## Non-Monotonic Logic

Marco Degano

Philosophical Logic 2025 1 December 2025

## Readings

### Suggested:

- ► Frank Veltman, lecture notes on counterfactuals (sec. 5). https://staff.fnwi.uva.nl/f.j.m.m.veltman/papers/ Notes\_Counterfactuals.pdf
- ► S. Kraus, D. Lehmann & M. Magidor (1990), *Nonmonotonic Reasoning, Preferential Models and Cumulative Logics*.

### Plan

- 1. Defeasible Reasoning
- 2. Cumulative Consequence Relation
- 3. Preferential Models
- 4. Applications

### Outline

- 1. Defeasible Reasoning
- 2. Cumulative Consequence Relation
- 3. Preferential Models
- 4. Applications

## Defeasible reasoning

### Birds fly.



- ► We often use **generalizations** that are *rationally compelling* but not deductively valid.
- ► We are talking about what *normally* or *typically* happens, not about exceptionless laws.
- Reasoning is defeasible when conclusions may have to be withdrawn in the light of further information, even if we keep our original premises.

## Broad cases of defeasible reasoning

- Everyday decision-making: Lights are off in a café, so you assume it's closed. Then you see people inside and an "Open" sign: you keep "lights looked off from outside", but give up "the café is closed".
- ► Social reasoning: Your friend normally replies within minutes. This time they don't answer for hours, so you think they're upset. Then you learn they were on a long flight: you keep that they usually reply quickly, but drop "they're upset with me".
- ▶ **Default categorization (Tweety):** From "Tweety is a bird" we infer "Tweety flies". Learning "Tweety is a penguin", we keep that Tweety is a bird and that birds normally fly, but reject "Tweety flies".

## Monotonicity: formal notion

Let  $\models$  be a (single-conclusion) consequence relation between *sets* of formulas and formulas.

 $\models$  is monotonic if for all sets  $\Gamma, \Delta$  and all  $\varphi$ :

if 
$$\Gamma \models \varphi$$
, then for every  $\Delta \supseteq \Gamma$ ,  $\Delta \models \varphi$ 

Adding premises never invalidates an earlier consequence.

Classical consequence |= (and standard proof systems for classical logic) satisfy monotonicity.

A consequence relation  $\triangleright$  is non-monotonic if there exist  $\Gamma \subseteq \Delta$  and  $\varphi$  such that

$$\Gamma \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \varphi \hspace{0.5em} \hspace{0.$$

New information can defeat previous conclusions.

## Non-monotonic reasoning in practice

We use a non-monotonic consequence symbol  $\sim$  (read  $\alpha \sim \beta := \text{ if } \alpha, \text{ then normally } \beta$ ):

$$Bird(x) \sim Flies(x)$$

From the knowledge base  $K = \{Bird(Tweety)\}\$  we may infer

$$K \sim \mathsf{Flies}(\mathsf{Tweety}).$$

► After adding Penguin(Tweety), and a more specific default

Penguin(
$$x$$
)  $\sim \neg \mathsf{Flies}(x)$ ,

the enlarged  $K' = K \cup \{Penguin(Tweety)\}\$  no longer supports Flies(Tweety). Instead we get

$$K' \sim \neg \mathsf{Flies}(\mathsf{Tweety}).$$

This is exactly a failure of monotonicity.

### From Aristotle to Al

Defeasible Reasoning



John McCarthy (1927 - 2011)



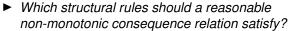
Raymond Reiter (1939-2002)

- Aristotle already distinguished strict demonstration from more tentative, practical reasoning based on generalizations.
- In modern logic, **non-monotonic logic** is the study of formal systems intended to capture defeasible reasoning patterns.
- In AI and knowledge representation, many formalisms were proposed:
  - Negation as failure in logic programming.
  - Circumscription (McCarthy).
  - Default logic (Reiter).
- These different formalisms all induce some non-monotonic consequence relation  $\sim$  on formulas.

# Today's focus: the KLM perspective Kraus, Lehmann & Magidor (KLM):

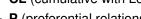
▶ Treat  $\triangleright$  itself as primitive:  $\varphi \triangleright \psi$  reads

"from  $\varphi$ , it *normally* follows  $\psi$ ."



- CL (cumulative with Loop).
- P (preferential relations).
- R (rational relations, later work).
- and some further, stronger systems.
- Prove representation theorems: each such family corresponds to a natural class of models (with preferences between worlds/states).





Daniel Lehmann

Sarit Kraus



Menachem Magidor

Today: focus on **cumulative and preferential relations** and system C and P.

### Outline

- 1. Defeasible Reasoning
- 2. Cumulative Consequence Relation
- 3. Preferential Models
- 4. Applications

## Cumulative consequence C

A consequence relation  $\sim$  is **cumulative** iff it satisfies the rules:

- 1. Reflexivity:  $\varphi \sim \varphi$
- 2. Left logical equivalence: if  $\varphi \models \psi$  and  $\psi \models \varphi$ , and  $\varphi \not\sim \chi$ , then  $\psi \not\sim \chi$ .
- 3. Right weakening: if  $\varphi \models \psi$  and  $\chi \triangleright \varphi$ , then  $\chi \triangleright \psi$ .
- 4. Cut: if  $\varphi \wedge \psi \triangleright \chi$  and  $\varphi \triangleright \psi$ , then  $\varphi \triangleright \chi$ .
- 5. Cautious monotonicity: if  $\varphi \sim \psi$  and  $\varphi \sim \chi$ , then  $\varphi \wedge \psi \sim \chi$ .
- ➤ System C is meant to be the *minimal* structural core of reasonable non-monotonic consequence.
- ightharpoonup is closed under classical equivalence and consequence.
- ► Plausible conclusions can be "re-used" (Cut).
- ► Learning something you *already* inferred as plausible never harms (CMon).

## Some derived rules (on blackboard)

In system C we can derive, among others, the following rules.

### ► Equivalence

$$\frac{\alpha \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \beta \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \alpha \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \gamma}{\beta \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \gamma}$$

If  $\alpha$  and  $\beta$  are plausible consequences of each other, then they have the same plausible consequences.

#### ► Another rule

$$\frac{\alpha \vee \beta \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \alpha \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \gamma}{\alpha \vee \beta \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \gamma}$$

If from  $\alpha \vee \beta$  we plausibly recover  $\alpha$ , and from  $\alpha$  we plausibly get  $\gamma$ , then already  $\alpha \vee \beta$  plausibly entails  $\gamma$ .

### Cumulative models

We give a semantics for system C. Fix a set W of worlds.

#### A cumulative model

$$\mathcal{M}_C = \langle S, \prec, V \rangle$$

#### has:

- ► A set S of states (sets of worlds, possible "epistemic states" of an agent).
- ▶ A labelling function  $\ell: S \to \mathcal{P}(W) \setminus \{\emptyset\}$ :

 $\ell(s)$  is the set of worlds compatible with state s.

ightharpoonup A binary relation  $\prec$  on S expressing preference / normality:

$$s' \prec s = "s'$$
 is more normal (preferred) than  $s$ ".

► A valuation *V* assigning truth values to atoms at worlds:

$$V(w,p) \in \{0,1\}$$
 for each world  $w \in W$  and atomic  $p$ .

Classical satisfaction  $\mathcal{M}_C$ ,  $w \models \alpha$  is defined in the usual way.

### Cumulative models

#### Definition (State satisfaction)

Let  $\mathcal{M}_C = \langle S, \ell, \prec, V \rangle$ . For a state  $s \in S$  and formula  $\alpha$ :

$$s \models \alpha \quad \text{iff} \quad \forall w \in \ell(s) \ (\mathcal{M}_C, w \models \alpha).$$

So s satisfies  $\alpha$  iff all its worlds do.

$$\llbracket \alpha \rrbracket^{\mathcal{M}_C} = \{ s \in S \mid s \models \alpha \}$$

for the set of states that satisfy  $\alpha$ .

### Cumulative models

#### Definition (Smoothness for states)

A subset  $A \subseteq S$  is **smooth** (with respect to  $\prec$ ) iff for every  $s \in A$ :

- $\blacktriangleright$  either s is  $\prec$ -minimal in A,
- $\blacktriangleright$  or there is some  $\prec$ -minimal  $s' \in A$  with  $s' \prec s$ .

 $\mathcal{M}_C$  is a **cumulative model** iff for every formula  $\alpha$ , the set  $\llbracket \alpha \rrbracket^{\mathcal{M}_C}$  is smooth.

Smoothness is the analogue of the "no infinite descent" / limit assumption, now formulated for sets of states.

## Example: a simple cumulative model

	p	q
$\overline{w_1}$	1	1
$w_2$	1	0
$w_3$	0	1

Let  $W = \{w_1, w_2, w_3\}$  and let V be the valuation given by the table. Define a structure

$$\mathcal{M}_C = \langle S, \ell, \prec, V \rangle$$

by:

$$S = \{s_0, s_1\}, \qquad \ell(s_0) = \{w_1, w_2\}, \qquad \ell(s_1) = \{w_3\}, \qquad s_0 \prec s_1.$$

$$[p]^{\mathcal{M}_C} = \{s_0\}, \qquad [q]^{\mathcal{M}_C} = \{s_1\}, \qquad [T]^{\mathcal{M}_C} = \{s_0, s_1\}.$$

In each of these sets the  $\prec$ -minimal elements are well behaved (e.g.  $s_0$  is the unique minimal element of  $[\![\top]\!]^{\mathcal{M}_C}$ ), so  $\mathcal{M}_C$  satisfies the smoothness condition and is therefore a cumulative model.

How to make smoothness fail here? Add  $s_1 \prec s_0$ .

## Consequence in cumulative models

#### Given a cumulative model

$$\mathcal{M}_C = \langle S, \ell, \prec, V \rangle$$

we define a model-relative consequence relation  $\sim_{\mathcal{M}_C}$ .

#### Definition (Cumulative consequence in $\mathcal{M}_C$ )

$$\alpha \sim_{\mathcal{M}_C} \beta$$
 iff for every  $\prec$ -minimal  $s \in [\![\alpha]\!]^{\mathcal{M}_C}, s \models \beta$ .

- ▶ Collect all states where  $\alpha$  holds:  $\llbracket \alpha \rrbracket^{\mathcal{M}_C}$ .
- ▶ Restrict to the *best* (most normal)  $\alpha$ -states (the  $\prec$ -minimal ones).
- ▶ If in all these best  $\alpha$ -states every compatible world satisfies  $\beta$ , then  $\beta$  is a *cumulative consequence* of  $\alpha$ .

### Example: consequence in a cumulative model Take atoms b, f, r.

Let  $W = \{w_1, w_2, w_3\}$  and let V be the valuation given by the table. Define a cumulative model

$$\mathcal{M}_C = \langle S, \ell, \prec, V \rangle$$

by

$$S = \{s_1, s_2, s_3\}, \quad \ell(s_1) = \{w_1\}, \ \ell(s_2) = \{w_2, w_3\}, \ \ell(s_3) = \{w_3\}$$

and a preference relation

$$s_1 \prec s_2 \qquad s_1 \prec s_3$$

with no further <-links.

$$b \hspace{0.2em}\not\sim_{\mathcal{M}_C} f \hspace{1cm} b \wedge r \hspace{0.2em}\not\sim_{\mathcal{M}_C} f \hspace{1cm} \neg f \hspace{0.2em}\not\sim_{\mathcal{M}_C} \neg r$$

## Representation theorem for system C

### Theorem (KLM representation for C)

A consequence relation  $\triangleright$  on formulas is **cumulative** iff there exists a cumulative model

$$\mathcal{M}_C = \langle S, \ell, \prec, V \rangle$$

such that, for all formulas  $\alpha, \beta$ ,

$$\alpha \sim \beta$$
 iff  $\alpha \sim_{\mathcal{M}_C} \beta$ 

- ► The proof rules of system C exactly capture reasoning in cumulative models.
- Preferential models are a special case:
  - S is a set of states, each labelled by a *single* world,
  - $\prec$  is a strict partial order on S.

### Outline

- 1. Defeasible Reasoning
- 2. Cumulative Consequence Relation
- 3. Preferential Models
- 4. Applications

Applications

## The move to preferential models

- ▶ In a **cumulative model**, a state  $s \in S$  is an *information state*:  $\ell(s) \subseteq W$  is the set of worlds compatible with what is currently taken for granted.
- ► States can therefore be "coarse-grained": one state may keep several classical possibilities open at once.
- ► A preferential model is the special case where every state is a singleton:  $\ell(s) = \{w\}$ . We *collapse* states with worlds and order the worlds directly.
- ▶ This makes the normality ordering more concrete: we compare "how things might be" world by world, just as in Lewis-Stalnaker similarity semantics for counterfactuals.
- Trade-off:

Defeasible Reasoning

- we gain simplicity and a tighter fit with the counterfactual picture
- but lose generality: not every cumulative consequence relation can be represented by an ordering on single worlds.

## Example

Take atoms p, q and worlds:  $w_1 : p \land q \qquad w_2 : p \land \neg q$ 

Consider states:

$$s_{ exttt{coarse}}:=\{w_1,w_2\}$$
 ("we know  $p$ , but  $q$  is still open") 
$$s_1:=\{w_1\}\qquad s_2:=\{w_2\}$$

We might prefer the "generic" state where only p is settled:

$$s_{\text{coarse}} \prec s_1, \qquad s_{\text{coarse}} \prec s_2.$$

In a pure **preferential** model, we only see  $w_1$  and  $w_2$ . There is no separate node for the information state

"p is known, q is undetermined".

So we lose the ability to order and reason about *information states* as such. We only order 'fully specified' ways the world might be.

The cumulative model above is actually 'representable' by a preferential model. How? Can you think of a case which cannot be represented in preferential models?

## Example

Take atoms p, q, r and worlds:

$$w_1: p \wedge q \wedge r$$
  $w_2: \neg p \wedge q \wedge r$   $w_3: p \wedge \neg q \wedge \neg r$ 

$$w_3: p \land \neg q \land \neg r$$

Preferential Models

Consider states:

$$s_p := \{w_1\}$$
  $s_q := \{w_2\}$   $s_{\text{coarse}} := \{w_2, w_3\}$ 

We order states by:

$$s_{\text{coarse}} \prec s_p, \qquad s_{\text{coarse}} \prec s_q$$

$$p \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} r \hspace{0.5em} (p \vee q) \not\hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} r$$

There is no way, just by ordering worlds, to make the "most normal"  $p \vee q$ -situation" behave like our coarse state  $\{w_2, w_3\}$  that mixes r and  $\neg r$  while still treating p and q separately as r-supporting.

### A preferential model is a triple

$$\mathcal{M}_P = \langle W, \prec, V \rangle$$

#### such that:

- ightharpoonup  $\prec$  is a strict partial order on W
- lacktriangle for every formula  $\alpha$ , the truth set  $[\![\alpha]\!]^{\mathcal{M}_P} = \{w \in W \mid w \models \alpha\}$  is **smooth** with respect to  $\prec$ .

#### Definition (Smoothness for worlds)

A subset  $A \subseteq W$  is **smooth** (with respect to  $\prec$ ) iff for every  $w \in A$ :

- $\triangleright$  either w is  $\prec$ -minimal in A,
- ightharpoonup or there is some  $\prec$ -minimal  $w' \in A$  with  $w' \prec w$ .

## Preferential consequence

### Definition (Preferential consequence)

Given a preferential model  $\mathcal{M}_P = \langle W, \prec, V \rangle$ , define:

$$\alpha \succ_{\mathcal{M}_P} \beta$$
 iff for every  $\prec$ -minimal  $w \in [\![\alpha]\!]^{\mathcal{M}_P}, \ w \models \beta.$ 

$$\llbracket \alpha \rrbracket^{\mathcal{M}_P} = \{ w \in W \mid \mathcal{M}_P, w \models \alpha \}$$

#### So we:

- $\blacktriangleright$  look at all  $\alpha$ -worlds.
- ightharpoonup pick the most normal ones (the  $\prec$ -minimal  $\alpha$ -worlds),
- $\blacktriangleright$  and require all of them to satisfy  $\beta$ .

Applications

## Knowledge bases and preferential entailment

Cumulative Consequence Relation

A (default) knowledge base K is a set of conditionals

$$\alpha \sim \beta$$

#### Definition (Satisfaction of a knowledge base)

Let K be a set of conditionals  $\alpha \succ \beta$ . A preferential model  $\mathcal{M}_P = \langle W, \prec, V \rangle$  satisfies K iff for every  $\alpha \succ \beta \in K$  we have

$$\alpha \sim_{\mathcal{M}_P} \beta$$

#### Definition (Preferential entailment)

Let K be a set of conditionals. We write

$$K \models_{\mathsf{pref}} \alpha \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \beta$$

iff for every preferential model  $\mathcal{M}_P$ : if  $\mathcal{M}_P$  satisfies K, then

$$\alpha \sim_{\mathcal{M}_P} \beta$$

## Example: a bird-penguin knowledge base

Atoms: b (bird), f (flies), p (penguin).

#### Consider the knowledge base

$$K = \{p \triangleright b, p \triangleright \neg f, b \triangleright f\}$$

- $\blacktriangleright$  Any preferential model  $\mathcal{M}_P$  that satisfies K must make the most normal p-worlds non-flying birds.
- ▶ In all such models we also have:

$$p \wedge b \sim_{\mathcal{M}_P} \neg f$$

SO

$$K \models_{\mathsf{pref}} p \wedge b \not\sim \neg f$$

- $ightharpoonup K 
  ot \models_{\mathsf{pref}} p \mid_{\sim} f \text{ [countermodel on blackboard]}$
- $ightharpoonup K \models_{\mathsf{pref}} f \hspace{-0.2em}\sim\hspace{-0.9em} \neg p \hspace{-0.5em} [\mathsf{proof} \hspace{-0.5em} \mathsf{on} \hspace{-0.5em} \mathsf{blackboard}]$

## System P (preferential consequence)

System  ${\bf P}$  is system  ${\bf C}$  plus one extra rule,  ${\bf Or}$ .

- 1. Reflexivity  $\varphi \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \varphi$
- 2. Left logical equivalence

If  $\varphi \models \psi$  and  $\psi \models \varphi$ , and  $\varphi \triangleright \chi$ , then  $\psi \triangleright \chi$ .

3. Right weakening

If  $\varphi \models \psi$  and  $\chi \triangleright \varphi$ , then  $\chi \triangleright \psi$ .

4. Cut

If  $\varphi \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \psi$  and  $\varphi \wedge \psi \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \chi$ , then  $\varphi \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \chi$ .

5. Cautious monotonicity

If  $\varphi \sim \psi$  and  $\varphi \sim \chi$ , then  $\varphi \wedge \psi \sim \chi$ .

6. **Or** 

If  $\varphi \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \chi$  and  $\psi \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \chi$ , then  $\varphi \vee \psi \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \chi$ .

## Reading the Or rule

**Or:** if  $\varphi \triangleright \chi$  and  $\psi \triangleright \chi$ , then  $\varphi \lor \psi \triangleright \chi$ .

If both  $\varphi$  and  $\psi$  individually are good enough reasons for  $\chi$ , then so is their disjunction.

- (1) a. If John comes to the party, it will normally be great.
  - b. If Cathy comes to the party, it will normally be great.
  - c. So if either John or Cathy comes, it will normally be great.

## Some derived rules in P (blackboard)

In P we can derive several rules:

S-rule

$$\frac{\varphi \wedge \psi \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \chi}{\varphi \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} (\psi \supset \chi)}$$

Union

$$\frac{\varphi \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \psi \hspace{0.2em} \chi \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \gamma}{\varphi \vee \chi \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \psi \vee \gamma}$$

## Representation theorem

### Theorem (KLM representation for P)

A consequence relation  $\sim$  on formulas satisfies all rules of system  ${\bf P}$  iff there exists a preferential model

$$\mathcal{M}_P = \langle W, \prec, V \rangle$$

such that, for all formulas  $\varphi, \psi$ ,

$$\varphi \sim \psi \quad \text{iff} \quad \varphi \sim_{\mathcal{M}_P} \psi$$

For a knowledge base K, the *closure* of K under the rules of  $\mathbf{P}$  coincides with preferential entailment:

$$K \vdash_{\mathbf{P}} \varphi \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \psi \hspace{0.5em} \text{ iff } \hspace{0.5em} K \models_{\mathsf{pref}} \varphi \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \psi$$

So to know what K preferentially entails, we can also reason syntactically with  ${\bf P}$  instead of quantifying over all preferential models.

## Duplicate labels and the representation theorem

Fix a propositional language with atoms p, q and this valuation:

$$\begin{array}{c|cccc} & p & q \\ \hline w_0 & 0 & 0 \\ w_1 & 1 & 0 \\ w_2 & 1 & 1 \\ w_3 & 1 & 1 \\ \end{array}$$

We define a preferential model  $\mathcal{M}_P = \langle W, \prec, V \rangle$  by:

$$W=\{w_0,w_1,w_2,w_3\}$$
 
$$w_0 \prec w_2, \qquad w_1 \prec w_3 \qquad \text{and no other } \prec\text{-links}.$$

- ► This is a perfectly good preferential model (smoothness holds, ≺ is a strict partial order).
- It defines a consequence relation  $\sim_W$  that satisfies all rules of system P, so  $\sim_W$  is a preferential consequence relation.
- ► However, there is **no** preferential model with *unique labels* (i.e. with taking worlds-as-valuations identifying  $w_2$  with  $w_3$ ) that induces exactly the same consequence relation  $\sim_W$ .

## Representation theorems

The **representation theorem** is a global, *structural* result about *consequence relations*:

- ► Fix a proof system (e.g. C or P).
- ► Consider all binary relations |~ on formulas.
- ► The theorem says:
  - $\mid\sim$  satisfies the rules of the system  $\iff$   $\mid\sim$  =  $\mid\sim_{\mathcal{M}}$  for some model  $\mathcal{M}$
- So it classifies which abstract consequence relations are exactly those induced by a given class of models.

**Soundness and Completeness** (for a fixed entailment notion) is a more familiar, *formula-level* result:

- ► Fix a semantics (e.g. preferential models) and a proof system (e.g. P).
- For a knowledge base K and a conditional  $\alpha \sim \beta$ :

$$K \vdash_{\mathbf{P}} \alpha \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \beta \hspace{0.5em}\iff \hspace{0.5em} K \models_{\mathsf{nref}} \alpha \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \beta$$

► This says: whatever is valid in *all* models is derivable, and vice versa.

### An underivable rule

Notably, P does *not* validate the following rule:<sup>1</sup>

#### **Rational Monotonicity:**

if 
$$\varphi \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \chi$$
 and  $\varphi \hspace{0.2em}\not\sim\hspace{-0.9em}\mid\hspace{0.58em} \neg \psi$ , then  $\varphi \wedge \psi \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \chi$ .

Is this a good rule for non-monotonic reasoning?

We have already encountered this rule before. How was it called and to what example was related?

Similarity analysis of counterfactuals: Strengthening with a Possibility axiom scheme and *The Verdi-Bizet-Satie* example.

 $<sup>^1</sup>$ Adding this rule to system  ${f P}$  yields system  ${f R}$  (Lehmann & Magidor 1992). Semantically, R corresponds to ranked models: there is a total preorder  $\leq$  on W (a ranking of worlds) whose strict part is  $\prec$ . Thus any two worlds are comparable in rank, i.e. for all  $v, w \in W$  we have v < w or w < v(or both).

### Preferential Models and Counterfactuals

We can make the connection with the similarity framework we introduced for counterfactuals, with the following additional assumptions:

- ▶ The limit assumption:  $\prec$  is a well-founded partial order on W.
- ▶ **Absoluteness:** for every  $u, w \in W : \prec_u = \prec_w [\prec_w \text{ is independent of } w]$
- ▶ Universality: for every  $w \in W$ ,  $W_w = W$  [the ordering is on W]

Recall the original clause for counterfactuals.

 $M\models (\varphi\leadsto \psi)$  iff for every world  $w\in W$ ,  $M,u\models \psi$  for every closest  $[\![\varphi]\!]$ -world u to w.

With Universality and Absoluteness, we can simplify it as follows:

 $M \models (\varphi \leadsto \psi) \text{ iff } M, u \models \psi \text{ for every } \prec\text{-minimal } \llbracket \varphi \rrbracket\text{-world } u.$ 

And then we rewrite  $M \models (\varphi \leadsto \psi)$  as  $\varphi \models_{\mathcal{M}_P} \psi$ .

## Object language vs metalanguage

► The conditional → (for counterfactuals) is an object-language connective:

$$\varphi \leadsto \psi$$

is a formula that can itself be embedded.

► The non-monotonic consequence symbol  $\sim$  is a *metalanguage* relation:

$$\varphi \sim \psi$$

is not a formula of L, but a statement about L.

► In the KLM approach, the central object of study *is* this relation  $\sim$  and its structural properties.

#### Outline

- 1. Defeasible Reasoning
- 2. Cumulative Consequence Relation
- 3. Preferential Models
- 4. Applications

# The frame problem

Consider designing how a robot should reason about actions and change.

- (2)If the daylight sensor is low, turn on the light. a.
  - If the temperature is low, turn on the heating. h.

In classical logic we might have:

DaylightLow → LightOn TempLow → HeatingOn

But what about persistence?

- ► After turning the light on, should the robot keep believing that the light is *still* on at the next time step?
- ▶ Writing explicit "frame axioms" saying that everything stays the same unless affected by an action quickly leads to an explosion of axioms.

# A non-monotonic take on the frame problem

Idea: use default persistence rules instead of explicit frame axioms.

Introduce discrete time steps  $t, t+1, \ldots$  and write  $F_t$  for "F holds at time t".

For each F (light on, heating on, . . . ) we have a default:

$$F_t \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} F_{t+1}$$

"if F holds at time t, then normally F still holds at time t+1".

- ► This replaces a huge family of classical frame axioms by a small, uniform schema of non-monotonic persistence rules.
- ► The rules are defeasible: e.g. if we also know that the action at t is TurnOffLight, then the default  $LightOn_t \sim LightOn_{t+1}$  is defeated.
- Non-monotonicity is crucial: new information about actions can make the robot *retract* a previous persistence conclusion without changing the earlier facts.

## Grice and implicatures

- (3)A: Does John speak English? B: Well, he knows the colours.
  - From B's answer we *normally* infer that John does not speak English (or at least not very well): if B could honestly say "Yes". that would be more informative.
  - ▶ This is a **conversational implicature**: a defeasible inference driven by the assumption that speakers respect conversational maxims.
  - ► The inference is **non-monotonic**: it can be cancelled without contradiction, e.g.
    - B: Well, he knows the colours. In fact, his English is pretty good.

### Pronoun resolution as default reasoning

- (4) John met Bill at the station. *He* greeted *him*.
  - ▶ By default, we resolve pronouns in line with simple preferences (e.g. subject  $\rightarrow$  "he", object  $\rightarrow$  "him"):

John = 
$$he$$
, Bill =  $him$ .

- ▶ This preferred interpretation is a **default**: it reflects what normally happens, given the syntax and discourse structure.
- But it is defeasible. Additional material can force a different resolution, e.g.
- (5) John met Bill at the station. *He* greeted *him*. Then John greeted him as well.
  - ▶ Now we are pushed to reinterpret the first sentence so that *he* = Bill, him = John.

#### Temporal anaphora and non-monotonicity

- In narrative discourse with simple past tense, there is a strong **default**: events are understood as occurring in the order in which they are mentioned (forward-moving timeline).
- ► This default is **defeasible**: world knowledge or discourse relations can override it.
- (6)John fell. Mary pushed him.
  - ► The default forward-reading would place John's falling *before* Mary's pushing.
  - ▶ But our knowledge about causation and the more natural discourse relation forces a different ordering: Mary pushed John before he fell.
  - ► Asher & Lascarides (2003) give a systematic non-monotonic account of such temporal and discourse inferences (and related phenomena like lexical disambiguation).

### Exercise: Soundness of rules in C

- ► Show that the rules of system C are sound with respect to cumulative and preferential models.
- ► Show that the **Or** rule is sound with respect to preferential models.
- ► Show that the Or rule is not valid in cumulative models.

## Exercise: Equivalence relation

We define an equivalence relation on formulas by:

$$\alpha \sim \beta$$
 :  $\iff$   $\alpha \sim \beta$  and  $\beta \sim \alpha$ 

Assume  $\mid \sim$  is a cumulative consequence relation (i.e., it satisfies the rules of System C).

Show that:

$$\alpha \sim \beta$$
 iff  $\forall \gamma \ (\alpha \sim \gamma \Leftrightarrow \beta \sim \gamma)$ 

#### Exercise: Underivable rules in P

#### Check that the following rules *cannot* be derived in P:

- 1. If  $\varphi \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} (\psi \supset \chi)$ , then  $\varphi \wedge \psi \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.58em} \chi$ .
- 2. If  $\varphi \lor \psi \sim \chi$ , then  $\varphi \sim \chi$  or  $\psi \sim \chi$ .

#### Exercise: The **And** rule

#### And

If  $\varphi \sim \psi$  and  $\varphi \sim \chi$ , then  $\varphi \sim \psi \wedge \chi$ .

- ▶ Show that the **And** rule is a derived rule of system C.
- ► Show semantically that **And** is valid in cumulative models, without using the representation theorem for C.

#### Exercise: rules for P

Show that the following system, with  $\pmb{\mathsf{And}}$  in place of  $\pmb{\mathsf{Cut}}$ , is an equivalent axiomatization of system  $\pmb{\mathsf{P}}$ :

- 1. Reflexivity  $\varphi \sim \varphi$ .
- 2. Left logical equivalence If  $\varphi \models \psi$  and  $\psi \models \varphi$ , and  $\varphi \not \sim \chi$ , then  $\psi \not \sim \chi$ .
- 3. Right weakening If  $\varphi \models \psi$  and  $\chi \triangleright \varphi$ , then  $\chi \triangleright \psi$ .
- 4. **And** If  $\varphi \sim \psi$  and  $\varphi \sim \chi$ , then  $\varphi \sim \psi \wedge \chi$ .
- 5. Cautious monotonicity If  $\varphi \sim \psi$  and  $\varphi \sim \chi$ , then  $\varphi \wedge \psi \sim \chi$ .
- 6. Or If  $\varphi \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.5em} \chi$  and  $\psi \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.5em} \chi$ , then  $\varphi \lor \psi \hspace{0.2em}\sim\hspace{-0.9em}\mid\hspace{0.5em} \chi$ .

\*Show that system C with **And** in place of **Cut** is is strictly weaker than C (define a non-monotonic consequence relation satisfying **Ref**, **LLE**, **RW**, **CMon**, and **And**, but not **Cut**.)

## Exercise: The Penguin Triangle

$$p =$$
 "penguin"  $b =$  "bird"  $f =$  "flies"

#### Suppose K contains:

1. 
$$p \sim b$$
 (penguins are normally birds)

2. 
$$p \sim \neg f$$
 (penguins normally do not fly)

3. 
$$b \sim f$$
 (birds normally fly)

Show semantically or by taking the closure of K under P:

- 1.  $b \sim \neg p$
- 2.  $b \lor p \succ f$
- 3.  $b \lor p \sim \neg p$

and explain informally why these are acceptable or problematic.