



POTSDAM INSTITUTE FOR
CLIMATE IMPACT RESEARCH



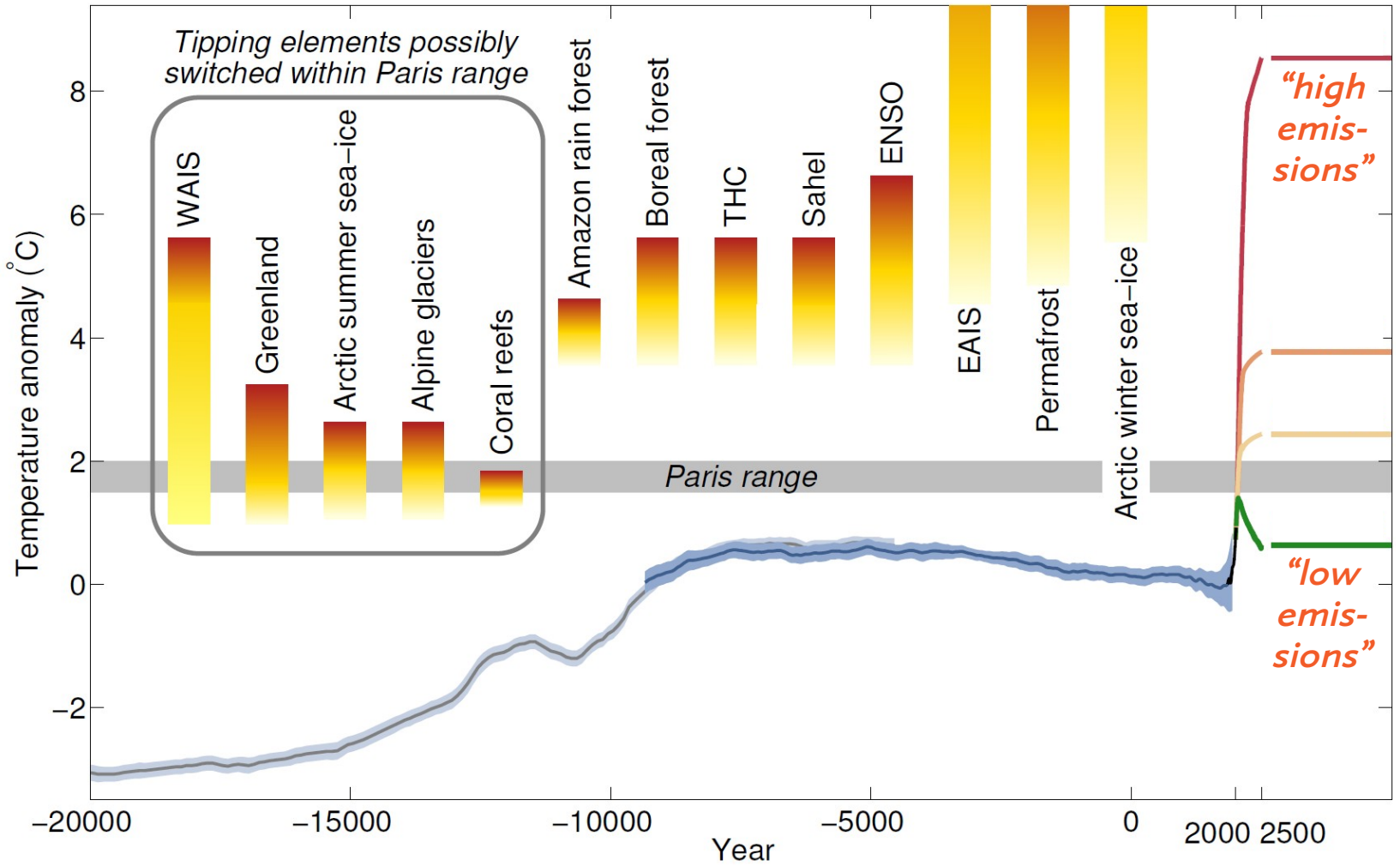
MEASURES OF INDIVIDUAL AND GROUPWISE EX-POST AND EX-ANTE RESPONSIBILITY IN EXTENSIVE-FORM GAMES WITH UNQUANTIFIABLE UNCERTAINTY

work in progress by **Jobst Heitzig & Sarah Hiller**

Formal Ethics 2019
Ghent, June 2019

Foretaste

CLIMATE TIPPING ELEMENTS, PARIS AGREEMENT

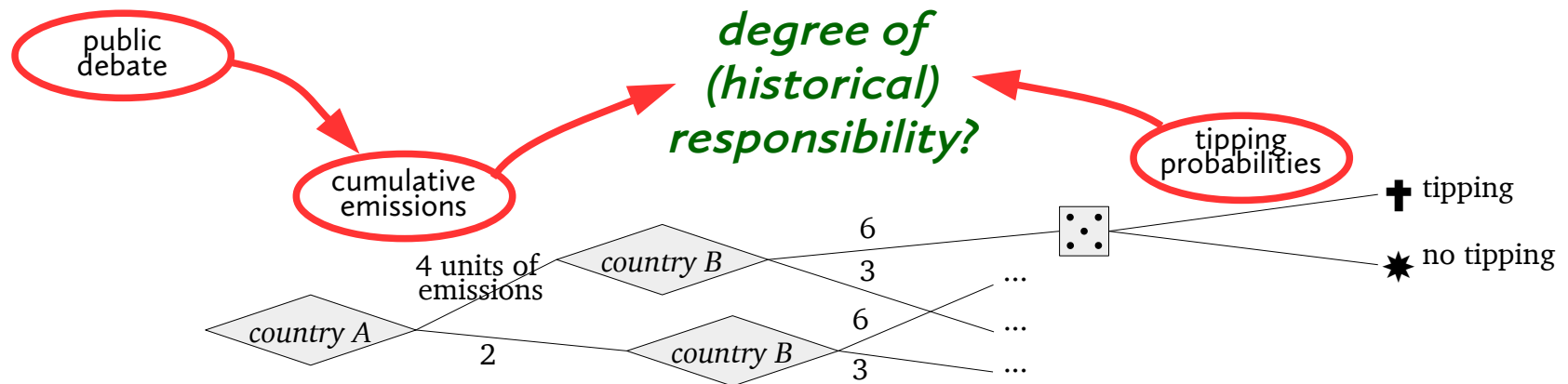


SchellInhuber et al., Nature Climate Change 2016



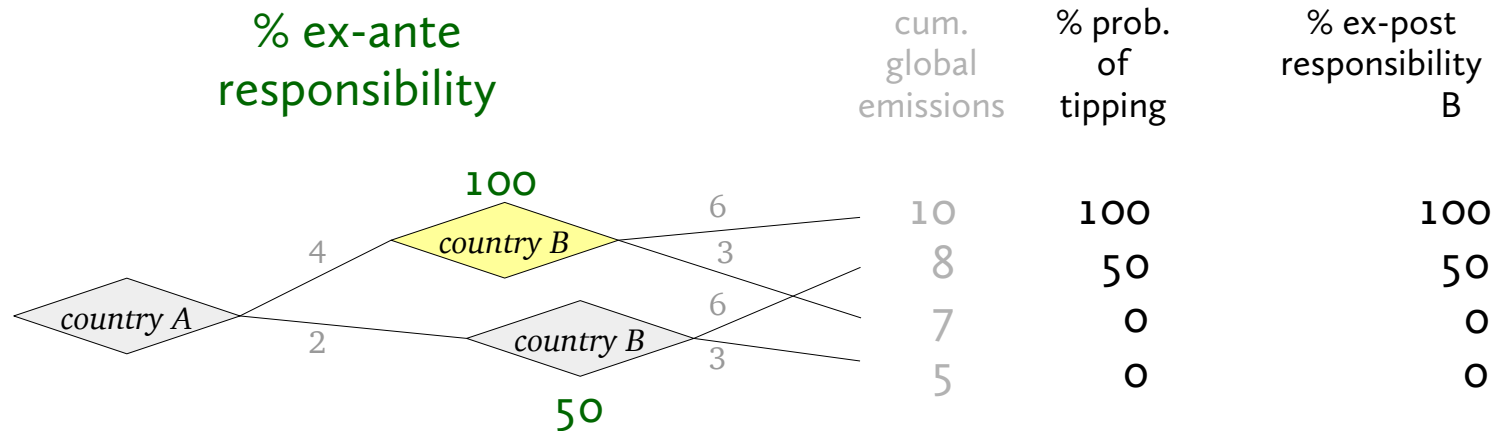
SIMPLISTIC EXAMPLE: TRIGGERING A CLIMATE TIPPING

Country A chooses high or low greenhouse gas emissions,
then country B chooses high or low emissions,
then unwanted climate tipping either occurs or not



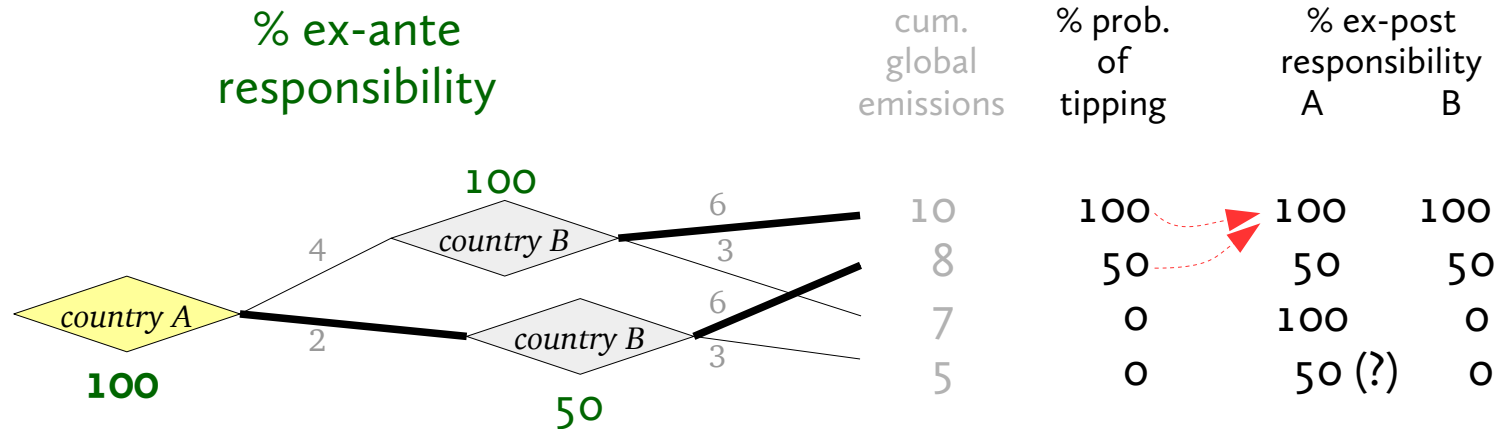
USE PROBABILITIES; BE OPTIMISTIC ABOUT INFLUENCE

Country A chooses high or low greenhouse gas emissions,
then country B chooses high or low emissions,
 then unwanted climate tipping either occurs or not



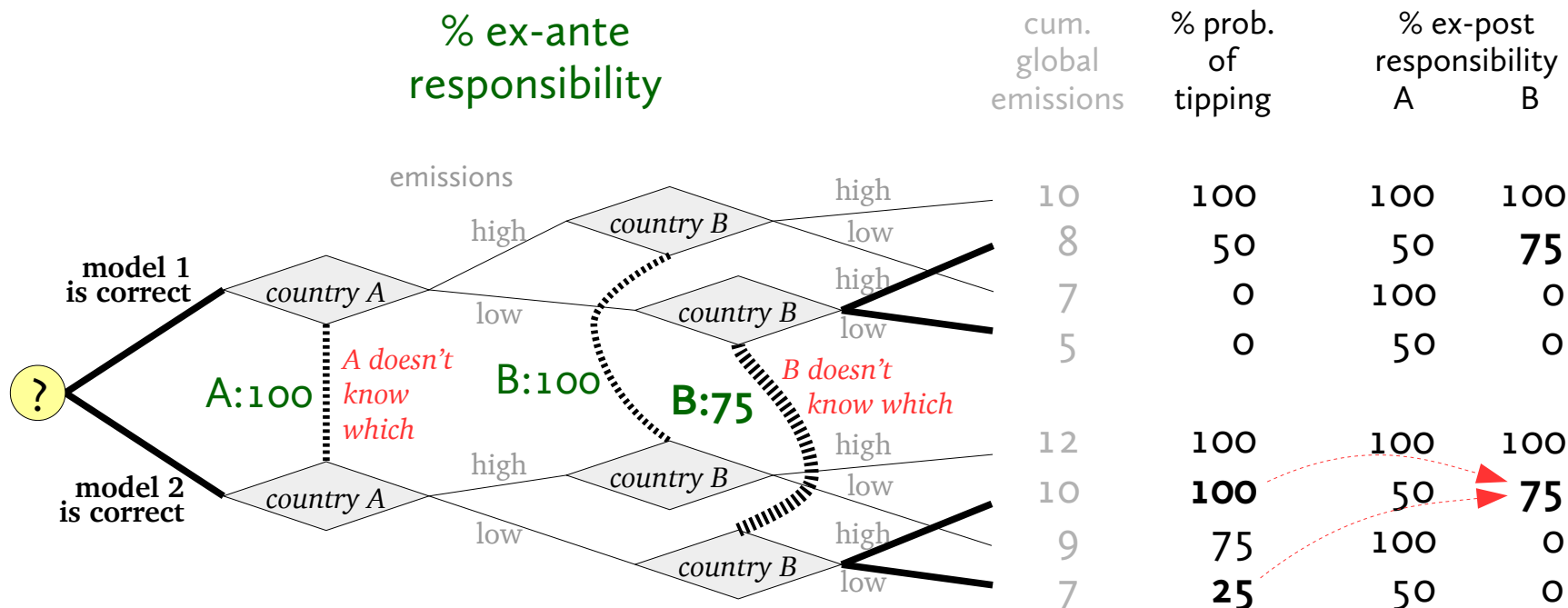
USE PROBABILITIES; BE OPTIMISTIC ABOUT INFLUENCE

Country A chooses high or low greenhouse gas emissions,
 then country B chooses high or low emissions,
 then unwanted climate tipping either occurs or not



BE PESSIMISTIC ABOUT UNKNOWNNS

Country A chooses high or low greenhouse gas emissions,
 then country B chooses high or low emissions,
 then unwanted climate tipping either occurs or not
 + uncertain probabilities



THE *MAX-DIFF* -FORMULA FOR *EX-ANTE* RESPONSIBILITY

Verbal definitions:

A *scenario S for agent group G at node v* is a choice of branch for each unquant. uncertainty node and all other agents' information states in the branch rooted at v.

A *strategy s for agent group G at node v* is a choice of action for each of G's information states in the branch rooted at v.

The *conditional value of G's strategy s at v given S* is a certain strictly increasing function f (e.g., $f(P) = P$ or $f(P) = \text{logit } P$) of the probability P , evaluated at v , of a good outcome, conditional on S and s .

G's conditional influence at v given S is the difference between the largest & smallest conditional values of all of G's strategies in S at v .

G's degree of ex-ante responsibility at v is its maximum conditional influence at v over all possible scenarios at v .

Formally:

$$\text{ear}(G, v) = \max \left\{ \max \left\{ f(P(\text{good}|v, S, s)) : \text{strat. } s \text{ for } G \text{ at } v \right\} - \min \left\{ f(P(\text{good}|v, S, s)) : \text{strat. } s \text{ for } G \text{ at } v \right\} : \text{scenario } S \text{ for } G \text{ at } v \right\}$$



Introduction

RANDOM CITATION:

“Whether humans are responsible for the bulk of climate change is going to be left to the scientists, but it's all of our responsibility to leave this planet in better shape for the future generations than we found it.”

(Mike Huckabee, US Republican)

SORT OF QUESTIONS WE CARE ABOUT HERE:

- Am I ethically **responsible for** climate change / Katrina / Bob's homelessness / etc.? And in what sense? And to what degree?
- Do I have ethical **responsibility to** mitigate climate change / compensate victims / etc.? And to what degree? And to what length or amount?
- How can responsibility be **quantified** in view of **many agents** & different forms of **uncertainty**?

Strategy:

study examples → formulate theses → suggest formulae

VAGUE INITIAL WORKING DEFINITION OF “MORE OR LESS RESPONSIBLE FOR/TO...”

Ex-post (=backward-looking) responsibility:

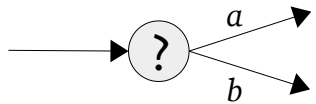
An agent i is the more responsible for a (typically ethically undesired) factual outcome q the more i could have exerted influence to prevent q .

Ex-ante (=forward-looking) responsibility:

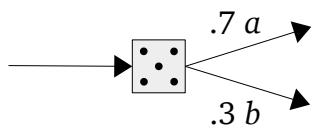
Given that outcome q is considered ethically undesirable, agent i has the more responsibility to help prevent q the more potential influence i can exert to prevent q .

Toolbox & First Example

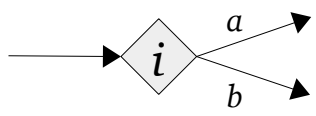
TREE-SHAPED MODELS OF MULTIAGENT DECISIONS UNDER UNCERTAINTY



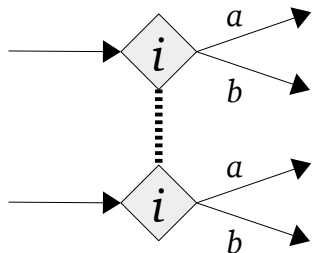
unquantified uncertainty node with branch labels
 modeller believes these are all possible subbranches but has no idea whatsoever about their probabilities



stochastic node with probabilities and branch labels
 modeller believes these are all possible subbranches and that these are their probabilities
 (different stochastic nodes are considered to be independent random events)



decision node with action labels for agent i
 modeller believes agent i has exactly these options and will decide among them at free will
 and that the probabilities of her choosing each option are unknown



information state with two decision nodes for agent i
 modeller believes agent i will not know in which of the two nodes she is

→ * **good outcome node**

→ † **bad outcome node**

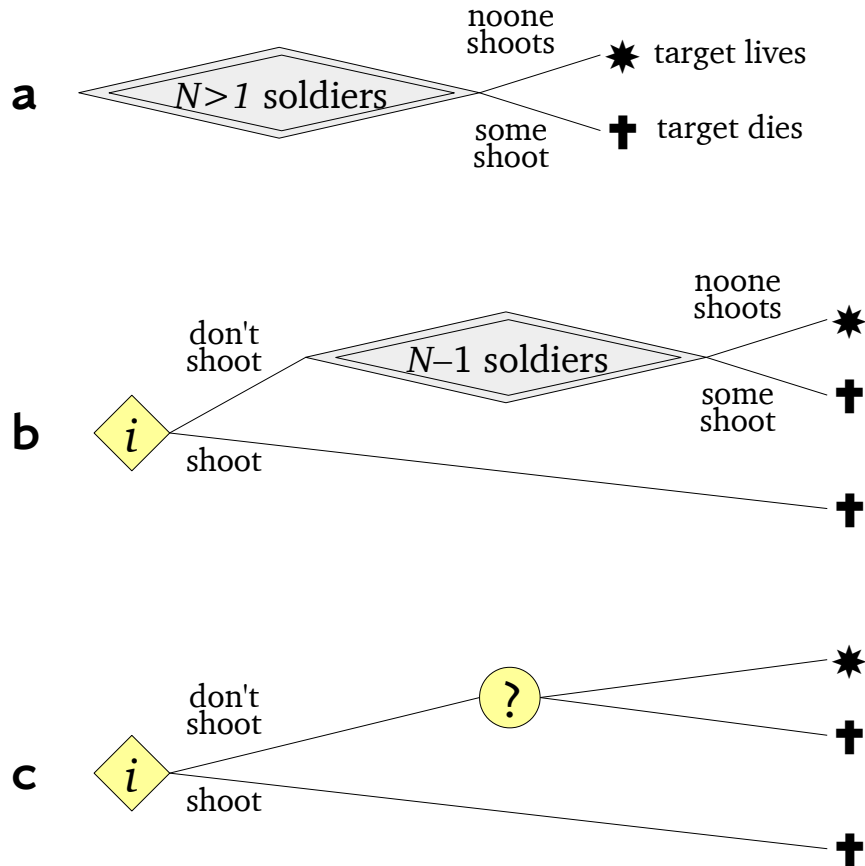


summarizing decision node with options for agent group G
 modeller believes group G has exactly these options and modeller is not interested in more detail

(enriched version of extensive-form game trees)

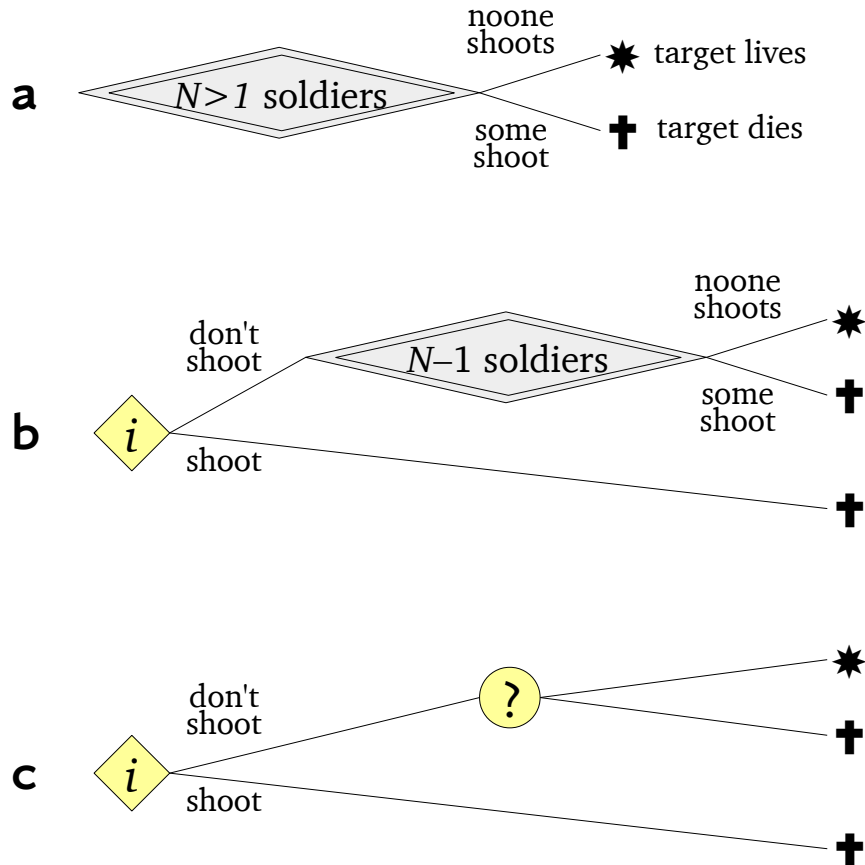
EXTREME EXAMPLE 1: KILLING BY FIRE-SQUAD

Different models of the situation:



EXTREME EXAMPLE 1: KILLING BY FIRE-SQUAD

Different models of the situation:



Thesis 1:

a,b,c are all *equivalent* w.r.t. the assessment of i 's responsibilities! In particular, the number N is irrelevant (if > 1)

Thesis 2:

i has, ethically, *full "ex-ante" responsibility* for what the result will be.

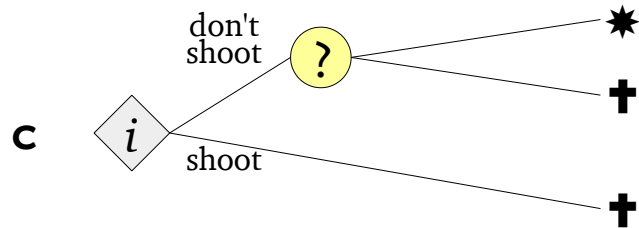
Thesis 3:

If i shoots, she has also *full "ex-post" responsibility* for the result. If not, she has none.

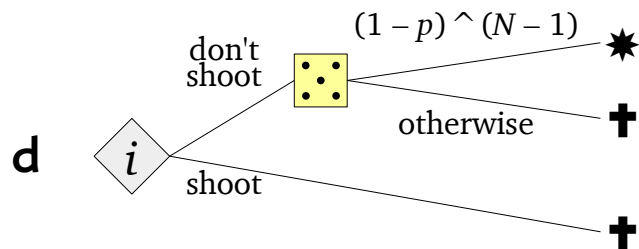
EXAMPLE 1 (CONTD.):

“ETHICAL” VS. “PSYCHOLOGICAL” ASSESSMENT

“Ethical observer's” model: *i* cannot know with what probability the others will shoot



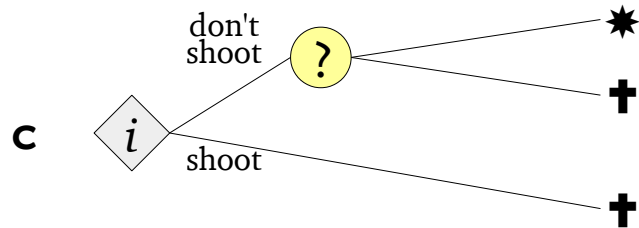
i might *prefer* the “psychological” model: each other soldier shoots with **probability** p , thus it's very likely that target dies anyway



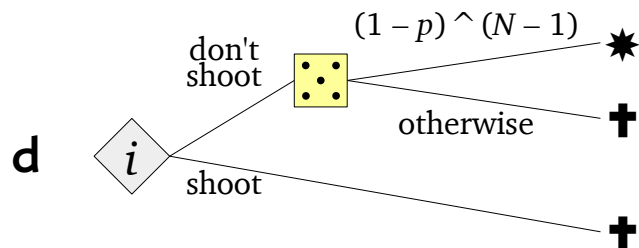
EXAMPLE 1 (CONTD.):

“ETHICAL” VS. “PSYCHOLOGICAL” ASSESSMENT

“Ethical observer's” model: *i* cannot know with what probability the others will shoot



i might *prefer* the “psychological” model: each other soldier shoots with **probability** p , thus it's very likely that target dies anyway



Thesis 4:

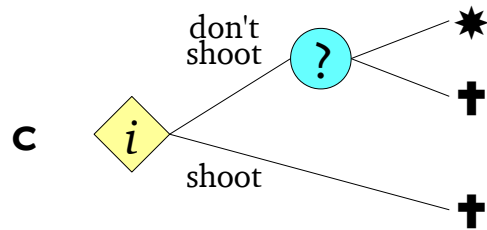
i cannot rightfully claim to know the probabilities that the others shoot, hence the *proper model* is **c** and not **d**.

Thesis 5:

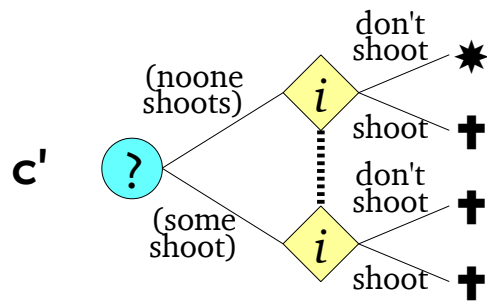
In other situations, where **d** would be the proper model, *i* would only have *partial responsibility* of degree $(1 - p)^{N - 1} \ll 1$ (resulting in exponentially fast “diffusion of responsibility”)

EXAMPLE 1 (CONTD.): WHAT ROLE DOES TIMING PLAY?

i decides “first”:

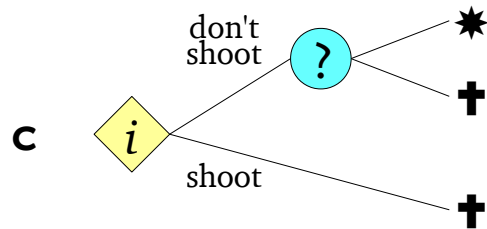


(some) others decide “first”,
but i doesn't know how
(*dashed line*):

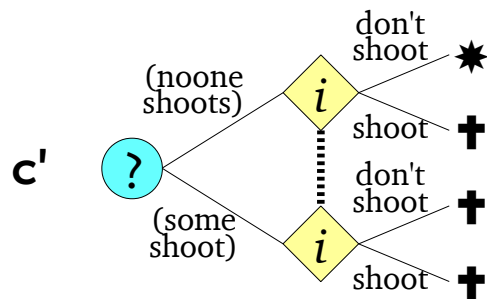


EXAMPLE 1 (CONTD.): WHAT ROLE DOES TIMING PLAY?

i decides “first”:



(some) others decide “first”,
but i doesn't know how
(*dashed line*):



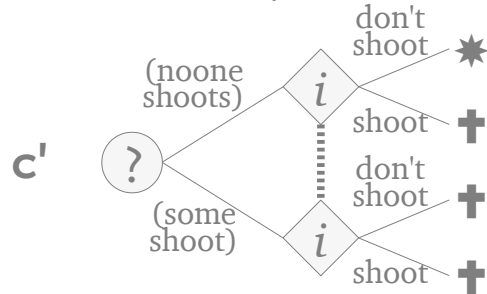
Thesis 6:

Such timing issues
have no effect on *what i can know*,
hence they are irrelevant here.
c, c' are in this sense equivalent.

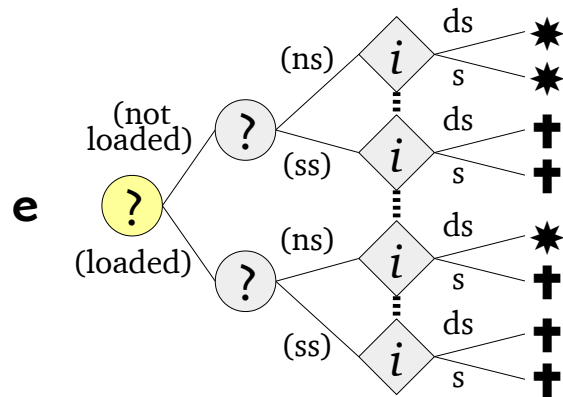
EXAMPLE 1 (CONTD.):

WHAT IF THE GUN MIGHT NOT BE LOADED?

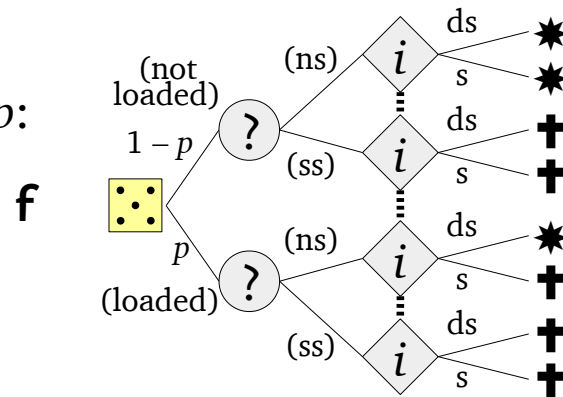
Gun certainly loaded:



Loaded with unknown probability:



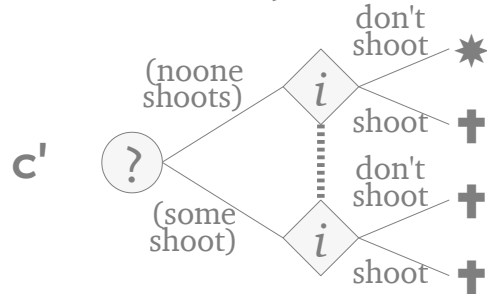
Loaded with known prob. p :



EXAMPLE 1 (CONTD.):

WHAT IF THE GUN MIGHT NOT BE LOADED?

Gun certainly loaded:



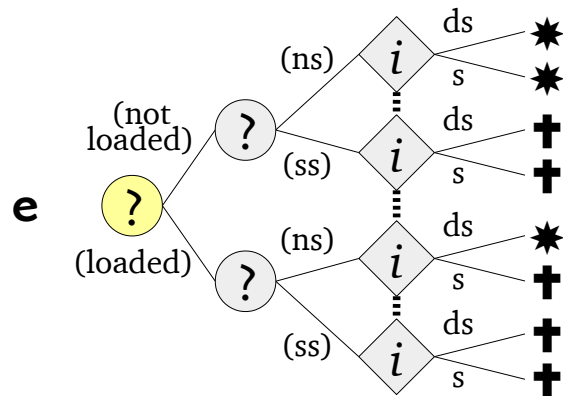
Thesis 7:

In **e**, *i* must consider the possibility that the gun might be loaded and has thus *full ex-ante responsibility*.

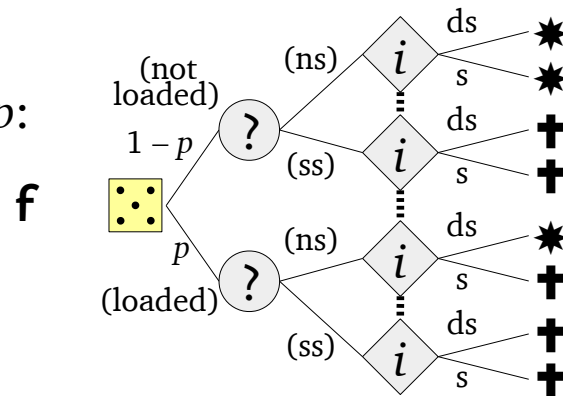
Thesis 8:

In **f**, *i* knows to have limited influence → *partial ex-ante responsibility of degree p*

Loaded with unknown probability:

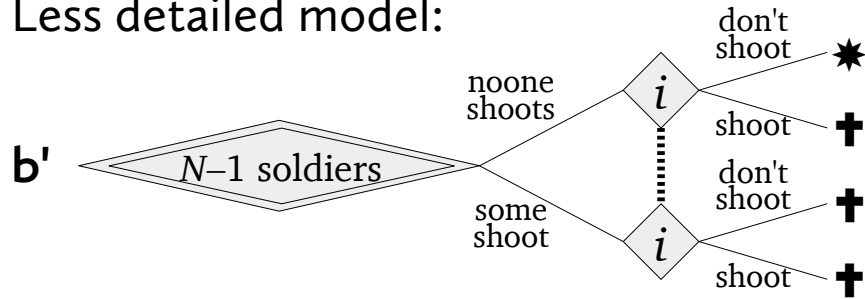


Loaded with known prob. p :

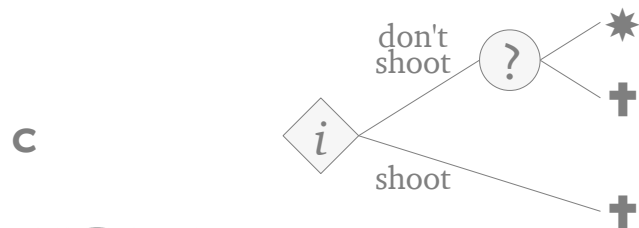
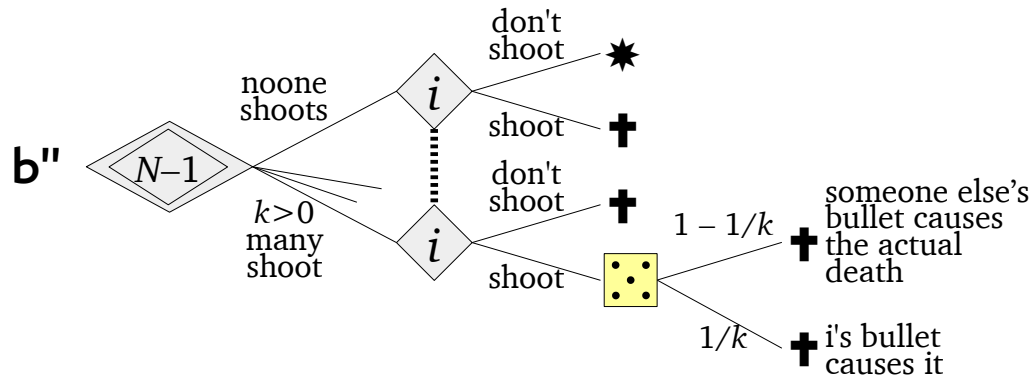


EXAMPLE 1 (CONTD.): WHOSE BULLET WAS IT?

Less detailed model:

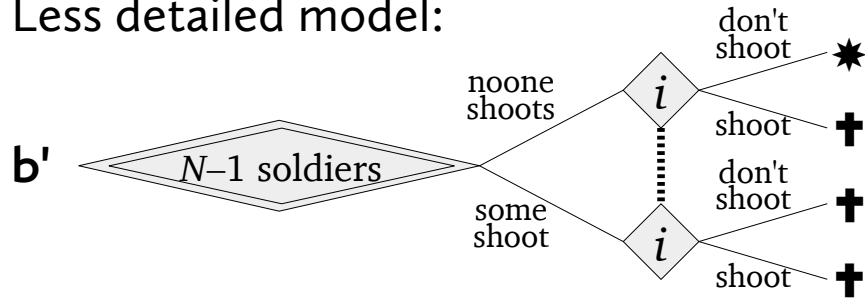


More detailed model:

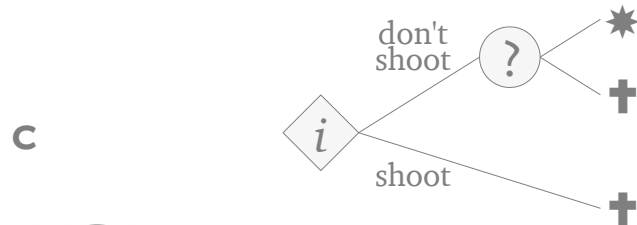
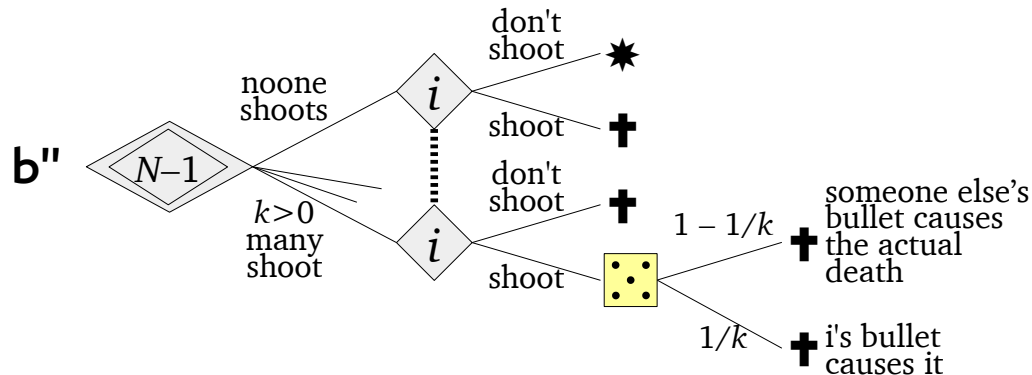


EXAMPLE 1 (CONTD.): WHOSE BULLET WAS IT?

Less detailed model:



More detailed model:



Thesis 9:

If each shot bullet *would* have killed, it doesn't matter which bullet *did* kill.

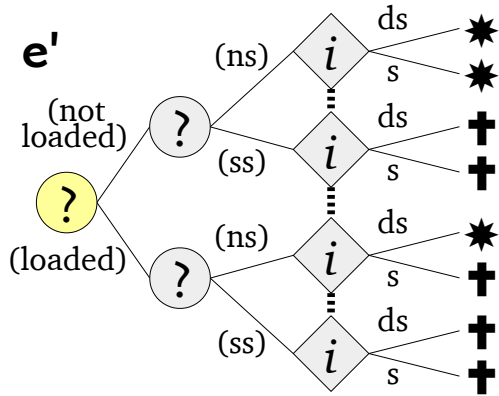
b'' is still *equivalent* to **b'** (and **b,c**) w.r.t. the assessment of i 's ex-post responsibility.

(So if i shoots, she has full ex-post responsibility for the actual result.)

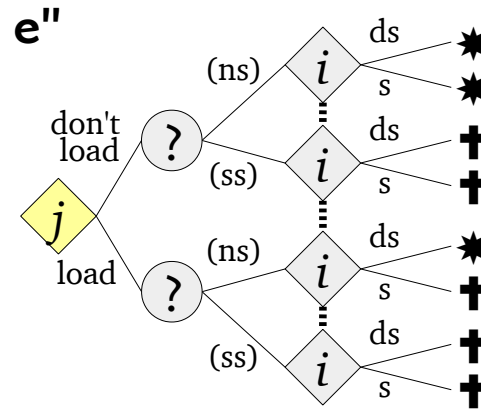
→ *Factual causes* matter less than *potential consequences and their probabilities*

“If you hadn't loaded the gun...” (STILL EXAMPLE 1)

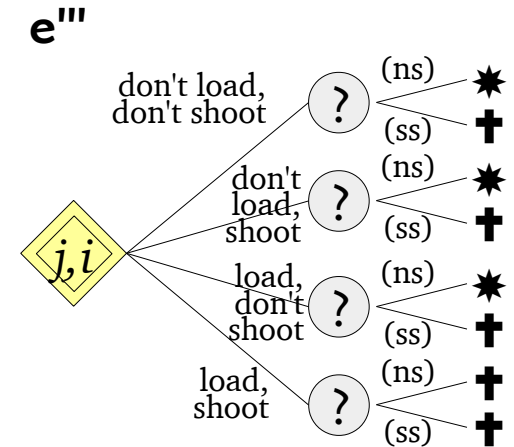
We had this model:



More specific model:



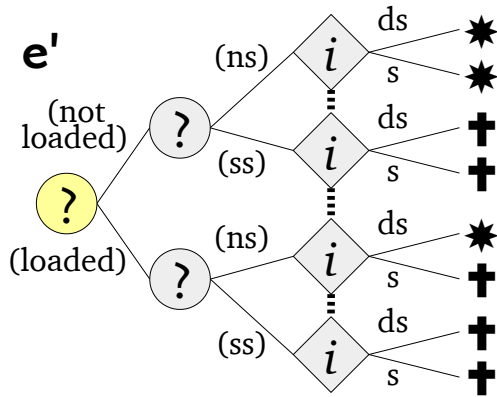
Less detailed model, treating $\{i,j\}$ as a group:



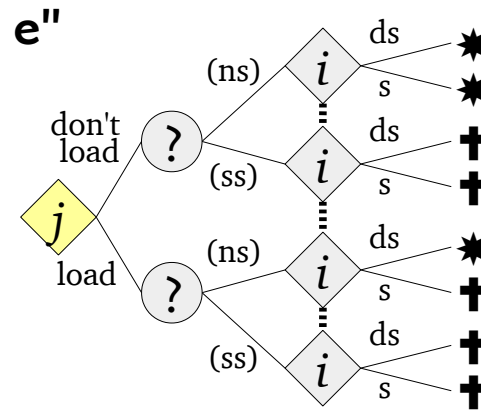
Is j less responsible than i?

“If you hadn't loaded the gun...” (STILL EXAMPLE 1)

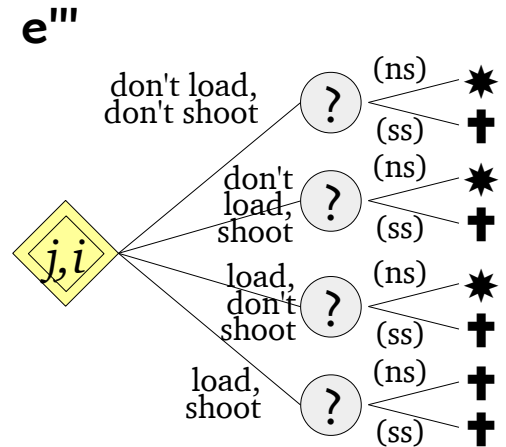
We had this model:



More specific model:



Less detailed model, treating $\{i,j\}$ as a group:



Is j less responsible than i?

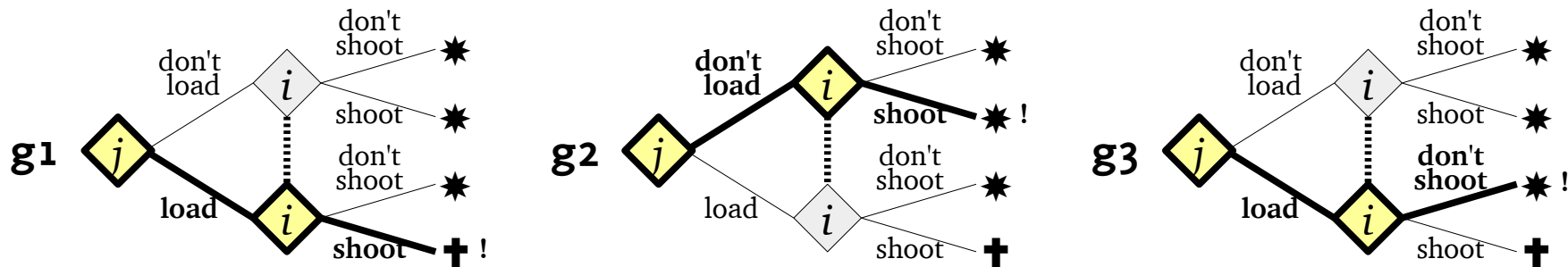
Thesis 10: *i must consider the possibility* that *j* has loaded the gun, but *likewise j must consider the possibility* that *i* will shoot, hence *both have full ex-ante responsibility* (i.e., “degree 1”).

Also the group $\{i,j\}$ has resp. degree $1 < 1+1$, hence responsibility is *not additive*.

(Note: some authors focussing on attribution consider a group only responsible if no smaller group is responsible)

FACTUAL VS. COUNTERFACTUAL EX-POST RESPONSIBILITY

Different *realizations* of the same, simplified model:

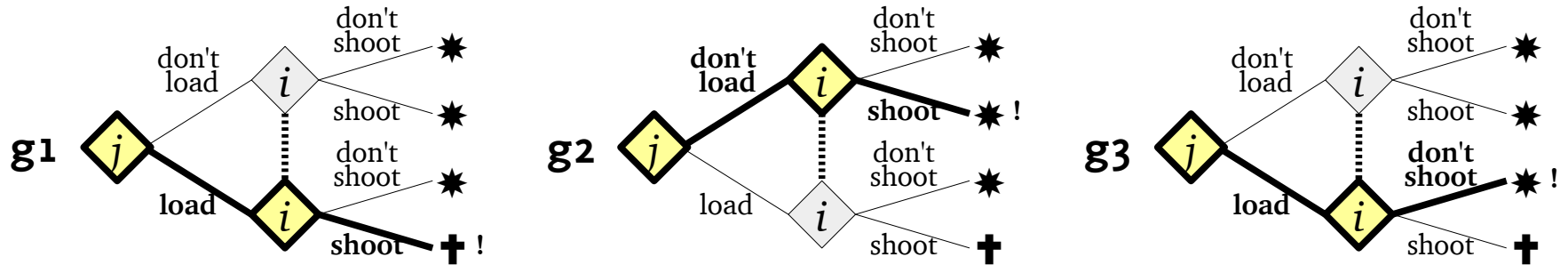


Both i and j are *fully ex-ante* responsible to avoid the target's death. In g_1 , both i and j are *fully ex-post* responsible for the factual death.

Should either of i or j be blamed in g_2 or g_3 ?

FACTUAL VS. COUNTERFACTUAL EX-POST RESPONSIBILITY

Different *realizations* of the same, simplified model:



Both i and j are *fully ex-ante* responsible to avoid the target's death.
In g_1 , both i and j are *fully ex-post* responsible for the factual death.

Should either of i or j be blamed in g_2 or g_3 ?

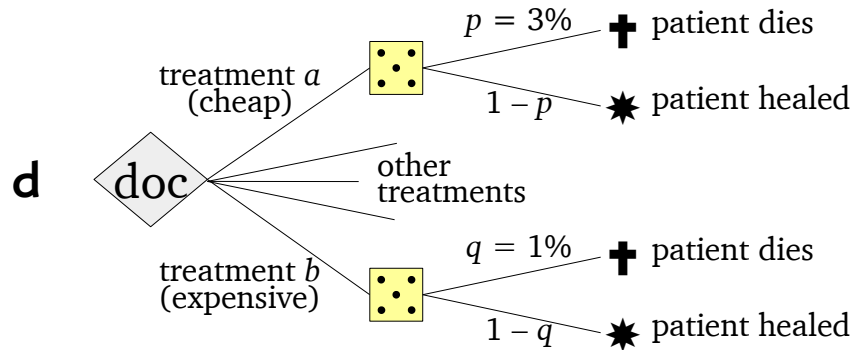
Thesis 11:

In g_2 , since the factual outcome is good, *i is not factually responsible* for a bad one, still *i 's action did make a bad outcome possible*, so i should be blamed for this: *i is fully counterfactually responsible regarding the target's possible death.*

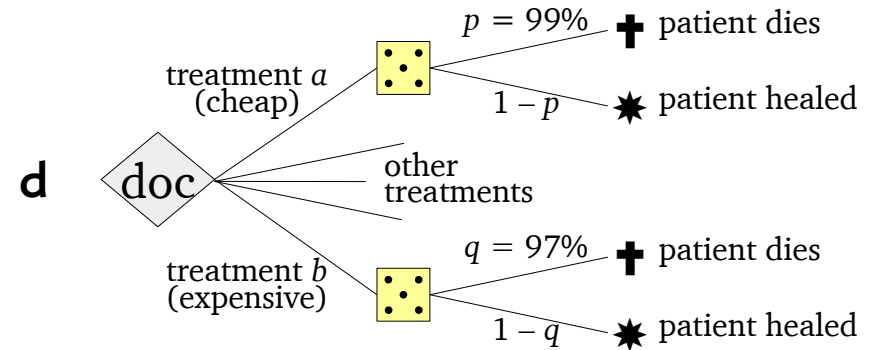
(and similarly j in g_3)

EXAMPLE 2: RISK OF DEATH IN MEDICAL TREATMENT

Low overall risk of death:

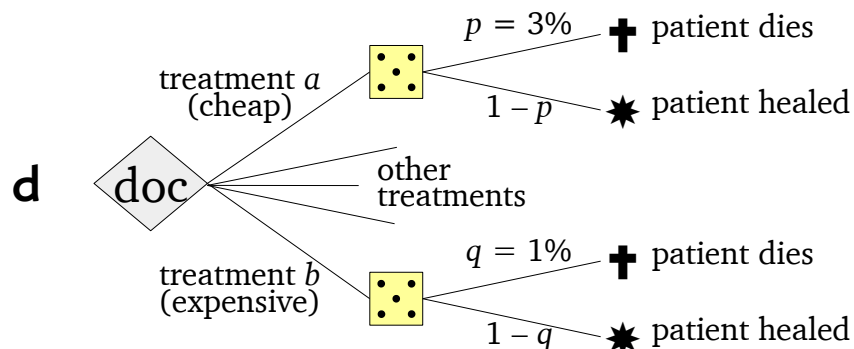


High overall risk of death:

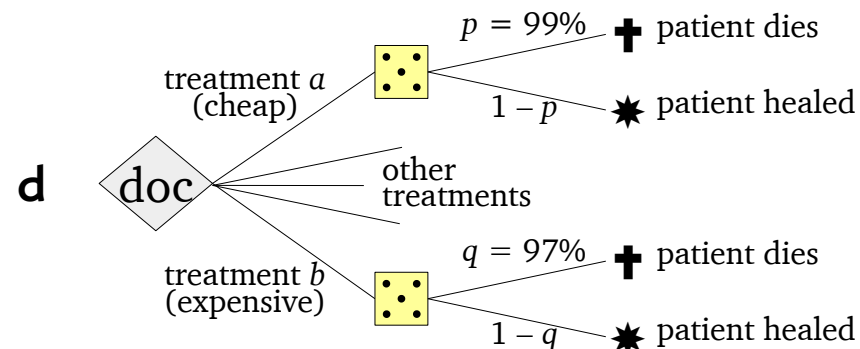


EXAMPLE 2: RISK OF DEATH IN MEDICAL TREATMENT

Low overall risk of death:



High overall risk of death:



Thesis 12:

Doctor's degree of ex-ante responsibility regarding the patient's possible death is a certain *function* $f(p,q)$ of the two probabilities p and q , which increases with larger p and decreases with larger q (assuming $p > q$).

Maybe simply $f(p,q) = p - q = 3\% - 1\% = 99\% - 97\% = 2\%$?

Or should $f(3\%,1\%)$, $f(51\%,49\%)$, and $f(99\%,97\%)$ all differ, and if so, how?

Suggested general principles & some candidate formulae

“[...] it seems rational to [...] concentrate [...] on the actual decision in light of the probabilities.”
(Nagel 1979: Moral luck)

PRINCIPLES FOR MEASURING RESPONSIBILITY (I)

derived from the preceding theses (nos. in brackets)

Unquantified uncertainty:

- *Others' unknown choices* (if at free will) should be treated as *unquantifiable uncertainty*, and not as stochasticity with some assumed probabilities. (1.+4.)
- Facing unquantified uncertainty, one must consider that *the “best” branch may have probability one.* (2.)
- Facing unquantified uncertainty, one must also consider that *the “worst” branch may have probability one.* (7.)

Stochastic uncertainty:

- Actual stochasticity may reduce the *degree* of responsibility in dependence of the resulting *probabilities.* (5.,8.,12.)

PRINCIPLES FOR MEASURING RESPONSIBILITY (II)

Three forms of responsibility:

“Forward-looking”: *options* (potential actions) may lead to

- *ex-ante* responsibility. (2.)

“Backward-looking”: *choices* (factual actions) may lead to

- *factual ex-post* responsibility or
- *counterfactual ex-post* responsibility. (3.+11.)

PRINCIPLES FOR MEASURING RESPONSIBILITY (III)

Relevant information:

- Consider *all potential consequences of all options*, not only the factual consequences or their factual causes. (9.,10.,11.)
- The factual outcome only determines the type of ex-post responsibility (factual or counterfactual), not its degree. (11.)
- Consider what agents *can be expected to know/can rightfully claim to know* at the time they act, not what they choose to believe. (4.)
- Timing issues that do not affect this knowledge are irrelevant. (6.)

PRINCIPLES FOR MEASURING RESPONSIBILITY (IV)

Sharing/division of responsibility:

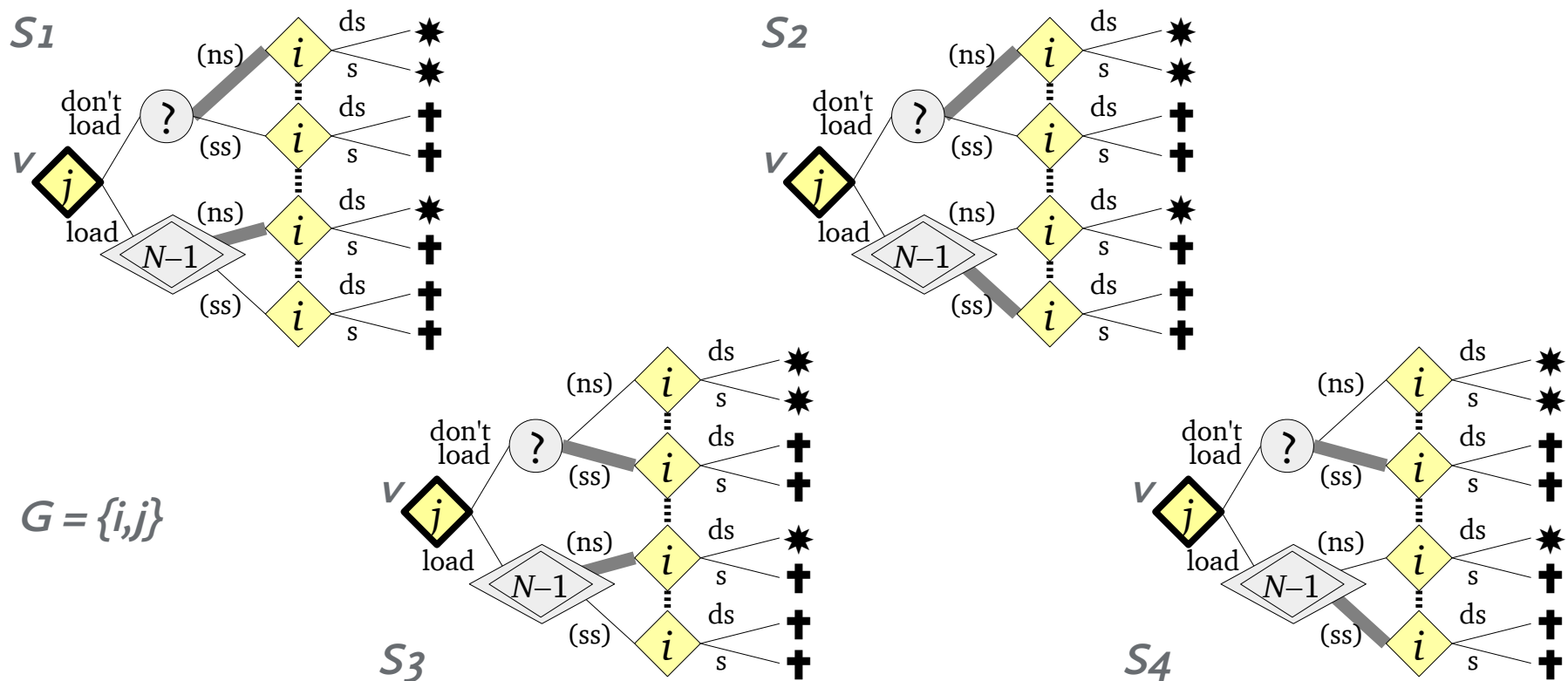
- Several agents may *all* individually be *fully ex-ante* responsible *and fully ex-post* responsible for the same outcome. (10.)
- A *group's* collective degree of responsibility (from 0 to 1) can be *smaller* than the sum of its members' individual degrees of responsibility. (10.)

(It is unclear at this point whether it can also be *larger*)

THE *MAX-DIFF*-FORMULA FOR *EX-ANTE* RESPONSIBILITY

Verbal definitions:

A *scenario S* for agent group *G* at node *v* is a choice of branch for each unquant. uncertainty node and all other agents' information states in the branch rooted at *v*.



THE *MAX-DIFF* -FORMULA FOR *EX-ANTE* RESPONSIBILITY

Verbal definitions:

A *scenario S for agent group G at node v* is a choice of branch for each unquant. uncertainty node and all other agents' information states in the branch rooted at v .

A *strategy s for agent group G at node v* is a choice of action for each of G 's information states in the branch rooted at v .

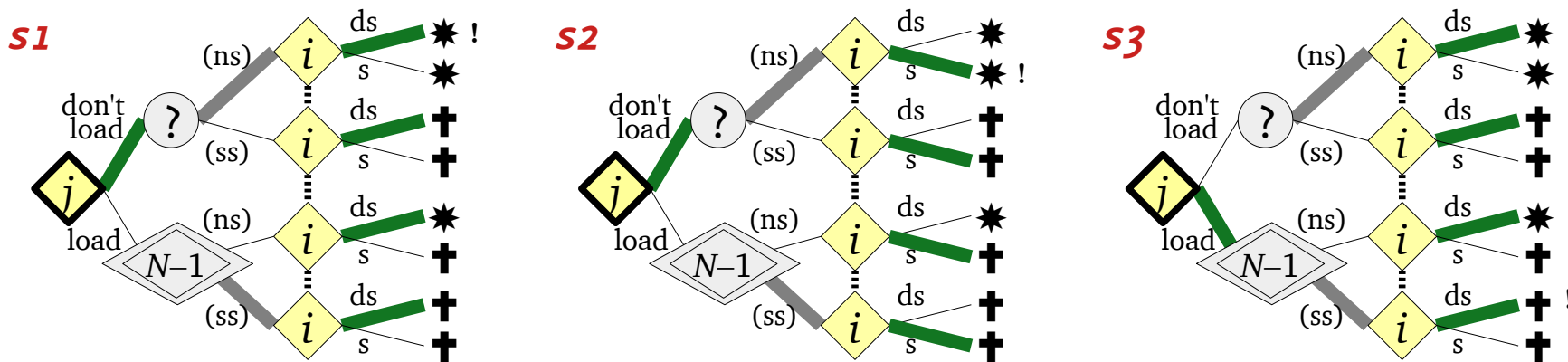
THE *MAX-DIFF* -FORMULA FOR *EX-ANTE* RESPONSIBILITY

Verbal definitions:

A *scenario S* for agent group *G* at node *v* is a choice of branch for each unquant. uncertainty node and all other agents' information states in the branch rooted at *v*.

A *strategy s* for agent group *G* at node *v* is a choice of action for each of *G*'s information states in the branch rooted at *v*.

$G = \{i, j\}$, scenario S_2 :



$s_4 = (\text{load}, \text{shoot})$

THE *MAX-DIFF* -FORMULA FOR *EX-ANTE* RESPONSIBILITY

Verbal definitions:

A *scenario S for agent group G at node v* is a choice of branch for each unquant. uncertainty node and all other agents' information states in the branch rooted at v .

A *strategy s for agent group G at node v* is a choice of action for each of G 's information states in the branch rooted at v .

The *conditional value of G's strategy s at v given S* is a certain strictly increasing function f (e.g., $f(P) = P$ or $f(P) = \text{logit } P$) of the probability P , evaluated at v , of a good outcome, conditional on S and s .



THE *MAX-DIFF* -FORMULA FOR *EX-ANTE* RESPONSIBILITY

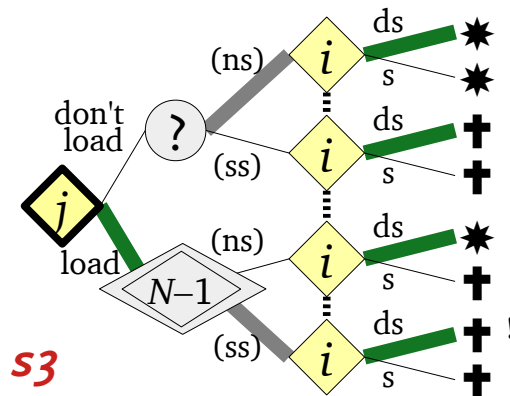
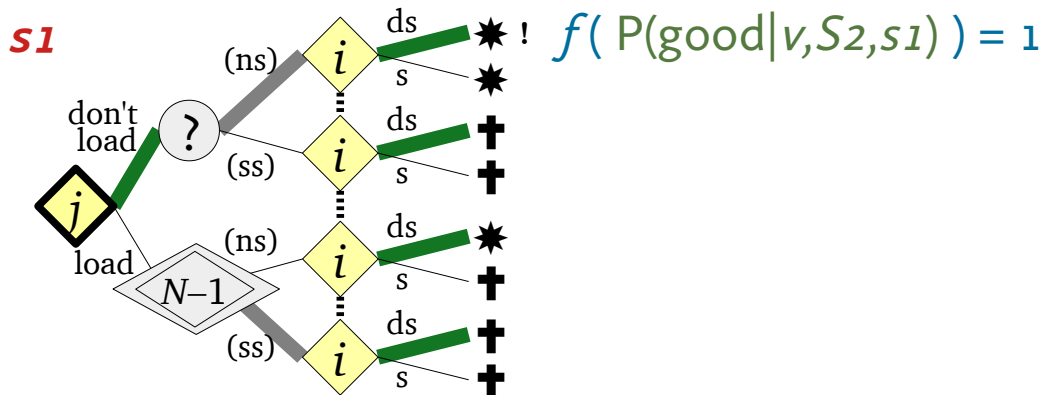
Verbal definitions:

A *scenario S* for agent group *G* at node *v* is a choice of branch for each unquant. uncertainty node and all other agents' information states in the branch rooted at *v*.

A *strategy s* for agent group *G* at node *v* is a choice of action for each of *G*'s information states in the branch rooted at *v*.

The *conditional value of G's strategy s at v given S* is a certain strictly increasing function *f* (e.g., $f(P) = P$ or $f(P) = \text{logit } P$) of the probability *P*, evaluated at *v*, of a good outcome, conditional on *S* and *s*.

$G = \{i, j\}$, scenario S_2 , using $f(P) = P$:



THE *MAX-DIFF* -FORMULA FOR *EX-ANTE* RESPONSIBILITY

Verbal definitions:

A *scenario S for agent group G at node v* is a choice of branch for each unquant. uncertainty node and all other agents' information states in the branch rooted at v.

A *strategy s for agent group G at node v* is a choice of action for each of G's information states in the branch rooted at v.

The *conditional value of G's strategy s at v given S* is a certain strictly increasing function f (e.g., $f(P) = P$ or $f(P) = \text{logit } P$) of the probability P , evaluated at v , of a good outcome, conditional on S and s .

G's conditional influence at v given S is the difference between the largest & smallest conditional values of all of G's strategies in S at v .

$$ci(G, v | S) = \max \{ f(P(\text{good} | v, S, s)) : \text{strat. } s \text{ for } G \text{ in } S \text{ at } v \} \\ - \min \{ f(P(\text{good} | v, S, s)) : \text{strat. } s \text{ for } G \text{ in } S \text{ at } v \}$$



THE *MAX-DIFF* -FORMULA FOR *EX-ANTE* RESPONSIBILITY

Verbal definitions:

A *scenario S for agent group G at node v* is a choice of branch for each unquant. uncertainty node and all other agents' information states in the branch rooted at v .

A *strategy s for agent group G at node v* is a choice of action for each of G 's information states in the branch rooted at v .

The *conditional value of G's strategy s at v given S* is a certain strictly increasing function f (e.g., $f(P) = P$ or $f(P) = \text{logit } P$) of the probability P , evaluated at v , of a good outcome, conditional on S and s .

G's conditional influence at v given S is the difference between the largest & smallest conditional values of all of G 's strategies in S at v .

G's degree of ex-ante responsibility at v is its maximum conditional influence at v over all possible scenarios at v .



THE *MAX-DIFF* -FORMULA FOR *EX-ANTE* RESPONSIBILITY

Verbal definitions:

A *scenario S for agent group G at node v* is a choice of branch for each unquant. uncertainty node and all other agents' information states in the branch rooted at v.

A *strategy s for agent group G at node v* is a choice of action for each of G's information states in the branch rooted at v.

The *conditional value of G's strategy s at v given S* is a certain strictly increasing function f (e.g., $f(P) = P$ or $f(P) = \text{logit } P$) of the probability P , evaluated at v , of a good outcome, conditional on S and s .

G's conditional influence at v given S is the difference between the largest & smallest conditional values of all of G's strategies in S at v .

G's degree of ex-ante responsibility at v is its maximum conditional influence at v over all possible scenarios at v .

Formally:

$$\text{ear}(G, v) = \max \left\{ \max \left\{ f(P(\text{good}|v, S, s)) : \text{strat. } s \text{ for } G \text{ at } v \right\} - \min \left\{ f(P(\text{good}|v, S, s)) : \text{strat. } s \text{ for } G \text{ at } v \right\} : \text{scenario } S \text{ for } G \text{ at } v \right\}$$



EX-POST RESPONSIBILITY VERSION 1: THE *SUM-MAX-DIFF-MAX*-FORMULA

Verbal definitions:

G's reachable target in S at v

is the maximum conditional value of all of G 's strategies in S at v .

G's shortfall in S at decision node v due to action a

is the difference between G 's reachable targets in S at nodes v and $v:a$.

G's incremental ex-post-degree of responsibility due to a

is its maximum shortfall at $node(a)$ due to a over all possible scenarios at $node(a)$.

G's total ex-post-degree of responsibility (version 1)

is the sum of its increm. ex-post-degrees over all actions actually taken by G .

In short:

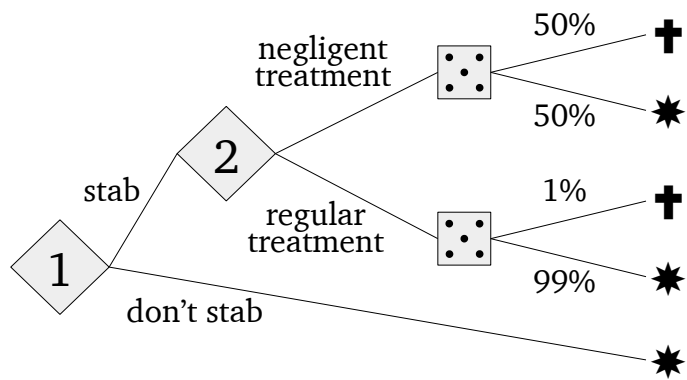
$rt(G, S, v) = \max \{ f(P(\text{good} | v, S, s)) : \text{strategy } s \text{ for } G \text{ in } S \text{ at } v \}$

$epr_1(G) = \text{sum} \{ \max \{ rt(G, S, v) - rt(G, S, v:a) : S \text{ at } node(a) \} : \text{action } a \text{ taken by } G \}$

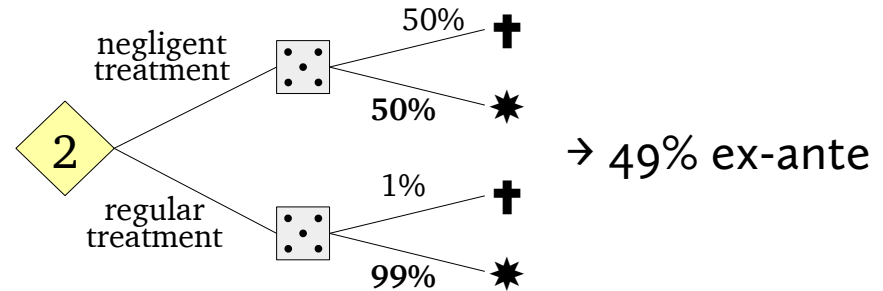


Further Examples

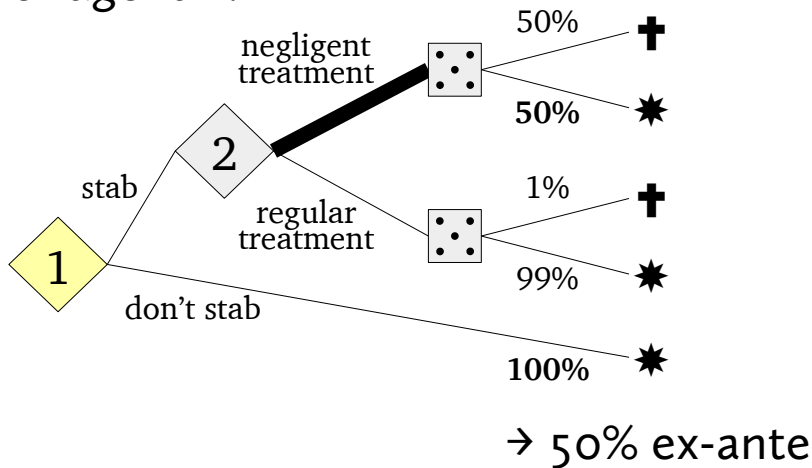
EXAMPLE à la Canavotto & Giordani (yesterday):



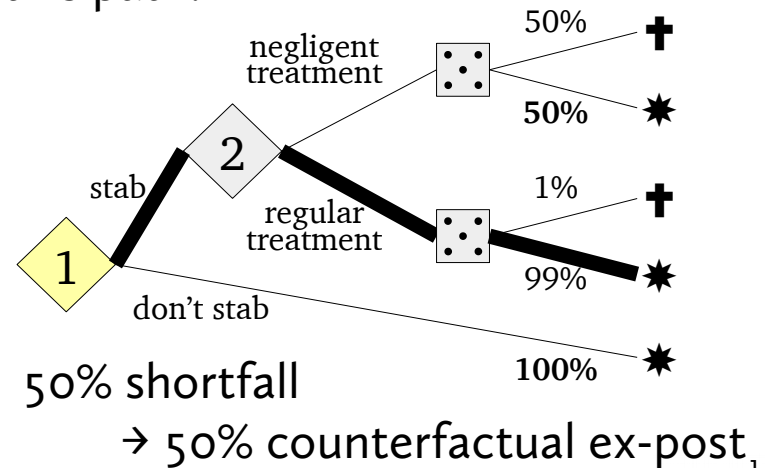
For agent 2, there is only one scenario:



Influence-maximizing scenario for agent 1:

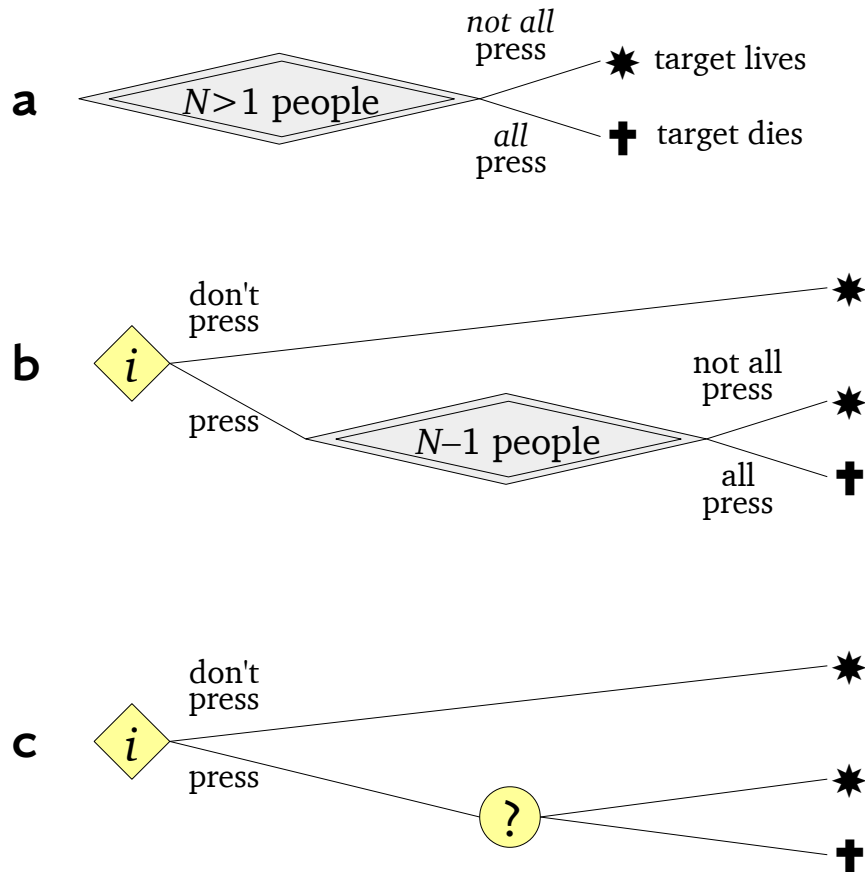


Ex-post resp. of agent 1 after this path:



AN ELECTRIC CHAIR

Target dies iff *all* N persons press a button.



Our formula implies:
a, b, c are all *equivalent*
w.r.t. the assessment
of i 's responsibilities;
the number N is irrelevant.

i has *full ex-ante responsibility*
(like any other subgroup of the N).

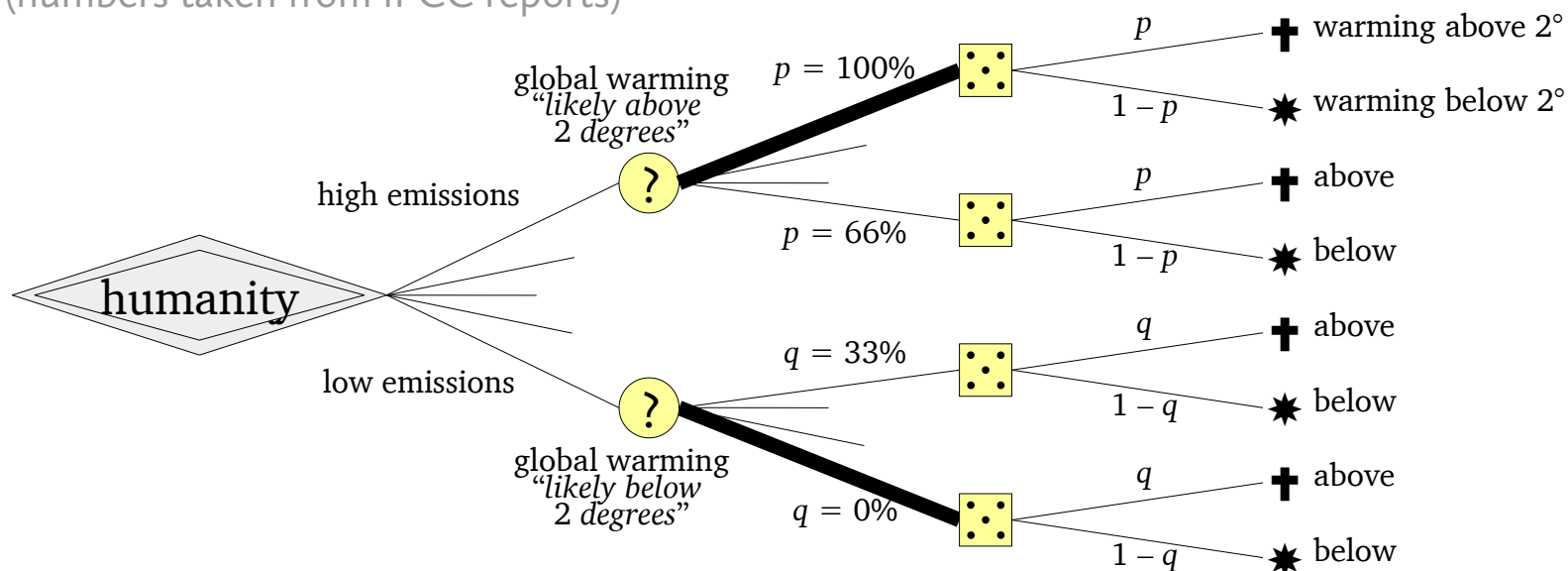
i [or some group G] has
either full or no ex-post resp.,
depending on whether
 i [or at least one member of G]
has pressed her button.

These implications seem OK...

INTERMEDIATE LEVELS OF UNCERTAINTY

Exact probabilities unknown, but known to be within a certain range
→ Represent by combination of unquantified uncertainty and stochasticity nodes

E.g.: staying below two degrees of global warming
(numbers taken from IPCC reports)



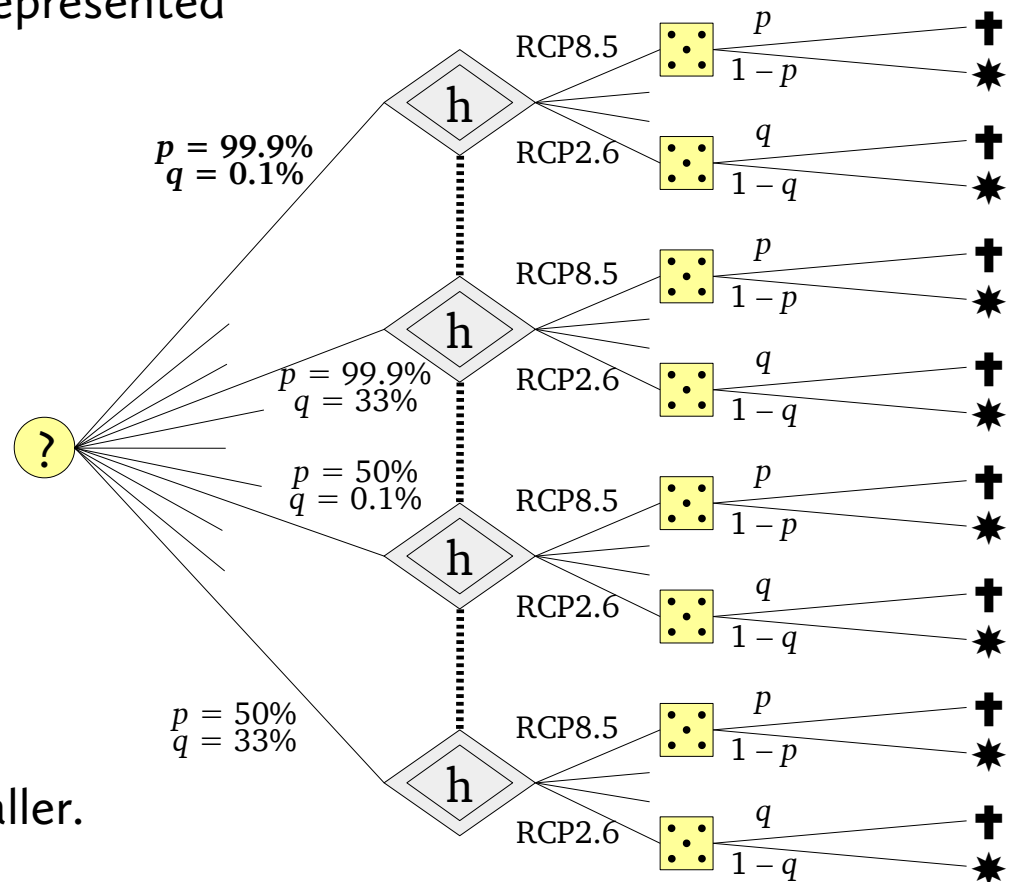
→ ex-ante resp. is $100\% - 0\% = 100\%$ (and not only $66\% - 33\% = 33\%$)

INDEPENDENT VS. DEPENDENT UNCERTAINTY

The situation may be more clearly represented by this (formally equivalent) model:

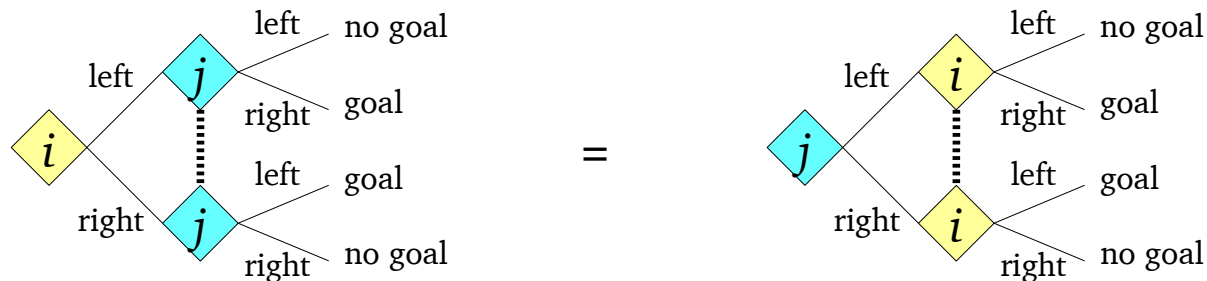
This representation shows that the model is adequate if the *uncertainty* about the effect of RCP8.5 and the uncertainty about the effect of RCP2.6 are *independent*.

If the science says they are somehow correlated, some of the many branches may not exist and hence the responsibility may be smaller.



UNAVOIDABLE EX-POST RESPONSIBILITY

Penalty kick: penalty-taker (i) kicks into one corner, goalie (j) jumps into another?



The complete symmetry of the model implies:

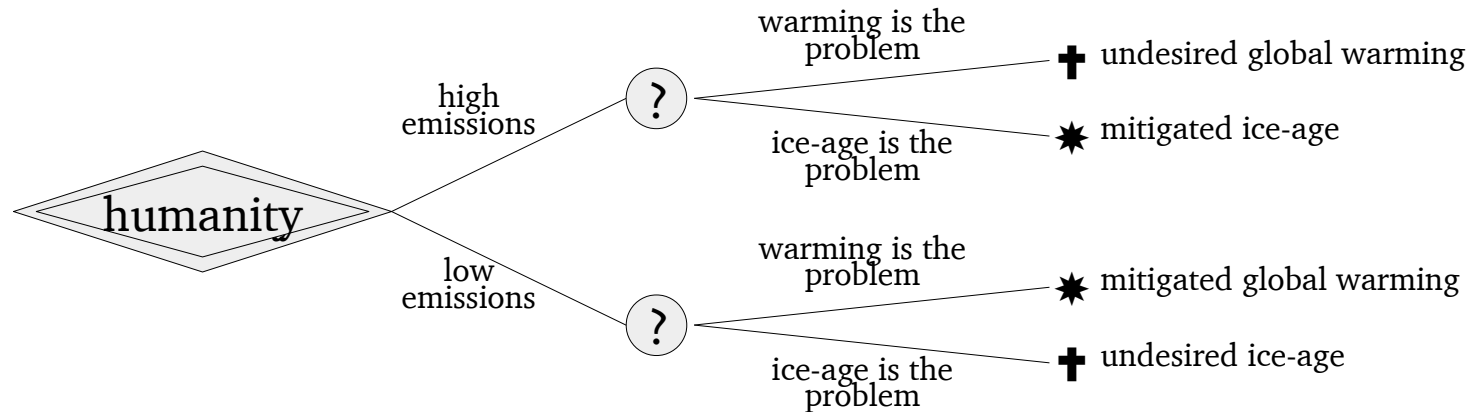
- both have the same degree of ex-ante responsibility
- both have the same degree of ex-post resp., no matter what they do
- the degree itself also does not depend on what they do!

In our theory:

*No matter what they do, both have full ex-post responsibility
(and thus both might be held accountable for the outcome).*

SOMETIMES THERE IS NO “RIGHT THING” TO DO?

1980ies: Unclear whether GHG emissions support undesired global warming or help preventing an undesired imminent ice-age (private comm. with B. Hoskins)



If this model represents their knowledge at the time, our theory implies:
They had to expect to be ex-post (either factually or counterfactually) responsible anyway, no matter what they would do.

Is this reasonable?

EX-POST RESPONSIBILITY VERSION 2: THE *SUM-DIFF-MAXIMIN* -FORMULA

Verbal definitions:

G's guaranteed value of strategy s at node v
is the minimum conditional value of s over all scenarios S of G at v .

G's precautionary target at v
is the maximum guaranteed value of all of G 's strategies at v .

G's shortfall at v due to action a
is the difference between G 's precautionary targets at nodes v and $v:a$.

G's total ex-post-degree of responsibility (version 2)
is the sum of its shortfalls over all node-action pairs actually traversed by G .

In short:

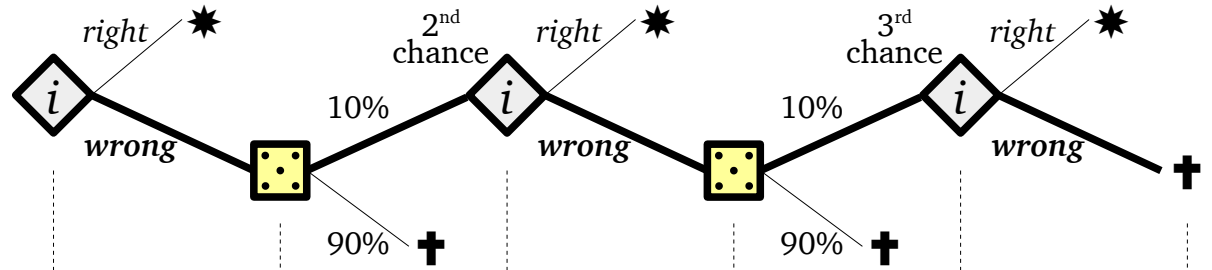
$pt(G,v) = \max \{ \min \{ f(P(\text{good}|v,S,s)) : \text{scen. } S \text{ of } G \text{ at } v \} : \text{strat. } s \text{ for } G \text{ in } S \text{ at } v \}$

$epr_2(G) = \text{sum} \{ pt(G,v) - pt(G,v:a) : \text{node-action pair } (v,a) \text{ traversed by } G \}$



REPEATED FAILURE CAN INCREASE EX-POST RESPONSIBILITY BEYOND 100%!

Unexpected 2nd and 3rd chances, all failed:



Quantities:

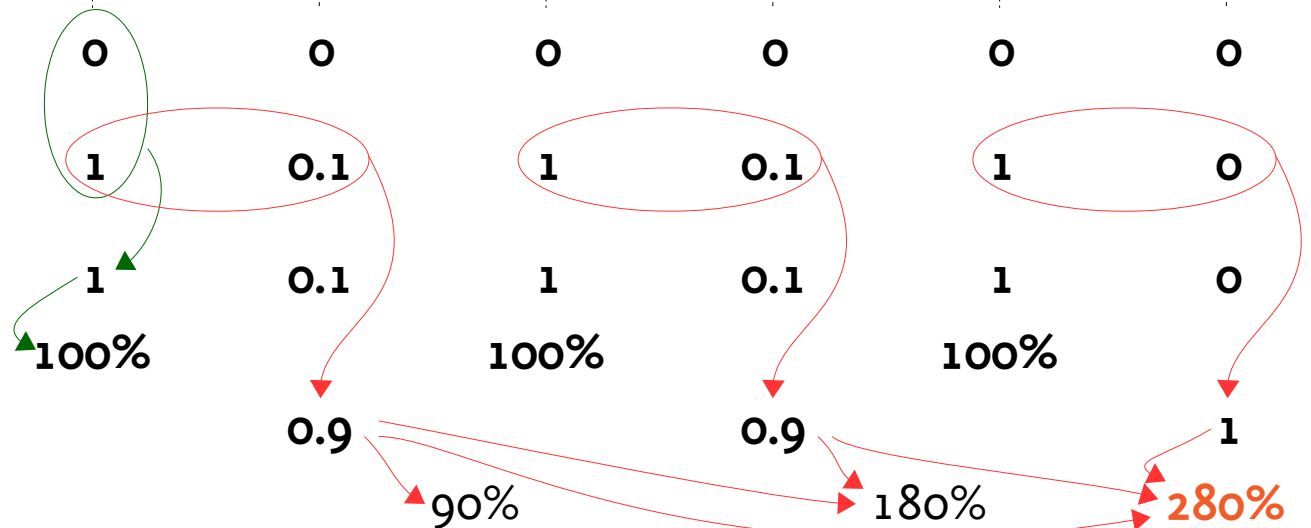
smallest prob. of ★
 reachable target =
 largest prob. of ★

influence = diff.

ex-ante resp.

shortfall = diff.

ex-post resp.



SOME FURTHER PROPERTIES

- responsibility **increases** with
 - more options or repeated chances to act
 - higher uncertainty, earlier chances to act
- responsibility **does not decrease** with
 - more other players (no division or “diffusion of resp.”)
 - luck after a wrong decision

SOME NEXT STEPS

- Read all the stuff you guys have already done on this (Sorry for not citing anything – we were ignorant of much of the existing formal literature)
- Figure out which formula for ex-post makes more sense
- Relate to logics (use game trees in semantics, include “degree-of-responsibility” quantifiers, etc.?)
- Analyse natural-language use of “responsibility”
- Apply to somewhat realistic model of climate change
- ... *Suggestions?*

Thank you for your attention!

heitzig@pik-potsdam.de