# On Making Counterfactual Assumptions:

## Veltman (2005) and Beyond

William B. Starr

*Dept. of Philosophy, Rutgers University*
*26 Nichol Ave.*
*New Brunswick, NJ 08904*

## 1  Background & Motivation

****All formal definitions can be found together in the last section****

### 1.1  What the Stalnaker-Lewis Theory Can't Explain

- Tichý (1976: 271) proposed the following discourse as a counterexample to Stalnaker (1968) & Lewis (1973)'s theories of counterfactuals:

  (1) a. Invariably, if it is raining, Jones wears his hat
      b. If it is not raining, Jones wears his hat at random
      c. Today, it is raining and so Jones is wearing his hat
      d. But, even if it had not been raining, Jones would have been wearing his hat

- Given (1a-c), (1d) is judged to be incoherent/false/bad
- Why is this a counterexample to the Stalnaker-Lewis theory?

  ○ According the the Stalnaker-Lewis theory a counterfactual is true in the actual world $w$ iff the consequent is true in every world in which the antecedent is true and in other respects is maximally similar to $w$

  ○ Given the constraints placed on the actual world $w$ in (1), which are the most $w$-simliar worlds where it is not raining?

  ○ Tichý argues that the most $w$-similar worlds are ones where Jones is wearing his hat, since, as (1c) says, Jones is wearing his hat in $w$

  ○ But then Stalnaker-Lewis predict that (1d) is true and so (1) should be fine!

  ○ Note that this counterexample assumes that particular matters of fact are taken into account when determining $w$-similarity

- Lewis (1979) responded with a system of weights that govern similarity

---

*Email:* wstarr@rutgers.edu.
*URL:* http://eden.rutgers.edu/~wbstarr.

- Among them was the following:

  "It is of little or no importance to secure approximate similarity of particular fact." (Lewis 1979: 472)

- Assuming Lewis' metric can be successfully implemented, this appears to solve the problem with (1), since the counterexample relied on similarity with respect to a particular fact
- Unfortunately for Lewis, there are examples which appear to require similarity to depend on particular facts: [1]

  (2) a. Jones always flips a coin before he opens the curtain to see what the weather is like
      b. If it's not raining and the coin comes up heads, he wears his hat
      c. If it's not raining and the coin comes up tails, he doesn't wear his hat
      d. Invariably, if it is raining he wears his hat
      e. Today, the coin came up heads and it is raining, so Jones is wearing his hat
      f. But, even if it hadn't been raining, Jones would have been wearing his hat

- Given (2a-e), (2f) is judged to be acceptable
- But, given Lewis' (1979) metric (2f) is predicted to be false!

  ○ (2f)'s truth requires that the closest non-raining worlds be ones where Jones is still wearing his hat, but it is a particular matter of fact that Jones is wearing his hat. So it seems that particular matters of fact do count towards similarity after all!

- There are simpler examples in the literature that illustrate this fact:

  (3) *If we had bought one more artichoke this morning, we would have had one for everyone at dinner tonight* (Sanford 1989: 173)

      ▷ This example requires that we hold fixed the particular fact of how many guests showed up for dinner tonight

  (4) (Your friend invites you to bet heads on a fair coin-toss. You decline. The coin is tossed and comes up heads.) *See, if you had bet heads you would have won!* (Slote 1978: 27 fn33) [2]

      ▷ This example requires that we hold fixed the particular fact that the coin came up heads

---

[1]  Veltman (2005: 164) attributes this example to Frank Mulkens.
Slote (1978: 27 fn33) attributes a parallel example to the late philosopher Sydney Morgenbesser.
[2]  Slote attributes this example to Sydney Morgenbesser.

## 1.2   What Went Wrong? Where are we Going?

- Veltman (2005) offers an alternative diagnosis of why (1d) is bad, while (2f) is acceptable

- When we make the counterfactual assumption *had it not rained*, we don't keep assuming that Jones wore his hat, since Jones was wearing his hat *because* it was raining

  ○ When you make a counterfactual assumption of $\phi$, you not only suspend your belief that $\neg\phi$ but also any other beliefs that depended on $\neg\phi$

- Veltman (2005: 164) puts it nicely:

   Similarity of particular fact is important, but only for facts that do not depend on other facts. Facts stand and fall together. In making a counterfactual assumption, we are prepared to give up everything that depends on something that we must give up to maintain consistency. But, we want to keep in as many independent facts as we can.

- Note that this constitutes a significant revision of **Ramsey's Test**, which motivated the Stalnaker-Lewis theory:

   This is how to evaluate a conditional: first, add the antecedent (hypothetically) to your stock of beliefs; second, make whatever adjustments are required to maintain consistency (without modifying the hypothetical belief in the antecedent); finally, consider whether or not the consequent is true.
   (Stalnaker 1968: 106)

 ▷ You also want to give up the beliefs that depend on the beliefs you have to give up to maintain consistency

- Kratzer (1989) developed **lumping semantics** to account for this kind of phenomena, but her approach faces technical difficulties. (Kanazawa *et al.* 2005; Kratzer 1990, 2002, 2005)

- Veltman's (2005) theory provides an alternative formalization that still captures this intuitive diagnosis

- Veltman's theory is developed in the framework of **Update Semantics** (Veltman 1996), where the meaning of a formula is taken to be an operation on cognitive states

 ▷ This more dynamic & procedural perspective on meaning makes the relationship between recipes like Ramsey's Test and the content of the formal theory utterly transparent. *Notice how untrue this is of the Stalnaker-Lewis theory!*

 ▷ It also allows one to interpret counterfactuals as a simple sequence of operations: *if had been* $\phi$ creates a subordinating info state $S'$ & then *it would have been* $\psi$ is tested in $S'$. This allows the semantics to operate on surface strings. Other theories require LFs that are *very* far-removed from what we see & hear. (see Kratzer 1986)

## 2 Technical Preliminaries

### 2.1 Possible Worlds Semantics

- Veltman begins by introducing the basic ideas of possible worlds semantics
- **Def.$[\![\cdot]\!]$** shows how to reinterpret the connectives of classical logic as set-theoretic operations on possible worlds
- When $\mathcal{A} = \{\mathsf{A}, \mathsf{B}, \mathsf{C}\}$, $W$ can be depicted as follows:

| $W$ | A | B | C |
|-----|---|---|---|
| $w_0$ | 0 | 0 | 0 |
| $w_1$ | 0 | 0 | 1 |
| $w_2$ | 0 | 1 | 0 |
| $w_3$ | 0 | 1 | 1 |
| $w_4$ | 1 | 0 | 0 |
| $w_5$ | 1 | 0 | 1 |
| $w_6$ | 1 | 1 | 0 |
| $w_7$ | 1 | 1 | 1 |

- $[\![\mathsf{A}]\!] = \{w_0, w_1, w_2, w_3\}$

- $[\![\neg\mathsf{A}]\!] = \{w_4, w_5, w_6, w_7\}$

- $[\![\mathsf{A} \wedge \mathsf{B}]\!] = \{w_6, w_7\}$

- $[\![(\mathsf{A} \wedge \mathsf{B}) \rightarrow \mathsf{C}]\!] = \{w_0, w_1, w_2, w_3, w_4, w_5, w_7\}$

- But notice also the presence of **situations**

  - These are incomplete worlds which only give truth values to certain formulae
  - They allow us to talk about commonalities among worlds & ignore irrelevant differences between worlds

### 2.2 Update Semantics

To capture the meaning of counterfactuals Veltman employs a **dynamic** notion of meaning, according to which the meaning of a sentence is the change it brings about in the cognitive state of anyone who wants to incorporate the information conveyed by it. He models this as follows.

- The meaning $[\phi]$ of a sentence $\phi$ is an operation on cognitive states
- Cognitive **states** are defined in **Def.$\mathcal{S}$** [3]
- Ordinary sentences will update $F_S$, while sentences stating laws will update both $F_S$ & $U_S$

  - After all, laws inform you both about what particular facts hold about the actual world, and what worlds are generally possible

- To represent laws formally, Veltman introduces an operator $\square$:

  - $\square\phi$ means *it is a law that $\phi$*

- **Def.$[\cdot]$** captures these two kinds of updates with two different clauses

---

[3] Veltman's division of information relevant to the interpretation of counterfactuals into laws and relevant facts resembles greatly the **cotenability** theory offered by Goodman (1947).

- These concepts allow the definition of another important notion in updates semantics: **support**
- The idea is that a state $S$ supports $\phi$ when updating with $\phi$ adds no information over & above what is already in $S$

- As for logical consequence, the idea is that some premises entail a conclusion if updating an arbitrary state with that sequence invariably leads to a state that supports the conclusion

### 2.2.1   Two Examples

Throughout the remainder of the presentation, we'll use the following translations:

$\triangleright$   *It is raining*                                      $\rightsquigarrow$   R

$\triangleright$   *Jones* $\begin{Bmatrix} \text{is wearing} \\ \text{wears} \end{Bmatrix}$ *his hat*   $\rightsquigarrow$   W

$\triangleright$   *The coin came up heads*         $\rightsquigarrow$   H

### Example 1

- Let's consider the state $S_2 = \mathbf{1}[\neg\mathsf{R}][\Box(\mathsf{H} \to \mathsf{W})]$, and how it comes about

| **1** | R | W | H |
|---|---|---|---|
| $\boldsymbol{w_0}$ | **0** | **0** | **0** |
| $\boldsymbol{w_1}$ | **0** | **0** | **1** |
| $\boldsymbol{w_2}$ | **0** | **1** | **0** |
| $\boldsymbol{w_3}$ | **0** | **1** | **1** |
| $\boldsymbol{w_4}$ | **1** | **0** | **0** |
| $\boldsymbol{w_5}$ | **1** | **0** | **1** |
| $\boldsymbol{w_6}$ | **1** | **1** | **0** |
| $\boldsymbol{w_7}$ | **1** | **1** | **1** |

$\xrightarrow{[\neg\mathsf{R}]}$

| $S_1$ | R | W | H |
|---|---|---|---|
| $\boldsymbol{w_0}$ | **0** | **0** | **0** |
| $\boldsymbol{w_1}$ | **0** | **0** | **1** |
| $\boldsymbol{w_2}$ | **0** | **1** | **0** |
| $\boldsymbol{w_3}$ | **0** | **1** | **1** |
| $w_4$ | 1 | 0 | 0 |
| $w_5$ | 1 | 0 | 1 |
| $w_6$ | 1 | 1 | 0 |
| $w_7$ | 1 | 1 | 1 |

$\xRightarrow{[\Box(\mathsf{H} \to \mathsf{W})]}$

| $S_1$ | R | W | H |
|---|---|---|---|
| $\boldsymbol{w_0}$ | **0** | **0** | **0** |
| ~~$\boldsymbol{w_1}$~~ | ~~**0**~~ | ~~**0**~~ | ~~**1**~~ |
| $\boldsymbol{w_2}$ | **0** | **1** | **0** |
| $\boldsymbol{w_3}$ | **0** | **1** | **1** |
| $w_4$ | 1 | 0 | 0 |
| ~~$w_5$~~ | ~~1~~ | ~~0~~ | ~~1~~ |
| $w_6$ | 1 | 1 | 0 |
| $w_7$ | 1 | 1 | 1 |

- ○ **Bold** worlds belong to $F_S$; all non-cancelled worlds are in $U_S$
    (NB: $F_S \subseteq U_S$)
- ○ We begin with $\mathbf{1} = \langle W, W \rangle$ and find $S_1 = \mathbf{1}[\neg\mathsf{R}]$:
    - $\triangleright$ To find $S_1$ intersect $F_\mathbf{1}$ with $[\![\neg\mathsf{R}]\!]$, i.e. throw the R-worlds out of $W$

- ○ Now find $S_1[\Box(\mathsf{H} \to \mathsf{W})]$:
    - $\triangleright$ Intersect $F_{S_1}$ **and** $U_{S_1}$ with $[\![(\mathsf{H} \to \mathsf{W})]\!]$, i.e. completely throw the $(\mathsf{H} \wedge \neg\mathsf{W})$-worlds out of the state (cross them out)

- Thinking graphically, $\Box$ can be thought of as a quantifier over rows of the table; it says: *cross out every row that falsifies the formula in my scope*

**Example 2**

- Consider the state $S_3 = \mathbf{1}[\Box(\mathsf{R} \to \mathsf{W})][\mathsf{R} \wedge \mathsf{W}]$, and how it comes about

| **1** | R | *W* | H |
|---|---|---|---|
| $\boldsymbol{w_0}$ | **0** | **0** | **0** |
| $\boldsymbol{w_1}$ | **0** | **0** | **1** |
| $\boldsymbol{w_2}$ | **0** | **1** | **0** |
| $\boldsymbol{w_3}$ | **0** | **1** | **1** |
| $\boldsymbol{w_4}$ | **1** | **0** | **0** |
| $\boldsymbol{w_5}$ | **1** | **0** | **1** |
| $\boldsymbol{w_6}$ | **1** | **1** | **0** |
| $\boldsymbol{w_7}$ | **1** | **1** | **1** |

$\xrightarrow{[\Box(\mathsf{R} \to \mathsf{W})]}$

| $S_1$ | R | W | H |
|---|---|---|---|
| $\boldsymbol{w_0}$ | **0** | **0** | **0** |
| $\boldsymbol{w_1}$ | **0** | **0** | **1** |
| $\boldsymbol{w_2}$ | **0** | **1** | **0** |
| $\boldsymbol{w_3}$ | **0** | **1** | **1** |
| ~~$\boldsymbol{w_4}$~~ | ~~1~~ | ~~0~~ | ~~0~~ |
| ~~$\boldsymbol{w_5}$~~ | ~~1~~ | ~~0~~ | ~~1~~ |
| $\boldsymbol{w_6}$ | **1** | **1** | **0** |
| $\boldsymbol{w_7}$ | **1** | **1** | **1** |

$\xrightarrow{[\mathsf{R} \wedge \mathsf{W}]}$

| $S_3$ | R | W | H |
|---|---|---|---|
| $w_0$ | 0 | 0 | 0 |
| $w_1$ | 0 | 0 | 1 |
| $w_2$ | 0 | 1 | 0 |
| $w_3$ | 0 | 1 | 1 |
| ~~$w_4$~~ | ~~1~~ | ~~0~~ | ~~0~~ |
| ~~$w_5$~~ | ~~1~~ | ~~0~~ | ~~1~~ |
| $w_6$ | 1 | 1 | 0 |
| $w_7$ | 1 | 1 | 1 |

# 3  Interpreting Counterfactuals Dynamically

Now it's time to look at Veltman's approach to interpreting counterfactuals

- Examples like (1) suggest that Ramsey's Test should be replaced with the following:

  **Veltman's Test**  To assess a counterfactual *if it had been $\phi$, it would have been $\psi$* carry out this process:
  - ① Remove $\neg\phi$ & anything that depends on it from your belief state
  - ② Update your beliefs hypothetically with $\phi$
  - ③ Test wether or not your beliefs now support $\psi$

- Veltman implements this test formally as follows:

  - Step ① is achieved by **retracting** $[\![\neg\phi]\!]$ from $S$, i.e. forming a state $S \downarrow [\![\neg\phi]\!]$
  - Step ② is achieved by updating $S \downarrow [\![\neg\phi]\!]$ with $\phi$: $(S \downarrow [\![\neg\phi]\!])[\phi]$
  - Step ③ is achieved by testing whether or not $(S \downarrow [\![\neg\phi]\!])[\phi] \models \psi$

- Carrying out steps ① & ② is what Veltman calls **making a counterfactual assumption**, and corresponds to the interpretation of the antecedent (**Def.4.3**), while step ③ corresponds to interpreting the consequent in the state derived in the first two steps

- The technical heart of the theory resides in correctly defining retraction

- So, before looking at the formal definition, let's consider what it must **do** to get examples (1) & (2) right

## 3.1  How to Get (1) & (2) Right

### 3.1.1  Example (1)

- Our goal is to predict that (1d) is not good in (1):

  (1) a. Invariably, if it is raining, Jones wears his hat
      b. If it is not raining, Jones wears his hat at random
      c. Today, it is raining and so Jones is wearing his hat
      d. But, even if it had not been raining, Jones would have been wearing his hat

- Given our translations in §2.2.1, (1) $\rightsquigarrow$ (1′):

  (1′) a. $\Box(\mathsf{R} \rightarrow \mathsf{W})$
       b. (Not included in Veltman's analysis)
       c. $\mathsf{R} \wedge \mathsf{W}$
       d. *If it had been* $\neg\mathsf{R}$, *it would have been* $\mathsf{W}$

- (1d′) is interpreted in the state $S_3 = \mathbf{1}[\Box(\mathsf{R} \rightarrow \mathsf{W})][\mathsf{R} \wedge \mathsf{W}]$, which we calculated earlier (§2.2.1):

| $S_3$ | R | W | H |
|-------|---|---|---|
| $w_0$ | 0 | 0 | 0 |
| $w_1$ | 0 | 0 | 1 |
| $w_2$ | 0 | 1 | 0 |
| $w_3$ | 0 | 1 | 1 |
| ~~$w_4$~~ | ~~0~~ | ~~0~~ | ~~0~~ |
| ~~$w_5$~~ | ~~0~~ | ~~0~~ | ~~0~~ |
| $\mathbf{w_6}$ | **1** | **1** | **0** |
| $\mathbf{w_7}$ | **1** | **1** | **1** |

- Recall that the plan is to assess a counterfactual by testing whether or not $(S \downarrow \llbracket \neg\phi \rrbracket)[\phi] \models \psi$, but we want to be clear what operation $\downarrow$ must be in order for this definition to cohere with the data

- So, we're looking for a definition of $\downarrow$ that allows us to show:

  **Fact 1**  $(S_3 \downarrow \llbracket \neg\neg\mathsf{R} \rrbracket)[\neg\mathsf{R}] \not\models \mathsf{W}$

- $S_4 := (S_3 \downarrow \llbracket \neg\neg\mathsf{R} \rrbracket)[\neg\mathsf{R}]$
- From Fact 1 we derive two constraints $F_{S_4}$:

 (i)  $w_6, w_7 \notin F_{S_4}$ (since $S_4$ arises from an update with $\neg\mathsf{R}$)
(ii)  $w_0 \in F_{S_4}$ or $w_1 \in F_{S_4}$, or both (since $S_4 \not\models \mathsf{W}$)

- The largest set consistent with these facts is shown below:

| $S_4$ | R | W | H |
|---|---|---|---|
| $\boldsymbol{w_0}$ | **0** | **0** | **0** |
| $\boldsymbol{w_1}$ | **0** | **0** | **1** |
| $\boldsymbol{w_2}$ | **0** | **1** | **0** |
| $\boldsymbol{w_3}$ | **0** | **1** | **1** |
| ~~$w_4$~~ | ~~1~~ | ~~0~~ | ~~0~~ |
| ~~$w_5$~~ | ~~1~~ | ~~0~~ | ~~1~~ |
| $w_6$ | 1 | 1 | 0 |
| $w_7$ | 1 | 1 | 1 |

- Now, we'll see what constraints example (2) places on ↓
- We can then see which definitions meet both sets of constraints, making more clear how & Veltman arrived at the definition that he did

### 3.1.2   Example (2)

- Our goal is to predict that (2f) is good in (2):
  (2) a. Jones always flips a coin before he opens the curtain to see what the weather is like
      b. If it's not raining and the coin comes up heads, he wears his hat
      c. If it's not raining and the coin comes up tails, he doesn't wear his hat
      d. Invariably, if it is raining he wears his hat
      e. Today, the coin came up heads and it is raining, so Jones is wearing his hat
      f. But, even if it hadn't been raining, Jones would have been wearing his hat

- Given our translations in §2.2.1, (2) $\rightsquigarrow$ (2′):

  (2′) a. (Not included in analysis)
       b. $\Box((R \vee H) \leftrightarrow W)$
       c. (Incorporated into (2b′))
       d. (Incorporated into (2b′))
       e. $(H \wedge R) \wedge W$
       f. *If it had been* ¬R, *it would have been* W

- So we want to a definition of ↓ that allows us to show the following, where $S_5 := \mathbf{1}[\Box((R \vee H) \leftrightarrow W)][H \wedge R]$:

  **Fact 2** $(S_5 \downarrow [\![\neg\neg R]\!])[\neg R] \models W$

- To see what this amounts to, we first find $S_5$:

| **1** | R | *W* | H |
|---|---|---|---|
| $w_0$ | 0 | 0 | 0 |
| $w_1$ | 0 | 0 | 1 |
| $w_2$ | 0 | 1 | 0 |
| $w_3$ | 0 | 1 | 1 |
| $w_4$ | 1 | 0 | 0 |
| $w_5$ | 1 | 0 | 1 |
| $w_6$ | 1 | 1 | 0 |
| $w_7$ | 1 | 1 | 1 |

$\xrightarrow{\Box((R \vee H) \leftrightarrow W)}$

| | R | W | H |
|---|---|---|---|
| $w_0$ | 0 | 0 | 0 |
| ~~$w_1$~~ | ~~0~~ | ~~0~~ | ~~1~~ |
| ~~$w_2$~~ | ~~0~~ | ~~1~~ | ~~0~~ |
| $w_3$ | 0 | 1 | 1 |
| ~~$w_4$~~ | ~~1~~ | ~~0~~ | ~~0~~ |
| ~~$w_5$~~ | ~~1~~ | ~~0~~ | ~~1~~ |
| $w_6$ | 1 | 1 | 0 |
| $w_7$ | 1 | 1 | 1 |

$\xrightarrow{[H \wedge R]}$

| $S_5$ | R | W | H |
|---|---|---|---|
| $w_0$ | 0 | 0 | 0 |
| ~~$w_1$~~ | ~~0~~ | ~~0~~ | ~~1~~ |
| ~~$w_2$~~ | ~~0~~ | ~~1~~ | ~~0~~ |
| $w_3$ | 0 | 1 | 1 |
| ~~$w_4$~~ | ~~1~~ | ~~0~~ | ~~0~~ |
| ~~$w_5$~~ | ~~1~~ | ~~0~~ | ~~1~~ |
| $w_6$ | 1 | 1 | 0 |
| $w_7$ | 1 | 1 | 1 |

- What follows about $S_6 := (S_5 \downarrow [\![\neg\neg R]\!])[\neg R]$ from Fact 2?

(i) $w_6, w_7 \notin F_{S_6}$

(ii) $w_0 \notin F_{S_6}$

- This narrows the possibilities down to one:

| $S_6$ | R | W | H |
|---|---|---|---|
| $w_0$ | 0 | 0 | 0 |
| ~~$w_1$~~ | ~~0~~ | ~~0~~ | ~~1~~ |
| ~~$w_2$~~ | ~~0~~ | ~~1~~ | ~~0~~ |
| $w_3$ | 0 | 1 | 1 |
| ~~$w_4$~~ | ~~1~~ | ~~0~~ | ~~0~~ |
| ~~$w_5$~~ | ~~1~~ | ~~0~~ | ~~1~~ |
| $w_6$ | 1 | 1 | 0 |
| $w_7$ | 1 | 1 | 1 |

- We can now say precisely what needs to happen to get examples (1) & (2) right

### 3.1.3 Putting the Pieces Together

- In the previous two sections we determined what a satisfactory definition of retraction must do (in graphical format):

**Condition 1**

| $S_3$ | R | W | H |
|---|---|---|---|
| $w_0$ | 0 | 0 | 0 |
| $w_1$ | 0 | 0 | 1 |
| $w_2$ | 0 | 1 | 0 |
| $w_3$ | 0 | 1 | 1 |
| ~~$w_4$~~ | ~~X~~ | ~~X~~ | ~~X~~ |
| ~~$w_5$~~ | ~~X~~ | ~~X~~ | ~~X~~ |
| $\boldsymbol{w_6}$ | 1 | 1 | 0 |
| $\boldsymbol{w_7}$ | 1 | 1 | 1 |

$\xrightarrow{(S_3\downarrow[\![\neg\neg\mathsf{R}]\!])[\neg\mathsf{R}]}$

| $S_4$ | R | W | H |
|---|---|---|---|
| $\boldsymbol{w_0}$ | **0** | **0** | **0** |
| $\boldsymbol{w_1}$ | **0** | **0** | **1** |
| $\boldsymbol{w_2}$ | **0** | **1** | **0** |
| $\boldsymbol{w_3}$ | **0** | **1** | **1** |
| ~~$w_4$~~ | ~~X~~ | ~~X~~ | ~~X~~ |
| ~~$w_5$~~ | ~~X~~ | ~~X~~ | ~~X~~ |
| $w_6$ | 1 | 1 | 0 |
| $w_7$ | 1 | 1 | 1 |

**Condition 2**

| $S_5$ | R | W | H |
|---|---|---|---|
| $w_0$ | 0 | 0 | 0 |
| ~~$w_1$~~ | ~~X~~ | ~~X~~ | ~~X~~ |
| ~~$w_2$~~ | ~~X~~ | ~~X~~ | ~~X~~ |
| $w_3$ | 0 | 1 | 1 |
| ~~$w_4$~~ | ~~X~~ | ~~X~~ | ~~X~~ |
| ~~$w_5$~~ | ~~X~~ | ~~X~~ | ~~X~~ |
| $w_6$ | 1 | 1 | 0 |
| $\boldsymbol{w_7}$ | **1** | **1** | **1** |

$\xrightarrow{(S_5\downarrow[\![\neg\neg\mathsf{R}]\!])[\neg\mathsf{R}]}$

| $S_6$ | R | W | H |
|---|---|---|---|
| $w_0$ | 0 | 0 | 0 |
| ~~$w_1$~~ | ~~X~~ | ~~X~~ | ~~X~~ |
| ~~$w_2$~~ | ~~X~~ | ~~X~~ | ~~X~~ |
| $\boldsymbol{w_3}$ | **0** | **1** | **1** |
| ~~$w_4$~~ | ~~X~~ | ~~X~~ | ~~X~~ |
| ~~$w_5$~~ | ~~X~~ | ~~X~~ | ~~X~~ |
| $w_6$ | 1 | 1 | 0 |
| $w_7$ | 1 | 1 | 1 |

- While Condition 1 suggests that $S_3\downarrow[\![\neg\neg\mathsf{R}]\!]$ is the set of all $\neg\mathsf{R}$-worlds in $F_{S_3}$ this strategy does not work for Condition 2, since $w_0$ can't be in $F_{S_6}$
- Veltman (2005: 168) proposes the following solution:

  - To retract $[\![\neg\phi]\!]$ from $S$ the following must be done for every $w\in F_S$ and every **basis** $s'$ for $w$:
    - ① If $s'$ **forces** $[\![\neg\phi]\!]$, make minimal adjustments to $s'$ s.t. it doesn't
    - ② Call each result of step ① $s$
    - ③ Each world in $U_S$ extending such an $s$ belongs to $F_{S\downarrow[\![\neg\phi]\!]}$

  - $S\downarrow P$ is defined more formally in **Definition 4**
  - *Basis* & *Forcing* defined formally in **Definition 3**
    - ▷ Intuitively, $s$ is a basis for $w$ if $s$ is among the minimal situations that can be used to uniquely identify $w$ in $U_S$
    - ▷ Intuitively, $s$ forces $P$ within $U_S$ just in case every world $w$ in $U_S$ that $s$ is in is also in $P$

- Let's see that Veltman's proposal meets Condition 1:
  - ○ There are two worlds in $F_{S_3}$: $w_6$ & $w_7$

    **Case 1** ($w_6$)
    - ① One basis: $\{\langle \mathsf{R}, 1\rangle, \langle \mathsf{H}, 0\rangle\}$; maximal subset that doesn't force $[\![\neg\neg\mathsf{R}]\!]$: $\{\langle \mathsf{H}, 0\rangle\}$
    - ③ $w_0, w_2$ & $w_6$ extend $\langle \mathsf{H}, 0\rangle$, so $\{w_0, w_2, w_6\} \subseteq F_{S_4}$

    **Case 2** ($w_7$)
    - ① One basis: $\{\langle \mathsf{R}, 1\rangle, \langle \mathsf{H}, 1\rangle\}$; maximal subset that doesn't force $[\![\neg\neg\mathsf{R}]\!]$: $\{\langle \mathsf{H}, 1\rangle\}$
    - ③ $w_1, w_3$ & $w_7$ extend $\langle \mathsf{H}, 1\rangle$, so now $\{w_0, w_1, w_2, w_3, w_6, w_7\} \subseteq F_{S_4}$
  - ○ So $S_3 \downarrow [\![\neg\neg\mathsf{R}]\!] = \langle U_S, \{w_0, w_1, w_2, w_3, w_6, w_7\}\rangle$
  - ○ Thus, $(S_3 \downarrow [\![\neg\neg\mathsf{R}]\!])[\neg\mathsf{R}] = \langle U_S, \{w_0, w_1, w_2, w_3\}\rangle = S_4$
  - ○ Perfect!

- Let's see that Veltman's proposal meets Condition 2:
  - ○ There is one world in $F_{S_5}$: $w_7$
    - ① One basis: $\{\langle \mathsf{R}, 1\rangle, \langle \mathsf{H}, 1\rangle\}$; maximal subset that doesn't force $[\![\neg\neg\mathsf{R}]\!]$: $\{\langle \mathsf{H}, 1\rangle\}$
    - ③ $w_3$ & $w_7$ extend $\langle \mathsf{H}, 1\rangle$, so $\{w_3, w_7\} = F_{S_5}$
  - ○ So $S_5 \downarrow [\![\neg\neg\mathsf{R}]\!] = \langle U_S, \{w_3, w_7\}\rangle$
  - ○ Thus, $(S_5 \downarrow [\![\neg\neg\mathsf{R}]\!])[\neg\mathsf{R}] = \langle U_S, \{w_3\}\rangle = S_6$
  - ○ $\sqrt{\phantom{--}}$ !

- To see how the formal definitions mirror this intuitive version of the theory we'll work through example (1) in full detail

- I've written a program in PLT Scheme that grinds out predictions for Veltman's theory; I've used it to verify that his predictions for (2) are indeed correct, although working through a full derivation would be a good exercise for those wishing to understand the theory better

### 3.1.4 Example (1) in Detail

- Recall that $(1) \rightsquigarrow (1')$:

  $(1')$a. $\Box(\mathsf{R} \to \mathsf{W})$
  
      b. (Not included in Veltman's analysis)
  
      c. $\mathsf{R} \wedge \mathsf{W}$
  
      d. *If it had been* $\neg\mathsf{R}$, *it would have been* $\mathsf{W}$

- We want to show:

  **Fact 1**
  $\mathbf{1}[\Box(\mathsf{R} \to \mathsf{W})][\mathsf{R} \wedge \mathsf{W}][\textit{if it had been } \neg\mathsf{R}] \not\models \mathsf{W}$

- For this example we only really need two atomic sentences, so let's recalculate $\mathbf{1}[\Box(\mathsf{R} \to \mathsf{W})][\mathsf{R} \wedge \mathsf{W}]$:

| **1** | R | W |
|---|---|---|
| $\boldsymbol{w_0}$ | **0** | **0** |
| $\boldsymbol{w_1}$ | **0** | **1** |
| $\boldsymbol{w_2}$ | **1** | **0** |
| $\boldsymbol{w_3}$ | **1** | **1** |

$\xrightarrow{[\Box(\mathsf{R} \to \mathsf{W})]}$

| $S_1$ | R | W |
|---|---|---|
| $\boldsymbol{w_0}$ | **0** | **0** |
| $\boldsymbol{w_1}$ | **0** | **1** |
| ~~$\boldsymbol{w_2}$~~ | ~~**1**~~ | ~~**0**~~ |
| $\boldsymbol{w_3}$ | **1** | **1** |

$\xrightarrow{[\mathsf{R} \wedge \mathsf{W}]}$

| $S_2$ | R | W |
|---|---|---|
| $w_0$ | 0 | 0 |
| $w_1$ | 0 | 1 |
| ~~$w_2$~~ | ~~1~~ | ~~0~~ |
| $w_3$ | 1 | 1 |

- Find $S_3[\textit{if it had been } \neg\mathsf{R}]$:

$$
\begin{aligned}
S_3[\textit{if it had been } \neg\mathsf{R}] &= S_3 \downarrow [\![\neg\neg\mathsf{R}]\!][\neg\mathsf{R}] && (\textbf{Def 4.3}) \\
&= \langle U_{S_3 \downarrow [\![\neg\neg\mathsf{R}]\!]}, F_{S_3 \downarrow [\![\neg\neg\mathsf{R}]\!]} \rangle[\neg\mathsf{R}] && (\textbf{Def 4.2}) \\
&= \langle U_{S_3}, F_{S_3 \downarrow [\![\neg\neg\mathsf{R}]\!]} \rangle[\neg\mathsf{R}] && (\textbf{Def 4.2.a}) \quad (5)
\end{aligned}
$$

  ○ Find $F_{S_3 \downarrow [\![\neg\neg\mathsf{R}]\!]}$:

$$
\begin{aligned}
F_{S_3 \downarrow [\![\neg\neg\mathsf{R}]\!]} &= \{w \mid w \in U_{S_3} \ \& \ \exists w' \in F_{S_3}, \exists s \in w' \downarrow [\![\neg\neg\mathsf{R}]\!] : s \subseteq w\} && (\textbf{Def 4.2.b}) \\
&= \{w \mid w \in U_{S_3} \ \& \ \exists w' \in F_{S_3}, \exists s \in w' \downarrow \{w_2, w_3\} : s \subseteq w\} && (\textbf{Def } [\![\neg\phi]\!]) \\
&= \{w \mid w \in U_{S_3} \ \& \ \exists s \in w_3 \downarrow \{w_2, w_3\} : s \subseteq w\} && (6) \\
& && (F_{S_3} = \{w_3\})
\end{aligned}
$$

  ▷ Find $w_3 \downarrow \{w_2, w_3\}$:

$$
\begin{aligned}
w_3 \downarrow \{w_2, w_3\} = \{s \mid s \subseteq w_3 \ \& \ \exists s' : s' \text{ is a basis for } w_3 \ \& && (\textbf{Def 4.1}) \\
s \text{ is a maximal subset of } s' \\
\text{not forcing } \{w_2, w_3\} \quad \}
\end{aligned}
$$

▶ Since $s \subseteq w_3$ the only four situations that could be in $w_3 \downarrow \{w_2, w_3\}$ are $s_1 = \{\langle \mathsf{R}, 1\rangle\}$, $s_2 = \{\langle \mathsf{R}, 1\rangle, \langle \mathsf{W}, 1\rangle\}$, $s_3 = \{\langle \mathsf{W}, 1\rangle\}$ & $\varnothing$.

⋆ $s_1 \in w_3 \downarrow \{w_2, w_3\}$? No; $s_1$ forces $\{w_2, w_3\}$, since the only world $w \in U_{S_3}$ s.t. $s_1 \subseteq w$ is $w_3$ and $w_3 \in \{w_2, w_3\}$.

⋆ $s_2 \in w_3 \downarrow \{w_2, w_3\}$? No, for the same reason as the last case.

⋆ $s_3 \in w_3 \downarrow \{w_2, w_3\}$? $s_3$ does not force $\{w_2, w_3\}$. $\exists s'$: $s'$ is a basis for $w_3$ & $s_3 \subseteq s'$? No. $w_3$ only has one basis, $\{\langle \mathsf{R}, 1\rangle\}$, and $s_3 \not\subseteq \{\langle \mathsf{R}, 1\rangle\}$

⋆ $\varnothing \in w_3 \downarrow \{w_2, w_3\}$? Yes. $\varnothing$ does not force $\{w_2, w_3\}$, since $\exists w \in U_{S_3}$ s.t. $\varnothing \subset w$ while $\varnothing \notin \{w_2, w_3\}$; namely $w_1$ & $w_2$. Also, $\{\langle \mathsf{R}, 1\rangle\}$ is the unique basis for $w_3$, so we find the set of all subsets of $\{\langle \mathsf{R}, 1\rangle\}$ that do not force $\{w_2, w_3\}$. $\{\langle \mathsf{R}, 1\rangle\}$ forces $\{w_2, w_3\}$, so the only such subset is $\varnothing$, hence it is the maximal such subset.

Therefore $w_3 \downarrow \{w_2, w_3\} = \{\varnothing\}$

Substituting this result back into (6) we get the following:

$$
\begin{aligned}
F_{S_3 \downarrow [\![ \neg\neg\mathsf{R} ]\!]} &= \{w \mid w \in U_{S_3} \ \& \ \exists s \in w_3 \downarrow \{w_2, w_3\} : s \subseteq w\} \\
&= \{w \mid w \in U_{S_3} \ \& \ \exists s \in \{\varnothing\} : s \subseteq w\} \\
&= \{w \mid w \in U_{S_3} \ \& \ \varnothing \subseteq w\} \\
&= U_{S_3} \quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad (7)
\end{aligned}
$$

Substituting this result back into (5) gives us:

$$
\begin{aligned}
S_3[\textit{if it had been } \neg\mathsf{R}] &= \langle U_{S_3}, F_{S_3 \downarrow [\![ \neg\neg\mathsf{R} ]\!]} \rangle [\neg\mathsf{R}] \\
&= \langle U_{S_3}, U_{S_3} \rangle [\neg\mathsf{R}] \\
&= \langle U_{S_3}, U_{S_3} \cap \{w_0, w_1\} \rangle \quad (\mathbf{Def.} S[\phi], \mathbf{Def.} [\![ \neg\phi ]\!]) \\
&= \langle U_{S_3}, \{w_0, w_1\} \rangle \quad\quad\quad\quad\quad\quad\quad\quad\quad (8)
\end{aligned}
$$

○ This is exactly what we want, since it allows us to show that $S_3[\textit{if it had been } \neg\mathsf{R}] \not\models \mathsf{W}$:

$$
\begin{aligned}
\langle U_{S_3}, \{w_0, w_1\} \rangle [\mathsf{W}] &= \langle U_{S_3}, \{w_0, w_1\} \cap \{w_1, w_3\} \rangle \\
&= \langle U_{S_3}, \{w_1\} \rangle \\
&\neq \langle U_{S_3}, \{w_0, w_1\} \rangle \quad\quad\quad\quad\quad (9)
\end{aligned}
$$

By $\mathbf{Def.}{\models}$, (9) implies that $\langle U_{S_3}, \{w_0, w_1\} \rangle \not\models \mathsf{W}$. From (8) we know that this means $S_3[\textit{if it had been } \neg\mathsf{R}] \not\models \mathsf{W}$. Unpacking $S_3$ we get:
**Fact 1  1$[\Box(\mathsf{R} \to \mathsf{W})][\mathsf{R} \wedge \mathsf{W}][\textit{if it had been } \neg\mathsf{R}] \not\models \mathsf{W}$**

# 4 Propositions, Support & Truth

- In the previous sections all we assumed about the interpretation of counterfactuals was:

  **Counterfactual Support Constraint**
    $S \models$ *if it had been* $\phi$, *it would have been* $\psi \iff S[\textit{if it had been } \phi] \models \psi$

- But this can't be all there is to the story
- If we allowed discourse *after* a counterfactual to be interpreted in the state $S[\textit{if it had been } \phi]$ we'd make the prediction that discourses like the following are good:

  (10)
      a. If Oswald hadn't killed Kennedy, someone else would have
      b. In fact, Oswald didn't kill Kennedy

- It will always be the case that $S[\textit{if it had been } \phi] \models \phi$, but this is a problem since you can't follow up a counterfactual by asserting its antecedent
- Similarly you can't generally follow up a counterfactual by asserting its consequent, but if we interpret discourse coming after a counterfactual in the state $S[\textit{if it had been } \phi]$ then we predict the opposite
- Accordingly, it looks as if only the consequent is interpreted in $S[\textit{if it had been } \phi]$, but subsequent discourse is interpreted in the original input state $S$
- This idea is captured formally in **Definition 5**, where counterfactuals are treated as **tests** [4]
- This implementation gives rise to some pressing questions:
  - Is anaphora out of counterfactuals possible?
  - Don't counterfactuals express propositions?
    - ▷ If so, which ones?
    - ▷ If not, how can counterfactuals be informative?
- Later in the course we'll consider the first question more; for now, let's focus on the second one

---

[4] Note that the notion of test here is different than in DPL. In DPL 'tests' are more like 'filters': they can eliminate assignments from the info state, but don't add any. Tests here, on the other hand, neither add nor subtract information; passing the test means proceeding with the exact same state that was the input and failing means transitioning to the absurd state. This difference is relevant here, since it is obvious how tests in the DPL sense can be informative.

## 4.1 How are Counterfactual Tests Informative?

- According to Definition 5, interpreting a counterfactual simply amounts to performing a test on your state of information
- When combined with Veltman's approach to interpretation, the resulting picture faces a puzzle:
  - Information states are conceived of as **individualistic** representations; they are psychological entities[5]
  - So, you interpret my utterances by evaluating them in your information state
  - But wait, if you don't already accept the counterfactual I am attempting to communicate, then aren't you going to just think it is false, since your info state will fail the test counterfactuals trigger?

---

[5] Of course there will be mutual public constraints on information states; such constraints seem to be largely governed by pragmatic reasoning. (see Thomason *et al.* to appear)

- Veltman's solution to this puzzle is two-fold:

(i) First, he maintains that when the laws are fixed, as in chess or classical mechanics, it is possible to give normal update & propositional clauses that are correct; namely those in **Definition 6**

(ii) When the laws are not fixed, we rely on pragmatic strategies to extract information from discourse involving counterfactuals:

> Given [Definition 5], sentences of the form ⌜*if had been $\phi$, would have been $\psi$*⌝ do not convey new information — not directly at least. They provide an invitation to perform a test. By asserting ⌜*if had been $\phi$, would have been $\psi$*⌝, a speaker makes a kind of comment: 'Given the general laws and the facts I am acquainted with, the sentence $\psi$ is supported by the state I get in when I assume that $\phi$ had been the case'. The addressee is supposed to determine whether the same holds on account of his or her own information. If not, a discussion will arise, and in the course of this discussion both the speaker and the hearer may learn some new laws and facts, which could affect the outcome of the test. (Veltman 2005: 171)

> ▷ Note that if your state does not pass the test it is indeterminate how you should change your state to pass it: you can either assume new laws or some new facts. This suggests that one must move to a pragmatic story for this kind of phenomena once info states get split into law and fact bases. It is only through the course of interrogation that an addressee can determine how to coordinate his info state with the speaker.

- Why doesn't Definition 6 work when the laws are not fixed?
  - Counterfactuals often communicate information about laws, but this is not captured in Definition 6; neither in the propositions they express (**Def 6.3**) nor in the updates they trigger (**Def 6.4**)
- Is their some intuitive justification for having things turn out like this, or is it just an unfortunate consequence of splitting info states into information pertaining to laws and information pertaining to facts?

## 5  Beyond Veltman (2005)

**The Duchess**
  The duchess has been murdered, and you are supposed to find the murderer.
  At some point only the butler and the gardener are left as suspects. At this
  point you believe:

  (11) *If the butler did not kill her, the gardener did*

  Still, somewhat later — after you found convincing evidence showing that
  the butler did it, and that the garderer had nothing to do with it — you
  get in a state in which you will *reject* the sentence:

  (12) *If the butler had not killed her, the gardener would have [killed her]*

- As Veltman (2005: 174-5) shows, his theory predicts that (12) is to be re-
  jected
- But, several authors have argued that (12) has a reading on which it is *true*
  - This is the so-called **epistemic reading**; the idea is that (12) communi-
    cates the idea that there was a stage in the evolution of your epistemic
    state where you would have concluded from the information that the but-
    ler didn't do it — even though you've found that he did — that the
    gardener did it
- I don't get any such reading, particularly when the elliptical *killed her* is
  filled in (though I'm not sure why)
- Veltman claims that there is no epistemic reading; Schulz disagrees. What
  to do?
  ▷ Let's look at more data


**The Hamburger** (Hansson 1989)
  Suppose that one night you approach a small town of which you know that
  it has exactly two snackbars. Just before entering the town you meet a
  man eating a hamburger. You have good reason to accept the following
  conditional:

  (13) *If snackbar A is closed, snackbar B is open*

  Suppose now that after entering the town & seeing the man with the burger,
  you see that A is in fact open. Can you then think to yourself:

  (14) *[Even] if snackbar A had been closed, snackbar B would have been
        open*

- This example seems much better to me, but still borderline.

## 5.1   Another Kind of Case

**King Ludwig** (Kratzer 1989: 640)

King Ludwig of Bavaria likes to spend his weekends at Leoni Castle. Whenever the Royal Bavarian flag is up and the lights are on, the King is in the castle. At the moment the lights are on, the flag is down, and the King is away. But:

(15) *If the flag had been up, then the King would have been in the castle*

- It seems as if (15) is true, but, as Veltman (2005: 177-6) shows, his theory predicts it to be unacceptable
- Although this is a problem for the theory, there are examples with the same abstract structure that require the prediction that Veltman's theory makes:

**Ann, Billie & Carol**

Consider the case of three sisters who own just one bed, large enough for two of them but too small for all three. Every night at least one of them has to sleep on the floor. Whenever Ann sleeps in the bed and Billie sleeps in the bed, Carol sleeps on the floor. At the moment Billie is sleeping in bed, Ann is sleeping on the floor, and Carol is sleeping in bed.

(16) *If Ann had been in bed, Carol would have been sleeping on the floor*

- (16) seems to be false; why wouldn't Billie be on the floor?
- In parallel we might ask why the lights wouldn't have been off and the King still away.

# 6 Formal Definitions

**Worlds** $W = \{w \mid w : \mathcal{A} \mapsto \{0, 1\}\}$
(A world $w$ is a function that gives every atomic formula a truth-value)

**Situations** $N = \{s \mid s \subseteq w \;\&\; w \in W\}$
(A situation $s$ is a partial world; it assigns only some atomic formulae truth-values)

**Proposition** A proposition is a set of worlds, intuitively the set of worlds where a sentence is true.

**Definition of $[\![\cdot]\!]$** (Propositions)

$$
\begin{aligned}
[\![\mathsf{A}]\!] &= \{w \in W \mid w(\mathsf{A}) = 1\}, \text{where } \mathsf{A} \in \mathcal{A} \\
[\![\neg\phi]\!] &= W \sim [\![\phi]\!] \\
[\![\phi \wedge \psi]\!] &= [\![\phi]\!] \cap [\![\psi]\!] \\
[\![\phi \vee \psi]\!] &= [\![\phi]\!] \cup [\![\psi]\!] \\
[\![\phi \to \psi]\!] &= (W \sim [\![\phi]\!]) \cup [\![\psi]\!]
\end{aligned}
$$

**Definition of $S$** (States)
- $S = \langle U_S, F_S \rangle$, where either of the following conditions holds:
  (1) $\varnothing \neq F_S \subseteq U_S \subseteq W$
  (2) $F_S = U_S = \varnothing$
- $w \in F_S$ just in case, for all the agent knows, $w$ might be the actual world (‘$F$’ for *fact*)
- $U_S$ is the set of worlds consistent with the general laws the agent knows (‘$U$’ for *universe*)
- $\mathbf{1} := \langle W, W \rangle$; $\mathbf{1}$ is the **minimal state**
- $\mathbf{0} := \langle \varnothing, \varnothing \rangle$; $\mathbf{0}$ is the **absurd state**

**Definition of $[\,\cdot\,]$** (Basic Interpretation)

$$
S[\phi] = \begin{cases} \langle U_S, F_S \cap [\![\phi]\!] \rangle & \text{if } F_S \cap [\![\phi]\!] \neq \varnothing \\ \mathbf{0} & \text{otherwise} \end{cases}
$$

$$
S[\Box\phi] = \begin{cases} \langle U_S \cap [\![\phi]\!], F_S \cap [\![\phi]\!] \rangle & \text{if } F_S \cap [\![\phi]\!] \neq \varnothing \\ \mathbf{0} & \text{otherwise} \end{cases}
$$

**Definition of $\models$** (Support, Logical Consequence)

$$
\begin{aligned}
S \models \phi &\iff S[\phi] = S &&(S \text{ \textbf{supports} } \phi) \\
\phi_1, \ldots, \phi_n \models \psi &\iff \forall S : S[\phi_1] \ldots [\phi_n] \models \psi &&(\phi_1, \ldots, \phi_n \text{ \textbf{logically entail} } \psi)
\end{aligned}
$$

**Definition 3** (Forcing, Determination, Basis)

(a) The situation $s$ **forces the proposition** $P$ **within** $U_S$ iff for every $w \in U_S$ such that $s \subseteq w$ it holds that $w \in P$.

(b) The situation $s$ **determines the world** $w$ iff $s$ forces $\{w\}$ within $U_S$.

(c) The situation $s$ is a **basis for the world** $w$ iff $s$ is a **minimal** situation determining $w$ within $U_S$.

**Definition 4** (Retraction, Counterfactual Assumption)

(1) Suppose $w \in U_S$, and $P \subseteq W$. Then:

$$w \downarrow P = \{s \mid s \subseteq w \ \& \ \exists s' : s' \text{ is a basis for } w \ \& \ s \text{ is a maximal subset of } s' \text{ not forcing } P\}$$

(2) $S \downarrow P = \langle U_{S \downarrow P}, F_{S \downarrow P} \rangle$ (the **retraction of** $P$ **from** $S$), where:

(a) $U_{S \downarrow P} = U_S$

(b) $F_{S \downarrow P} = \{w \mid w \in U_S \ \& \ \exists w' \in F_S, \exists s \in w' \downarrow P : s \subseteq w\}$

(3) $S[\textit{if it had been the case that } \phi] = (S \downarrow [\![\neg \phi]\!])[\phi]$

**Intuitive Gloss of Retraction**

○ To retract $[\![\neg \phi]\!]$ from $S$ the following must be done for every $w \in F_S$ and every **basis** $s'$ for $w$:

① If $s'$ **forces** $[\![\neg \phi]\!]$, make minimal adjustments to $s'$ s.t. it doesn't

② Call each result of step ① $s$

③ Each world in $U_S$ extending such an $s$ belongs to $F_{S \downarrow [\![\neg \phi]\!]}$

**Definition 5** (Counterfactuals as Tests)

$$S[\textit{if it had been } \phi, \textit{ it would have been } \psi] = \begin{cases} S & \text{if } S[\textit{if it had been } \phi] \models \psi \\ \mathbf{0} & \text{otherwise} \end{cases}$$

**Definition 6** (Counterfactual Updates & Propositions)

When $U_S$ is fixed, normal update & proposition clauses can be given for counterfactuals

(1) $\exists v \in F_S : \langle U_S, \{v\}\rangle[\textit{if it had been } \phi] \models \psi \implies$

$S[\textit{if had been } \phi, \textit{would have been } \psi] = \langle U_S, \{v \in F_S \mid \langle U, \{v\}\rangle[\textit{if had been } \phi] \models \psi\}\rangle$

(2) Otherwise, $S[\textit{if had been } \phi, \textit{would have been } \psi] = \mathbf{0}$

(3) $[\![\textit{if had been } \phi, \textit{would have been } \psi]\!] = \{w \in W \mid \langle U, \{w\}\rangle[\textit{if had been } \phi] \models \psi\}$

(4) $S[\textit{if had been } \phi, \textit{would have been } \psi] = \langle U_S, F_S \cap [\![\textit{if had been } \phi, \textit{would have been } \psi]\!]\rangle$

# References

GOODMAN, N. (1947). 'The Problem of Counterfactual Conditionals'. *The Journal of Philosophy*, **44**: 113–118.

HANSSON, S. O. (1989). 'New Operators for Theory Change'. *Theoria*, **55**: 114–132.

KANAZAWA, M., KAUFMANN, S. & PETERS, S. (2005). 'On the Lumping Semantics of Counterfactuals'. *Journal of Semantics*, **22**: 129–151.

KRATZER, A. (1986). 'Conditionals'. In *Proceedings from the 22nd Regional Meeting of the Chicago Linguistic Society*, Chicago: University of Chicago. URL `http://semanticsarchive.net/Archive/ThkMjYxN/Conditionals.pdf`

KRATZER, A. (1989). 'An Investigation of the Lumps of Thought'. *Linguistics and Philosophy*, **12**(5): 607–653.

KRATZER, A. (1990). 'How Specific is a Fact?' In *Proceedings of the 1990 Conference on Theories of Partial In- formation*, Center for Cognitive Science, University of Texas at Austin.

KRATZER, A. (2002). 'Facts: Particulars or Information Units?' *Linguistics and Philosophy*, **25**(5–6): 655–670.

KRATZER, A. (2005). 'Constraining Premise Sets for Counterfactuals'. *Journal of Semantics*, **22**: 153–158.

LEWIS, D. K. (1973). *Counterfactuals*. Cambridge, Massachusetts: Harvard University Press.

LEWIS, D. K. (1979). 'Counterfactual Dependence and Time's Arrow'. *Noûs*, **13**: 455–476.

SANFORD, D. H. (1989). *If P then Q: Conditionals and the Foundations of Reasoning*. London: Routledge.

SLOTE, M. (1978). 'Time in Counterfactuals'. *Philosophical Review*, **7**(1): 3–27.

STALNAKER, R. C. (1968). 'A Theory of Conditionals'. In N. Rescher (ed.) *Studies in Logical Theory*, 98–112, Oxford: Basil Blackwell Publishers.

THOMASON, R., STONE, M. & DEVAULT, D. (to appear). 'Enlightened Update: A Computational Architecture for Presupposition and other Pragmatic Phenomena'. In D. Byron, C. Roberts & S. Schwenter (eds.) *Presupposition Accomodation*, Ohio State University. URL `http://www.cs.rutgers.edu/~mdstone/pubs/osu06.pdf`

TICHÝ, P. (1976). 'A Counterexample to the Stalnaker-Lewis Analysis of Counterfactuals'. *Philosophical Studies*, **29**: 271–273.

VELTMAN, F. (1996). 'Defaults in Update Semantics'. *Journal of Philosophical Logic*, **25**(3): 221–261.

VELTMAN, F. (2005). 'Making Counterfactual Assumptions'. *Journal of Semantics*, **22**: 159–180. URL `http://staff.science.uva.nl/~veltman/papers/FVeltman-mca.pdf`