

THE WHOLE TRUTH ABOUT THE LAW – REASONING ABOUT EXCEPTIONS IN LEGAL AI

Stephan Leuenberger / Burkhard Schafer

Lecturer, University of Glasgow, Dept. of Philosophy, School of Humanities

University Gardens, Glasgow G12 8QQ, UK

Stephan.Leuenberger@glasgow.ac.uk; <http://www.gla.ac.uk/schools/humanities/staff/stephanleuenberger>

Professor of Computational Legal Theory, The University of Edinburgh, SCRIPT Centre for IT and IP law

Old College, South Bridge, Edinburgh, EH8 9YL, UK

B.schafer@ed.ac.uk <http://www.law.ed.ac.uk/people/burkhardshafer>

Keywords: *Defeasibility, Holton, That's it, legal reasoning*

Abstract: *Defeasible reasoning plays an important part in understanding and modelling legal argumentation. The most commonly used approaches in AI and Law however do not capture legal disputes that are themselves about the legal understanding of defeasibility, argument types that nonetheless play an important role in judicial review or in appeals. We introduce a theory of reasoning about exceptions (or the lack of them) that has been developed by Richard Holton in an attempt to clarify our understanding of the status of ethical norms. We show its potential to add to our theoretical machinery for the analysis of legal reasoning, but also suggest some necessary refinements.*

1. Introduction and motivation¹

An important skill for a practicing lawyer is to identify, discover or invent exceptions to superficially rigid looking legal rules. The ability to «distinguish» a precedent and to argue for a novel exception to the ratio that it (seemed to have) established is an important mechanism for the evolution of the law and its ability to adopt to changing circumstances – though it can also lead to the famous «loopholes» or «technicalities» of popular discourse.

This indicates an inherent tension in the law: On the one hand, we want rules that are simple, clear, come with a high degree of predictability in their application, enabled as a principle of fairness to treat like cases alike, and also are authoritative in the sense that they can be traced back to an appropriate legislative source. On the other hand, we also want to do justice in the individual case and also achieve wider social goals in rapidly changing environments without having to legislate constantly. *Fiat Justitia et pereat mundus* has never been a tenable proposition for a legal system.

Legislators and judges through the centuries have aimed to maintain an uneasy equilibrium between these conflicting demands, accepting on the one hand the role and importance of general rules given by the sovereign, but at the same time recognise the importance for flexibility and discretionary decision making in the everyday application of the law in courts. Law, then is (also) a system of rules, but of rules that allow for exceptions and the identification of as-yet unrecognised defences.²

For legal theorists, especially legal theorists interested in formal accounts of the validity of legal argumentation, the resulting practice of legal decision making provides a conundrum.

¹ Work on this paper was supported by the Arts and Humanities Research Council (grant number AH/M009610/1).

² For a comprehensive informal discussion of defeasibility in law from a jurisprudential perspective see D'ALMEIDA, *Allowing for Exceptions: A Theory of Defences and Defeasibility in Law*. Oxford University Press, Oxford 2015.

The following argument (A) seems to be a valid justification of the imposition of punishment as part of a court decision:

- 1) As a valid rule of our legal system, killing another human being is prohibited and punished with 20 years imprisonment.
 - 2) Jane has been found by this court to have killed another human being.
- Therefore
- C) Jane has done something prohibited and is to be punished with 20 years imprisonment.

We can recognise this as a valid legal justification, and it is therefore tempting to reconstruct it as a simple syllogism, with 1) as a general premise of the form of a universally quantified conditional:

- 1b) Valid rule: For all x, if x kills a human being, then the punishment is 20 years.³

However, even though we recognise the legal argument A as a valid justification, we also know that 1b) is strictly speaking false, as there are numerous exceptions to this rule, from self-defence to insanity to duress. We also know that if the facts of the case had been different, or even maybe if Jane had a more competent advocate arguing one of these defences, the outcome might have been different. The task then becomes to give an account of A that on the one hand preserves its validity, and at the same time does not commit us to accept 1) as literally true.

One strategy to achieve this is to think of 1) and legal rules of its kind as essentially incomplete.⁴ Sometimes, this incompleteness is due to rhetorical or pragmatic constraints. We could for instance expand 1b) to a more realistic

- 1c) For all x, if x kills a human being AND does not act in self-defence AND does not act under duress AND is not insane, then the punishment is 20 years.

We would then also have to amend argument A by additional explicit premises that state that in this case, there was no self-defence, duress or insanity. This fix however comes with a number of problems of its own. First, our reconstruction of the argument now deviates considerably from what we found in our hypothetical court judgement. In particular, it adds three factual claims (no self-defence, duress, insanity) which not only are not stated as such in the decision, but also are unsupported by evidence. Furthermore, at least for the common law of England and Scotland, duress and insanity are exceptions to the homicide rule that were carved out by the judiciary, when prompted to do so in specific cases, where the application of the rule would have been repugnant for our collective sense of justice. Indeed, many exceptions in common law (and, so we would argue, in civilian jurisdiction as well) have entered the legal realm when a rule that had previously been deemed to hold without exceptions was causing unjust results when applied under conditions not foreseen, or foreseeable, to the legislator. Through a form of «sceptical meta-induction», we can therefore infer that it remains possible that also at a future time, under the facts of a future case, judges will again «recognise» a new exception that we are not yet aware of. It seems therefore not only impractical, but impossible in principle to give a complete account of 1).

The importance of incomplete rules for legal reasoning has been recognised within the law and AI community early on. Drawing in particular on the work by Pollock from general argumentation theory, researchers in legal AI developed formal representations of legal reasoning that could at the same time prove inferences such as A as valid, but without the need to reformulate the universal premise and list all possible exceptions explicitly

³ For the purpose of this paper, we do not distinguish practical syllogisms, that is syllogisms where the universal premise is a normative statement and the conclusion a mandate for action, from a syllogism simpliciter. For our purpose this distinction does not matter and could confuse the discussion.

⁴ See e.g. PRAKKEN, *Incomplete arguments in legal discourse: a case study*. Proceedings of the JURIX 2002 conference. IOS, Amsterdam 2002, pp. 93–102.

in its antecedents.⁵ This very active research field has created numerous formal approaches, often differing in technical details. What they tend to have in common though is to think of legal reasoning as an argumentation dialogue where arguments try to defeat other arguments or try to defend themselves against attackers. Such attacks are directed against premises of a specific kind – those that are defeasible. One way to defeat such a premise is to argue that the case under consideration is subject to a valid, if unstated, exception. Rather than reformulating the antecedents of the legal rule, a distinction is made in the formal representation of the implication itself, the «then» part of the rule.

In particular, these approaches distinguish between a strict implication, formally represented as « $A \rightarrow B$ » and interpreted like the implication of classical logic, and a weak or defeasible implication, represented as « $A \sim B$ » for « A weakly or defeasibly implies B ». Informally, we can think of the former as saying: «If A , then definitely B ». The latter could be rendered in a number of ways, each with subtly different shades of meaning (and possibly as a result subtly different logical behaviour),⁶ e.g. «If A , then normally B », «If A , typically B », «If A , then unless stated otherwise, B » or «If A , then assume B until proven otherwise».

This approach not only fulfils the adequacy criteria noted above – it accounts for the validity of the legal argument in the presence of implicit knowledge about possible exceptions – by using formal methods developed in general argumentation theory, it also embeds legal reasoning in a general theory of rational persuasion. The solution is not an ad hoc modification aimed at proving the rationality of legal decision making against counterexamples, but is motivated by observations across a wide range of disciplines and also everyday reasoning tasks. To use a canonical example:

If asked: «why do you think that this animal can fly», stating «Because it is a bird, and birds fly» is in many contexts a perfectly appropriate response. However, the implicit premise «All birds fly» is not generally true. If we learn that the animal in question is a penguin or ostrich, we may have to revisit our belief.

Epistemic contexts like this seem to be particularly suitable for an analysis that emphasises the temporary nature of our knowledge and beliefs, and unsurprisingly, evidential reasoning about the facts of a case has become one of the focal points of this type of analysis.⁷

It is however that very strength of the argumentation schema approach that also creates some of its limitations. Certain features that are crucial for legal reasoning qua legal reasoning cannot easily be captured in this way.

First, law very often self-reflexively turns on itself. Not only do we *use* strict and defeasible rules in legal reasoning, whether or not a specific law should be interpreted as applying strictly or allows for exceptions can in turn become subject of a legal debate that uses the rules of statutory interpretation as a form of argument. The «*Expressio unius est exclusio alterius*» rule of statutory interpretation for instance states that if the legislator has explicitly enumerated certain exceptions, we have to presume that these are the only exceptions permissible.⁸

Arguments that try to establish if a specific norm is best understood as introducing a strict implication or a weak implication, and if the latter, what exactly the meaning and scope of the « \sim » is, are difficult to represent adequately in argumentation logics. In logical terms, this is because they evoke the meta-theory of the logic

⁵ See e.g. PRAKKEN, Logical Tools for Modelling Legal Argument. A Study of Defeasible Reasoning in Law. Kluwer Dordrecht, 1997; PRAKKEN/REED/WALTON, Argumentation schemes and generalisations in reasoning about evidence. Proceedings of the Ninth International Conference on Artificial Intelligence and Law, Edinburgh 2003. ACM, New York 2003, pp. 32–34; BENCH-CAPON, Argument in artificial intelligence and law. Artificial Intelligence and Law 1997, Volume 5, pp. 249–261; GOVERNATORI, On the relationship between Carneades and defeasible logic. Proceedings of the 13th International Conference on Artificial Intelligence and Law, Pittsburgh 2011, ACM, New York 2011, pp. 31–40.

⁶ See also PRAKKEN/SARTOR, The three faces of defeasibility in the law. Ratio Juris 2004, Volume 17, pp. 118–139.

⁷ See e.g. PRAKKEN, Analysing reasoning about evidence with formal models of argumentation. Law, Probability & Risk 2004, Volume 4, pp. 33–50; BEX/PRAKKEN/REED/WALTON, Towards a formal account of reasoning about evidence: argumentation schemes and generalisations. Artificial Intelligence and Law 2003, Volume 11(2-3), pp. 125–165.

⁸ For an example see *Andrus v. Glover Const. Co.*, 446 U.S. 608, 616–617, 27 May 1980 (citing *Continental Casualty Co. v. United States*, 314 U.S. 527, 533, 5 January 1942)

under consideration, or maybe even the informal rules of translating a text into its logical form. The problem here is that we only know that « \leftrightarrow » is supposed to represent something like the natural language «it generally follows» when we look at the semantic meta-language of the formalism, or the informal explanation and motivation. Distinguishing object- and meta-language issues is of course a common and necessary technical device in formal logic. Unfortunately, in legal discourse object and meta-language issues often appear side by side, and any more or less artificial separation between the two means that some aspects of a legal decision can't be expressed any longer, or not expressed adequately.

The effect can be mitigated, of course. We can for instance introduce new argumentation schemata that take the rules of statutory interpretation as their antecedent. In the case of the «*exclusio alterius*» rule, we could then use this new premise to «undercut» an argument that reasons for making an exception. However, what gets lost in this analysis is that the issue is strictly about meaning and appropriateness of the « \leftrightarrow » in the legal rule under discussion. If it were possible to enrich our language so that it can talk directly about the meaning and extension of defeasibility (while avoiding the inconsistencies that often arise when object and meta language are conflated) then an entire class of legal arguments could be represented more faithfully.

The second aspect is best explained by comparing the penguin argument that tries to establish a fact about the world with the homicide example that establishes reasons for a specific course of legal action. In the former, defeasibility is closely connected to the concept of falsifiability. All our knowledge is provisional, it is always possible that we have to revise past inferences in the light of new observations. We can more or less arbitrarily terminate an argument at a given point in time to ask which position has been best defended, but nothing prevents us in principle to continue the process of revision indefinitely. Here, the logic of science differs crucially from the logic of law. Law seeks finality, the resolution of a conflict once and for all. In legal doctrine, this is expressed e.g. through the doctrine of *stare decisis*. Even were it the case that in a future decision, a new exception to the homicide prohibition is successfully argued, this will not normally lead to older cases that had been decided under the stricter reading of the rule being reopened.⁹ Similarly, once a decision in a court case is reached, the ability for the parties to raise new material facts or proffer new legal arguments will be severely limited, if permitted at all. If in our example, the defendant had failed to raise the issue of self-defence in the trial of first instance, then she will in Scots law be barred from raising it on appeal. Similarly, if her solicitor failed to convince the appeal court that the law should allow for a new exception, he will not be able to reopen the debate merely because he can think of a new argument for the exception at a later stage.

We argue that this finality of legal reasoning is so important to understand the epistemology of the law, and so deeply enshrined in procedural rules across legal systems from a huge variety of cultures, that it merits to be explicitly expressed in a formal representation of legal argumentation. In our example above, that means that the reasoning of the court is best understood if another explicit premise is added that states, roughly, «that that is it» – all arguments have been heard, all relevant exceptions been raised. In the approaches to defeasibility referenced above, this moment of closure is again only «visible» when one looks at the semantic meta-theory of the respective logic. We suggest that we can gain valuable insights about the nature of legal reasoning when in addition, it can be expressed explicitly in the object language of the argument itself.

In what follows, we will propose a way to express claims about the presence or absence of exceptions in the very premises that the argument at hand is using. It takes as its starting point Richard Holton's «that's it» clauses, building on its underlying intuitions while changing the formal representation.

2. «That's it, folks»

Let us return briefly to the homicide argument above.

⁹ In the US, this is made explicit in art 416.1489 of the Code of Federal Regulations.

So far we have shown how we can express some of the underlying issue through a logic with defeasibility.

1) As a valid rule of our legal system, killing another human being is *normally* prohibited and punished with 20 years imprisonment.

2) Jane has been found by this court to have killed another human being.

Therefore

C) Jane has done something prohibited and is to be punished with 20 years imprisonment.

Formally

1* $\forall x \text{ KilledHuman}(x) \sim \text{Punished}(x)$

2* $\text{KilledHuman}(\text{Jane})$

C* $\text{Punished}(\text{Jane})$

In the defeasible systems discussed above, 3 is derived from the premises because there was no argumentative move by the opposition that defeated the premise. To understand this though, we have to look at the meta-theory of the proposed logic that explicates the meaning of « \sim »

In a series of influential papers, Richard Holton has proposed a different reconstruction.¹⁰ According to him, ethical or legal arguments are subject to a *that's it* premise that states that locally, no relevant exceptions apply. Using such a *that's it* premise, we can use the following as a first stab at rendering our initial argument formally:

i. $\forall x(\text{KilledHuman}(x) \wedge \dots \wedge F_n(x) \wedge \text{That's it} \rightarrow \text{Punished}(x))$

ii. $\text{KilledHuman}(\text{Jane}) \wedge \dots F_n(\text{Jane})$

iii. That's it

(C) $\forall a$

We have now added a premise to the effect that not only do the stated exceptions to the general rule not apply in the case under consideration – what is expressed by the conjunction ... $F_n(x)$ – but that no further exceptions need to be considered. This is a direct formal representation of the «*exclusio alterius*» rule of statutory interpretation. Intuitively, the universal premise is now true. In every case of a potentially falsifying instance of (1), there is a defeating reason. That is, if a killing is not wrong, then there is a reason for why it is not – that it was an act of self-defence, for example. So in that case *That's it* was falsely asserted, and it follows logically that the relevant instance of (i) is true.

What does it take for the third premise of a *That's it* argument, the sentence *That's it* itself, to be true? Holton states the truth-conditions by quantifying over *That's it* arguments. Specifically, he introduces a relation of supersession between such arguments, and takes *That's it*, as it occurs in an argument, to be true just in case that argument is not superseded by any sound *That's it* argument.

What does it take, in turn, for a superseding argument to be sound? All its singular premises have to be true, of course. In addition, *That's it* has to be true as it occurs in the superseding argument, and likewise the universal premise, which contains *That's it*. This means that the truth of *That's it* as it occurs in one argument depends on the truth of *That's it* as it occurs in another argument. So Holton's account is not a recursive definition of *That's it* in terms of the semantic values of other expressions.

In many ways the intuition behind this analysis closely resemble the intuition that also informs the existing defeasible logics in legal AI: there too a proposition wins if there is no remaining successful attack move available (Holton's superseding arguments). The main difference is that this idea is now expressed in the object language.

What is the logical structure of *That's it* arguments? We will call «simple arguments» arguments the type of problematic argument with false universal premise with which this paper started. Every simple argument is

¹⁰ HOLTON, Principles and Particularism, Proceedings of the Aristotelian Society 2002, Supplementary Volume 67, pp. 191–209;
HOLTON, The Exception Proves the Rule. Journal of Political Philosophy 2010, Volume 18, pp. 369–388.

uniquely specified by the name of the agent or action (in our case, Jane), the predicates F_1, \dots, F_n occurring in the singular premise (In our example, killing someone) and a predicate V occurring in the conclusion (in our case, the 20 years punishment). Let $s(a, F_1, \dots, F_n, V)$ refer to the simple argument from (1) and (2) to (C) displayed above, with «V» standing for the Verdict in the conclusion. Likewise, every *that's it* argument is uniquely specified by such a name and such predicates. We shall use $t(a, F_1, \dots, F_n, V)$ to refer to the *That's it* argument from (i), (ii) and (iii) to (C) displayed above.

With this notation, Holton's account of supersession among *That's it* arguments can be defined as follows:

Definition 1 : $t(a, F_1, \dots, F_n, V)$ supersedes $t(b, G_1, \dots, G_m, V')$ if

(i) $F_1(a) \wedge \dots \wedge F_n(a)$ entails $G_1(b) \wedge \dots \wedge G_m(b)$

(ii) $G_1(b) \wedge \dots \wedge G_m(b)$ does not entail $F_1(a) \wedge \dots \wedge F_n(a)$

(iii) $V(a)$ and $V(b)$ are incompatible

(Typically, the second condition is only satisfied if $a = b$.)

Supersession is thus both a matter of the logical relations between the predicates involved in the arguments, and of the properties had by the individuals involved. We obtain a plausible example by taking a to be an individual who killed in self-defence, $b = a$, and F_1, F_2, V and V' to be «killed», «acted in self-defence», «is guilty», and «is not guilty», respectively.

Holton then specifies the truth-conditions of *That's it* in terms of supersession. On his account, briefly sketched above already, *That's it* means something different in different arguments. Specifically, *That's it* says that «*this* argument is not superseded by any sound argument», where the demonstrative «*this*» picks out the very argument in which it occurs.

While Holton uses such a demonstrative to state the truth-conditions, this is not essential, since we can pick out the relevant argument by descriptive means. The notation we introduced above helps us do that. The account can then be formulated as follows:

Constraint 1 : An occurrence of *That's it* in the argument $t(a, F_1, \dots, F_n, V)$ is true if there is no sound argument superseding $t(a, F_1, \dots, F_n, V)$.

Since definitions need to be non-circular, on the orthodox view, we are calling this biconditional a «constraint». It certainly does constrain the truth-values that occurrences of *That's it* receive in arguments, given an interpretation of the rest of the language. What we do not know, before further investigation, is whether there is always a unique assignment of truth-values to *That's it* that satisfies the constraint. If (and only if) it did, the constraint would arguably deserve to be called a «definition».

At this point, we depart from Holton's formal account while trying to remain true to the spirit of his argument. To ensure that premise 3), the *that's it* proposition, is not rendered false just because there have been some cases where self-defence was accepted as an exception, but expresses instead the idea that *in the case under consideration*, self-defence has not been raised, we have to be able to quantify inside the *That's it* proposition. To enable us to do so, we shall introduce a family of *That's it* predicates. Given a stock of basic predicates in the language, these can be taken to be complex predicates. Specifically, whenever F_1, \dots, F_n as well as V are basic predicates, then $T_v^{F_1, \dots, F_n}$ is a complex predicate. To a first approximation, $T_v^{F_1, \dots, F_n}$ is true of an agent x if and only if F_1, \dots, F_n cover between them all the facts legally or procedurally relevant to whether V applies to x . In terms of argument supersession, it is natural to modify Constraint 1 as follows:

Constraint 2. $T_v^{F_1, \dots, F_n}$ is true of x if there is no sound argument superseding $t(x, F_1, \dots, F_n, V)$.

This account allows us to quantify into the *That's it* clause. On this revised account, a *That's it* argument has the following syntactic form:

(1) $F_1(a) \wedge \dots \wedge F_n(a)$

(2) $\forall x (F_1(x) \wedge \dots \wedge F_n(x) \wedge T_v^{F_1, \dots, F_n}(x) \rightarrow V(x))$

- (3) $T_V^{F_1, \dots, F_n}(a)$
 (C) $V(a)$

From now on, we will use $\langle t(a, F_1, \dots, F_n, V) \rangle$ to refer to such an argument, rather than one in which *That's it* contains no quantifiable variable.

Such arguments are easily seen to be valid:

$$F_1(a) \wedge \dots \wedge F_n(a) \wedge T_V^{F_1, \dots, F_n}(a) \rightarrow V(a)$$

follows from (2) by universal instantiation, and together with (1) and (3) entails (C).

The modification avoids the problem of false universal premises. The instance of the scheme concerning Jane's killing looks as follows:

- (1) Jane killed someone.
 (2) For all x (x killed someone $\wedge T_{\text{guilty}}^{\text{killed}}(x) \rightarrow x$ is guilty)
 (3) $T_{\text{guilty}}^{\text{killed}}(\text{Jane})$
 (C) Jane is guilty.

The problem with the earlier version of this argument was that if the *that's it* premise (3) is true, then premise (2) is false. Or in other words, as soon as an exception to the general rule has been argued successfully *anywhere*, it also has to be considered in the specific case at hand. This allows us to understand the evolution of the law and its tendency to add new exceptions to established rules over time. But crucially it does not account for the observation of stasis that we described above, and that legal decisions remain valid even if an exception is argued in a later case, or if an available exception was not argued. With other words, it did not really resolve the problem that often, legal decisions seem to be based on a prima facie false universal rule. In this version, in contrast, premise (2) may well be true even if there are other cases, such as say the case of a John Doe who killed in self-defence:

John is not a counterinstance any longer because $T_{\text{guilty}}^{\text{killed}}(\text{John})$ is in this case false – the sound argument $t(\text{John}, \text{killed}, \text{in self-defence}, \text{not guilty})$ supersedes $t(\text{John}, \text{killed}, \text{guilty})$.

3. Discussion

We started our analysis with a familiar observation: legal arguments often seemed to be based on general premises that are false at face value. This problem has long been recognised in legal AI research, with argumentation frameworks in the tradition of Pollock and Dung being widely used to capture the way in which claims in a legal discourse can be defeated by attacking their support through claiming a relevant exception.

While a significant class of legal arguments can be represented in this way, these approaches struggle when the very question of legal interpretation, the validity and legitimacy of exceptions and the degree if any, of defeasibility is raised within a court judgement. The reason for this gap, so we argued, is that these argumentation frameworks gain much of their explanatory value in their respective meta-theories only. We understand *how* they capture the defeasible nature of certain legal rules when looking at the meta-definition of what makes an argument successful, but we cannot replicate this insight within the legal discourse directly. Legal argumentation however is in a crucial sense self-referential. The formal rules of substantive law always co-exist with the procedural rules of correct legal interpretation, which can at any time themselves become the object of an argument. Court decisions can, at the same time, *use* a legal rule to justify a conclusion and *mention* that very same rule to argue for a specific interpretation of it.

The existence of these procedural rules and rules of interpretation and their interaction with substantive arguments is so central to the nature of legal reasoning, not just in some legal systems but as a universal feature of modern law, that the ability to represent them adds a valuable tool in our analytical toolkit. One role of some of these procedural rules, again as a universal or near universal feature of legal systems, is that they impose a finality on the reasoning process that is missing in other fields of defeasible reasoning. At some point, there

are rules that say «and that's been it», even though in theory, other valid defeaters may exist. This notion of finality is again so central to the structure of legal systems that it is desirable to be able to represent, analyse and discuss it.

We argued therefore for an expansion of the formal repertoire that re-introduces some of the intuitions that inform argumentation frameworks into the object language. We identified Richard Holton's «that's it» propositions as a prima facie promising candidate for such a task, that closely mirrors some important procedural rules of legal interpretation. We argued however that to reap the full benefit of this approach, a richer representation that allows *that's it* to be in the scope of a quantifier was needed.

The resulting logic is capable of handling interesting and important examples of legal reasoning. In some respects however, it is in the version introduced above not capable of correctly analysing a class of legal arguments that can be represented convincingly in argumentation frameworks. Our approach allows us to explain why it can be the case that John is acquitted or murder while Jane is convicted, even though the legal rule that justifies both decisions simply states that killing is prohibited: John successfully pleaded an exception, self-defence, that Jane did not. However, what about the situation where John did act in self-defence (so that the exception is in principle available to him) and he also raises it in a procedurally correct way, and his defence is nonetheless unsuccessful because his actions were deemed excessive? Intuitively, we have then an exception from an exception which, if argued by the prosecution, reinstates the original guilty verdict.

Our proposal as formulated cannot adequately represent this situation. That it is not just a trivial extension of the general idea of defeasibility should not come as a surprise though. As the discussion between PRAKKEN and HORTY has shown for the way in which argumentation frameworks handle the same issue, it is not entirely clear, or at least debatable, what exactly happens if an initial rule is reinstated when a possible defeater is itself defeated.¹¹ We agree with PRAKKEN that reinstatement is a real phenomenon and that the way in which argumentation systems represent them is sound. For our purpose, this means that we have to refine the proposal above in some minor ways, though the technical details will have to be left for a future paper.

¹¹ HORTY, Argument construction and reinstatement in logics for defeasible reasoning. *Artificial intelligence and Law* 2001, Volume 9, pp. 1–28; PRAKKEN, Intuitions and the modelling of defeasible reasoning: some case studies. arXiv preprint cs/0207031 2002, <https://arxiv.org/abs/cs/0207031> (accessed 12 January 2016).