

# 13. Abductive Reasoning in Dynamic Epistemic Logic

Angel Nepomuceno–Fernández, Fernando Soler–Toscano, Fernando R. Velázquez–Quesada

This chapter proposes a study of abductive reasoning addressing it as an epistemic process that involves both an agent's information and the actions that modify this information. More precisely, this proposal presents and discusses definitions of an abductive problem and an abductive solution in terms of an agent's information (her knowledge and beliefs) and the involved epistemic actions (observation and belief revision). The discussion is then formalized with tools from dynamic epistemic logic; under such framework, the properties of the given definitions are studied, an epistemic action representing the application of an abductive step is introduced, and an illustrative example is provided. A number of the most interesting properties of abductive reasoning (those highlighted by Peirce) are shown to be better modeled within this approach.

13.1	<b>Classical Abduction</b> .....	270
13.2	<b>A Dynamic Epistemic Perspective</b> .....	272
13.2.1	What Is an Abductive Problem? .....	272
13.2.2	What Is an Abductive Solution? .....	273
13.2.3	How is the Best Explanation Selected? .....	273
13.2.4	How is the Best Explanation Incorporated Into the Agent's Information? .....	274
13.2.5	Abduction in a Picture .....	274
13.3	<b>Representing Knowledge and Beliefs</b> .....	275
13.3.1	Language and Models .....	275
13.3.2	Operations on Models .....	276
13.4	<b>Abductive Problem and Solution</b> .....	278
13.4.1	Abductive Problem .....	278
13.4.2	Classifying Problems .....	279
13.4.3	Abductive Solutions .....	279
13.4.4	Classifying Solutions .....	280
13.5	<b>Selecting the Best Explanation</b> .....	281
13.5.1	Ordering Explanations .....	282
13.6	<b>Integrating the Best Solution</b> .....	284
13.6.1	Abduction in a Picture, Once Again .....	285
13.6.2	Further Classification .....	285
13.6.3	Properties in a Picture .....	287
13.7	<b>Working with the Explanations</b> .....	287
13.7.1	A Modality .....	288
13.8	<b>A Brief Exploration to Nonideal Agents</b> .....	289
13.8.1	Considering Inference .....	290
13.8.2	Different Reasoning Abilities .....	290
13.9	<b>Conclusions</b> .....	290
	<b>References</b> .....	292

Within logic, abductive reasoning has been studied mainly from a purely syntactic perspective. Definitions of an abductive problem and its solution(s) are given in terms of a theory and a formula, and therefore most of the formal logical work on the subject has focused on:

1. Discussing what a theory and a formula should satisfy in order to constitute an abductive problem, and what a formula should satisfy in order to be an abductive solution [13.1]; see also Chap. 10
2. Proposing algorithms to find abductive solutions [13.2–6]
3. Analyzing the structural properties of abductive consequence relations [13.7–9].

In all these studies, which follow the so-called Aliseda–Kakas/Kowalski–Magnani/Meheus (AKM)-schema of abduction Chap. 10, explanationism and consequentialism are considered, but the epistemic character of abductive reasoning seems to have been

pushed into the background. Such character is considered crucial in this chapter, as it will be discussed.

This chapter's main proposal is an *epistemic and dynamic* approach to abductive reasoning. The proposal is close to the ideas of [13.10–13] in that it stresses the key role that agents play within the abductive reasoning scenario; after all, at the heart, abduction deals with agents and their (individual or collective) information. In this sense, this collaboration is closer to the Gabbay–Woods (GW)-schema [13.14, 15], see also Chap. 14, which is based on the concept of *ignorance problem* that arises when a cognitive agent has a cognitive target that cannot be attained from what she currently knows, and thus highlights the distinctive epistemic feature of abduction that is key to this chapter's considerations. Even so, this presentation goes one step further, as it fully adopts a dynamic perspective by making explicit the actions involved in the abductive process; after all, abduction studies the way agents react epistemically (as individuals or groupwise) to new observations.

More precisely, this proposal argues (Sect. 13.2) that abductive reasoning can be better understood as a *process* that involves *an agent's information*. To this end, it presents definitions of an abductive problem and an abductive solution in terms of an agent's knowledge and her beliefs as well as a *subjective* criteria for selecting *the agent's* best explanation, and outlines a policy through which the chosen abductive solution can be integrated into the agent's information. Then, the discussed ideas and definitions are formalized using tools from *dynamic epistemic logic* (DEL). This choice is not accidental: classical epistemic logic (EL [13.16, 17]) with its possible worlds semantic model is a powerful framework that allows to represent an agent's knowledge and beliefs not only about propositional facts but also about her own information. Its dynamic extension, DEL [13.18, 19], allows the representation of diverse epistemic actions (as diverse forms of announcements and different policies for belief revision) that make such

information change. Section 13.3 introduces the needed tools, and then the ideas and definitions discussed in Sect. 13.2 are formalized in Sects. 13.4, 13.5, 13.6, and 13.7. The chapter closes with a brief exploration (Sect. 13.8) of the epistemic and dynamic aspects of abductive reasoning that are brought to light when non-ideal agents are considered.

*Abductive reasoning* The concept of abductive reasoning has been discussed in various fields, and this has led to different ideas of what abduction should consist of (see [13.20], among others). For example, while certain authors claim that there is an abductive problem only when neither the observed  $\chi$  nor its negation follows from a theory [13.2], others say that there is also an abductive problem when, though  $\chi$  does not follow, its negation does [13.1], a situation that has been typically called a belief revision problem. There are also several opinions of what an abductive solution is. Most of the work on strategies for finding abductive solutions focuses on formulas that are already part of the system (the aforementioned [13.2–6]), while some others take a broader view, allowing not only changes in the underlying logical consequence relation [13.21] but also the creation and modification of concepts [13.22].

The present proposal focuses on a simple account: Abductive reasoning will be understood as a reasoning process that goes from a single unjustified fact to its abductive explanations, where an explanation is a formula of the system that satisfies certain properties. Still, similar epistemic and dynamic approaches can be made to other interpretations of abduction, as those that involve the creation of new concepts or changes in awareness [13.23, 24].

*Abductive reasoning in dynamic epistemic logic* This contribution is a revised version of a proposal whose different parts have been presented in diverse venues. While Sects. 13.2, 13.4, and 13.6 are based on [13.25], Sects. 13.5 and 13.7 are based on [13.26] and Sect. 13.8 is based on [13.27].

## 13.1 Classical Abduction

After Peirce's formulation of abductive reasoning (see [13.28] and Chap. 10), he immediately adds [13.29, p. 231] that:

“[The abductive solution] cannot be abductively inferred, or if you prefer the expression, cannot be abductively conjectured, until its entire content is already present in the premises, *If [the abductive solution] were true, [the abductive problem] would be a matter of course.*”

According to these ideas, abduction is a process that is triggered when a surprising fact is observed by an epistemic agent. Although the process returns an explicative hypothesis, the genuine result of an abductive inference is the plausibility of such hypothesis. The truth of the obtained hypothesis is thereby conjectured as plausible, which makes abduction an inferential process of a nonmonotonic character whose conclusion is rather a provisional proposal that could be revised in the light of new information.

When formalized within logical frameworks, the key concepts in abductive reasoning have traditionally taken the following form (Chap. 10). First, it is said that an abductive problem arises when there is a formula that does not follow from the current theory.

### Definition 13.1 Abductive problem

Let  $\Phi$  and  $\chi$  be a theory and a formula, respectively, in some language  $\mathcal{L}$ . Let  $\vdash$  be a consequence relation on  $\mathcal{L}$ :

- The pair  $(\Phi, \chi)$  constitutes a (*novel*) *abductive problem* when neither  $\chi$  nor  $\neg\chi$  are consequences of  $\Phi$ , that is, when

$$\Phi \not\vdash \chi \quad \text{and} \quad \Phi \not\vdash \neg\chi.$$

- The pair  $(\Phi, \chi)$  constitutes an *anomalous abductive problem* when, though  $\chi$  is not a consequence of  $\Phi$ ,  $\neg\chi$  is, that is, when

$$\Phi \not\vdash \chi \quad \text{and} \quad \Phi \vdash \neg\chi.$$

It is typically assumed that the theory  $\Phi$  is a set of formulas closed under logical consequence, and that  $\vdash$  is a truth-preserving consequence relation.

Consider a novel abductive problem. The observation of a  $\chi$  about which the theory  $\Phi$  does not have any opinion shows that  $\Phi$  is incomplete. Further information that completes  $\Phi$  making  $\chi$  a consequence of it solves the problem, as now the theory is strong enough to *explain*  $\chi$ . Consider now an anomalous abductive problem. The observation of a  $\chi$  whose negation is entailed by the theory shows that the theory contains a mistake. Now two steps are needed. First, perform a *theory revision* that stops  $\neg\chi$  from being a consequence of  $\Phi$ ; this turns the anomalous problem into a novel one, and now the search for further information that completes the theory, making  $\chi$  a consequence of it, can be performed. Here are the formal definitions.

### Definition 13.2 Abductive solution

- Given a *novel* abductive problem  $(\Phi, \chi)$ , the formula  $\eta$  is said to be an *abductive solution* when

$$\Phi \cup \{\eta\} \vdash \chi.$$

- Given an *anomalous* abductive problem  $(\Phi, \chi)$ , the formula  $\eta$  is an *abductive solution* when it is possible to perform a theory revision to get a *novel* problem  $(\Phi', \chi)$  for which  $\eta$  is a solution.

This definition of an abductive solution is often considered as too weak:  $\eta$  can take many trivial forms, as anything that contradicts  $\Phi$  (then everything, including  $\chi$ , follows from  $\Phi \cup \{\eta\}$ ) and even  $\chi$  itself (clearly,  $\Phi \cup \{\chi\} \vdash \chi$ ). Further conditions can be imposed in order to define more satisfactory solutions; here are some of them [13.1] (Chap. 10).

### Definition 13.3 Classification of abductive solutions

Let  $(\Phi, \chi)$  be an abductive problem. An abductive solution  $\eta$  is

consistent	iff	$\Phi, \eta \not\vdash \perp$
explanatory	iff	$\eta \not\vdash \chi$
minimal	iff	for every other solution $\zeta$ , $\eta \vdash \zeta$ implies $\zeta \vdash \eta$

The *consistency* requirement discards solutions that are inconsistent with the theory, something a reasonable explanation should not do. In a similar way, the *explanatory* requirement discards those explanations that would justify the problem by themselves, since it is preferred that the explanation only complements the current theory. Finally, the *minimality* requirement works as Occam's razor, looking for the simplest explanation: A solution is minimal when it is in fact logically equivalent to any other solution it implies. For further details on these definitions, the reader is referred to Chap. 10.

## 13.2 A Dynamic Epistemic Perspective

The present contribution proposes an approach to abductive reasoning from an epistemic and dynamic perspective. Instead of understanding abductive reasoning as a process that *modifies a theory* whenever *there is a formula that is not entailed by the theory under some particular consequence relation*, as the traditional definition of an abductive problem does, the proposed approach understands abductive reasoning as a process that *changes an agent's information* whenever, *due to some epistemic action, the agent has come to know or believe a fact that she could not have predicted otherwise*.

Such an epistemic and dynamic approach is natural. First, abduction, as other forms of nonmonotonic reasoning (e.g., belief revision, default reasoning), is classified as form of a common-sense reasoning rather than a mathematical one, and most of its classic examples involve *real* agents and their information (e.g., Mary observes that the light does not go on; Karen observes that the lawn is wet; Holmes observes that Mr. Wilson's right cuff is very shiny). Thus, even though abductive reasoning has been linked to scientific theories (as interpreted in philosophy of science), in its most basic forms it deals with an agent's (or a set of agents') information. Second, abductive reasoning implies a change in the agent's information (Mary assumes that the electricity supply has failed; Karen assumes it has rained; Holmes assumes Mr. Wilson has done a lot of writing lately), and thus it is essential to distinguish the different stages during the abductive process: the stage before the observation, the stage after the observation has raised the abductive problem (and thus the one when the agent starts looking for an explanation), and the stage in which the explanation that has been chosen is incorporated into the agent's information. This describes, of course, a dynamic process.

There is a final issue that is crucial for an epistemic approach to abductive reasoning. From this contribution's perspective, abductive reasoning involves not one epistemic attitude (as is typically assumed in most approaches) but rather (at least) two: that of those propositions about which the agent has full certainty; and that of those propositions that she considers very likely but she still cannot be certain about. The reason is that an agent typically tries to explain facts she has come to *know* due to some observation, but the chosen solution, being a *hypothesis* that might be dropped in the light of further observations, should not attain the full certainty status. The use of different epistemic notions also gives more flexibility to deal with a wider variety of abductive problems and abductive solutions,

making the analysis closer, from the authors' perspective, to Peirce's original formulation.

All in all, the abductive process can be studied by asking four questions:

1. What is an abductive problem?
2. What is an abductive solution?
3. How is the *best* solution(s) selected?
4. How does the agent assimilate the chosen solution(s)?

In the following, answers to these questions are discussed.

### 13.2.1 What Is an Abductive Problem?

There are, from an epistemic and dynamic perspective, two important concepts in the definition of an abductive problem. The first is what a formula  $\chi$  should satisfy in order to become an abductive problem. The second is the action that triggers the abductive problem, that is, the action that turns a formula  $\chi$  into an abductive problem.

For the former concept, a formula is typically said to be an abductive problem when it is surprising. There are different ways to define *a surprising observation of  $\chi$*  (some of them in a DEL setting [13.30]). Most of the approaches that define this notion in terms of what the agent knows (believes) understand a surprise as something that does not follow from such knowledge (beliefs). In other words, it is said that a given  $\chi$  is surprising whenever the agent does not know (believe) it, or, more radically, whenever the agent knows (believes)  $\neg\chi$ .

Now, note how in the context of abductive reasoning it is not reasonable to define a surprising observation in terms of what the agent knows (believes) *after such epistemic action*. The reason is that, after observing  $\chi$ , an agent would typically come to know (believe) it. Thus, if the mentioned definitions are followed focusing on the agent's information after the observation, no  $\chi$  would be surprising and there would be no abductive problems at all! It is more reasonable to define a surprising observation not in terms of what the agent knows (believes) as a result of the observation, but rather in terms of what she knew (believed) *before* it. More precisely, it will be said that a known (believed)  $\chi$  is surprising with respect to an agent whenever she could not have come to know (believe) it.

Of course, the meaning of the sentence *the agent could have come to know (believe)  $\chi$*  still needs to be clarified. This is a crucial notion, as it will indicate not only when a formula  $\chi$  is an abductive problem (the

agent could not have come to know (believe)  $\chi$ ), but also what a formula  $\eta$  needs in order to be an abductive solution (with the help of  $\eta$ , the agent could have come to know (believe)  $\chi$ ). Here the ability to come to know (believe) a given formula will be understood as the ability to infer it, and the simplest way to state this idea is the following: An agent could have come to know (believe)  $\chi$  if and only if there is an implication  $\eta \rightarrow \chi$  such that the agent knew both the implication and its antecedent. Other formulations that do not use the material implication  $\rightarrow$  are also possible (e.g., the agent may know both  $\neg\eta \vee \chi$  and  $\eta$  to come to know  $\chi$ ), but in the semantic model this contribution uses (Sect. 13.3), they are logically equivalent to the proposed one.

With respect to the action that triggers an abductive problem  $\chi$ , this action is typically assumed to be the observation of  $\chi$  itself. Here a more general idea will be considered: The action that triggers the abductive problem will be simply the observation of *some* formula  $\psi$ . Thus, though  $\psi$  should indeed be related to  $\chi$  (after all,  $\chi$  is an abductive problem because the agent comes to know  $\chi$  by observing  $\psi$ ), the agent will not be restricted to look for explanations of the formula that has been observed: She will also be able to look for explanations of *any* formula  $\chi$  she has come to know (believe) through the observation but could not have come to know (believe) by herself before. Note how other actions are also reasonable, as the agent might want to explain a belief she attained after a belief revision (Sect. 13.4.1).

Here is the intuitive definition of an abductive problem in full detail:

“Let  $s_1$  represent the epistemic state of an agent, and let  $s_2$  be the epistemic state that results from the agent observing some given  $\psi$ . A formula  $\chi$  constitutes an abductive problem for the agent at  $s_2$  whenever  $\chi$  is known and there is no implication  $\eta \rightarrow \chi$  such that the agent knew both the implication and its antecedent at  $s_1$ .”

It is important to emphasize how an abductive problem has been defined *with respect to an agent and stage* (i.e., *some epistemic situation*). Thus, whether a formula is an abductive problem depends on the formula *but also on the information* of that given agent at that given stage. The definition is given purely in terms of the agent’s knowledge, but it can also be given purely in terms of her beliefs, or even in terms of both, as it will be seen later.

The presented definition could seem very restrictive. Even if the reader agrees with the basic idea ( $\chi$  is an abductive problem for a given agent whenever she knows  $\chi$  but she could not have come to know (believe) it), she/he does not need to agree with the way key parts of it are understood. Nevertheless, as stated

in the introduction, this contribution does not intend on providing a full account of the many different understandings of what abductive reasoning does. Rather, its aim is to show how an epistemic and dynamic perspective can shed a new light on the way abductive reasoning is understood, even when assuming its simplest interpretation.

### 13.2.2 What Is an Abductive Solution?

In this proposal’s setting, an abductive solution for a given  $\chi$  will be defined in terms of what the agent could have been able to infer *before* the observation that raised the problem. As mentioned before, it will be said that  $\eta$  is a solution for the abductive problem  $\chi$  when the agent could have come to know (believe)  $\chi$  with the help of  $\eta$ . In this simple case in which the ability to come to know (believe) a given formula is understood as the ability to infer the formula by means of a simple modus ponens step, the following definition is obtained:

“A formula  $\eta$  constitutes an abductive solution for the abductive problem  $\chi$  at some given state  $s_2$  if the agent knew  $\eta \rightarrow \chi$  at the previous state  $s_1$ . Thus, the set of solutions for an abductive problem  $\chi$  is the set of antecedents of implications which have  $\chi$  as consequent and were known before the observation that triggered the abductive problem.”

Note how abductive solutions are looked for not when the agent has come to know (believe)  $\chi$ , but rather at the stage immediately before it. Thus,  $\eta$  is a solution when, had it been known (believed) before, would have allowed the agent to come to know (believe) (to predict/expect)  $\chi$ .

### 13.2.3 How is the Best Explanation Selected?

Although there are several notions of explanation for modeling the behavior of why-questions in scientific contexts (e.g., the law model, the statistical relevance model, or the genetic model), most of these consider a consequence (entailment) relation; *explanation* and *consequence* go typically hand in hand. However, finding suitable and reasonable criteria for selecting *the best* explanation has constituted a fundamental problem in abductive reasoning [13.31–33], and in fact many authors consider it to be the heart of the subject. Many approaches are based on logical criteria, but beyond requisites to avoid triviality and certain restrictions to the syntactic form, the definition of suitable criteria is still an open problem. Some approaches have suggested the use of *contextual aspects*, such as an ordering among formulas or among full theories. In particular,

for the latter, a typical option is the use of *preferential models* based on qualitative properties that are beyond the pure causal or deductive relationship between the abductive problem and its abductive solution. However, these preference criteria are seen as an external device which works on top of the deductive part of the explanatory mechanism, and as such they have been criticized because they seem to fall outside the logical framework.

Approaching abductive reasoning from an epistemic point of view provides a different perspective. It has been discussed already how the explanation an agent will choose for a given abductive problem does not depend on how the problematic formula could have been predicted, but rather on how *the agent* could have predicted it. In general, different agents have different information, and thus they might disagree in what each one calls *the best explanation* (and even in what each one calls *explanation* at all). This suggests that, instead of looking for criteria to select *the best explanation*, the goal should be a criterion to select *the agent's best explanation*. Now, once the agent has a set of formulas that explain the abductive problem from her point of view, how can she choose the best? This proposal's answer makes use of the fact that the considered agents have not only knowledge but also beliefs: Among all these explanations, some are more plausible than others *from her point of view*. These are precisely the ones the agent will choose when trying to explain a surprising observation: The best explanation can be defined in terms of a preference ordering among the agent's epistemic possibilities. It could be argued that this criterion is not *logical in the classic sense* because it is not based exclusively on the *deductive* relationship between the observed fact and different ways in which it could have been derived. Nevertheless, it is *logical in a broader sense* since it does depend on the agent's information: her knowledge and, crucially, her beliefs.

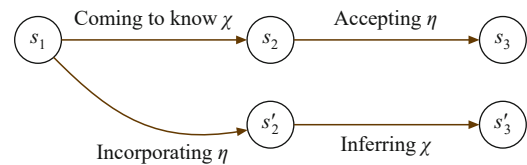
### 13.2.4 How is the Best Explanation Incorporated Into the Agent's Information?

Once the best explanation has been selected, it has to be incorporated into the agent's information. One of

the features that distinguishes abductive reasoning from deductive reasoning is its nonmonotonic nature: The chosen explanation does not need to be true, and in fact can be discarded in the light of further information. This indicates that an abductive solution cannot be assimilated as knowledge. Nevertheless, an epistemic agent has not only this hard form of information which is not subjected to modifications; she also has a soft form that can be revised as many times as it is needed: *beliefs*. Therefore, once the best abductive solution  $\eta$  has been chosen, the agent's information can be changed, leading her to *believe* that  $\eta$  is the case.

### 13.2.5 Abduction in a Picture

It is interesting to notice how the stated definitions of abductive problem and abductive solution rely on some form of *counterfactivity*, as in Peirce's original formulation (and also as discussed in [13.15]): A given  $\eta$  is a solution of a problem  $\chi$  if it would have allowed the agent to predict  $\chi$ . This can be better described with the following diagram.



The upper path is the real one: By means of an observation, the agent goes from the epistemic state  $s_1$  to the epistemic state  $s_2$  in which she knows  $\chi$ , and by accepting the abductive solution  $\eta$  she goes further to  $s_3$ . The existence of this path, the fact that  $\chi$  is an abductive problem and  $\eta$  is one of its abductive solutions, indicates that, at  $s_1$ , the lower path would have been possible: Incorporating  $\eta$  to the agent's information would have taken her to an epistemic state  $s'_2$  where she would have been able to infer  $\chi$ . Of course,  $s'_3$  is not identical to  $s_3$ : In  $s'_3$  both  $\eta$  and  $\chi$  are *equally reliable* because the second is inferred from the first, but in  $s_3$ ,  $\eta$  is less reliable than  $\chi$  since although the second is obtained via an observation, the first is just a hypothesis that is subject to revision in the light of further information.

### 13.3 Representing Knowledge and Beliefs

As mentioned, the most natural framework for formalizing the discussed ideas is that of DEL, the *dynamic* extension of *epistemic logic*. In particular, with the *plausibility models* of [13.34] it is possible to represent an agent’s knowledge and beliefs as well as acts of observation and belief revision, all of which are crucial to the stated understanding of the abductive process. This section introduces these needed tools; the discussed definitions will be formalized in Sect. 13.4.

#### 13.3.1 Language and Models

##### Definition 13.4 Language

Given a set of atomic propositions  $P$ , formulas  $\varphi$  of the language  $\mathcal{L}$  are given by

$$\varphi ::= p \mid \neg\varphi \mid \varphi \vee \varphi \mid \langle \leq \rangle \varphi \mid \langle \sim \rangle \varphi,$$

where  $p \in P$ . Formulas of the form  $\langle \leq \rangle \varphi$  are read as *there is a world at least as plausible (as the current one) where  $\varphi$  holds*, and those of the form  $\langle \sim \rangle \varphi$  are read as *there is a world epistemically indistinguishable (from the current one) where  $\varphi$  holds*. Other Boolean connectives ( $\wedge, \rightarrow, \leftrightarrow$ ) as well as the universal modalities,  $[\leq]$  and  $[\sim]$ , are defined as usual ( $[\leq]\varphi := \neg\langle \leq \rangle\neg\varphi$  and  $[\sim]\varphi := \neg\langle \sim \rangle\neg\varphi$  for the latter).

The modalities  $\langle \leq \rangle$  and  $\langle \sim \rangle$ , respectively, make it possible to define the notions of belief and knowledge within  $\mathcal{L}$ . The language’s semantic model, a *plausibility model*, is defined as follows.

##### Definition 13.5 Plausibility model

Let  $P$  be a set of atomic propositions. A *plausibility model* is a tuple  $M = \langle W, \leq, V \rangle$ , where:

1.  $W$  is a nonempty set of *possible worlds*
2.  $\leq \subseteq (W \times W)$  is a locally connected and conversely well-founded preorder, the agent’s *plausibility relation*, representing the plausibility order of the worlds from her point of view ( $w \leq u$  is read as *u is at least as plausible as w*)
3.  $V : W \rightarrow \wp(P)$  is an *atomic valuation function*, indicating the atoms in  $P$  that are true at each possible world.

A *pointed plausibility model*  $(M, w)$  is a plausibility model with a distinguished world  $w \in W$ .

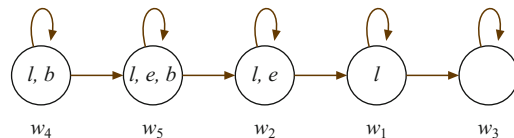
Before proceeding, recall that a relation  $R \subseteq (W \times W)$  is *locally connected* when every two elements that are  $R$ -comparable to a third are also  $R$ -comparable. It is *conversely well-founded* when there is no infinite  $\bar{R}$ -

ascending chain of elements in  $W$ , where  $\bar{R}$ , the *strict* version of  $R$ , is defined as  $\bar{R}wu$  iff  $Rwu$  and not  $Ruw$ . Finally, it is a *preorder* when it is reflexive and transitive.

The key idea behind plausibility models is that an agent’s beliefs can be defined as what is true *in the most plausible worlds from the agent’s perspective*, and modalities for the plausibility relation  $\leq$  will allow this definition to be formed. In order to define the agent’s knowledge, the approach is to assume that two worlds are epistemically indistinguishable for the agent if and only if she considers one of them at least as plausible as the other (if and only if they are comparable via  $\leq$ ). The *epistemic indistinguishability relation*  $\sim$  can therefore be defined as the union of  $\leq$  and its converse, that is, as  $\sim := \leq \cup \geq$ . Thus,  $\sim$  is the symmetric closure of  $\leq$  and hence  $\leq \subseteq \sim$ . Moreover, since  $\leq$  is reflexive and transitive,  $\sim$  is an *equivalence* relation. This epistemic indistinguishability relation  $\sim$  should not be confused with the *equal plausibility* relation, denoted by  $\simeq$ , and defined as the *intersection* of  $\leq$  and  $\geq$ , that is,  $\simeq := \leq \cap \geq$ . For further details and discussion on these models, their requirements and their properties, the reader is referred to [13.34, 35].

##### Example 13.1

The following diagram represents a plausibility model  $M$  based on the atomic propositions  $P := \{l, e, b\}$ . Circles represent possible worlds (named  $w_1$  up to  $w_5$ ), and each one of them includes exactly the atomic propositions that are true at that world (e.g., at  $w_2$ , the atomic propositions  $l$  and  $e$  are true, but  $b$  is false). Arrows represent the plausibility relation, with transitive arcs omitted (so  $w_4 \leq w_5 \leq w_2 \leq w_1 \leq w_3$ , but also  $w_4 \leq w_2, w_4 \leq w_1, w_4 \leq w_3$  and so on). Moreover,  $\sim$  is then the full Cartesian product, that is, for every worlds  $u$  and  $v$  in the model,  $u \sim v$ .



For the semantic interpretation, the two modalities  $\langle \leq \rangle$  and  $\langle \sim \rangle$  are interpreted with the help of their respective relations in the standard modal way.

##### Definition 13.6 Semantic interpretation

Let  $(M, w)$  be a pointed plausibility model with  $M = \langle W, \leq, V \rangle$ . Atomic propositions and Boolean operators

are interpreted as usual. For the remaining cases,

$$\begin{aligned} (M, w) \Vdash \langle \leq \rangle \varphi &\text{ iff } \exists u \in W \text{ s.t. } w \leq u \text{ \& } (M, u) \Vdash \varphi \\ (M, w) \Vdash \langle \sim \rangle \varphi &\text{ iff } \exists u \in W \text{ s.t. } w \sim u \text{ \& } (M, u) \Vdash \varphi. \end{aligned}$$

### Defining Knowledge and Beliefs

The notion of knowledge in plausibility models is defined by means of the epistemic indistinguishability relation in the standard way: The agent knows  $\varphi$  at some world  $w$  if and only if  $\varphi$  is the case in *every world she considers to be epistemically possible from  $w$* . (This makes knowledge a very strong notion, corresponding to an “absolutely unrevisable belief” [13.34]). The modality  $[\sim]$  can be used to this end. For the notion of beliefs, the idea is, as stated before, that the agent believes  $\varphi$  at a given  $w$  if and only if  $\varphi$  is the case *in the most plausible worlds from  $w$* . Thanks to the properties of the plausibility relation (a locally connected and conversely well-founded preorder),  $\varphi$  is true in the most plausible (i. e., the  $\leq$ -maximal) worlds from  $w$  if and only if, in accordance with the plausibility order, from some moment onward there are only  $\varphi$ -worlds (see [13.34, 36, 37] for the technical details). The modalities  $\langle \leq \rangle$  and  $[\leq]$  can be used to this end. Summarizing,

$$\begin{aligned} \text{The agent knows } \varphi & \quad K\varphi := [\sim]\varphi \\ \text{The agent believes } \varphi & \quad B\varphi := \langle \leq \rangle [\leq]\varphi \end{aligned}$$

Observe how, since  $\leq \subseteq \sim$ , the formula  $K\varphi \rightarrow B\varphi$  is valid (but its converse is not).

The dual of these notions, epistemic possibility and most likely possibility, can be defined as the correspondent modal duals

$$\hat{K}\varphi := \langle \sim \rangle \varphi \quad \hat{B}\varphi := [\leq] \langle \leq \rangle \varphi.$$

#### Example 13.2

Consider the plausibility model  $M$  of Example 13.1, and take  $w_2$  as the evaluation point. Since  $w_2 \sim u$  holds for every possible world  $u$  in the model, every world is epistemically possible from  $w_2$ 's perspective. But every world in the model satisfies  $b \rightarrow l$  (the implication is true at  $w_2, w_1$ , and  $w_3$  because the antecedent  $b$  is false, and true at  $w_4$  and  $w_5$  because the consequent  $l$  is true), so  $[\sim](b \rightarrow l)$ , that is,  $K(b \rightarrow l)$  is true at  $w_2$ : *The agent knows  $b \rightarrow l$  at  $w_2$* . On the other hand,  $\neg l$  is not true in every world, but it is true in  $w_3$ , the most plausible one from  $w_2$ 's perspective, so  $\langle \leq \rangle [\leq] \neg l$ , that is,  $B\neg l$ , is true at  $w_2$ : *The agent believes  $\neg l$  at  $w_2$* . Moreover, observe how  $b$  is neither known (it is not true in every

epistemically indistinguishable world) nor believed (it is not true in the most plausible worlds) at  $w_2$ . Still, it is true in some epistemic possibilities from  $w_2$  (e.g.,  $w_5$ ); hence,  $\langle \sim \rangle b$  (i. e.,  $\hat{K}b$ ) holds at  $w_2$ : *At that world, the agent considers  $b$  possible*.

A more detailed description of this framework, a number of the epistemic notions that can be defined within it, its technical details and its axiom system can be found in [13.34].

Following the DEL idea, actions that modify an agent's information can be represented as operations that transform the underlying semantic model. In the rest of this section, operations that can be applied over plausibility models will be recalled, and extensions of the language that allow to describe the changes such operations bring about will be provided. These will be used in Sect. 13.4 to represent and describe abductive reasoning.

### 13.3.2 Operations on Models

#### Update, Also Known as Observation

The most natural operation over Kripke-like semantic models is that of *update*. This operation reduces the domain of the model, and is typically given in terms of the formula the worlds should satisfy in order to survive the operation.

#### Definition 13.7 Update operation

Let the tuple  $M = \langle W, \leq, V \rangle$  be a plausibility model and let  $\psi$  be a formula in  $\mathcal{L}$ . The *update* operation yields the plausibility model  $M_{\psi!} = \langle W', \leq', V' \rangle$  where  $W' := \{w \in W \mid (M, w) \Vdash \psi\}$ ,  $\leq' := \leq \cap (W' \times W')$  and, for every  $w \in W'$ ,  $V'(w) := V(w)$ .

This operation reduces the domain of the model (preserving only those worlds that satisfy the given  $\psi$ ) and restricts the plausibility relation and the atomic valuation function accordingly. Since a submodel is obtained, the operation preserves the (universal) properties of the plausibility relation and hence it preserves plausibility models: If  $M$  is a plausibility model, then so is  $M_{\psi!}$ .

In order to describe the effects of an update within the language, existential modalities of the form  $\langle \psi! \rangle$  are used, for every formula  $\psi$ . Here is their semantic interpretation

$$\begin{aligned} (M, w) \Vdash \langle \psi! \rangle \varphi &\text{ iff } (M, w) \Vdash \varphi \\ &\text{ and } (M_{\psi!}, w) \Vdash \varphi. \end{aligned}$$

In words, an update formula  $\langle \psi! \rangle \varphi$  holds at  $(M, w)$  if and only if  $\varphi$  is the case (i. e., the evaluation point

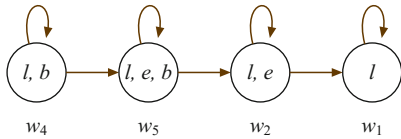


will survive the operation) and, after the update,  $\varphi$  is the case. The universal modality  $[\psi!]$  is defined as the modal dual of  $\langle\psi!\rangle$ , that is,  $[\psi!]\varphi := \neg\langle\psi!\rangle\neg\varphi$ .

In addition to being the most natural operation over Kripke-like models, an update also has a straightforward epistemic interpretation: it works as an act of a public announcement [13.38, 39] or, as it will be called here, an act of *observation*. When the agent observes a given  $\psi$ , she can discard those epistemically possible worlds that fail to satisfy this formula, thereby obtaining a model with only worlds that satisfied  $\psi$  before the operation. More details on this operation and its modalities (including an axiom system) can be found in the papers [13.38, 39] or in the textbooks [13.18, 19].

### Example 13.3

Consider the model  $M$  in Example 13.1 again. Suppose the agent observes  $l$ ; this can be modeled as an update with  $l$ , which yields the following model  $M_l$



The most plausible world in  $M$  has been discarded in  $M_l$ . As explained in Example 13.2, the agent believes  $\neg l$  in  $M$ , but after the observation this is not the case anymore:  $\neg l$  does not hold in the unique most plausible world of the new model  $M_l$ . In fact,  $\neg l$  does not hold in any epistemically possible world, and thus after the observation the agent knows  $l$ ; in symbols

$$(M_l, w_2) \models Kl, \text{ that is, } (M, w_2) \models [l!]Kl.$$

### Upgrade, Also Known as Belief Revision

Another natural operation over plausibility-like models is the rearrangement of worlds within an epistemic partition. Of course, there are several ways in which a new order can be defined. The following rearrangement, taken from [13.40], is one of the many possibilities.

#### Definition 13.8 Upgrade operation

Let the tuple  $M = \langle W, \leq, V \rangle$  be a plausibility model and let  $\psi$  be a formula in  $\mathcal{L}$ . The *upgrade* operation produces the plausibility model  $M_{\psi\uparrow} = \langle W, \leq', V \rangle$ , which differs from  $M$  just in the plausibility order, given now by

$$\leq' := \{(w, u) \mid w \leq u \text{ and } (M, u) \models \psi\} \cup \{(w, u) \mid w \leq u \text{ and } (M, w) \models \neg\psi\}$$

$$\{(w, u) \mid w \sim u, (M, w) \models \neg\psi \text{ and } (M, u) \models \psi\}.$$

The new plausibility relation states that after an upgrade with  $\psi$ , all  $\psi$ -worlds become more plausible than all  $\neg\psi$ -worlds, and within the two zones the old ordering remains [13.40]. More precisely, a world  $u$  will be at least as plausible as a world  $w$ ,  $w \leq' u$ , if and only if they already are of that order and  $u$  satisfies  $\psi$ , or they already are of that order and  $w$  satisfies  $\neg\psi$ , or they are comparable,  $w$  satisfies  $\neg\psi$  and  $u$  satisfies  $\psi$ . This operation preserves the properties of the plausibility relation and hence preserves plausibility models, as shown in [13.35].

In order to describe effects of this operation within the language, an existential modality  $\langle\psi\uparrow\rangle$  is introduced for every formula  $\psi$ ,

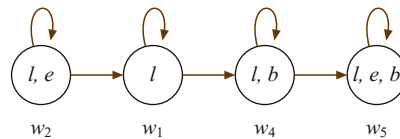
$$(M, w) \models \langle\psi\uparrow\rangle\varphi \text{ iff } (M_{\psi\uparrow, w}) \models \varphi.$$

In words, an upgrade formula  $\langle\psi\uparrow\rangle\varphi$  holds at  $(M, w)$  if and only if  $\varphi$  is the case after an upgrade with  $\psi$ . The universal modality  $[\psi\uparrow]$  is defined as the modal dual of  $\langle\psi\uparrow\rangle$ , as in the update case.

This operation also has a very natural epistemic interpretation. The plausibility relation defines the agent's beliefs, and hence any changes in the relation can be interpreted as changes in the agent's beliefs [13.34, 40, 41]. In particular, an act of *revising* beliefs after a reliable and yet fallible source has suggested  $\psi$  can be represented by an operation that puts  $\psi$ -worlds at the top of the plausibility order. Moreover, each one of the different methods to obtain a relation with former  $\psi$ -worlds at the top can be seen as a different policy for revising beliefs. Details on the operation and its modalities (including an axiom system) can be found in the papers [13.34, 40] or in the textbook [13.19].

### Example 13.4

Consider the model  $M_l$  in Example 13.3, that is, the model that results from the agent observing  $l$  at the initial model  $M$  in Example 13.1. Suppose the agent performs a belief revision toward  $b$ ; this can be modeled as an upgrade with  $b$ , which yields the following model  $(M_l)_{b\uparrow}$ :



The ordering of the worlds has changed, making those worlds that satisfy  $b$  ( $w_4$  and  $w_5$ ) more plausible than

those that do not ( $w_2$  and  $w_1$ ), keeping the old ordering with these two zones ( $w_5$  strictly above  $w_4$  and  $w_1$  strictly above  $w_2$ ). In  $M_{!}$  the agent believed  $\neg b \wedge \neg e$ , as such formula was the case in the model's unique most plausible world  $w_1$ , but this is not the case anymore in  $(M_{!})_{b\uparrow}$ : The unique most plausible world,  $w_5$ , satisfies

$b \wedge e$ , and thus the formula is part of the agent's beliefs. In symbols,

$$\begin{aligned} ((M_{!})_{b\uparrow}, w_2) \Vdash B(b \wedge e), \\ \text{that is, } (M_{!}, w_2) \Vdash [! \uparrow]B(b \wedge e). \end{aligned}$$

## 13.4 Abductive Problem and Solution

Given the intuitive definitions discussed and the formal tools introduced, it is now time to formalize the ideas.

### 13.4.1 Abductive Problem

First, the definition of what an abductive problem is.

#### Definition 13.9 Abductive problem

Let  $(M, w)$  be a pointed plausibility model, and consider  $(M_{\psi!}, w)$ , the pointed plausibility model that results from observing a given  $\psi$  at  $(M, w)$ .

A formula  $\chi$  is an abductive problem at  $(M_{\psi!}, w)$  if and only if it is known at such stage but it was not known before, that is, if and only if

$$(M_{\psi!}, w) \Vdash K\chi \quad \text{and} \quad (M, w) \Vdash \neg K\chi.$$

Equivalently, a formula  $\chi$  can become an abductive problem at  $(M, w)$  if and only if it is not known at such stage but will be known after observing  $\psi$ , that is, if and only if

$$(M, w) \Vdash \neg K\chi \wedge [\psi!]K\chi.$$

Note again how the definition of an abductive problem is relative to an agent's information at some given stage (the one represented by the pointed model).

There are two points worth emphasizing. First, note again how the definition distinguishes between the formula that becomes the abductive problem,  $\chi$ , and the formula whose observation triggers the abductive problem,  $\psi$ . Although these two formulas are typically understood to be the same ( $\chi$  becomes an abductive problem after being observed), the choice in this contribution is to distinguish between them. One reason for this is technical: Here the idea is that the agent will look for explanations of formulas that she could not have known before the observation but knows afterward. However, stating this as *the agent knows  $\chi$  after observing it* is restrictive in the DEL setting as not every formula satisfies this condition. This is because the underlying EL framework is powerful enough to talk

about the knowledge an agent has not only about facts but also about her own knowledge, and so there are formulas expressing situations such as *it is raining and you do not know it* ( $r \wedge \neg Kr$ ), which can be observed but are not known afterward (now you know that it is raining!). Another reason is, as stated earlier, generality: The described agent will be able to look for explanations not only of the formulas she can observe, but also of those that she can come to know through an observation. Still, this choice does not imply that the observed formula and the one that becomes an abductive problem are unrelated: In order for the agent to know  $\chi$  after observing  $\psi$ , she must have known  $\psi \rightarrow [\psi!]\chi$  before the action. This is nothing but the *reduction axiom* for the knowledge modality in *public announcement logic*

$$[\psi!]K\chi \leftrightarrow (\psi \rightarrow K(\psi \rightarrow [\psi!]\chi)).$$

Second, the requirements Definition 13.9 asks for  $\chi$  to be an abductive problem are not exactly the ones stated in the previous section: The sentence *there is no implication  $\eta \rightarrow \chi$  such that, before  $\chi$  became an abductive problem, the agent knew both the implication and its antecedent* has been replaced by *the agent did not know  $\chi$  before  $\chi$  became an abductive problem*. The reason is that, in DEL, the agent's knowledge and beliefs are closed under logical consequence (still, small variations of the EL framework allows the representation of nonideal agents and their abductive reasoning; see Sect. 13.8), and in such setting the two statements are equivalent: If there is an  $\eta$  such that the agent knew  $\eta \rightarrow \chi$  and  $\eta$  before  $\chi$  became an abductive problem, then clearly she knew  $\chi$  too, and if she knew  $\chi$ , then there was a  $\eta$  such that  $\eta \rightarrow \chi$  and  $\eta$  were both known, namely  $\chi$  itself. In fact, the restatement of the requirement emphasizes that it is the observation of  $\psi$  what causes the agent to know  $\chi$  and hence what creates the abductive problem.

It is worthwhile to highlight how, although the definition of an abductive problem was given in terms of the agent's knowledge, it can also be given in terms of her beliefs: It also makes sense for her to look for explanations of what she has come to believe!

“A formula  $\chi$  is said to be an abductive problem at  $(M_{\psi \uparrow}, w)$  if and only if  $(M_{\psi \uparrow}, w) \Vdash B\chi$  and  $(M, w) \Vdash \neg B\chi$ .”

With this definition, a formula is a problem if it is believed now but was not believe before a belief revision with  $\psi$ . But not only that. The agent can also face abductive problems that combine knowledge and beliefs. For example, she can face an abductive problem with  $\chi$  if she does not know the formula at some stage but believes it after a belief revision with  $\psi$ :

“A formula  $\chi$  is said to be an abductive problem at  $(M_{\psi \uparrow}, w)$  if and only if  $(M_{\psi \uparrow}, w) \Vdash B\chi$  and  $(M, w) \Vdash \neg K\chi$ .”

The stated definition allows to describe several forms of abductive problems, all of which differ in the strength of the attachment of the agent to the problematic  $\chi$  (known, strongly believed, safely believed, believed, etc.) *after* the epistemic action (update, upgrade) and the strength of her attachment to the formula *before* the action.

### 13.4.2 Classifying Problems

As mentioned, some approaches classify an abductive problem  $\chi$  according to whether  $\chi$  or  $\neg\chi$  follows from the theory: If neither  $\chi$  nor  $\neg\chi$  follows, then  $\chi$  is called a *novel* abductive problem; if  $\chi$  does not follow but  $\neg\chi$  does, then  $\chi$  is called an *anomalous* abductive problem. Given the requirement *the agent did not know  $\chi$  before  $\chi$  became an abductive problem* ( $\neg K\chi$ ) in Definition 13.9, one could suggest *the agent knew  $\neg\chi$*  ( $K\neg\chi$ ) as an alternative, but since the definition also asks for  $\chi$  to be known after the observation in order to be an abductive problem, such suggestion turns out to be too strong for propositional formulas: If  $\neg\chi$  is propositional and the agent knows it at some stage, then every epistemic possibility satisfies  $\neg\chi$ . Thus, since no *epistemic* action can change the (propositional) formula's truth value, the only way for the agent to know  $\chi$  afterward is for the action to eliminate every epistemic possibility, making  $K\varphi$  true for *every* formula  $\varphi$  and thus turning the agent inconsistent. But even though it is not possible to classify abductive problems in terms of the knowledge the agent had about the formula before the observation, it is still possible (and more reasonable) to classify them by using weaker notions, such as beliefs. Here is one possibility.

**Definition 13.10** *Expected, novel and anomalous problems*

Suppose  $\chi$  is an abductive problem at  $(M_{\psi \uparrow}, w)$ . Then  $\chi$  is said to be:

- An *expected* abductive problem if and only if  $(M, w) \Vdash B\chi$
- An *novel* abductive problem if and only if  $(M, w) \Vdash \neg B\chi \wedge \neg B\neg\chi$
- An *anomalous* abductive problem if and only if  $(M, w) \Vdash B\neg\chi$ .

Many people would not call the first case an abductive problem: The observation is a confirmation rather than a surprise, and thus it does not need to trigger any further epistemic action. Nevertheless, the case shows how this proposal allows for such situations to be considered. In fact, the classification can be refined by considering further attitudes, such as the *safe beliefs* of [13.34] or the *strong beliefs* of [13.42] (both definable in  $\mathcal{L}$ ).

### 13.4.3 Abductive Solutions

An abductive solution is now to be defined. Here is a version that uses only the notion of knowledge.

**Definition 13.11** *Abductive solution*

Let  $(M, w)$  be a pointed plausibility model, and consider  $(M_{\psi \uparrow}, w)$ , the pointed plausibility model that results from observing  $\psi$  at  $(M, w)$ .

If at  $(M_{\psi \uparrow}, w)$  the formula  $\chi$  is an abductive problem, then  $\eta$  is an abductive solution if and only if the agent knew that  $\eta$  implied  $\chi$  before the observation, that is, if and only if

$$(M, w) \Vdash K(\eta \rightarrow \chi).$$

Equivalently, if at  $(M, w)$  the formula  $\chi$  *can become* an abductive problem, then  $\eta$  will be an abductive solution if and only if the agent knows that  $\eta$  implies  $\chi$ , that is, if and only if

$$(M, w) \Vdash K(\eta \rightarrow \chi).$$

Just as in the case of abductive problem, it is also possible to define an abductive solution in terms of weaker notions as beliefs. For example, while a very strict agent would accept  $\eta$  as explanation only when  $\eta \rightarrow \chi$  was known, a less strict agent could accept it when such implication was only believed.

It is worth emphasizing that, in the stated definition, a solution for a problem  $\chi$  (at some  $M_{\psi \uparrow}$ ) is a formula  $\eta$  such that  $\eta \rightarrow \chi$  is known not when the abductive problem has arisen (at  $M_{\psi \uparrow}$ ) but rather at the stage immediately before (at  $M$ ). This is because an explanation is a piece of information that would have allowed the agent to predict the surprising observation. In fact, if

an abductive solution for a problem  $\chi$  were defined as a formula  $\eta$  such that  $\eta \rightarrow \chi$  is known once  $\chi$  is an abductive problem (at  $M_{\psi!}$ ), then every formula  $\varphi$  would be a solution since (at  $M_{\psi!}$ )  $K\chi$  would be the case (because  $\chi$  is an abductive problem) and hence so would be  $K(\varphi \rightarrow \chi)$  for every formula  $\varphi$ .

Observe also how, again in the stated definition, if  $\eta$  is a solution for the abductive problem  $\chi$  (at some  $M_{\psi!}$ ), then  $\eta$  could not be known before the observation that triggered the problem (at  $M$ ). Otherwise, both  $K(\eta \rightarrow \chi)$  and  $K\eta$  would be the case at such stage ( $M$ ) and hence, by the closure under logical consequence of knowledge in EL, so would be  $K\chi$ , contradicting the fact that  $\chi$  is an abductive problem.

#### Proposition 13.1

Let  $\chi$  be an abductive problem and  $\eta$  be one of its abductive solutions, both at  $(M_{\psi!}, w)$ . Then,  $(M, w) \Vdash \neg K\eta$ .

### 13.4.4 Classifying Solutions

It is common in the literature to classify abductive solutions according to their properties (Chap. 10). For example (Definitions 13.2 and 13.3; again, see Chap. 10), given a surprising observation  $\chi$ , an abductive solution  $\eta$  is said to be:

- *Plain* when it is a solution
- *Consistent* when it does not contradict the agent's information
- *Explanatory* when it does not explain  $\chi$  by itself.

Similar properties can be described in the present setting. To begin with, the *plain* property simply states that  $\eta$  is an abductive solution; a definition that has been already provided (Definition 13.11).

For the *consistency* property, the intuitive idea is for the solution to be compatible with the agent's information. To this end, consider the following definition.

#### Definition 13.12 Consistent solution

Let  $\chi$  be an abductive problem and  $\eta$  be one of its abductive solutions, both at  $(M_{\psi!}, w)$ . It is said that  $\eta$  is a *consistent* solution if and only if the agent considers it possible at  $(M_{\psi!}, w)$ , that is, if and only if

$$(M_{\psi!}, w) \Vdash \hat{K}\eta.$$

Thus, a solution is consistent when it is epistemically possible. Note how this requirement is given in terms of the agent's information *after* the epistemic action that triggered the abductive problem, and not

before it. In fact, there are formulas that, in a given situation, are solutions according to the stated definition, and yet not epistemically possible once the abductive problem has been raised.

#### Fact 13.1

Not every abductive solution is consistent.

*Proof:* Let  $\eta$  and  $\chi$  be propositional formulas, and take a model  $M$  in which the agent considers at least one  $(\neg\eta \wedge \neg\chi)$ -world to be epistemically possible, with the rest of the epistemic possibilities being  $(\neg\eta \wedge \chi)$ -worlds. After observing  $\chi$ ,  $\neg\chi$ -worlds will be discarded and there will be only  $(\neg\eta \wedge \chi)$ -worlds left, thus making  $\chi$  itself an abductive problem (it is not known at  $M$  but it will be known at  $M_{\chi!}$ ) and  $\eta$  an abductive solution (every epistemic possibility at  $M$  satisfies  $\eta \rightarrow \chi$ , so the agent knows this implication). Nevertheless, there are no  $\eta$ -worlds at  $M_{\chi!}$ , and therefore  $\hat{K}\eta$  is false at such stage. ■

The *explanatory* property is interesting. The idea in the classic setting is to avoid solutions that imply the problematic  $\chi$  per se, such as  $\chi$  itself or any formula logically equivalent to it. In the current epistemic setting, this idea can be understood in a different way: A solution  $\eta$  is explanatory when the acceptance of  $\eta$  (which, as discussed, will be modeled via belief revision; see Sect. 13.6) changes the agent's information, that is, when the agent's information is different from  $(M_{\psi!}, w)$  to  $((M_{\psi!})_{\eta\uparrow}, w)$  (the model that results after integrating the solution  $\eta$ ). This assertion could be formalized by stating that the agent's information is the same in two pointed models if and only if the agent has the same knowledge in both, but this would be insufficient: The model operation representing an act of belief revision (the upgrade of Definition 13.8) is devised to change only the agent's beliefs (although certain knowledge, such as knowledge about beliefs, might also change). A second attempt would be to state that the agent's information is the same in two pointed models if and only if they coincide in the agent's knowledge and beliefs, but the mentioned operation can change a model without changing the agent's beliefs.

Within the current *modal* epistemic logic framework, a more natural way of specifying the idea of an agent having the same information in two models is via the notion of bisimulation.

#### Definition 13.13 Bisimulation

Let  $P$  be a set of atomic propositions and let  $M = \langle W, \leq, V \rangle$  and  $M' = \langle W', \leq', V' \rangle$  be two plausibility models based on this set. A nonempty relation  $Z \subseteq (W \times W')$  is called a *bisimulation* between  $M$  and  $M'$  (notation:  $M \leftrightarrow_Z M'$ ) if and only if, for every  $(w, w') \in Z$ :

- $V(w) = V'(w')$ , that is,  $w$  and  $w'$  satisfy the same atomic propositions
- If there is a  $u \in W$  such that  $w \leq u$ , then there is a  $u' \in W'$  such that  $w' \leq' u'$  and  $Zuu'$
- If there is a  $u' \in W'$  such that  $w' \leq' u'$ , then there is a  $u \in W$  such that  $w \leq u$  and  $Zuu'$ .

Two models  $M$  and  $M'$  are *bisimilar* (notation:  $M \leftrightarrow M'$ ) when there is a bisimulation between them, and two pointed models  $(M, w)$  and  $(M', w')$  are bisimilar (notation:  $(M, w) \leftrightarrow (M', w')$ ) when there is a bisimulation between  $M$  and  $M'$  containing the pair  $(w, w')$ .

This notion is significant because, under image-finiteness (a plausibility model is *image-finite* if and only if every world can  $\leq$ -see only a finite number of worlds), it characterizes modal equivalence, that is, it characterises models that satisfy exactly the same formulas in the modal language.

**Theorem 13.1**

Let  $P$  be a set of atomic propositions and let  $M = \langle W, \leq, V \rangle$  and  $M' = \langle W', \leq', V' \rangle$  be two image-finite plausibility models. Then  $(M, w) \leftrightarrow (M', w')$  if and only if, for every formula  $\varphi \in \mathcal{L}$ ,  $(M, w) \Vdash \varphi$  iff  $(M', w') \Vdash \varphi$ .

Now it is possible to state a formal definition of what it means for a solution to be explanatory.

**Definition 13.14 Explanatory solution**

Let  $\chi$  be an abductive problem and  $\eta$  be one of its abductive solutions, both at  $(M_{\psi!}, w)$ . It is said that  $\eta$  is an *explanatory* solution if and only if its acceptance changes the agent’s information, that is, if and

only if there is *no* bisimulation between  $(M_{\psi!}, w)$  and  $((M_{\psi!})_{\eta\uparrow}, w)$ .

This definition, devised in order to avoid solutions that explain the abductive problem per se, has pleasant side effects. In the abductive reasoning literature, a solution is called *trivial* when it is logically equivalent to the abductive problem  $\chi$  (i.e., when it is not explanatory) or when it is a contradiction (to the agent’s knowledge, or a logical contradiction). Under the given definition, every trivial solution is not explanatory: *Accepting any such solution will not change the agent’s information*. The reason is that, in both cases, the upgrade operation *will not make any change in the model*: In the first case because, after the observation, the agent knows the abductive problem formula, and hence every epistemically possible world satisfies it (as well as every formula logically equivalent to the problem); in the second case because *no* epistemically possible world satisfies it. In this way, this framework characterizes trivial solutions not in terms of their form, as is typically done, but rather in terms of their effect: *Accepting them will not give the agent any new information*.

In particular, this shows how the act of incorporating a contradictory explanation will not make the agent *collapse* and turn into someone that knows and believes everything, as happens in traditional approaches; thus, a logic of formal inconsistency (e.g., [13.43]; see also Chap. 15) is not strictly necessary. This is a consequence of two simple but powerful ideas:

1. Distinguishing an agent’s different epistemic attitudes
2. Assimilating an abductive solution not as knowledge, but rather as a belief.

## 13.5 Selecting the Best Explanation

Finding suitable and reasonable criteria for selecting the best explanation is a fundamental problem in abductive reasoning [13.32, 33], and in fact many authors consider this to be the heart of the subject. The so-called *thesis of purpose*, stated in [13.33], establishes that the aim of scientific abduction is:

1. To generate new hypotheses
2. To select hypotheses for further examination and testing.

Hence a central issue in scientific abduction is to provide methods for selecting. Because the true state of the world is unknown, selecting the best explanation requires more than just consistency with the available

information, and there are many proposals of what these extra criteria should be.

Some approaches are based on probabilistic measurements [13.44–46]. Even Sherlock Holmes advised that, in order to evaluate explanations, one should “balance probabilities and choose the most likely” (*The Hound of the Baskervilles*), but unfortunately explanations rarely come equipped with probabilities.

In abductive logic programming, a common strategy is to look for abductive solutions at the *dead ends* of prolog proofs [13.47]. Sound and complete procedures can be defined also by using stable models and answer sets [13.48, 49]. Apart from selection criteria based on consistency and integrity constraints, it is common to

start with a set of *abducible* predicates and select explanations built only from ground atoms using them (see Chap. 10 for more details on abductive logic programming).

There are also approaches that use logical criteria, but beyond the already mentioned requisites to avoid triviality, the definition of suitable criteria is still an open problem. One of the most pursued ideas is that of minimality, a concept that can be understood syntactically (e.g., [13.3] and [13.5] look for literals), semantically (a minimal explanation is equivalent to any other explanation it implies [13.1]), with respect to the set of possible explanations (the best explanation is the weakest, i.e., the one that is implied by the rest of them), and even with respect to the current information (the best explanation is the one that disrupt less the current information).

In fact, most logical criteria are based on restrictions on the logical form of the solutions but, as mentioned in [13.1], finer criteria to select between two equally valid solutions require contextual aspects. With this idea in mind some approaches have proposed to use an ordering among formulas [13.10, 50, 51] or among full theories (i.e., possible worlds [13.52, 53]). In particular, for the latter, a common option is the use of *preferential models* (e.g., [13.54]) in which preferential criteria for selecting the best explanation are regarded as qualitative properties that are beyond the pure causal or deductive relationship between the abductive problem and its abductive solution. But these preference criteria are normally treated as an external device, which works on top of the logical or deductive part of the explanatory mechanism, and thus it has been criticized because it seems to fall outside a logical framework.

The epistemic approach of this proposal provides with an interesting alternative. The concepts of an abductive problem and an abductive solution have been defined in terms of the agent's epistemic attitudes, so it is natural to use such attitudes as a criterion for selecting the best explanation. Consider, for instance, the following elaboration of an example presented in Chap. 10.

“Mary and Gaby arrive late to Mary's apartment; the light switch is pressed but the light does not turn on. Knowing that the apartment is old, Mary assumes a failure in the electric line as the explanation for the light not turning on. Gaby, on the other hand, does not have any information about the apartment, so she explains the light not turning on by assuming that the bulb is burned out.”

After pressing the switch, both Mary and Gaby observe that the light does not turn on. There are several explanations for this: It is possible that the electric line failed, as Mary assumed, but it can also be the case

that the bulb is burned out, as Gaby thinks, and it is even possible that the switch is faulty. Then, why do they choose a different explanation? The reason is that, though they both observe that the light does not turn on, they have different background information: Mary knows that the apartment is old, and hence she considers a failure in the electric line more likely than any other explanation, but Gaby does not have that piece of information, so for her a burned out bulb explains the lack of light better.

The example shows that, even when facing the same surprising observation (the light does not turn on), agents with different knowledge and beliefs may choose a different best explanation: While Mary assumes that the electric line has failed, Gaby thinks that the bulb is burned out. Both explanations are equally *logical* since either a failure on the electric line or else a burned out bulb is enough to explain why the light does not turn on. What makes Mary to choose the first and Gaby the second is that they have different knowledge and different beliefs. This suggests first, that, instead of looking for criteria to select *the* best explanation, the goal should be a criteria to select *the agent's* best explanation.

But there is more. The explanation an agent will choose for a given abductive problem depends not only on how the problematic formula could have been predicted, but also on what the agent herself knows and what she considers more likely to be the case. It could be argued that this criterion is not *logical in the classical sense* because it is not based exclusively on the *deductive* relationship between the observed fact and the different ways in which it could have been derived. Nevertheless, it is *logical in a broader sense* since it does depend on the agent's information: her knowledge and her beliefs. In particular, in the plausibility models framework, the agent's knowledge and beliefs are defined in terms of a plausibility relation among epistemic possibilities, so it is natural to use precisely this relation as a criterion for selecting each agent's best explanation(s).

This section presents a straightforward use of this idea. It discusses how the plausibility order among epistemic possibilities can be lifted to a plausibility order among formulas, thus providing a natural criterion to select the agent's best explanation. A generalization of this idea that works instead with all explanations will be discussed later (Sect. 13.7).

### 13.5.1 Ordering Explanations

A plausibility model provides an ordering among possible worlds. This order can be lifted to get an ordering among set of worlds, that is, an ordering among formulas of the language (with each formula seen as the set

of those worlds that make it true). The different ways in which such ordering can be defined has been studied in *preference logic* (see [13.55–57] or, for a more detailed exposition [13.58, Chap. 3.3]); this section recalls the main ideas, showing how they can be applied to the task of selecting the best explanation in abductive reasoning.

In general, an ordering among objects can be lifted to an ordering among sets of such objects in different ways. For example, one can say that the lifted ordering puts the set of objects satisfying the property  $\psi$  (the set of  $\psi$ -objects) over the set of objects satisfying the property  $\varphi$  (the set of  $\varphi$ -objects) when *there is a  $\psi$ -object that the original ordering among objects places above some  $\varphi$ -object* (a  $\exists\exists$  preference of  $\psi$  over  $\varphi$ ; see below). But one can be more drastic and say that the set of  $\psi$ -objects is above the set of  $\varphi$ -ones when the original ordering places *every  $\psi$ -object above every  $\varphi$ -one* (a  $\forall\forall$  preference of  $\psi$  over  $\varphi$ ). This quantification combination gives raise to the following possibilities

$\varphi \leq_{\exists\exists} \psi$	iff	<i>there is a <math>\varphi</math>-object <math>w</math> and there is a <math>\psi</math>-object <math>u</math> such that <math>w \leq u</math></i>
$\varphi \leq_{\forall\exists} \psi$	iff	<i>for every <math>\varphi</math>-object <math>w</math> there is a <math>\psi</math>-object <math>u</math> such that <math>w \leq u</math></i>
$\varphi \leq_{\forall\forall} \psi$	iff	<i><math>w \leq u</math> for every <math>\varphi</math>-object <math>w</math> and every <math>\psi</math>-object <math>u</math></i>
$\varphi \leq_{\exists\forall} \psi$	iff	<i>there is a <math>\varphi</math>-object <math>w</math> such that <math>w \leq u</math> for every <math>\psi</math>-object <math>u</math></i>

The first two orderings can be defined within the language  $\mathcal{L}$

$$\begin{aligned}\varphi \leq_{\exists\exists} \psi &:= \langle \sim \rangle (\varphi \wedge \langle \leq \rangle \psi) \\ \varphi \leq_{\forall\exists} \psi &:= [\sim] (\varphi \rightarrow \langle \leq \rangle \psi).\end{aligned}$$

The first formula indicates that there is a  $\psi$ -world that is at least as plausible as a  $\varphi$ -one,  $\varphi \leq_{\exists\exists} \psi$ , exactly when there is an epistemic possibility that satisfies  $\varphi$  and that can see an at least as plausible  $\psi$ -world. The second one only changes the first quantification (turning, accordingly, the conjunction into an implication): For *every  $\varphi$ -world* there is a  $\psi$ -world that is at least as plausible.

The last two orderings are not immediate. Given the formulas for the previous two orderings, one could propose  $[\sim](\varphi \rightarrow \langle \leq \rangle \psi)$  for the  $\forall\forall$  case, but this formula is not correct: It states that every world that is at least as plausible as any  $\varphi$ -world satisfies  $\psi$ , but it does not guarantee that *every  $\psi$ -world is indeed above every  $\varphi$ -world*:

1. There might be a  $\psi$ -world incomparable to some  $\varphi$ -one, and even if all worlds are comparable

2. There might be a  $\psi$ -world strictly below a  $\varphi$ -one ( $<$ , the strict version of  $\leq$ , is defined as  $w < u$  if and only if  $w \leq u$  and *not*  $u \leq w$ ).

The plausibility order is locally connected (i. e., inside each epistemic partition, every world is comparable to each other) so (1) cannot occur. Thus, a formula defining  $\leq_{\forall\forall}$  only needs to guarantee that no  $\psi$ -world is *strictly* below a  $\varphi$ -one; in other words, it needs to express that, given any  $\psi$ -world, every world that is strictly more plausible satisfies  $\neg\varphi$ . Such formula can be easily stated in a language that extends  $\mathcal{L}$  with a standard modality for the relation  $<$

$$\varphi \leq_{\forall\forall} \psi := [\sim](\psi \rightarrow \langle < \rangle \neg\varphi).$$

Finally, the  $\exists\forall$  ordering presents a similar situation. Following the first two cases one could propose  $\langle \sim \rangle (\varphi \wedge \langle < \rangle \psi)$ , but such formula is not appropriate, even in the current full-comparability case: It holds even when there are  $\psi$ -worlds below the chosen  $\varphi$ -one. In order to guarantee the existence of a  $\varphi$ -world that is at most as plausible as every  $\psi$ -world, the formula should state that every world that is strictly *less* plausible than the  $\varphi$ -world satisfies  $\neg\psi$ . Extending the language again, this time with a modality for  $>$ , makes such formula straightforward

$$\varphi \leq_{\exists\forall} \psi := \langle \sim \rangle (\varphi \wedge \langle > \rangle \neg\psi).$$

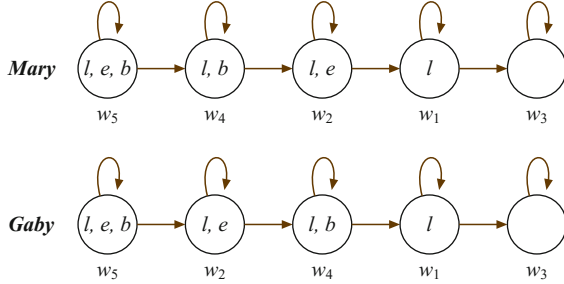
All in all, the important fact is that among these four orderings on sets of worlds (i. e., formulas), two are definable within  $\mathcal{L}$  and the other two only need simple extensions. This shows how the plausibility order among worlds that defines the agent's knowledge and beliefs (Sect. 13.3.1) also defines plausibility orderings among formulas (sets of worlds), and hence provides a criterion for selecting the best abductive solution *for a given agent*. It will now be shown how this criterion can be used, and how it leads to situations in which agents with different knowledge and beliefs choose different best explanations.

### Example 13.5

Recall Mary and Gaby's example. Both observe that after pressing the switch the light does not turn on, but each one of them chooses a different explanation: While Mary assumes that the electric line failed, Gaby thinks that the bulb is burned out. As it has been argued, the reason why they choose different explanations is that they have different knowledge and beliefs. Here is a formalization of the situation.

The following plausibility models show Mary and Gaby's knowledge and beliefs before pressing the

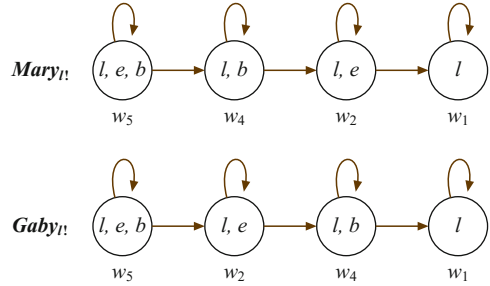
switch. They involve tree atomic propositions:  $l$  standing for lack of light,  $e$  standing for a failure in the electric line and  $b$  standing for a burned out bulb. Again, each possible world has indicated within it exactly those atomic propositions that are true in each one of them, and the arrows represent the plausibility relation (transitive arrows are omitted).



Observe how both Mary and Gaby know that both a failure on the electric line and a burned out bulb imply lack of light (both  $e \rightarrow l$  and  $b \rightarrow l$  hold in every world). In fact, the only difference in the models is the plausibility order between worlds  $w_2$  and  $w_4$ . Mary knows that the apartment is old so she considers a failure on the line ( $e$ ) more likely than a burned out bulb ( $b$ ), and hence the situation where the electric line fails but the bulb is not burned out ( $w_2$ ) is more likely than its opposite ( $w_4$ ). Gaby, on the other hand, does not know anything about the apartment, and hence for her a burned out bulb with a working electric line ( $w_4$ ) is more plausible than a working bulb and a failing electric line ( $w_2$ ). It is also assumed that, for both of them, the most likely possibility is the one in which everything works correctly ( $w_1$ ) and the least plausible case is the one in which everything fails ( $w_5$ ).

After they both observe that pressing the switch does not turn on the light, the unique world where  $l$  is

not the case,  $w_3$ , is eliminated, thus producing the following models.



As a result of the observation, Mary and Gaby know that there is no light ( $\neg l$  holds in both models), something that they did not know before. Thus, following Definition 13.9, both have an abductive problem with  $l$ .

According to Definition 13.11, both  $e$  and  $b$  are abductive solutions for the abductive problem  $l$  for both Mary and Gaby: Both formulas are the antecedent of implications that have  $l$  as a consequent and that were known before the observation. So, how can each girl choose her own best explanation? For Mary, the unique ordering that puts  $b$  above  $e$  is the weakest one,  $\exists\exists$  (there is a  $b$ -world,  $w_4$ , at least as plausible as a  $e$ -one,  $w_5$ ). Nevertheless, from her point of view,  $e$  is above  $b$  not only in the weak  $\exists\exists$  way ( $w_2$  is at least as plausible as  $w_4$ ) but also in the stronger  $\forall\exists$  way (every  $b$ -world has a  $e$ -world that is at least as plausible as it). Thus, one can say that  $e$  is a more plausible explanation from Mary's perspective. In Gaby's case something analogous happens:  $b$  is above  $e$  not only in the weak  $\exists\exists$  way ( $w_4$  is at least as plausible as  $w_2$ ) but also in the strong  $\forall\exists$  way. Hence, it can be said that, for Gaby,  $b$  is the best explanation.

### 13.6 Integrating the Best Solution

Once the agent has selected the best explanation for her, she can incorporate it into her information. As discussed in Sect. 13.2, even though the nonmonotonic nature of abductive reasoning indicates that an abductive solution should not be assimilated as knowledge, the richness of the present framework allows the possibility to integrate it as a part of the agent's beliefs. Here is a modality describing such action.

**Definition 13.15 Modality for abductive reasoning**  
 Let  $(M, w)$  be a pointed plausibility model and consider again  $(M_{\psi!}, w)$ , the pointed plausibility model

that results from observing  $\psi$  at  $(M, w)$ . Every pair of formulas  $\eta$  and  $\chi$  in  $\mathcal{L}$  define an existential modality  $\langle \text{Abd}_{\eta}^{\chi} \rangle \varphi$ , read as *the agent can perform an abductive step for  $\chi$  with  $\eta$  after which  $\varphi$  is the case*, and whose semantic interpretation is as follows

$$\begin{aligned}
 (M_{\psi!}, w) \models \langle \text{Abd}_{\eta}^{\chi} \rangle \varphi & \\
 \text{iff} & \\
 (1) (M_{\psi!}, w) \models K\chi \text{ and } (M, w) \models \neg K\chi, & \\
 (2) (M, w) \models K(\eta \rightarrow \chi), \text{ and} & \\
 (3) ((M_{\psi!})_{\eta \uparrow}, w) \models \varphi. &
 \end{aligned}$$



Equivalently,  $\langle \text{Abd}_\eta^x \rangle \varphi$ 's semantic interpretation can be defined as

$$\begin{aligned} (M_{\psi!}, w) \models \langle \text{Abd}_\eta^x \rangle \varphi \\ \text{iff} \\ (M, w) \models \neg K\chi \wedge K(\eta \rightarrow \chi) \wedge [\psi!](K\chi \wedge [\eta \uparrow]\varphi). \end{aligned}$$

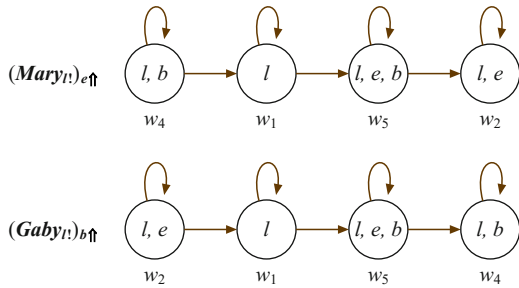
The definition states that  $\langle \text{Abd}_\eta^x \rangle \varphi$  is true at  $(M_{\psi!}, w)$  if and only if:

1.  $\chi$  is an abductive problem at  $(M_{\psi!}, w)$
2.  $\eta$  is an abductive solution also at  $(M_{\psi!}, w)$
3. An upgrade (Definition 13.8) with  $\eta$  will make  $\varphi$  true.

The last part makes precise the idea of how an agent should incorporate the selected explanation: It cannot be incorporated as knowledge, but it can be incorporated as a belief.

**Example 13.6**

Returning to Example 13.5, once Mary and Gaby have selected their respective best explanation, they can perform an abductive step. In Mary's case, worlds satisfying  $e$  ( $w_5$  and  $w_2$ ) will become more plausible than worlds that do not satisfy it ( $w_4$  and  $w_1$ ); in Gaby's case, worlds satisfying  $b$  ( $w_5$  and  $w_4$ ) will become more plausible than worlds that do not satisfy it ( $w_2$  and  $w_1$ ). Applying these upgrades to the models  $Mary_{\eta!}$  and  $Gaby_{\eta!}$  produces the following models.

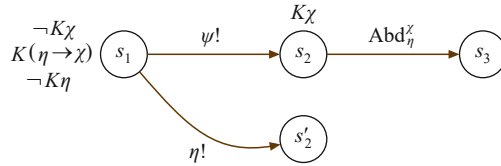


As a result of the abductive step, each agent believes her own explanation: Mary believes that the electric line has failed ( $e$  is true in her unique most plausible world  $w_2$ ), and Gaby believes that the bulb is burned out ( $b$  is true in her unique most plausible world  $w_4$ ). That is, for every  $w \in \{w_1, w_2, w_4, w_5\}$ ,

$$\begin{aligned} (Mary_{\eta!}, w) \models \langle \text{Abd}_e^l \rangle Be \\ (Gaby_{\eta!}, w) \models \langle \text{Abd}_b^l \rangle Bb. \end{aligned}$$

**13.6.1 Abduction in a Picture, Once Again**

The definitions that have been provided allow more precision in the diagram of abductive reasoning presented in Sect. 13.2.5. Here is the updated version for the case in which the definitions are given just in terms of the agent's knowledge. Note how the *inferring*  $\chi$  step has been dropped, as it is not needed in an omniscient setting such as DEL. Again, circles represent the agent's epistemic states (i. e., full plausibility models) and arrows are labeled with the operations that modify the agent's information.



Again, the upper path represents what really happened. After observing  $\psi$ , the agent reaches the epistemic state  $s_2$  in which she knows  $\chi$ . But before the observation, at  $s_1$ , she did not know  $\chi$ , and thus this formula is an abductive problem at  $s_2$ . Observe how  $\eta \rightarrow \chi$  was known at  $s_1$ : hence,  $\eta$  is an abductive solution at  $s_2$  and the agent can perform an abductive step with it to reach state  $s_3$ . This abductive solution  $\eta$  would have helped the agent to infer (and hence to come to know)  $\chi$ , and the lower path represents this alternative situation. In general, it cannot be guaranteed that the agent would have known  $\chi$  (or even  $\eta$ ) at state  $s'_2$ : these formulas could have had epistemic modalities, and hence the observation could have changed their truth value. However, if both formulas are propositional,  $K\chi$  and  $K\eta$  hold at  $s'_2$ .

**13.6.2 Further Classification**

Section 13.4.4 presented an epistemic version of the the common classification of abductive solutions. But the current DEL setting allows further possibilities and hence a finer classification. For example, here are two straightforward ideas. First, a solution  $\eta$  has been defined as the antecedent of an implication that has  $\chi$  as a consequent and that was known *before* the epistemic action that triggered the problem. Nevertheless, given that both formulas might contain epistemic operators, the agent can go from knowing the implication to not knowing it. Second, it has been stated that the agent incorporates the selected explanation via a belief revision (i. e., an upgrade). Nevertheless, since the solution might contain epistemic operators, the upgrade does not guarantee that the agent will believe the solution after the operation.

**Definition 13.16 Adequate solution and successful solution**

Let the formula  $\eta$  be an abductive solution for the abductive problem  $\chi$  at  $(M_{\psi!}, w)$ . Then:

- $\eta$  is an *adequate* solution if and only if the agent still knows  $\eta \rightarrow \chi$  at  $(M_{\psi!}, w)$ , that is, if and only if  $(M_{\psi!}, w) \Vdash K(\eta \rightarrow \chi)$ .
- $\eta$  is a *successful* solution if and only if it is believed after the abductive step, that is, if and only if  $(M_{\psi!}, w) \Vdash \langle \text{Abd}_{\eta}^{\chi} \rangle B\eta$ .

Here it is a result about the adequacy property.

**Proposition 13.2**

Every abductive solution is *adequate*.

*Proof:* More precisely, suppose that at  $(M_{\psi!}, w)$  the formula  $\chi$  is an abductive problem and  $\eta$  is one of its abductive solutions. Since  $\chi$  is an abductive problem,  $(M_{\psi!}, w) \Vdash K\chi$  and hence  $(M_{\psi!}, w) \Vdash K(\eta \rightarrow \chi)$ . ■

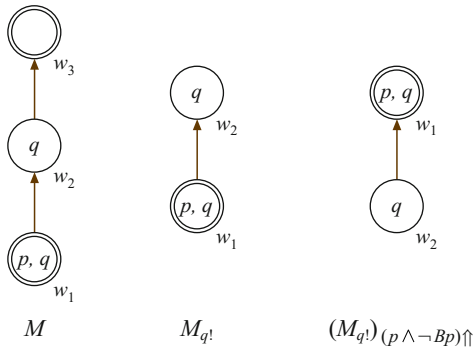
Given this result, this property is of little interest in the current setting. However, it becomes interesting in settings with nonomniscient agents. In such frameworks, it is possible for the agent not to know  $\eta \rightarrow \chi$  even when she knows  $\chi$  and she knew  $\eta \rightarrow \chi$  before.

Here it is another result, now about the property of being a successful solution.

**Fact 13.2**

Not every abductive solution is successful.

*Proof:* Let  $P = \{p, q\}$  be the set of atomic propositions, and consider the pointed plausibility models below (reflexive and transitive arrows omitted) in which the evaluation points are double circled.



Observe how  $q$  is an abductive problem at  $M_{q!}$  since it is not known at  $M$  (there is an epistemically possible world where  $q$  fails, namely,  $w_3$ ) but it is known at  $M_{q!}$ . Observe also how  $p \wedge \neg Bp$  is an abductive solution since  $K((p \wedge \neg Bp) \rightarrow q)$  holds at  $M$  (it is true at  $w_1$  and  $w_2$  because  $q$  is true in those worlds, and also true at  $w_3$  because  $p \wedge \neg Bp$  fails in this world). Furthermore,  $p \wedge \neg Bp$  is a consistent solution since it is epistemically possible in  $M_{q!}$  ( $p$  and  $\neg Bp$  are both true at  $w_1$ , the latter because there is a most plausible world,  $w_2$ , where  $p$  is not the case, and hence the agent does not believe  $p$ ). Nevertheless, after an upgrade with  $p \wedge \neg Bp$  this very formula is not believed. It fails at the unique most plausible world  $w_1$  because  $\neg Bp$  fails at it: the most plausible world ( $w_1$  itself) satisfies  $p$  and hence the agent now believes  $p$ , that is,  $Bp$  is the case. ■

Nevertheless, if a propositional solution  $\eta$  is also consistent, then it is successful.

**Proposition 13.3**

Suppose that at  $(M_{\psi!}, w)$  the formula  $\eta$  is an abductive solution for the abductive problem  $\chi$ . If  $\eta$  is a propositional and consistent solution, then it is successful.

*Proof:* If  $\eta$  is a consistent solution, then at  $(M_{\psi!}, w)$  there is at least one epistemically possible  $\eta$ -world. Therefore, an upgrade with  $\eta$  will put worlds that satisfied  $\eta$  in  $(M_{\psi!}, w)$  on top of the plausibility order. Now,  $\eta$  is propositional, and hence its truth value depends only on the valuation of each possible world; since the upgrade operation does not affect the valuation, then any world satisfying  $\eta$  in  $M_{\psi!}$  will still satisfy it in  $(M_{\psi!})_{\eta\uparrow}$ . Hence, after the operation, the most plausible worlds will satisfy  $\eta$ , and thus  $((M_{\psi!})_{\eta\uparrow}, w) \Vdash B\eta$  will be the case. This, together with the fact that at  $M_{\psi!}$  the formula  $\chi$  is an abductive problem and the formula  $\eta$  is an abductive solution, yield  $(M_{\psi!}, w) \Vdash \langle \text{Abd}_{\eta}^{\chi} \rangle B\eta$ . ■

It has been already stated that a solution is explanatory when it changes the agent's information. A further classification of abductive solutions can be provided according to *how much* they change the agent's information, that is, according to the attitude of the agent toward the solution *before* it was incorporated.

**Definition 13.17**

Suppose that  $\chi$  is an abductive problem at  $(M_{\psi!}, w)$ . An explanatory abductive solution  $\eta$  is:

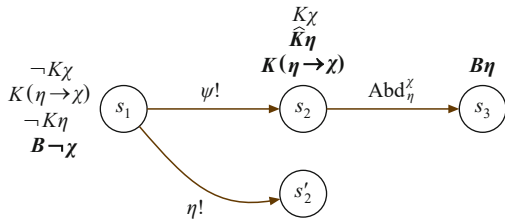
- *Weakly explanatory* when  $(M_{\psi!}, w) \Vdash B\eta$

- Neutral when  $(M_{\psi!}, w) \models \neg B\eta \wedge \neg B\neg\eta$
- Strongly explanatory when  $(M_{\psi!}, w) \models B\neg\eta$ .

Again, there are more possibilities if further epistemic attitudes are considered.

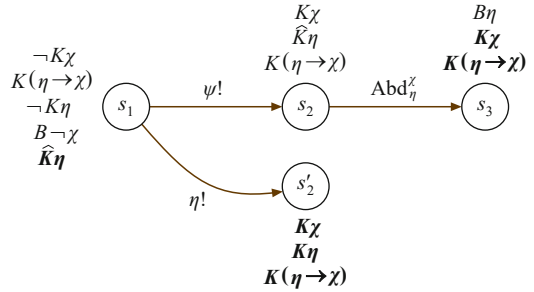
### 13.6.3 Properties in a Picture

Consider an anomalous abductive problem  $\chi$  (i. e.,  $B\neg\chi$  holds at  $s_1$ ) whose abductive solution  $\eta$  is consistent ( $\hat{K}\eta$  holds at  $s_2$ ) and successful ( $B\eta$  holds at  $s_3$ ), recalling also that every solution is adequate (so  $K(\eta \rightarrow \chi)$  holds at  $s_2$ ). This extends the diagram of Sect. 13.6.1 in the following way.



Moreover, consider the case in which both  $\chi$  and  $\eta$  are propositional, the typical case in abductive reasoning in which the agent looks for explanations of facts, and not of her own (or, in a multiagent setting, of other agents') epistemic state. First, in such case,  $\eta$  should be an epistemic possibility not only at  $s_2$  but also at  $s_1$ . But not only that; it is possible now to state the effects of the abductive step at  $s_2$  (the agent will believe  $\eta$  and will still know  $\eta \rightarrow \chi$ ) and of the hypothetical announcement of

$\eta$  at  $s_1$  (she would have known both  $\eta$  and  $\chi$ , and she would have still known  $\eta \rightarrow \chi$ ). Therefore,



This diagram beautifully illustrates what lies behind this proposal's understanding of abductive reasoning. In the propositional case, if  $\eta$  is a consistent and successful abductive solution for the abductive problem  $\chi$ , then, after abductive reasoning, the agent will know  $\chi$  and will believe  $\eta$ . In fact, when the observed formula  $\psi$  is actually the same  $\chi$  that becomes an abductive problem, the epistemic effect of abductive reasoning, from knowledge to beliefs, can be described with the following validity [13.59],

$$K(\eta \rightarrow \chi) \rightarrow [\chi!](K\chi \rightarrow \langle \text{Abd}_\eta^\chi \rangle B\eta) .$$

What makes  $\eta$  a reasonable solution is the existence of an *alternative reality* in which she observed  $\eta$  and, thanks to that, came to know  $\chi$ . Similar diagrams can be obtained for the cases in which the definitions of an abductive problem and an abductive solution are given in terms of epistemic attitudes other than knowledge.

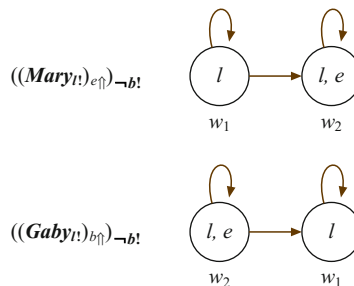
## 13.7 Working with the Explanations

The reason why abductive solutions are incorporated as beliefs and not as knowledge is because the selected explanation (in fact, *any* explanation) is just a hypothesis, subject to change in light of further information. Consider the following continuation of the Mary and Gaby's situation.

### Example 13.7

After their respective abductive steps (models  $(Mary!_l)_{e\uparrow}$  and  $(Gaby!_l)_{b\uparrow}$  of Example 13.6), Mary and Gaby take a closer look at the bulb and observe that it is not burned out ( $\neg b$ ). Semantically this is simply an observation operation that eliminates  $w_4$  and  $w_5$ , exactly those epistemic possibilities where the bulb is burned out (i. e., where  $b$  holds). The resulting models

are the following.



This observation does not affect Mary's explanation: She still believes that the electric line has failed ( $e$  is true in her unique most plausible world  $w_2$ ). But Gaby's

case is different: She does not have an explanation for  $l$  anymore. Although she knows it ( $Kl$  holds at the model on the bottom, that is,  $l$  is true in every epistemic possibility), she neither knows nor believes the antecedent of a known implication with  $l$  as a consequent (besides, of course, the trivial ones); she needs to perform a further abductive step in order to explain it.

There is, nevertheless, a way to avoid the extra abductive reasoning step. Recall that after applying the defined upgrade operation (Definition 13.8), all the worlds satisfying the given formula become more plausible than the ones that do not satisfy it, *and within the two zones the old ordering remains*. If the *lifted* worlds are not those that satisfy the agent's most plausible explanation but rather those that satisfy *at least one* of her explanations, the resulting model will have two layers: the lower one with worlds that do not satisfy any explanation, and the upper one with worlds that satisfy at least one. But inside the upper layer the old ordering will remain. In other words, the most plausible worlds in the resulting model (i. e., the most plausible ones in the upper layer) will be the ones that satisfy at least one explanation and that were already more plausible than the rest. Thus, with respect to the most plausible explanation, the same result is achieved: After such upgrade, roughly, the agent will believe the explanation that was the most plausible for her.

The difference with respect to the approach of the previous section is that the worlds that appear below the most plausible ones are not arbitrary. Worlds on the second best layer satisfy already some explanation; an explanation that was not chosen because it was not the most plausible one. Then, if further observations make the original best explanation obsolete, once that the correspondent (and now also obsolete) worlds have been discarded, the ones that will be at the top of the plausibility ordering will be the previously second best. Thus, an explanation will be already present and no further abductive steps will be needed.

### 13.7.1 A Modality

The idea just described is formalized now by introducing a modality that, given an abductive problem  $\chi$ , upgrades those worlds that satisfy at least one of its abductive explanations.

#### *Definition 13.18 Modality for formula-based abduction*

Let  $(M, w)$  be a pointed plausibility model and consider again  $(M_{\psi!}, w)$ , the pointed plausibility model that results from observing  $\psi$  at  $(M, w)$ . Every formula  $\chi$  in

$\mathcal{L}$  defines an existential modality of the form  $\langle \text{Abd } \chi \rangle \varphi$ , read as *the agent can perform a complete abductive step for  $\chi$  after which  $\varphi$  is the case*, and whose semantic interpretation is as follows

$$\begin{aligned} (M_{\psi!}, w) \Vdash \langle \text{Abd } \chi \rangle \varphi \\ \text{iff} \\ (1) (M_{\psi!}, w) \Vdash K\chi \text{ and } (M, w) \Vdash \neg K\chi, \\ (2) ((M_{\psi!})_{(\bigvee \Sigma_{\chi})\uparrow}, w) \Vdash \varphi, \end{aligned}$$

where  $\Sigma_{\chi}$  is the set of abductive solutions for  $\chi$ , that is,

$$\Sigma_{\chi} := \{\eta \mid (M, w) \Vdash K(\eta \rightarrow \chi)\}.$$

Equivalently,  $\langle \text{Abd } \chi \rangle \varphi$ 's semantic interpretation can be defined as

$$\begin{aligned} (M_{\psi!}, w) \Vdash \langle \text{Abd } \chi \rangle \varphi \\ \text{iff} \\ (M, w) \Vdash \neg K\chi \wedge [\psi!](K\chi \wedge [\bigvee \Sigma_{\chi} \uparrow]\varphi). \end{aligned}$$

The correspondent universal modality,  $[\text{Abd } \chi]$ , is defined as usual.

The definition states that  $\langle \text{Abd } \chi \rangle \varphi$  is true at  $(M_{\psi!}, w)$  if and only if (1)  $\chi$  is an abductive problem at  $(M_{\psi!}, w)$ , and (2) an upgrade with  $\bigvee \Sigma_{\chi}$  will make  $\varphi$  true. The last part makes precise the idea of working with all the solutions:  $\Sigma_{\chi}$  contains all abductive solutions for  $\chi$ , so  $\bigvee \Sigma_{\chi}$  is a disjunction characterising those worlds that satisfy at least one of them and hence an upgrade with it will move such worlds to the topmost layer. But inside this layer, the former plausibility order will persist, and hence worlds at the top of it will be precisely those that satisfy at least one solution for  $\chi$  and, among them, were already the most plausible ones.

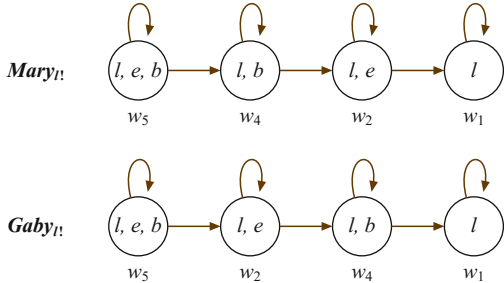
#### *Remark 13.1*

The set  $\Sigma_{\chi}$  contains, among others,  $\chi$ ,  $\chi \wedge \chi$ , and so on, and hence  $\bigvee \Sigma_{\chi}$  is an infinite disjunction. Syntactic restrictions can be imposed in order to avoid such situations (e.g., asking for solutions that are also minimal conjunctions of literals). Another possibility, closer to the semantic spirit of this approach, is to work with finite plausibility models, and then look for solutions among the formulas that characterize each possible world.

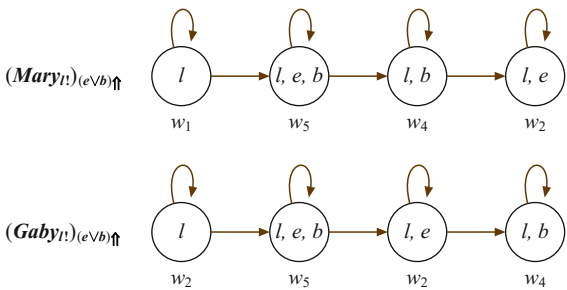
The following example shows how this new operation allows the agent to have ready another explanation in case the initially *best* one turns out to be incorrect.

**Example 13.8**

Let us go back to Mary and Gaby’s example all the way to the stage after which they have observed that the light does not turn on (models  $Mary_{!1}$  and  $Gaby_{!1}$  of Example 13.5, repeated here).



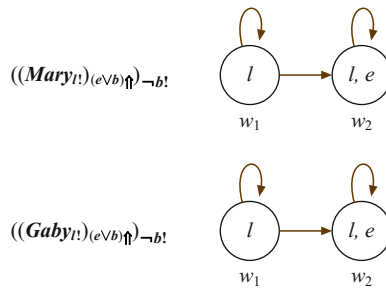
Suppose that, instead of selecting their respective most plausible explanation and assimilating it (as they did in Example 13.5), Mary and Gaby work with all their explanations: Instead of an upgrade with  $e$  for Mary and an upgrade with  $b$  for Gaby, both of them perform an upgrade with  $e \vee b$ . This produces the following models.



The worlds satisfying  $e \vee b$  ( $w_2$ ,  $w_4$ , and  $w_5$ ) have been upgraded. As a result of this, both Mary and Gaby have an explanation for  $l$ , but each one of them has *her own* explanation: While Mary believes that the electric line

has failed ( $e$  is the case in the most plausible world at the model on the top), Gaby believes that the bulb is burned out ( $b$  holds in the most plausible world at the model on the bottom).

So far the result of the upgrade is, with respect to Mary and Gaby beliefs, exactly the same as with the previous proposal where only worlds that satisfy the most plausible explanation are upgraded (in both cases,  $w_2$  and  $w_4$  are Mary’s and Gaby’s most plausible worlds, respectively). But note what happens now when they both observe that the bulb is in fact not burned out ( $\neg b$ ): Such action produces the following situation.



Again, the observation does not affect Mary’s explanation ( $e$  still holds in the most plausible world at model on the top), but it does change Gaby’s since her previous explanation  $b$  is not possible anymore. The difference is that now she does not need to perform an extra abductive step because she has already another explanation: She now believes that the electric line has failed ( $e$  holds in the most plausible world at model on the bottom).

Thus, after an upgrade with all explanations, what the agent will be lead to believe depends on her plausibility order, just as with the first proposal. Nevertheless, if further information invalidates such best explanation, the agent will believe the next to best one without the need of further abductive steps.

**13.8 A Brief Exploration to Nonideal Agents**

As most (if not all) proposals for representing a given phenomena, the presented epistemic and dynamic approach to abduction has made some assumptions for the sake of simplicity. One of the most important of these is the fact that agents whose information is represented within the plausibility framework are *ideal*: Their knowledge and beliefs are closed under logical consequence. This supposition is not exclusive of this approach; the classic logical definitions of abductive reasoning assume not only that the given set of formulas  $\Phi$ , the theory, is closed under logical con-

sequence, but also that  $\vdash$  is the logical consequence relation.

The present proposal highlights the epistemic nature of abductive reasoning, and so it is natural to ask how such reasoning process works for a different kind of agents, in particular, for those whose information does not need to have *ideal* properties and thus are, in that sense, closer to real computational agents with limited resources (and also closer to us human beings). This final section briefly discusses some ideas; further developments in this direction can be found in [13.60].

### 13.8.1 Considering Inference

Suppose Karl is in his dining room and sees smoke coming out of the kitchen. This seems unjustified at first, but then he realises that the chicken he placed on the fire has been there for a long time. Initially Karl did not have any explanation for the smoke, but after a moment he realized that such event was actually not surprising at all.

This case is different from the discussed ones because Karl is not an ideal agent: He does not have at hand all logical consequences of his information, and therefore he did not realize that the information he had before seeing the smoke was enough to predict it (i. e., to infer that there would be smoke). Described in more technical terms, seeing the smoke raised an abductive problem for Karl, but such problem arose because he did not have, at the time of the observation, all the logical consequences of the information he actually had (otherwise there would have been no abductive problem at all). Accordingly, in such case the abductive solution is not necessarily a piece of information that would have allowed Karl to predict the smoke; it might be a simple inference step that made *explicit* what was only *implicit* before.

This shows not only how agents whose information is not closed under logical consequence can face at least a new kind of abductive problem, but also how such problems give rise to a different kind of solutions.

### 13.8.2 Different Reasoning Abilities

In the previous example, the abductive solution was a simple inference step because Karl had the needed

reasoning tools to infer *there is smoke in the kitchen* from *the chicken has been on the fire for a long time*. But what if that was not the case? That is, what if, besides not having at hand all the logical consequences of his information, Karl did not have the required reasoning tools to infer some of them?

In such new situation, Karl faces again an abductive problem, but this time of a different nature. The surprising observation could have been predicted in the sense that it is a logical consequence of Karl's information *the chicken has been on the fire for a long time*, just as in the initial version of this example. The difference is that such observation is not something that Karl could have predicted by himself: He did not have the needed tools. One can say that, even though *there is smoke in the kitchen* is *objectively* derivable from the initial information, it is not *subjectively* derivable in the sense that Karl could not have done it. To put it in other words, besides not having at hand all the logical consequences of her actual information, Karl might not even be able to reach them.

Accordingly, the simple inference step of before cannot be a solution to the problem now, as Karl does not have the needed tools to perform it. One possible solution is, as in the traditional case, a piece of information that would have allowed Karl to predict the smoke from some other previously known fact, but a more interesting one is some reasoning tool that would have helped him to predict the fire from the known fact *the chicken has been on the fire for a long time*.

New cases arise when further kinds of agents are considered. A systematic study of such cases can be found in [13.61].

## 13.9 Conclusions

This chapter has proposed an epistemic and dynamic approach to abductive reasoning, understanding this form of reasoning as a process that:

1. Is triggered by an epistemic action through which the agent comes to know or believe certain  $\chi$  that otherwise she could not have been able to know or believe
2. Looks for explanations for  $\chi$  in the set of formulas that could have helped the agent to come to know or believe  $\chi$
3. Incorporates the chosen explanation as a part of the agent's beliefs.

Besides providing formal definitions of what an abductive problem and an abductive solution are in terms

of an agent's knowledge and beliefs, the present proposal has discussed:

1. A classification of abductive problems in terms of both how convinced the agent is of the problematic formula after the observation (she *knows* it, or just *believes* it) and how plausible the formula was *before* the epistemic action that triggered the problem
2. A classification of abductive solutions based not only on their deductive relation with the abductive problem or their syntactic form, but also in terms of both their plausibility *before* the problem was raised and the way it will affect the agent's information once they are incorporated
3. A new perspective that looks not for *the best* explanation but rather for *the agent's best* explanation,

and the possibility to carry out this search in terms of which explanations are more likely from the agent's point of view, that is, in terms of the agent's beliefs

4. The possibility of integrating the chosen solution into the agent's information as part of her beliefs, which allows not only to identify trivial solutions because of their effect rather than their form, but also to revise and eventually discard solutions that become obsolete in the light of further information.

Crucial for all these contributions has been the use of plausibility models and, in general, the DEL guidelines, which puts emphasis in the representation of both epistemic attitudes and the actions that affect them.

It is worthwhile to compare, albeit briefly, the present proposal to other epistemic approaches to abductive reasoning. Besides immediate differences in the respective semantic models (while other approaches follow the Alchourrón–Gärdenfors–Makinson (AGM) belief revision, using a set of formulas for representing the agent's information, here possible worlds are used), there are two main points that distinguish the presented ideas from other proposals. First, here several epistemic attitudes are taken into account, thus making a clear difference between what the agent holds with full certainty (knowledge) and what she considers very likely but still cannot guarantee (beliefs); this allows to distinguish between the certainty of both the previous information and the surprising observation, and the mere plausibility of the chosen solution (recall the validity  $K(\eta \rightarrow \chi) \rightarrow [\chi!](K\chi \rightarrow \langle \text{Abd}_\eta^x \rangle B\eta)$ , briefly discussed at the end of Sect. 13.6). Second, this approach goes one step further by making explicit the different stages of the abductive process, thus making also explicit the epistemic actions involved. This highlights the importance of actions such as *belief revision*, commonly understood in epistemic approaches to abduction as the one *triggered* by the abductive problem [13.12, 62], and also such as *observation*, understood here as the one that *triggers* the abductive process.

This chapter presents only the first steps toward a proper study of abductive reasoning from an epistemic and dynamic perspective, and several of the current proposals can be refined. For example, the specific definition of an abductive problem (Definition 13.9) relies on the fact that, within the DEL framework, agents are logically omniscient. As it has been hinted at in Sect. 13.8, in a nonomniscient DEL setting [13.35, 63] the ideas discussed in Sect. 13.2 would produce a different formal definition (which, incidentally, would allow to classify abductive problems and abductive solutions according to some *derivability* criteria). Moreover, it would be possible to analyze the full abductive picture presented in Sect. 13.2.1, which requires inference steps

in the alternative reality path. These extensions are relevant: They would allow a better understanding of the abductive process as performed by *real* agents.

But it is also possible to do more than just follow the traditional research lines in abductive reasoning, and here are two interesting possibilities (whose development exceeds the limits of this chapter). First, the DEL framework allows multiagent scenarios in which abductive problems would arise in the context of a community of agents. In such setting, further to the public observation and revision used here, actions that affect the knowledge and beliefs of different agents in different ways are possible. For example, an agent may be privately informed about  $\psi$ : If this raises an abductive problem  $\chi$  for her and another agent has private information about  $\eta \rightarrow \chi$ , they can interact to obtain the abductive solution  $\eta$ . Second, the DEL framework deals with high-order knowledge, thus allowing to study cases in which an agent, instead of looking for an explanation of a fact, looks for an explanation of her own epistemic state. Interestingly, explanations might involve epistemic actions as well as the lack of them.

According to those considerations, this logical approach takes into account the dynamics aspects of logical information processing, and one of them is abductive inference, one of the most important forms of inference in scientific practices. The aforementioned multiagent scenarios allow to model concrete practices, particularly those that develop a methodology based on observation, verification, and systematic formulation of provisional hypotheses, such as in empirical sciences, social sciences, and clinical diagnosis. The epistemological repercussions of this DEL approach is given by the conceptual resources that it offers, useful to model several aspects of explanatory processes. If known theories of belief revision, at the last resort, say nothing about context of discovery, by means of DEL the accessibility of this context to rational epistemological and logical analysis is extended, further on classical logical treatment of abduction. From the perspective of game theoretic semantics, for example, now it is easier to determine what rules are strategic and what are operatories when abductive steps were given. But applications should also be considered to tackle certain philosophical problems. For example, abductive scenarios within multiagent settings can be used to study the implications of different forms of communication within scientific communities.

**Acknowledgments.** The first author acknowledges the support of the project *Logics of discovery, heuristics and creativity in the sciences* (PAPIIT, IN400514-3), granted by the National Autonomous University of Mexico (UNAM).

## References

- 13.1 A. Aliseda: *Abductive Reasoning. Logical Investigations into Discovery and Explanation*, Synthese Library, Vol. 330 (Springer, Dordrecht 2006)
- 13.2 A.C. Kakas, R.A. Kowalski, F. Toni: Abductive logic programming, *J. Logic Comput.* **2**(6), 719–770 (1992)
- 13.3 M.C. Mayer, F. Pirri: First order abduction via tableau and sequent calculi, *Log. J. IGPL* **1**(1), 99–117 (1993)
- 13.4 M.C. Mayer, F. Pirri: Propositional abduction in modal logic, *Logic J. IGPL* **3**(6), 907–919 (1995)
- 13.5 A.L. Reyes-Cabello, A. Aliseda, Á. Nepomuceno-Fernández: Towards abductive reasoning in first-order logic, *Logic J. IGPL* **14**(2), 287–304 (2006)
- 13.6 S. Klarman, U. Eudriss, S. Schlobar: ABox abduction in the description logic ACC, *J. Autom. Reason.* **46**, 43–80 (2011)
- 13.7 J. Lobo, C. Uzcátegui: Abductive consequence relations, *Artif. Intell.* **89**(1/2), 149–171 (1997)
- 13.8 A. Aliseda: Mathematical reasoning vs. abductive reasoning: A structural approach, *Synthese* **134**(1/2), 25–44 (2003)
- 13.9 B. Walliser, D. Zwirn, H. Zwirn: Abductive logics in a belief revision framework, *J. Log. Lang. Info.* **14**(1), 87–117 (2004)
- 13.10 H.J. Levesque: A knowledge-level account of abduction, *Proc. 11th Intl. Joint Conf. on Artif. Intell.*, ed. by N.S. Sridharan (Morgan Kaufmann, Burlington 1989), pp. 1061–1067, Detroit 1989
- 13.11 C. Boutilier, V. Becher: Abduction as belief revision, *Artif. Intell.* **77**(1), 43–94 (1995)
- 13.12 A. Aliseda: Abduction as epistemic change: A Peircean model in artificial intelligence. In: *Abduction and Induction: Essays on Their Relation and Integration*, Applied Logic, ed. by P.A. Flach, A.C. Kakas (Kluwer, Dordrecht 2000) pp. 45–58
- 13.13 L. Magnani: *Abductive Cognition: The Epistemological and Eco-Cognitive Dimensions of Hypothetical Reasoning*, Cognitive Systems Monographs, Vol. 3 (Springer, Heidelberg 2009)
- 13.14 D. Gabbay, J. Woods (Eds.): *The Reach of Abduction: Insight and Trial, A Practical Logic of Cognitive Systems*, Vol. 2 (Elsevier, Amsterdam 2005)
- 13.15 J. Woods: Cognitive economics and the logic of abduction, *Rev. Symb. Log.* **5**(1), 148–161 (2012)
- 13.16 J. Hintikka: *Knowledge and Belief: An Introduction to the Logic of the Two Notions* (Cornell Univ. Press, Ithaca 1962)
- 13.17 R. Fagin, J.Y. Halpern, Y. Moses, M.Y. Vardi: *Reasoning About Knowledge* (MIT Press, Cambridge 1995)
- 13.18 H. van Ditmarsch, W. van der Hoek, B. Kooi: *Dynamic Epistemic Logic*, Synthese Library, Vol. 337 (Springer, Dordrecht 2007)
- 13.19 J. van Benthem: *Logical Dynamics of Information and Interaction* (Cambridge Univ. Press, Cambridge 2011)
- 13.20 P.A. Flach, A.C. Kakas: *Abduction and Induction: Essays on their Relation and Integration*, Applied Logic (Kluwer, Dordrecht 2000)
- 13.21 F. Soler-Toscano, D. Fernández-Duque, Á. Nepomuceno-Fernández: A modal framework for modeling abductive reasoning, *Log. J. IGPL* **20**(2), 438–444 (2012)
- 13.22 M.E. Quilici-Gonzalez, W.F.G. Haselager: Creativity: Surprise and abductive reasoning, *Semiotica* **153**(1–4), 325–342 (2005)
- 13.23 J. van Benthem, F.R. Velázquez-Quesada: The dynamics of awareness, *Synthese (Knowl., Rationality and Action)* **177**, 5–27 (2010)
- 13.24 B. Hill: Awareness dynamics, *J. Phil. Log.* **39**(2), 113–137 (2010)
- 13.25 F.R. Velázquez-Quesada, F. Soler-Toscano, Á. Nepomuceno-Fernández: An epistemic and dynamic approach to abductive reasoning: Abductive problem and abductive solution, *J. Appl. Log.* **11**(4), 505–522 (2013)
- 13.26 Á. Nepomuceno-Fernández, F. Soler-Toscano, F.R. Velázquez-Quesada: An epistemic and dynamic approach to abductive reasoning: Selecting the best explanation, *Log. J. IGPL* **21**(6), 943–961 (2013)
- 13.27 F. Soler-Toscano, F.R. Velázquez-Quesada: A dynamic-epistemic approach to abductive reasoning. In: *Logic of Knowledge. Theory and Applications*, Dialogues and the Games of Logic. A Philosophical Perspective, Vol. 3, ed. by C. Barés Gómez, S. Magnier, F.J. Salguero (College Publications, London 2012) pp. 47–78
- 13.28 C.S. Peirce: *The Essential Peirce. Selected Philosophical Writings (1893–1913)*, Vol. 2 (Indiana Univ., Bloomington, Indianapolis 1998), ed. by N. Houser
- 13.29 C.S. Peirce: *The Essential Peirce. Selected Philosophical Writings (1867–1893)*, Vol. 1 (Indiana Univ., Bloomington, Indianapolis 1992), ed. by N. Houser, C. Kloesel
- 13.30 E. Lorini, C. Castelfranchi: The cognitive structure of surprise: Looking for basic principles, *Topoi* **26**(1), 133–149 (2007)
- 13.31 G.H. Harman: The inference to the best explanation, *Phil. Rev.* **74**(1), 88–95 (1965)
- 13.32 P. Lipton: *Inference to the Best Explanation* (Routledge, London, New York 2004)
- 13.33 J. Hintikka: What is abduction? The fundamental problem of contemporary epistemology, *Trans. C.S. Peirce Soc.* **34**(3), 503–533 (1998)
- 13.34 A. Baltag, S. Smets: A qualitative theory of dynamic interactive belief revision. In: *Logic and the Foundations of Game and Decision Theory (LOFT)*, Texts in Logic and Games, Vol. 3, ed. by G. Bonanno, W. van der Hoek, M. Wooldridge (Amsterdam Univ. Press, Amsterdam 2008) pp. 13–60
- 13.35 F.R. Velázquez-Quesada: Dynamic epistemic logic for implicit and explicit beliefs, *J. Log. Lang. Info.* **23**(2), 107–140 (2014)
- 13.36 C. Boutilier: Unifying default reasoning and belief revision in a modal framework, *Artif. Intell.* **68**(1), 33–85 (1994)



- 13.37 R. Stalnaker: On logics of knowledge and belief, *Phil. Stud.* **128**(1), 169–199 (2006)
- 13.38 J.A. Plaza: Logics of public communications, *Proc. 4th Intl. Symp. Methodol. Intell. Sys.*, ed. by M.L. Emrich, M.S. Pfeifer, M. Hadzikadic, Z.W. Ras (North-Holland, Amsterdam 1989) pp. 201–216
- 13.39 J. Gerbrandy, W. Groeneveld: Reasoning about information change, *J. Log. Lang. Info.* **6**(2), 147–196 (1997)
- 13.40 J. van Benthem: Dynamic logic for belief revision, *J. Appl. Non-Class. Log.* **17**(2), 129–155 (2007)
- 13.41 H. van Ditmarsch: Prolegomena to dynamic logic for belief revision, *Synthese* **147**(2), 229–275 (2005)
- 13.42 A. Baltag, S. Smets: Learning by questions and answers from belief-revision cycles to doxastic fixed points. In: *Logic, Language, Information and Computation*, ed. by H. Ono, M. Kanazawa, R. de Queiroz (Springer, Berlin, Heidelberg 2009) pp. 124–139
- 13.43 W.A. Carnielli: Surviving abduction, *Log. J. IGPL* **14**(2), 237–256 (2006)
- 13.44 J. Pearl: *Probabilistic Reasoning in Intelligent Systems – Networks of Plausible Inference* (Morgan Kaufmann, San Francisco 1989)
- 13.45 D. Poole: Probabilistic horn abduction and bayesian networks, *Artif. Intell.* **64**(1), 81–129 (1993)
- 13.46 D. Dubois, A. Gilio, G. Kern-Isberner: Probabilistic abduction without priors, *Intl. J. Approx. Reason.* **47**(3), 333–351 (2008)
- 13.47 M. Denecker, D. De Schreye: SLDNFA: An abductive procedure for normal abductive programs, *Proc. Intl. Joint Conf. Symp. Log. Program.*, ed. by K.R. Apt (MIT Press, Washington 1992) pp. 686–700
- 13.48 A.C. Kakas, P. Mancarella: Generalized stable models: A semantics for abduction, *Proc. 9th Eur. Conf. Artif. Intell. ECAI '90*, ed. by L.C. Aiello (Pitman, Stockholm 1990) pp. 385–391
- 13.49 F. Lin, J.-H. You: Abduction in logic programming: A new definition and an abductive procedure based on rewriting, *Proc. 17th Int. Joint Conf. Artif. Intell., IJCAI*, ed. by B. Nebel (Morgan Kaufmann, Seattle 2001) pp. 655–666
- 13.50 M.C. Mayer, F. Pirri: Abduction is not deduction-in-reverse, *Log. J. IGPL* **4**(1), 95–108 (1996)
- 13.51 P. Gärdenfors, D. Makinson: Nonmonotonic inference based on expectations, *Artif. Intell.* **65**(2), 197–245 (1994)
- 13.52 R. Pino-Pérez, C. Uzcátegui: Jumping to explanations versus jumping to conclusions, *Artif. Intell.* **111**(1/2), 131–169 (1999)
- 13.53 R. Pino-Pérez, C. Uzcátegui: Preferences and explanations, *Artif. Intell.* **149**(1), 1–30 (2003)
- 13.54 D. Makinson: Bridges between classical and non-monotonic logic, *Log. J. IGPL* **11**(1), 69–96 (2003)
- 13.55 J. van Benthem, S. van Otterloo, O. Roy: Preference logic, conditionals and solution concepts in games. In: *Modality Matters: Twenty-Five Essays in Honour of Krister Segerberg*, Uppsala Philosophical Studies, ed. by H. Lagerlund, S. Lindström, R. Sliwinski (Univ. Uppsala, Uppsala 2006) pp. 61–76
- 13.56 P. Girard: *Modal Logic for Belief and Preference Change*, Ph.D. Thesis (Stanford Univ., Stanford 2008)
- 13.57 J. van Benthem, P. Girard, O. Roy: Everything else being equal: A modal logic for ceteris paribus preferences, *J. Phil. Log.* **38**(1), 83–125 (2009)
- 13.58 F. Liu: *Reasoning about Preference Dynamics*, Synthese Library, Vol. 354 (Springer, Heidelberg 2011)
- 13.59 F.R. Velázquez-Quesada: Reasoning processes as epistemic dynamics, *Axiomathes* **25**(1), 41–60 (2015)
- 13.60 F. Soler-Toscano, F.R. Velázquez-Quesada: Generation and selection of abductive explanations for non-omniscient agents, *J. Log. Lang. Info.* **23**(2), 141–168 (2014)
- 13.61 F. Soler-Toscano, F.R. Velázquez-Quesada: Abduction for (non-omniscient) agents, *Workshop Proc. MALLOW 2010*, Vol. 627, ed. by O. Boissier, A. El Fallah Seghrouchni, S. Hassas, N. Maudet (CEUR, Lyon 2010), [www.ceur-ws.org/vol-627/lrba\\_4.pdf](http://www.ceur-ws.org/vol-627/lrba_4.pdf)
- 13.62 J. van Benthem: Abduction at the interface of logic and philosophy of science, *Theoria* **22**(3), 271–273 (2009)
- 13.63 F.R. Velázquez-Quesada: Explicit and implicit knowledge in neighbourhood models. In: *Logic, Rationality, and Interaction – Proc. 4th Int. Workshop LORI 2013, Hangzhou*, Lecture Notes in Computer Science, Vol. 8196, ed. by D. Grossi, O. Roy, H. Huang (Springer, Berlin, Heidelberg 2013) pp. 239–252