the employees are completely trustworthy,

## Default / prototypical facts

- Company presidents typically have secretaries intercepting their phone calls.
- Cars have four wheels.
- Companies generally do not allow employees that work together to be married.

## Intentional facts

- John believes that Henry is trying to blackmail him.
- Jane does not want Jim to think that she loves John.

## Others ...

---

4.

# Resolution

# Goal

---

Deductive reasoning in language as close as possible to full FOL

$\neg, \wedge, \vee, \exists, \forall$

Knowledge Level:

given KB,  $\alpha$ , determine if  $\text{KB} \models \alpha$ .

or given an open  $\alpha[x_1, x_2, \dots, x_n]$ , find  $t_1, t_2, \dots, t_n$  such that  $\text{KB} \models \alpha[t_1, t_2, \dots, t_n]$

When KB is finite  $\{\alpha_1, \alpha_2, \dots, \alpha_k\}$

$\text{KB} \models \alpha$

iff  $\models [(\alpha_1 \wedge \alpha_2 \wedge \dots \wedge \alpha_k) \supset \alpha]$

iff  $\text{KB} \cup \{\neg\alpha\}$  is unsatisfiable

iff  $\text{KB} \cup \{\neg\alpha\} \models \text{FALSE}$

where FALSE is something like  $\exists x.(x \neq x)$

So want a procedure to test for validity, or satisfiability, or for entailing FALSE.

Will now consider such a procedure (first without quantifiers)

# Clausal representation

---

Formula = set of clauses

Clause = set of literals

Literal = atomic sentence or its negation

positive literal and negative literal

Notation:

If  $p$  is a literal, then  $\bar{p}$  is its complement

$$\bar{\bar{p}} \Rightarrow p \quad \overline{\neg p} \Rightarrow p$$

To distinguish clauses from formulas:

[ and ] for clauses:  $[p, \bar{r}, s]$     { and } for formulas:  $\{ [p, \bar{r}, s], [p, r, s], [\bar{p}] \}$

$[\ ]$  is the empty clause

$\{ \}$  is the empty formula

So  $\{ \}$  is different from  $\{ [\ ] \}$ !

Interpretation:

Formula understood as conjunction of clauses

$\{ [p, \neg q], [r], [s] \}$

$[ \ ]$

Clause understood as disjunction of literals

represents

represents

Literals understood normally

$((p \vee \neg q) \wedge r \wedge s)$

FALSE



# CNF and DNF

---

Every propositional wff  $\alpha$  can be converted into a formula  $\alpha'$  in Conjunctive Normal Form (CNF) in such a way that  $\models \alpha \equiv \alpha'$ .

1. eliminate  $\supset$  and  $\equiv$  using  $(\alpha \supset \beta) \rightsquigarrow (\neg\alpha \vee \beta)$  etc.
2. push  $\neg$  inward using  $\neg(\alpha \wedge \beta) \rightsquigarrow (\neg\alpha \vee \neg\beta)$  etc.
3. distribute  $\vee$  over  $\wedge$  using  $((\alpha \wedge \beta) \vee \gamma) \rightsquigarrow ((\alpha \vee \gamma) \wedge (\beta \vee \gamma))$
4. collect terms using  $(\alpha \vee \alpha) \rightsquigarrow \alpha$  etc.

Result is a conjunction of disjunction of literals.

an analogous procedure produces DNF,  
a disjunction of conjunction of literals

We can identify CNF wffs with clausal formulas

$$(p \vee \neg q \vee r) \wedge (s \vee \neg r) \rightsquigarrow \{ [p, \neg q, r], [s, \neg r] \}$$

So: given a finite KB, to find out if  $\text{KB} \models \alpha$ , it will be sufficient to

1. put  $(\text{KB} \wedge \neg\alpha)$  into CNF, as above
2. determine the satisfiability of the clauses

# Resolution rule of inference

---

Given two clauses, infer a new clause:

From clause  $\{ p \} \cup C_1$ ,  
and  $\{ \neg p \} \cup C_2$ ,  
infer clause  $C_1 \cup C_2$ .

$C_1 \cup C_2$  is called a resolvent of input clauses with respect to  $p$ .

Example:

clauses  $[w, r, q]$  and  $[w, s, \neg r]$  have  $[w, q, s]$  as resolvent wrt  $r$ .

Special Case:

$[p]$  and  $[\neg p]$  resolve to  $[\ ]$  (the  $C_1$  and  $C_2$  are empty)

A derivation of a clause  $c$  from a set  $S$  of clauses is a sequence  $c_1, c_2, \dots, c_n$  of clauses, where  $c_n = c$ , and for each  $c_i$ , either

1.  $c_i \in S$ , or
2.  $c_i$  is a resolvent of two earlier clauses in the derivation

Write:  $S \rightarrow c$  if there is a derivation

# Rationale

---

Resolution is a symbol-level rule of inference, but has a connection to knowledge-level logical interpretations

Claim: Resolvent is entailed by input clauses.

Suppose  $\mathcal{I} \models (p \vee \alpha)$  and  $\mathcal{I} \models (\neg p \vee \beta)$

Case 1:  $\mathcal{I} \models p$

then  $\mathcal{I} \models \beta$ , so  $\mathcal{I} \models (\alpha \vee \beta)$ .

Case 2:  $\mathcal{I} \not\models p$

then  $\mathcal{I} \models \alpha$ , so  $\mathcal{I} \models (\alpha \vee \beta)$ .

Either way,  $\mathcal{I} \models (\alpha \vee \beta)$ .

So:  $\{(p \vee \alpha), (\neg p \vee \beta)\} \models (\alpha \vee \beta)$ .

Special case:

$[p]$  and  $[\neg p]$  resolve to  $[\ ]$ ,

so  $\{[p], [\neg p]\} \models \text{FALSE}$

that is:  $\{[p], [\neg p]\}$  is unsatisfiable

# Derivations and entailment

---

Can extend the previous argument to derivations:

If  $S \rightarrow c$  then  $S \models c$

Proof: by induction on the length of the derivation.

Show (by looking at the two cases) that  $S \models c_i$ .

But the converse does not hold in general

Can have  $S \models c$  without having  $S \rightarrow c$ .

Example:  $\{\neg p\} \models [\neg p, \neg q]$  i.e.  $\neg p \models (\neg p \vee \neg q)$   
but no derivation

However.... Resolution *is* refutation complete!

**Theorem:**  $S \rightarrow []$  iff  $S \models []$

sound and complete  
when restricted to  $[]$

Result will carry over to quantified clauses (later)

So for any set  $S$  of clauses:  $S$  is unsatisfiable iff  $S \rightarrow []$ .

Provides method for determining satisfiability: search all derivations for  $[]$ .

So provides a method for determining all entailments

# A procedure for entailment

---

To determine if  $KB \models \alpha$ ,

- put  $KB, \neg\alpha$  into CNF to get  $S$ , as before
- check if  $S \rightarrow []$ .

If  $KB = \{\}$ , then we are testing the validity of  $\alpha$

## Non-deterministic procedure

1. Check if  $[]$  is in  $S$ .  
If yes, then return **UNSATISFIABLE**
2. Check if there are two clauses in  $S$  such that they resolve to produce a clause that is not already in  $S$ .  
If no, then return **SATISFIABLE**
3. Add the new clause to  $S$  and go to 1.

Note: need only convert  $KB$  to CNF once

- can handle multiple queries with same  $KB$
- after addition of new fact  $\alpha$ , can simply add new clauses  $\alpha'$  to  $KB$

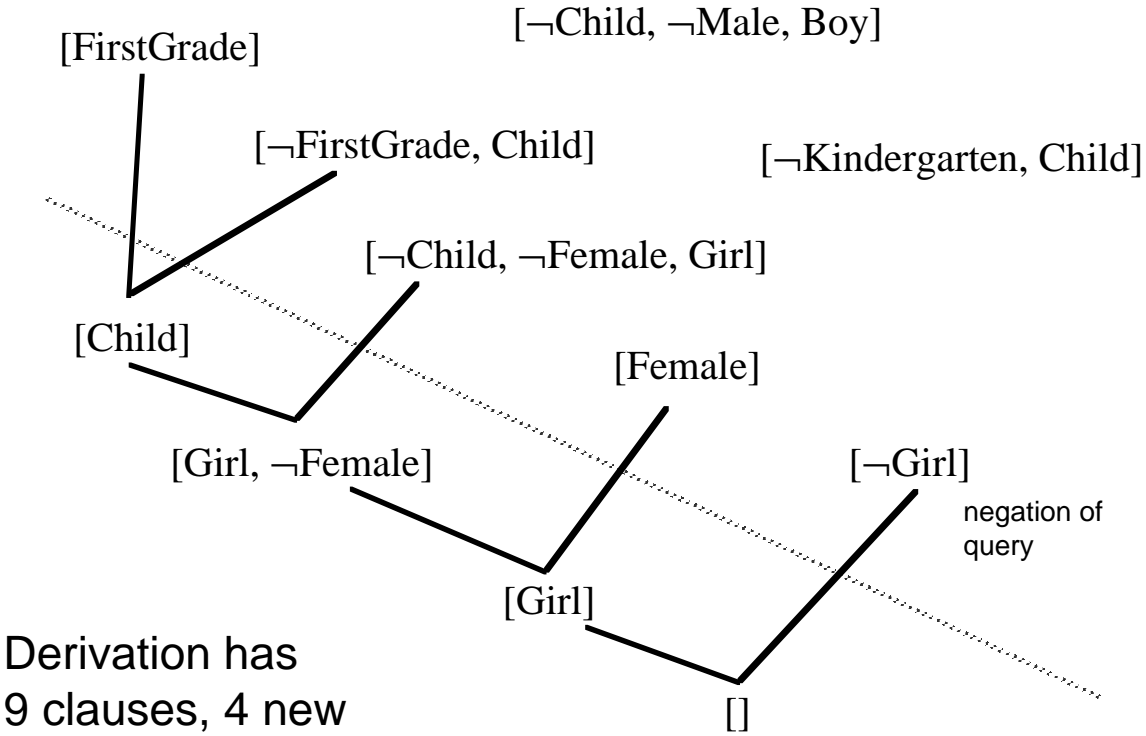
So: good idea to keep  $KB$  in CNF

# Example 1

KB

- FirstGrade
- FirstGrade  $\supset$  Child
- Child  $\wedge$  Male  $\supset$  Boy
- Kindergarten  $\supset$  Child
- Child  $\wedge$  Female  $\supset$  Girl
- Female

Show that KB  $\models$  Girl



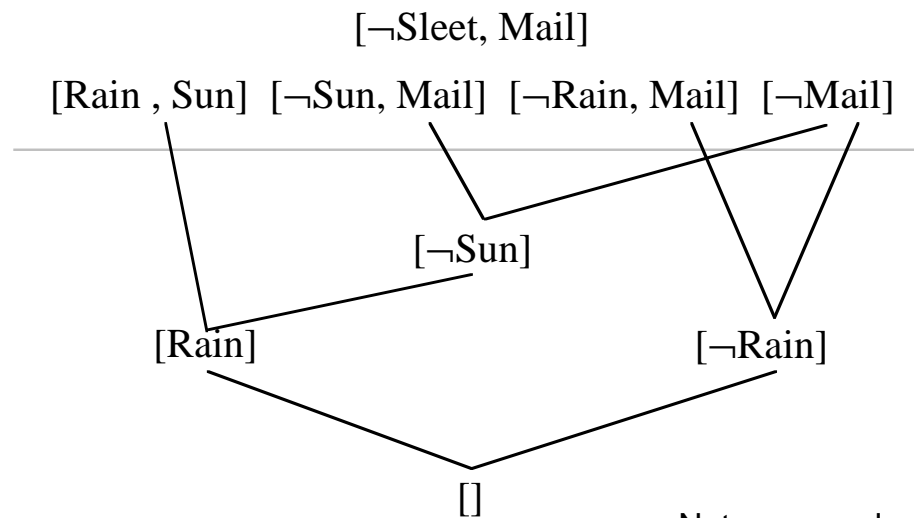
Derivation has  
9 clauses, 4 new

# Example 2

KB

$(\text{Rain} \vee \text{Sun})$   
 $(\text{Sun} \supset \text{Mail})$   
 $((\text{Rain} \vee \text{Sleet}) \supset \text{Mail})$

Show  $\text{KB} \models \text{Mail}$



Note: every clause not in  $S$  has 2 parents

Similarly  $\text{KB} \not\models \text{Rain}$

Can enumerate all resolvents given  $\neg\text{Rain}$ , and  $[\ ]$  will not be generated

# Quantifiers

---

Clausal form as before, but atom is  $P(t_1, t_2, \dots, t_n)$ , where  $t_i$  may contain variables

Interpretation as before, but variables are understood *universally*

Example:  $\{ [P(x), \neg R(a, f(b, x))], [Q(x, y)] \}$

interpreted as

$\forall x \forall y \{ [R(a, f(b, x)) \supset P(x)] \wedge Q(x, y) \}$

Substitutions:  $\theta = \{v_1/t_1, v_2/t_2, \dots, v_n/t_n\}$

Notation: If  $\rho$  is a literal and  $\theta$  is a substitution, then  $\rho\theta$  is the result of the substitution (and similarly,  $c\theta$  where  $c$  is a clause)

Example:  $\theta = \{x/a, y/g(x, b, z)\}$

$P(x, z, f(x, y))\theta = P(a, z, f(a, g(x, b, z)))$

A literal is ground if it contains no variables.

A literal  $\rho$  is an instance of  $\rho'$ , if for some  $\theta$ ,  $\rho = \rho'\theta$ .



# Generalizing CNF

---

Resolution will generalize to handling variables

Ignore = for now

But to convert wffs to CNF, we need three additional steps:

1. eliminate  $\supset$  and  $\equiv$

2. push  $\neg$  inward using also  $\neg\forall x.\alpha \rightsquigarrow \exists x.\neg\alpha$  etc.

3. standardize variables: each quantifier gets its own variable

e.g.  $\exists x[P(x)] \wedge Q(x) \rightsquigarrow \exists z[P(z)] \wedge Q(x)$  where  $z$  is a new variable

4. eliminate all existentials (*discussed later*)

5. move universals to the front using  $(\forall x\alpha) \wedge \beta \rightsquigarrow \forall x(\alpha \wedge \beta)$

where  $\beta$  does not use  $x$

6. distribute  $\vee$  over  $\wedge$

7. collect terms

Get universally quantified conjunction of disjunction of literals

then drop all the quantifiers...

# First-order resolution

---

Main idea: a literal (with variables) stands for all its instances; so allow all such inferences

So given  $[P(x,a), \neg Q(x)]$  and  $[\neg P(b,y), \neg R(b,f(y))]$ ,  
want to infer  $[\neg Q(b), \neg R(b,f(a))]$  among others

since  $[P(x,a), \neg Q(x)]$  has  $[P(b,a), \neg Q(b)]$  and  
 $[\neg P(b,y), \neg R(b,f(y))]$  has  $[\neg P(b,a), \neg R(b,f(a))]$

Resolution:

Given clauses:  $\{\rho_1\} \cup C_1$  and  $\{\bar{\rho}_2\} \cup C_2$ .

Rename variables, so that distinct in two clauses.

For any  $\theta$  such that  $\rho_1\theta = \bar{\rho}_2\theta$ , can infer  $(C_1 \cup C_2)\theta$ .

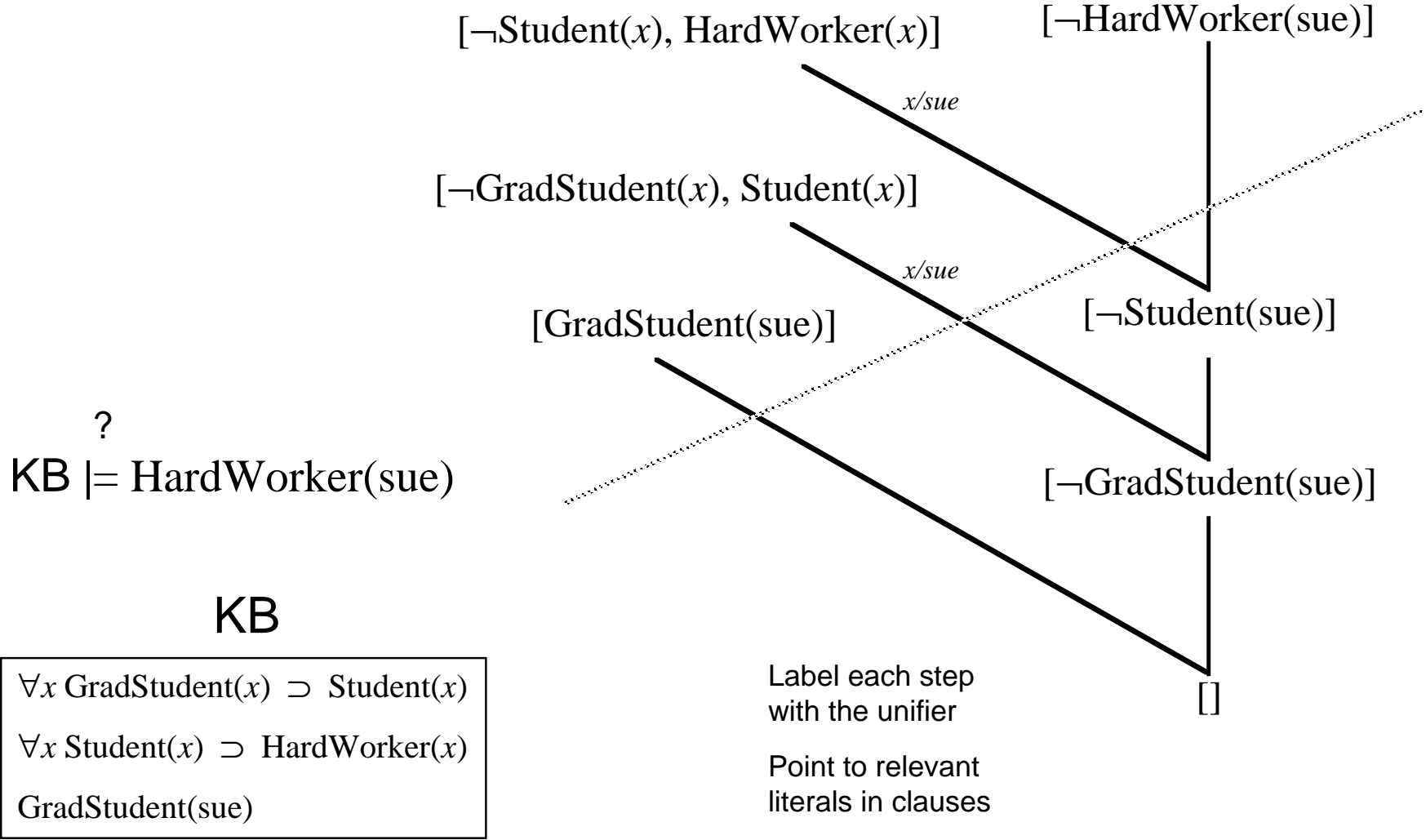
We say that  $\rho_1$  unifies with  $\bar{\rho}_2$  and that  $\theta$  is a unifier of the two literals

Resolution derivation: as before

**Theorem:**  $S \rightarrow []$  iff  $S \models []$  iff  $S$  is unsatisfiable

Note: There are pathological examples where a slightly more general definition of Resolution is required. We ignore them for now...

# Example 3



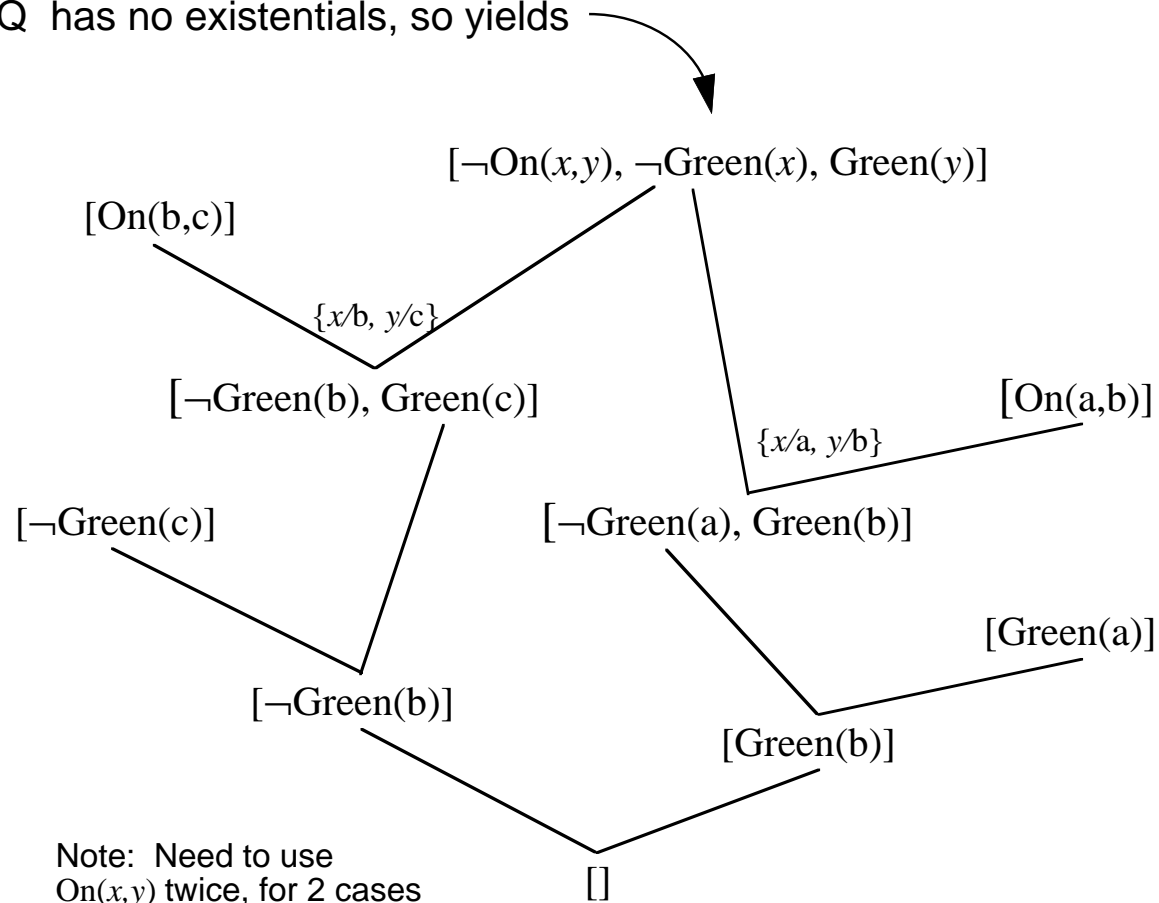
# The 3 block example

KB = {On(a,b), On(b,c), Green(a),  $\neg$ Green(c)}

already in CNF

Query =  $\exists x \exists y [\text{On}(x,y) \wedge \text{Green}(x) \wedge \neg \text{Green}(y)]$

Note:  $\neg Q$  has no existentials, so yields



# Arithmetic

KB:  $\text{Plus}(\text{zero}, x, x)$   
 $\text{Plus}(x, y, z) \supset \text{Plus}(\text{succ}(x), y, \text{succ}(z))$

Q:  $\exists u \text{Plus}(2, 3, u)$

For readability,  
we use

- 0 for zero,
- 1 for succ(zero),
- 2 for succ(succ(zero))

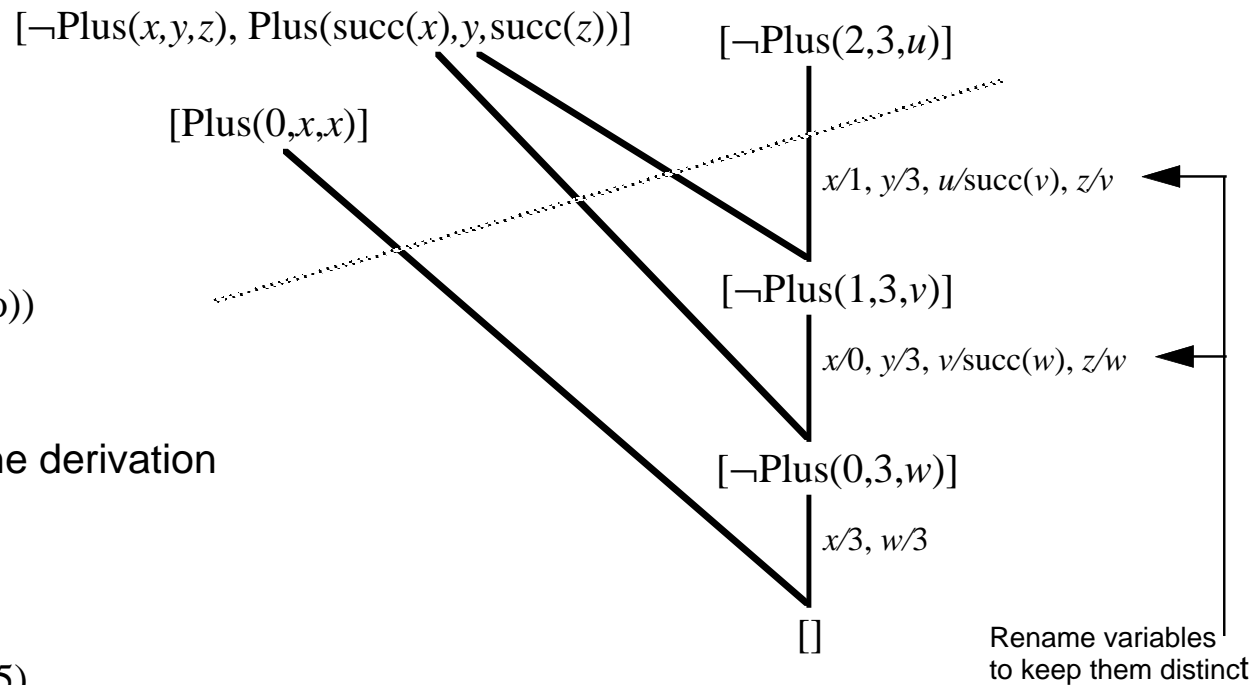
etc.

Can find the answer in the derivation

$u/\text{succ}(\text{succ}(3))$

that is:  $u/5$

Can also derive  $\text{Plus}(2, 3, 5)$



# Answer predicates

In full FOL, we have the possibility of deriving  $\exists xP(x)$  without being able to derive  $P(t)$  for any  $t$ .

e.g. the three-blocks problem

$$\exists x\exists y[\text{On}(x,y) \wedge \text{Green}(x) \wedge \neg\text{Green}(y)]$$

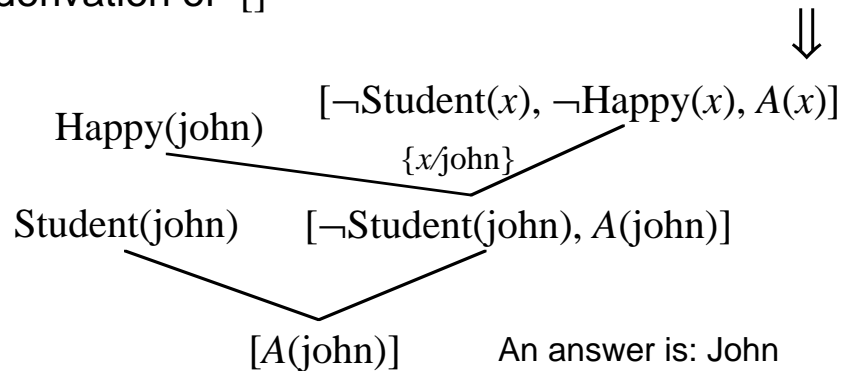
but cannot derive which block is which

## Solution: answer-extraction process

- replace query  $\exists xP(x)$  by  $\exists x[P(x) \wedge \neg A(x)]$   
where  $A$  is a new predicate symbol called the answer predicate
- instead of deriving  $\square$ , derive any clause containing just the answer predicate
- can always convert to and from a derivation of  $\square$

KB: Student(john)  
Student(jane)  
Happy(john)

Q:  $\exists x[\text{Student}(x) \wedge \text{Happy}(x)]$



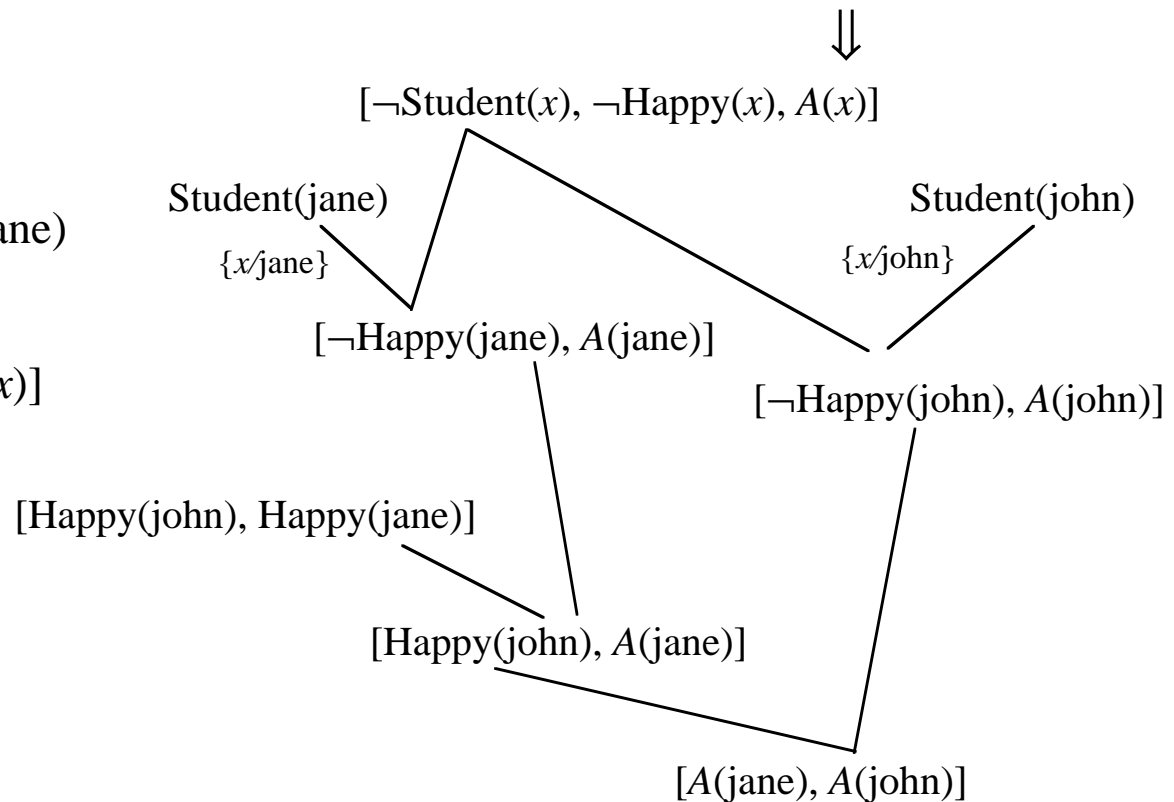
# Disjunctive answers

KB:

Student(john)  
 Student(jane)  
 Happy(john)  $\vee$  Happy(jane)

Query:

$\exists x[\text{Student}(x) \wedge \text{Happy}(x)]$



An answer is: either Jane or John

Note:

- can have variables in answer
- need to watch for Skolem symbols... (next)

# Skolemization

---

So far, converting wff to CNF ignored existentials

e.g.  $\exists x \forall y \exists z P(x, y, z)$

Idea: names for individuals claimed to exist, called Skolem constant and function symbols

there exists an  $x$ , call it  $a$

for each  $y$ , there is a  $z$ , call it  $f(y)$

get  $\forall y P(a, y, f(y))$

So replace  $\forall x_1 (\dots \forall x_2 (\dots \forall x_n (\dots \exists y [\dots y \dots] \dots) \dots) \dots)$   
by  $\forall x_1 (\dots \forall x_2 (\dots \forall x_n (\dots [\dots f(x_1, x_2, \dots, x_n) \dots] \dots) \dots) \dots)$

$f$  is a new function symbol that appears nowhere else

Skolemization does not preserve equivalence

e.g.  $\not\models \exists x P(x) \equiv P(a)$

But it does preserve satisfiability

$\alpha$  is satisfiable iff  $\alpha'$  is satisfiable (where  $\alpha'$  is the result of Skolemization)

sufficient for resolution!



# Variable dependence

---

Show that  $\exists x \forall y R(x,y) \models \forall y \exists x R(x,y)$

show  $\{\exists x \forall y R(x,y), \neg \forall y \exists x R(x,y)\}$  unsatisfiable

$$\exists x \forall y R(x,y) \rightsquigarrow \forall y R(a,y)$$

$$\neg \forall y \exists x R(x,y) \rightsquigarrow \exists y \forall x \neg R(x,y) \rightsquigarrow \forall x \neg R(x,b)$$

then  $\{ [R(a,y)], [\neg R(x,b)] \} \rightarrow []$  with  $\{x/a, y/b\}$ .

Show that  $\forall y \exists x R(x,y) \not\models \exists x \forall y R(x,y)$

show  $\{\forall y \exists x R(x,y), \neg \exists x \forall y R(x,y)\}$  satisfiable

$$\forall y \exists x R(x,y) \rightsquigarrow \forall y R(f(y),y)$$

$$\neg \exists x \forall y R(x,y) \rightsquigarrow \forall x \exists y \neg R(x,y) \rightsquigarrow \forall x \neg R(x,g(x))$$

then get  $\{ [R(f(y),y)], [\neg R(x,g(x))] \}$

where the two literals do not unify

Note: important to get dependence of variables correct

$R(f(y),y)$  vs.  $R(a,y)$  in the above

# A problem

---

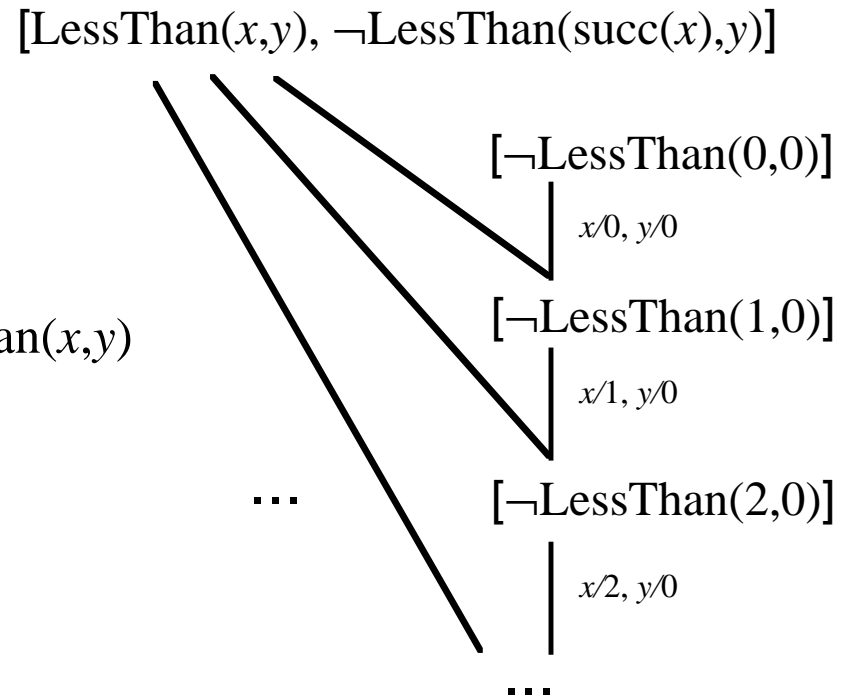
KB:

$\text{LessThan}(\text{succ}(x),y) \supset \text{LessThan}(x,y)$

Query:

$\text{LessThan}(\text{zero},\text{zero})$

Should fail since  $\text{KB} \not\models Q$



Infinite branch of resolvents

cannot use a simple depth-first procedure to search for []

# Undecidability

---

Is there a way to detect when this happens?

No! FOL is very powerful

can be used as a full programming language

just as there is no way to detect in general when  
a program is looping

There can be no procedure that does this:

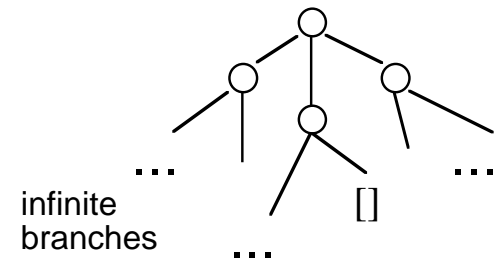
```
Proc[Clauses] =  
  If Clauses are unsatisfiable  
  then return YES  
  else return NO
```

However: Resolution is complete

some branch will contain [], for unsatisfiable clauses

So breadth-first search guaranteed to find []

search may not terminate on satisfiable clauses



# Overly specific unifiers

---

In general, no way to guarantee efficiency, or even termination

later: put control into users' hands

One thing that can be done:

reduce redundancy in search, by keeping search as general as possible

Example

$\dots, P(g(x), f(x), z) \quad [\neg P(y, f(w), a), \dots$

unified by

$\theta_1 = \{x/b, y/g(b), z/a, w/b\}$  gives  $P(g(b), f(b), a)$

and by

$\theta_2 = \{x/f(z), y/g(f(z)), z/a, w/f(z)\}$  gives  $P(g(f(z)), f(f(z)), a)$ .

Might not be able to derive the empty clause from clauses having overly specific substitutions

wastes time in search!

# Most general unifiers

---

$\theta$  is a most general unifier (MGU) of literals  $\rho_1$  and  $\rho_2$  iff

1.  $\theta$  unifies  $\rho_1$  and  $\rho_2$
2. for any other unifier  $\theta'$ , there is a another substitution  $\theta^*$  such that  $\theta' = \theta\theta^*$

Note: composition  $\theta\theta^*$  requires applying  $\theta^*$  to terms in  $\theta$

for previous example, an MGU is

$$\theta = \{x/w, y/g(w), z/a\}$$

for which

$$\theta_1 = \theta\{w/b\}$$

$$\theta_2 = \theta\{w/f(z)\}$$

**Theorem:** Can limit search to most general unifiers only without loss of completeness (with certain caveats)

# Computing MGUs

---

Computing an MGU, given a set of literals  $\{\rho_i\}$

usually only have two literals

1. Start with  $\theta := \{\}$ .
2. If all the  $\rho_i\theta$  are identical, then done;  
otherwise, get disagreement set,  $DS$   
e.g  $P(a, f(a, g(z)), \dots) P(a, f(a, u), \dots)$   
disagreement set,  $DS = \{u, g(z)\}$
3. Find a variable  $v \in DS$ , and a term  $t \in DS$  not containing  $v$ .  
If not, fail.
4.  $\theta := \theta\{v/t\}$
5. Go to 2

Note: there is a better *linear* algorithm

# Herbrand Theorem

---

Some 1st-order cases can be handled by converting them to a propositional form

Given a set of clauses  $S$

- the Herbrand universe of  $S$  is the set of all terms formed using only the function symbols in  $S$  (at least one)

e.g., if  $S$  uses (unary)  $f$ , and  $c, d$ ,  $U = \{c, d, f(c), f(d), f(f(c)), f(f(d)), f(f(f(c))), \dots\}$

- the Herbrand base of  $S$  is the set of all  $c\theta$  such that  $c \in S$  and  $\theta$  replaces the variables in  $c$  by terms from the Herbrand universe

Theorem:  $S$  is satisfiable iff Herbrand base is

(applies to Horn clauses also)

Herbrand base has no variables, and so is essentially *propositional*, though usually infinite

- finite, when Herbrand universe is finite  
can use propositional methods (guaranteed to terminate)
- sometimes other “type” restrictions can be used to keep the Herbrand base finite  
include  $f(t)$  only if  $t$  is the correct type

# Resolution is difficult!

---

First-order resolution is not guaranteed to terminate.

What can be said about the propositional case?

Shown by Haken in 1985 that there are unsatisfiable clauses  $\{c_1, c_2, \dots, c_n\}$  such that the *shortest* derivation of [] contains on the order of  $2^n$  clauses

Even if we could always find a derivation immediately, the most clever search procedure will still require *exponential* time on some problems

Problem just with resolution?

Probably not.

Determining if a set of clauses is satisfiable was shown by Cook in 1972 to be NP-complete

No easier than an extremely large variety of computational tasks

Roughly: any search task where what is searched for can be verified in polynomial time can be recast as a satisfiability problem

- » satisfiability
- » does graph of cities allow for a full tour of size  $\leq k$  miles?
- » can N queens be put on an N×N chessboard all safely?      and many, many more....

Satisfiability is believed by most people to be unsolvable in polynomial time



# SAT solvers

---

In the propositional case, procedures have been proposed for determining the satisfiability of a set of clauses that appear to work much better in practice than Resolution.

The most popular is called DP (or DPLL) based on ideas by Davis, Putnam, Loveland and Logemann. (See book for details.)

These procedures are called SAT solvers as they are mostly used to find a satisfying interpretation for clauses that are satisfiable.

related to constraint satisfaction programs (CSP)

Typically they have the property that if they *fail* to find a satisfying interpretation, a Resolution derivation of [ ] can be reconstructed from a trace of their execution.

so worst-case exponential behaviour, via Haken's theorem!

One interesting counter-example to this is the procedure GSAT, which has different limitations. (Again, see the book.)

# Implications for KR

---

Problem: want to produce entailments of KB as needed for immediate action

full theorem-proving may be too difficult for KR!

need to consider other options ...

- giving control to user e.g. procedural representations (later)
- less expressive languages e.g. Horn clauses (and a major theme later)

In some applications, it is reasonable to wait

e.g. mathematical theorem proving, where we care about specific formulas

Best to hope for in general: reduce redundancy

main example: MGU, as before

but many other strategies (as we will see)

ATP: automated theorem proving

- area of AI that studies strategies for automatically proving difficult theorems
- main application: mathematics, but relevance also to KR

# Strategies

---

## 1. Clause elimination

- pure clause

contains literal  $p$  such that  $p$  does not appear in any other clause  
clause cannot lead to  $\square$

- tautology

clause with a literal and its negation  
any path to  $\square$  can bypass tautology

- subsumed clause

a clause such that one with a subset of its literals is already present  
path to  $\square$  need only pass through short clause  
can be generalized to allow substitutions

## 2. Ordering strategies

many possible ways to order search, but best and simplest is

- unit preference

prefer to resolve unit clauses first

Why? Given unit clause and another clause, resolvent is a smaller one  $\Rightarrow \square$

# Strategies 2

---

## 3. Set of support

KB is usually satisfiable, so not very useful to resolve among clauses with only ancestors in KB

contradiction arises from interaction with  $\neg Q$

always resolve with at least one clause that has an ancestor in  $\neg Q$   
preserves completeness (sometimes)

## 4. Connection graph

pre-compute all possible unifications

build a graph with edges between any two unifiable literals of opposite polarity

label edge with MGU

Resolution procedure:

repeatedly:     select link  
                  compute resolvent  
                  inherit links from parents after substitution

Resolution as search: find sequence of links  $L_1, L_2, \dots$  producing []

# Strategies 3

---

## 5. Special treatment for equality

instead of using axioms for =

reflexivity, symmetry, transitivity, substitution of equals for equals

use new inference rule: paramodulation

from  $\{(t=s)\} \cup C_1$  and  $\{P(\dots t' \dots)\} \cup C_2$

where  $t\theta = t'\theta$

infer  $\{P(\dots s \dots)\}\theta \cup C_1\theta \cup C_2\theta$ .

collapses many resolution steps into one

see also: theory resolution (later)

## 6. Sorted logic

terms get sorts:

$x$ : Male    mother:[Person  $\rightarrow$  Female]

keep taxonomy of sorts

only unify  $P(s)$  with  $P(t)$  when sorts are compatible

assumes only “meaningful” paths will lead to []

# Finally...

---

## 7. Directional connectives

given  $[\neg p, q]$ , can interpret as either

from  $p$ , infer  $q$  (forward)

to prove  $q$ , prove  $p$  (backward)

procedural reading of  $\supset$

In 1st case: would only resolve  $[\neg p, q]$  with  $[p, \dots]$  producing  $[q, \dots]$

In 2nd case: would only resolve  $[\neg p, q]$  with  $[\neg q, \dots]$  producing  $[\neg p, \dots]$

### Intended application:

forward:  $\text{Battleship}(x) \supset \text{Gray}(x)$

do not want to try to prove something is gray  
by trying to prove that it is a battleship

backward:  $\text{Person}(x) \supset \text{Has}(x, \text{spleen})$

do not want to conclude the spleen property for  
each individual inferred to be a person

This is the starting point for the procedural representations (later)

---

5.

# Reasoning with Horn Clauses

# Horn clauses

---

Clauses are used two ways:

- as disjunctions: (rain  $\vee$  sleet)
- as implications: ( $\neg$ child  $\vee$   $\neg$ male  $\vee$  boy)

Here focus on 2nd use

Horn clause = at most one +ve literal in clause

- positive / definite clause = exactly one +ve literal

e.g.  $[\neg p_1, \neg p_2, \dots, \neg p_n, q]$

- negative clause = no +ve literals

e.g.  $[\neg p_1, \neg p_2, \dots, \neg p_n]$  and also  $[\ ]$

Note:  $[\neg p_1, \neg p_2, \dots, \neg p_n, q]$  is a representation for  
 $(\neg p_1 \vee \neg p_2 \vee \dots \vee \neg p_n \vee q)$  or  $[(p_1 \wedge p_2 \wedge \dots \wedge p_n) \supset q]$

so can read as: If  $p_1$  and  $p_2$  and ... and  $p_n$  then  $q$

and write as:  $p_1 \wedge p_2 \wedge \dots \wedge p_n \Rightarrow q$  or  $q \Leftarrow p_1 \wedge p_2 \wedge \dots \wedge p_n$

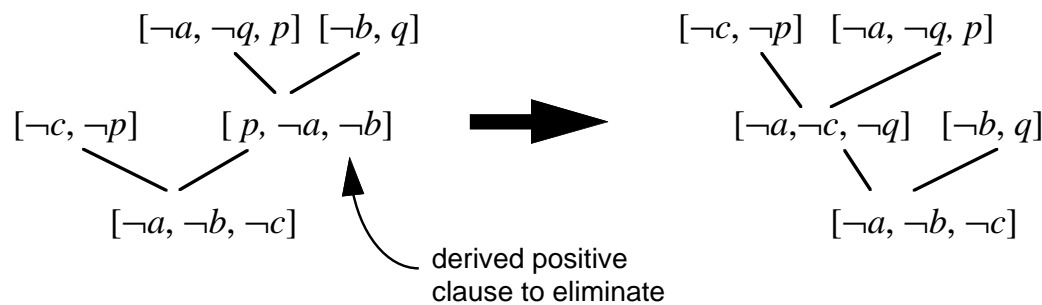


# Resolution with Horn clauses

Only two possibilities:



It is possible to rearrange derivations of negative clauses so that all new derived clauses are negative



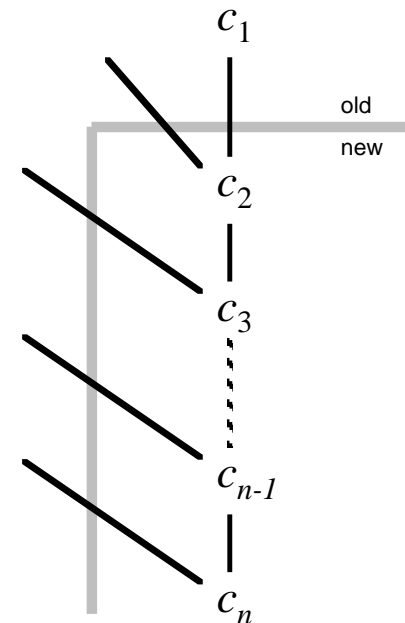
# Further restricting resolution

Can also change derivations such that each derived clause is a resolvent of the previous derived one (negative) and some positive clause in the original set of clauses

- Since each derived clause is negative, one parent must be positive (and so from original set) and one parent must be negative.
- Chain backwards from the final negative clause until both parents are from the original set of clauses
- Eliminate all other clauses not on this direct path

This is a recurring pattern in derivations

- See previously:
  - example 1, example 3, arithmetic example
- But not:
  - example 2, the 3 block example



# SLD Resolution

---

An SLD-derivation of a clause  $c$  from a set of clauses  $S$  is a sequence of clause  $c_1, c_2, \dots, c_n$  such that  $c_n = c$ , and

1.  $c_1 \in S$
2.  $c_{i+1}$  is a resolvent of  $c_i$  and a clause in  $S$

Write:  $S \xrightarrow{\text{SLD}} c$

SLD means S(elected) literals  
L(inear) form  
D(efinite) clauses

Note: SLD derivation is just a special form of derivation and where we leave out the elements of  $S$  (except  $c_1$ )

In general, cannot restrict ourselves to just using SLD-Resolution

Proof:  $S = \{[p, q], [p, \neg q], [\neg p, q], [\neg p, \neg q]\}$ . Then  $S \rightarrow []$ .

Need to resolve some  $[\rho]$  and  $[\bar{\rho}]$  to get  $[\ ]$ .

But  $S$  does not contain any unit clauses.

So will need to derive both  $[\rho]$  and  $[\bar{\rho}]$  and then resolve them together.

# Completeness of SLD

---

However, for Horn clauses, we can restrict ourselves to SLD-Resolution

**Theorem:** SLD-Resolution is refutation complete for Horn clauses:  $H \rightarrow []$  iff  $H \xrightarrow{\text{SLD}} []$

So:  $H$  is unsatisfiable iff  $H \xrightarrow{\text{SLD}} []$

This will considerably simplify the search for derivations

**Note:** in Horn version of SLD-Resolution, each clause in the  $c_1, c_2, \dots, c_n$ , will be negative

So clauses  $H$  must contain at least one negative clause,  $c_1$  and this will be the only negative clause of  $H$  used.

Typical case:

- KB is a collection of positive Horn clauses
- Negation of query is the negative clause

# Example 1 (again)

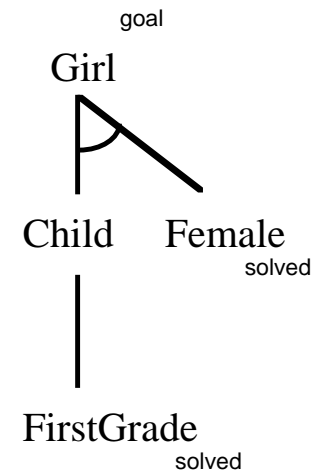
## KB

FirstGrade  
FirstGrade  $\supset$  Child  
Child  $\wedge$  Male  $\supset$  Boy  
Kindergarten  $\supset$  Child  
Child  $\wedge$  Female  $\supset$  Girl  
Female

## SLD derivation

[ $\neg$ Girl]  
|  
[ $\neg$ Child,  $\neg$ Female]  
|  
[ $\neg$ Child]  
|  
[ $\neg$ FirstGrade]  
|  
[]

## alternate representation



Show  $\text{KB} \cup \{\neg\text{Girl}\}$  unsatisfiable

A goal tree whose nodes are atoms, whose root is the atom to prove, and whose leaves are in the KB

# Prolog

---

## Horn clauses form the basis of Prolog

Append(nil,y,y)

Append(x,y,z)  $\Rightarrow$  Append(cons(w,x),y,cons(w,z))

With SLD derivation, can  
always extract answer from proof

$H \models \exists x \alpha(x)$

iff

for some term  $t$ ,  $H \models \alpha(t)$

Different answers can be found  
by finding other derivations

What is the result of appending [c] to the list [a,b] ?

Append(cons(a,cons(b,nil)), cons(c,nil),  $u$ )      goal

|  $u / \text{cons}(a,u')$

Append(cons(b,nil), cons(c,nil),  $u'$ )

|  $u' / \text{cons}(b,u'')$

Append(nil, cons(c,nil),  $u''$ )

solved:  $u'' / \text{cons}(c,\text{nil})$

So goal succeeds with  $u = \text{cons}(a,\text{cons}(b,\text{cons}(c,\text{nil})))$   
that is: Append([a b],[c],[a b c])

# Back-chaining procedure

---

```
Solve[ $q_1, q_2, \dots, q_n$ ] = /* to establish conjunction of  $q_i$  */
  If  $n=0$  then return YES; /* empty clause detected */
  For each  $d \in \text{KB}$  do
    If  $d = [q_1, \neg p_1, \neg p_2, \dots, \neg p_m]$  /* match first  $q$  */
      and /* replace  $q$  by -ve lits */
      Solve[ $p_1, p_2, \dots, p_m, q_2, \dots, q_n$ ] /* recursively */
    then return YES
  end for; /* can't find a clause to eliminate  $q$  */
  Return NO
```

## Depth-first, left-right, back-chaining

- depth-first because attempt  $p_i$  before trying  $q_i$
- left-right because try  $q_i$  in order, 1,2, 3, ...
- back-chaining because search from goal  $q$  to facts in KB  $p$

## This is the execution strategy of Prolog

First-order case requires unification *etc.*

# Problems with back-chaining

Can go into infinite loop

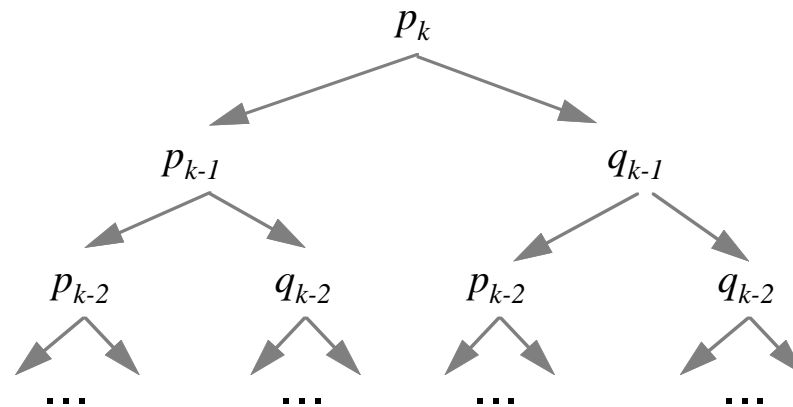
tautologous clause:  $[p, \neg p]$  (corresponds to Prolog program with  $p :- p$ ).

Previous back-chaining algorithm is inefficient

Example: Consider  $2n$  atoms,  $p_0, \dots, p_{n-1}, q_0, \dots, q_{n-1}$  and  $4n-4$  clauses

$(p_{i-1} \Rightarrow p_i), (q_{i-1} \Rightarrow p_i), (p_{i-1} \Rightarrow q_i), (q_{i-1} \Rightarrow q_i)$ .

With goal  $p_k$  the execution tree is like this



Solve $[p_k]$  eventually fails after  $2^k$  steps!

Is this problem inherent in Horn clauses?



# Forward-chaining

---

Simple procedure to determine if Horn KB  $\models q$ .

main idea: mark atoms as solved

1. If  $q$  is marked as solved, then return **YES**
2. Is there a  $\{p_1, \neg p_2, \dots, \neg p_n\} \in \text{KB}$  such that  $p_2, \dots, p_n$  are marked as solved, but the positive lit  $p_1$  is not marked as solved?
  - no: return **NO**
  - yes: mark  $p_1$  as solved, and go to 1.

FirstGrade example:

Marks: FirstGrade, Child, Female, Girl then done!

Note: FirstGrade gets marked since all the negative atoms in the clause (none) are marked

Observe:

- only letters in KB can be marked, so at most a linear number of iterations
- not goal-directed, so not always desirable
- a similar procedure with better data structures will run in *linear* time overall

# First-order undecidability

---

Even with just Horn clauses, in the first-order case we still have the possibility of generating an infinite branch of resolvents.

KB:

$\text{LessThan}(\text{succ}(x),y) \Rightarrow \text{LessThan}(x,y)$

Query:

$\text{LessThan}(\text{zero},\text{zero})$

As with full Resolution,  
there is no way to detect  
when this will happen

There is no procedure that will test for the  
satisfiability of first-order Horn clauses

the question is undecidable

$[\neg\text{LessThan}(0,0)]$

$\downarrow_{x/0, y/0}$

$[\neg\text{LessThan}(1,0)]$

$\downarrow_{x/1, y/0}$

$[\neg\text{LessThan}(2,0)]$

$\downarrow_{x/2, y/0}$

...

As with non-Horn clauses, the best that we can do is to give control of the deduction to the *user*

to some extent this is what is done in Prolog,  
but we will see more in “Procedural Control”

---

6.

# Procedural Control of Reasoning

# Declarative / procedural

---

Theorem proving (like resolution) is a general domain-independent method of reasoning

Does not require the user to know how knowledge will be used

will try all logically permissible uses

Sometimes we have ideas about how to use knowledge, how to search for derivations

do not want to use arbitrary or stupid order

Want to communicate to theorem-proving procedure some *guidance* based on properties of the domain

- perhaps specific method to use
- perhaps merely method to avoid

Example: directional connectives

In general: control of reasoning

# DB + rules

---

Can often separate (Horn) clauses into two components:

Example:

MotherOf(jane,billy)

FatherOf(john,billy)

FatherOf(sam, john)

...

ParentOf(x,y)  $\leftarrow$  MotherOf(x,y)

ParentOf(x,y)  $\leftarrow$  FatherOf(x,y)

ChildOf(x,y)  $\leftarrow$  ParentOf(y,x)

AncestorOf(x,y)  $\leftarrow$  ...

...

a database of facts

- basic facts of the domain
- usually ground atomic wffs

collection of rules

- extends the predicate vocabulary
- usually universally quantified conditionals

Both retrieved by unification matching

Control issue: how to use the rules

# Rule formulation

---

Consider AncestorOf in terms of ParentOf

Three logically equivalent versions:

1.  $\text{AncestorOf}(x,y) \Leftarrow \text{ParentOf}(x,y)$   
 $\text{AncestorOf}(x,y) \Leftarrow \text{ParentOf}(x,z) \wedge \text{AncestorOf}(z,y)$
2.  $\text{AncestorOf}(x,y) \Leftarrow \text{ParentOf}(x,y)$   
 $\text{AncestorOf}(x,y) \Leftarrow \text{ParentOf}(z,y) \wedge \text{AncestorOf}(x,z)$
3.  $\text{AncestorOf}(x,y) \Leftarrow \text{ParentOf}(x,y)$   
 $\text{AncestorOf}(x,y) \Leftarrow \text{AncestorOf}(x,z) \wedge \text{AncestorOf}(z,y)$

Back-chaining goal of  $\text{AncestorOf}(\text{sam},\text{sue})$  will ultimately reduce to set of  $\text{ParentOf}(-,-)$  goals

1. get  $\text{ParentOf}(\text{sam},z)$ : find child of Sam searching *downwards*
2. get  $\text{ParentOf}(z,\text{sue})$ : find parent of Sue searching *upwards*
3. get  $\text{ParentOf}(-,-)$ : find parent relations searching *in both directions*

Search strategies are not equivalent

if more than 2 children per parent, (2) is best

# Algorithm design

---

Example: Fibonacci numbers

1, 1, 2, 3, 5, 8, 13, 21, ...

Version 1:

Fibo(0, 1)

Fibo(1, 1)

$\text{Fibo}(s(s(n)), x) \Leftarrow \text{Fibo}(n, y) \wedge \text{Fibo}(s(n), z) \wedge \text{Plus}(y, z, x)$

Requires *exponential* number of Plus subgoals

Version 2:

$\text{Fibo}(n, x) \Leftarrow \text{F}(n, 1, 0, x)$

$\text{F}(0, c, p, c)$

$\text{F}(s(n), c, p, x) \Leftarrow \text{Plus}(p, c, s) \wedge \text{F}(n, s, c, x)$

Requires only *linear* number of Plus subgoals

# Ordering goals

---

Example:

$\text{AmericanCousinOf}(x,y) \Leftarrow \text{American}(x) \wedge \text{CousinOf}(x,y)$

In back-chaining, can try to solve either subgoal first

Not much difference for  $\text{AmericanCousinOf}(\text{fred}, \text{sally})$ , but big difference for  $\text{AmericanCousinOf}(x, \text{sally})$

1. find an American and then check to see if she is a cousin of Sally
2. find a cousin of Sally and then check to see if she is an American

So want to be able to order goals

better to generate cousins and test for American

In Prolog: order clauses, and literals in them

Notation:  $G :- G_1, G_2, \dots, G_n$  stands for

$$G \Leftarrow G_1 \wedge G_2 \wedge \dots \wedge G_n$$

but goals are attempted in presented order



# Commit

---

Need to allow for backtracking in goals

$\text{AmericanCousinOf}(x,y) \text{ :- CousinOf}(x,y), \text{American}(x)$

for goal  $\text{AmericanCousinOf}(x,\text{sally})$ , may need to try to solve the goal  $\text{American}(x)$  for many values of  $x$

But sometimes, given clause of the form

$G \text{ :- } T, S$

goal  $T$  is needed only as a *test* for the applicability of subgoal  $S$

- if  $T$  succeeds, commit to  $S$  as the *only* way of achieving goal  $G$ .
- if  $S$  fails, then  $G$  is considered to have failed
  - do not look for other ways of solving  $T$
  - do not look for other clauses with  $G$  as head

In Prolog: use of cut symbol

Notation:  $G \text{ :- } T_1, T_2, \dots, T_m, !, G_1, G_2, \dots, G_n$

attempt goals in order, but if all  $T_i$  succeed, then commit to  $G_i$

# If-then-else

---

Sometimes inconvenient to separate clauses in terms of unification:

$G(\text{zero}, -) :- \text{method 1}$   
 $G(\text{succ}(n), -) :- \text{method 2}$

For example, may split based on computed property:

$\text{Expt}(a, n, x) :- \text{Even}(n), \dots$  (*what to do when  $n$  is even*)  
 $\text{Expt}(a, n, x) :- \text{Even}(s(n)), \dots$  (*what to do when  $n$  is odd*)

want: check for even numbers only once

Solution: use ! to do if-then-else

$G :- P, !, Q.$   
 $G :- R.$

To achieve  $G$ : if  $P$  then use  $Q$  else use  $R$

Example:

$\text{Expt}(a, n, x) :- n = 0, !, x = 1.$   
 $\text{Expt}(a, n, x) :- \text{Even}(n), !, \text{ (for even } n)$   
 $\text{Expt}(a, n, x) :- \text{ (for odd } n)$

Note: it would be correct to write

$\text{Expt}(a, 0, x) :- !, x = 1.$

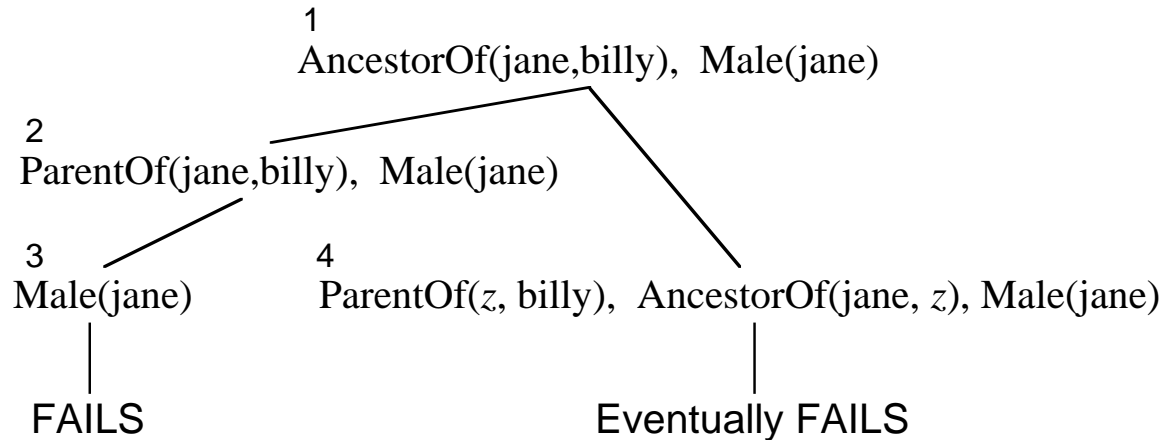
but not

$\text{Expt}(a, 0, 1) :- !.$

# Controlling backtracking

---

Consider solving a goal like



So goal should really be:  $\text{AncestorOf}(\text{jane}, \text{billy}), !, \text{Male}(\text{jane})$

Similarly:

$\text{Member}(x, l) \leftarrow \text{FirstElement}(x, l)$

$\text{Member}(x, l) \leftarrow \text{Rest}(l, l') \wedge \text{Member}(x, l')$

If only to be used for testing, want

$\text{Member}(x, l) \text{ :- } \text{FirstElement}(x, l), !, .$

On failure, do not try  
to find another  $x$  later  
in the rest of the list

# Negation as failure

---

Procedurally: we can distinguish between the following:

can solve goal  $\neg G$  vs. cannot solve goal  $G$

Use **not**( $G$ ) to mean the goal that succeeds if  $G$  fails, and fails if  $G$  succeeds

Roughly: **not**( $G$ ) :-  $G$ , !, fail.                    /\* fail if  $G$  succeeds \*/  
                  **not**( $G$ ).                                        /\* otherwise succeed \*/

Only terminates when failure is *finite* (no more resolvents)

Useful when DB + rules is complete

NoChildren( $x$ ) :- **not**(ParentOf( $x,y$ ))

or when method already exists for complement

Composite( $n$ ) :-  $n > 1$ , **not**(PrimeNum( $n$ ))

Declaratively: same reading as  $\neg$ , but not when *new* variables in  $G$

[**not**(ParentOf( $x,y$ ))  $\supset$  NoChildren( $x$ )]    ✓  
vs. [ $\neg$ ParentOf( $x,y$ )  $\supset$  NoChildren( $x$ )]    ✗

# Dynamic DB

---

Sometimes useful to think of DB as a snapshot of the world that can be changed dynamically

assertions and deletions to the DB

then useful to consider 3 procedural interpretations for rules like

$$\text{ParentOf}(x,y) \Leftarrow \text{MotherOf}(x,y)$$

1. If-needed: Whenever have a goal matching  $\text{ParentOf}(x,y)$ , can solve it by solving  $\text{MotherOf}(x,y)$   
ordinary back-chaining, as in Prolog
2. If-added: Whenever something matching  $\text{MotherOf}(x,y)$  is added to the DB, also add  $\text{ParentOf}(x,y)$   
forward-chaining
3. If-removed: Whenever something matching  $\text{ParentOf}(x,y)$  is removed from the DB, also remove  $\text{MotherOf}(x,y)$ , if this was the reason  
keeping track of dependencies in DB

Interpretations (2) and (3) suggest demons

procedures that monitor DB and fire when certain conditions are met

# The Planner language

---

Main ideas:

1. DB of facts

(Mother susan john) (Person john)

2. If-needed, if-added, if-removed procedures consisting of

- body: program to execute
- pattern for invocation (Mother  $x$   $y$ )

3. Each program statement can succeed or fail

- **(goal  $p$ )**, **(assert  $p$ )**, **(erase  $p$ )**,
- **(and  $s$  ...  $s$ )**, statements with backtracking
- **(not  $s$ )**, negation as failure
- **(for  $p$   $s$ )**, do  $s$  for every way  $p$  succeeds
- **(finalize  $s$ )**, like cut
- a lot more, including all of Lisp

Shift from proving conditions  
to making conditions hold!

examples: **(proc if-needed (cleartable)**  
          **(for (on  $x$  table)**  
            **(and (erase (on  $x$  table)) (goal (putaway  $x$ ))))**  
**(proc if-removed (on  $x$   $y$ ) (print  $x$  " is no longer on "  $y$ ))**

---

7.

# Rules in Production Systems

# Direction of reasoning

---

A conditional like  $P \Rightarrow Q$  can be understood as transforming

- assertions of  $P$  to assertions of  $Q$
- goals of  $Q$  to goals of  $P$

Can represent the two cases explicitly:  $(\text{assert } P) \Rightarrow (\text{assert } Q)$   
 $(\text{goal } Q) \Rightarrow (\text{goal } P)$

and then distinguish between

1. goal vs. data directed reasoning
  - goal: from  $Q$  towards  $P$
  - data: from  $P$  towards  $Q$
2. forward vs. backward-chaining
  - forward: along the  $\Rightarrow$
  - backward: against the  $\Rightarrow$

Possible to have

- **(proc if-added (mygoal  $Q$ ) ... (mygoal  $P$ ))**
- **(proc if-needed (myassert  $P$ )... (myassert  $Q$ ))**

How to do data-directed reasoning in Prolog

Now: a formalism with forward-chaining



# Production systems

---

Idea: working memory + production rule set

Working memory: like DB, but volatile

Production rule: **IF** *conditions* **THEN** *actions*

condition: tests on WM

action: changes to WM

Basic operation: cycle of

1. recognize

find conflict set: rules whose conditions are satisfied by current WM

2. resolve

determine which of the rules will fire

3. act

perform required changes to WM

Stop when no rules fire

# Working memory

---

Set of working memory elements (WME)

Each WME is of the form  $(type\ attr_1\ val_1\ attr_2\ val_2\ \dots\ attr_n\ val_n)$

where  $type, attr_i, val_i$  are all atoms

Examples: (person age 27 home Toronto)  
(goal task openDoor importance 5)  
(student name JohnSmith dept CS)

Understood as  $\exists x[type(x) \wedge attr_1(x)=val_1 \wedge \dots \wedge attr_n(x)=val_n]$

- individual is not explicitly named
- order of attributes is not significant

Can handle n-ary relations as usual

(myAssertion relation OlderThan firstArg John secondArg Mary)

# Rule conditions

---

Conditions: tested conjunctively

a condition is  $p$  or  $\neg p$ , where  $p$  is a pattern of the form

$(type\ attr_1\ spec_1\ \dots\ attr_k\ spec_k)$

where each specification must be one of

- an atom
- an expression within [ ]
- a variable
- a test, within { }
- the  $\wedge$ ,  $\vee$ ,  $\neg$  of a specification

Examples:

(person age [ $n+4$ ] occupation  $x$ )  
- (person age { $< 23 \wedge > 6$ })

A rule is applicable if there are values of the variables to satisfy all the conditions

- for a pattern, need WME of the correct type and for each  $attr$  in pattern,  $val$  must match  $spec$
- for  $\neg p$ , there must be no WME that matches  $p$  ∴ negation as failure

# Rule actions

---

Actions: performed sequentially

An action is of the form

- **ADD** *pattern*
- **REMOVE** *index*
- **MODIFY** *index* (*attr spec*)

where

- index *i* refers to the WME that matched *i*-th pattern (inapplicable to *-p*)
- variables and expressions refer to values obtained in the matching

Examples:

**IF** (Student name *x*)  
**THEN**   **ADD** (Person name *x*)  
                  ordinary forward chaining

**IF** (Person age *x*) (Birthday)  
**THEN**   **REMOVE** 2  
                  **MODIFY** 1 (age [*x*+1])  
                                  database update

**IF** (starting)  
**THEN**   **REMOVE** 1  
                  **ADD** (phase val 1)   control information

# Example 1

---

Placing bricks in order of size

largest in place 1, next in place 2, etc.

Initial working memory

(counter index 1) (brick name A size 10 place heap) (brick name B size 30 place heap) (brick name C size 20 place heap)
--

Production rules:

**IF** (brick place heap name  $n$  size  $s$ )  
-(brick place heap size  $\{> s\}$ )  
-(brick place hand)

put the largest  
brick in your hand

**THEN MODIFY 1** (place hand)

**IF** (brick place hand) (counter index  $i$ )  
**THEN MODIFY 1** (place  $i$ )  
**MODIFY 2** (index  $[i+1]$ )

put a brick in your  
hand at the next spot

# Trace

---

Only one rule can fire at a time, so no conflict resolution is required

The following modifications to WM

1. (brick name B size 30 place hand)
2. (brick name B size 30 place 1)  
(counter index 2)
3. (brick name C size 20 place hand)
4. (brick name C size 20 place 2)  
(counter index 3)
5. (brick name A size 10 place hand)
6. (brick name A size 10 place 3)  
(counter index 4)

So the final working memory is

(counter index 4) (brick name A size 10 place 3) (brick name B size 30 place 1) (brick name C size 20 place 2)
---

## Example 2

---

How many days are there in a year?

Start with: (want-days year  $n$ )

End with: (has-days days  $m$ )

1. **IF** (want-days year  $n$ )  
**THEN REMOVE 1**  
**ADD** (year mod4 [ $n \bmod 4$ ]  
mod100 [ $n \bmod 100$ ]  
mod400 [ $n \bmod 400$ ])
2. **IF** (year mod400 0)  
**THEN REMOVE 1 ADD** (has-days days 366)
3. **IF** (year mod100 0 mod400 { $\neq 0$ })  
**THEN REMOVE 1 ADD** (has-days days 365)
4. **IF** (year mod4 0 mod100 { $\neq 0$ })  
**THEN REMOVE 1 ADD** (has-days days 366)
5. **IF** (year mod4 { $\neq 0$ })  
**THEN REMOVE 1 ADD** (has-days days 365)

# Applications

---

## 1. Psychological modeling

**IF** (goal is get-unit-digit)  
(minuend unit  $d$ )  
(subtrahend unit  $\{> d\}$ )

fine-grained modeling of symbol  
manipulation performed by people  
during problem solving

**THEN REMOVE 1**  
**ADD** (goal is borrow-from-tens)

## 2. Expert systems

rules used by experts in a problem area to perform complex tasks  
*(examples later)*

### Claimed advantages:

- modularity: each rule acts independently of the others
- fine-grained control: no complex goal or control stack
- transparency: can recast rules in English to provide explanation of behaviour

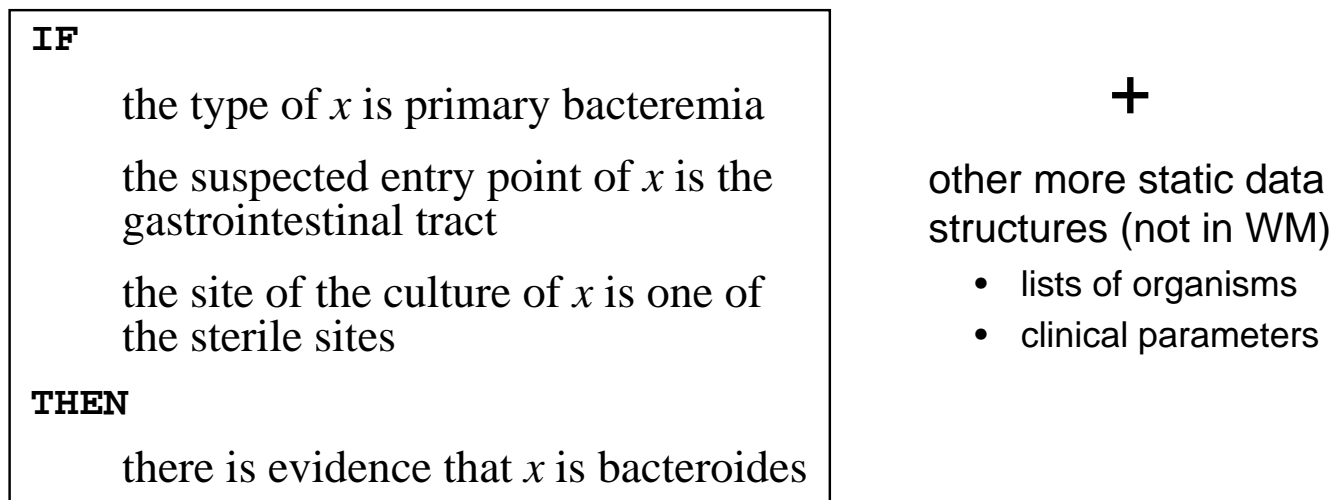


# MYCIN

---

System developed at Stanford to aid physicians in treating bacterial infections

Approximately 500 rules for recognizing about 100 causes of infection



## Certainty factors

numbers from 0 to 1 attached to conclusions to rank order alternatives

**AND** – take min      **OR** – take max

# XCON

---

System developed at CMU (as R1) and used extensively at DEC (now owned by Compaq) to configure early Vax computers

Nearly 10,000 rules for several hundred component types

Major stimulus for commercial interest in rule-based expert systems ★

**IF**

the context is doing layout and assigning a power supply

an sbi module of any type has been put in a cabinet

the position of the sbi module is known

there is space available for the power supply

there is no available power supply

the voltage and the frequency of the components are known

**THEN**

add an appropriate power supply

# Context switching

---

XCON and others use rules of the form

```
IF the current context is  $x$   
THEN deactivate  $x$   
         activate context  $y$ 
```

organized to fire when no other rules apply

Useful for grouping rules

```
IF (control phase 1) AND ...  
THEN ...  
...  
IF (control phase 1) AND ...  
THEN ... MODIFY 1 (phase 2) ...
```

Allows emulation of  
control structures.

```
IF (control phase 2) AND ...  
THEN ...  
...  
IF (control phase 2) AND ...  
THEN ... MODIFY 1 (phase 3) ...
```

But still difficult for  
*complex* control

# Conflict resolution

---

Sometimes with data-directed reasoning, we want to fire *all* applicable rules

With goal-directed reasoning, we may want a *single* rule to fire

- arbitrary
- first rule in order of presentation (as in Prolog)
- specificity, as in
  - IF** (bird) **THEN ADD** (can-fly)
  - IF** (bird weight {> 100}) **THEN ADD** (cannot-fly)
  - IF** (bird) (penguin) **THEN ADD** (cannot-fly)
- recency
  - fire on rule that uses most recent WME
  - fire on least recently used rule
- refractoriness
  - never use same rule for same value of variables (called rule instance)
  - only use a rule/WME pair once (will need a “refresh” otherwise)

# Conflict combinations

---

## OPS5:

1. discard rule instances that have already been used
2. order remaining instances in terms of recency of WME matching 1st condition (and then of 2nd condition, etc.)
3. if still no single rule, order rules by number of conditions
4. select arbitrarily among those remaining

## SOAR:

system that attempts to find a way to move from a start state to a goal state by applying productions

selecting what rule to fire

≡

deciding what to do next

if unable to decide, SOAR sets up the selection as a new (meta-)goal to solve, and the process iterates

# Rete procedure

Early systems spent 90% of their time matching, even with indexing and hashing.

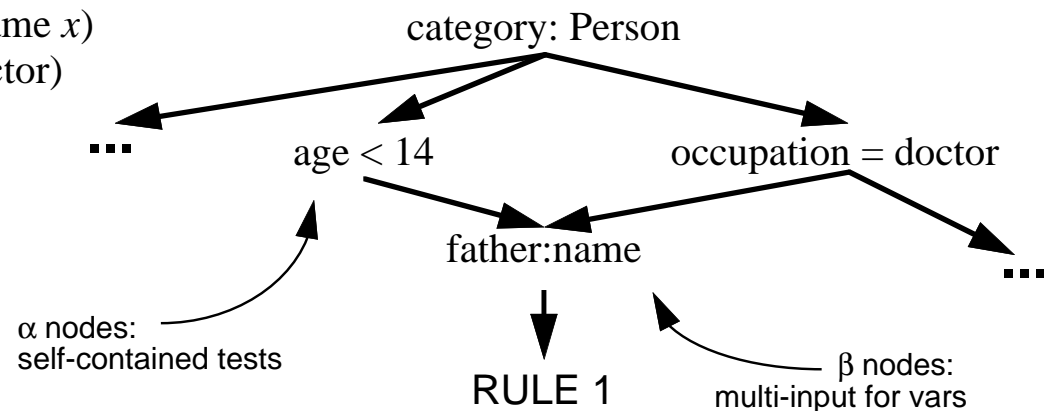
But:

- WM is modified only slightly on each cycle
- many rules share conditions

So:

- incrementally pass WME through network of tests
- tokens that make it through satisfy all conditions and produce conflict set
- can calculate new conflict set in terms of old one and change to WM

**IF** (Person father y age {< 14} name x)  
(Person name y occupation doctor)  
**THEN** ...



---

8.

# Object-Oriented Representation

# Organizing procedures

---

With the move to put control of inference into the user's hands, we're focusing on more procedural representations

knowing facts by executing code

Even production systems are essentially programming languages.

Note also that everything so far is *flat*, i.e., sentence-like representations

- information about an object is scattered in axioms
- procedure fragments and rules have a similar problem

With enough procedures / sentences in a KB, it could be critical to *organize* them

- production systems might have rule sets, organized by context of application
- but this is not a natural, *representational* motivation for grouping



# Object-centered representation

---

Most obvious organizational technique depends on our ability to see the world in terms of objects

- physical objects:
  - a desk has a surface-material, # of drawers, width, length, height, color, procedure for unlocking, etc.
  - some variations: no drawers, multi-level surface, built-in walls (carrel)
- also, *situations* can be object-like:
  - a class: room, participants, teacher, day, time, seating arrangement, lighting, procedures for registering, grading, etc.
  - leg of a trip: destination, origin, conveyance, procedures for buying ticket, getting through customs, reserving hotel room, locating a car rental etc.

Suggests clustering procedures for determining properties, identifying parts, interacting with parts, as well as constraints between parts, all of *objects*

- legs of desk connect to and support surface
  - beginning of a travel leg and destination of prior one
- object-centered constraints

# Situation recognition

---

Focus on objects as an organizational / chunking mechanism to make some things easier to find

Suggests a different kind of reasoning than that covered so far

basic idea originally proposed by Marvin Minsky

- recognize (guess) situation; activate relevant object representations
- use those object representations to set up expectations
  - some for verification; some make it easier to interpret new details
- flesh out situation once you've recognized

Wide applicability, but typical applications include

- relationship recognition e.g., story understanding
- data monitoring
- propagation and enforcement of constraints for planning tasks
  - this latter is most doable and understandable,  
so we will concentrate on it

# Basic frame language

---

Let's call our object structures frames

note wide variety of interpretations in literature

Two types:

- individual frames  
represent a single object like a person, part of a trip
- generic frames  
represent categories of objects, like students

An individual frame is a named list of buckets called slots. What goes in the bucket is called a filler of the slot. It looks like this:

*(frame-name*  
    <*slot-name1* *filler1*>  
    <*slot-name2* *filler2* > ...)  
*)*                    where frame names and slot names are atomic,  
                                and fillers are either numbers, strings or the  
                                names of other individual frames.

Notation: individual frames: toronto  
                  slot names:       :Population (note ":" at start)  
                  generic frames: CanadianCity

# Instances and specializations

---

Individual frames have a special slot called :INSTANCE-OF whose filler is the name of a generic frame:

```
(toronto
  <:INSTANCE-OF CanadianCity>
  <:Province ontario>
  <:Population 4.5M>...)
```

```
(tripLeg123-1
  <:INSTANCE-OF TripLeg>
  <:Destination toronto>...)
```

Generic frames have a syntax that is similar to that of individual frames, except that they have a slot called :IS-A whose filler is the name of another generic frame

```
(CanadianCity
  <:IS-A City>
  <:Province CanadianProvince>
  <:Country canada>...)
```

We say that the frame toronto is an instance of the frame CanadianCity and that the frame CanadianCity is a specialization of the frame City

# Procedures and defaults

---

Slots in generic frames can have associated procedures

1. computing a filler (when no slot filler is given)

(Table

<:Clearance [**IF-NEEDED** computeClearanceFromLegs]> ...)

2. propagating constraints (when a slot filler is given)

(Lecture

<:DayOfWeek WeekDay>

<:Date [**IF-ADDED** computeDayOfWeek]> ...)

If we create an instance of Table, the :Clearance will be calculated as needed. Similarly, the filler for :DayOfWeek will be calculated when :Date is filled.

For instances of CanadianCity, the :Country slot will be filled automatically. But we can also have

(city135

<:**INSTANCE-OF** CanadianCity>

<:Country holland>)

The filler canada in CanadianCity is considered a default value.

# IS-A and inheritance

---

Specialization relationships imply that procedures and fillers of more general frame are applicable to more specific frame:  
inheritance.

For example, instances of MahoganyCoffeeTable will inherit the procedure from Table (via CoffeeTable)

(CoffeeTable  
  <:**IS-A** Table> ...)  
(MahoganyCoffeeTable  
  <:**IS-A** CoffeeTable> ...)

Similarly, default values are inheritable, so that Clyde inherits a colour from RoyalElephant, not Elephant

(Elephant  
  <:**IS-A** Mammal>  
  <:Colour gray> ...)  
(RoyalElephant  
  <:**IS-A** Elephant>  
  <:Colour white>)  
(clyde  
  <:**INSTANCE-OF** RoyalElephant>)

# Reasoning with frames

---

Basic (local) reasoning goes like this:

1. user instantiates a frame, i.e., declares that an object or situation exists
2. slot fillers are inherited where possible
3. inherited **IF-ADDED** procedures are run, causing more frames to be instantiated and slots to be filled.

If the user or any procedure requires the filler of a slot then:

1. if there is a filler, it is used
2. otherwise, an inherited **IF-NEEDED** procedure is run, potentially causing additional actions

Globally:

- make frames be major situations or object-types you need to flesh out
- express constraints between slots as **IF-NEEDED** and **IF-ADDED** procedures
- fill in default values when known

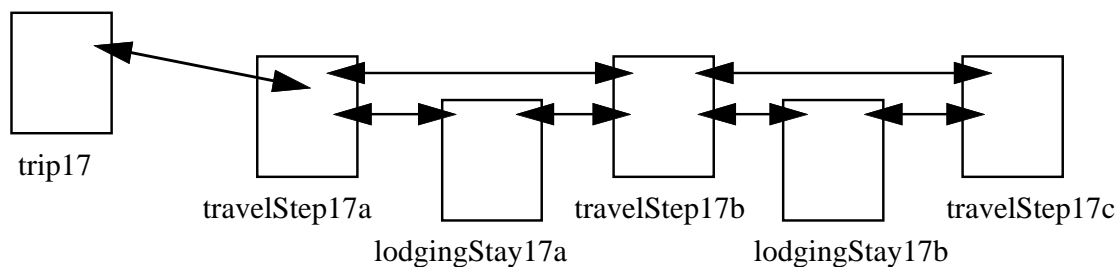
⇒ like a fancy, semi-symbolic spreadsheet

# Planning a trip

A simple example: a frame system to assist in travel planning  
(and possibly documentation – automatically generate forms)

Basic structure (main frame types):

- a Trip will be a sequence of TravelSteps  
these will be linked together by slots
- a TravelStep will usually terminate in a LodgingStay (except the last, or one with two travels on one day)
  - a LodgingStay will point to its arriving TravelStep and departing TravelStep
  - TravelSteps will indicate the LodgingStays of their origin and destination



```
(trip17
  <:INSTANCE-OF Trip>
  <:FirstStep travelStep17a>
  <:Traveler ronB> ...)
```



# Parts of a trip

---

TravelSteps and LodgingStays share some properties (e.g., :BeginDate, :EndDate, :Cost, :PaymentMethod), so we might create a more general category as the parent frame for both of them:

(Trip

<:FirstStep TravelStep>  
<:Traveler Person>  
<:BeginDate Date>  
<:TotalCost Price> ...)

(TripPart

<:BeginDate>  
<:EndDate>  
<:Cost>  
<:PaymentMethod> ...)

(TravelStep

<:**IS-A** TripPart>  
<:Means>  
<:Origin> <:Destination>  
<:NextStep> <:PreviousStep>  
<:DepartureTime> <:ArrivalTime>  
<:OriginLodgingStay>  
<:DestinationLodgingStay> ...)

(LodgingStay

<:**IS-A** TripPart>  
<:ArrivingTravelStep>  
<:DepartingTravelStep>  
<:City>  
<:LodgingPlace> ...)

# Travel defaults and procedures

---

## Embellish frames with defaults and procedures

(TravelStep  
 <:Means airplane> ...)

(TripPart  
 <:PaymentMethod visaCard> ...)


(TravelStep  
 <:Origin [**IF-NEEDED** {if no SELF:PreviousStep then newark}]]>)

(Trip  
 <:TotalCost  
 [**IF-NEEDED**  
 { x←SELF:FirstStep;  
 result←0;  
 repeat  
 { if exists x:NextStep  
 then  
 { result←result + x:Cost +  
 x:DestinationLodgingStay:Cost;  
 x←x:NextStep }  
 else return result+x:Cost } }]]>)

Program notation (for an imaginary language):

- SELF is the current frame being processed
- if  $x$  refers to an individual frame, and  $y$  to a slot, then  $xy$  refers to the filler of the slot

assume this  
is 0 if there is  
no LodgingStay



## More attached procedures

---

```
(TravelStep
  <:NextStep
    [IF-ADDED
      {if SELF:EndDate ≠ SELF:NextStep:BeginDate
        then
          SELF:DestinationLodgingStay ←
            SELF:NextStep:OriginLodgingStay ←
              create new LodgingStay
                with :BeginDate = SELF:EndDate
                and with :EndDate = SELF:NextStep:BeginDate
                and with :ArrivingTravelStep = SELF
                and with :DepartingTravelStep = SELF:NextStep
              ...}}>
  ...)
```

Note: default :City of LodgingStay, etc. can also be calculated:

```
(LodgingStay
  <:City [IF-NEEDED {SELF:ArrivingTravelStep:Destination}]...> ...)
```

# Frames in action

---

Propose a trip to Toronto on Dec. 21, returning Dec. 22

(trip18

<:INSTANCE-OF Trip>  
<:FirstStep travelStep18a>)

the first thing to do is to create  
the trip and the first step

(travelStep18a

<:INSTANCE-OF TravelStep>  
<:BeginDate 12/21/98>  
<:EndDate 12/21/98>  
<:Means>  
<:Origin>  
<:Destination toronto>  
<:NextStep> <:PreviousStep>  
<:DepartureTime> <:ArrivalTime>)

the next thing to do is to create  
the second step and link it to the first  
by changing the :NextStep

(travelStep18b

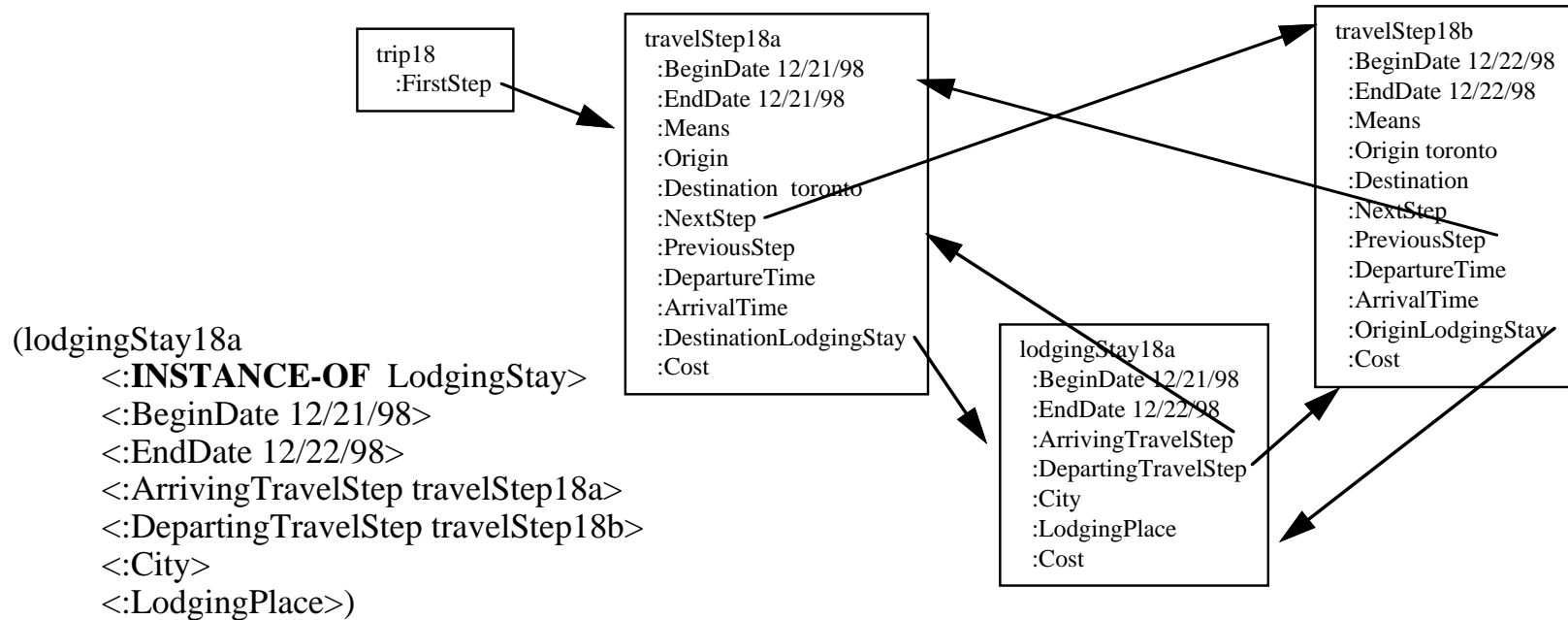
<:INSTANCE-OF TravelStep>  
<:BeginDate 12/22/98>  
<:EndDate 12/22/98>  
<:Means>  
<:Origin toronto>  
<:Destination>  
<:NextStep>  
<:PreviousStep travelStep18a>  
<:DepartureTime> <:ArrivalTime>)

(travelStep18a

<:NextStep travelStep18b>)

# Triggering procedures

**IF-ADDED** on :NextStep then creates a LodgingStay:

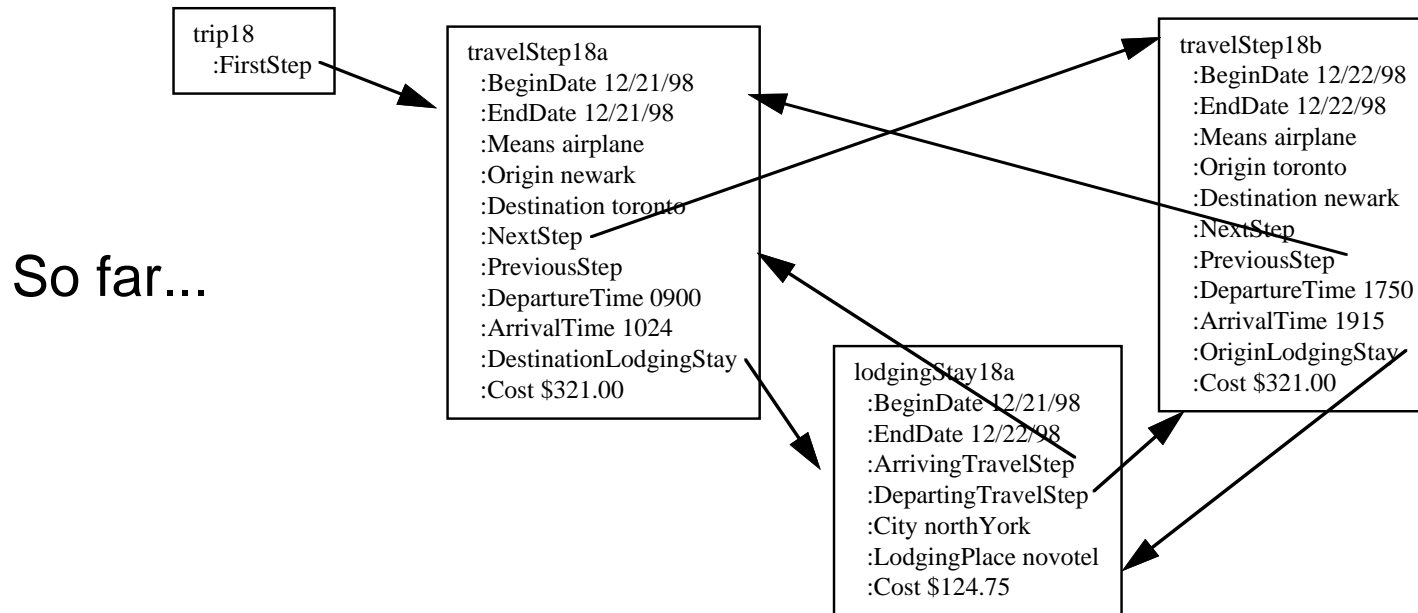


If requested, **IF-NEEDED** can provide :City for lodgingStay18a (toronto)

which could then be overridden by hand, if necessary  
 (e.g. usually stay in North York, not Toronto)

Similarly, apply default for :Means and default calc for :Origin

# Finding the cost of the trip



Finally, we can use `:TotalCost` **IF-NEEDED** procedure (see above) to calculate the total cost of the trip:

- `result ← 0, x ← travelStep18a, x:NextStep = travelStep18b`
- `result ← 0 + $321.00 + $124.75; x ← travelStep18b, x:NextStep = NIL`
- `return: result = $445.75 + $321.00 = $766.75`

# Using the formalism

---

Main purpose of the above: embellish a sketchy description with defaults, implied values

- maintain consistency
- use computed values to
  1. allow derived properties to look explicit
  2. avoid up front, potentially unneeded computation

## Monitoring

- hook to a DB, watch for changes in values
- like an ES somewhat, but monitors are more object-centered, inherited

## Scripts for story understanding

generate expectations (e.g., restaurant)

## Real, Minsky-like commonsense reasoning

- local cues  $\Rightarrow$  potentially relevant frames  $\Rightarrow$  further expectations
- look to match expectations ; mismatch  $\Rightarrow$  “differential diagnosis”

# Extensions

---

## 1. Types of procedures

- **IF-REMOVED**

e.g., remove TravelStep  $\Rightarrow$  remove LodgingStay

- “servants” and “demons”

flexible “pushing” and “pulling” of data

## 2. Slots

- multiple fillers

- “facets” – more than just defaults and fillers

- [**REQUIRE** <class>] (or procedure)

- **PREFER** – useful if conflicting fillers

## 3. Metaframes

(CanadianCity <:**INSTANCE-OF** GeographicalCityType> ...)

(GeographicalCityType <:**IS-A** CityType>

<:AveragePopulation NonNegativeNumber> ...)

## 4. Frames as actions (“scripts”)



# Object-oriented programming

---

Somewhat in the manner of production systems, specifying problems with frames can easily slide into a style of *programming*, rather than a declarative object-oriented modeling of the world

- note that direction of procedures (pushing/pulling) is explicitly specified  
not declarative

This drifts close to conventional object-oriented programming (developed concurrently).

- same advantages:
  - definition by specialization
  - localization of control
  - encapsulation
  - etc.
- main difference:
  - frames: centralized, conventional control regime (instantiate/ inherit/trigger)
  - object-oriented programming: objects acting as small, independent agents sending each other messages

---

9.

# Structured Descriptions

# From sentences to objects

---

As we saw with frames, it useful to shift the focus away from the true *sentences* of an application towards the categories of *objects* in the application and their properties.

In frame systems, this was done *procedurally*, and we concentrated on hierarchies of frames as a way of organizing collections of procedures.

In this section, we look at the categories of objects themselves:

- objects are members of multiple categories  
e.g. a doctor, a wife, a mother of two
- categories of objects can be more or less specific than others  
e.g. a doctor, a professional, a surgeon
- categories of objects can have parts, sometimes in multiples  
e.g. books have titles, tables have legs
- the relation among the parts of an object can be critical in its being a member of a category  
e.g. a stack vs. a pile of bricks

# Noun phrases

---

In FOL, all categories and properties of objects are represented by atomic predicates.

- In some cases, these correspond to simple *nouns* in English such as Person or City.
- In other cases, the predicates seem to be more like *noun phrases* such as MarriedPerson or CanadianCity or AnimalWithFourLegs.

Intuitively, these predicates have an internal structure and connections to other predicates.

e.g. A married person must be a person.

These connections hold by *definition* (by virtue of what the predicates themselves mean), not by virtue of the facts we believe about the world.

In FOL, there is no way to break apart a predicate to see how it is formed from other predicates.

Here we will examine a logic that allows us to have both atomic and non-atomic predicates: a description logic

# Concepts, roles, constants

---

In a description logic, there are sentences that will be true or false (as in FOL).

In addition, there are three sorts of expressions that act like nouns and noun phrases in English:

- concepts are like category nouns      Dog, Teenager, GraduateStudent
- roles are like relational nouns      :Age, :Parent, :AreaOfStudy
- constants are like proper nouns      johnSmith, chair128

These correspond to unary predicates, binary predicates and constants (respectively) in FOL.

See also: generic frames, slots, and individual frames.

However: roles can have multiple fillers.

However, unlike in FOL, concepts need not be atomic and can have semantic relationships to each other.

roles will remain atomic (for now)

# The symbols of DL

---

## Three types of non-logical symbols:

- atomic concepts:

Dog, Teenager, GraduateStudent

We include a distinguished concept: Thing

- roles: (all are atomic)

:Age, :Parent, :AreaOfStudy

- constants:

johnSmith, chair128

## Four types of logical symbols:

- punctuation: [, ], (, )

- positive integers: 1, 2, 3, ...

- concept-forming operators: ALL, EXISTS, FILLS, AND

- connectives:  $\equiv$ ,  $\doteq$ , and  $\rightarrow$

# The syntax of DL

---

The set of concepts is the least set satisfying:

- Every atomic concept is a concept.
- If  $r$  is a role and  $d$  is a concept, then  $[\text{ALL } r d]$  is a concept.
- If  $r$  is a role and  $n$  is an integer, then  $[\text{EXISTS } n r]$  is a concept.
- If  $r$  is a role and  $c$  is a constant, then  $[\text{FILLS } r c]$  is a concept.
- If  $d_1, \dots, d_k$  are concepts, then so is  $[\text{AND } d_1, \dots, d_k]$ .

Three types of sentences in DL:

- If  $d$  and  $e$  are concepts, then  $(d \equiv e)$  is a sentence.
- if  $d$  and  $e$  are concepts, then  $(d \dot{=} e)$  is a sentence.
- If  $d$  is a concept and  $c$  is a constant, then  $(c \rightarrow d)$  is a sentence.

# The meaning of concepts

---

Constants stand for individuals, concepts for sets of individuals, and roles for binary relations.

The meaning of a complex concept is derived from the meaning of its parts the same way a noun phrases is:

- [EXISTS  $n$   $r$ ] describes those individuals that stand in relation  $r$  to at least  $n$  other individuals
- [FILLS  $r$   $c$ ] describes those individuals that stand in the relation  $r$  to the individual denoted by  $c$
- [ALL  $r$   $d$ ] describes those individuals that stand in relation  $r$  only to individuals that are described by  $d$
- [AND  $d_1 \dots d_k$ ] describes those individuals that are described by all of the  $d_i$ .

For example:

“a company with at least 7 directors,  
whose managers are all women with  
PhDs, and whose min salary is \$24/hr”

[AND Company  
[EXISTS 7 :Director]  
[ALL :Manager [AND Woman  
[FILLS :Degree PhD]]]  
[FILLS :MinSalary \$24.00/hour]]



# A DL knowledge base

---

A DL knowledge base is a set of DL sentences serving mainly to

- give names to definitions

e.g. (FatherOfDaughters  $\doteq$   
[AND Male [EXISTS 1 :Child]  
[ALL :Child Female]] )

“A FatherOfDaughters is precisely a male with at least one child and all of whose children are female”

- give names to partial definitions

e.g. (Dog  $\sqsubseteq$  [AND Mammal Pet  
CarnivorousAnimal  
[FILLS :VoiceCall barking]])

“A dog is among other things a mammal that is a pet and a carnivorous animal whose voice call includes barking”

gives necessary but not sufficient conditions

- assert properties of individuals

e.g. (joe  $\rightarrow$   
[AND FatherOfDaughters Surgeon]])

“Joe is a FatherOfDaughters and a Surgeon”

Other types of DL sentences are typically not used in a KB.

e.g. ([AND Rational Animal]  $\doteq$  [AND Featherless Biped])

# Formal semantics

---

Interpretation  $\mathcal{I} = \langle D, I \rangle$  as in FOL, where

- for every constant  $c$ ,  $I[c] \in D$
- for every atomic concept  $a$ ,  $I[a] \subseteq D$
- for every role  $r$ ,  $I[r] \subseteq D \times D$

We then extend the interpretation to all concepts as subsets of the domain as follows:

- $I[\text{Thing}] = D$
- $I[[\text{ALL } r \ d]] = \{x \in D \mid \text{for any } y, \text{ if } \langle x, y \rangle \in I[r] \text{ then } y \in I[d]\}$
- $I[[\text{EXISTS } n \ r]] = \{x \in D \mid \text{there are at least } n \ y \text{ such that } \langle x, y \rangle \in I[r]\}$
- $I[[\text{FILLS } r \ c]] = \{x \in D \mid \langle x, I[c] \rangle \in I[r]\}$
- $I[[\text{AND } d_1 \ \dots \ d_k]] = I[d_1] \cap \dots \cap I[d_k]$

A sentence of DL will then be true or false as follows:

- $\mathcal{I} \models (d \sqsubseteq e)$  iff  $I[d] \subseteq I[e]$
- $\mathcal{I} \models (d \doteq e)$  iff  $I[d] = I[e]$
- $\mathcal{I} \models (c \rightarrow e)$  iff  $I[c] \in I[e]$

# Entailment and reasoning

---

Entailment in DL is defined as in FOL:

A set of DL sentences  $S$  entails a sentence  $\alpha$  (which we write  $S \models \alpha$ ) iff  
for every  $\mathcal{I}$ , if  $\mathcal{I} \models S$  then  $\mathcal{I} \models \alpha$

A sentence is valid iff it is entailed by the empty set.

Given a KB consisting of DL sentences, there are two basic sorts of reasoning we consider:

1. determining if  $KB \models (c \rightarrow e)$

whether a named individual satisfies a certain description

2. determining if  $KB \models (d \sqsubseteq e)$

whether one description is subsumed by another

the other case,  $KB \models (d \dot{\sqsubseteq} e)$  reduces to

$KB \models (d \sqsubseteq e)$  and  $KB \models (d \sqsupseteq e)$

# Entailment vs. validity

---

In some cases, an entailment will hold because the sentence in question is valid.

- ([AND Doctor Female]  $\models$  Doctor)
- ([FILLS :Child sue]  $\models$  [EXISTS 1 :Child])
- (john  $\rightarrow$  [ALL :Hobby Thing])

But in most other cases, the entailment depends on the sentences in the KB.

For example,

([AND Surgeon Female]  $\models$  Doctor)

is not valid.

But it is entailed by a KB that contains

(Surgeon  $\stackrel{\bullet}{\models}$  [AND Specialist [FILLS :Specialty surgery]])

(Specialist  $\models$  Doctor)

# Computing subsumption

---

We begin with computing subsumption, that is, determining whether or not  $\text{KB} \models (d \sqsubseteq e)$ .

and therefore  
whether  $d \dot{\sqsubseteq} e$

Some simplifications to the KB:

- we can remove the  $(c \rightarrow d)$  assertions from the KB
- we can replace  $(d \sqsubseteq e)$  in KB by  $(d \dot{\sqsubseteq} [\text{AND } e a])$ , where  $a$  is a new atomic concept
- we assume that in the KB for each  $(d \dot{\sqsubseteq} e)$ , the  $d$  is atomic and appears only once on the LHS
- we assume that the definitions in the KB are acyclic  
vs. cyclic  $(d \dot{\sqsubseteq} [\text{AND } e f]), (e \dot{\sqsubseteq} [\text{AND } d g])$

Under these assumptions, it is sufficient to do the following:

- normalization: using the definitions in the KB, put  $d$  and  $e$  into a special normal form,  $d'$  and  $e'$
- structure matching: determine if each part of  $e'$  is matched by a part of  $d'$ .

# Normalization

Repeatedly apply the following operations to the two concepts:

- expand a definition: replace an atomic concept by its KB definition

- flatten an AND concept:

$$[\text{AND } \dots [\text{AND } d e f] \dots] \Rightarrow [\text{AND } \dots d e f \dots]$$

- combine the ALL operations with the same role:

$$[\text{AND } \dots [\text{ALL } r d] \dots [\text{ALL } r e] \dots] \Rightarrow [\text{AND } \dots [\text{ALL } r [\text{AND } d e]] \dots]$$

- combine the EXISTS operations with the same role:

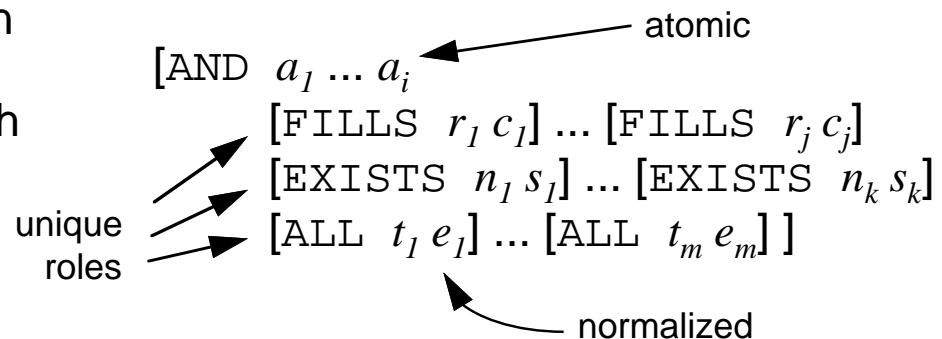
$$[\text{AND } \dots [\text{EXISTS } n_1 r] \dots [\text{EXISTS } n_2 r] \dots] \Rightarrow$$

$$[\text{AND } \dots [\text{EXISTS } n r] \dots] \quad (\text{where } n = \text{Max}(n_1, n_2))$$

- remove a vacuous concept: Thing, [ALL r Thing], [AND]

- remove a duplicate expression

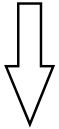
In the end, we end up with a normalized concept of the following form



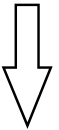
# Normalization example

---

[AND Person  
[ALL :Friend Doctor]  
[EXISTS 1 :Accountant]  
[ALL :Accountant [EXISTS 1 :Degree]]  
[ALL :Friend Rich]  
[ALL :Accountant [AND Lawyer [EXISTS 2 :Degree]]]]



[AND Person  
[EXISTS 1 :Accountant]  
[ALL :Friend [AND Rich Doctor]]  
[ALL :Accountant [AND Lawyer [EXISTS 1 :Degree] [EXISTS 2 :Degree]]]]



[AND Person  
[EXISTS 1 :Accountant]  
[ALL :Friend [AND Rich Doctor]]  
[ALL :Accountant [AND Lawyer [EXISTS 2 :Degree]]]]

# Structure matching

---

Once we have replaced atomic concepts by their definitions, we no longer need to use the KB.

To see if a normalized concept  $[\text{AND } e_1 \dots e_m]$  subsumes a normalized concept  $[\text{AND } d_1 \dots d_n]$ , we do the following:

For each component  $e_j$ , check that there is a matching component  $d_i$ , where

- if  $e_j$  is atomic or  $[\text{FILLS } r \ c]$ , then  $d_i$  must be identical to it;
- if  $e_j = [\text{EXISTS } 1 \ r]$ , then  $d_i$  must be  $[\text{EXISTS } n \ r]$  or  $[\text{FILLS } r \ c]$ ;
- if  $e_j = [\text{EXISTS } n \ r]$  where  $n > 1$ , then  $d_i$  must be of the form  $[\text{EXISTS } m \ r]$  where  $m \geq n$ ;
- if  $e_j = [\text{ALL } r \ e']$ , then  $d_i$  must be  $[\text{ALL } r \ d']$ , where recursively  $e'$  subsumes  $d'$ .

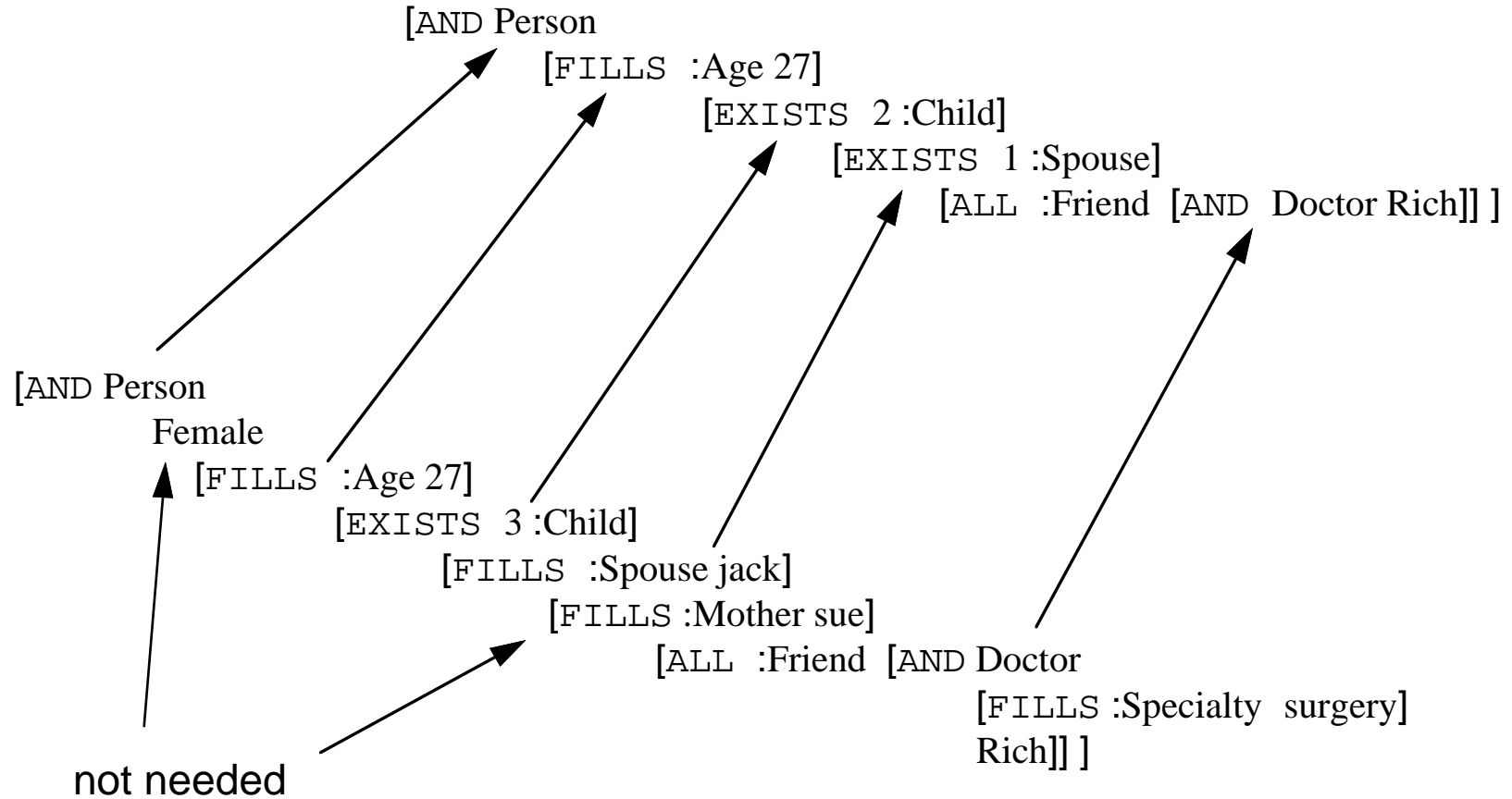
In other words, for every part of the more general concept, there must be a corresponding part in the more specific one.

It can be shown that this procedure is sound and complete:  
it returns YES iff  $\text{KB} \models (d \sqsubseteq e)$ .



# Structure matching example

---



# Computing satisfaction

---

To determine if  $\text{KB} \models (c \rightarrow e)$ , we use the following procedure:

1. find the most specific concept  $d$  such that  $\text{KB} \models (c \rightarrow d)$
2. determine whether or not  $\text{KB} \models (d \sqsubseteq e)$ , as before.

To a first approximation, the  $d$  we need is the AND of every  $d_i$  such that  $(c \rightarrow d_i) \in \text{KB}$ .

Suppose the KB contains

However, this can  
miss some inferences!

(joe  $\rightarrow$  Person)  
(canCorp  $\rightarrow$  [AND Company  
[ALL :Manager Canadian]  
[FILLS :Manager joe]])

then the KB  $\models$  (joe  $\rightarrow$  Canadian).

To find the  $d$ , a more complex procedure is used that *propagates* constraints from one individual (canCorp) to another (joe).

The individuals we need to consider need not be named by constants; they can be individuals that arise from EXISTS (like Skolem constants).

# Taxonomies

---

Two common sorts of queries in a DL system:

- given a query concept  $q$ , find all constants  $c$  such that  $\text{KB} \models (c \rightarrow q)$   
e.g.  $q$  is [AND Stock FallingPrice MyHolding]      might want to trigger a procedure for each such  $c$
- given a query constant  $c$ , find all *atomic* concepts  $a$  such that  $\text{KB} \models (c \rightarrow a)$

We can exploit the fact that concepts tend to be structured hierarchically to answer queries like these more efficiently.

Taxonomies arise naturally out of a DL KB:

- the nodes are the atomic concepts that appear on the LHS of a sentence  $(a \sqsubseteq d)$  or  $(a \doteq d)$  in the KB
- there is an edge from  $a_i$  to  $a_j$  if  $(a_i \sqsubseteq a_j)$  is entailed and there is no distinct  $a_k$  such that  $(a_i \sqsubseteq a_k)$  and  $(a_k \sqsubseteq a_j)$ .

can link every constant  $c$  to the most specific atomic concepts  $a$  in the taxonomy such that  $\text{KB} \models (c \rightarrow a)$

Positioning a new atom in a taxonomy is called classification

# Computing classification

---

Consider adding  $(a_{new} \dot{\equiv} d)$  to the KB.

- find  $S$ , the most specific subsumers of  $d$ : the atoms  $a$  such that  $\text{KB} \models (d \sqsupseteq a)$ , but nothing below  $a$  see below
- find  $G$ , the most general subsumees of  $d$ : the atoms  $a$  such that  $\text{KB} \models (a \sqsupseteq d)$ , but nothing above  $a$ 
  - if  $S \cap G$  is not empty, then  $a_{new}$  is not new
- remove any links from atoms in  $G$  to atoms in  $S$
- add links from all the atoms in  $G$  to  $a_{new}$  and from  $a_{new}$  to all the atoms in  $S$
- reorganize the constants:
  - for each constant  $c$  such that  $\text{KB} \models (c \rightarrow a)$  for all  $a \in S$ , but  $\text{KB} \not\models (c \rightarrow a)$  for no  $a \in G$ , and where  $\text{KB} \models (c \rightarrow d)$ , remove links from  $c$  to  $S$  and put a single link from  $c$  to  $a_{new}$ .

Adding  $(a_{new} \sqsupseteq d)$  is similar, but with no subsumees.

# Subsumers and subsumees

---

Calculating the most specific subsumers of a concept  $d$ :

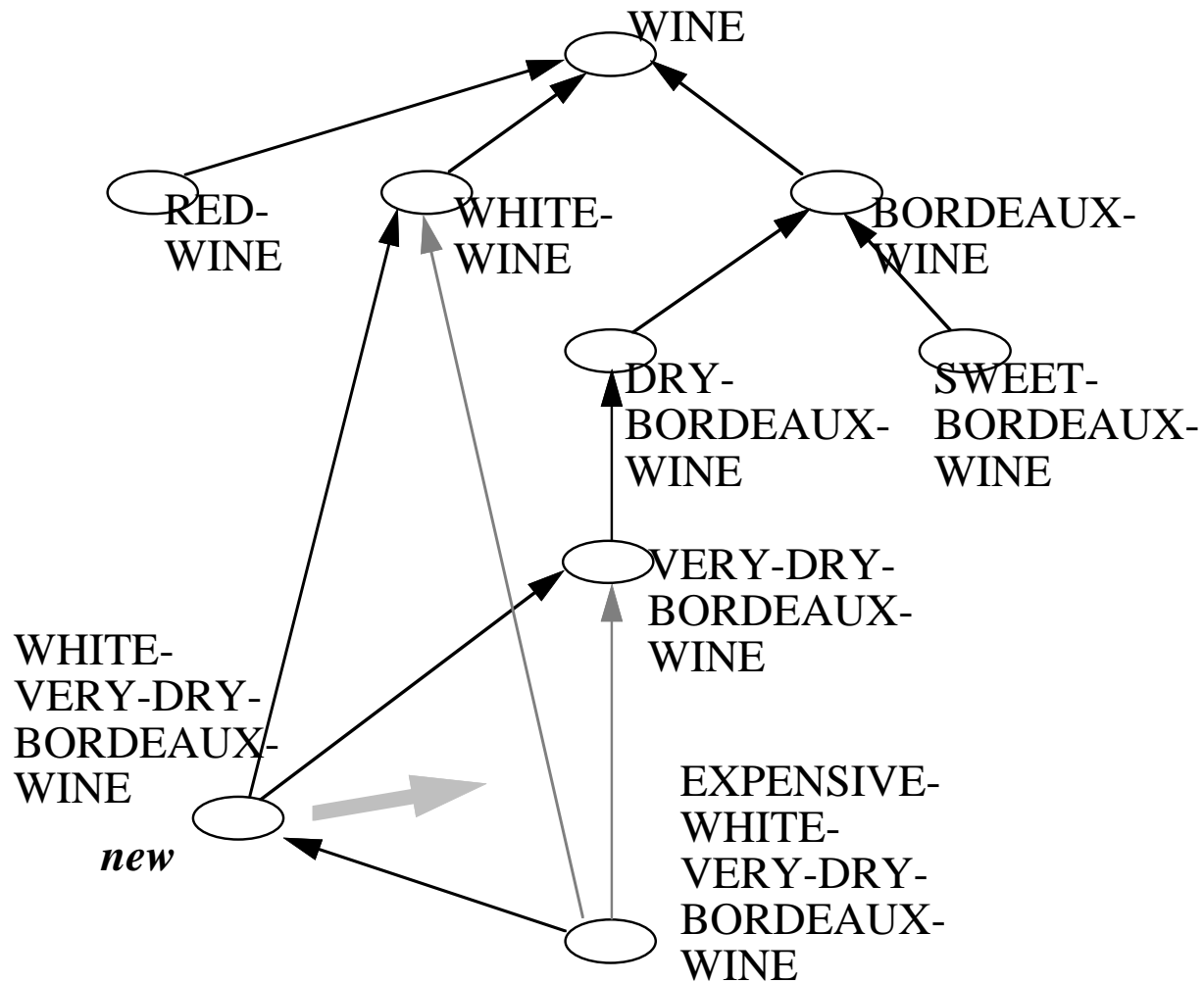
- Start with  $S = \{\text{Thing}\}$ .
- Repeatedly do the following:
  - Suppose that some  $a \in S$  has at least one child  $a'$  just below it in the taxonomy such that  $\text{KB} \models (d \sqsubseteq a')$ .
  - Then remove  $a$  from  $S$  and replace it by all such children  $a'$ .

Calculating the most general subsumees of a concept  $d$ :

- Start with  $G =$  the most specific subsumers.
- Repeatedly do the following:
  - Suppose that for some  $a \in G$ ,  $\text{KB} \not\models (a \sqsubseteq d)$ .
  - Then remove  $a$  from  $G$  and replace it by all of its children (or delete it, if there are none).
- Repeatedly delete any element of  $G$  that has a parent subsumed by  $d$ .

# An example of classification

---



# Using the taxonomic structure

---

Note that classification uses the structure of the taxonomy:

If there is an  $a'$  just below  $a$  in the taxonomy such that  $\text{KB} \not\models (d \sqsubseteq a')$ , we never look below this  $a'$ . If this concept is sufficiently high in the taxonomy (e.g. just below Thing), an entire subtree will be ignored.

Queries can also exploit the structure:

For example, to find the constants described by a concept  $q$ , we simply classify  $q$  and then look for constants in the part of the taxonomy subtended by  $q$ . The rest of the taxonomy not below  $q$  is ignored.

This natural structure allows us to build and use very large knowledge bases.

- the time taken will grow linearly with the *depth* of the taxonomy
- we would expect the depth of the taxonomy to grow *logarithmically* with the size of the KB
- under these assumptions, we can handle a KB with thousands or even millions of concepts and constants.

# Taxonomies vs. frame hierarchies

---

The taxonomies in DL look like the **IS-A** hierarchies with frames.

There is a big difference, however:

- in frame systems, the KB designer gets to decide what the fillers of the :IS-A slot will be; the :IS-A hierarchy is constructed manually
- in DL, the taxonomy is completely determined by the meaning of the concepts and the subsumption relation over concepts

For example, a concept such as

[AND Fish [FILLS :Size large]]

must appear in the taxonomy below Fish even if it was first constructed to be given the name Whale. It cannot simply be positioned below Mammal.

To correct our mistake, we need to associate the name with a different concept:

[AND Mammal [FILLS :Size large] ...]



# Inheritance and propagation

---

As in frame hierarchies, atomic concepts in DL inherit properties from concepts higher up in the taxonomy.

For example, if a Doctor has a medical degree, and Surgeon is below Doctor, then a Surgeon must have a medical degree.

This follows from the logic of concepts:

If  $KB \models (\text{Doctor} \sqsubseteq [\text{EXISTS } 1 : \text{MedicalDegree}])$   
and  $KB \models (\text{Surgeon} \sqsubseteq \text{Doctor})$   
then  $KB \models (\text{Surgeon} \sqsubseteq [\text{EXISTS } 1 : \text{MedicalDegree}])$

This is a simple form of *strict* inheritance (*cf.* next chapter)

Also, as noted in computing satisfaction (e.g. with joe and canCorp), adding an assertion like  $(c \rightarrow e)$  to a KB can cause other assertions  $(c' \rightarrow e')$  to be entailed for other individuals.

This type of propagation is most interesting in applications where membership in classes is monitored and changes are significant.

# Extensions to the language

---

A number of extensions to the DL language have been considered in the literature:

- upper bounds on the number of fillers

[AND [EXISTS 2 :Child] [AT-MOST 3 :Child]]

opens the possibility of inconsistent concepts

- sets of individuals: [ALL :Child [ONE-OF wally theodore]]
- relating the role fillers: [SAME-AS :President :CEO]
- qualified number restriction: [EXISTS 2 :Child Female] vs.  
[AND [EXISTS 2 :Child] [ALL :Child Female]]
- complex (non-atomic) roles: [EXISTS 2 [RESTR :Child Female]]  
[ALL [RESTR :Child Female] Married] vs.  
[ALL :Child [AND Female Married]]

Each of these extensions adds extra complexity to the problem of calculating subsumption.

This topic will be explored for RESTR in Chapter 16.

# Some applications

---

Like production systems, description logics have been used in a number of sorts of applications:

- interface to a DB
  - relational DB, but DL can provide a nice higher level view of the data based on objects
- working memory for a production system
  - instead of a having rules to reason about a taxonomy and inheritance of properties, this part of the reasoning can come from a DL system
- assertion and classification for monitoring
  - incremental change to KB can be monitored with certain atomic concepts declared “critical”
- contradiction detection in configuration
  - for a DL that allows contradictory concepts, can alert the user when these are detected. This works well for incremental construction of a concept representing e.g. a configuration of a computer.

---

10.

# Inheritance

# Hierarchy and inheritance

---

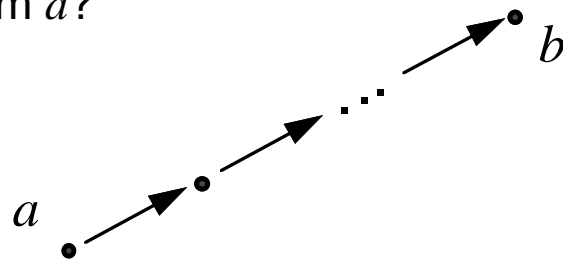
As we noticed with both frames and description logics, hierarchy or taxonomy is a natural way to view the world

importance of *abstraction* in remembering and reasoning

- groups of things share properties in the world
- do not have to repeat representations
  - e.g. sufficient to say that “elephants are mammals” to know a lot about them

Inheritance is the result of transitivity reasoning over paths in a network

- for strict networks, *modus ponens* (if-then reasoning) in graphical form
- “does  $a$  inherit from  $b$ ?” is the same as “is  $b$  in the transitive closure of **:IS-A** (or subsumption) from  $a$ ?”



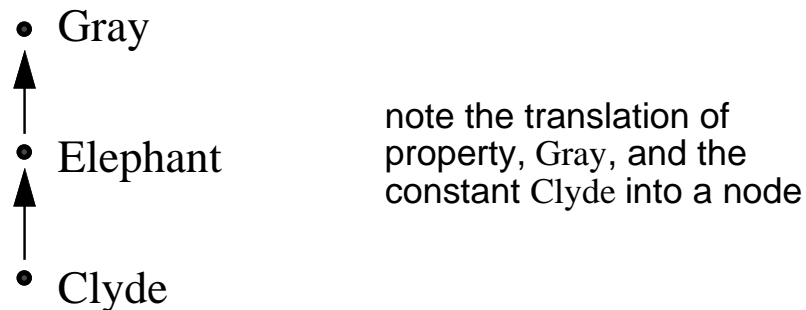
graphically, is there a path of **:IS-A** connections from  $a$  to  $b$ ?

# Path-based reasoning

---

## Focus just on inheritance and transitivity

- many interesting considerations in looking just at where information comes from in a network representation
- abstract frames/descriptions, and properties into nodes in graphs, and just look at reasoning with paths and the conclusions they lead us to

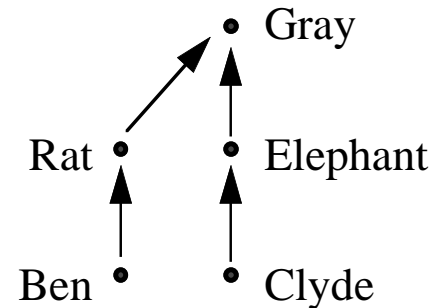


- edges in the network: Clyde·Elephant, Elephant·Gray
- paths included in this network: edges plus {Clyde·Elephant·Gray}  
in general, a path is a sequence of 1 or more edges
- conclusions supported by the paths:  
Clyde → Elephant; Elephant → Gray; Clyde → Gray

# Inheritance networks

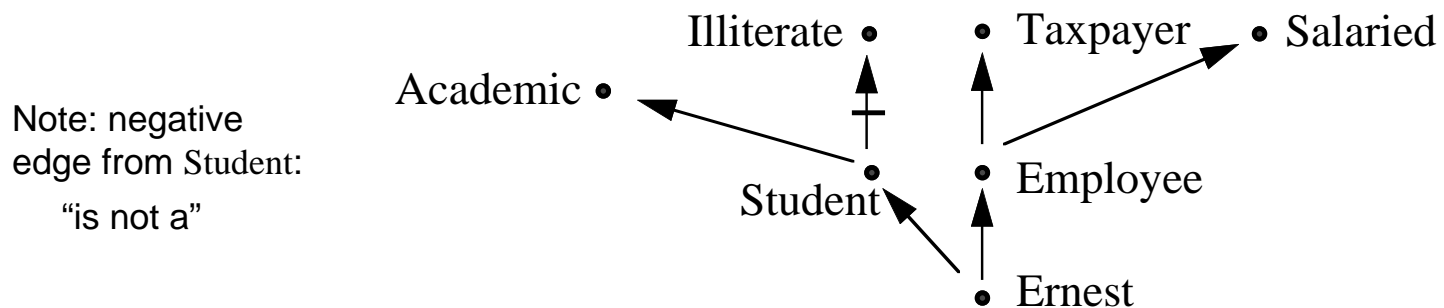
## (1) Strict inheritance in trees

- as in description logics
- conclusions produced by complete transitive closure on all paths (any traversal procedure will do); all reachable nodes are implied



## (2) Strict inheritance in DAGs

- as in DL's with multiple AND parents (= multiple inheritance)
- same as above: all conclusions you can reach by any paths are supported

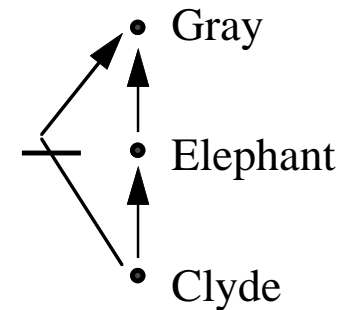


# Inheritance with defeasibility

## (3) Defeasible inheritance

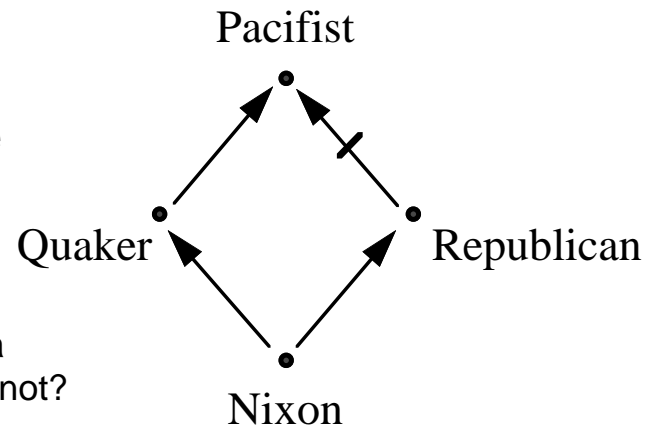
- as in frame systems
- inherited properties do not always hold, and can be *overridden* (defeated)
- conclusions determined by searching upward from “focus node” and selecting first version of property you want

while elephants in general are gray, Clyde is not



## A key problem: *ambiguity*

- *credulous* accounts choose arbitrarily
- *skeptical* accounts are more conservative



Is Nixon a pacifist or not?

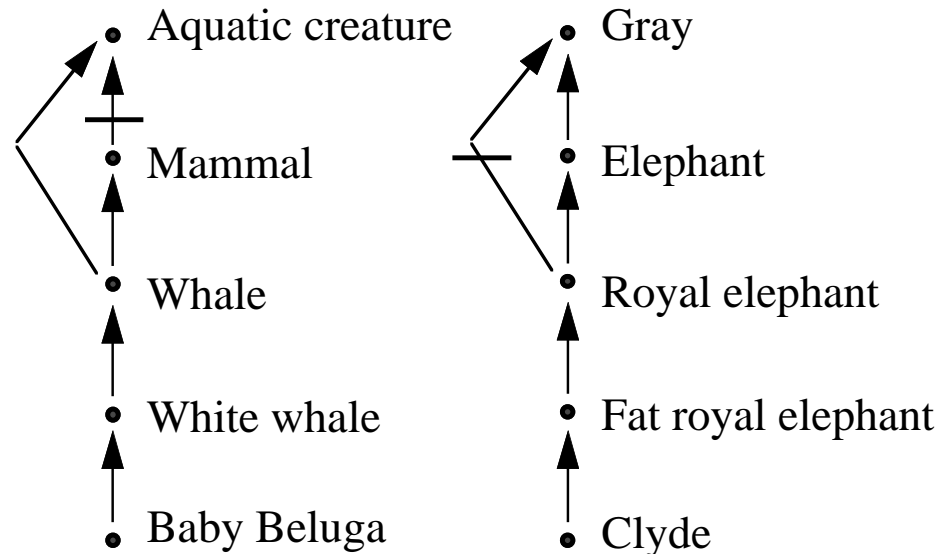


# Shortest path heuristic

## Defeasible inheritance in DAGs

- links have *polarity* (positive or negative)
- use shortest path heuristic to determine which polarity counts

Intuition: inherit from the most specific subsuming class



- as a result, not all paths count in generating conclusions
  - some are “preempted”
  - but some are “admissible”

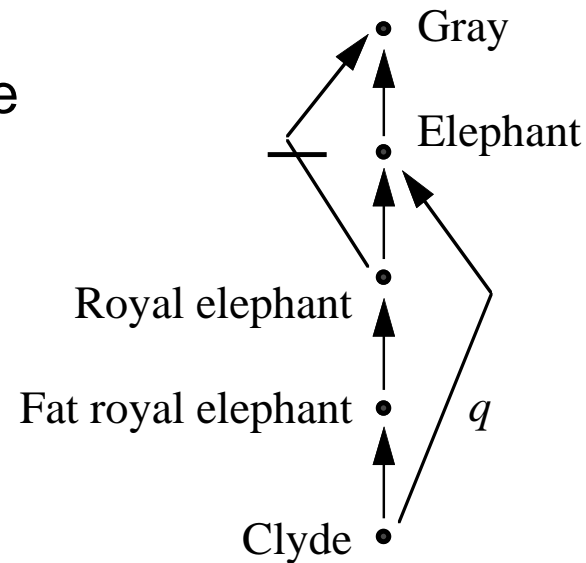
think of paths as *arguments* in support of conclusions

⇒ the inheritance problem = what are the admissible conclusions?

# Problems with shortest path

1. Shortest path heuristic produces incorrect answers in the presence of redundant edges (which are already implied!)

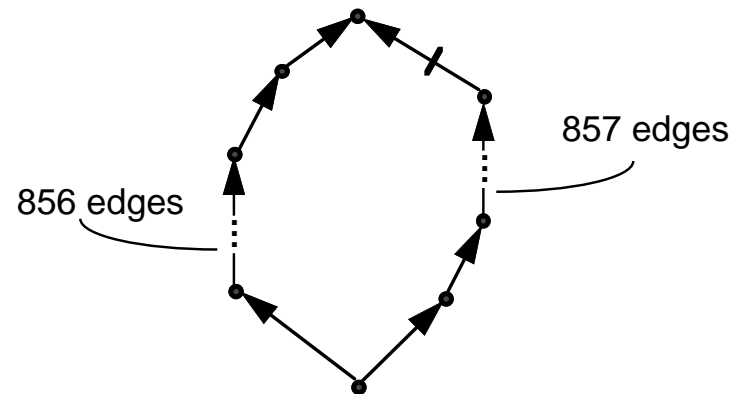
the redundant edge  $q$ ,  
expressing that Clyde is an  
Elephant changes polarity of  
conclusion about color



2. Anomalous behavior with ambiguity

adding 2 edges to the  
left side changes the  
conclusion!

Why should length be a factor?  
This network should be ambiguous...



# Specificity criteria

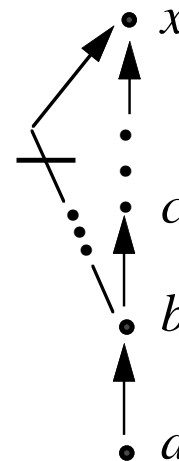
---

Shortest path is a specificity criterion (sometimes called a preemption strategy) which allows us to make admissibility choices among competing paths

- It's not the only possible one
- Consider "*inferential distance*":  
not linear distance, but topologically based
  - a node  $a$  is nearer to node  $b$  than to node  $c$  if there is a path from  $a$  to  $c$  through  $b$
  - idea: conclusions from  $b$  preempt those from  $c$

This handles  $\text{Clyde} \rightarrow \neg\text{Gray}$  just fine,  
as well as redundant links

- But what if path from  $b$  to  $c$  has some of its edges preempted? what if some are redundant?



# A formalization (Stein)

---

An inheritance hierarchy  $\Gamma = \langle V, E \rangle$  is a directed, acyclic graph (DAG) with positive and negative edges, intended to denote “(normally) is-a” and “(normally) is-not-a”, respectively.

- positive edges are written  $a \cdot x$
- negative edges are written  $a \cdot \neg x$

A sequence of edges is a path:

- a positive path is a sequence of one or more positive edges  $a \cdot \dots \cdot x$
- a negative path is a sequence of positive edges followed by a single negative edge  $a \cdot \dots \cdot v \cdot \neg x$

Note: there are no paths with more than 1 negative edge.

Also: there might be 0 positive edges.

A path (or argument) supports a conclusion:

- $a \cdot \dots \cdot x$  supports the conclusion  $a \rightarrow x$  ( $a$  is an  $x$ )
- $a \cdot \dots \cdot \neg x$  supports  $a \not\rightarrow x$  ( $a$  is not an  $x$ )

Note: a conclusion may be supported by many arguments

However: not all arguments are equally believable...

# Support and admissibility

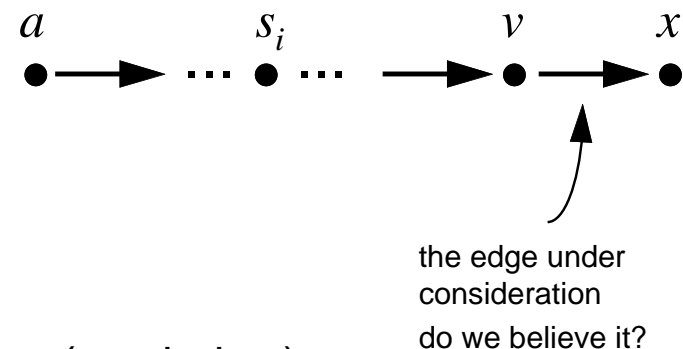
$\Gamma$  supports a path  $a \cdot s_1 \cdot \dots \cdot s_n \cdot (\neg)x$  if the corresponding set of edges  $\{a \cdot s_1, \dots, s_n \cdot (\neg)x\}$  is in  $E$ , and the path is admissible according to specificity (see below).

the hierarchy supports a conclusion  $a \rightarrow x$  (or  $a \not\rightarrow x$ )  
if it supports some corresponding path

A path is admissible if every edge in it is admissible.

An edge  $v \cdot x$  is admissible in  $\Gamma$  wrt  $a$  if there is a positive path  $a \cdot s_1 \dots s_n \cdot v$  ( $n \geq 0$ ) in  $E$  and

1. each edge in  $a \cdot s_1 \dots s_n \cdot v$  is admissible in  $\Gamma$  wrt  $a$  (recursively);
2. no edge in  $a \cdot s_1 \dots s_n \cdot v$  is redundant in  $\Gamma$  wrt  $a$  (see below);
3. no intermediate node  $a, s_1, \dots, s_n$  is a preemptor of  $v \cdot x$  wrt  $a$  (see below).

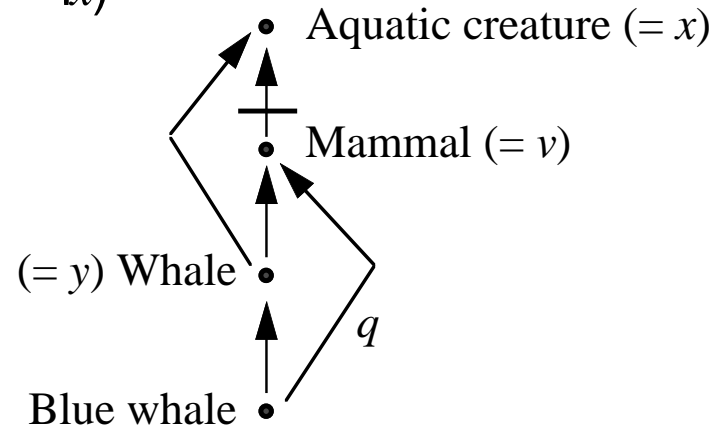


A negative edge  $v \cdot \neg x$  is handled analogously.

# Preemption and redundancy

A node  $y$  along path  $a \cdot \dots \cdot y \cdot \dots \cdot v$  is a preemptor of the edge  $v \cdot x$  wrt  $a$  if  $y \cdot \neg x \in E$  (and analogously for  $v \cdot \neg x$ )

for example, in this figure  
the node Whale preempts  
the negative edge from  
Mammal to Aquatic creature  
wrt both Whale and Blue whale



A positive edge  $b \cdot w$  is redundant in  $\Gamma$  wrt node  $a$  if there is some positive path  $b \cdot t_1 \cdot \dots \cdot t_m \cdot w \in E$  ( $m \geq 1$ ), for which

1. each edge in  $b \cdot t_1 \cdot \dots \cdot t_m$  is admissible in  $\Gamma$  wrt  $a$ ;
2. there are no  $c$  and  $i$  such that  $c \cdot \neg t_i$  is admissible in  $\Gamma$  wrt  $a$ ;
3. there is no  $c$  such that  $c \cdot \neg w$  is admissible in  $\Gamma$  wrt  $a$ .

The edge labelled  $q$  above is redundant

The definition for a negative edge  $b \cdot \neg w$  is analogous

# Credulous extensions

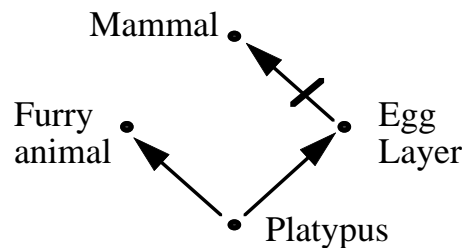
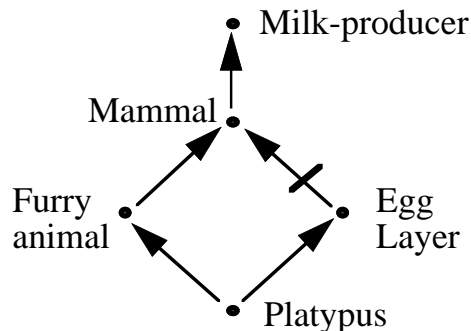
$\Gamma$  is *a*-connected iff for every node  $x$  in  $\Gamma$ , there is a path from  $a$  to  $x$ , and for every edge  $v \cdot (\neg) x$  in  $\Gamma$ , there is a *positive* path from  $a$  to  $v$ .

In other words, every node and edge is reachable from  $a$

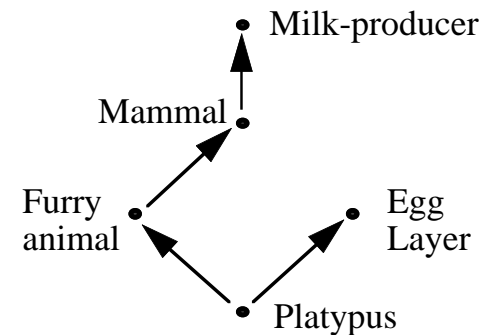
$\Gamma$  is (potentially) ambiguous wrt a node  $a$  if there is some node  $x \in V$  such that both  $a \cdot s_1 \dots s_n \cdot x$  and  $a \cdot t_1 \dots t_m \cdot \neg x$  are paths in  $\Gamma$

A credulous extension of  $\Gamma$  wrt node  $a$  is a maximal unambiguous *a*-connected subhierarchy of  $\Gamma$  wrt  $a$

If  $X$  is a credulous extension of  $\Gamma$ , then adding an edge of  $\Gamma$  to  $X$  makes  $X$  either ambiguous or not *a*-connected



**Extension 1**



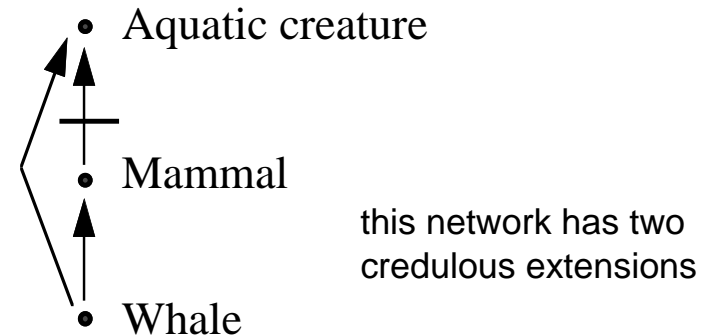
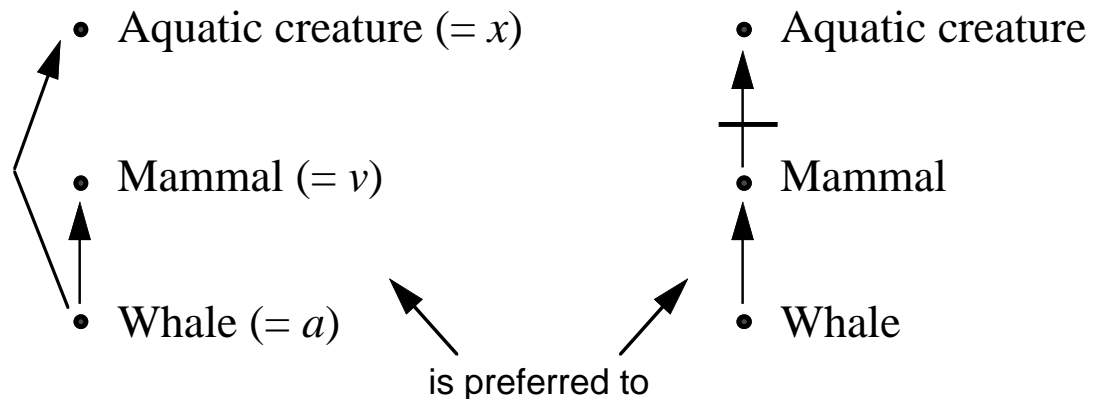
**Extension 2**

# Preferred extensions

Credulous extensions do not incorporate any notion of admissibility or preemption.

Let  $X$  and  $Y$  be credulous extensions of  $\Gamma$  wrt node  $a$ .  $X$  is preferred to  $Y$  iff there are nodes  $v$  and  $x$  such that:

- $X$  and  $Y$  agree on all edges whose endpoints precede  $v$  topologically,
- there is an edge  $v \cdot x$  (or  $v \cdot \neg x$ ) that is *inadmissible* in  $\Gamma$ ,
- this edge is in  $Y$ , but not in  $X$ .



A credulous extension is a preferred extension if there is no other extension that is preferred to it.



# Subtleties

---

## What to believe?

- “credulous” reasoning: choose a preferred extension and believe all the conclusions supported
- “skeptical” reasoning: believe the conclusions from any path that is supported by all preferred extensions
- “ideally skeptical” reasoning: believe the conclusions that are supported by all preferred extensions

note: ideally skeptical reasoning cannot be computed in a path-based way (conclusions may be supported by different paths in each extension)

## We’ve been doing “upwards” reasoning

- start at a node and see what can be inherited from its ancestor nodes
- there are many variations on this definition; none has emerged as the agreed upon, or “correct” one
- an alternative looks from the top and sees what propagates down  
upwards is more efficient

---

11.

# Defaults

# Strictness of FOL

---

To reason from  $P(a)$  to  $Q(a)$ , need either

- facts about  $a$  itself
- universals, e.g.  $\forall x(P(x) \supset Q(x))$ 
  - something that applies to all instances
  - all or nothing!

But most of what we learn about the world is in terms of generics

e.g., encyclopedia entries for ferris wheels, violins, turtles, wildflowers

Properties are not strict for all instances, because

- genetic / manufacturing varieties
  - early ferris wheels
- cases in exceptional circumstances
  - dried wildflowers
- borderline cases
  - toy violins
- imagined cases
  - flying turtles

*etc.*

# Generics vs. universals

---

✓ Violins have four strings.

vs.

✗ All violins have four strings.

vs.

? All violins that are not  $E_1$  or  $E_2$  or ... have four strings.

(exceptions usually cannot be enumerated)

Similarly, for general properties of individuals

- Alexander the great: ruthlessness
- Ecuador: exports
- pneumonia: treatment

Goal: be able to say a  $P$  is a  $Q$  in general, but not necessarily

It is reasonable to conclude  $Q(a)$  given  $P(a)$ ,  
unless there is a good reason not to

Here: qualitative version (no numbers)

# Varieties of defaults (I)

---

## General statements

- prototypical: The prototypical  $P$  is a  $Q$ .  
Owls hunt at night.
- normal: Under typical circumstances,  $P$ 's are  $Q$ 's.  
People work close to where they live.
- statistical: Most  $P$ 's are  $Q$ 's.  
The people in the waiting room are growing impatient.

## Lack of information to the contrary

- group confidence: All known  $P$ 's are  $Q$ 's.  
Natural languages are easy for children to learn.
- familiarity: If a  $P$  was not a  $Q$ , you would know it.
  - an older brother
  - very unusual individual, situation or event

# Varieties of defaults (II)

---

## Conventional

- conversational: Unless I tell you otherwise, a  $P$  is a  $Q$   
“There is a gas station two blocks east.”  
the default: the gas station is open.
- representational: Unless otherwise indicated, a  $P$  is a  $Q$   
the speed limit in a city

## Persistence

- inertia: A  $P$  is a  $Q$  if it used to be a  $Q$ .
  - colours of objects
  - locations of parked cars (for a while!)

Here: we will use “Birds fly” as a typical default.

# Closed-world assumption

---

Reiter's observation:

There are usually many more -ve facts than +ve facts!

Example: airline flight guide provides

DirectConnect(cleveland,toronto)      DirectConnect(toronto,northBay)  
DirectConnect(toronto,winnipeg)      ...

but not:  $\neg$ DirectConnect(cleveland,northBay)

Conversational default, called CWA:

only +ve facts will be given, relative to some vocabulary

But note:  $KB \not\models$  -ve facts (would have to answer: "I don't know")

Proposal: a new version of entailment:  $KB \models_c \alpha$  iff  $KB \cup Negs \models \alpha$

where  $Negs = \{\neg p \mid p \text{ atomic and } KB \not\models p\}$

Note: relation to negation as failure

a common pattern:

$KB' = KB \cup \Delta$

Gives:  $KB \models_c$  +ve facts and -ve facts

# Properties of CWA

---

For every  $\alpha$  (without quantifiers),  $KB \models_c \alpha$  or  $KB \models_c \neg\alpha$

Why? Inductive argument:

- immediately true for atomic sentences
- push  $\neg$  in, e.g.  $KB \models \neg\neg\alpha$  iff  $KB \models \alpha$
- $KB \models (\alpha \wedge \beta)$  iff  $KB \models \alpha$  and  $KB \models \beta$
- Say  $KB \not\models_c (\alpha \vee \beta)$ . Then  $KB \not\models_c \alpha$  and  $KB \not\models_c \beta$ .  
So by induction,  $KB \models_c \neg\alpha$  and  $KB \models_c \neg\beta$ . Thus,  $KB \models_c \neg(\alpha \vee \beta)$ .

CWA is an assumption about complete knowledge

never any unknowns, relative to vocabulary

In general, a KB has incomplete knowledge,

e.g. Let KB be  $(p \vee q)$ . Then  $KB \models (p \vee q)$ ,  
but  $KB \not\models p$ ,  $KB \not\models \neg p$ ,  $KB \not\models q$ ,  $KB \not\models \neg q$

With CWA, have: If  $KB \models_c (\alpha \vee \beta)$ , then  $KB \models_c \alpha$  or  $KB \models_c \beta$ .

similar argument to above



# Query evaluation

---

With CWA can reduce queries (without quantifiers) to the atomic case:

$KB \models_c (\alpha \wedge \beta)$  iff  $KB \models_c \alpha$  and  $KB \models_c \beta$

$KB \models_c (\alpha \vee \beta)$  iff  $KB \models_c \alpha$  or  $KB \models_c \beta$

$KB \models_c \neg(\alpha \wedge \beta)$  iff  $KB \models_c \neg\alpha$  or  $KB \models_c \neg\beta$

$KB \models_c \neg(\alpha \vee \beta)$  iff  $KB \models_c \neg\alpha$  and  $KB \models_c \neg\beta$

$KB \models_c \neg\neg\alpha$  iff  $KB \models_c \alpha$

reduces to:  $KB \models_c \rho$ , where  $\rho$  is a literal

If  $KB \cup Negs$  is consistent, get  $KB \models_c \neg\alpha$  iff  $KB \not\models_c \alpha$

reduces to:  $KB \models_c p$ , where  $p$  is atomic

If atoms stored as a table, deciding if  $KB \models_c \alpha$  is like DB-retrieval:

- reduce query to set of atomic queries
- solve atomic queries by table lookup

Different from ordinary logic reasoning (e.g. no reasoning by cases)

# Consistency of CWA

---

If KB is a set of atoms, then  $KB \cup Negs$  is always consistent

Also works if KB has conjunctions and if KB has only negative disjunctions

If KB contains  $(\neg p \vee \neg q)$ . Add both  $\neg p, \neg q$ .

Problem when  $KB \models (\alpha \vee \beta)$ , but  $KB \not\models \alpha$  and  $KB \not\models \beta$

e.g.  $KB = (p \vee q)$   $Negs = \{\neg p, \neg q\}$

$KB \cup Negs$  is inconsistent and so for every  $\alpha$ ,  $KB \models_c \alpha$  !

Solution: only apply CWA to atoms that are “uncontroversial”

One approach: GCWA

$Negs = \{\neg p \mid \text{If } KB \models (p \vee q_1 \vee \dots \vee q_n) \text{ then } KB \models (q_1 \vee \dots \vee q_n)\}$

When KB is consistent, get:

- $KB \cup Negs$  consistent
- everything derivable is also derivable by CWA

# Quantifiers and equality

---

So far, results do not extend to wffs with quantifiers

can have  $\text{KB} \not\models_c \forall x.\alpha$  and  $\text{KB} \not\models_c \neg\forall x.\alpha$

e.g. just because for every  $t$ , we have  $\text{KB} \models_c \neg\text{DirectConnect}(\text{myHome}, t)$   
does not mean that  $\text{KB} \models_c \forall x[\neg\text{DirectConnect}(\text{myHome}, x)]$

But may want to treat KB as providing complete information about what individuals exist

Define:  $\text{KB} \models_{cd} \alpha$  iff  $\text{KB} \cup \text{Negs} \cup \text{Dc} \models \alpha$  where the  $c_i$  are all the constants appearing in KB (assumed finite)

where  $\text{Dc}$  is domain closure:  $\forall x[x=c_1 \vee \dots \vee x=c_n]$ ,

Get:  $\text{KB} \models_{cd} \exists x.\alpha$  iff  $\text{KB} \models_{cd} \alpha[x/c]$ , for some  $c$  appearing in the KB  
 $\text{KB} \models_{cd} \forall x.\alpha$  iff  $\text{KB} \models_{cd} \alpha[x/c]$ , for all  $c$  appearing in the KB

Then add:  $\text{Un}$  is unique names:  $(c_i \neq c_j)$ , for  $i \neq j$

Get:  $\text{KB} \models_{cdu} (c = d)$  iff  $c$  and  $d$  are the same constant

→ full recursive query evaluation

# Non-monotonicity

---

Ordinary entailment is monotonic

If  $KB \models \alpha$ , then  $KB^* \models \alpha$ , for any  $KB \subseteq KB^*$

But CWA entailment is *not* monotonic

Can have  $KB \models_c \alpha$ ,  $KB \subseteq KB'$ , but  $KB' \not\models_c \alpha$

e.g.  $\{p\} \models_c \neg q$ , but  $\{p, q\} \not\models_c \neg q$

Suggests study of non-monotonic reasoning

- start with explicit beliefs
- generate implicit beliefs non-monotonically, taking *defaults* into account
- implicit beliefs may not be uniquely determined (vs. monotonic case)

Will consider three approaches:


- minimal entailment: interpretations that minimize abnormality
- default logic: KB as facts + default rules of inference
- autoepistemic logic: facts that refer to what is/is not believed

# Minimizing abnormality

---

CWA makes the extension of all predicates as small as possible  
by adding negated literals

Generalize: do this only for selected predicates  
Ab predicates used to talk about typical cases

Example KB: Bird(chilly),  $\neg$ Flies(chilly),  
Bird(tweety), (chilly  $\neq$  tweety),  
 $\forall x[\text{Bird}(x) \wedge \neg \text{Ab}(x) \supset \text{Flies}(x)]$   All birds that  
are normal fly

Would like to conclude by default Flies(tweety), but  $\text{KB} \not\models \text{Flies}(tweety)$   
because there is an interpretation  $\mathcal{I}$  where  $\mathcal{I}[\text{tweety}] \in \mathcal{I}[\text{Ab}]$

Solution: consider only interpretations where  
 $\mathcal{I}[\text{Ab}]$  is as small as possible, relative to KB

for example: KB requires that  $\mathcal{I}[\text{chilly}] \in \mathcal{I}[\text{Ab}]$

this is sometimes  
called “circumscription”  
since we circumscribe  
the Ab predicate

Generalizes to many  $\text{Ab}_i$  predicates

# Minimal entailment

---

Given two interps over the same domain,  $\mathcal{I}_1$  and  $\mathcal{I}_2$

$\mathcal{I}_1 \leq \mathcal{I}_2$  iff  $I_1[Ab] \subseteq I_2[Ab]$  for every Ab predicate

$\mathcal{I}_1 < \mathcal{I}_2$  iff  $\mathcal{I}_1 \leq \mathcal{I}_2$  but not  $\mathcal{I}_2 \leq \mathcal{I}_1$  read:  $\mathcal{I}_1$  is more normal than  $\mathcal{I}_2$

Define a new version of entailment,  $\models_{\leq}$  by

$KB \models_{\leq} \alpha$  iff for every  $\mathcal{I}$ , if  $\mathcal{I} \models KB$  and no  $\mathcal{I}^* < \mathcal{I}$  s.t.  $\mathcal{I}^* \models KB$   
then  $\mathcal{I} \models \alpha$ .

So  $\alpha$  must be true in all interps satisfying KB that are *minimal* in abnormalities

Get:  $KB \models_{\leq} \text{Flies}(\text{tweety})$

because if interp satisfies KB and is minimal, only  $I[\text{chilly}]$  will be in  $I[Ab]$

Minimization need not produce a *unique* interpretation:

$\text{Bird}(a), \text{Bird}(b), [\neg \text{Flies}(a) \vee \neg \text{Flies}(b)]$  yields two minimal interpretations

$KB \not\models_{\leq} \text{Flies}(a), KB \not\models_{\leq} \text{Flies}(b), KB \models_{\leq} \text{Flies}(a) \vee \text{Flies}(b)$

Different from the CWA: no inconsistency!

But stronger than GCWA: conclude a or b flies

# Fixed and variable predicates

---

Imagine KB as before +  $\forall x[\text{Penguin}(x) \supset \text{Bird}(x) \wedge \neg\text{Flies}(x)]$

Get:  $\text{KB} \models \forall x[\text{Penguin}(x) \supset \text{Ab}(x)]$

So minimizing Ab also minimizes penguins:  $\text{KB} \models_{\leq} \forall x \neg\text{Penguin}(x)$

McCarthy's definition: Let **P** and **Q** be sets of predicates

$\mathcal{S}_1 \leq \mathcal{S}_2$  iff same domain and

1.  $I_1[P] \subseteq I_2[P]$ , for every  $P \in \mathbf{P}$       Ab predicates
2.  $I_1[Q] = I_2[Q]$ , for every  $Q \notin \mathbf{Q}$       fixed predicates

so only predicates in **Q** are allowed to vary

Get definition of  $\models_{\leq}$  that is parameterized by what is minimized *and* what is allowed to vary

Previous example: minimize Ab, but allow only Flies to vary.

Problems:

- need to decide what to allow to vary
- cannot conclude  $\neg\text{Penguin}(\text{tweety})$  by default!

only get default ( $\neg\text{Penguin}(\text{tweety}) \supset \text{Flies}(\text{tweety})$ )

# Default logic

---

Beliefs as deductive theory

explicit beliefs = axioms

implicit beliefs = theorems = least set closed under inference rules

e.g. If we can prove  $\alpha$  and  $(\alpha \supset \beta)$ , then infer  $\beta$

Would like to generalize to default rules:

If can prove  $Bird(x)$ , but *cannot* prove  $\neg Flies(x)$ , then infer  $Flies(x)$ .

Problem: how to characterize theorems

cannot write a derivation, since do not know when to apply default rules

no guarantee of unique set of theorems

If cannot infer  $p$ , infer  $q$  + If cannot infer  $q$ , infer  $p$  ??

Solution: default logic

no notion of theorem

instead, have extensions: sets of sentences that are “reasonable” beliefs, given explicit facts and default rules



# Extensions

---

Default logic KB uses two components:  $KB = \langle F, D \rangle$

- $F$  is a set of sentences (facts)
- $D$  is a set of default rules: triples  $\langle \alpha : \beta / \gamma \rangle$  read as

If you can infer  $\alpha$ , and  $\beta$  is *consistent*, then infer  $\gamma$

$\alpha$ : the prerequisite,  $\beta$ : the justification,  $\gamma$ : the conclusion

e.g.  $\langle \text{Bird}(\text{tweety}) : \text{Flies}(\text{tweety}) / \text{Flies}(\text{tweety}) \rangle$

treat  $\langle \text{Bird}(x) : \text{Flies}(x) / \text{Flies}(x) \rangle$  as set of rules

Default rules where  $\beta = \gamma$  are called normal and write as  $\langle \alpha \Rightarrow \beta \rangle$

will see later a reason for wanting non-normal ones

A set of sentences  $E$  is an extension of  $\langle F, D \rangle$  iff for every sentence  $\pi$ ,  $E$  satisfies the following:

$$\pi \in E \text{ iff } F \cup \Delta \models \pi, \text{ where } \Delta = \{ \gamma \mid \langle \alpha : \beta / \gamma \rangle \in D, \alpha \in E, \neg \beta \notin E \}$$

So, an extension  $E$  is the set of entailments of  $F \cup \{ \gamma \}$ , where the  $\gamma$  are assumptions from  $D$ .

to check if  $E$  is an extension, guess at  $\Delta$  and show that it satisfies the above constraint

# Example

---

Suppose KB has

$F = \text{Bird}(\text{chilly}), \neg\text{Flies}(\text{chilly}), \text{Bird}(\text{tweety})$

$D = \langle \text{Bird}(x) \Rightarrow \text{Flies}(x) \rangle$

then there is a unique extension, where  $\Delta = \text{Flies}(\text{tweety})$

- This is an extension since *tweety* is the only *t* for this  $\Delta$  such that  $\text{Bird}(t) \in E$  and  $\neg\text{Flies}(t) \notin E$ .
- No other extension, since this applies no matter what  $\text{Flies}(t)$  assumptions are in  $\Delta$ .

But in general can have multiple extensions:

$F = \{ \text{Republican}(\text{dick}), \text{Quaker}(\text{dick}) \}$       $D = \{ \langle \text{Republican}(x) \Rightarrow \neg\text{Pacifist}(x) \rangle, \langle \text{Quaker}(x) \Rightarrow \text{Pacifist}(x) \rangle \}$

Two extensions:  $E_1$  has  $\Delta = \neg\text{Pacifist}(\text{dick})$ ;      $E_2$  has  $\Delta = \text{Pacifist}(\text{dick})$

Which to believe?

credulous: choose an extension arbitrarily

skeptical: believe what is common to all extensions

Can sometimes use non-normal defaults to avoid conflicts in defaults

$\langle \text{Quaker}(x) : \text{Pacifist}(x) \wedge \neg\text{Republican}(x) / \text{Pacifist}(x) \rangle$

but then need to consider all possible interactions in defaults!

# Unsupported conclusions

---

Extension tries to eliminate facts that do not result from either  $F$  or  $D$ .

e.g., we do not want  $\text{Yellow}(\text{tweety})$  and its entailments in the extension

But the definition has a problem:

Suppose  $F = \{\}$  and  $D = \langle p : \text{True} / p \rangle$ .

Then  $E = \text{entailments of } \{p\}$  is an extension

since  $p \in E$  and  $\neg\text{True} \notin E$ , for above default

However, no good reason to believe  $p$ !

Only support for  $p$  is default rule, which requires  $p$  itself as a prerequisite

So default should have no effect. Want one extension:  $E = \text{entailments of } \{\}$

Reiter's definition:

For any set  $S$ , let  $\Gamma(S)$  be the least set containing  $F$ , closed under entailment, and satisfying

if  $\langle \alpha : \beta / \gamma \rangle \in D$ ,  $\alpha \in \Gamma(S)$ , and  $\neg\beta \notin S$ , then  $\gamma \in \Gamma(S)$ .

A set  $E$  is an extension of  $\langle F, D \rangle$  iff  $E = \Gamma(E)$ .

← note: not  $\Gamma(S)$

called a fixed point of the  $\Gamma$  operator

# Autoepistemic logic

---

One disadvantage of default logic is that rules cannot be combined or reasoned about

$$\langle \alpha : \beta / \gamma \rangle \rightsquigarrow \langle \alpha : \beta / (\gamma \vee \delta) \rangle$$

Solution: express defaults as *sentences* in an extended language that talks about belief explicitly

for any sentence  $\alpha$ , we have another sentence  $\mathbf{B}\alpha$

$\mathbf{B}\alpha$  says "I believe  $\alpha$ ": autoepistemic logic

e.g.  $\forall x[\text{Bird}(x) \wedge \neg \mathbf{B}\neg \text{Flies}(x) \supset \text{Flies}(x)]$

All birds fly except those that I believe to not fly =

Any bird not believed to be flightless flies.

No longer expressing defaults using formulas of FOL.

# Semantics of belief

---

These are not sentences of FOL, so what semantics and entailment?

- modal logic of belief provide semantics
- for here: treat  $\mathbf{B}\alpha$  as if it were an new atomic wff
- still get entailment:  $\forall x[\text{Bird}(x) \wedge \neg\mathbf{B}\neg\text{Flies}(x) \supset \text{Flies}(x) \vee \text{Run}(x)]$

Main property for set of implicit beliefs,  $E$ :

1. If  $E \models \alpha$  then  $\alpha \in E$ . (closed under entailment)
2. If  $\alpha \in E$  then  $\mathbf{B}\alpha \in E$ . (positive introspection)
3. If  $\alpha \notin E$  then  $\neg\mathbf{B}\alpha \in E$ . (negative introspection)

Any such set of sentences is called stable

Note: if  $E$  contains  $p$  but does not contain  $q$ , it will contain  $\mathbf{B}p$ ,  $\mathbf{B}\mathbf{B}p$ ,  $\mathbf{B}\mathbf{B}\mathbf{B}p$ ,  $\neg\mathbf{B}q$ ,  $\mathbf{B}\neg\mathbf{B}q$ ,  $\mathbf{B}(\mathbf{B}p \wedge \neg\mathbf{B}q)$ , etc.

# Stable expansions

---

Given KB, possibly containing **B** operators, our implicit beliefs should be a stable set that is minimal.

Moore's definition: A set of sentences  $E$  is called a stable expansion of KB iff it satisfies the following:

$$\pi \in E \text{ iff } \text{KB} \cup \Delta \models \pi, \text{ where } \Delta = \{\mathbf{B}\alpha \mid \alpha \in E\} \cup \{\neg\mathbf{B}\alpha \mid \alpha \notin E\}$$

fixed point of another operator

analogous to the extensions of default logic

Example: for  $\text{KB} = \{ \text{Bird}(\text{chilly}), \neg\text{Flies}(\text{chilly}), \text{Bird}(\text{tweety}), \forall x[\text{Bird}(x) \wedge \neg\mathbf{B}\neg\text{Flies}(x) \supset \text{Flies}(x)] \}$

get a unique stable expansion containing  $\text{Flies}(\text{tweety})$

As in default logic, stable expansions are not uniquely determined

$$\text{KB} = \{(\neg\mathbf{B}p \supset q), (\neg\mathbf{B}q \supset p)\}$$

2 stable expansions  
(one with  $p$ , one with  $q$ )

$$\text{KB} = \{(\neg\mathbf{B}p \supset p)\} \quad (\text{self-defeating default})$$

no stable expansions!  
so what to believe?

# Enumerating stable expansions

---

Define: A wff is objective if it has no **B** operators

When a KB is propositional, and **B** operators only dominate objective wffs, we can enumerate all stable expansions using the following:

1. Suppose  $\mathbf{B}\alpha_1, \mathbf{B}\alpha_2, \dots, \mathbf{B}\alpha_n$  are all the **B** wffs in KB.
2. Replace some of these by True and the rest by  $\neg$ True in KB and simplify.  
Call the result  $\text{KB}^\circ$  (it's objective).  
at most  $2^n$  possible replacements
3. Check that for each  $\alpha_i$ ,
  - if  $\mathbf{B}\alpha_i$  was replaced by True, then  $\text{KB}^\circ \models \alpha_i$
  - if  $\mathbf{B}\alpha_i$  was replaced by  $\neg$ True, then  $\text{KB}^\circ \not\models \alpha_i$
4. If yes, then  $\text{KB}^\circ$  determines a stable expansion.  
entailments of  $\text{KB}^\circ$  are the objective part

# Example enumeration

---

For  $KB = \{ \text{Bird}(\text{chilly}), \neg\text{Flies}(\text{chilly}), \text{Bird}(\text{tweety}),$   
 $[\text{Bird}(\text{tweety}) \wedge \neg\mathbf{B}\neg\text{Flies}(\text{tweety}) \supset \text{Flies}(\text{tweety})],$   
 $[\text{Bird}(\text{chilly}) \wedge \neg\mathbf{B}\neg\text{Flies}(\text{chilly}) \supset \text{Flies}(\text{chilly})] \}$

Two  $\mathbf{B}$  wffs:  $\mathbf{B}\neg\text{Flies}(\text{tweety})$  and  $\mathbf{B}\neg\text{Flies}(\text{chilly})$ ,  
so four replacements to try.

Only one satisfies the required constraint:

$\mathbf{B}\neg\text{Flies}(\text{tweety}) \rightarrow \neg\text{True},$   
 $\mathbf{B}\neg\text{Flies}(\text{chilly}) \rightarrow \text{True}$

Resulting  $KB^\circ$  has

$(\text{Bird}(\text{tweety}) \supset \text{Flies}(\text{tweety}))$

and so entails

$\text{Flies}(\text{tweety})$

as desired.



# More ungroundedness

---

Definition of stable expansion may not be strong enough

$KB = \{(\mathbf{B}p \supset p)\}$  has 2 stable expansions:

- one without  $p$  and with  $\neg\mathbf{B}p$   
corresponds to  $KB^\circ = \{\}$
- one with  $p$  and  $\mathbf{B}p$ .  
corresponds to  $KB^\circ = \{p\}$

But why should  $p$  be believed?

only justification for having  $p$  is having  $\mathbf{B}p$ !

similar to problem with default logic extension

Konolige's definition:

A grounded stable expansion is a stable expansion that is minimal wrt to the set of sentences without  $\mathbf{B}$  operators.

rules out second stable expansion

Other examples suggest that an even stronger definition is required!

can get an equivalence with Reiter's definition of extension in default logic

---

12.

# Vagueness, Uncertainty and Degrees of Belief

# Noncategorical statements

---

Ordinary commonsense knowledge quickly moves away from categorical statements like “a  $P$  is *always (unequivocably)* a  $Q$ ”

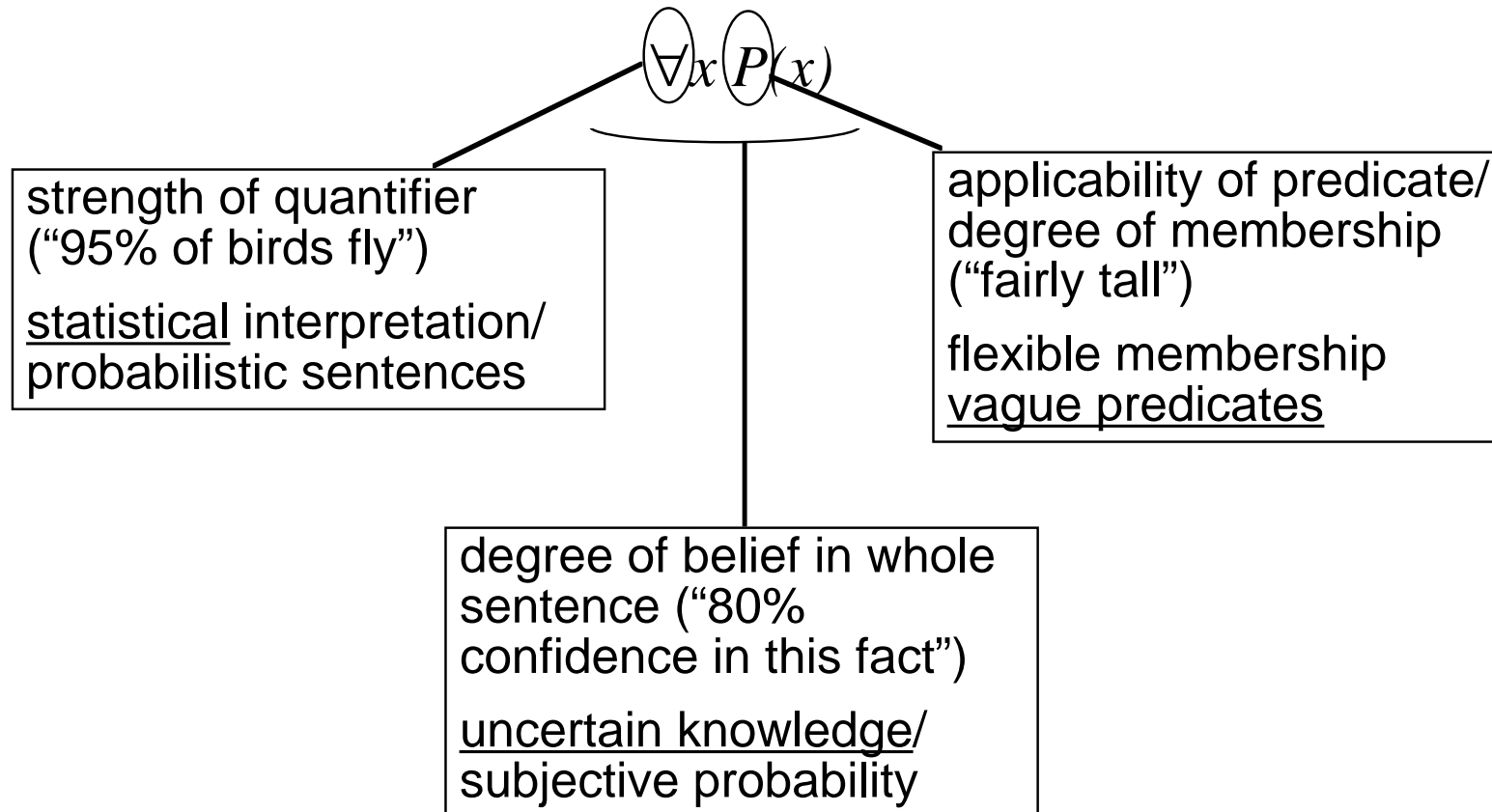
There are many ways in which we can come to less than categorical information

- things are *usually (almost never, occasionally, seldomly, rarely, almost always)* a certain way
- judgments about how good an example something is  
e.g., barely rich, a poor example of a chair, not very tall
- imprecision of sensors  
e.g., the best you can do is to get within +/-10%
- reliability of sources of information  
e.g., “most of the time he’s right on the money”
- strength/confidence/trust in generic information or deductive rules

Conclusions will not “follow” in the usual sense

# Weakening a universal

There are at least 3 ways a universal like  $\forall x P(x)$  can be made to be less categorical:



# Objective probability

---

## Statistical (frequency) view of sentences

objective: does not depend on who is assessing the probability

## Always applied to collections

can not assign probabilities to (random) events that are not members of any obvious repeatable sequence:

- ok for “the probability that I will pick a red face card from the deck”
- not ok for “the probability that the Blue Jays will win the World Series this Fall”
- “the probability that Tweety flies is between .9 and .95” is always false (either Tweety flies or not)

Can use probabilities to correspond to English words like “rarely,” “likely,” “usually”

generalized quantifiers: “most,” “many,” “few”

For most  $x$ ,  $Q(x)$  vs. For all  $x$ ,  $Q(x)$

# The basic postulates

---

Numbers between 0 and 1 representing frequency of an event in a (large enough) random sample

extremes: 0 = never happens; 1 = always happens

Start with set  $U$  of all possible occurrences. An event  $a$  is any subset of  $U$ . A probability measure is any function  $Pr$  from events to  $[0,1]$  satisfying:

- $Pr(U) = 1$ .
- If  $a_1, \dots, a_n$  are disjoint events, then  $Pr(\cup a_i) = \sum Pr(a_i)$

Conditioning: the probability of one event may depend on its interaction with others

$$Pr(a/b) = \text{probability of } a, \text{ given } b = Pr(a \cap b) / Pr(b)$$

Conditional independence:

event  $a$  is judged independent of event  $b$  conditional on background knowledge  $s$  if knowing that  $b$  happened does not affect the probability of  $a$

$$Pr(a/s) = Pr(a/b,s) \quad (\text{note: CI is symmetric})$$

Note: without independence,  $Pr(a/s)$  and  $Pr(a/b,s)$  can be very different.

# Some useful consequences

---

Conjunction:

$$\mathbf{Pr}(ab) = \mathbf{Pr}(a/b) \cdot \mathbf{Pr}(b)$$

conditionally independent:  $\mathbf{Pr}(ab) = \mathbf{Pr}(a) \cdot \mathbf{Pr}(b)$

Negation:

$$\mathbf{Pr}(\neg s) = 1 - \mathbf{Pr}(s)$$

$$\mathbf{Pr}(\neg s/d) = 1 - \mathbf{Pr}(s/d)$$

If  $b_1, b_2, \dots, b_n$  are pairwise disjoint and exhaust all possibilities, then

$$\mathbf{Pr}(a) = \sum \mathbf{Pr}(ab_i) = \sum \mathbf{Pr}(a / b_i) \cdot \mathbf{Pr}(b_i)$$

$$\mathbf{Pr}(a / c) = \sum \mathbf{Pr}(ab_i / c)$$

Bayes' rule:

$$\mathbf{Pr}(a/b) = \mathbf{Pr}(a) \cdot \mathbf{Pr}(b/a) / \mathbf{Pr}(b)$$

if  $a$  is a disease and  $b$  is a symptom, it is usually easier to estimate numbers on RHS of equation (see below, for subjective probabilities)

# Subjective probability

---

It is reasonable to have non-categorical beliefs even in categorical sentences

- confidence/certainty in a sentence
- “your” probability = *subjective*

Similar to defaults

- move from statistical/group observations to belief about individuals
- but not categorical: how certain am I that Tweety flies?

“Prior probability”  $Pr(x/s)$  ( $s$  = prior state of information or background knowledge)

“Posterior probability”  $Pr(x/E,s)$  ( $E$  = new evidence)

Need to *combine evidence* from various sources

how to derive new beliefs from prior beliefs and new evidence?

want explanations; probability is just a summary



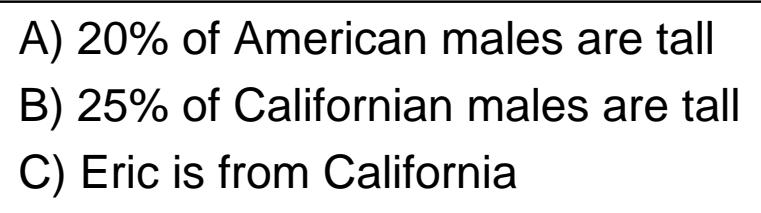
# From statistics to belief

---

Would like to go from statistical information (e.g., the probability that a bird chosen at random will fly) to a degree of belief (e.g., how certain are we that this particular bird, Tweety, flies)

Traditional approach is to find a reference class for which we have statistical information and use the statistics for that class to compute an appropriate degree of belief for an individual

Imagine trying to assign a degree of belief to the proposition “Eric (an American male) is tall” given facts like these

- 
- A) 20% of American males are tall
  - B) 25% of Californian males are tall
  - C) Eric is from California

This is called direct inference

Problem: individuals belong to many classes

- with just A  $\rightarrow$  .2
- A,B,C - prefer more specific  $\rightarrow$  .25
- A,C - no statistics for more specific class  $\rightarrow$  .2?
- B - are Californians a representative sample?

# Basic Bayesian approach

---

Would like a more principled way of calculating subjective probabilities

Assume we have  $n$  atomic propositions  $p_1, \dots, p_n$  we care about. A logical interpretation  $I$  can be thought of as a specification of which  $p_i$  are true and which are false.

Notation: for  $n=4$ , we use  $\langle \neg p_1, p_2, p_3, \neg p_4 \rangle$  to mean the interpretation where only  $p_2$  and  $p_3$  are true.

A joint probability distribution  $J$ , is a function from interpretations to  $[0,1]$  satisfying  $\sum J(I) = 1$  (where  $J(I)$  is the degree of belief in the world being as per  $I$ ).

The degree of belief in any sentence  $\alpha$ :  $Pr(\alpha) = \sum_{I \models \alpha} J(I)$

$$\begin{aligned} \text{Example: } Pr(p_2 \wedge \neg p_4) &= J(\langle \neg p_1, p_2, p_3, \neg p_4 \rangle) + \\ &J(\langle \neg p_1, p_2, \neg p_3, \neg p_4 \rangle) + \\ &J(\langle p_1, p_2, p_3, \neg p_4 \rangle) + \\ &J(\langle p_1, p_2, \neg p_3, \neg p_4 \rangle). \end{aligned}$$

# Problem with the approach

---

To calculate the probabilities of arbitrary sentences involving the  $p_i$ , we would need to know the full joint distribution function.

For  $n$  atomic sentences, this requires knowing  $2^n$  numbers

impractical for all but very small problems

Would like to make plausible assumptions to cut down on what needs to be known.

In the simplest case, all the atomic sentences are independent.

This gives us that

$$J(\langle P_1, \dots, P_n \rangle) = \Pr(P_1 \wedge \dots \wedge P_n) = \prod \Pr(P_i) \quad (\text{where } P_i \text{ is either } p_i \text{ or } \neg p_i)$$

and so only  $n$  numbers are needed.

But this assumption is too strong. A better assumption:

the probability of each  $P_i$  only depends on a small number of  $P_j$   
and the dependence is acyclic.

# Belief networks

Represent all the atoms in a belief network (or Bayes' network).

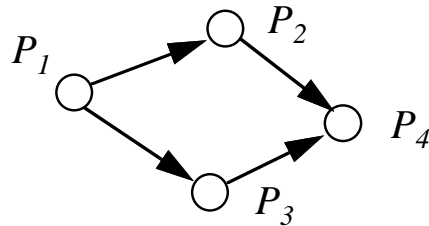
Assume:

$$J(\langle P_1, \dots, P_n \rangle) = \prod \Pr(P_i | c(P_i))$$

where  $\Pr(c(P_i)) > 0$

$c(P)$  = parents of node  $P$

Example:



$$J(\langle P_1, P_2, P_3, P_4 \rangle) = \Pr(P_1) \cdot \Pr(P_2 | P_1) \cdot \Pr(P_3 | P_1) \cdot \Pr(P_4 | P_2, P_3).$$

$$\begin{aligned} \text{So: } J(p_1, \bar{p}_2, p_3, \bar{p}_4) &= \Pr(p_1) \cdot \Pr(\bar{p}_2 | p_1) \cdot \Pr(p_3 | p_1) \cdot \Pr(\bar{p}_4 | \bar{p}_2, p_3) \\ &= \Pr(p_1) \cdot [1 - \Pr(p_2 | p_1)] \cdot \Pr(p_3 | p_1) \cdot [1 - \Pr(p_4 | \bar{p}_2, p_3)] \end{aligned}$$

To fully specify the joint distribution (and therefore probabilities over any subset of the variables), we only need  $\Pr(P | c(P))$  for every node  $P$ .

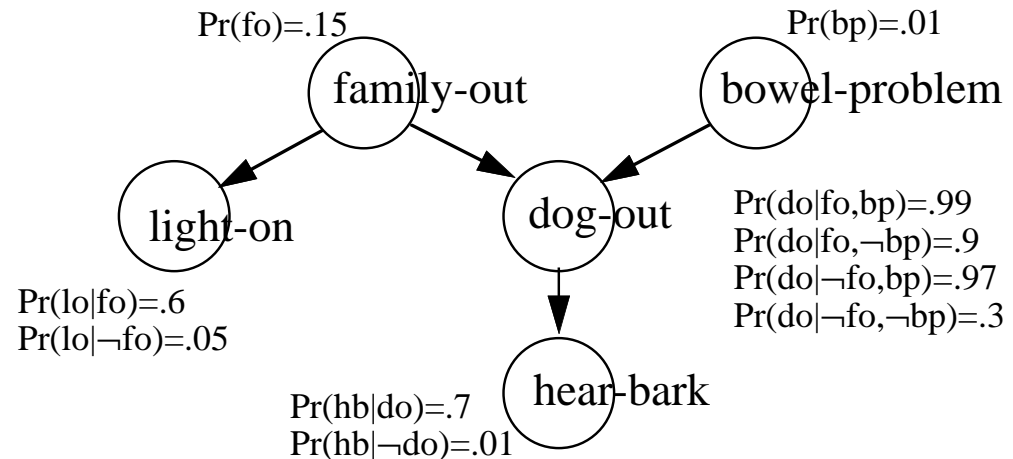
If node  $P$  has parents  $Q_1, \dots, Q_m$ , then we need to know the values of  $\Pr(p | q_1, q_2, \dots, q_m)$ ,  $\Pr(p | \bar{q}_1, q_2, \dots, q_m)$ ,  $\Pr(p | q_1, \bar{q}_2, \dots, q_m)$ , ...,  $\Pr(p | \bar{q}_1, \bar{q}_2, \dots, \bar{q}_m)$ .

$n \cdot 2^m$  numbers  $\ll 2^n$  numbers !

# Using belief networks

Assign a node to each variable in the domain and draw arrows toward each node  $P$  from a select set  $c(P)$  of nodes perceived to be “direct causes” of  $P$ .

arcs can often be interpreted as causal connections



From the DAG, we get that

$$J(\langle FO, LO, BP, DO, HB \rangle) = \mathbf{Pr}(FO) \times \mathbf{Pr}(LO \mid FO) \times \mathbf{Pr}(BP) \times \mathbf{Pr}(DO \mid FO, BP) \times \mathbf{Pr}(HB \mid DO)$$

Using this formula and the 10 numbers above, we can calculate the full joint distribution

# Example calculation

---

Suppose we want to calculate  $Pr(\text{fo} \mid \text{lo}, \neg\text{hb})$

$$Pr(\text{fo} \mid \text{lo}, \neg\text{hb}) = Pr(\text{fo}, \text{lo}, \neg\text{hb}) / Pr(\text{lo}, \neg\text{hb}) \quad \text{where}$$

$$Pr(\text{fo}, \text{lo}, \neg\text{hb}) = \sum J(\langle \text{fo}, \text{lo}, \text{BP}, \text{DO}, \neg\text{hb} \rangle) \quad \text{first 4 values below}$$

$$Pr(\text{lo}, \neg\text{hb}) = \sum J(\langle \text{FO}, \text{lo}, \text{BP}, \text{DO}, \neg\text{hb} \rangle) \quad \text{all 8 values below}$$

$$J(\langle \text{fo}, \text{lo}, \text{bp}, \text{do}, \neg\text{hb} \rangle) = .15 \cdot .6 \cdot .01 \cdot .99 \cdot .3 = .0002673 +$$

$$J(\langle \text{fo}, \text{lo}, \text{bp}, \neg\text{do}, \neg\text{hb} \rangle) = .15 \cdot .6 \cdot .01 \cdot .01 \cdot .99 = .00000891 +$$

$$J(\langle \text{fo}, \text{lo}, \neg\text{bp}, \text{do}, \neg\text{hb} \rangle) = .15 \cdot .6 \cdot .99 \cdot .9 \cdot .3 = .024057 +$$

$$J(\langle \text{fo}, \text{lo}, \neg\text{bp}, \neg\text{do}, \neg\text{hb} \rangle) = .15 \cdot .6 \cdot .99 \cdot .1 \cdot .99 = .0088209 +$$

$$J(\langle \neg\text{fo}, \text{lo}, \text{bp}, \text{do}, \neg\text{hb} \rangle) = .85 \cdot .05 \cdot .01 \cdot .97 \cdot .3 = .000123675$$

$$J(\langle \neg\text{fo}, \text{lo}, \text{bp}, \neg\text{do}, \neg\text{hb} \rangle) = .85 \cdot .05 \cdot .01 \cdot .03 \cdot .99 = .0000126225 +$$

$$J(\langle \neg\text{fo}, \text{lo}, \neg\text{bp}, \text{do}, \neg\text{hb} \rangle) = .85 \cdot .05 \cdot .99 \cdot .3 \cdot .3 = .00378675$$

$$J(\langle \neg\text{fo}, \text{lo}, \neg\text{bp}, \neg\text{do}, \neg\text{hb} \rangle) = .85 \cdot .05 \cdot .99 \cdot .7 \cdot .99 = .029157975$$

$$Pr(\text{fo} \mid \text{lo}, \neg\text{hb}) = .03316 / .06624 = .5$$

# Bypassing the full calculation

---

Often it is possible to calculate some probability values without first calculating the full joint distribution

Example: what is  $\Pr(\text{fo} \mid \text{lo})$ ?

by Bayes rule:  $\Pr(\text{fo} \mid \text{lo}) = \Pr(\text{lo} \mid \text{fo}) \cdot \Pr(\text{fo}) / \Pr(\text{lo})$

but:  $\Pr(\text{lo}) = \Pr(\text{lo} \mid \text{fo}) \cdot \Pr(\text{fo}) + \Pr(\text{lo} \mid \bar{\text{fo}}) \cdot \Pr(\bar{\text{fo}})$

But in general, the problem is NP-hard

- the problem is even hard to approximate in general
- much of the attention on belief networks involves special-purpose procedures that work well for restricted topologies

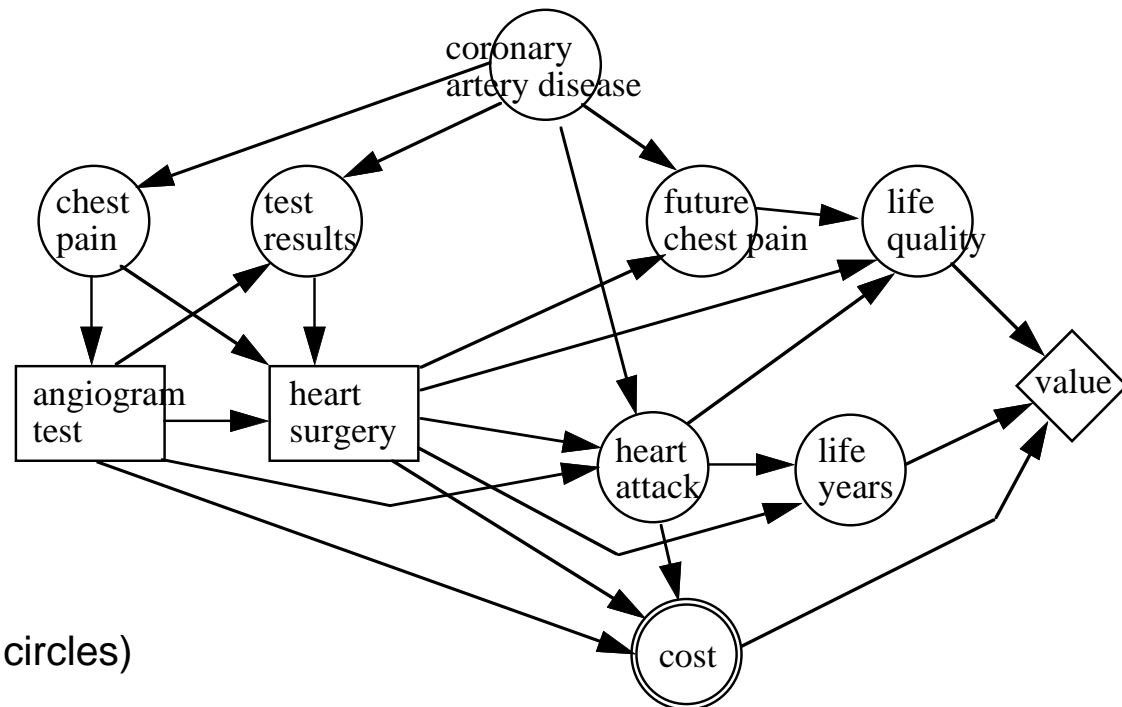
# Influence diagrams

## Graphical knowledge representation for decision problems

- nodes represent propositions or quantities of interest, including decision variables, states of the world, and preference values
- arcs represent influence or relevance (probabilistic or deterministic relationships between the variables)

### Node types

- chance nodes (circles)
- value nodes (diamonds)
- decision nodes (rectangles)
- deterministic nodes (double circles)





# Dempster-Shafer theory

## Another attempt at evidence-pooling

for cases where there is uncertainty about probability

## Uses two-part measure: *belief* and *plausibility*

these are lower and upper bounds on probabilities of a proposition

Relational DB example

Name	Age
a	[ 22,26]
b	[ 20,22]
c	[ 30,35]
d	[ 20,22]
e	[ 28,30]

→ {20,21,22} is the set of possibilities of Age(d), or the *possibility distribution* of Age(d)

Set membership questions like  $\text{Age}(x) \in Q$  cease to be applicable; more natural to ask about the possibility of  $Q$  given the table above of  $\text{Age}(x)$

if  $Q=[20,25]$ , it is *possible* that  $\text{Age}(a) \in Q$ , *not possible* that  $\text{Age}(c) \in Q$ ,  
*certain* that  $\text{Age}(d) \in Q$

What is the probability that the age of someone is in the range [20,25]?

belief=2/5; plausibility=3/5. So answer is [.4,.6].

## DS combination rule → multiple sources

# Vague predicates

---

Not every predicate fits every object exactly (nor fails completely)

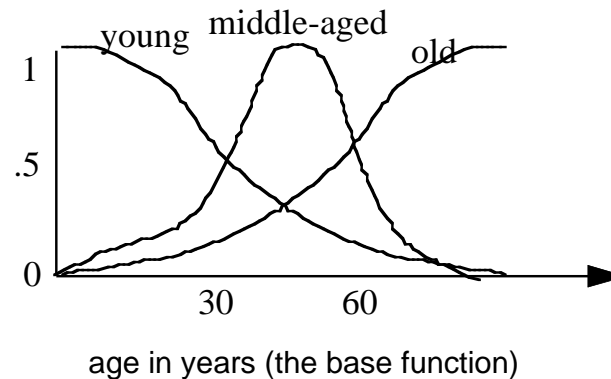
- Categories with degrees of membership  
e.g., fast, old, distant
- Problem: reference sets
  - big fly vs. big elephant

We call predicates that are thought of as holding to a degree vague predicates (or fuzzy predicates).

For each vague predicate, there is a precise base function in terms of which it is understood.

- tall: height
- rich: net worth
- bald: percent hair cover

A degree curve maps the base function to  $[0,1]$ .



# Conjunction and disjunction

---

As with probabilities, we need boolean combinations of properties

Negation is as with probability:

degree of membership in  $\neg P = 1 - \text{degree of membership in } P$

But handle conjunction with MIN and disjunction with MAX!

Example:

Suppose an individual has very high (.95) degree of membership in predicates Tall, Coordinated, Strong, ... for 20 predicates.

Then want to say very high (.95) degree of membership in (Tall  $\wedge$  Coordinated  $\wedge$  Strong  $\wedge$  ...)

as opposed to

Suppose there is a very high (.95) probability of being Tall, of being Coordinated, of being Strong, ... for 20 predicates.

The probability of being all of them at the same time (Tall  $\wedge$  Coordinated  $\wedge$  Strong  $\wedge$  ...) can be low.

Other operators: “very” = square; “somewhat” = square root

# Rules with vague predicates

---

Imagine degrees of fraud = {high, somewhat high, medium, somewhat low, low}, based on a numeric universe of discourse (to some maximum amount)

Construct a set of rules that indicate degrees of fraud based on authorizations and difference in amount of recorded accountability and actual stock:

- 1) If        number of authorizations is often  
   then     fraud is somewhat high
- 2) If        amount is larger than usual  
   then     high fraud

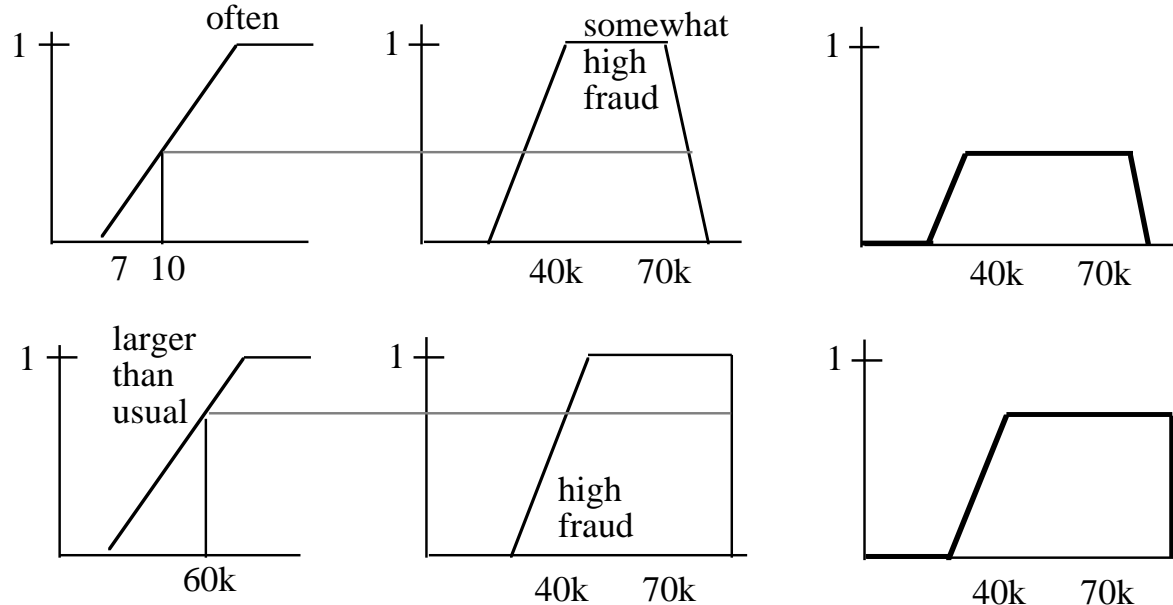
Want to estimate the amount of fraud given inputs

10 authorizations,  
amount of \$60K

# Applying rules

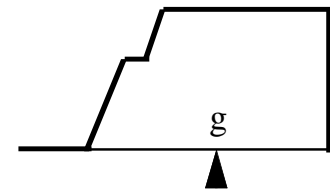
Use degree curves for “somewhat high”, “larger than usual” etc.

Can combine with rules in a way that allows conclusion of rule to apply to the degree that the condition of the rule applied.



Given: 10 authorizations  
amount of 60k

conclusion:  
g = center  
of gravity



---

13.

# Explanation and Diagnosis

# Abductive reasoning

---

So far: reasoning has been primarily *deductive*:

- given KB, is  $\alpha$  an implicit belief?
- given KB, for what  $x$  is  $\alpha[x]$  an implicit belief?

Even default / probabilistic reasoning has a similar form

Now consider a new type of question:

Given KB, and an  $\alpha$  that I do *not* believe,

what would be sufficient to make me believe that  $\alpha$  was true?

or what else would I have to believe for  $\alpha$  to become an implicit belief?

or what would *explain*  $\alpha$  being true?

Deduction: given  $(p \supset q)$ , from  $p$ , deduce  $q$

Abduction: given  $(p \supset q)$ , from  $q$ , abduce  $p$

$p$  is sufficient for  $q$  or one way for  $q$  to be true is for  $p$  to be true

Also induction: given  $p(t_1), q(t_1), \dots, p(t_n), q(t_n)$ , induce  $\forall x (p(x) \supset q(x))$

Can be used for causal reasoning: (*cause*  $\supset$  *effect*)

# Diagnosis

---

One simple version of diagnosis uses abductive reasoning

KB has facts about symptoms and diseases

including:  $(Disease \wedge Hedges \supset Symptoms)$

Goal: find disease(s) that best explain observed symptoms

Observe: we typically do not have knowledge of the form

$(Symptom \wedge \dots \supset Disease)$

so reasoning is not deductive

Example:

<p>(tennis-elbow <math>\supset</math> sore-elbow) (tennis-elbow <math>\supset</math> tennis-player) (arthritis <math>\wedge</math> untreated <math>\supset</math> sore-joints) (sore-joints <math>\supset</math> sore-elbow <math>\wedge</math> sore-hip)</p>
---

Explain: sore-elbow

Want: tennis-elbow,  
(arthritis  $\wedge$  untreated),  
...

Non-uniqueness: multiple equally good explanations

+ logical equivalences:  $(untreated \wedge \neg\neg arthritis)$



# Adequacy criteria

---

Given KB, and  $\beta$  to be explained, we want an  $\alpha$  such that

1.  $\alpha$  is sufficient to account for  $\beta$

$$\text{KB} \cup \{\alpha\} \models \beta \quad \text{or} \quad \text{KB} \models (\alpha \supset \beta)$$

2.  $\alpha$  is not ruled out by KB

$$\text{KB} \cup \{\alpha\} \text{ is consistent} \quad \text{or} \quad \text{KB} \not\models \neg\alpha$$

otherwise  $(p \wedge \neg p)$  would count as an explanation

3.  $\alpha$  is as simple as possible

parsimonious : as few *terms* as possible  
explanations should not unnecessarily  
strong or unnecessarily weak

e.g.  $\text{KB} = \{(p \supset q), \neg r\}$  and  $\beta = q$

$\alpha = (p \wedge s \wedge \neg t)$  is too strong

$\alpha = (p \vee r)$  is too weak

4.  $\alpha$  is in the appropriate vocabulary

atomic sentences of  $\alpha$  should be drawn  
from **H**, possible hypotheses in terms of  
which explanations are to be phrased  
e.g. diseases, original causes

e.g. sore-elbow explains sore-elbow  
trivial explanation

sore-joints explains sore-elbow  
may or may not be suitable

Call such  $\alpha$  an explanation of  $\beta$  wrt KB

# Some simplifications

---

From criteria of previous slide, we can simplify explanations in the propositional case, as follows:

- To explain an arbitrary wff  $\beta$ , it is sufficient to choose a new letter  $p$ , add  $(p \equiv \beta)$  to KB, and then explain  $p$ .

$$\text{KB} \models (E \supset \beta) \quad \text{iff} \quad \text{KB} \cup \{(p \equiv \beta)\} \models (E \supset p)$$

- Any explanation will be (equivalent to) a conjunction of literals (that is, the negation of a clause)

Why? If  $\alpha$  is a purported explanation, and  $\text{DNF}[\alpha] = (d_1 \vee d_2 \vee \dots \vee d_n)$  then each  $d_i$  is also an explanation that is no less simple than  $\alpha$

A simplest explanation is then the negation of a clause with a *minimal* set of literals

So: to explain a literal  $\rho$ , it will be sufficient to find the minimal clauses  $C$  (in the desired vocabulary) such that

1.  $\text{KB} \models (\neg C \supset \rho)$  or  $\text{KB} \models (C \cup \{\rho\})$                       sufficient
2.  $\text{KB} \not\models C$     consistent

# Prime implicates

---

A clause  $C$  is a prime implicate of a KB iff

1.  $\text{KB} \models C$
2. For no  $C^* \subset C$ ,  $\text{KB} \models C^*$


Note: For any clause  $C$ , if  $\text{KB} \models C$ , then some subset of  $C$  is a prime implicate

Example:  $\text{KB} = \{(p \wedge q \wedge r \supset g), (\neg p \wedge q \supset g), (\neg q \wedge r \supset g)\}$

Prime implicates:

$(p \vee \neg q \vee g),$   
 $(\neg r \vee g),$  and  
 $(p \vee \neg p), (g \vee \neg g), \dots$

Note: tautology  $(a \vee \neg a)$  is always a prime implicate unless  $\text{KB} \models a$  or  $\text{KB} \models \neg a$



For explanations:

- want minimal  $C$  such that  $\text{KB} \models (C \cup \{\rho\})$  and  $\text{KB} \not\models C$
- so: find prime implicates  $C$  such that  $\rho \in C$ ;  
then  $\neg(C - \rho)$  must be an explanation for  $\rho$

Example: explanations for  $g$  in example above

- 3 prime implicates contain  $g$ , so get 3 explanations:  $(\neg p \wedge q)$ ,  $r$ , and  $g$

# Computing explanations

---

Given KB, to compute explanations of literal  $\rho$  in vocabulary  $\mathbf{H}$ :

calculate the set  $\{\neg(C - \rho) \mid C \text{ is a prime implicate and } \rho \in C\}$   
prime implicates containing  $\rho$

But how to compute prime implicates?

Can prove: Resolution is complete for non-tautologous prime implicates

$KB \models C$  iff  $KB \rightarrow C$       completeness for  $[\ ]$  is a special case!

So: assuming KB is in CNF, generate *all* resolvents in language  $\mathbf{H}$ , and retain those containing  $\rho$  that are minimal

Could pre-compute all prime implicates, but there may be *exponentially* many, even for a Horn KB

Example: atoms:  $p_i, q_i, E_i, O_i, 0 \leq i < n + E_n, O_n$

wffs:  $E_i \wedge p_i \supset O_{i+1}, E_i \wedge q_i \supset E_{i+1},$   
 $O_i \wedge p_i \supset E_{i+1}, O_i \wedge q_i \supset O_{i+1},$   
 $E_0, \neg O_0$

explain:  $E_n$

# Circuit example

## Components

$\text{Gate}(x) \equiv \text{Andgate}(x) \vee \text{Orgate}(x) \vee \text{Xorgate}(x)$

$\text{Andgate}(a1), \text{Andgate}(a2),$

$\text{Orgate}(o1),$

$\text{Xorgate}(b1), \text{Xorgate}(b2)$

$\text{Fulladder}(f)$  the whole circuit

## Connectivity

$\text{in1}(b1) = \text{in1}(f), \text{in2}(b1) = \text{in2}(f)$

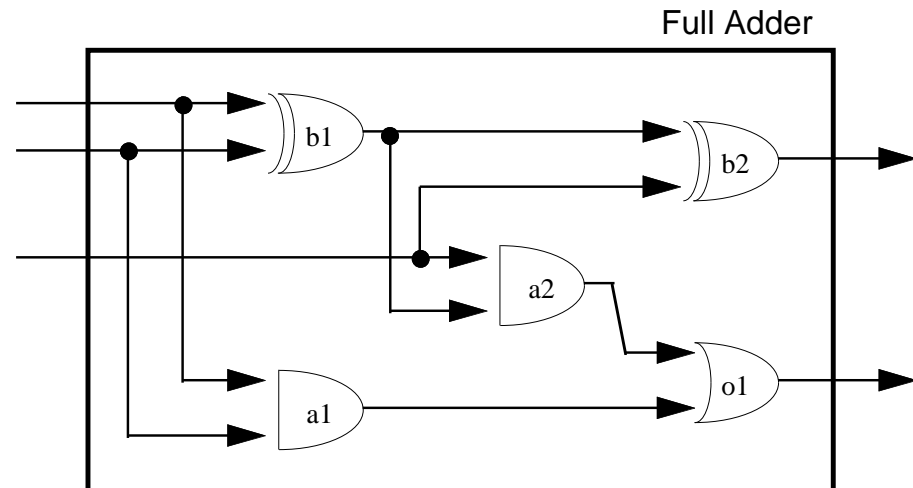
$\text{in1}(b2) = \text{out}(b1), \text{in2}(b2) = \text{in3}(f)$

$\text{in1}(a1) = \text{in1}(f), \text{in2}(a1) = \text{in2}(f)$

$\text{in1}(a2) = \text{in3}(f), \text{in2}(a2) = \text{out}(b1)$

$\text{in1}(o1) = \text{out}(a2), \text{in2}(o1) = \text{out}(a1)$

$\text{out1}(f) = \text{out}(b2), \text{out2}(f) = \text{out}(o1)$



# Circuit behaviour

---

## Truth tables for logical gates

$\text{and}(0,0) = 0, \text{ and}(0,1) = 0, \dots$        $\text{or}(0,0) = 0, \text{ or}(0,1) = 1, \dots$   
 $\text{xor}(0,0) = 0, \text{ xor}(0,1) = 1, \dots$

## Normal behaviour

$\text{Andgate}(x) \wedge \neg \text{Ab}(x) \supset \text{out}(x) = \text{and}(\text{in1}(x), \text{in2}(x))$

$\text{Orgate}(x) \wedge \neg \text{Ab}(x) \supset \text{out}(x) = \text{or}(\text{in1}(x), \text{in2}(x))$

$\text{Xorgate}(x) \wedge \neg \text{Ab}(x) \supset \text{out}(x) = \text{xor}(\text{in1}(x), \text{in2}(x))$

## Abnormal behaviour: fault models

### Examples

$[\text{Orgate}(x) \vee \text{Xorgate}(x)] \wedge \text{Ab}(x) \supset \text{out}(x) = \text{in2}(x)$       (short circuit)

### Other possibilities ...

- some abnormal behaviours may be inexplicable
- some may be compatible with normal behaviour on certain inputs

# Abductive diagnosis

---

Given KB as above + input settings

e.g.  $KB \cup \{in1(f) = 1, in2(f) = 0, in3(f) = 1\}$

we want to explain observations at outputs

e.g.  $(out1(f) = 1 \wedge out2(f) = 0)$

in the language of Ab

We want conjunction of Ab literals  $\alpha$  such that

$$KB \cup Settings \cup \{\alpha\} \models Observations$$

Compute by “propositionalizing”:

For the above,  $x$  ranges over 5 components and  $u, v$  range over 0 and 1.

Easiest to do by preparing a table ranging over all Ab literals, and seeing which conjunctions entail the observations.

# Table for abductive diagnosis

	Ab(b1)	Ab(b2)	Ab(a1)	Ab(a2)	Ab(o1)	Entails observation?
1.	Y	Y	Y	Y	Y	N
2.	Y	Y	Y	Y	N	N
3.	Y	Y	Y	N	Y	N
4.	Y	Y	Y	N	N	N
5.	Y	Y	N	Y	Y	Y
6.	Y	Y	N	Y	N	N
7.	Y	Y	N	N	Y	Y
8.	Y	Y	N	N	N	Y
9.	Y	N	Y	Y	Y	N
10.	Y	N	Y	Y	N	N
11.	Y	N	Y	N	Y	N
12.	Y	N	Y	N	N	N
13.	Y	N	N	Y	Y	Y
14.	Y	N	N	Y	N	N
15.	Y	N	N	N	Y	Y
...						
32.	N	N	N	N	N	N



# Example diagnosis

---

Using the table, we look for minimal sets of literals.

For example, from line (5), we have that

$$Ab(b1) \wedge Ab(b2) \wedge \neg Ab(a1) \wedge Ab(a2) \wedge Ab(o1)$$

entails the observations. However, lines (5), (7), (13) and (15) together lead us to a smaller set of literals (the first explanation below).

The explanations are

1.  $Ab(b1) \wedge \neg Ab(a1) \wedge Ab(o1)$
2.  $Ab(b1) \wedge \neg Ab(a1) \wedge \neg Ab(a2)$
3.  $Ab(b2) \wedge \neg Ab(a1) \wedge Ab(o1)$

Note: not all components are mentioned since for these settings, get the same observations whether or not they are working

but for this fault model only

Can narrow down diagnosis by looking at a number of different settings  
differential diagnosis

# Diagnosis revisited

---

Abductive definition has limitations

- often only care about what is not working
- may not be able to characterize all possible failure modes
- want to prefer diagnoses that claim as few broken components as possible

## Consistency-based diagnosis:

Assume KB uses the predicate  $Ab$  as before, but perhaps only characterizes the normal behaviour

e.g.  $Andgate(x) \wedge \neg Ab(x) \supset out(x) = and(in1(x), in2(x))$

Want a minimal set of components  $D$ , such that

$\{Ab(c) \mid c \in D\} \cup \{\neg Ab(c) \mid c \notin D\}$

can use table as before  
with last column changed  
to “consistency”

is consistent with  $KB \cup Settings \cup Observations$

In previous example, get 3 diagnoses:  $\{b1\}$ ,  $\{b2, a2\}$  and  $\{b2, o1\}$

Note: more complex to handle non-minimal diagnoses

# Some complications

---

## 1. negative evidence

- allow for missing observations  
e.g. ensure that  $KB \cup \{\alpha\} \not\models \text{fever}$

## 2. variables and quantification

- same definition, modulo “simplicity”, (but how to use Resolution?)
- useful to handle open wffs also  
 $KB \cup \{x = 3\} \models P(x)$  handles WH-questions

## 3. probabilities

- not all simplest explanations are equally likely
- also: replace  $(Disease \wedge \dots \supset Symptom)$  by a probabilistic version

## 4. defaults

- instead of requiring  $KB \cup \{\alpha\} \models \beta$ , would prefer that given  $\alpha$ , it is reasonable to believe  $\beta$   
e.g. being a bird explains being able to fly

# Other applications

---

## 1. object recognition

what scene would account for image elements observed?

what objects would account for collection of properties discovered?

## 2. plan recognition

what high-level goals of an agent would account for the actions observed?

## 3. hypothetical reasoning

instead of asking: what would I have to be told to believe  $\beta$ ?

ask instead: what would I learn if I was told that  $\alpha$ ?

Dual of explanation: want  $\beta$  such that

Solution: you learn  $\beta$  on being told  $\alpha$

iff

$\neg\beta$  is an explanation for  $\neg\alpha$

can use the abduction procedure

$KB \cup \{\alpha\} \models \beta$

$KB \not\models \beta$

simplicity, parsimony

using correct vocabulary

---

14.

# Actions

# Situation calculus

---

The situation calculus is a dialect of FOL for representing dynamically changing worlds in which all changes are the result of named actions.

There are two distinguished sorts of terms:

- actions, such as
  - $\text{put}(x,y)$  put object  $x$  on top of object  $y$
  - $\text{walk}(loc)$  walk to location  $loc$
  - $\text{pickup}(r,x)$  robot  $r$  picks up object  $x$
- situations, denoting possible world histories. A distinguished constant  $S_0$  and function symbol  $do$  are used
  - $S_0$  the initial situation, before any actions have been performed
  - $do(a,s)$  the situation that results from doing action  $a$  in situation  $s$

for example:  $do(\text{put}(A,B),do(\text{put}(B,C),S_0))$

the situation that results from putting A on B after putting B on C in the initial situation

# Fluents

---

Predicates or functions whose values may vary from situation to situation are called fluents.

These are written using predicate or function symbols whose last argument is a situation

for example:  $\text{Holding}(r, x, s)$ : robot  $r$  is holding object  $x$  in situation  $s$

can have:  $\neg\text{Holding}(r, x, s) \wedge \text{Holding}(r, x, \text{do}(\text{pickup}(r,x),s))$

the robot is not holding the object  $x$  in situation  $s$ , but is holding it in the situation that results from picking it up

**Note:** there is no distinguished “current” situation. A sentence can talk about many different situations, past, present, or future.

A distinguished predicate symbol  $\text{Poss}(a,s)$  is used to state that  $a$  may be performed in situation  $s$

for example:  $\text{Poss}(\text{pickup}(r,x), S_0)$

it is possible for the robot  $r$  to pickup object  $x$  in the initial situation

This is the entire language.

# Preconditions and effects

---

It is necessary to include in a KB not only facts about the initial situation, but also about world dynamics: what the actions do.

Actions typically have preconditions: what needs to be true for the action to be performed

- $Poss(\text{pickup}(r,x), s) \equiv \forall z. \neg \text{Holding}(r,z,s) \wedge \neg \text{Heavy}(x) \wedge \text{NextTo}(r,x,s)$

a robot can pickup an object iff it is not holding anything, the object is not too heavy, and the robot is next to the object

Note: free variables assumed to be universally quantified

- $Poss(\text{repair}(r,x), s) \equiv \text{HasGlue}(r,s) \wedge \text{Broken}(x,s)$

it is possible to repair an object iff the object is broken and the robot has glue

Actions typically have effects: the fluents that change as the result of performing the action

- $\text{Fragile}(x) \supset \text{Broken}(x, do(\text{drop}(r,x),s))$

dropping a fragile object causes it to break

- $\neg \text{Broken}(x, do(\text{repair}(r,x),s))$

repairing an object causes it to be unbroken



# The frame problem

---

To really know how the world works, it is also necessary to know what fluents are *unaffected* by performing an action.

- $\text{Colour}(x,c,s) \supset \text{Colour}(x, c, \text{do}(\text{drop}(r,x),s))$   
dropping an object does not change its colour
- $\neg\text{Broken}(x,s) \wedge [ x \neq y \vee \neg\text{Fragile}(x) ] \supset \neg\text{Broken}(x, \text{do}(\text{drop}(r,y),s)$   
not breaking things

These are sometimes called frame axioms.

**Problem:** need to know a vast number of such axioms. ( Few actions affect the value of a given fluent; most leave it invariant. )

an object's colour is unaffected by picking things up, opening a door, using the phone, turning on a light, electing a new Prime Minister of Canada, *etc.*

The frame problem:

- in building KB, need to think of these  $\sim 2 \times A \times F$  facts about what does not change
- the system needs to reason efficiently with them

# What counts as a solution?

---

- Suppose the person responsible for building a KB has written down *all* the effect axioms
  - for each fluent  $F$  and action  $A$  that can cause the truth value of  $F$  to change, an axiom of the form  $[R(s) \supset \pm F(do(A,s))]$ , where  $R(s)$  is some condition on  $s$
- We want a systematic procedure for generating all the frame axioms from these effect axioms
- If possible, we also want a *parsimonious* representation for them (since in their simplest form, there are too many)

Why do we want such a solution?

- frame axioms are necessary to reason about actions and are not entailed by the other axioms
  - convenience for the KB builder
  - for theorizing about actions
- |  |                                      |
|--|--------------------------------------|
|  | – modularity: only add effect axioms |
|  | – accuracy: no inadvertent omissions |

# The projection task

---

What can we do with the situation calculus?

We will see later that it can be used for planning.

A simpler job we can handle directly is called the projection task.

Given a sequence of actions, determine what would be true in the situation that results from performing that sequence.

This can be formalized as follows:

Suppose that  $R(s)$  is a formula with a free situation variable  $s$ .

To find out if  $R(s)$  would be true after performing  $\langle a_1, \dots, a_n \rangle$  in the initial situation, we determine whether or not

$$KB \models R(\text{do}(a_n, \text{do}(a_{n-1}, \dots, \text{do}(a_1, S_0) \dots)))$$

For example, using the effect and frame axioms from before, it follows that  $\neg \text{Broken}(B, s)$  would hold after doing the sequence

$$\langle \text{pickup}(A), \text{pickup}(B), \text{drop}(B), \text{repair}(B), \text{drop}(A) \rangle$$

# The legality task

---

The projection task above asks if a condition would hold after performing a sequence of actions, but not whether that sequence can in fact be properly executed.

We call a situation legal if it is the initial situation or the result of performing an action whose preconditions are satisfied starting in a legal situation.

The legality task is the task of determining whether a sequence of actions leads to a legal situation.

This can be formalized as follows:

To find out if the sequence  $\langle a_1, \dots, a_n \rangle$  can be legally performed in the initial situation, we determine whether or not

$$KB \models Poss(a_i, do(a_{i-1}, \dots, do(a_1, S_0) \dots))$$

for every  $i$  such that  $1 \leq i \leq n$ .

# Limitations of the situation calculus

---

This version of the situation calculus has a number of limitations:

- no time: cannot talk about how long actions take, or when they occur
- only known actions: no hidden exogenous actions, no unnamed events
- no concurrency: cannot talk about doing two actions at once
- only discrete situations: no continuous actions, like pushing an object from A to B.
- only hypotheticals: cannot say that an action has occurred or will occur
- only primitive actions: no actions made up of other parts, like conditionals or iterations

We will deal with the last of these below.

First we consider a simple solution to the frame problem ...

# Normal form for effect axioms

---

Suppose there are two positive effect axioms for the fluent *Broken*:

$$\text{Fragile}(x) \supset \text{Broken}(x, \text{do}(\text{drop}(r, x), s))$$

$$\text{NextTo}(b, x, s) \supset \text{Broken}(x, \text{do}(\text{explode}(b), s))$$

These can be rewritten as

$$\begin{aligned} \exists r \{a = \text{drop}(r, x) \wedge \text{Fragile}(x)\} \vee \exists b \{a = \text{explode}(b) \wedge \text{NextTo}(b, x, s)\} \\ \supset \text{Broken}(x, \text{do}(a, s)) \end{aligned}$$

Similarly, consider the negative effect axiom:

$$\neg \text{Broken}(x, \text{do}(\text{repair}(r, x), s))$$

which can be rewritten as

$$\exists r \{a = \text{repair}(r, x)\} \supset \neg \text{Broken}(x, \text{do}(a, s))$$

In general, for any fluent  $F$ , we can rewrite all the effect axioms as two formulas of the form

$$P_F(\mathbf{x}, a, s) \supset F(\mathbf{x}, \text{do}(a, s)) \quad (1)$$

$$N_F(\mathbf{x}, a, s) \supset \neg F(\mathbf{x}, \text{do}(a, s)) \quad (2)$$

where  $P_F(\mathbf{x}, a, s)$  and  $N_F(\mathbf{x}, a, s)$  are formulas whose free variables are among the  $x_i$ ,  $a$ , and  $s$ .

# Explanation closure

---

Now make a completeness assumption regarding these effect axioms:

assume that (1) and (2) characterize *all* the conditions under which an action  $a$  changes the value of fluent  $F$ .

This can be formalized by explanation closure axioms:

$$\neg F(\mathbf{x}, s) \wedge F(\mathbf{x}, do(a,s)) \supset P_F(\mathbf{x}, a, s) \quad (3)$$

if  $F$  was false and was made true by doing action  $a$   
then condition  $P_F$  must have been true

$$F(\mathbf{x}, s) \wedge \neg F(\mathbf{x}, do(a,s)) \supset N_F(\mathbf{x}, a, s) \quad (4)$$

if  $F$  was true and was made false by doing action  $a$   
then condition  $N_F$  must have been true

These explanation closure axioms are in fact disguised versions of frame axioms!

$$\neg F(\mathbf{x}, s) \wedge \neg P_F(\mathbf{x}, a, s) \supset \neg F(\mathbf{x}, do(a,s))$$

$$F(\mathbf{x}, s) \wedge \neg N_F(\mathbf{x}, a, s) \supset F(\mathbf{x}, do(a,s))$$

# Successor state axioms

---

Further assume that our KB entails the following

- integrity of the effect axioms:  $\neg \exists \mathbf{x}, a, s. P_F(\mathbf{x}, a, s) \wedge N_F(\mathbf{x}, a, s)$
- unique names for actions:

$$A(x_1, \dots, x_n) = A(y_1, \dots, y_n) \supset (x_1 = y_1) \wedge \dots \wedge (x_n = y_n)$$

$$A(x_1, \dots, x_n) \neq B(y_1, \dots, y_m) \quad \text{where } A \text{ and } B \text{ are distinct}$$

Then it can be shown that KB entails that (1), (2), (3), and (4) together are logically equivalent to

$$F(\mathbf{x}, do(a, s)) \equiv P_F(\mathbf{x}, a, s) \vee (F(\mathbf{x}, s) \wedge \neg N_F(\mathbf{x}, a, s))$$

This is called the successor state axiom for  $F$ .

For example, the successor state axiom for the *Broken* fluent is:

$$\begin{aligned} \text{Broken}(x, do(a, s)) \equiv & \\ & \exists r \{a = \text{drop}(r, x) \wedge \text{Fragile}(x)\} \\ & \vee \exists b \{a = \text{explode}(b) \wedge \text{NextTo}(b, x, s)\} \\ & \vee \text{Broken}(x, s) \wedge \neg \exists r \{a = \text{repair}(r, x)\} \end{aligned}$$

An object  $x$  is broken after doing action  $a$   
iff  
 $a$  is a dropping action and  $x$  is fragile,  
 or  $a$  is a bomb exploding  
     where  $x$  is next to the bomb,  
 or  $x$  was already broken and  
      $a$  is not the action of repairing it

Note universal quantification: for *any* action  $a$  ...



# A simple solution to the frame problem

---

This simple solution to the frame problem (due to Ray Reiter) yields the following axioms:

- one successor state axiom per fluent
- one precondition axiom per action
- unique name axioms for actions

Moreover, we do not get fewer axioms at the expense of prohibitively long ones

the length of a successor state axioms is roughly proportional to the number of actions which affect the truth value of the fluent

The conciseness and perspicuity of the solution relies on

- quantification over actions
- the assumption that relatively few actions affect each fluent
- the completeness assumption (for effects)

Moreover, the solution depends on the fact that actions always have deterministic effects.

## Limitation: primitive actions

---

As yet we have no way of handling in the situation calculus complex actions made up of other actions such as

- conditionals: If the car is in the driveway then drive else walk
- iterations: while there is a block on the table, remove one
- nondeterministic choice: pickup up some block and put it on the floor

and others

Would like to *define* such actions in terms of the primitive actions, and inherit their solution to the frame problem

Need a compositional treatment of the frame problem for complex actions

Results in a novel programming language for discrete event simulation and high-level robot control

# The Do formula

---

For each complex action  $A$ , it is possible to define a formula of the situation calculus,  $Do(A, s, s')$ , that says that action  $A$  when started in situation  $s$  may legally terminate in situation  $s'$ .

Primitive actions:  $Do(A, s, s') = Poss(A, s) \wedge s' = do(A, s)$

Sequence:  $Do([A; B], s, s') = \exists s''. Do(A, s, s'') \wedge Do(B, s'', s')$

Conditionals:  $Do([if \phi then A else B], s, s') =$   
 $\phi(s) \wedge Do(A, s, s') \vee \neg\phi(s) \wedge Do(B, s, s')$

Nondeterministic branch:  $Do([A / B], s, s') = Do(A, s, s') \vee Do(B, s, s')$

Nondeterministic choice:  $Do([\pi x. A], s, s') = \exists x. Do(A, s, s')$

*etc.*

**Note:** programming language constructs with a purely logical situation calculus interpretation

# GOLOG

---

GOLOG (Algol in logic) is a programming language that generalizes conventional imperative programming languages

- the usual imperative constructs + concurrency, nondeterminism, more...
- bottoms out *not* on operations on internal states (assignment statements, pointer updates) but on *primitive actions* in the world (e.g. pickup a block)
- what the primitive actions do is user-specified by precondition and successor state axioms

What does it mean to “execute” a GOLOG program?

- find a sequence of primitive actions such that performing them starting in some initial situation  $s$  would lead to a situation  $s'$  where the formula  $Do(A, s, s')$  holds
- give the sequence of actions to a robot for actual execution in the world

**Note:** to find such a sequence, it will be necessary to reason about the primitive actions

$A ; \textit{if Holding}(x) \textit{ then } B \textit{ else } C$

to decide between  $B$  and  $C$  we need to determine if the fluent *Holding* would be true after doing  $A$

# GOLOG example

---

Primitive actions:  $\text{pickup}(x)$ ,  $\text{putonfloor}(x)$ ,  $\text{putontable}(x)$

Fluents:  $\text{Holding}(x,s)$ ,  $\text{OnTable}(x,s)$ ,  $\text{OnFloor}(x,s)$

Action preconditions:  $\text{Poss}(\text{pickup}(x), s) \equiv \forall z. \neg \text{Holding}(z, s)$

$\text{Poss}(\text{putonfloor}(x), s) \equiv \text{Holding}(x, s)$

$\text{Poss}(\text{putontable}(x), s) \equiv \text{Holding}(x, s)$

Successor state axioms:

$\text{Holding}(x, \text{do}(a,s)) \equiv a = \text{pickup}(x) \vee$

$\text{Holding}(x,s) \wedge a \neq \text{putontable}(x) \wedge a \neq \text{putonfloor}(x)$

$\text{OnTable}(x, \text{do}(a,s)) \equiv a = \text{putontable}(x) \vee \text{OnTable}(x,s) \wedge a \neq \text{pickup}(x)$

$\text{OnFloor}(x, \text{do}(a,s)) \equiv a = \text{putonfloor}(x) \vee \text{OnFloor}(x,s) \wedge a \neq \text{pickup}(x)$

Initial situation:  $\forall x. \neg \text{Holding}(x, S_0)$

$\text{OnTable}(x, S_0) \equiv x = A \vee x = B$

Complex actions:

*proc* ClearTable : **while**  $\exists b. \text{OnTable}(b)$  **do**  $\pi b [\text{OnTable}(b)? ; \text{RemoveBlock}(b)]$

*proc* RemoveBlock( $x$ ) :  $\text{pickup}(x) ; \text{putonfloor}(x)$

# Running GOLOG

---

To find a sequence of actions constituting a legal execution of a GOLOG program, we can use Resolution with answer extraction.

For the above example, we have

$$KB \models \exists s. Do(\text{ClearTable}, S_0, s)$$

The result of this evaluation yields

$$s = do(\text{putonfloor}(B), do(\text{pickup}(B), do(\text{putonfloor}(A), do(\text{pickup}(A), S_0))))))$$

and so a correct sequence is

$$\langle \text{pickup}(A), \text{putonfloor}(A), \text{pickup}(B), \text{putonfloor}(B) \rangle$$

When what is known about the actions and initial state can be expressed as Horn clauses, the evaluation can be done in Prolog.

The GOLOG interpreter in Prolog has clauses like

```
do(A, S1, do(A, S1)) :- prim_action(A), poss(A, S1).
```

```
do(seq(A, B), S1, S2) :- do(A, S1, S3), do(B, S3, S2).
```

This provides a convenient way of controlling a robot at a high level.

---

15.

# Planning

# Planning

---

So far, in looking at actions, we have considered how an agent could figure out what to do given a high-level program or complex action to execute.

Now, we consider a related but more general reasoning problem: figure out what to do to make an arbitrary condition true. This is called planning.

- the condition to be achieved is called the goal
- the sequence of actions that will make the goal true is called the plan

Plans can be at differing levels of detail, depending on how we formalize the actions involved

- “do errands” vs. “get in car at 1:32 PM, put key in ignition, turn key clockwise, change gears,…”

In practice, planning involves anticipating what the world will be like, but also observing the world and replanning as necessary...



# Using the situation calculus

---

The situation calculus can be used to represent what is known about the current state of the world and the available actions.

The planning problem can then be formulated as follows:

Given a formula  $Goal(s)$ , find a sequence of actions  $a$  such that

$$KB \models Goal(do(\mathbf{a}, S_0)) \wedge Legal(do(\mathbf{a}, S_0))$$

where  $do(\langle a_1, \dots, a_n \rangle, S_0)$  is an abbreviation for

$$do(a_n, do(a_{n-1}, \dots, do(a_2, do(a_1, S_0)) \dots))$$

and where  $Legal(\langle a_1, \dots, a_n \rangle, S_0)$  is an abbreviation for

$$Poss(a_1, S_0) \wedge Poss(a_2, do(a_1, S_0)) \wedge \dots \wedge Poss(a_n, do(\langle a_1, \dots, a_{n-1} \rangle, S_0))$$

So: given a goal formula, we want a sequence of actions such that

- the goal formula holds in the situation that results from executing the actions, and
- it is possible to execute each action in the appropriate situation

# Planning by answer extraction

---

Having formulated planning in this way, we can use Resolution with answer extraction to find a sequence of actions:

$$KB \models \exists s. Goal(s) \wedge Legal(s)$$

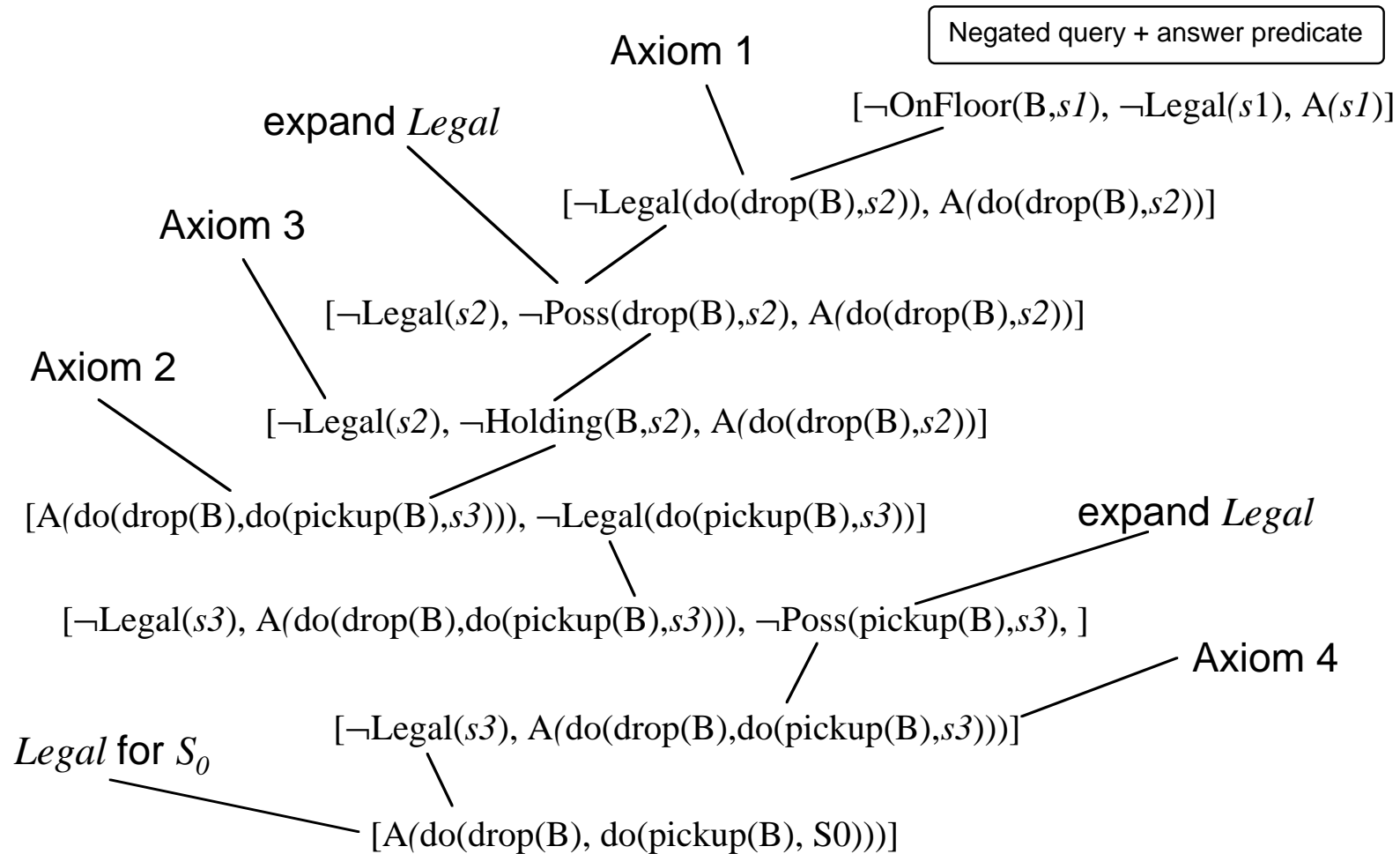
We can see how this will work using a simplified version of a previous example:

An object is on the table that we would like to have on the floor. Dropping it will put it on the floor, and we can drop it, provided we are holding it. To hold it, we need to pick it up, and we can always do so.

- Effects:       $OnFloor(x, do(drop(x),s))$   
                   $Holding(x, do(pickup(x),s))$   
                  Note: ignoring frame problem
- Preconds:     $Holding(x, s) \supset Poss(drop(x), s)$   
                   $Poss(pickup(x), s)$
- Initial state:  $OnTable(B, S_0)$
- The goal:      $OnFloor(B, s)$

KB

# Deriving a plan



Here is the plan: in the initial situation, pickup block B, and in the resulting situation, drop B.

# Using Prolog

---

Because all the required facts here can be expressed as Horn clauses, we can use Prolog directly to synthesize a plan:

```
onfloor(X,do(drop(X),S)).
holding(X,do(pickup(X),S)).
poss(drop(X),S) :- holding(X,S).
poss(pickup(X),S).
ontable(b,s0).
legal(s0).
legal(do(A,S)) :- poss(A,S), legal(S).
```

With the Prolog goal `?- onfloor(b,S), legal(S).`

we get the solution `S = do(drop(b),do(pickup(b),s0))`

But planning problems are rarely this easy!

Full Resolution theorem-proving can be problematic for a complex set of axioms dealing with actions and situations explicitly...

# The STRIPS representation

---

STRIPS is an alternative representation to the pure situation calculus for planning.

from work on a robot called Shaky at SRI International in the 60's.

In STRIPS, we do not represent histories of the world, as in the situation calculus.

Instead, we deal with a single world state at a time, represented by a database of ground atomic wffs (e.g.,  $\text{In}(\text{robot}, \text{room}_1)$ )

This is like the database of facts used in procedural representations and the working memory of production systems

Similarly, we do not represent actions as part of the world model (cannot reason about them directly), as in the situation calculus.

Instead, actions are represented by operators that syntactically transform world models

An operator takes a DB and transforms it to a new DB

# STRIPS operators

---

Operators have pre- and post-conditions

- precondition = formulas that need to be true at start
- “delete list” = formulas to be removed from DB
- “add list” = formulas to be added to DB

Example:  $\text{PushThru}(o, d, r_1, r_2)$

“the robot pushes object  $o$  through door  $d$  from room  $r_1$  to room  $r_2$ ”

- precondition:  $\text{InRoom}(\text{robot}, r_1), \text{InRoom}(o, r_1), \text{Connects}(d, r_1, r_2)$
- delete list:  $\text{InRoom}(\text{robot}, r_1), \text{InRoom}(o, r_1)$
- add list:  $\text{InRoom}(\text{robot}, r_2), \text{InRoom}(o, r_2)$

STRIPS problem space =  $\left\{ \begin{array}{l} \text{initial world model, } \text{DB}_0 \text{ (list of ground atoms)} \\ \text{set of operators (with preconds and effects)} \\ \text{goal statement (list of atoms)} \end{array} \right.$

desired plan: sequence of ground operators

# STRIPS Example

---

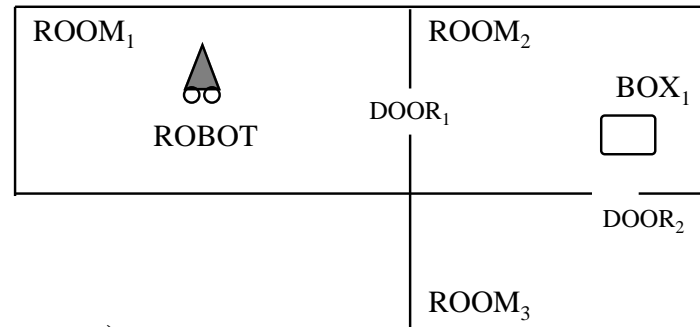
In addition to PushThru, consider

GoThru( $d, r_1, r_2$ ):

precondition:  $\text{InRoom}(\text{robot}, r_1), \text{Connects}(d, r_1, r_2)$

delete list:  $\text{InRoom}(\text{robot}, r_1)$

add list:  $\text{InRoom}(\text{robot}, r_2)$



DB<sub>0</sub>:

$\text{InRoom}(\text{robot}, \text{room}_1) \text{ InRoom}(\text{box}_1, \text{room}_2)$

$\text{Connects}(\text{door}_1, \text{room}_1, \text{room}_2) \text{ Box}(\text{box}_1)$

$\text{Connects}(\text{door}_2, \text{room}_2, \text{room}_3) \dots$

Goal: [  $\text{Box}(x) \wedge \text{InRoom}(x, \text{room}_1)$  ]

# Progressive planning

---

Here is one procedure for planning with a STRIPS like representation:

**Input** : a world model and a goal (ignoring variables)  
**Output** : a plan or fail.

ProgPlan[DB,Goal] =

  If Goal is satisfied in DB, then return empty plan

  For each operator  $o$  such that precondition( $o$ ) is satisfied in the current DB:

    Let  $DB' = DB + \text{addlist}(o) - \text{dellist}(o)$

    Let plan = ProgPlan[ $DB'$ ,Goal]

    If plan  $\neq$  fail, then return [act( $o$ ) ; plan]

  End for

  Return fail

This depth-first planner searches forward from the given  $DB_0$  for a sequence of operators that eventually satisfies the goal

$DB'$  is the progressed world state



# Regressive planning

---

Here is another procedure for planning with a STRIPS like representation:

**Input** : a world model and a goal (ignoring variables)  
**Output** : a plan or fail.

```
RegrPlan[DB,Goal] =  
  If Goal is satisfied in DB, then return empty plan  
  For each operator  $o$  such that  $\text{dellist}(o) \cap \text{Goal} = \{\}$ :  
    Let  $\text{Goal}' = \text{Goal} + \text{precond}(o) - \text{addlist}(o)$   
    Let  $\text{plan} = \text{RegrPlan}[\text{DB}, \text{Goal}']$   
    If  $\text{plan} \neq \text{fail}$ , then return  $[\text{plan} ; \text{act}(o)]$   
  End for  
  Return fail
```

This depth-first planner searches backward for a sequence of operators that will reduce the goal to something satisfied in  $\text{DB}_0$

Goal' is the regressed goal

# Computational aspects

---

Even without variables, STRIPS planning is NP-hard.

Many methods have been proposed to avoid redundant search

e.g. partial-order planners, macro operators

One approach: application dependent control

Consider this range of GOLOG programs:

< any deterministic program >



*while*  $\neg$ Goal *do*  $\pi a . a$

fully specific about sequence  
of actions required

easy to execute

any sequence such that *Goal*  
holds at end

as hard as planning!

pick an action



In between, the two extremes we can give domain-dependent guidance to a planner:

*while*  $\neg$ Goal *do*  $\pi a . [Acceptable(a)? ; a]$

where *Acceptable* is formalized separately

This is called forward filtering .

# Hierarchical planning

---

The basic mechanisms of planning so far still preserve all detail needed to solve a problem

- attention to too much detail can derail a planner to the point of uselessness
- would be better to first search through an *abstraction space*, where unimportant details were suppressed
- when solution in abstraction space is found, account for remaining details

## ABSTRIPS

precondition wffs in abstraction space will have fewer literals than those in ground space

e.g., PushThru operator

- high abstraction: applicable whenever an object is pushable and a door exists
- lower: robot and obj in same room, connected by a door to target room
- lower: door must be open
- original rep: robot next to box, near door

predetermined partial order of predicates with “criticality” level

# Reactive systems

---

Some suggest that explicit, symbolic production of formal plans is something to be avoided (especially considering computational complexity)

even propositional case is intractable; first-order case is undecidable

Just “react”: observe conditions in the world and decide (or look up) *what to do next*

can be more robust in face of unexpected changes in the environment

⇒ reactive systems

“Universal plans”: large lookup table (or boolean circuit) that tells you exactly what to do based on current conditions in the world

Reactive systems have impressive performance on certain low-level problems (e.g. learning to walk), and can even look “intelligent”

but what are the limitations? ...

---

16.

# The Tradeoff between Expressiveness and Tractability

# Limit expressive power?

---

Defaults, probabilities, *etc.* can all be thought of as extensions to FOL, with obvious applications

Why not strive for the *union* of all such extensions? all of English?

Problem: automated reasoning

Lesson here:

reasoning procedures required for more expressive languages  
may not work very well in practice

Tradeoff: expressiveness vs. tractability

Overview:

- a Description Logic example
- limited languages
- the problem with cases
- vivid reasoning as an extreme case
- less vivid reasoning
- hybrid reasoning systems

# Simple description logic

---

Consider the language FL defined by:

$\langle \text{concept} \rangle ::= \text{atom}$	$\langle \text{role} \rangle ::= \text{atom}$
[AND $\langle \text{concept} \rangle \dots \langle \text{concept} \rangle$ ]	[RESTR $\langle \text{role} \rangle \langle \text{concept} \rangle$ ]
[ALL $\langle \text{role} \rangle \langle \text{concept} \rangle$ ]	
[SOME $\langle \text{role} \rangle$ ] (= [EXISTS 1 $\langle \text{role} \rangle$ ])	

Example: [ALL :Child [AND Female Student]]

an individual whose children are female students

[ALL [RESTR :Child Female] Student]

an individual whose female children are students

there may or may not be male children and they may or may not be students

Interpretation  $\mathcal{I} = \langle D, I \rangle$  as before, but with

$$I[[\text{RESTR } r \ c]] = \{ (x,y) \mid (x,y) \in I[r] \text{ and } y \in I[c] \}$$

So [RESTR :Child Female] is the :Child relation restricted to females = :Daughter

Subsumption defined as usual

# Computing subsumption

---

First for  $FL^- = FL$  without the `RESTR` operator

- put the concepts into normalized form
- to see if  $C$  subsumes  $D$  make sure that

$$\begin{aligned} & [\text{AND } p_1 \dots p_k \\ & \quad [\text{SOME } r_1] \dots [\text{SOME } r_m] \\ & \quad [\text{ALL } s_1 c_1] \dots [\text{ALL } s_n c_n] \end{aligned}$$

1. for every  $p \in C$ ,  $p \in D$
2. for every  $[\text{SOME } r] \in C$ ,  $[\text{SOME } r] \in D$
3. for every  $[\text{ALL } s c] \in C$ , find an  $[\text{ALL } s d] \in D$  such that  $c$  subsumes  $d$ .

Can prove that this method is sound and complete relative to definition based on interpretations

Running time:

- normalization is  $O(n^2)$
- structural matching: for each part of  $C$ , find a part of  $D$ . Again  $O(n^2)$

What about all of  $FL$ , including `RESTR`?



# Subsumption in FL

---

- cannot settle for part-by-part matching

[ALL [RESTR :Friend [AND Male Doctor]] [AND Tall Rich]]

subsumes

[AND [ALL [RESTR :Friend Male] [AND Tall Bachelor]]  
[ALL [RESTR :Friend Doctor] [AND Rich Surgeon]]]

- complex interactions

[SOME [RESTR  $r$  [AND  $a$   $b$ ]]]

subsumes

[AND [SOME [RESTR  $r$  [AND  $c$   $d$ ]]] [ALL [RESTR  $r$   $c$ ] [AND  $a$   $e$ ]]  
[ALL [RESTR  $r$  [AND  $d$   $e$ ]  $b$ ]]

In general: FL is powerful enough to encode *all* of propositional logic.

There is a mapping  $\Omega$  from CNF wffs to FL where

$\models (\alpha \supset \beta)$  iff  $\Omega(\alpha)$  is subsumed by  $\Omega(\beta)$

But  $\models (\alpha \supset (p \wedge \neg p))$  iff  $\alpha$  is unsatisfiable

Conclusion: there is no good algorithm for FL unless P=NP

# Moral

---

Even small doses of expressive power can come at a significant computational price

Questions:

- what properties of a representation language control its difficulty?
- how far can expressiveness be pushed without losing good algorithms
- when is easy reasoning adequate for KR purposes?

These questions remain unanswered, but some progress:

- need for case analyses is a major factor
- tradeoff for DL languages is reasonably well understood
- best addressed (perhaps) by looking at working systems

Useful approach:

- find reasoning tasks that are tractable
- analyze difficulty in extending them

# Limited languages

---

Many reasoning problems that can be formulated in terms of FOL entailment ( $KB \models \alpha$ ) admit very specialized methods because of the restricted form of either KB or  $\alpha$

although problem could be solved using full resolution, there is no need

## Example 1: Horn clauses

- SLD resolution provides more focussed search
- in propositional case, a linear procedure is available

## Example 2: Description logics

Can do DL subsumption using Resolution

Introduce predicate symbols for concepts, and “meaning postulates” like

$$\begin{aligned} \forall x[P(x) \equiv \forall y(\text{Friend}(x,y) \supset \text{Rich}(y)) \\ \wedge \forall y(\text{Child}(x,y) \supset \\ \forall z(\text{Friend}(y,z) \supset \text{Happy}(z)))] \end{aligned} \quad \begin{aligned} [\text{AND} [\text{ALL} : \text{Friend Rich}] \\ [\text{ALL} : \text{Child} \\ [\text{ALL} : \text{Friend Happy}]]] \end{aligned}$$

Then ask if  $MP \models \forall x[P(x) \supset Q(x)]$

# Equations

---

## Example 3: linear equations

Let  $E$  be the usual axioms for arithmetic:

$$\forall x \forall y (x+y = y+x), \quad \forall x (x+0 = x), \quad \dots \quad \text{Peano axioms}$$

Then we get the following:

$$E \models (x+2y=4 \wedge x-y=1) \supset (x=2 \wedge y=1)$$

Can “solve” linear equations using Resolution!

But there is a much better way:

Gauss-Jordan method with back substitution

- subtract (2) from (1):  $3y = 3$
- divide by 3:  $y = 1$
- substitute in (1):  $x = 2$

In general, a set of linear equations can be solved in  $O(n^3)$  operations

This idea obviously generalizes!

always advantageous to use a specialized procedure when it is available,  
rather than a general method like Resolution

# When is reasoning hard?

---

Suppose that instead of linear equations, we have something like

$$(x+2y=4 \vee 3x-y=7) \wedge x-y=1$$

Can still show using Resolution:  $y > 0$

To use GJ method, we need to split cases:

$$\begin{array}{l} x+2y=4 \wedge x-y=1 \quad \rightsquigarrow \quad y=1 \\ 3x-y=7 \wedge x-y=1 \quad \rightsquigarrow \quad y=2 \end{array} \quad \therefore y > 0$$

What if 2 disjunctions?  $(eqnA_1 \vee eqnB_1) \wedge (eqnA_2 \vee eqnB_2)$

there are *four* cases to consider with GJ method

What if  $n$  binary disjunctions?  $(eqnA_1 \vee eqnB_1) \wedge \dots \wedge (eqnA_n \vee eqnB_n)$

there are  $2^n$  cases to consider with GJ method

with  $n=30$ , would need to solve  $10^9$  systems of equations!

Conclusion: case analysis is still a *big* problem.

Question: can we avoid case analyses??

# Expressiveness of FOL

---

Ability to represent incomplete knowledge

$P(a) \vee P(b)$             but which?

$\exists x P(x)$              $P(a) \vee P(b) \vee P(c) \vee \dots$

and even

$c \neq 3$              $c=1 \vee c=2 \vee c=4 \vee \dots$

Reasoning with facts like these requires somehow “covering” all the implicit cases

languages that admit efficient reasoning do *not* allow this type of knowledge to be represented

- Horn clauses,
- description logics,
- linear equations, ...

only limited forms of disjunction, quantification etc.

# Complete knowledge

---

One way to ensure tractability:

somehow restrict contents of KB so that reasoning by cases is not required

But is complete knowledge enough for tractability?

suppose  $KB \models \alpha$  or  $KB \models \neg\alpha$ , as in the CWA

Get: queries reduce to  $KB \models \rho$ , literals

But: it can still be hard to answer for literals

Example:  $KB = \{(p \vee q), (\neg p \vee q), (\neg p \vee \neg q)\}$

Have:  $KB \models \neg p \wedge q$  complete!

But to find literals may require case analysis

So complete knowledge is not enough to avoid case analyses if the knowledge is “hidden” in the KB.

Need a form of complete knowledge that is more explicit...

# Vivid knowledge

---

Note: If KB is complete and consistent, then it is satisfied by a *unique* interpretation  $I$

Why? define  $I$  by  $I \models p$  iff  $\text{KB} \models p$

ignoring  
quantifiers  
for now

Then for any  $I^*$ , if  $I^* \models \text{KB}$  then  $I^*$  agrees with  $I$  on all atoms  $p$

Get:  $\text{KB} \models \alpha$  iff  $I \models \alpha$

entailments of KB are sentences that are true at  $I$

explains why queries reduce to atomic case

$(\alpha \vee \beta)$  is true iff  $\alpha$  is true or  $\beta$  is true, etc.

if we have the  $I$ , we can easily determine what is or is not entailed

Problem: KB can be complete and consistent, but unique interpretation may be hard to find

Solution: a KB is vivid if it is a complete and consistent set of literals (for some language)

e.g.  $\text{KB} = \{\neg p, q\}$

specifies  $I$  directly



# Quantifiers

---

As with the CWA, we can generalize the notion of vivid to accommodate queries with quantifiers

A first-order KB is vivid iff for some finite set of positive function-free ground literals  $KB^+$ ,  $KB = KB^+ \cup Negs \cup Dc \cup Un$ .

Get a simple recursive algorithm for  $KB \models \alpha$ :

$KB \models \exists x.\alpha$  iff  $KB \models \alpha[x/c]$ , for some  $c \in KB^+$

$KB \models (\alpha \vee \beta)$  iff  $KB \models \alpha$  or  $KB \models \beta$

$KB \models \neg\alpha$  iff  $KB \not\models \alpha$

$KB \models (c = d)$  iff  $c$  and  $d$  are the same constant

$KB \models p$  iff  $p \in KB^+$

This is just database retrieval

- useful to store  $KB^+$  as a collection of relations
- only  $KB^+$  is needed to answer queries, but *Negs*, *Dc*, and *Un* are required to *justify* the correctness of the procedure

# Analogues

---

Can think of a vivid KB as an analogue of the world

there is a 1-1 correspondence between

- objects in the world and constants in the KB<sup>+</sup>
- relationships in the world and syntactic relationships in the KB<sup>+</sup>

for example, if constants  $c_1$  and  $c_2$  stand for objects in the world  $o_1$  and  $o_2$

there is a relationship  $R$  holding between objects  $o_1$  and  $o_2$  in the world

iff

constants  $c_1$  and  $c_2$  appear as a tuple in the relation represented by  $R$

Not true in general

for example, if  $\text{KB} = \{P(a)\}$  then it only uses 1 constant, but could be talking about a world where there are 5 individuals of which 4 satisfy  $P$

Result: certain reasoning operations are easy

- how many objects satisfy  $P$  (by counting)
- changes to the world (by changes to KB<sup>+</sup>)

# Beyond vivid

---

Requirement of vividness is very strict.

Want weaker alternatives with good reasoning properties

## Extension 1

ignoring  
quantifiers  
again

Suppose KB is a finite set of literals

- not necessarily a complete set (no CWA)
- assume consistent, else trivial

Cannot reduce  $\text{KB} \models \alpha$  to literal queries

if  $\text{KB} = \{p\}$  then  $\text{KB} \models (p \wedge q \vee p \wedge \neg q)$  but  $\text{KB} \not\models p \wedge q$  and  $\text{KB} \not\models p \wedge \neg q$

But: assume  $\alpha$  is small. Can put into CNF

$$\alpha \rightsquigarrow (c_1 \wedge \dots \wedge c_n)$$

- $\text{KB} \models \alpha$  iff  $\text{KB} \models c_i$ , for every clause in CNF of  $\alpha$
- $\text{KB} \models c$  iff  $c$  has complimentary literals – tautology  
or  $\text{KB} \cap c$  is not empty

## Extension 2

---

Imagine KB vivid as before + new definitions:

$\forall xyz[R(x,y,z) \equiv \dots \text{wff in vivid language } \dots]$

Example: have vivid KB using predicate ParentOf

add:  $\forall xy[\text{MotherOf}(x,y) \equiv \text{ParentOf}(x,y) \wedge \text{Female}(x)]$

To answer query containing  $R(t_1, t_2, t_3)$ , simply macro expand it with definition and continue

- can handle arbitrary logical operators in definition since they become part of *query*, not KB
- can generalize to handle predicates not only in vivid KB, provided that they bottom out to KB<sup>+</sup>

$\forall xy[\text{AncestorOf}(x,y) \equiv \text{ParentOf}(x,y) \vee$   
 $\exists z \text{ParentOf}(x,z) \wedge \text{AncestorOf}(z,y)]$

- clear relation to Prolog  
a version of logic programming based on inductive definitions,  
not Horn clauses

# Other extensions

---

Vivification: given non-vivid KB, attempt to make vivid e.g. by eliminating disjunctions *etc.*

for example,

- use taxonomies to choose between disjuncts

Flipper is a whale or a dolphin.

- use intervals to encompass disjuncts

The picnic will be on June 2, 3, or 4th.

- use defaults to choose between disjuncts

Serge works in Toronto or Montreal.

Problem: what to do with function symbols, when Herbrand universe is not finite?

partial Herbrand base?

# Hybrid reasoning

---

Want to be able to incorporate a number of special-purpose efficient reasoners into a single scheme such as Resolution

Resolution will be the glue that holds the reasoners together

Simple form: semantic attachment

- attach procedures to functions and predicates  
e.g. numbers: procedures on plus, LessThan, ...
- ground terms and atomic sentences can be *evaluated* prior to Resolution
  - $P(\text{factorial}(4), \text{times}(2,3)) \rightsquigarrow P(24, 6)$
  - $\text{LessThan}(\text{quotient}(36,6), 5) \vee \alpha \rightsquigarrow \alpha$
- much better than reasoning directly with axioms

More complex form: theory resolution

- build theory into unification process (the way paramodulation builds in =)
- extended notion of complimentary literals

$\{\alpha, \text{LessThan}(2,x)\}$  and  $\{\text{LessThan}(x,1), \beta\}$  resolve to  $\{\alpha,\beta\}$

# Using descriptions

---

Imagine that predicates are defined elsewhere as concepts in a description logic

Married  $\doteq$  [AND ...]                      Bachelor  $\doteq$  [AT-MOST ...]

then  $\{P(x), \text{Married}(x)\}$  and  $\{\text{Bachelor}(\text{john}), Q(y)\}$  resolve to  $\{P(\text{john}), Q(y)\}$

Can use description logic procedure to decide if two predicates are complimentary

instead of explicit meaning postulates

Residues: for “almost” complimentary literals

$\{P(x), \text{Male}(x)\}$  and  $\{\neg\text{Bachelor}(\text{john}), Q(y)\}$

resolve to

$\{P(\text{john}), Q(y), \text{Married}(\text{john})\}$

since the two literals are contradictory *unless* John is married

Main issue: what resolvents are necessary to get the same conclusions as from meaning postulates?

residues are necessary for completeness

---

**THE END**