



HAL
open science

The Surprise Deception Paradox

Benjamin Icard

► To cite this version:

| Benjamin Icard. The Surprise Deception Paradox. 2024. hal-04647404v4

HAL Id: hal-04647404

<https://hal.science/hal-04647404v4>

Preprint submitted on 13 Aug 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

The Surprise Deception Paradox

Benjamin Icard

LIP6, Sorbonne Université, CNRS, France

Abstract

This article tackles an epistemic puzzle formulated by R. Smullyan that we call the ‘*Surprise Deception Paradox*’. On the morning of April 1st 1925, his brother announced that he would deceive him during the day, but apparently nothing happened. Since R. Smullyan waited all day to be deceived by some action, he was actually deceived, but by the lack of an action, that is to say by omission. Afterwards, Smullyan felt immediately puzzled: because he expected to be deceived, he was not deceived; but since he was not deceived the way he expected, he was actually deceived. We use dynamic belief revision logic to look more clearly into this puzzle. We argue that Smullyan’s reasoning is not a self-referential paradox but shares common features with the more famous Surprise Examination Paradox. In Smullyan’s riddle, we show that a misleading default mechanism makes R. Smullyan surprised by the deception he has been preyed to. We also use this solution to discuss whether such defaults, compared to other forms of truth-telling deception, may qualify as lies or not.

Keywords — Deception, Omission, Default Inference, Dynamic Belief Revision Logic, Doxastic Paradox, Surprise, Veridical Deception, Lying

1 Introduction

In *What is the Name of this Book?*, Raymond Smullyan tells a personal anecdote concerning his first introduction to logic [49]. While he was suffering from flu on April 1st 1925, his older brother — EMILE —, said to him: “*Well, RAYMOND, today is April Fool’s Day, and I will deceive you as you have never been deceived before!*”. After this announcement, RAYMOND waited all day long to be deceived but (apparently) nothing happened. Late at night, though he was no longer expecting to be deceived (or so he thought), he still felt highly concerned by his brother’s announcement. Noticing his perplexity, EMILE asks him in a Socratic way:

EMILE: “*So, you expected me to deceive you, didn’t you?*”

RAYMOND: *Yes.*

EMILE: *But I didn’t, did I?*

RAYMOND: *No.*

EMILE: *But you expected me to, didn’t you?*

RAYMOND: *Yes.*

EMILE: *So I deceived you, didn’t I?*

RAYMOND: *Yes.”*

After this explanation, RAYMOND lays in his bed and starts reasoning. He rapidly falls into puzzling thoughts. On the one hand, supposing that he wasn’t deceived, then he did not get what he expected, because he expected to be deceived after his brother’s

announcement, and hence he was actually deceived. But on the other hand, supposing that he was actually deceived, then he exactly did get what he expected, and hence he was not *deceived* stricto sensu.

EMILE’s trick may be decomposed into three distinct stages: the *ex ante*, *ex interim* and *ex post* stages. At the *ex ante stage*, EMILE announces that he will deceive RAYMOND later on the day. At this moment, RAYMOND is not yet deceived but he learns he will be. At the *ex interim stage*, deception actually occurs: RAYMOND believes that EMILE will deceive him by doing some specific action (commission) but in fact, he deceives him by doing none (omission). At the *ex post stage*, RAYMOND is no longer deceived but EMILE’s later explanation makes him puzzled and surprised.

In this paper, we analyze the dynamics of EMILE’s deception and show that this latter paradox is specious. However, by comparing EMILE’s announcement with the teacher’s announcement in the surprise examination paradox [e.g. 41, 47, 48], we argue that EMILE’s announcement should not succeed in principle. His deception succeeds, and RAYMOND feels surprised, only because EMILE’s announcement leads RAYMOND to draw the misleading default interpretation that he will be deceived only by commission. Leaving this mechanism aside, EMILE’s announcement should theoretically fail, as in Gerbrandy’s formalizations of the three-day version of the surprise exam paradox [29, 5]. RAYMOND is surprised at the end of the day because he did not expect at all that he could be deceived both by commission and by omission. This is a source of strong surprise.

In Section 2, we adapt extant dynamic belief revision logic [6, 11] to analyze Smullyan’s story in details. Our language $\mathcal{L}_{(\mathbf{B}, \uparrow)}$ is a propositional doxastic syntax enriched with a belief upgrade operator taken from [6, 7]. This upgrade operator is used to provide abbreviations for two defined operators inspired by the quantification over actions in arbitrary announcement logic [26, 33]. Those defined operators aim to express RAYMOND’s mistaken *default interpretation* of EMILE’s announcement as meaning that he will be deceived by some action (*deception by commission*), when in fact he will be deceived by none (*deception by omission*) (subsection 2.1). This notion of default interpretation motivates the use of doxastic plausibility models to provide $\mathcal{L}_{(\mathbf{B}, \uparrow)}$ with semantic clauses (subsection 2.2). That being done, we rely on this apparatus to define various notions of deception involved in Smullyan’s story, such as deception as a state versus as an attitude, deception on a content versus simpliciter, etc. (subsection 2.3). Finally, $\mathcal{L}_{(\mathbf{B}, \uparrow)}$ helps to set out several facts of Smullyan’s tale before analyzing them in more details (subsection 2.4).

In Section 3, we use $\mathcal{L}_{(\mathbf{B}, \uparrow)}$ to describe the stages of the April Fool’s deception. We distinguish EMILE’s *ex ante announcement* in which EMILE informs RAYMOND that he will deceive him during the day (subsection 3.1), from the *ex interim deception* coming later on that day (subsection 3.2). In this case, we focus on the two distinct states of deception RAYMOND goes through successively (subsection 3.2.1) and then express EMILE’s two distinct attitudes leading to those states (subsection 3.2.2). Finally, we model EMILE’s *ex post explanation* that describe his whole deceptive plot (subsection 3.3).

Section 4 delves deeper into the *ex post* stage of deception by analyzing RAYMOND’s puzzling thoughts and surprise after EMILE’s explanation. We first point out that RAYMOND’s reasoning is self-referential but not paradoxical as it stands (subsection 4.1). Then, we compare EMILE’s announcement with the teacher’s announcement in the Surprise Examination Paradox to argue that EMILE’s announcement of a surprising deception should fail in principle (subsection 4.2). Basically, we rely on Gerbrandy’s formal-

ization of surprise [29] to provide a radical upgrade reading of EMILE’s announcement in line with the three-day version of the surprise exam paradox (subsection 4.2.1). Using plausibility models as in [5], we argue that EMILE’s announcement does succeed only as a result of RAYMOND being deceived by misleading default interpretation (subsection 4.2.2). Finally, we argue that RAYMOND is in fact surprised in two different ways by EMILE on April 1st 1925 (subsection 4.2.3).

We conclude in Section 5 by showing that Smullyan’s riddle is informative on the way one can deceive by saying the truth, as in so-called “*veridical deception*”. Classically, lying is considered the standard form of veridical deception since, traditionnally, a speaker can lie by saying the truth, provided that the speaker is insincere when making her utterance. We discuss EMILE’s announcement in light of this standard definition of lying and by taking into account debates related to other forms of veridical deception, such as false implicatures [e.g. 23, 1, 39].

2 A Language for Analysis

2.1 A Dynamic Doxastic Syntax $\mathcal{L}_{(\mathbf{B}, \uparrow)}$

We present $\mathcal{L}_{(\mathbf{B}, \uparrow)}$ which relies on extant dynamic belief revision theory [e.g. 11, 6, 7]. The syntax of $\mathcal{L}_{(\mathbf{B}, \uparrow)}$ contains a static modal operator for (conditional) belief, as well as a dynamic modal operator for upgrading beliefs. The set of all $\mathcal{L}_{(\mathbf{B}, \uparrow)}$ -formulas is given by the following Backus-Naur Form:

$$\langle \mathbf{Formulas} \rangle \quad \varphi, \psi \quad ::= \quad \top \mid p \mid \neg\varphi \mid (\varphi \wedge \psi) \mid \mathbf{B}^\psi\varphi \mid [\uparrow\varphi]\psi,$$

Atomic formulas are abbreviated by p . Complex formulas, such as φ or ψ , can be either tautologies (using the common abbreviation \top), atomic formulas, negation of a formula, conjunction of formulas, static formulas $\mathbf{B}^\psi\varphi$ (for beliefs in φ conditional on ψ), or dynamic formulas of the form $[\uparrow\varphi]\psi$ (for radical upgrade with a formula φ leading to ψ). Additional connectives ($\perp, \vee, \rightarrow, \leftrightarrow$) are defined as usual by using the primitive symbols of language $\mathcal{L}_{(\mathbf{B}, \uparrow)}$. As long as this does not introduce ambiguity, we now omit parentheses when writing formulas for easier reading.

In the specific context of Smullyan’s story, the informal reading of operator $\mathbf{B}^\psi\varphi$ is more precisely “RAYMOND *believes that φ conditional on formula ψ* ”,¹ while the reading of the dynamic operator $[\uparrow\varphi]\psi$ is: “*after RAYMOND commits a radical upgrade with φ , ψ is the case*”. Operator \mathbf{B} stands for “RAYMOND *plainly believes that φ* ” and can be defined from $\mathbf{B}^\top\varphi$ corresponding to “RAYMOND *conditionally believes that φ on \top* ”: $\mathbf{B}\varphi := \mathbf{B}^\top\varphi$. The connection between Smullyan’s story and the complex formulas of $\mathcal{L}_{(\mathbf{B}, \uparrow)}$ is made more precise in sections 3 and 4.

2.2 Doxastic Plausibility Models for $\mathcal{L}_{(\mathbf{B}, \uparrow)}$

As indicated previously, we choose to give a semantic interpretation to $\mathcal{L}_{(\mathbf{B}, \uparrow)}$ by using doxastic plausibility models defined by [e.g. 11, 6, 7, 12], instead of more classical Kripke models used in public announcement logics for instance [30, 4, 43]. This choice is motivated by the fact that, after hearing EMILE’s announcement of deception on April 1st 1925, RAYMOND performs the *default interpretation* that he will be deceived by some action, as in so-called *deception by commission*, instead of no action, as *deception by*

¹Following [9], we can think of conditional beliefs $\mathbf{B}^\psi\varphi$ as ‘contingency’ plans for belief change: in case RAYMOND will find out that ψ was the case, he will believe that φ was the case.

omission. From a more qualitative perspective, this default interpretation corresponds to the *most plausible way* people would interpret EMILE’s early announcement: “*Well, RAYMOND, today is April Fool’s Day, and I will deceive you as you have never been deceived before!*”. In order to capture this notion of default interpretation based on high plausibility, we use doxastic plausibility models to interpret language $\mathcal{L}_{(\mathbf{B}, \uparrow)}$ semantically.

A plausibility model for $\mathcal{L}_{(\mathbf{B}, \uparrow)}$ is a relational structure $\mathbb{S} = \langle \mathcal{S}, \leq, \| \cdot \| \rangle$ with:

- ❖ \mathcal{S} which is a finite non-empty set of “possible states” s (or “worlds”).
- ❖ \leq which is a “plausibility relation” for RAYMOND such that $\leq \subseteq \mathcal{S} \times \mathcal{S}$. This relation \leq is a total preorder, that is to say a transitive, strongly connected and thus reflexive relation over the set \mathcal{S} .
- ❖ $\| \cdot \|: At \rightarrow \wp(\mathcal{S})$ which is a standard “valuation map” where At is the set of all atomic formulas p , and $\wp(\mathcal{S})$ is the set of subsets of \mathcal{S} .

The conventional reading of the plausibility order is the following one: when $s \leq t$ (for all $s, t \in \mathcal{S}$), it means that RAYMOND considers the state t to be “at least as plausible as” the state s . The truth of a formula φ at the actual state s in the doxastic plausibility model \mathbb{S} , denoted $\mathbb{S}, s \models \varphi$, is defined inductively by extending the valuation map to all the complex formulas of $\mathcal{L}_{(\mathbf{B}, \uparrow)}$:

- $\mathbb{S}, s \models \top$ Always.
- $\mathbb{S}, s \models p$ iff $s \in \| p \|$.
- $\mathbb{S}, s \models \neg \varphi$ iff $\mathbb{S}, s \not\models \varphi$.
- $\mathbb{S}, s \models (\varphi \wedge \psi)$ iff $\mathbb{S}, s \models \varphi$ and $\mathbb{S}, s \models \psi$.
- $\mathbb{S}, s \models \mathbf{B}^\psi \varphi$ iff for all $t \in \text{Max}_{\leq}^{\|\psi\|} : \mathbb{S}, t \models \varphi$.

Here $\text{Max}_{\leq}^{\|\psi\|}$ is the set of states that satisfy ψ and that, among all ψ -states, are maximal for the ordering \leq : $\text{Max}_{\leq}^{\|\psi\|} = \{u \in \|\psi\| \mid \text{for all } v \in \|\psi\| \ v \leq u\}$, where $\|\psi\|$ is the extension of $\| \cdot \|$ to arbitrary formulas $\psi \in \mathcal{L}_{(\mathbf{B}, \uparrow)}$, such that, given \mathbb{S} : $\|\psi\| = \{t \in \mathcal{S} \mid \mathbb{S}, t \models \psi\}$. In other words, in the semantic clause corresponding to $\mathbf{B}^\psi \varphi$, the ψ -states on which RAYMOND conditions his beliefs in φ are the ψ -states of the ordering that RAYMOND considers to be the most plausible ones *and* that also satisfy φ . As we said, plain beliefs are recovered from conditional beliefs by conditioning on \top , and the clause can be reduced to:

- $\mathbb{S}, s \models \mathbf{B}\varphi$ iff for all $t \in \text{Max}_{\leq}^{\|\top\|} : \mathbb{S}, t \models \varphi$.

Following [11] and [7], we define a belief radical upgrade with a formula φ , written $[\uparrow \varphi]$, as a mapping of the following kind:

$$[\uparrow \varphi]: \mathbb{S} \mapsto \mathbb{S}^{[\uparrow \varphi]}$$

Here \mathbb{S} is the initial plausibility model and $\mathbb{S}^{[\uparrow \varphi]}$ is the transformed model obtained once the intended operation $[\uparrow \varphi]$ has been performed on \mathbb{S} . The semantic clause for belief radical upgrade is defined by:

- $\mathbb{S}, s \models [\uparrow \varphi]\psi$ iff $\mathbb{S}^{[\uparrow \varphi]}, s \models \psi$.

The doxastic plausibility model $\mathbb{S}^{[\uparrow \varphi]}$ is defined from \mathbb{S} in the following way:

$$\mathbb{S}^{[\uparrow \varphi]} = \langle \mathcal{S}^{[\uparrow \varphi]}, \leq^{[\uparrow \varphi]}, \| \cdot \|^{[\uparrow \varphi]} \rangle$$

where:

$$\begin{aligned}
\mathcal{S}^{[\uparrow\varphi]} &:= \mathcal{S} \\
\leq^{[\uparrow\varphi]} &:= (\leq \cap (\mathcal{S} \times \|\varphi\|)) \cup (\leq \cap (\|\neg\varphi\| \times \mathcal{S})) \cup (\|\neg\varphi\| \times \|\varphi\|) \\
\|\cdot\|^{[\uparrow\varphi]} &:= \|\cdot\|
\end{aligned}$$

Here $\mathcal{S}^{[\uparrow\varphi]}$ is identical to \mathcal{S} from the initial model \mathcal{S} , while on the atomic fragment of $\mathcal{L}_{(\mathcal{B}, \uparrow)}$, $\|\cdot\|^{[\uparrow\varphi]}$ is similar to $\|\cdot\|$. The special feature of $\mathcal{S}^{[\uparrow\varphi]}$ is the plausibility order $\leq^{[\uparrow\varphi]}$. The reordering of states defined by $\leq^{[\uparrow\varphi]}$ ensures that the states where φ is true are promoted in plausibility. More precisely, the first part ($\leq \cap (\mathcal{S} \times \|\varphi\|)$) of $\leq^{[\uparrow\varphi]}$ states that the relative ordering of worlds where φ is true is the same as in the original order \leq . The second part ($\leq \cap (\|\neg\varphi\| \times \mathcal{S})$) states that the relative ordering of worlds where φ is false is the same as in the original order \leq . Finally, the third part ($\|\neg\varphi\| \times \|\varphi\|$) states that the worlds where φ is true become equally or more plausible than the worlds where φ is false. So, in the upgraded model $\mathcal{S}^{[\uparrow\varphi]}$, although all φ -states are made *equally or more plausible than* all $\neg\varphi$ -states, everything else is left the same *within* the φ and $\neg\varphi$ zones themselves.

For clarity's sake, let us put more emphasis on the dynamic attitude corresponding to operator $[\uparrow\varphi]\psi$. The upgrade part of this operator, that is to say $[\uparrow\varphi]$, indicates that the agent does not know the source of φ to be *absolutely reliable*, thus *completely honest* and *truthful*. But the agent believes the source to be *highly reliable*, thus *strongly honest* and *truthful*.² The whole complex formula $[\uparrow\varphi]\psi$ according to which ψ is the case after the belief upgrade $[\uparrow\varphi]$, can be seen as a *default rule of interpretation* based on the intuition that ψ is *more plausible* than $\neg\psi$ in all the cases in which φ obtains. This notion of *higher plausibility* which is the core intuition governing non-monotonic reasoning has been captured in [51] by a formal rule written $(\varphi \Rightarrow \psi)$ meaning that formulas “ φ are normally ψ ” or, to put it another way, that worlds satisfying the conjunction $(\varphi \wedge \psi)$ are *more plausible* than worlds satisfying the conjunction $(\varphi \wedge \neg\psi)$. For this reason, we assume that the default rule $(\varphi \Rightarrow \psi)$ can be reasonably approximated by the formula $[\uparrow\varphi]\mathcal{B}\psi$ in language $\mathcal{L}_{(\mathcal{B}, \uparrow)}$.

Various other dynamic revision operators exist in the literature [e.g. 11, 7, 12], based on previous work on belief revision policies [e.g. 45, 27]. Within all this diversity, two well-known representatives are hard updates such as e.g. public announcements, usually written “!”, in which sources are considered completely truthful and receive absolute trust [30, 4, 43], or softer upgrades such as e.g. “conservative upgrades”, usually written “ \uparrow ”, in which sources are taken with much weaker confidence as only barely truthful [7, 12]. Those policies can also be read as types of default interpretations but regarding the trust RAYMOND should put into EMILE on April 1st 1925, radical upgrades are the most appropriate match: high confidence, though not perfect confidence. From a theoretical perspective, RAYMOND cannot have perfect confidence in EMILE: the day is April Fool's Day and EMILE explicitly states that he will deceive him during that day. On a more practical note, RAYMOND's behaviour all along the day shows that he has, and keeps having, strong confidence in EMILE, — even at the end of the day when nothing has happened and even after his brother's late explanation. So, weaker confidence operators,

²According to van Benthem's terminology, the agent does not have a *hard* doxastic attitude towards the source of information, but a *softer* one. The agent is not sure that the source always tells the truth but she is *highly confident* that the source does.

such as conservative upgrades, would not match adequately with RAYMOND’s effective behaviour.

As proved in [7] and [11], $\mathcal{L}_{(\mathbf{B}, \uparrow)}$ is axiomatized completely by a set of axioms and inferences rules for propositional logic, belief operator \mathbf{B} and the dynamic upgrade operator \uparrow . Here we simply put emphasis on two axioms of central importance for analyzing RAYMOND’s paradoxical thoughts, in addition to other more classical principles (see subsection 4.1 in that respect). For $\varphi, \psi \in \mathcal{L}_{(\mathbf{B}, \uparrow)}$, the two reduction axioms we use are:

$$\begin{array}{ll} (i) & \vdash \quad [\uparrow \varphi] \neg \psi \quad \leftrightarrow \quad \neg [\uparrow \varphi] \psi \\ (ii) & \vdash \quad [\uparrow \varphi] (\psi \wedge \phi) \quad \leftrightarrow \quad ([\uparrow \varphi] \psi \wedge [\uparrow \varphi] \phi) \end{array}$$

2.3 Expressing Deception with $\mathcal{L}_{(\mathbf{B}, \uparrow)}$

We first define the distinction between *action of commission* leading to a formula ψ versus *omission of an action* leading to ψ to draw further distinctions regarding RAYMOND’s deception on April 1st 1925.³ Broadly speaking, we distinguish the fact for RAYMOND to be in states of deception on April 1st 1925 from EMILE’s deceptive attitudes leading to RAYMOND to fall in such states.

2.3.1 Commission versus Omission of an Action

Definition 1 (commission such that ψ) We let $[\uparrow^+] \psi$ stands for the action of commission leading to formula ψ on April 1st 1925. Given some $\psi \in \mathcal{L}_{(\mathbf{B}, \uparrow)}$, this corresponds to the following abbreviation:

$$[\uparrow^+] \psi \quad := \quad [\uparrow \varphi] \psi \text{ for some } \varphi \in \mathcal{L}_{(\mathbf{B}, \uparrow)}$$

The abbreviation $[\uparrow^+] \psi$ provided in Definition 1 states that formula ψ obtains as the result of a radical upgrade with some formula φ . In the context of Smullyan’s story, this abbreviation $[\uparrow^+] \psi$ can be read as “*after RAYMOND commits some radical upgrade, ψ is the case*”.

Definition 2 (omission such that ψ) We let $[\uparrow^-] \psi$ stands for the omission of an action leading to formula ψ on April 1st 1925. Given some $\psi \in \mathcal{L}_{(\mathbf{B}, \uparrow)}$, this corresponds to the abbreviation:

$$[\uparrow^-] \psi \quad := \quad [\uparrow \varphi] \psi \text{ for no } \varphi \in \mathcal{L}_{(\mathbf{B}, \uparrow)}$$

By contrast with Definition 1, abbreviation of Definition 2 states that formula ψ now obtains as the result of the lack of any radical upgrade with some φ . In the context of Smullyan’s story, the abbreviation $[\uparrow^-] \psi$ can be read as “*after RAYMOND omits to perform any radical upgrade, ψ is the case*”.

Observation 1 (commission vs omission such that ψ) Notice that given a formula $\psi \in \mathcal{L}_{(\mathbf{B}, \uparrow)}$, a doxastic plausibility model $\mathbb{S} = \langle \mathcal{S}, \leq, \|\cdot\| \rangle$ for $\mathcal{L}_{(\mathbf{B}, \uparrow)}$ with a state $\mathbf{s} \in \mathcal{S}$, we have $\mathbb{S}, \mathbf{s} \models [\uparrow^-] \psi$ iff $\mathbb{S}, \mathbf{s} \not\models [\uparrow^+] \psi$. Indeed, suppose that we have $\mathbb{S}, \mathbf{s} \models [\uparrow^-] \psi$. So, by Definition 2, we have $\mathbb{S}, \mathbf{s} \models [\uparrow^-] \psi$ iff $\mathbb{S}, \mathbf{s} \models [\uparrow \varphi] \psi$ for no $\varphi \in \mathcal{L}_{(\mathbf{B}, \uparrow)}$ iff $\mathbb{S}, \mathbf{s} \not\models [\uparrow \varphi] \psi$ for all $\varphi \in \mathcal{L}_{(\mathbf{B}, \uparrow)}$, so that we have in particular $\mathbb{S}, \mathbf{s} \not\models [\uparrow \varphi] \psi$ for some $\varphi \in \mathcal{L}_{(\mathbf{B}, \uparrow)}$, i.e. $\mathbb{S}, \mathbf{s} \not\models [\uparrow^+] \psi$ by Definition 1.

³Notice that the nature of omission regarding causality [e.g. 46, 18] raises deontological issues linked to the moral responsibility of agents [19, 54] in computational ethics and related logics [15, 16, 21]. These questions of crucial importance go beyond the scope of this article which is more concerned with modelling attitudes of deception based on commission versus omission and epistemic paradoxes involved in Smullyan’s recollection.

Consistent with conceptual intuitions, Observation 1 shows that the omission of an action resulting in ψ amounts to the *absence*, or *lack*, of an action of commission resulting in ψ .

2.3.2 Deception as a State

We start out by defining what it means for RAYMOND to be in a doxastic state of deception, first concerning a specific formula ψ and then simpliciter.

Definition 3 (*deceived on ψ*) We write deceived_ψ the fact for RAYMOND to be deceived on some formula ψ on April 1st 1925. The corresponding abbreviation is:

$$\text{deceived}_\psi := (\psi \wedge \mathbf{B}\neg\psi) \text{ for some } \psi \in \mathcal{L}_{(\mathbf{B}, \uparrow)}$$

Outside of Smullyan's story, Definition 3 would be considered inadequate for the reason that it does not express the presence of the deceiver and sets aside his or her intentions to deceive the addressee. But firstly, the truly central dimensions of EMILE's deceptive plan are the dimensions of commission and omission, not EMILE's intention to deceive RAYMOND, which is naturally assumed. Secondly, EMILE's intention to deceive does not play any particular role in the doxastic paradoxes involved in Smullyan's story (see subsections 4.1 and 4.2). For those reasons, we assume that EMILE has the intention to deceive RAYMOND, but we do not formalize it explicitly in $\mathcal{L}_{(\mathbf{B}, \uparrow)}$ or by introducing further definitions.

Definition 4 (*deceived simpliciter*) We write deceived the fact for RAYMOND to be deceived simpliciter on April 1st 1925. The corresponding abbreviation is:

$$\text{deceived} := \text{deceived}_\psi \text{ for some } \psi \in \mathcal{L}_{(\mathbf{B}, \uparrow)}$$

Definition 4 simply states that being deceived on some formula $\psi \in \mathcal{L}_{(\mathbf{B}, \uparrow)}$ is sufficient for RAYMOND to be deceived *stricto sensu* by EMILE on April 1st 1925.

2.3.3 Deception as an Attitude

After having defined deception as a state, we define deception as an attitude here. For capturing the subtleties of EMILE's deceptive trick, we need to be able to capture the distinction between deception by an action, i.e. *deception by commission*, versus deception by the lack of an action, i.e. *deception by omission*.

Definition 5 (*deception by commission on ψ*) We abbreviate by d_ψ^+ the action of commission leading RAYMOND to be deceived on formula ψ . That is:

$$d_\psi^+ := [\uparrow^+] \text{deceived}_\psi \text{ for some } \psi \in \mathcal{L}_{(\mathbf{B}, \uparrow)}$$

Based on Definition 3, Definition 5 states that RAYMOND is deceived by commission on some formula $\psi \in \mathcal{L}_{(\mathbf{B}, \uparrow)}$ in case there is some action such that after this action, RAYMOND is deceived on ψ , since ψ is true but RAYMOND believes ψ to be false.

Definition 6 (*deception by omission on ψ*) We use d_ψ^- to abbreviate the omission of an action leading RAYMOND to be deceived on formula ψ . That is:

$$d_\psi^- := [\uparrow^-] \text{deceived}_\psi \text{ for some } \psi \in \mathcal{L}_{(\mathbf{B}, \uparrow)}$$

Again, based on Definition 3, Definition 6 states that RAYMOND is deceived by omission on some formula $\psi \in \mathcal{L}_{(\mathbf{B}, \uparrow)}$ in case there is no action such that, as the result of this lack of action, RAYMOND is deceived on ψ , since ψ is true but RAYMOND believes ψ to be false.

Observation 2 (*deception by commission vs omission such that ψ*) Following Observation 2, given a formula $\psi \in \mathcal{L}_{(\mathbf{B}, \uparrow)}$, a doxastic plausibility model $\mathbb{S} = \langle \mathcal{S}, \leq, \parallel \cdot \parallel \rangle$

for $\mathcal{L}_{(\mathcal{B}, \uparrow)}$ with a state $s \in \mathcal{S}$, we have $\mathbb{S}, s \models d_{\psi}^{-}$ iff $\mathbb{S}, s \not\models d_{\psi}^{+}$. Again, Observation 2 is consistent with conceptual intuitions since the fact that EMILE will deceive RAYMOND by omission on a formula ψ on April 1st 1925 implies that he will not deceive him by commission on ψ , and the other way around.

Definition 7 (deception on ψ) We abbreviate by d_{ψ} the fact that EMILE deceive RAYMOND on formula $\psi \in \mathcal{L}_{(\mathcal{B}, \uparrow)}$ on April 1st 1925, either by commission or by omission.

$$d_{\psi} := (d_{\psi}^{+} \vee d_{\psi}^{-}) \text{ for some } \psi \in \mathcal{L}_{(\mathcal{B}, \uparrow)}$$

As emphasized in Observation 2, d_{ψ}^{+} and d_{ψ}^{-} cannot be true together for a given formula ψ since d_{ψ}^{+} is semantically equivalent to $(d_{\psi}^{+} \wedge \neg d_{\psi}^{-})$, while d_{ψ}^{-} is semantically equivalent to $(\neg d_{\psi}^{+} \wedge d_{\psi}^{-})$. Naturally, this semantic incompatibility does not impose any syntactic restriction in Definition 7.

Definition 8 (deception simpliciter, versions of) Let us use abbreviations d^{+} , d^{-} and d for attitudes of deception by commission simpliciter, deception by omission simpliciter and deception simpliciter. That is:

- Deception by commission simpliciter: $d^{+} := d_{\psi}^{+}$ for some $\psi \in \mathcal{L}_{(\mathcal{B}, \uparrow)}$
- Deception by omission simpliciter: $d^{-} := d_{\psi}^{-}$ for some $\psi \in \mathcal{L}_{(\mathcal{B}, \uparrow)}$
- Deception simpliciter: $d := d_{\psi}$ for some $\psi \in \mathcal{L}_{(\mathcal{B}, \uparrow)}$

In definition 8, formula d means that, in case EMILE actually deceive RAYMOND, he deceives him either by commission (d^{+}) or by omission (d^{-}). Conceptually, formula d^{+} states that EMILE will deceive RAYMOND by *doing some specific action* on April 1st 1925. That is: EMILE will make *some move* such that RAYMOND will be *deceived on some fact* afterwards. Such a strategy is known as “*deception by commission*” in the literature [17]. By contrast, formula d^{-} states that EMILE will deceive RAYMOND by *not making any action* on April 1st 1925, that is: by the absence of an action. Now EMILE won’t perform any action but, in doing so, he will deceive his brother by depriving him of some informational facts. This strategy is also called “*deception by omission*” in the literature [17].

2.4 Contextualizing $\mathcal{L}_{(\mathcal{B}, \uparrow)}$ to Smullyan’s story

Formulas d^{+} and d^{-} should be seen as *types* of deceptive attitudes. Those formulas express the two kinds of deceptive means EMILE can use to deceive his brother after his morning announcement. Let us abbreviate by formula D EMILE’s morning announcement in $\mathcal{L}_{(\mathcal{B}, \uparrow)}$, more precisely the formula capturing its content:

Definition 9 (informal reading of Emile’s announcement) We write D the content of EMILE’s morning announcement on April 1st 1925. That is:

D : “Well, RAYMOND, today is April Fool’s Day, and I will deceive you as you have never been deceived before!”

In Definition 9, abbreviation D is only a first, and still informal, way of writing EMILE’s morning announcement.⁴ Intuitively, three distinct combinations of formulas d^{+} and d^{-} align with D being true: $(d^{+} \wedge \neg d^{-})$, $(\neg d^{+} \wedge d^{-})$ and $(d^{+} \wedge d^{-})$. The first combination $(d^{+} \wedge \neg d^{-})$ means that RAYMOND is deceived *only* by commission on April 1st 1925. The second combination $(\neg d^{+} \wedge d^{-})$ that he is deceived *only* by omission on that date. The last combination $(d^{+} \wedge d^{-})$ means that RAYMOND is deceived *both*

⁴In Definition 12, we provide a formal expression of this announcement, written \mathbf{D} , which clearly insists on the surprise dimension of EMILE’s announcement of deception, based on [29, 5].

by commission *and* by omission on April 1st 1925. Contrariwise, in case RAYMOND is *neither* deceived by commission *nor* by omission on April 1st 1925, then formulas \mathbf{d}^+ and \mathbf{d}^- are false, so we have the combination $(\neg\mathbf{d}^+ \wedge \neg\mathbf{d}^-)$.

Let us now point out facts that are objectively true in Smullyan's story. To begin with, RAYMOND's upgrade with EMILE's announcement D is an action. After this action, RAYMOND is indeed *deceived by commission* \mathbf{d}^+ but, in fact, he is not deceived by commission *only*, so the combination $(\mathbf{d}^+ \wedge \neg\mathbf{d}^-)$ is false: $\neg(\mathbf{d}^+ \wedge \neg\mathbf{d}^-)$. In addition to that, RAYMOND also turns out to be *deceived by omission* at the end of the day, so the combination $(\mathbf{d}^+ \wedge \mathbf{d}^-)$ is true. As a result, the combination stating that RAYMOND is deceived *only* by omission is false: $\neg(\neg\mathbf{d}^+ \wedge \mathbf{d}^-)$.

From a more pragmatic perspective, the announcement D triggers the conclusion that EMILE will deceive RAYMOND only by commission, which is encoded by the conjunctive formula $(\mathbf{d}^+ \wedge \neg\mathbf{d}^-)$. This is the conclusion that RAYMOND believe after EMILE's announcement D . This conclusion relies on the default inference one would make to interpret some announcement of deception in common circumstances. When hearing such an announcement, which can be written by $[\uparrow D]$, one would naturally infer that deception only by commission $(\mathbf{d}^+ \wedge \neg\mathbf{d}^-)$ is *more plausible* than the negation of it $\neg(\mathbf{d}^+ \wedge \neg\mathbf{d}^-)$. In other words, one would infer that the conjunction $D \wedge (\mathbf{d}^+ \wedge \neg\mathbf{d}^-)$ is a *more plausible* option than the conjunction $D \wedge \neg(\mathbf{d}^+ \wedge \neg\mathbf{d}^-)$. Using the classical writing of default rules given in [51], the default rule at stake in Smullyan's story can be written $D \Rightarrow (\mathbf{d}^+ \wedge \neg\mathbf{d}^-)$.

The issue is that the formula $(\mathbf{d}^+ \wedge \neg\mathbf{d}^-)$ is false in the context of Smullyan's tale: it is not the case that RAYMOND will be deceived only by commission after his brother's announcement. Actually, EMILE's explanation reveals that RAYMOND is deceived also by omission at the end of the day. Though being perfectly natural as a default rule, the inference $D \Rightarrow (\mathbf{d}^+ \wedge \neg\mathbf{d}^-)$ should not be used by RAYMOND on April 1st since using this rule leads him to draw the false conclusions that he will be deceived only by commission and thus, that he won't be deceived also by omission. We will show later that having those false beliefs is what causes his surprise afterwards. Before doing this, we take the opportunity to describe EMILE's misleading announcement in more details.

3 Smullyan's April Fool's Trick

3.1 Ex Ante Interpretation

Semantically, RAYMOND uses the default rule $D \Rightarrow (\mathbf{d}^+ \wedge \neg\mathbf{d}^-)$ to interpret EMILE's announcement as meaning that he will deceive him only by commission, i.e. by believing that $(\mathbf{d}^+ \wedge \neg\mathbf{d}^-)$: $\mathbf{B}(\mathbf{d}^+ \wedge \neg\mathbf{d}^-)$. Let us now define a reasonable approximation of this default interpretation in language $\mathcal{L}_{(\mathbf{B}, \uparrow)}$.

Definition 10 (default deception) Let us write **default** the abbreviation of the formula $[\uparrow D]\mathbf{B}(\mathbf{d}^+ \wedge \neg\mathbf{d}^-)$ that we previously accepted as a reasonable approximation of the default rule $D \Rightarrow (\mathbf{d}^+ \wedge \neg\mathbf{d}^-)$. That is:

$$\mathbf{default} \quad := \quad [\uparrow D]\mathbf{B}(\mathbf{d}^+ \wedge \neg\mathbf{d}^-)$$

3.2 Ex Interim Deception

3.2.1 Raymond's Two States of Deception

We now describe these two states of deception RAYMOND successively goes through during the day. On April 1st 1925, RAYMOND is actually **deceived two times** by his

brother. To begin with, EMILE's morning announcement that D brings about RAYMOND's *first state of deception*. By applying the default rule **default** to D , RAYMOND computes the wrong conclusion that he will be deceived only by commission on April 1st 1925: $(d^+ \wedge \neg d^-)$. Accordingly, RAYMOND believes a formula, i.e. $(d^+ \wedge \neg d^-)$, which is false. So based on definition 3, RAYMOND is deceived on formula $\neg(d^+ \wedge \neg d^-)$ after his brother's announcement, which can be described by the conjunction (before and after eliminating double negation):⁵

$$\begin{aligned} \text{deceived}_{\neg(d^+ \wedge \neg d^-)} &:= \neg(d^+ \wedge \neg d^-) \wedge \mathbf{B}\neg\neg(d^+ \wedge \neg d^-) \\ &:= \neg(d^+ \wedge \neg d^-) \wedge \mathbf{B}(d^+ \wedge \neg d^-) \end{aligned}$$

Late on that day, a *second state of deception* adds to this first state when EMILE explains his April Fool's plot. Indeed, EMILE reveals that RAYMOND was not deceived only by commission but that he was in fact deceived both by commission and by omission: $(d^+ \wedge d^-)$. Since RAYMOND believes this possibility to be false all day, his *second state of deception* consists in being deceived on this formula $(d^+ \wedge d^-)$:

$$\text{deceived}_{(d^+ \wedge d^-)} := (d^+ \wedge d^-) \wedge \mathbf{B}\neg(d^+ \wedge d^-)$$

To sum up, RAYMOND goes through *two states of deception* in Smullyan's tale. First, he is deceived on the fact that he will be deceived only by commission: $\text{deceived}_{\neg(d^+ \wedge \neg d^-)}$. Then, he is deceived on the fact that he will be deceived both by commission and by omission: $\text{deceived}_{(d^+ \wedge d^-)}$. RAYMOND's first state of deception is made possible by applying the default rule **default** to EMILE's announcement D . RAYMOND's second state of deception results from the fact that RAYMOND does not realize that he will be deceived also by omission after realizing that he has not been deceived (only) by commission. The next subsection is devoted to studying EMILE's attitudes that lead to those two states of deception.

3.2.2 Emile's Two Attitudes of Deception

Raymond's *first state of deception* is that he wrongly believes that he will be deceived only by commission on April 1st 1925: $\neg(d^+ \wedge \neg d^-) \wedge \mathbf{B}\neg\neg(d^+ \wedge \neg d^-)$, or after eliminating double negation: $\neg(d^+ \wedge \neg d^-) \wedge \mathbf{B}(d^+ \wedge \neg d^-)$. As we have said, this state results from the **default rule** he applies to EMILE's announcement. This announcement D is in fact the event that leads RAYMOND to be deceived a first time, since after this announcement D the following formula holds as an assumption:

$$[\uparrow D](\neg(d^+ \wedge \neg d^-) \wedge \mathbf{B}(d^+ \wedge \neg d^-))$$

By applying definition 3 and the definition of the commission operator \uparrow^+ :

$$[\uparrow^+] \text{deceived}_{\neg(d^+ \wedge \neg d^-)}$$

So by applying definition 5, we have:

$$\mathbf{d}_{\neg(d^+ \wedge \neg d^-)}^+$$

Late in the afternoon, comes RAYMOND's second state of deception since he still does not realize that he is actually deceived both by commission and by omission: $(d^+ \wedge d^-)$. Interestingly, this second state of deception does not result from EMILE making any action, but specifically from the lack of an action, or omission. Let us sum up the event that leads to RAYMOND's *second state of deception*. Late on the afternoon, the following formula holds:

⁵From now on, we eliminate double negations immediately when it is possible.

$$[\uparrow^-]((\mathbf{d}^+ \wedge \mathbf{d}^-) \wedge \mathbf{B}\neg(\mathbf{d}^+ \wedge \mathbf{d}^-))$$

By applying definition 3 and the definition of the omission operator \uparrow^- :

$$[\uparrow^-]\mathbf{deceived}_{(\mathbf{d}^+ \wedge \mathbf{d}^-)}$$

So by applying definition 6:

$$\mathbf{d}^-_{(\mathbf{d}^+ \wedge \mathbf{d}^-)}$$

We are now able to summarize EMILE’s deceptive plot on April 1st 1925 into a single formula, as given in Definition 11 hereafter.

Definition 11 (*deception plot*) We define EMILE’s deceptive plot on April 1st 1925, written *deception*, which relies on using both deception by commission and deception by omission as follows:

$$\mathbf{deception} := (\mathbf{d}^-_{\neg(\mathbf{d}^+ \wedge \neg \mathbf{d}^-)} \wedge \mathbf{d}^-_{(\mathbf{d}^+ \wedge \mathbf{d}^-)})$$

This definition is in line with the deceptive plot happening in Smullyan’s story since, according to the left conjunct \mathbf{d}^+ , EMILE deceives RAYMOND by commission, and, according to the right conjunct \mathbf{d}^- , EMILE also deceives RAYMOND by omission. Specifically, EMILE deceives RAYMOND by commission on formula $\neg(\mathbf{d}^+ \wedge \neg \mathbf{d}^-)$ and deceives RAYMOND by omission on formula $(\mathbf{d}^+ \wedge \mathbf{d}^-)$.

In the next section, we use dynamic plausibility models to show that EMILE’s explanation acts as a way of teaching RAYMOND the meaning of this formula *Deception*, that is to say his plan to deceive his brother on April 1st 1925.

3.3 Ex Post Explanation

We know that late on April Fool’s Day, RAYMOND has still not understood the trick. His mother intervenes and EMILE finally unveils the inner workings of his deceptive plot. As a matter of fact, EMILE’s explanation can be seen as a Socratic dialogue to teach him the meaning of sentence *Deception*. To show this, we now use plausibility models to represent the dynamics of this dialogue. The epistemic landscape of all the deceptive possibilities can be given by the set:

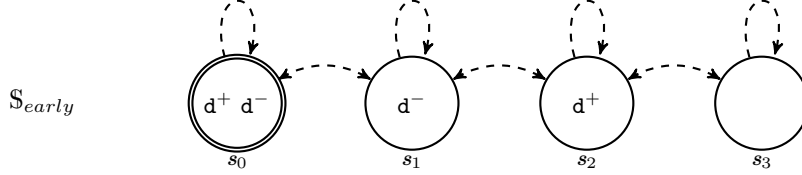
$$\{\mathbf{s}_0^{(\mathbf{d}^+ \wedge \mathbf{d}^-)}, \mathbf{s}_1^{(\neg \mathbf{d}^+ \wedge \mathbf{d}^-)}, \mathbf{s}_2^{(\mathbf{d}^+ \wedge \neg \mathbf{d}^-)}, \mathbf{s}_3^{(\neg \mathbf{d}^+ \wedge \neg \mathbf{d}^-)}\}$$

For simplicity’s sake, our modelling only represents the logical space resulting from complex formulas \mathbf{d}^+ and \mathbf{d}^- , setting aside the modelling of atomic formulas.⁶ The formula holding at each state is right-upperscripted. State \mathbf{s}_3 corresponds to EMILE not deceiving RAYMOND either by commission or by omission, \mathbf{s}_2 to deceiving him only by commission, \mathbf{s}_1 to deceiving only by omission and \mathbf{s}_0 to deceiving him by both options. The actual state is \mathbf{s}_0 since in Smullyan’s story, EMILE deceives RAYMOND by both commission and by omission. Graphically, to represent the fact that the state \mathbf{t} is *at least as plausible as* the state \mathbf{s} for RAYMOND (i.e. $\mathbf{s} \leq \mathbf{t}$), we draw a dashed right arrow “ \dashrightarrow ” from state \mathbf{s} to state \mathbf{t} . When states \mathbf{s} and \mathbf{t} are *equally plausible* for him (that

⁶States are only labelled with complex formulas \mathbf{d}^+ and \mathbf{d}^- since, when modelling EMILE’s explanation, we are only concerned with the evolution of RAYMOND’s beliefs regarding deception (only) by commission and deception (only) by omission, both represented by complex formulas \mathbf{d}^+ and \mathbf{d}^- , respectively. For similar reasons, we only represent \mathbf{d}^+ and \mathbf{d}^- when modelling the effect of EMILE’s announcement of surprise in subsection 4.2.2.

is $s \leq t$ and $t \leq s$), we draw a left-right dashed arrow “ $\leftarrow\text{----}\rightarrow$ ” between s and t . For easier readability, we omit representing the transitivity of plausibility orders in any of the models.

Early in the morning on April Fool’s Day, before any announcement, RAYMOND has no particular belief in the possibility of his brother deceiving him on that day. More precisely, being deceived and not being deceived are equally plausible options, and so are all the states of the epistemic landscape: s_0 , s_1 , s_2 and s_3 . The following model depicts this initial configuration, before any announcement:

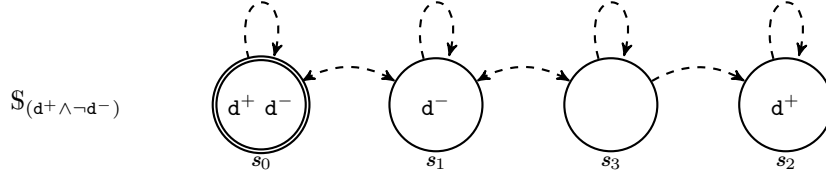


According to RAYMOND’s recollection, EMILE’s explanation starts as follows:

EMILE: “*So, you expected me to deceive you, didn’t you?*”

RAYMOND: *Yes.*”

Here EMILE’s question can be understood as “*So, you expected me to deceive you by some action, didn’t you?*”. Since this question is in the affirmative form, it can be formally expressed by RAYMOND upgrading his beliefs with formula $(d^+ \wedge \neg d^-)$ corresponding to the belief he reached through the default interpretation of EMILE’s announcement. The matching operation $[\uparrow (d^+ \wedge \neg d^-)]: S_{early} \rightarrow S_{(d^+ \wedge \neg d^-)}$ returns the following model:

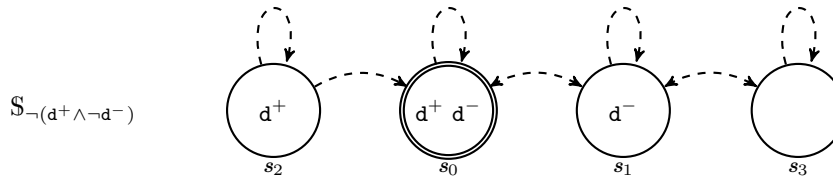


We can see that $S_{(d^+ \wedge \neg d^-)}, s_0 \models B(d^+ \wedge \neg d^-)$. Here RAYMOND acknowledges that he came to believe that $(d^+ \wedge \neg d^-)$ after his brother’s announcement that D . Immediately, though, EMILE informs him that this default interpretation was wrong:

EMILE: “*But I didn’t, did I?*”

RAYMOND: *No.*”

By those lines, EMILE invites his brother to retract his belief in formula $(d^+ \wedge \neg d^-)$. The corresponding dynamic operation is $[\uparrow \neg(d^+ \wedge \neg d^-)]: S_{(d^+ \wedge \neg d^-)} \rightarrow S_{\neg(d^+ \wedge \neg d^-)}$. Graphically, the model becomes:



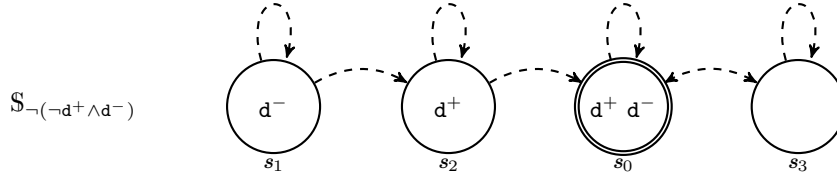
In this model, RAYMOND now believes that formula $(d^+ \wedge \neg d^-)$ is false, that is: $\mathbb{S}_{\neg(d^+ \wedge \neg d^-)}, s_0 \models \mathbf{B}\neg(d^+ \wedge \neg d^-)$. By doing so, he acknowledges that he was wrong to believe that $(d^+ \wedge \neg d^-)$ after his brother's announcement and thus, that he has been deceived on $(d^+ \wedge \neg d^-)$. At this step of EMILE's explanation, RAYMOND is actually uncertain whether he has been deceived at all since s_0 , s_1 and s_3 are the most plausible states of the ordering, yet equally plausible.

But EMILE's demonstration goes even further. He explains that not only has RAYMOND been wrong, and thus deceived, on formula $(d^+ \wedge \neg d^-)$ but that he has also been deceived on formula $(d^+ \wedge d^-)$. The reason is that RAYMOND kept dismissing the true possibility that he would be deceived both by commission d^+ and by omission d^- after realizing, as he just did, that he would not be deceived by commission *only*: $\neg(d^+ \wedge \neg d^-)$. EMILE starts explaining why in the following lines:

EMILE: “*But you expected me to, didn't you?*”

RAYMOND: *Yes.*”

Here EMILE's repeats the same question he asked before (“*So, you expected me to deceive you, didn't you?*”) using the contrastive word “*But*” to generate a different interpretation that we can translate by reformulating EMILE's question: “*But you kept dismissing that you could be deceived also by omission, didn't you?*”. This way, EMILE invites RAYMOND to realize that he kept rejecting this option even though he observed that he was not deceived only by commission in model $\mathbb{S}_{\neg(d^+ \wedge \neg d^-)}$. Rejecting deception by omission as a possibility corresponds to the upgrade $[\uparrow \neg(\neg d^+ \wedge d^-)]: \mathbb{S}_{\neg(d^+ \wedge \neg d^-)} \rightarrow \mathbb{S}_{\neg(\neg d^+ \wedge d^-)}$ which returns the following model:

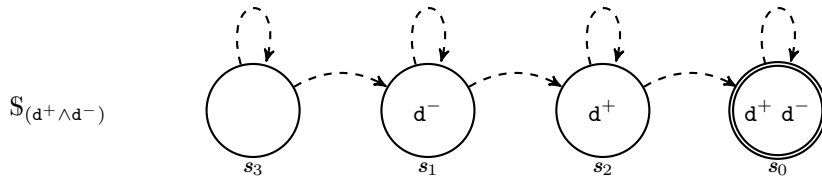


In model $\mathbb{S}_{\neg(\neg d^+ \wedge d^-)}$, RAYMOND is now uncertain whether he is deceived both by commission and by omission, as encoded by $(d^+ \wedge d^-)$ true in state s_0 , or not deceived at all, as encoded by $(\neg d^+ \wedge \neg d^-)$ true in state s_3 . But still, option $(d^+ \wedge d^-)$ is now a possibility. EMILE asks RAYMOND a last question to make him realize that option option $(d^+ \wedge d^-)$ is in fact the true one:

EMILE: “*So I deceived you, didn't I?*”

RAYMOND: “*Yes.*”

Between options $(d^+ \wedge d^-)$ and $(\neg d^+ \wedge \neg d^-)$, the only remaining option for RAYMOND to be deceived *per se* is of course to update his beliefs with formula $(d^+ \wedge d^-)$. The corresponding operation is $[\uparrow (d^+ \wedge d^-)]: \mathbb{S}_{\neg(\neg d^+ \wedge d^-)} \rightarrow \mathbb{S}_{(d^+ \wedge d^-)}$ returning the final model:



We have described the distinct steps by which RAYMOND learns *why* and *how* he was fooled by his brother on April 1st 1925. By going backward from last model $S_{\text{deceived}_{(d^+ \wedge d^-)}}$ to initial model S_{early} , we can reconstitute the dynamics of EMILE’s explanation, as described in Table 1.

Dialogical Explanation	Learning Step	Some Epistemic Effect
E: “ <i>So you expected...?</i> ” R: “ <i>Yes.</i> ”	(a) $[\uparrow (d^+ \wedge \neg d^-)]$	$B(d^+ \wedge \neg d^-)$
E: “ <i>But I didn’t...?</i> ” R: “ <i>No.</i> ”	(b) $[\uparrow \neg(d^+ \wedge \neg d^-)]$	$\text{deceived}_{\neg(d^+ \wedge \neg d^-)}$
E: “ <i>But you expected...?</i> ” R: “ <i>Yes.</i> ”	(c) $[\uparrow \neg(\neg d^+ \wedge d^-)]$	$B\neg(d^+ \wedge d^-)$
E: “ <i>So I deceived you...?</i> ” R: “ <i>Yes.</i> ”	(d) $[\uparrow (d^+ \wedge d^-)]$	$\text{deceived}_{(d^+ \wedge d^-)}$

Table 1: EMILE’s Explanation to RAYMOND.

Formulas (a) to (d) of Table 1 summarize the distinct steps through which RAYMOND learns EMILE’s deceptive strategy. At step (a), RAYMOND learns that after hearing announcement D , he made an upgrade with the default interpretation that he would be deceived only by commission: $[\uparrow (d^+ \wedge \neg d^-)]$. But at step (b), EMILE reveals that this default interpretation was a mistake: $[\uparrow \neg(d^+ \wedge \neg d^-)]$. By concatenating steps (a) and (b), RAYMOND can actually conclude that he has been deceived on formula $(d^+ \wedge \neg d^-)$, which is captured by $\text{deceived}_{\neg(d^+ \wedge \neg d^-)}$. But at step (c), EMILE informs RAYMOND that though realizing that he had not been deceived (only) by commission, he kept dismissing that he could be deceived only by omission: $[\uparrow \neg(\neg d^+ \wedge d^-)]$. At step (d), EMILE finally informs him that deception both by commission and by omission was the right conclusion to draw: $[\uparrow (d^+ \wedge d^-)]$. By concatenating steps (c) and (d), RAYMOND can conclude that he has also been deceived on formula $(d^+ \wedge d^-)$, which is now captured by $\text{deceived}_{(d^+ \wedge d^-)}$.

EMILE’s explanation can be conceived as a way of directing RAYMOND’s attention to the two states of deception he has been prey to: firstly, $\text{deceived}_{\neg(d^+ \wedge \neg d^-)}$; secondly, $\text{deceived}_{(d^+ \wedge d^-)}$. Through the learning process, RAYMOND understands *why* and *how* he has been deceived by his malicious brother. However, as we point out hereafter, this learning also throws RAYMOND into puzzling thoughts.

4 Some Puzzling Thoughts on Deception

In this section, we show that RAYMOND’s late reasoning is puzzling but not of a paradoxical nature. That being shown, we argue that EMILE’s announcement may raise some other paradoxical issues that share echoes to the so-called “Surprise Exam Paradox” in

the epistemological and epistemic logic literature. Finally, we put more emphasis on the kind of surprise RAYMOND experiences when EMILE finally unveils his trick.

4.1 Raymond's Self-Referential Thoughts

After his brother's explanation, RAYMOND lies in bed and starts reasoning about the deception he has been subjected to. He eventually falls into a paradox that can be summed up as follows:

On the one hand, supposing that he, RAYMOND, wasn't deceived, then he didn't get what he expected (because he expected to be deceived after his brother's announcement), and hence he was actually deceived. But on the other hand, supposing that he was deceived, then he exactly did get what he expected, and hence he was not deceived stricto sensu.

We argue that this paradox is specious by expressing RAYMOND's two lines of reasoning in $\mathcal{L}_{(\mathbf{B}, \uparrow)}$. More precisely, only his first line of his reasoning ("On the one hand...") seems to bring about a contradiction. His second line of reasoning ("On the other hand...") is not contradictory as it stands.

Let us start by examining this second line before considering the first. RAYMOND's second line of reasoning can be reformulated as follows: if RAYMOND supposes that he is deceived by his brother, he has to suppose that formula **deceived** is true after having processed EMILE's announcement D : $[\uparrow D]\mathbf{deceived}$. But then, RAYMOND *exactly did get what he expected* since he expected to be deceived after his brother's announcement: $[\uparrow D]\mathbf{B}(\mathbf{deceived})$. From those two premises, the derivation hereafter reflects RAYMOND's conclusion that he is not deceived after all: $[\uparrow D]\neg\mathbf{deceived}$. By applying the reduction axiom for negation, i.e. axiom (i), we obtain a straightforward contradiction: $\neg[\uparrow D]\mathbf{deceived}$. Let us examine this derivation and look for some mistaken step:

- | | | |
|----|--|---|
| 1. | $[\uparrow D]\mathbf{deceived}$ | [Hypothesis] |
| 2. | $[\uparrow D]\mathbf{B}(\mathbf{deceived})$ | [Hypothesis] |
| 3. | $[\uparrow D]\mathbf{deceived} \wedge [\uparrow D]\mathbf{B}(\mathbf{deceived})$ | 1-2, Conjunction |
| 4. | $[\uparrow D](\mathbf{deceived} \wedge \mathbf{B}(\mathbf{deceived}))$ | 3, Axiom (ii) for $\mathcal{L}_{(\mathbf{B}, \uparrow)}$ |
| 5. | $[\uparrow D]\neg\mathbf{deceived}_{\mathbf{deceived}}$ | 4, Failure of Definition 3 |
| 6. | $[\uparrow D]\neg\mathbf{deceived}$ | 5, Failure of Definition 4 |
| 7. | $\neg[\uparrow D]\mathbf{deceived}$ | 6, Axiom (i) for $\mathcal{L}_{(\mathbf{B}, \uparrow)}$. Contra 1 |

At step 6, it is logically incorrect to derive that $[\uparrow D]\neg\mathbf{deceived}$ by affirming that, after $[\uparrow D]$, if RAYMOND is not deceived on formula **deceived** itself, i.e. $\neg\mathbf{deceived}_{\mathbf{deceived}}$, then he is not deceived simpliciter, i.e. $\neg\mathbf{deceived}$. This inference is too strong since RAYMOND may perfectly be deceived simpliciter, i.e. **deceived**, because he is deceived on formulas distinction from formula **deceived** itself. This is exactly what happens in Smullyan's story in which RAYMOND is deceived on formulas $\neg(\mathbf{d}^+ \wedge \neg\mathbf{d}^-)$ and $(\mathbf{d}^+ \wedge \mathbf{d}^-)$, and thus deceived simpliciter. Therefore, the conclusion $\neg[\uparrow D]\mathbf{deceived}$ that RAYMOND obtains at step 7 is not warranted: he cannot logically conclude that after upgrading his beliefs with his brother's announcement that D , he is deceived (step 1) *only if* he is not deceived (step 7). So RAYMOND's second line of reasoning does not lead to a genuine contradiction.

Concerning RAYMOND's first line of reasoning, this line *does* lead to a genuine contradiction but only because one of RAYMOND's assumptions, viz. that he is not deceived after his brother's announcement, is false. RAYMOND's first line of reasoning can be reformulated as follows: on the one hand, if RAYMOND supposes that he is not deceived by

EMILE, he has to assume that formula $\neg\text{deceived}$ is true after his brother's announcement D : $[\uparrow D]\neg\text{deceived}$. But then, *he did not get what he expected* since he clearly expected to be deceived after his brother's announcement: $[\uparrow D]\mathbf{B}(\text{deceived})$. From those two premises, we derive a contradiction as follows:

- | | |
|---|---|
| 1. $[\uparrow D]\neg\text{deceived}$ | [Hypothesis] |
| 2. $[\uparrow D]\mathbf{B}(\text{deceived})$ | [Hypothesis] |
| 3. $[\uparrow D]\neg\text{deceived} \wedge [\uparrow D]\mathbf{B}(\text{deceived})$ | 1-2, Conjunction |
| 4. $[\uparrow D](\neg\text{deceived} \wedge \mathbf{B}(\text{deceived}))$ | 3, Axiom (ii) for $\mathcal{L}_{(\mathbf{B}, \uparrow)}$ |
| 5. $[\uparrow D]\text{deceived}_{\neg\text{deceived}}$ | 4, Applying Definition 3 |
| 6. $[\uparrow D]\text{deceived}$ | 5, Applying Definition 4 |
| 7. $\neg[\uparrow D]\text{deceived}$ | 1, Axiom (i) for $\mathcal{L}_{(\mathbf{B}, \uparrow)}$. Contra 6 |

Though RAYMOND's reasoning is perfectly sound here, premise 1, that is to say $[\uparrow D]\neg\text{deceived}$, is false. As argued in subsection 3.2.1, it is false that RAYMOND is not deceived after his brother's announcement D . In particular, we know that RAYMOND is deceived by the default interpretation he makes of EMILE's announcement D as meaning deception only by commission: $[\uparrow D]\text{deceived}_{\neg(d^+ \wedge \neg d^-)}$. Applying definition 4 on formula $\text{deceived}_{\neg(d^+ \wedge \neg d^-)}$, we have $[\uparrow D]\text{deceived}$, contrary to the assumption at step 1. In other words, the contradiction in RAYMOND's reasoning can be bypassed since it relies on the false assumption that $[\uparrow D]\neg\text{deceived}$.

To sum up, RAYMOND's evening thoughts are not of a paradoxical nature. In his first line of reasoning, the assumption that he is not deceived after his brother's announcement can be rejected and in RAYMOND's second line of reasoning, the step leading to a contradiction can be blocked. However, we will see now that paradox may strike back if we rephrase EMILE's announcement in a slightly different way putting more emphasis on another dimension of EMILE's announcement: the surprise EMILE announces to RAYMOND when saying that he will deceive him "*as [he has] never been deceived before!*". Such a surprise is noticeable in RAYMOND's late reasoning since he is surprised by the *type of deception* EMILE used: $(d^+ \wedge d^-)$. Let us now investigate on this dimension of surprise behind EMILE's announcement D .

4.2 Emile's Surprise Deception

4.2.1 Expressing Surprise in Smullyan's Story

We start out by reminding EMILE's early announcement: "*Well, Raymond, today is April Fool's Day, and I will deceive you as you have never been deceived before!*". We have used a specific formula, i.e. D (Definition 9), to abbreviate this announcement and, then, defined more complex formulas, i.e. d^+ and d^- (Definitions 5 to 8), to express deception by commission and deception by omission, respectively. Using abbreviation D was sufficient at this step to express EMILE's morning announcement of deception. That said, it left aside the dimension of surprise involved in EMILE's announcement, as conveyed by the formulation "... *as you have never been deceived before!*". To capture this surprise aspect of EMILE's announcement, we now propose to use another formula, written \mathbf{D} , encoding the specific fact that EMILE's deception will be indeed surprising for RAYMOND.

We may ask which explicit formal definition can be given to \mathbf{D} . Extant analyzes of the so-called "Surprise Examination Paradox" [e.g. 41, 47, 48] provide some path in that respect. In the surprise exam puzzle, a teacher announces to her students that she will give a (single) surprise exam the next week. This announcement can be compared

with EMILE’s announcement D along similar lines. They both announce a surprising event to come and we will see that announcing D raises similar issues as the teacher’s announcement does in the surprise exam case. To perform this comparison, we focus on the *three-day* version of the surprise exam paradox since, in Smullyan’s story, there are also three distinct ways in which RAYMOND can be deceived: deception only by commission ($d^+ \wedge \neg d^-$), deception only by omission ($\neg d^+ \wedge d^-$) and deception by both commission and omission ($d^+ \wedge d^-$). Based on those possibilities, we provide a formal expression of EMILE’s announcement based on the formalizations given by [28, 29] and [20] in the surprise exam case.⁷⁸

In the *three-day* version of the puzzle, the teacher makes the following announcement to the students: “*I will give you a single surprise exam either on Monday or on Tuesday or on Wednesday next week*”. Between the teacher and the students, it is common knowledge that an exam comes as a surprise if the students do not know the evening before the exam that it will happen the next day. But after this announcement, the students, who are assumed to be perfectly rational, start reasoning and finally conclude that such an exam is impossible in principle. Logically speaking, the exam would fail to be a surprise on any day of the week. The classical argument for this failure is based on so-called *backward induction* and runs as follows: if the test takes place on Wednesday, the students would know on Tuesday evening that it is on Wednesday because Wednesday is the last available working day of the week. Then, it would not be a surprise. But nor will it be a surprise if the test is on Tuesday because the students would know on Monday evening that the test is on Tuesday (for they know that the test is not on Wednesday due to the previous reasoning). The same argument applies such that the exam cannot be on Monday either. At the end of this backward reasoning, the students would conclude that the test cannot be a surprise on any day of the week. Unfortunately, the teacher unexpectedly gives the exam on Tuesday (for instance) and this is a complete surprise.

Given the various deceptive options EMILE has, a more explicit, yet informal, formulation of EMILE’s announcement D might be: “*I will surprise you by deceiving you only by commission, or only by omission, or both by commission and by omission*”. We know that Raymond misinterprets this explicit announcement in the following sense: “*I will surprise you by deceiving you only by commission*”. But from a strictly semantic point of view, all the three deceptive options, namely *deception only by commission*, *deception only by omission* and *deception both by commission and by omission*, are possible surprising options for RAYMOND and should be taken into account to capture the meaning of D .

This informal formulation of EMILE’s announcement D can result in many formal expressions depending on the way the concept of *surprise* is formalized. Following a basic intuition on which we elaborate later (see subsection 4.2.3), a state of surprise is generally triggered by a *mismatch of beliefs* between expectations and occurring events. Indeed, the agent is surprised by some events if those events come out to be true after

⁷Those formalizations for the surprise exam rely on using a knowledge modal operator, say K , for which we could provide the following interpretation in our setting: $\mathcal{S}, s \models K\varphi$ iff for all $t \in \mathcal{S}$, if $s \sim t$ then $\mathcal{S}, t \models \varphi$, with \sim a comparability relation over \mathcal{S} such that $\sim := \leq \cup \geq$ (where \geq is the converse of \leq). Notice that the way it is defined, \sim is a *symmetric* preorder over \mathcal{S} , and thus an *equivalence relation*, since for two states s and t , the conditions for relations $s \sim t$ and $s \sim t$ to obtain are identical. Accordingly, our knowledge operator K is S5, as in [20], while the operator used in [28, 29] is K45. But since our analysis relies on doxastic operator B in $\mathcal{L}_{(B, \uparrow)}$, this difference has no impact on our analysis.

⁸See also [10] and [13] for discussions of epistemic puzzles related to the Surprise Exam Paradox in dynamic logic settings, such as Fitch’s paradox, Moore’s paradox and Williamson’s Margin for Error paradox.

having remained unacknowledged by agents until they come out to be true. As a result, the agent experiences states of *surprise*, or even of *astonishment*, when she learns the existence of those facts and events. To see more clearly into this, let us first provide the formal abbreviation of the announcement D we favor in Definition 12.

Definition 12 (*surprise-deception reading of Emile’s announcement*) We write D the reading of EMILE’s morning announcement insisting on the surprise dimension of the announcement. That is:

$$\begin{aligned} D & := ((d^+ \wedge \neg d^-) \wedge \neg \mathbf{B}(d^+ \wedge \neg d^-)) \\ & \quad \vee ((\neg d^+ \wedge d^-) \wedge [\uparrow \neg(d^+ \wedge \neg d^-)] \neg \mathbf{B}(\neg d^+ \wedge d^-)) \\ & \quad \vee ((d^+ \wedge d^-) \wedge [\uparrow \neg(d^+ \wedge \neg d^-)][\uparrow \neg(\neg d^+ \wedge d^-)] \neg \mathbf{B}(d^+ \wedge d^-)) \end{aligned}$$

The abbreviation provided in Definition 12 is inspired by Gerbrandy’s formalization of the teacher’s announcement in the Surprise Exam case, using belief radical upgrades \uparrow as proposed in [5, 8] instead of announcement operators “!”.⁹ Firstly, formula D states that EMILE’s deception will be a surprise since, according to D , if RAYMOND is deceived only by commission, formula $(d^+ \wedge \neg d^-)$ is true but RAYMOND denies it to be true: $\neg \mathbf{B}(d^+ \wedge \neg d^-)$. Secondly, in case RAYMOND is deceived only by omission, now formula $(\neg d^+ \wedge d^-)$ is true. But again, it will be a surprise because, after having recognized that he would not be deceived only by commission: $[\uparrow \neg(d^+ \wedge \neg d^-)]$, RAYMOND denies $(\neg d^+ \wedge d^-)$ to be true: $\neg \mathbf{B}(\neg d^+ \wedge d^-)$. Thirdly, formula D states that if RAYMOND is deceived both by commission and by omission, i.e. $(d^+ \wedge d^-)$, it will be a surprise too because, after having recognized that $\neg(d^+ \wedge \neg d^-)$ and then that $\neg(\neg d^+ \wedge d^-)$: $[\uparrow \neg(d^+ \wedge \neg d^-)][\uparrow \neg(\neg d^+ \wedge d^-)]$, RAYMOND (still) denies $(d^+ \wedge d^-)$ to be true: $\neg \mathbf{B}(d^+ \wedge d^-)$. In other words, following D , no matter whether RAYMOND is deceived only by commission or only by omission, or by both, in all cases it will come as a surprise.

That being said, however, none of the deception types encoded in formula D will surprise RAYMOND. Again backward induction helps explain why. Suppose that RAYMOND has not been deceived only by commission, i.e. $(d^+ \wedge \neg d^-)$, or only by omission, i.e. $(\neg d^+ \wedge d^-)$. So the only remaining option for RAYMOND to be deceived is to be deceived both by commission and by omission: $(d^+ \wedge d^-)$. But then, RAYMOND cannot be surprised by $(d^+ \wedge d^-)$ since this is the only remaining option for him to be deceived after having rejected options $(d^+ \wedge \neg d^-)$ and $(\neg d^+ \wedge d^-)$. Suppose now that RAYMOND is deceived only by omission: $(\neg d^+ \wedge d^-)$. This implies that deception only by commission, i.e. option $(d^+ \wedge \neg d^-)$, has not taken place, so the only remaining possibility for RAYMOND to be deceived is $(\neg d^+ \wedge d^-)$, and here again he cannot be surprised by this option. Having rejected options $(d^+ \wedge d^-)$ and $(\neg d^+ \wedge d^-)$ as possible sources of surprise, suppose finally that RAYMOND is deceived only by commission: $(d^+ \wedge \neg d^-)$. Since $(d^+ \wedge \neg d^-)$ is the only remaining option for RAYMOND to be surprised, he won’t be surprised by it. Then, if RAYMOND’s beliefs remain consistent, he cannot be surprised at all after learning EMILE’s morning announcement D .

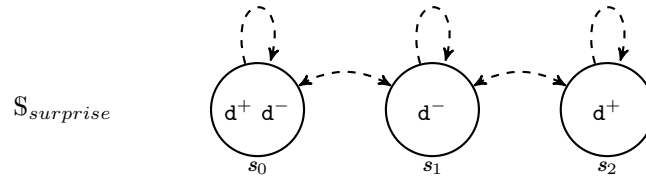
⁹A (three-day) Gerbrandy’s formalization of the teacher’s announcement of surprise with announcement operators would be: $(mon \wedge \neg \mathbf{K} mon) \vee (tue \wedge [!\neg mon] \neg \mathbf{K} tue) \vee (wed \wedge [!\neg mon][!\neg tue] \neg \mathbf{K} wed) \vee \mathbf{K} \perp$, such that “mon”, “tue” and “wed” stand for *monday*, *tuesday* and *wednesday*, \mathbf{K} is a K45 knowledge operator and “!” is an announcement operator. This formalization, written \mathbf{S} by Gerbrandy, means that the exam will come as a surprise for the reason that if the exam is on wednesday (*wed*), the agent does not know this ($\neg \mathbf{K} wed$) after having learned that it is not on tuesday ($[!\neg tue]$); if the exam is on tuesday (*tue*), the agent does not know this ($\neg \mathbf{K} tue$) after having learned that it is not on monday ($[!\neg mon]$); if the exam is on monday (*mon*), the agent does not know it either on monday ($\neg \mathbf{K} mon$); if the agent has contradictory information about the date of the exam (\perp), the latter will also come as a surprise (according to Gerbrandy). Using the reduction axioms of dynamic epistemic logic, Gerbrandy rewrite \mathbf{S} as follows: $(mon \wedge \neg \mathbf{K} mon) \vee (tue \wedge \neg \mathbf{K}(\neg mon \rightarrow tue)) \vee (wed \wedge \neg \mathbf{K}(\neg mon \wedge \neg tue \rightarrow wed)) \vee \mathbf{K} \perp$.

4.2.2 Emile’s Announcement of Surprise

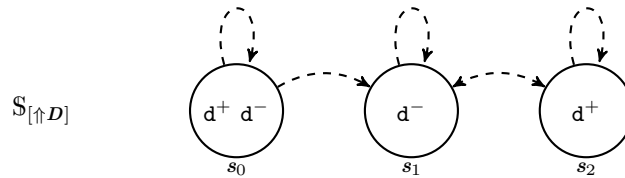
As pointed by Baltag & Smets in [5, 8], the most accepted reading of the surprise announcement is to understand the announcement as a self-referential statement, as in Gerbrandy’s third expression of the teacher’s announcement. In the context of Smullyan’s story, this amounts to the following reading of EMILE’s announcement: “*I will deceive you as you have never been deceived before, even after I am telling you all this*”. This self-referential statement can be interpreted as an infinite iteration of non-self-referential upgrades as follows:

$$[\uparrow D]; [\uparrow D]; [\uparrow D]; \dots$$

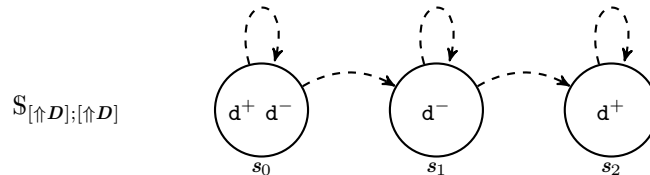
Informally, this sequence first says that EMILE’s deception will be a surprise even after EMILE announcing that it will be a surprise, then it says that even after this announcement, EMILE’s deception will be a surprise even after announcing it a second time, etc. To model this sequence, let us set aside the state in which RAYMOND is not deceived (s_3) and only consider the states where RAYMOND is actually deceived: s_0 , s_1 , s_2 . Initially, before EMILE’s announcement, all those states are considered to be equally plausible by RAYMOND, as shown in the following model:¹⁰



Let us consider the model $\mathbb{S}_{[\uparrow D]}$ obtained after a first iteration of $[\uparrow D]$ on $\mathbb{S}_{surprise}$. In $\mathbb{S}_{[\uparrow D]}$, all D -states become more plausible than all $\neg D$ -states, strictly reducing to state s_0 here.



Let us now consider the model $\mathbb{S}_{[\uparrow D];[\uparrow D]}$ obtained after the second iteration of $[\uparrow D]$, now applied on model $\mathbb{S}_{[\uparrow D]}$. In $\mathbb{S}_{[\uparrow D];[\uparrow D]}$, again, all D -states from $\mathbb{S}_{[\uparrow D]}$ are made more plausible than all $\neg D$ -states.



Now, a fixed-point has been reached with respect to the announcement of surprise D : any new update of $\mathbb{S}_{[\uparrow D];[\uparrow D]}$ with further operations $[\uparrow D]$ will leave model $\mathbb{S}_{[\uparrow D];[\uparrow D]}$

¹⁰As in subsection 3.3 when modelling EMILE’s explanation, we do not show the transitivity of plausibility orders in any of the models to enhance readability.

unchanged. The reason is that in model $\mathbb{S}_{[\uparrow D];[\uparrow D]}$, update $[\uparrow D]$ is still infinitely executable but the result is empty: D is false, RAYMOND believes it to be false and surprise cannot happen.

Based on this conclusion that D cannot be true, we may wonder why RAYMOND does not simply dismiss all the iterated upgrades $[\uparrow D]$ he has made and thus go back to his original plausibility order expressed in model $\mathbb{S}_{surprise}$. However, as argued by [5] in the surprise examination case, RAYMOND would then cancel all the reasons that led him to conclude that EMILE's announcement D was false. If he reverts to the initial model $\mathbb{S}_{surprise}$, then he no longer knows, more precisely believes, that EMILE's announcement is false, and then surprise is possible. That RAYMOND concludes EMILE's announcement to be false does not provide RAYMOND any warrant to dismiss all the iterated radical upgrades he made since, precisely, this conclusion is warranted by all those upgrades. In other words, RAYMOND should rationally stick to final model $\mathbb{S}_{[\uparrow D];[\uparrow D]}$ where surprise fails in principle.

So to end this section, let us now explain why, contrary to this conclusion, RAYMOND is actually surprised by EMILE. Subsection 3.3 where EMILE's explanation is modelled explains why and how RAYMOND is in fact surprised in a way consistent with Gerbrandy's reading of surprise D . As one can observe in model $\mathbb{S}_{(d^+ \wedge d^-)}$, since $(d^+ \wedge d^-)$ is true in the actual state s_0 :

$$\mathbb{S}_{(d^+ \wedge d^-)}, s_0 \models (d^+ \wedge d^-)$$

But we can also verify that the following dynamic formula is true:

$$\mathbb{S}_{(d^+ \wedge d^-)}, s_0 \models [\uparrow \neg(d^+ \wedge \neg d^-)][\uparrow \neg(\neg d^+ \wedge d^-)]B\neg(d^+ \wedge d^-)$$

So, we have by conjunction:

$$\mathbb{S}_{(d^+ \wedge d^-)}, s_0 \models (d^+ \wedge d^-) \wedge [\uparrow \neg(d^+ \wedge \neg d^-)][\uparrow \neg(\neg d^+ \wedge d^-)]B\neg(d^+ \wedge d^-)$$

Crucially, this conjunction admits formula D as a semantic consequence, by the consistency axiom for B :

$$(d^+ \wedge d^-) \wedge [\uparrow \neg(d^+ \wedge \neg d^-)][\uparrow \neg(\neg d^+ \wedge d^-)]B\neg(d^+ \wedge d^-) \models D$$

Then, applying the semantic deduction theorem, we have:

$$\mathbb{S}_{(d^+ \wedge d^-)}, s_0 \models D$$

In other words, following Gerbrandy's notion of surprise we encoded by D in our context, the source of RAYMOND's surprise is EMILE revealing that he was deceived both by commission and by omission. In other words, RAYMOND is surprised exactly by learning that he was deceived by option $(d^+ \wedge d^-)$ he had not anticipated at all. In the next subsection, however, we argue that this surprise à-la-Gerbrandy is not the only kind of surprise RAYMOND is preyed to in Smullyan's tale.

4.2.3 Raymond's Two States of Surprise

In line of one of the Gerbrandy's reading of surprise [29], RAYMOND is surprised when EMILE unveils his April Fool's deceptive trick at the end of the day. Here we rely on conceptual analysis to show that RAYMOND also encounters another form of surprise which is weaker than the one resulting from EMILE's explanation.

According to [e.g. 40, 36], a cognitive agent reaches a state of surprise when she is led to recognize an *inconsistency*, that is to say a *discrepancy* or a *mismatch*, between what she believes about the world and the actual state of the world. Following a basic intuition, an event which is unexpected is considered “surprising” and the more unexpected it is, the more surprising the event turns out to be [e.g. 42, 40]. In other words, an event that was expected but does not happen *is* surprising. But an event which actually happens but was totally unexpected is *even more* surprising.

In that sense, Lorini & Castelfranchi distinguish two main kinds of surprise. The first kind is called “*mismatch-based surprise*” and results from a “*conflict between a perceived fact and a scrutinized representation*”. In that case, the agent is surprised because she has some anticipatory representation of a fact or event, but she cannot make the incoming data fit with this anticipation. Accordingly, the intensity of the induced surprise depends of the probability, more crucially on the *implausibility*, which is assigned by the agent to the conflicting data she receives.

But Lorini & Castelfranchi contrast this first kind of surprise with a stronger kind named “*astonishment*” (or “*surprise in recognition*”). Contrary to the first kind, this type of surprise is of a *second-order nature* since now it is rooted on the *recognition* of the implausibility of a perceived fact compared to expectations: “[we] *perceive a certain fact and recognize the implausibility of this*”. Lorini & Castelfranchi also precise that this astonishment can depend upon two distinct mental processes. Firstly, we can be astonished by a fact/event φ because we assigned a high probability to $\neg\varphi$ and after perceiving φ we realize that we would not have expected that event φ [see 36, 3]. Second, we can be astonished by φ in case perceiving φ makes us infer the falsity (and incongruity) of our initial disbelief that φ or belief that $\neg\varphi$.¹¹

In Smullyan’s tale, both kinds of deeper surprise are involved because both forms of unexpectation are involved. In the previous subsection based on Gerbrandy’s reading of surprise, EMILE’s explanation leads RAYMOND to be *astonished* (or *surprised in recognition*) when EMILE tells him that he has been deceived both by commission and by omission. Since RAYMOND kept discarding this configuration as very implausible all along the day, RAYMOND is astonished when he recognized that EMILE used this option to trick him. But in fact, aside from this Gerbrandy’s reading, RAYMOND is also subjected to mismatch-based surprise on April 1st 1925. As we know, RAYMOND expects to be deceived only by commission after EMILE’s announcement and until his late explanation, but this does not happen. As a result of this absence, RAYMOND also encounters mismatch-based surprise before being astonished later.

Consistent with Lorini & Castelfranchi’s distinction is the fact that RAYMOND’s *astonishment* is stronger in terms of surprise compared to his state of *mismatch-based surprise*. Astonishment follows from his strong disbelief that he could be deceived both by commission and by omission. His belief that he would be deceived only by commission is less entrenched and only leads him to moderate surprise.

¹¹Lorini & Castelfranchi contrast those “*deeper and slower forms of surprise which are due to symbolic representations of expected events*”[see 36, 1] to a *first-hand surprise* which corresponds to a perceptual mismatch between a stimulus (namely between the agents can *see* or *hear* in their immediate environment) and what the agents expect from a sensory-motor perspective.

5 Conclusion

We used dynamic belief revision to analyze a puzzle in which Raymond Smullyan’s brother EMILE uses a deceptive stratagem to fool him on April 1st 1925. By a morning announcement, EMILE misleads RAYMOND’s expectations to be deceived (only) by commission. But in doing so, he also sets the stage for deceiving him by omission too. This stratagem leads RAYMOND to paradoxical thoughts and strong surprise. But contrary to what it seems, RAYMOND’s thoughts are not paradoxical: only the success of EMILE’s announcement is. Based on a formal comparison with the ‘Surprise Exam Paradox’, we finally observed that EMILE’s deception by commission and omission is the main source of RAYMOND’s surprise.

RAYMOND’s reasoning on surprise is flawless but it results in the false conclusion that EMILE will not surprise him. Of course, this should not invite RAYMOND to dismiss all the iterated radical upgrades he has made on April 1st 1925 [see 5], but one lesson for *future* RAYMOND is to be more cautious with EMILE’s new announcements. An interesting investigation in that respect would be to consider so-called “improvement operators” [34, 14] in the future. As observed for Moorean sentences [see e.g. 7], RAYMOND’s mistaken conclusion of no surprise is an illustration that the “*success postulate*” of the AGM framework [2] can be problematic. Still, formulas that do not give rise to any Moorean phenomena are successful in those settings. Improvement operators are less permissive in that respect since some versions authorize the plausibility of the information to increase iteratively without being accepted immediately, — no matter the logical form of the incoming formula. This specificity aligns with the realistic scenario where RAYMOND should now treat all EMILE’s announcements with cautious skepticism.

The theoretical failure of EMILE’s announcement also raises a more epistemological discussion on whether EMILE’s announcement qualifies as a lie or not. Traditionally, lying consists in making “*a believed-false statement to another person with the intention that the other person believe that statement to be true*” [37]. More precisely, a speaker X lies to a hearer Y on a proposition p if and only if (i) X believes that p is false (*untruthfulness clause*), (ii) X intends Y to believe that p is true (*intention-to-deceive clause*), and (iii) X tells Y that p (*addressed statement clause*).¹² Following this so-called ‘*subjective account*’ then, EMILE’s announcement cannot be considered a lie since condition (i) is not met: EMILE fails to be *untruthful* since he believes his morning utterance to be true. So, EMILE is not lying in the classical sense.

That said, EMILE’s announcement also a default interpretation leading RAYMOND to be deceived. Conceptually, such a misleading rule pertains to a broader category we may call *veridical deception*. In veridical deception, a non-cooperative speaker makes an utterance that is (objectively) true but whose pragmatic interpretation is false and intended to deceive some addressee. Veridical deception also includes *false implicatures*, *pretending to deceive* and *presupposition faking*. In *false implicatures*, speakers aim to make addressees believe a false information φ by telling them some true information ψ which conversationally implicates φ [e.g. 1, 38, 39]. With *pretending to deceive* [e.g. 53, 25], speakers make addressees disbelieve a true piece of information φ by arousing

¹²In more details, condition (i) requires that the speaker be insincere (or dishonest) with respect to p : she believes that p is false. Condition (ii) requires the speaker to expect her addressee to believe that p is true. Finally, condition (iii) requires the speaker to assert that p to some specific addressee for her utterance to count as a lie. Note that in this subjective definition, it is not required that content p be *objectively* false: according to condition (i), it is sufficient that the speaker *subjectively* believe that p is false. In that sense, one can perfectly lie while saying the truth.

their suspicion towards the source of φ . Finally, *presupposition faking* [e.g. 31, 53], consists in making addressees inadequately believe a true formula φ that the speakers fail to semantically or pragmatically account for.

Classically, lying is considered a strict semantic phenomenon, contrary to veridical deception which involves pragmatics. What matters is that the liar believe the semantic content of her utterance to be false (as in the *untruthfulness clause*), — no matter whether she believes the pragmatic interpretations of this utterance is false or believed to be false. By contrast, in veridical deception, and in EMILE’s misleading default in particular, the speaker believes the semantic content of her utterance to be true, although its pragmatic meaning is false, or believed to be false. Accordingly, veridical deception does not involve any untruthfulness clause.

The relationship between lying and veridical deception has been discussed concerning false implicatures in particular [e.g. 23, 1, 39, 22, 52, 55]. For most authors, such as [24] and [50], false implicatures are *not lies* since the semantic meaning of the speaker’s utterance is not defective, though its pragmatic meaning is. But for some others [e.g. 1, 39, 22, 52], untruthful implicatures are lies, consistent with the Gricean view that truthful statements which conversationally implicate believed-to-be-false statements *are* lies. Since the distinction between defaults and implicatures is debated [e.g. 35, 44, 3, 32], the question of whether misleading defaults are lies also remains open. Conducting empirical studies would be useful to help adjudicate on this issue. We leave this investigation for future work.

Acknowledgments

We are grateful to Alexandru Baltag, Timothée Bernard, Quentin Blomet, Denis Bonnay, Philippe Capet, Hans van Ditmarsch, Paul Égré, Raul Fervari, Jakob Süsskind, Peter van Emde Boas, and Sonja Smets for stimulating exchanges and comments on intermediary versions of the manuscript. This work was supported by the programs HYBRINFOX (ANR-21-ASIA-0003), FRONTCOG (ANR-17-EURE-0017), and THEMIS (DOS-0222794/00 and DOS-0222795/00).

References

- [1] Jonathan E Adler. Lying, deceiving, or falsely implicating. *The Journal of Philosophy*, 94(9):435–452, 1997.
- [2] Carlos E Alchourrón, Peter Gärdenfors, and David Makinson. On the logic of theory change: Partial meet contraction and revision functions. *The journal of symbolic logic*, 50(2):510–530, 1985.
- [3] Kent Bach. Semantic slack: What is said and more. *Foundations of speech act theory: Philosophical and linguistic perspectives*, pages 267–291, 1994.
- [4] Alexandru Baltag, Lawrence S. Moss, and Slawomir Solecki. The logic of public announcements, common knowledge, and private suspicions. *Proceeding TARK ’98 Proceedings of the 7th conference on Theoretical aspects of rationality and knowledge*, 1998.
- [5] Alexandru Baltag and S Smets. Surprise?! an answer to the hangman, or how to avoid unexpected exams! In *Logic and Interactive Rationality Seminar (LIRA), slides*, 2009.

- [6] Alexandru Baltag and Sonja Smets. Dynamic belief revision over multi-agent plausibility models. In *Proceedings of LOFT*, volume 6, pages 11–24, 2006.
- [7] Alexandru Baltag and Sonja Smets. A qualitative theory of dynamic interactive belief revision. *Logic and the foundations of game and decision theory (LOFT 7)*, 3:9–58, 2008.
- [8] Alexandru Baltag and Sonja Smets. Multi-agent belief dynamics. *course given at NASSLLI*, 2010.
- [9] Alexandru Baltag and Sonja Smets. Keep changing your beliefs, aiming for the truth. *Erkenntnis*, 75(2):255, 2011.
- [10] Johan van Benthem. What one may come to know. *Analysis*, 64(282):95–105, 2004.
- [11] Johan van Benthem. Dynamic logic for belief revision. *Journal of applied non-classical logics*, 17(2):129–155, 2007.
- [12] Johan van Benthem and Sonja Smets. Dynamic logics of belief change. In *Handbook of Logics for Knowledge and Belief*, pages 299–368. College Publications, 2015.
- [13] Denis Bonnay and Paul Égré. Knowing one’s limits: An analysis in centered dynamic epistemic logic. In Patrick Girard, Olivier Roy, and Mathieu Marion, editors, *Dynamic Formal Epistemology*, pages 103–126. Springer, 2011.
- [14] Richard Booth, Eduardo L Fermé, Sébastien Konieczny, and Ramón Pino Pérez. Credibility-limited improvement operators. In *ECAI*, volume 263, pages 123–128, 2014.
- [15] Gauvain Bourgne, Camilo Sarmiento, and Jean-Gabriel Ganascia. Ace modular framework for computational ethics: dealing with multiple actions, concurrency and omission. In *International Workshop on Computational Machine Ethics*, 2021.
- [16] Ilaria Canavotto. *Where responsibility takes you*. PhD thesis, Springer, 2020.
- [17] Roderick M. Chisholm and Thomas D. Feehan. The intent to deceive. *Journal of Philosophy*, 74(3):143–159, 1977.
- [18] Randolph Clarke. What is an omission? *Philosophical Issues*, 22:127–143, 2012.
- [19] Randolph K Clarke. *Omissions: Agency, metaphysics, and responsibility*. Oxford University Press, 2014.
- [20] Hans P. van Ditmarsch and Barteld Kooi. The secret of my success. *Synthese*, 153(2):339–339, 2006.
- [21] Hein Duijf. *The logic of responsibility voids*. Springer, 2022.
- [22] Marta Dynel. A web of deceit: A neo-gricean view on types of verbal deception. *International Review of Pragmatics*, 3(2):139–167, 2011.
- [23] Don Fallis. A conceptual analysis of disinformation. *iConference 2009 Papers*, 2009.
- [24] Don Fallis. What is lying? *The Journal of Philosophy*, 106(1):29–56, 2009.
- [25] Don Fallis. The varieties of disinformation. In *The Philosophy of Information Quality*, pages 135–161. Springer, 2014.

- [26] Tim French and Hans Van Ditmarsch. Undecidability for arbitrary public announcement logic. In *Advances in modal logic*, pages 23–42. College publications, 2008.
- [27] Peter Gärdenfors and Hans Rott. Belief revision. *Computational Complexity*, 63(6):35–132, 1995.
- [28] Jelle Gerbrandy. *Bisimulations on planet Kripke*. ILLC Dissertation Series, 1999.
- [29] Jelle Gerbrandy. The surprise examination in dynamic epistemic logic. *Synthese*, 155(1):21–33, 2007.
- [30] Jelle Gerbrandy and Willem Groeneveld. Reasoning about information change. *Journal of logic, language and information*, 6(2):147–169, 1997.
- [31] Peter Harder and Christian Erik J. Kock. *The theory of presupposition failure*. Akademisk forlag, 1976.
- [32] Laurence R Horn. The border wars: A neo-gricean perspective. *Where semantics meets pragmatics*, 16:21–48, 2006.
- [33] Aaron Hunter and François Schwarzentruber. Arbitrary announcements in propositional belief revision. In *DARe@IJCAI*, 2015.
- [34] Sébastien Konieczny and Ramón Pino Pérez. Improvement operators. *KR*, 8:177–187, 2008.
- [35] Stephen C Levinson. *Presumptive meanings: The theory of generalized conversational implicature*. MIT press, 2000.
- [36] Emiliano Lorini and Cristiano Castelfranchi. The unexpected aspects of surprise. *International Journal of Pattern Recognition and Artificial Intelligence*, 20(06):817–833, 2006.
- [37] James E Mahon. The definition of lying and deception. *The Stanford Encyclopedia of Philosophy*, 2015.
- [38] Jörg Meibauer. Lying and falsely implicating. *Journal of Pragmatics*, 37(9):1373–1399, 2005.
- [39] Jörg Meibauer. *Lying at the semantics-pragmatics interface*. Walter de Gruyter GmbH & Co KG, 2014.
- [40] Wulf-Uwe Meyer, Rainer Reisenzein, and Achim Schützwohl. Toward a process analysis of emotions: The case of surprise. *Motivation and Emotion*, 21(3):251–274, 1997.
- [41] Donald J O’Connor. Pragmatic paradoxes. *Mind*, 57(227):358–359, 1948.
- [42] Andrew Ortony and Derek Partridge. Surprisingness and expectation failure: what’s the difference? In *IJCAI*, pages 106–108, 1987.
- [43] Jan Plaza. Logics of public communications. *Synthese*, 158(2):165, 2007.
- [44] François Récanati. *Literal meaning*. Cambridge University Press, 2004.
- [45] Hans Rott. Conditionals and theory change: Revisions, expansions, and additions. *Synthese*, 81:91–113, 1989.

- [46] Gilbert Ryle. Negative 'actions'. *Hermathena*, pages 81–93, 1973.
- [47] Michael Scriven. Paradoxical announcements. *Mind*, 60(239):403–407, 1951.
- [48] Robert Shaw. The paradox of the unexpected examination. *Mind*, 67(267):382–384, 1958.
- [49] Raymond M Smullyan. *What is the Name of this Book? The Riddle of Dracula and Other Logical Puzzles: Mysteries, Paradoxes, Gödel's Discovery*. Prentice-Hall, 1978.
- [50] Roy Sorensen. Lying with conditionals. *The Philosophical Quarterly*, 62(249):820–832, 2012.
- [51] Frank Veltman. Defaults in update semantics. *Journal of philosophical logic*, 25(3):221–261, 1996.
- [52] Emanuel Viebahn. Non-literal lies. *Erkenntnis*, 82(6):1367–1380, 2017.
- [53] Jocelyne M. Vincent and Cristiano Castelfranchi. On the art of deception: how to lie while saying the truth. *Possibilities and limitations of pragmatics*, pages 749–777, 1981.
- [54] Elazar Weinryb. Omissions and responsibility. *The Philosophical Quarterly (1950-)*, 30(118):1–18, 1980.
- [55] Alex Wiegmann and Pascale Willemsen. How the truth can make a great lie: An empirical investigation of lying by falsely implicating. In *Proceedings of the 39th Annual Conference of the Cognitive Science Society*, pages 3516–3621, 2017.