



**HAL**  
open science

## Representing Pure Nash Equilibria in Argumentation

Bruno Yun, Srdjan Vesic, Nir Oren

► **To cite this version:**

Bruno Yun, Srdjan Vesic, Nir Oren. Representing Pure Nash Equilibria in Argumentation. *Argument and Computation*, 2020, 13 (2), pp.195-208. 10.3233/AAC-210007 . lirmm-03039438

**HAL Id: lirmm-03039438**

**<https://hal-lirmm.ccsd.cnrs.fr/lirmm-03039438v1>**

Submitted on 3 Dec 2020

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

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

# Representing Pure Nash Equilibria in Argumentation

Bruno YUN <sup>a,1</sup>, Srdjan VESIC <sup>b</sup> and Nir OREN <sup>a</sup>

<sup>a</sup>*Department of Computing Science, University of Aberdeen, Scotland*

<sup>b</sup>*CRIL Lens, University of Artois, France*

**Abstract.** In this paper we describe an argumentation-based representation of normal form games, and demonstrate how argumentation can be used to compute pure strategy Nash equilibria. Our approach builds on Modgil’s Extended Argumentation Frameworks. We demonstrate its correctness, prove several theoretical properties it satisfies, and outline how it can be used to explain why certain strategies are Nash equilibria to a non-expert human user.

**Keywords.** Argumentation, Game theory, Nash equilibrium, Pure strategy

## 1. Introduction

Game theory studies how multiple rational decision-makers should act given interactions between their strategies, and preferences over the resultant outcomes. Game theory has been applied to myriad fields [8]. Within game theory, decision-makers (referred to as players), their strategies, preferences and outcomes are represented within a game, and the solutions to a game identify some form of rational outcome. One such solution concept is that of a *dominant* strategy, where a player has a strategy or a set of strategies that will always result in the best outcome for them, regardless of what other players do. However, such dominant strategies often do not exist. In this work, we consider instead the notion of a *Nash equilibrium*, which identifies optimal strategies given that other players also pursue their own optimal strategies. Such Nash equilibria therefore represent a form of best response, and provide a well understood solution concept in game theory. However, finding Nash equilibria is computationally difficult, and it is sometimes difficult for a non-expert to understand why a given strategy is (or is not) a Nash equilibrium. We believe that by providing an argumentation-based representation of games, dialogues can be used to explain a Nash equilibrium to such non-experts. While work such as [6] has considered game theory in the context of ABA, to our knowledge, this work is the first to link abstract argumentation and Nash equilibria. We consider only so-called *pure strategies* for *normal form games* and intend to relax this restriction in future work.

The remainder of the paper is structured as follows. In Section 2, we provide a brief overview of argumentation and game-theory concepts necessary to understand our article. In Section 3, we describe how a normal form game can be encoded using argumen-

<sup>1</sup>Corresponding Author: E-mail: bruno.yun@abdn.ac.uk

tation. Section 4 examines some formal properties of our approach. Lastly, we discuss related and future work in Section 5 before concluding.

## 2. Background

We begin by providing the necessary background in game theory and argumentation required for the rest of the paper.

### 2.1. Game Theory

In this paper, we use the usual *normal form* for games [13].

**Definition 1. (Normal Game)** A (normal) game is  $G = (Ag, Ac, Av, Ou, Ef, \leq)$  where  $Ag = \{0, 1, \dots, n\}$  is a finite set of players;  $Ac$  is a finite set of strategies;  $Av = [Ac_0, \dots, Ac_n]$  with  $Ac_i \subseteq Ac$  denoting the strategies available to  $i$ ;  $Ou = \{o_0, \dots, o_m\}$  is a set of possible outcomes;  $Ef : Ac^n \rightarrow Ou^n$  captures the consequences of the joint strategies for each player; and  $\leq = [\leq_0, \dots, \leq_n]$  with  $\leq_i \subseteq Ou \times Ou$  denoting the preference relation for player  $i$ .

The notation  $o_k \leq_i o_l$  means that player  $i$  prefers outcome  $o_l$  to  $o_k$ . As commonly done, we write  $o_i <_i o_j$  iff  $o_i \leq_i o_j$  and  $o_j \not\leq_i o_i$ <sup>2</sup>. A *pure strategy profile*  $S$  is a tuple containing one strategy from each player in the game. The set of all such pure strategy profiles is  $S_G = \prod_{i \in Ag} Ac_i$ , and represents one joint strategy of all players. A *partial strategy profile* is a tuple containing a single strategy for a subset of the players. Given any pure strategy profile  $S = [s_0, \dots, s_n]$ , we write  $S_{-i}$  to denote the *partial strategy profile*  $[s_0, \dots, s_{i-1}, \emptyset, s_{i+1}, \dots, s_n]$ , where the strategy for player  $i$  is not specified. We then write  $S_{-i} \oplus s_i$  to denote strategy profile  $S$ . With a slight abuse of notation, for any  $S, S' \in S_G$  we write that  $S \leq_i S'$  iff  $Ef(S)_i \leq_i Ef(S')_i$ <sup>3</sup>.

**Example 1.** Let us consider the stag hunt game  $G = (\{0, 1\}, Ac, Av, Ou, Ef, \leq)$ , where  $Ac = \{stag, hare\}$ ,  $Av = [Ac, Ac]$ ,  $Ou = \{4, 3, 2, 1\}$ ,  $\leq$  is the standard less than relation over numbers. Table 1a graphically illustrates this game in normal form, and specifies  $Ef$ . For example, the tuple  $(1, 3)$  in the column “hare” and row “stag” means that  $Ef([stag, hare]) = (1, 3)$ . Given the pure strategy profile  $S = [stag, hare]$ ,  $S_{-0} = [\emptyset, hare]$  and  $S_{-0} \oplus hare = [hare, hare]$ . Here  $[stag, hare] \leq_0 [hare, hare]$  because  $(1, 3)_0 \leq_0 (2, 2)_0$  but  $[hare, hare] \leq_1 [stag, hare]$ .

In asking why a player should pursue a some strategy, we must take into account the strategies of others. A *Nash equilibrium* is the best response a player can make given optimal play by all other players.

**Definition 2.** Let  $G = (Ag, Ac, Av, Ou, Ef, \leq)$ , we say that  $S \in S_G$  is a *Nash equilibrium* if for every  $i \in Ag$  and for any strategy  $s \in Ac_i$ , it holds that  $S_{-i} \oplus s \leq_i S$ .

<sup>2</sup>We assume that our preferences are acyclic. I.e., if  $a <_i b <_i c$  then  $c \not\leq_i a$ .

<sup>3</sup>The notation  $Ef(S')_i$  means the  $i$ -th element of  $Ef(S')$ .

	Player 1			Player 1			
		<i>stag</i>	<i>hare</i>		<i>heads</i>	<i>tails</i>	
Player 0	<i>stag</i>	4	3	Player 0	<i>heads</i>	-1	1
	<i>hare</i>	1	2		<i>tails</i>	1	-1
		4	1			-1	1
		3	2			-1	1

(a) Stag Hunt

(b) Matching pennies.

**Table 1.** Two games in normal form.

A simple algorithm to identify all Nash equilibrium in the presence of pure strategies involves iterating through every player and identifying the best strategy profile (in terms of  $Ef$  for that player) given all other players' possible joint strategies. Any strategy profile which all players consider best is then a Nash equilibrium.

Given a game in normal form, the above algorithm involves – for a two player game – scanning down each column and marking the best strategy for the row player, and then doing the same for each row marking the best strategy for the column player. Each cell marked for both players is a Nash equilibrium. In the remainder of this paper, we show an argumentation-based alternative.

**Example 2** (Cont'd). *There are two Nash equilibria in the stag hunt game:  $[stag, stag]$  and  $[hare, hare]$ . The strategy profile  $[stag, stag]$  is a Nash equilibrium because  $[hare, stag] \leq_0 [stag, stag]$  and  $[stag, hare] \leq_1 [stag, stag]$ . Similarly,  $[hare, hare]$  is also a Nash equilibrium as  $[stag, hare] \leq_0 [hare, hare]$  and  $[hare, stag] \leq_1 [hare, hare]$ .*

## 2.2. Argumentation

We encode normal form games in terms of arguments and attacks by building on Modgil's Extended Argumentation Frameworks (EAF) [11].

**Definition 3.** *An Extended Argumentation Framework is a triple  $\langle \mathbb{A}, \mathbb{C}, \mathbb{D} \rangle$  where  $\mathbb{A}$  is a set of arguments,  $\mathbb{C} \subseteq \mathbb{A} \times \mathbb{A}$ ,  $\mathbb{D} \subseteq \mathbb{A} \times \mathbb{C}$  and if  $(z, (x, y)), (z', (y, x)) \in \mathbb{D}$  then  $(z, z'), (z', z) \in \mathbb{C}$ .*

**Definition 4** (Defeat). *Let  $\mathcal{AS} = (\mathbb{A}, \mathbb{C}, \mathbb{D})$  be an EAF,  $x, y \in \mathbb{A}$  and  $Y \subseteq \mathbb{A}$ . We say that  $y$  defeats  $x$  w.r.t.  $Y$ , denoted  $y \rightarrow_Y x$  iff  $(y, x) \in \mathbb{C}$  and there is no  $z \in Y$  s.t.  $(z, (y, x)) \in \mathbb{D}$ .*

**Definition 5** (Argumentation semantics). *Let  $\mathcal{AS} = (\mathbb{A}, \mathbb{C}, \mathbb{D})$  be an EAF and  $E \subseteq \mathbb{A}$ . We say that:*

- *$E$  is conflict-free iff for every  $x, y \in E$ , if  $(y, x) \in \mathbb{C}$  then  $(x, y) \notin \mathbb{C}$ , and there exists  $z \in E$  s.t.  $(z, (y, x)) \in \mathbb{D}$ .*
- *$x \in \mathbb{A}$  is acceptable w.r.t.  $E$  iff for every  $y \in \mathbb{A}$  s.t.  $y \rightarrow_E x$ , there exists  $z \in E$  s.t.  $z \rightarrow_E y$  and there exists  $R_E = \{x_1 \rightarrow_E y_1, \dots, x_n \rightarrow_E y_n\}$  s.t. for every  $i \in \{1, \dots, n\}$ ,  $x_i \in E$ ,  $z \rightarrow_E y \in R_E$  and for every  $x_j \rightarrow_E y_j \in R_E$ , for every  $y'$  s.t.  $(y', (x_j, y_j)) \in \mathbb{D}$ , there exists  $x' \rightarrow_E y' \in R_E$*
- *$E$  is an admissible extension iff every argument in  $E$  is acceptable w.r.t.  $E$*

- $E$  is a preferred extension iff  $E$  is a maximal (w.r.t.  $\subseteq$ ) admissible extension
- $E$  is a stable extension iff for every  $y \notin E$ , there exists  $x \in E$  such that  $x \rightarrow_E y$ .

We will use the notation  $Ext_s(\mathcal{AS})$  (resp.  $Ext_p(\mathcal{AS})$ ) to denote the set of all stable (resp. preferred) extensions.

### 3. Argumentation-based approach for games

We consider an argumentation framework with multi-level arguments. At the base level, we consider all possible strategy profiles as arguments. Since only a single strategy profile can ever occur (as players execute one set of strategies in the interaction), every argument at this level must attack every other argument. We refer to such arguments as *game-based arguments*, and note that they are equivalent to pure strategy profiles.

**Definition 6** (Game-based argument). *Let  $G = (Ag, Ac, Av, Ou, Ef, \leq)$  be a game, a game-based argument (w.r.t.  $G$ ) is a pure strategy profile  $S \in S_G$ .*

The set of all game-based arguments for a game  $G$  is denoted by  $\mathcal{A}_g(G)$ .

Next, we introduce *preference arguments*. Intuitively, these can be interpreted as statements of the form: “Given that the other players are performing a given set of strategies, the remaining player’s preferred strategy should be playing  $x$ ”.

**Definition 7** (Preference argument). *Let  $G = (Ag, Ac, Av, Ou, Ef, \leq)$  be a game,  $S \in S_G$  be a pure strategy profile and  $i \in Ag$ . A preference argument (w.r.t.  $G$ ) is a tuple  $(S_{-i}, s)$ , where  $s \in Ac_i$ .*

The set of preference arguments for a game  $G$  is denoted by  $\mathcal{A}_p(G)$ . A *cluster* of preference arguments is a maximal set of preference arguments sharing the same partial strategy profile.

Finally, we introduce *valuation arguments*, which can be interpreted as statements of the form: “Given that the other players are performing a given set of strategies, it is the case that the outcome of strategy  $s$  is better than the outcome of strategy  $s'$  for the remaining player”.

**Definition 8** (Valuation argument). *Let  $G = (Ag, Ac, Av, Ou, Ef, \leq)$  be a game,  $i \in Ag$ ,  $(S_{-i}, s), (S_{-i}, s') \in \mathcal{A}_p(G)$  be two preference arguments and  $S_{-i} \oplus s' <_i S_{-i} \oplus s$ . A valuation argument (w.r.t.  $G$ ) is the pair  $(S_{-i}, s' < s)$ .*

The set of valuation arguments for a game  $G$  is denoted by  $\mathcal{A}_v(G)$ .

**Example 3** (Cont’d). *The sets of game-based, preference and valuation arguments w.r.t.  $G$  are shown in Table 2. The argument  $a_1$  represents the case where player 0 chooses to hunt a stag and player 1 chooses to hunt a hare. The argument  $a_9$  represents the argument: “Given that player 0 chooses to hunt a hare, player 2’s preferred strategy should be to hunt a stag”. The argument  $a_{16}$  represents the argument: “Given that player 1 chooses to hunt a hare, the outcome of hunting a hare is better than the outcome of hunting a stag for player 0”.*

Game-based arguments	Preference arguments	Valuation arguments
$a_1 = [stag, hare]$	$a_5 = ([stag, \emptyset], stag)$	$a_{13} = ([stag, \emptyset], stag > hare)$
$a_2 = [stag, stag]$	$a_6 = ([stag, \emptyset], hare)$	$a_{14} = ([\emptyset, stag], stag > hare)$
$a_3 = [hare, stag]$	$a_7 = ([\emptyset, stag], stag)$	$a_{15} = ([hare, \emptyset], hare > stag)$
$a_4 = [hare, hare]$	$a_8 = ([\emptyset, stag], hare)$	$a_{16} = ([\emptyset, hare], hare > stag).$
	$a_9 = ([hare, \emptyset], stag)$	
	$a_{10} = ([hare, \emptyset], hare)$	
	$a_{11} = ([\emptyset, hare], stag)$	
	$a_{12} = ([\emptyset, hare], hare)$	

**Table 2.** Arguments for the stag hunt game

We now turn our attention to attacks. We note that preference and valuation arguments provide reasons why one argument should not attack another, and therefore introduce not only attacks between arguments, but also attacks on attacks.

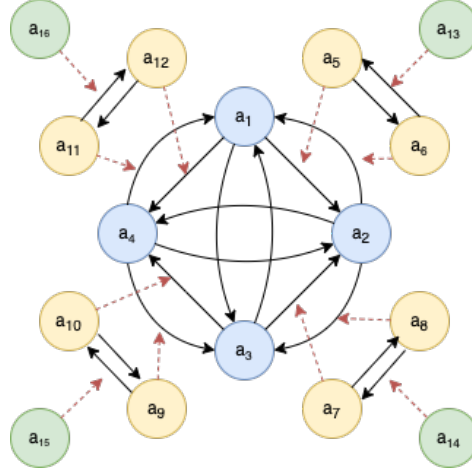
**Definition 9 (Attack).** For a game  $G = (Ag, Ac, Av, Ou, Ef, \leq)$ ,  $a_1, a_2 \in \mathcal{A}_g(G)$ ,  $a_3 = (S_1, s_2), a_4 = (S_3, s_4) \in \mathcal{A}_p(G)$  and  $a_5 = (S_5, s_6 > s_7) \in \mathcal{A}_v(G)$ . We say that:

- $a_1$  attacks  $a_2$ , denoted  $(a_1, a_2) \in \mathcal{C}_r(G)$ , iff  $a_1 \neq a_2$ .
- $a_3$  attacks  $a_4$ , denoted  $(a_3, a_4) \in \mathcal{C}_p(G)$ , iff  $S_1 = S_3$  and  $s_2 \neq s_4$ .
- $a_3$  attacks  $(a_1, a_2) \in \mathcal{C}_r(G)$ , denoted by  $(a_3, (a_1, a_2)) \in \mathcal{C}_u(G)$ , iff there exists  $s \in Ac$  such that  $S_1 \oplus s = a_1$  and  $S_1 \oplus s_2 = a_2$ .
- $a_5$  attacks  $(a_3, a_4) \in \mathcal{C}_p(G)$ , denoted by  $(a_5, (a_3, a_4)) \in \mathcal{C}_v(G)$ , iff  $S_5 = S_3$ ,  $s_6 = s_4$  and  $s_7 = s_2$ .

The first attack captured within Definition 9 is between every two distinct game-based arguments. As each player has to choose exactly one strategy, different strategy profiles are clearly incompatible. The second bullet point represents attacks between preference arguments. In the stag hunt example for instance,  $a_5$  attacks  $a_6$  (and vice-versa) because in the event of player 0 hunting a stag, player 1 can either hunt a stag or a hare. The third type of attack captures attacks from preference arguments to attacks between game-based arguments. Within the stag hunt,  $a_5$  attacks  $(a_1, a_2)$  because  $a_5$  states that it is preferable for player 1 to hunt a stag when player 0 is also hunting a stag. Note that in general, the preference argument  $(S_1, s_2)$  attacks *all* attacks against the game-based argument  $S_1 \oplus s_2$  coming from any other game-based arguments of the form  $S_1 \oplus s'$ , for any  $s' \in Ac$  such that  $s' \neq s_2$ . The last type of attack captures attacks from valuation arguments to attacks between preference arguments. Returning to the stag hunt,  $a_{13}$  attacks  $(a_6, a_5)$  as  $a_{13}$  states that the strategy “hunt a stag” is better than the strategy “hunt a hare” for player 1 when player 0 is hunting a stag.

The arguments and attacks induce a very specific type of extended argumentation framework, where object-level (game-based) arguments have their attacks attacked by meta-arguments (preference arguments) at level one, and where attacks between these meta-arguments are attacked by meta-arguments at level two (valuation arguments).

**Definition 10 (Argumentation framework).** Let  $G$  be a game. The argumentation framework corresponding to  $G$  is the tuple  $\mathcal{AS}_G = (\mathbb{A}, \mathbb{C}, \mathbb{D})$  where  $\mathbb{A} = \mathcal{A}_g(G) \cup \mathcal{A}_p(G) \cup \mathcal{A}_v(G)$ ,  $\mathbb{C} = \mathcal{C}_r(G) \cup \mathcal{C}_p(G)$  and  $\mathbb{D} = \mathcal{C}_u(G) \cup \mathcal{C}_v(G)$ .



**Figure 1.** Argumentation graph corresponding to stag hunt game

**Example 4** (Example 3 Contd). *Figure 1 represents the game-based, preference and valuation arguments of  $G$  using blue, yellow and green nodes respectively. The attacks between arguments ( $\mathbb{C}$ ) and on attacks ( $\mathbb{D}$ ) are represented using solid black arrows and dashed red arrows respectively.*

For our framework to be an EAF, it must satisfy some constraints, as described in [10], and we can easily show that this is the case.

**Proposition 1.** *Let  $G$  be a game and  $\mathcal{AS}_G = (\mathbb{A}, \mathbb{C}, \mathbb{D})$  be the corresponding argumentation framework, it holds that if  $(z, (x, y)), (z', (y, x)) \in \mathbb{D}$  then  $(z, z'), (z', z) \in \mathbb{C}$ .*

*Proof.* There are only two types of attacks on attacks: (1) attacks coming from valuation arguments to attacks between preference arguments and (2) attacks coming from preference arguments to attacks between game-based arguments. In the rest of this proof, we prove that Proposition 1 is satisfied for the two types of attacks on attacks.

- Considering (1), for a fixed partial strategy profile  $S_i$ , and fixed strategies  $s_j, s_k \in Ac$ , there is exactly one (or no) valuation argument of the form  $(S_i, s_j > s_k)$  or  $(S_i, s_k > s_j)$ . As a result, the condition in Proposition 1 is trivially satisfied for attacks coming from valuation arguments.
- We now study the case (2) and show that Proposition 1 is also satisfied for attacks coming from preference arguments on attacks between game-based arguments. Assume that  $(a_3, (x, y)), (a_4, (y, x)) \in \mathbb{D}$ , where  $a_3 = (S_1, s_2)$ ,  $a_4 = (S_1, s_4)$ ,  $x = S_1 \oplus s_4$  and  $y = S_1 \oplus s_2$ . By Definition 9,  $s_2 \neq s_4$  thus  $(a_3, a_4), (a_4, a_3) \in \mathcal{C}_p(G) \subseteq \mathbb{C}$ .

□

Since – given Proposition 1 – our argumentation system is an EAF, we can use EAF semantics to evaluate it.

**Example 5** (Example 4 Contd). *In our running example,  $a_5$  defeats  $a_6$  w.r.t.  $\mathbb{A}$  as  $(a_5, a_6) \in \mathbb{C}$  and there is no argument  $z \in \mathbb{A}$  such that  $(z, (a_5, a_6)) \in \mathbb{D}$ . However,  $a_6$  does not defeat  $a_5$  w.r.t.  $\mathbb{A}$  because  $(a_{13}, (a_6, a_5)) \in \mathbb{D}$ . All extensions contain arguments  $\{a_{16}, a_{15}, a_{14}, a_{13}, a_{12}, a_{10}, a_7, a_5\}$ , while one preferred extension contains  $\{a_2\}$  and the other contains  $\{a_4\}$ .*

#### 4. System Properties

Having described our system, we now consider its properties. The most important result we seek to show is the correspondence between argumentation semantics and Nash equilibria, and we begin by laying the groundwork for this. We then consider how many arguments will be generated for an arbitrary normal form game.

We begin by considering which preference arguments will appear in a preferred extension. This result is used in later proofs.

**Lemma 1.** *Let  $G = (Ag, Ac, Av, Ou, Ef, \leq)$  be a game, and  $\mathcal{AS}_G$  be the corresponding AS. For each preferred extension  $E$  of  $\mathcal{AS}_G$ , for each cluster  $C$  of preference arguments, there exists a unique argument  $c \in C$  such that  $c \in E$ .*

*Proof.* Assume a partial strategy profile  $S = [s_0, \dots, s_{i-1}, \emptyset, s_{i+1}, s_n]$  and the corresponding cluster of preference arguments  $C$ . Because our preferences are acyclic, we know that there exists a strategy  $s^*$  such that for every  $s \in Ac_i$ ,  $S \oplus s \leq_i S \oplus s^*$ . From the definition of the valuation argument, there are no valuation arguments attacking the attacks from the preference argument  $(S, s^*)$  to other preference arguments. As a result, we conclude that  $(S, s^*)$  is in a preferred extension  $E$  and that all the other preference in  $C$  are not  $E$ . Moreover, you need to choose one of such arguments from the cluster  $C$  for each preferred extension to satisfy the maximality condition of the semantics.  $\square$

Next, we show that if there is a preferred extension with game-based arguments, then each such extension has exactly one game-based argument.

**Lemma 2.** *If any preferred extension of  $\mathcal{AS}_G$  contains a game-based argument, then it contains exactly one game-based argument.*

*Proof.* Let  $E$  be a preferred extension containing game-based arguments. We prove by contradiction that it is not possible for  $E$  to have more than one game-based argument. Assume that  $E$  contains two game-based arguments  $a_1$  and  $a_2$ . By definition of the attack relation, there is a symmetric attack between  $a_1$  and  $a_2$ . Hence there must exist two preference arguments  $p_3$  and  $p_4$  with  $(p_3, (a_1, a_2)), (p_4, (a_2, a_1)) \in \mathbb{D}$  and  $(p_3, p_4), (p_4, p_3) \in \mathbb{C}$ . It is not possible for both  $(p_4, p_3)$  and  $(p_3, p_4)$  to be attacked by valuation arguments as this would require an inconsistency or cycle in  $\leq$ . By this observation,  $E$  contains only  $p_3$  or  $p_4$ . Hence,  $\{a_1, a_2\}$  is not conflict-free, contradiction.  $\square$

We now show that a game-based argument which is not a Nash equilibrium will not appear in any preferred extension of the associated argumentation system.



**Lemma 3.** *Let  $G = (Ag, Ac, Av, Ou, Ef, \leq)$  be a game, and  $\mathcal{AS}_G$  be the corresponding AS. If  $S \in S_G$  such that  $S$  is not a Nash equilibrium then for every preferred extension  $E$ ,  $S \notin E$ .*

*Proof.* Assume there is a non-Nash equilibrium game-based argument  $S' = [s'_0, \dots, s'_n]$  in a preferred extension  $E$ . Then, from Lemma 2,  $E$  does not contain any other game-based arguments. Since  $S'$  is not a Nash equilibrium, there exists  $i \in Ag$  and  $s \in Ac_i$  such that  $S'_{-i} \oplus s'_i <_i S'_{-i} \oplus s$ . In the rest of this proof, we consider the strategy  $s^*$  such that for every  $s \in Ac_i$ ,  $S'_{-i} \oplus s \leq_i S'_{-i} \oplus s^*$ . By definition, the attack from  $S'$  to  $S'_{-i} \oplus s^*$  is attacked by the preference argument  $(S'_{-i}, s^*)$ . Moreover, the preference argument  $(S'_{-i}, s^*)$  attacks all the other preference arguments  $(S'_{-i}, s')$ , where  $s' \in Ac_i$  and  $s' \neq s$ . By definition of the valuation arguments, none of the attacks from  $(S'_{-i}, s^*)$  to those other preference arguments is defeated. As a result, we conclude that there is preferred extension that contains  $(S'_{-i}, s^*)$ . Let  $s^+ = \{s \in Ac_i \mid S'_{-i} \oplus s \leq_i S'_{-i} \oplus s^* \text{ and } S'_{-i} \oplus s^* \leq_i S'_{-i} \oplus s\}$ , we can conclude that there is at least one argument  $(S'_{-i}, s_o), s_o \in s^+$  in  $E$  (Lemma 1) and  $(S'_{-i}, s_o)$  attacks the attack from  $S'$  to  $S'_{-i} \oplus s_o$ , contradiction.  $\square$

**Corollary 1.** *Let  $G = (Ag, Ac, Av, Ou, Ef, \leq)$  be a game, and  $\mathcal{AS}_G$  be the corresponding AS. If  $E$  is a preferred extension that contains a game-based argument  $S$ , then  $S$  is a Nash equilibrium.*

In the next proposition, we show that if there is only one preferred extension that contains a game-based argument, then there is an equivalence between preferred and stable extensions.

**Proposition 2.** *Let  $G$  be a game and  $\mathcal{AS}_G = (\mathbb{A}, \mathbb{C}, \mathbb{D})$  be the corresponding argumentation framework. If  $E \in Ext_p(\mathcal{AS}_G)$  and  $E \cap Ag(G) \neq \emptyset$  then  $E \in Ext_s(\mathcal{AS}_G)$ .*

*Proof.* We show that if a preferred extension possesses a game-based argument, then it is also a stable extension. Assume  $E$  contains a single game-based argument. By Lemma 2,  $E$  contains exactly one game-based argument. Therefore, all game-based arguments not in the extension are defeated by the game-based argument within the extension with respect to  $E$ , meaning that the game-based argument is a member (at the game-based level) of the stable extension.  $\square$

It may seem intuitive that the preferred and stable extension should coincide where multiple preferred extensions exist. However, this is not the case, as demonstrated by the following counter-example.

**Example 6.** *Consider the matching pennies game  $G = (Ag, Ac, Av, Ou, Ef, \leq)$  where  $Ag = \{0, 1\}$ ,  $Ac = \{\text{heads}, \text{tails}\}$ ,  $Av = [Ac, Ac]$ ,  $Ou = \{1, -1\}$ ,  $\leq$  is defined as the “less-than relation” for each player, and  $Ef$  is defined in Table 1b.*

*The set of arguments is  $\mathbb{A} = \{b_1, b_2, b_3, \dots, b_{16}\}$  and are listed in Table 3. There is only one preferred extension  $\{b_{16}, b_{15}, b_{14}, b_{13}, b_{12}, b_{10}, b_8, b_6\}$  but no stable extensions.*

Furthermore, even when multiple preferred extensions exist, these may not coincide with the stable extensions.

Game-based arguments	Preference arguments	Valuation arguments
$b_1 = [\text{heads}, \text{heads}]$	$b_5 = ([\text{heads}, \emptyset], \text{heads})$	$b_{13} = ([\text{heads}, \emptyset], \text{tails} > \text{heads})$
$b_2 = [\text{heads}, \text{tails}]$	$b_6 = ([\text{heads}, \emptyset], \text{tails})$	$b_{14} = ([\emptyset, \text{tails}], \text{tails} > \text{heads})$
$b_3 = [\text{tails}, \text{tails}]$	$b_7 = ([\emptyset, \text{tails}], \text{heads})$	$b_{15} = ([\text{tails}, \emptyset], \text{heads} > \text{tails})$
$b_4 = [\text{tails}, \text{heads}]$	$b_8 = ([\emptyset, \text{tails}], \text{tails})$	$b_{16} = ([\emptyset, \text{heads}], \text{heads} > \text{tails})$
	$b_9 = ([\text{tails}, \emptyset], \text{tails})$	
	$b_{10} = ([\text{tails}, \emptyset], \text{heads})$	
	$b_{11} = ([\emptyset, \text{heads}], \text{tails})$	
	$b_{12} = ([\emptyset, \text{heads}], \text{heads})$	

Table 3. Arguments for the matching pennies game

		Player 1		
		heads	tails	edge
Player 0	heads	1	-1	1
	tails	-1	1	1
	edge	-1	1	1

Table 4. Three strategy variant of the matching pennies game.

**Example 7.** Let us consider the following variant of the matching pennies game with three strategies for each player. We have  $G = (Ag, Ac, Av, Ou, Ef, \leq)$  where  $Ag = \{0, 1\}$ ,  $Ac = \{\text{heads}, \text{tails}, \text{edge}\}$ ,  $Av = [Ac, Ac]$ ,  $Ou = \{1, -1\}$ ,  $\leq$  is defined as the "less-than" relation for numbers for each player, and  $Ef$  is defined in Table 4. This variant of the game has eight distinct preferred extensions, but none contain any game-based arguments.

We now turn to our main result, namely the equivalence of the Nash equilibrium with the game-based arguments found in the preferred extensions.

**Proposition 3 (Equivalence).** Let  $G = (Ag, Ac, Av, Ou, Ef, \leq)$  be a game, and  $\mathcal{AS}_G$  be the argument framework for the game. A strategy profile  $S = [s_0, \dots, s_n] \in S_G$  is a Nash equilibrium iff there exists  $E \in \text{Ext}_p(\mathcal{AS}_G)$  such that  $S \in E$ .

*Proof.* We split this proof in two parts:

( $\Rightarrow$ ) We need to show that if  $S$  is a Nash equilibrium, then it is within a preferred extension of  $\mathcal{AS}_G$ . Let us consider the set of arguments  $E = \{S\} \cup \mathcal{A}_v(G) \cup \{(S_{-i}, s_i) \mid i \in Ag\}$ . We now show that  $E$  is a preferred extension of  $\mathcal{AS}_G$ . It is clear that  $E$  is conflict-free as for every  $x, y \in E$ ,  $(x, y) \notin \mathbb{C}$ . Every argument in  $\mathcal{A}_v(G)$  is acceptable w.r.t.  $E$  as valuation arguments are not attacked. Every argument  $a = (S_{-i}, s_i)$  is also acceptable w.r.t.  $E$  because for every  $s' \in Ac_i$  and  $s' \neq s_i$ , the attacks from  $a' = (S_{-i}, s')$  to  $a$ , is either not a defeat w.r.t.  $E$  (if there is a valuation argument that attacks  $(a', a)$ ) or it is a defeat but  $a'$  is defeated by  $a$  w.r.t.  $E$ . The argument  $S$  is also acceptable w.r.t.  $E$  because for every  $S' \in S_G$  and  $S' \neq S$ , the attack from  $S'$  to  $S$  is not a defeat w.r.t.  $E$  as the arguments  $(S_{-i}, s_i)$  are attacking those attacks. We conclude that the set  $E$  is admissible.

Following Lemma 2 and 1, we conclude that  $E$  is maximal for set inclusion as it contains all the valuation arguments, one preference argument per cluster and exactly one game-based argument.

( $\Leftarrow$ ) We need to show that if  $S$  is within a preferred extension, then  $S$  is a Nash equilibrium. This follows directly from the result from Corollary 1. □

Returning to the stable extensions, the following result shows that there is a one-to-one correspondence between the sets of Nash equilibria and the set of classes of stable extensions<sup>4</sup>, where each Nash equilibrium  $S$  corresponds to the class of stable extensions containing argument  $S$ .

**Corollary 2.** *Let  $G = (Ag, Ac, Av, Ou, Ef, \leq)$  be a game, and  $\mathcal{AS}_G$  be the corresponding EAF. There is a bijection between  $Y = \{S \in S_G \mid S \text{ is a Nash equilibrium}\}$  and  $\{\{E \in Ext_s(\mathcal{AS}_G) \mid S' \in E\} \mid S' \in Y\}$*

*Proof.* Follows directly from Proposition 3 and Proposition 2. □

Finally, we consider how many arguments an argumentation system containing representing a normal form game will contain.

**Proposition 4** (Number of arguments). *Let  $G = (Ag, Ac, Av, Ou, Ef, \leq)$  be a game s.t.  $|Ag| = n$  and  $m = \max_{i \in Ag} |Ac_i|$ , the number of arguments in  $\mathcal{AS}_G$  is in  $\mathcal{O}(m^{n+1} \cdot n)$ .*

*Proof.* The proof is split into three parts.

1. Suppose  $n$  players and  $m$  strategies per player. Each game-based argument corresponds to a pure strategy profile, i.e., there are  $m^n$  game-based arguments.
2. Consider the number of the preference arguments. There are  $m^{n-1} \cdot n$  partial strategy profiles. Roughly speaking, a preference argument is obtained from a partial strategy profile by replacing the empty set with a strategy. Hence, there are up to  $m^{n-1} \cdot n \cdot m = m^n \cdot n$  preference arguments.
3. We estimate the number of valuation arguments. Each valuation argument is obtained from one partial strategy profile and one pair of different strategies. There are  $m^{n-1} \cdot n$  partial strategy profiles and up to  $m \cdot (m - 1)$  pairs of different strategies. Furthermore, if a strategy  $x$  is preferred to strategy  $y$ , then  $y$  is not preferred to  $x$ . Thus, there are up to  $\frac{m \cdot (m-1)}{2}$  possible combinations to consider. Hence, the total number of valuation arguments is limited by  $\frac{m^{n-1} \cdot m \cdot (m-1) \cdot n}{2}$  which is in  $\mathcal{O}(m^{n+1} \cdot n)$ . Thus, the total number of arguments is in  $\mathcal{O}(m^n) + \mathcal{O}(m^n \cdot n) + \mathcal{O}(m^{n+1} \cdot n)$  which is in  $\mathcal{O}(m^{n+1} \cdot n)$ . □

We note that computing Nash equilibria is known to be computationally difficult, and the result regarding the number of arguments is therefore unsurprising.

<sup>4</sup>We say two stable extensions are equivalent iff they have the same game-based argument

## 5. Discussion, Related and Future Work

In this paper, we described how normal form games can be given an argumentation-based interpretation so as to allow – via argumentation semantics – for pure Nash equilibria to be computed. Intuitively, a Nash equilibrium identifies the best strategy a player can pursue given others’ strategies. However, explaining – to a non-expert – why some set of strategies forms a Nash equilibrium is often difficult, and our argument-based interpretation is the first step towards an explanatory dialogue for such explanation. Other work has shown the utility of providing such dialogue-based explanations [4, 7, 12]. In the current context, such an explanation could build on Modgil’s proof dialogues for extended argumentation frameworks [10], and could result in a dialogue as follows for the Stag hunt game shown in Figure 1.

**User** “Why should both players hunt a stag?” (why  $a_2$ ?)

**System** “It is the best response because  $a_2$  defeats all the other game-based arguments, namely  $a_1$ ,  $a_3$  and  $a_4$ ”.

Assume now that the user agrees that  $a_2$  defeats  $a_3$  and  $a_4$ ; hence they ask further about why  $a_2$  defeats  $a_1$ .

**User:** “Why should player 1 play stag if player 0 plays stag?” (why does  $a_2$  defeats  $a_1$ ?)

**System:** “Because playing stag gives a better outcome to player 1 if player 0 plays stag” ( $a_5$  defeats the attack ( $a_1, a_2$ ))

**User:** “Why does player 1 not prefer the outcome when hare is played”? (why not  $a_6$ )?

**System:** “Because of the valuation defined for player 1” ( $a_{13}$ )

**User:** “I understand.”

In the short term, we intend to formalise the dialogue and empirically evaluate its explanatory capability with human subjects. Other extensions which we intend to investigate include providing an argumentation semantics for mixed Nash equilibria (perhaps through the use of some form of ranking semantics [1, 3, 9]), and investigating other solution concepts (e.g., Pareto optimality) for more complex types of games. Finally, there are clear links between game theory and group-based practical reasoning. Building on work such as [2, 15], we intend to investigate how an argument-based formulation to practical reasoning underpinned by game theory can be created.

Several other authors have investigated some links between game theory and argumentation. For example, in his seminal paper, Dung [5] noted that the stable extension corresponds to the stable solution of an cooperative  $n$ –person game, but did not seem to deal with non-cooperative games as we do here. Game theory was also used to describe argument strength by Matt and Toni [9], and Rahwan and Larson [14] investigated the links between argumentation and game theory from a mechanism design point of view. Perhaps most closely related to the current work is Fan and Toni’s work [6] exploring the links between dialogue and assumption-based argumentation (ABA). Here, the authors showed how admissible sets of arguments obtained from their ABA constructs are equivalent to Nash equilibria. In contrast to the current work, they only considered two player games and utilised structured argumentation, allowing them to describe a proof dialogue with associated strategies.

## 6. Conclusions

In this paper, we provided an argumentation-based interpretation of pure strategies in normal form games, demonstrating how argumentation semantics can be aligned with the Nash equilibrium as a solution concept, and examining some of the argumentation system's properties.

We believe that this work has significant application potential in the context of argument-based explanation. At the same time, we recognise that there are significant open avenues for research in this area, but believe that the current work is an important step in investigating the linkages between the two domains.

## References

- [1] L. Amgoud, J. Ben-Naim, D. Doder, and S. Vesic. Ranking Arguments With Compensation-Based Semantics. In *KR*, 2016.
- [2] K. Atkinson and T. J. M. Bench-Capon. Argument schemes for reasoning about the actions of others. In *Proc. COMMA*, volume 287, pages 71–82, 2016.
- [3] E. Bonzon, J. Delobelle, S. Konieczny, and N. Maudet. A comparative study of ranking-based semantics for abstract argumentation. In *Proc. AAI-16*, pages 914–920, 2016.
- [4] M. Caminada, R. Kutlák, N. Oren, and W. W. Vasconcelos. Scrutable plan enactment via argumentation and natural language generation. In *AAMAS*, 2014.
- [5] P. M. Dung. On the acceptability of arguments and its fundamental role in non-monotonic reasoning, logic programming and n-person games. *Artificial Intelligence*, 77(2):321–357, Sept. 1995.
- [6] X. Fan and F. Toni. On the Interplay between Games, Argumentation and Dialogues. In *Proc. AAMAS-16*, pages 260–268, May 2016.
- [7] C. Kristijonas, S. Ken, and T. Francesca. Explanation for Case-Based Reasoning via Abstract Argumentation. *Frontiers in Artificial Intelligence and Applications*, pages 243–254, 2016.
- [8] A. Matsumoto and F. Szidarovszky. *Game Theory and Its Applications*. Springer Japan, 2016.
- [9] P.-A. Matt and F. Toni. A Game-Theoretic Measure of Argument Strength for Abstract Argumentation. In *Logics in Artificial Intelligence*, LNCS, pages 285–297, 2008.
- [10] S. Modgil. Labellings and games for extended argumentation frameworks. In *Proc. IJCAI-09*, pages 873–878, 2009.
- [11] S. Modgil. Reasoning about preferences in argumentation frameworks. *Artificial Intelligence*, 173(9-10):901–934, June 2009.
- [12] N. Oren, K. van Deemter, and W. W. Vasconcelos. Argument-Based Plan Explanation. In M. Vallati and D. Kitchin, editors, *Knowledge Engineering Tools and Techniques for AI Planning*, pages 173–188. Springer International Publishing, 2020.
- [13] M. Osborne. *Introduction to Game Theory: International Edition*. OUP, 2009.
- [14] I. Rahwan and K. Larson. Argumentation and Game Theory. In G. Simari and I. Rahwan, editors, *Argumentation in Artificial Intelligence*, pages 321–339. Springer US, Boston, MA, 2009.
- [15] Z. Shams, M. D. Vos, N. Oren, and J. Padget. Argumentation-based reasoning about plans, maintenance goals, and norms. *ACM Trans. Auton. Adapt. Syst.*, 14(3), 2020.