

Arms Races and Conflict: Experimental Evidence

KLAUS ABBINK, LU DONG, LINGBO HUANG*

June 2020

Abstract

We study escalation and aggression in an experimental first-strike game in which two participants play multiple rounds of a money-earning task. In each round, both players can spend money to accumulate weapons. The player with more weapons can spend money to strike against the other player, which almost totally eliminates the victim's earnings potential and removes their capacity to strike. Weapons can serve as a means of deterrence. In four treatments, we find that deterrence is strengthened if weapon stocking cannot be observed, that a balance of power is effective in maintaining peace, and that mutually beneficial trade decreases the risk of confrontation, but not necessarily the likelihood of costly arms races.

Keywords

Mutually assured destruction, balance of power, arms races, deterrence, trade, laboratory experiment

JEL Classification Codes

C72, C91, F51, N40

* Abbink: Monash Business School, Clayton, VIC, Australia, klaus.abbink@monash.edu. Dong: Economics Experimental Laboratory, Nanjing Audit University, Nanjing 211815, China, lu.dong@outlook.com. Huang: Economics Experimental Laboratory, Nanjing Audit University, Nanjing 211815, China, lingbo.huang@outlook.com. Financial support from the Australian Research Council (DP1411900) and National Natural Science Foundation of China (Grant No. 71873068) is gratefully acknowledged.

North Korean Leader Kim Jong Un just stated that the “Nuclear Button is on his desk at all times.” Will someone from his depleted and food starved regime please inform him that I too have a Nuclear Button, but it is a much bigger & more powerful one than his, and my Button works!

—Donald J. Trump, Jan. 2, 2018, via Twitter

1. Introduction

During the Cold War, the once-allied US and USSR entered an unprecedented arms race. At its peak in 1986, the two superpowers had more than 60,000 nuclear warheads targeted at each other (Norris and Kristensen 2010). Though the notion of complete human extinction by nuclear war (“overkill”) is largely a myth, the devastation caused in a nuclear confrontation would still be immense, with death tolls at least in the hundreds of millions, if not billions. Civilizations would be set back to primitive states. Fortunately, the Cold War ended without cataclysm, but the players’ nuclear arsenals are still there, albeit much diminished, and new players like Iran and North Korea are entering the field.

The arms race created what is known as the *Hobbesian trap*: the dilemma that occurs if fear of an attack leads to more armament, which leads to more fear, and so on, until one or both sides may feel tempted to launch a pre-emptive strike. Even if *no* side wishes to destroy the other, the pre-emptive strikes occurs out of bilateral fear of an imminent attack. Thomas Schelling (1960) offered a classical analogy to a homeowner confronting a burglar, both carrying a gun, with each being tempted to shoot the other before being shot. One way out of this trap is the doctrine of *Mutually Assured Destruction* (MAD), which guided the relations between the superpowers during much of the Cold War. It requires both sides to be strong enough that they possess not only first-strike but also second-strike capabilities. If neither side is able to eliminate the opponent’s arsenal in a single strike, a first strike will trigger retaliation with similar force. This threat of a second strike deters both sides from attacking and thus maintains a tense peace.

The MAD doctrine is built primarily on two pillars. The first is *deterrence*. If a military power has an arsenal large enough to sustain second-strike capability, then it would be suicidal for any rival to attack it. MAD requires mutual recognition of this capability. During the Cold War, this state was achieved with the introduction of submarine-based intercontinental missiles. Their mobility would prevent a first-striker from locating and destroying the opponent’s capacity for a second strike. The second building block of the MAD doctrine is the maintenance of a *balance of power*. If power disparity is significant, the powerful party is likely to launch a first strike to avoid future escalation. However, once a power balance has been established, it is often unstable, since each side strives to outstrip the other (or at least suspects as much of its rival). Both sides therefore try to develop more and more effective weaponry, leading to expensive arms races.

The maintenance of the MAD doctrine is both extremely costly and risky. It is costly because it requires enormous investment—at its height, in 1967, 9.1% of US GDP was spent on defense

(World Bank 2018). This expenditure is, by the very design of MAD, unproductive—huge arsenals are piled up in order never to be used. In fact, many commentators see the Soviet Union’s inability to sustain the levels of defense spending required to match Ronald Reagan’s expansion of the military budget as a cause of its collapse (Fitzgerald 2000). In addition, the MAD strategy is inherently risky. Global nuclear war could be triggered by misperception, miscommunication, or simple false alarms.¹

MAD helped preserve a tense peace during the Cold War. However, this does not mean that arms races are necessary or inevitable. Clearly, they are not desirable. They are even to some extent paradoxical, since it is hard to imagine a motive for an annihilatory attack. Even the winner of a nuclear war would gain little from obliterating entire nations and turning their territory into radioactive rubble. Moreover, in so doing an attacker would forego any opportunity for productive interaction through mutually beneficial trade. So why, at great cost, protect yourself against something no rational actor would do?

This argument, of course, requires that all actors always act rationally, and that maximizing their own country’s benefit is their main motivation. This may not be universally true. It is sufficient that one, or both, actors fear that the other may act for reasons incomprehensible to them. Abbink and de Haan (2014) analyze experimentally how fear can motivate destruction in their first-strike game, in which two players are locked in a situation of mutual threat for multiple rounds. In every round, one player can decide to strike against the other. A player who strikes first is protected against future strikes by their opponent; in addition, the victim loses almost all payoff. Though no material benefit is gained from striking, and therefore purely money-motivated players would never attack, the authors observe that fear can lead up to 80% of subjects to strike (on a more optimistic note, they also find that trust can build up peaceful relations).

It is difficult to study nuclear arms races using observational data, due to the low number of observations and the highly specific contexts in which they occur (see section 4 for an overview of the International Relations literature). In this paper we present a novel experimental paradigm for studying arms races. We extend the first-strike game to incorporate the possibility of arms races between two players. In each round players can invest in arms; think of them as “rockets.” In contrast to the first-strike game, a player can successfully strike *only if* they have more rockets than the opponent. In our model rockets are entirely unproductive and strikes are costly, thus in the subgame perfect equilibrium we should observe no arms investments and no strikes. However, the motives of fear that drove destruction in Abbink and de Haan (2014) are still present. It is thus not unreasonable to expect that experimental subjects may wish to protect

¹ In 1983, Stanislav Petrov, a lieutenant-colonel in the Soviet intelligence service, saw five attacking US missiles on the Soviet early warning system he monitored. Petrov decided that it was a false alarm and did not pass it on, which averted full-on nuclear war. He later reported that his decision was based on “a guess” (Hoffman 1999).

themselves through armament, or strike if they have the opportunity, and hence an arms race is born out of mutual fear and the desire for self-protection (Jervis 1976).²

In different treatments we study the effect of the aforementioned building blocks of MAD. We examine the contribution of (visible) deterrence by contrasting our baseline game with a treatment in which one's strength is not known to one's opponent (the *Hidden* treatment). It is therefore not possible to calibrate an arsenal to deter the opponent from striking, and therefore there is less reason to arm. A player might nonetheless want to arm as a kind of *insurance*, albeit incomplete, against attacks. Such attacks could still be possible, as a highly invested player might be tempted to strike expecting that the other player is weak, and thereby make future investment unnecessary. However, striking is always risky in this treatment, since it is possible that the other player has heavily invested as well. With the opponent kept in the dark, deterrence can only work indirectly, through the *potential* possession of retaliatory forces, the acquisition of which can be rationalized by the insurance motive.

We study the importance of a balance of power with the help of an *Asymmetry* treatment. In this treatment one of the players has a head start of half a rocket, while the other starts with no rockets. A balance of power can therefore never be achieved: one player is always ahead of the other. Again, the effect of this treatment on war and peace is, at the outset, ambiguous. The initially stronger player may be more likely to strike, perhaps immediately, to rid himself of any threat once and for all. On the other hand, one can also imagine very peaceful low-armament outcomes, since the stronger player may not perceive the weaker one as a serious threat: after one initial investment, the stronger player has full control. He can strike if the weaker invests too, or desist from rocket purchases otherwise.

We further introduce a *Trade* treatment to test the effect of economic exchange on arms races and conflict. The role of trade as a peacekeeper has long been subject to debate (e.g., Mansfield 1994; Barbieri 2002; Martin, Mayer, and Thoenig 2008; Li and Reuveny 2011). It seems reasonable to assume that countries with strong, mutually beneficial economic ties are less likely to attack each other, as they would forego future gains. However, the two most prominent historic examples seem to suggest otherwise. The Cold War never escalated to a military confrontation between the blocs, despite very little economic interdependence. Before World War I, on the other hand, there was ample trade between all European powers. In fact, some contemporaries proclaimed war to be a thing of the past for that very reason (Walsh 1910). They were clearly wrong. To test the effect of trade in our framework, subjects play a bargaining game that determines their earnings, in contrast to the other treatments where both players earn their payoffs independently. If one player strikes, bargaining opportunities vanish.

Our results can be summarized as follows. We observe much lower strike rates when the opponent's strength is unknown. This does not, however, significantly lower investments in

² Subjects may strike to avoid future costs of armament if they believe their opponents to be relentless in attempting to catch up with them (although neither the belief nor the opponent's buying rockets may be rational). Similar incentives have been the main focus of Garfinkel and Skaperdas (2000).

rockets. The balance of power is vital for peace: in its absence, strike rates dramatically increase. Finally, economic interaction promotes peace in our context.

Our study is, to the best of our knowledge, the first to study arms races in an experimental game, with its parallels to the threat of nuclear attacks particularly novel. However, it does not stand in isolation. Our game builds on the study by Abbink and de Haan (2014) mentioned earlier. It also shares some characteristics of centipede games (e.g., McKelvey and Palfrey 1992; Nagel and Tang 1998; Palacios-Huerta and Volij 2009; Levitt, List, and Sadoff 2011), as both feature different kinds of coordination failure. Like centipede games, our games involve a temptation to deviate early and thus forego the higher joint earnings available if players cooperate throughout the game. However, our payoff structure and equilibrium are very different. First, deviating/striking immediately is the only subgame perfect equilibrium in centipede games, whereas in our games cooperating till the end is the only subgame perfect equilibrium. Second, in centipede games, an early deviation eliminates both players' potential future earnings, whereas in our games it only affects the victim's. Our game is also reminiscent of stag hunt games (e.g., Rankin, Van Huyck, and Battalio 2000; Battalio, Samuelson, and Van Huyck 2001), in the sense that the efficient equilibrium is strategically risky. It also bears some resemblance to deterministic contests such as all-pay contests (Siegel 2009) and Bertrand competitions (Dufwenberg and Gneezy 2000; Abbink and Brandts 2008; Boone et al. 2012; Cracau and Sadrieh 2016), though payoff and equilibrium, again, are very different. In particular, the tension of conflict in our game is *not* driven by potential material benefit (in fact it is costly to strike), but is hypothesized to be driven by fear and perceived need for self-defense.³

But before we turn to our experimental game, we shall note that while we have chosen a stylized setting to study arms races and strikes, some of the aspects might be more realistic in some scenarios than in others. For instance, across all treatments reported in this paper, second strikes capacity is only granted once one player has at least as many rockets as the other, and will be automatically executed if this player is attacked. One might imagine a fuzzier design in which even a very weak player might be able to counterstrike with a small probability (such as by sponsoring terrorism) might better fit the real world scenario. We believe our framework is tractable enough to accommodate this and many other ideas surrounding various versions of “realistic” conflict that researchers have in mind. More on tractability is to be discussed in Section 5.

The remainder of the paper is organized as follows. Section 2 describes the experimental design. Section 3 reports the results. Section 4 discusses the contribution of the present experiment

³ Escalation has been studied in other, less familiar contexts. For example, in Chuah, Hoffmann, and Lerner (2014), players decide whether to escalate or acquiesce in a multi-stage bargaining game. The bargaining structure is that of a chicken game where both players' aggressions would lead to a head-on conflict. The motivation for aggression, as in the chicken game, is greed (monetary benefit) rather than fear. Lacomba et al. (2014) study a contest game in which the winner can appropriate the loser's resources. When the game is repeated, appropriation effort and resulting conflict can be escalated. However, here again the aggressor's motivation is to obtain the opponent's resources. In another repeated competitive setting, Bolle, Tan, and Zizzo (2014) consider a two-player vendetta game and attribute the escalation of retaliatory behavior to negative emotions.

through the lens of the literature. Section 5 concludes and discusses the possibility of studying other interesting aspects of arms races under our experimental framework.

2. Experimental Design

The experiment comprises an extensive-form game of altogether 12 full rounds and one partial one. We opted for a setup with a known length of play rather than an indefinitely repeated game because of its clear subgame perfect prediction and easy experimental implementation.⁴

2.1 The Game

The basic setup for the game is as follows. In each of the first 12 rounds, each player makes three different decisions on the same computer screen (see Figure 1):

Lottery Decision: In each round, both subjects earn a payoff, paralleling income from engaging in productive activities. How this occurs is not essential for our research purpose. To avoid experimenter demand effects, we decided to make subjects perform a meaningful task, rather than just giving them a sum of money each round. We opted for lottery choice tasks, as they are standard tools in experiments. Each player chooses between two lotteries which are played out independently with the player privately informed of the outcome. The lottery pairs comprise two series of Holt and Laury (2002)-style lottery pair sequences. The lottery outcome ranges from AU\$1.70 to AU\$3.40. The sequence in which the lottery pairs were presented was randomized for each participant.⁵

Rocket-buying Decision: Each player decides whether or not to buy a rocket, or, in the language of the experimental instructions, spend AU\$1.00 to buy a “token.” (In each round, players are only allowed to buy one rocket.) Any rocket bought in the current round will be effective in the next. Expenses on rockets are nonrefundable and the stack of rockets grows from round to round.⁶

Deactivation Decision: Each player decides whether or not to strike against their opponent. In the experimental language, this means spending AU\$1.50 to “deactivate” the other player. If the player with more rockets presses the deactivation button during a round, the other will be deactivated. However, if the player with fewer rockets presses the deactivation button, they

⁴ An alternative design is to model the repeated interaction between the two players as an indefinitely repeated game. Though fully solving the theoretical prediction of this game could be challenging, from an experimental viewpoint, running such a game would be interesting as it may better capture the uncertain nature of real-world confrontations. For example, war may occur abruptly due to human or machine errors.

⁵ We were not particularly interested in eliciting risk preferences. For that purpose we would have had to tailor the experiment to this aim. In lottery choice experiments, one lottery is typically chosen randomly to be paid, to avoid subjects perceiving the choices as a compound lottery. In our game, payoffs had to be cumulative. Further, we randomized the sequence in which the lottery pairs were presented for each subject, to avoid possible interference between the dynamics of the game and different expected payoffs across the lottery sequences. These design choices were necessary for our experimental game, but would be considered fatal flaws in an experiment on risk preferences.

⁶ Players are endowed with a startup fund of AU\$5 (in addition to a show-up fee of AU\$5) at the beginning of the experiment. This startup fund allows them to buy rockets in the first round.

will only deactivate themselves; if both have an equal number of rockets, both will be deactivated when either presses the deactivation button. Note that players who press the deactivation button still need to pay the deactivation cost (cutting into the show-up fee) even if they themselves are deactivated in that round. The deactivated person faces severe payoff losses: (1) they will lose all earnings from previous lotteries; 2) they will only earn 10% of the value of future lotteries; and 3) they will not be able to deactivate the other side in the future. In this sense, deactivation is a small-scale equivalent of an all-out nuclear first strike: the targeted country is devastated beyond recovery, and no longer poses any threat to the attacker. Note that the constellation with equal armament corresponds to the MAD world, in which a first strike immediately triggers a counterstrike.

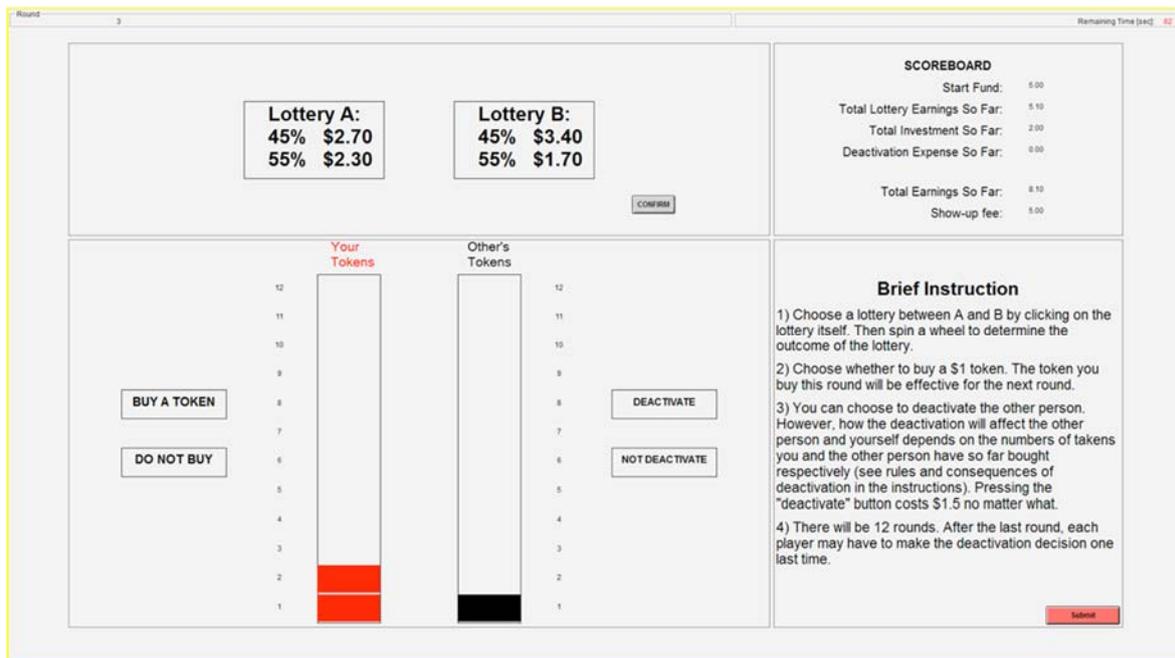


Figure 1: a screenshot of the game interface (Base treatment)

Notes: This figure is a screenshot of the game interface (in the Base treatment). In the upper left panel, after a player chose a lottery, the payoff would be realized immediately. In the lower left panel, a player made two decisions: whether to buy a rocket (token) and whether to deactivate the other player. The upper right panel listed a player's accumulated earnings and expenses up to the current round. Only after all three decisions were made could a player proceed to the next round.

To keep things neat, we designed a situation in which already a minimal superiority in armaments creates a first-strike capability. This gives us the best chance to observe many cases in which first strike opportunities appear and disappear, though in reality, an exact balance is neither possible nor required. Iran or North Korea will never be able to match the United States' nuclear arsenal warhead by warhead. Nevertheless, their nuclear ambitions cause great concern. While they could not win a nuclear war against the US, they can threaten a devastating strike

against her allies, Israel and South Korea, should the US attempt to invade them.⁷ Balance of power effectively means that both sides maintain second-strike capability.

In the 13th round, players only need to make the deactivation decision. This is to allow any rocket bought in the 12th round to be effective.

2.2 Deterrence

As discussed in the introduction, MAD is built on deterrence, which requires that the opposition knows about one's own second-strike capability. However, Cold War-era behavior was paradoxical in this regard, as both blocs were extremely secretive about their capabilities, and espionage was big business. It is in fact not obvious which is the more effective way to establish deterrence: be open about your strengths (which then, of course, also reveals your weaknesses), or keep your rival in the dark. Our game allows us to address this question with a very simple treatment variation. In the *Base* treatment, players see their own and their adversary's current rocket (token) balances on the game interface in every round. Thus, both players always know the consequence of a strike: either the player with fewer rockets or both players (in the case of equal rocket stocks) will be neutralized. In the *Hidden* treatment, we hide the adversary's current rocket balance to remove signaling as an element of the deterrence strategy, such that the consequence of deactivation is uncertain. Our two treatments thus represent idealized extreme cases of openness and secrecy. In the Base treatment the players know everything about their opponent's strength; in the Hidden treatment they know nothing.⁸

As visible deterrence is not possible in the Hidden treatment, it thus seems reasonable to expect fewer deactivations (in cases where rocket asymmetry arises) and less rocket investment in the Hidden treatment. However, players may still be "deterred" by the mere possibility that the opponent may be fully invested and attacking her may lead to their own destruction. In this case, it does not matter whether the opponent's rockets actually exist. Nonetheless, while one's own actual investment cannot deter an opponent, it may serve as some *insurance* against irrational aggression. A fully invested player would survive an attack by a less invested one. However, the insurance is not complete, because an attack by a fully invested opponent would still be fatal. It should be noted that the same insurance motive can also be present in the Base treatment. Therefore, the comparison to the Hidden treatment reveals to what extent escalation is due to visible deterrence.

2.3 The balance of power

In the Base treatment, peace can be sustained under the balance of terror, i.e. both sides can create a credible deterrent to each other by escalating at the same rate, but neither can deactivate

⁷ One could call this strategy 'MAD by proxy'—the nuclear newcomer cannot threaten the hegemon directly but can destroy a small close ally of it. In our study we model the direct confrontation of two powers.

⁸ Note that in the two treatments players can press the deactivation button in the very first round when no rocket has yet been bought (and therefore the only possible outcome is mutual deactivation). As will be clear later, this design feature is meant to keep treatment comparisons easy because in other treatments one side does possess some initial endowment of rockets.

others without destroying themselves. To study whether the balance of power is conducive to increased escalation and reduced deactivation, we introduce the *Asymmetry* treatment, in which we disrupt the power balance by endowing one side with half a rocket at the beginning of the game. The rocket escalation can no longer create a credible deterrent simultaneously to both sides; in any round, there is always at least a half-rocket difference, which is sufficient for the player with more rockets to gain the capability to deactivate the other side without deactivating herself.

Will the perpetual power asymmetry drive or arrest arms races? This is hard to say before the fact. The initially stronger player may be tempted to strike immediately to eliminate all future threats and make any investment in rockets unnecessary. This option may indeed be ideal if the player is motivated by (1) the fear that the initially weaker player will invest in the first round, striving to overtake and strike against them, (2) the goal of securing their payoff while minimizing costs, and (3) near-complete indifference to the opponent's suffering. However, unless the stronger player is very certain about (1), there might be an alternative strategy that satisfies (2) and does not require (3) in all its cruelty. The stronger player can invest in the first round and then observe his opponent's behavior—if they have invested, too, then strike, if not, hold back. In the latter case, the stronger player has even saved some money, since investing in a single rocket costs AU\$1, less than the deactivation cost of AU\$1.50. The stronger player also remains completely safe, since they retain the ability to respond to any future armament by their opponent with a deactivation. Even if the weaker player has invested, the loss for the stronger player is minor, consisting of one additional rocket (AU\$1) relative to the payoff he would have earned had he struck immediately. Therefore, as long as the weaker player knows their place, a power imbalance need not lead to war and destruction. This strategy might be viewed as “superiority” which was eventually pursued by the US in the Cold War.

2.4 Trade opportunities

In a situation marked by confrontation, trust-building is difficult and loss of a balance of power can form an inexorable driver of head-on conflict. Human beings seem to have overcome this trap globally in past decades. One explanation is the increasing volume of worldwide trade. Statistics suggest that bilateral peace between nations is positively correlated with the scale of their economic interdependence (Oneal, Russett, and Berbaum 2003; Li and Reuveny 2011), even when their military power is imbalanced. But is this relationship causal? For example, lowered tolerance for violence, the spread of democracy, and the growth of international organizations may all have contributed to the lock-in effect between peace and trade (Kant 1795; Pinker 2011).

If rocket escalation cannot sustain peace in the *Asymmetry* treatment, do trade opportunities provide an alternative path? To study whether trade opportunities help reduce the deactivation inclination absent a balance of power, in the *Trade* treatment, we add to the *Asymmetry* treatment a trading stage where the two players can bargain over a fixed amount of money. We could have added the trading stage to the *Base* treatment, but we believe it is more informative to establish the effect of trade in reducing the deactivation rate in the situation where

deactivations are more common (we designed this treatment after we had conducted the others, and knew that the Asymmetry treatment had produced the most deactivations). Specifically, in each round, after each player learns their lottery payoff (with payoffs ranging from AU\$1.70 to AU\$3.40), they proceed to bargain over AU\$7. The size of the bargaining pie ensures that there are always potential gains for both sides. The bargaining happens in real time and is semi-structured. Each player can move a cursor on a scale to indicate their demand out of the AU\$7, and they can learn the other player’s demand, which is updated in real time. After 60 seconds, if the sum of their demands does not exceed AU\$7 (trade success), they receive their respective demands; otherwise (trade failure), they only receive their own lottery payoff.⁹ Other features of the Asymmetry treatment are all retained: one of the two players has a head start of half a rocket. In each round players must decide whether to buy a rocket and whether to deactivate the other side. Once at least one side is deactivated, there will be *no* trading opportunities in future rounds; the deactivator will earn the standard lottery payoff in future rounds, while the deactivated player loses all previous earnings and will only earn 10% of future lottery payoffs.

The effect of bargaining on peace and escalations will depend on whether the players find mutually agreeable bargaining outcomes. If they do, then the costs of deactivating are now much higher even for a fearful player. These costs now comprise not only direct payment for deactivation, but also the opportunity costs of lost future trade. The effect on escalations is less clear. It is possible that a player may invest not only to deter their opponent, but also to create a threat that may boost their bargaining power—the deactivation button can turn even a poor proposal into an offer the opponent cannot refuse.

Table 1 provides a summary of all treatments.

Table 1: Experiment Design

Treatment	Num. of Obs.	Initial Tokens	Other’s Tokens	Trade Option
Base	$29 \times 2 = 58$	Both start with 0	Informed	No
Hidden	$29 \times 2 = 58$	Both start with 0	Not informed	No
Asymmetry	$30 \times 2 = 60$	One starts with 0.5	Informed	No
Trade	$26 \times 2 = 52$	One starts with 0.5	Informed	Yes

2.5 Game theoretic prediction

Here, we use the concept of subgame perfect equilibrium (SPE) to derive game theory predictions about deactivation and rocket-buying decisions. The Online Appendix contains a full analysis of the SPE for all treatments.

Starting from the last (13th) round, where the only decision to make is whether to deactivate the other side, it is clear that neither side will press the deactivation button whether or not

⁹ This is essentially Nash bargaining with each side’s lottery payoff as disagreement points. We chose this bargaining protocol over more structured ones such as Ultimatum or alternating offer bargaining primarily for its simplicity and clarity for the subjects. We also believe that this protocol has sufficiently captured the mutual benefits of trading opportunities to allow us to study their effects on escalation and deactivation.

rocket asymmetry exists, because deactivation costs money (AU\$1.50) and returns zero profit. Moving backward to the 12th round, as both sides anticipate no deactivation in the final round, they do not have to buy rockets. Given that there is no need to buy rockets and therefore no money to be saved through deactivation, they will also not deactivate in this round. By this chain of backward induction, the only SPE of the game is for both sides to neither escalate nor deactivate in any round. Moreover, neither deterrence nor balance of power affects this equilibrium, in theory.

In the Trade treatment, the characterization of a SPE also involves strategies in the bargaining stages. As with all Nash bargaining games, there exists a continuum of bargaining equilibria, since players can divide the pies in any way as long as the bargaining payoffs are no lower than their lottery payoff in a round. However, this multiplicity does not change the SPEs with respect to deactivation and investments. On a SPE path, no deactivation is carried out, since this would be costly even if all future bargaining should fail (in such a case, future payoffs are the same whether or not the opponent is deactivated, but the direct cost of deactivation makes this outcome less remunerative). Any threat created by rockets would therefore not be credible, and hence no investments should be made.

2.6 Experimental procedure

The experiment was programmed in z-Tree (Fischbacher 2007) and conducted at the Monash Laboratory for Experimental Economics (MonLEE). We recruited 228 participants from a university-wide undergraduate student pool. Participants were randomly seated in a partitioned computer terminal upon arrival. The experimental instructions (see Appendix A) were provided to them in written form and were read aloud by the experimenter at the start of each session. Participants were then asked to complete a comprehension quiz, which was designed to ensure that every participant understood the instructions. This was particularly important in our game since mistakes made early in the experiment (like inadvertent deactivation) cannot be rectified later. An average of 25 minutes per session was dedicated to ensuring comprehension. At the end of the experiment, participants completed a survey concerning demographics and strategies used in the game. Participants were then paid privately and instructed to leave the laboratory one at a time. A typical session lasted about one hour with average earnings of AU\$31.50 (approximately US\$23.80 or €20.40).

3. Results

We present the results in three parts: section 3.1 looks at the role of deterrence by comparing the Base and Hidden treatments; section 3.2 studies the role of power balance by comparing the Base and Asymmetry treatments; and section 3.3 investigates the role of trade by comparing the Asymmetry and Trade treatments. In each part, we first look at deactivation decisions (conflict), followed by rocket decisions (arms race). We mainly present aggregated results for each treatment; detailed group-level results can be found in Figures B1 to B4 in Appendix B.

Figure 2 shows the cumulative deactivation rate over the course of the game. Figures 3a and 3b show the average rocket investment over rounds. Recall that the unique subgame perfect equilibrium is zero rocket investment and zero deactivation in all treatments; but even a casual look at the figures suggests otherwise.

3.1 Deterrence

Deactivation occurred in 13 out of 29 groups (44.8%) in the Base treatment, but only in 3 out of 29 groups (10.3%) in the Hidden treatment, a significant difference ($p = 0.002$, two-tailed Z test). When did the deactivation occur? As expected, in the Base treatment participants only pressed the deactivation button when they had more rockets than the adversary in a round. Of 17 groups where rocket asymmetry arose,¹⁰ nine deactivated their adversary immediately upon achieving a one-rocket advantage, four did so after some delay,¹¹ and only four groups remained peaceful. Out of these four peaceful groups, two had rocket asymmetry only in the last round, where deactivation only means money burning, and the other two groups reached rocket symmetry in late rounds. In contrast, in the Hidden treatment, all three instances of deactivation resulted in deactivation of both players, i.e., deactivation occurred when both had an equal number of rockets.¹² Of 16 groups in which rocket asymmetry had arisen at any point (though players were unaware of this), all remained in peace throughout the game. These results are consistent with the hypothesis of fewer deactivations when one's rocket stock is unknown to one's opponent.

¹⁰ In the other 12 groups, all players bought rockets *every* round so that token asymmetry never arose. See next paragraph for more detail.

¹¹ Three participants actually deactivated their adversary in the last round (two of whom had obtained the rocket advantage in an earlier round). Their behavior is similar to subjects in other games who are willing to burn other's money and may be interpreted as spite or nastiness (Abbink and Sadrieh 2009). We looked up these players' answers regarding their deactivation decision strategies in the post-experiment survey. None of them explicitly expressed pleasure at hurting others (although one deactivator might hint at this motive in saying "although the other person was smart to invest until the very last second, he/she put him/herself at risk to have all money being lost. This risk cost the person their earnings.").

¹² One of the three instances of deactivation occurred in the first round. This participant admitted that it was a mistake in the post-experiment survey. The other two instances occurred in round 6 and round 7 where both players in the group had bought rockets in every previous round. One interpretation is that these deactivators might think that deactivating the other side was more likely to succeed now than in previous rounds.

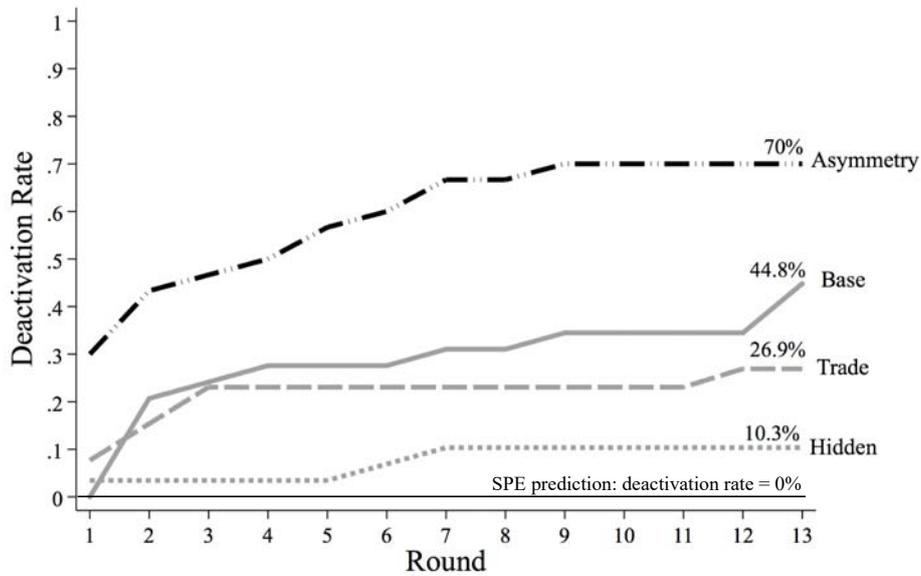


Figure 2: Cumulative deactivation rate over rounds

Notes: Deactivation occurred in 13 out of 29 groups in the Base treatment, 3 out of 29 groups in the Hidden treatment, 21 out of 30 groups in the Asymmetry treatment and 7 out of 26 groups in the Trade treatment. At round 13, deactivation was allowed to make tokens bought at round 12 effective.

Turning to the rocket-buying decisions, more groups in the Base treatment escalated neck-and-neck up to the end of the game than in the Hidden treatment (Base: 12 out of 29 groups; Hidden: 7 out of 29 groups). This difference, though not statistically significant, is directionally consistent with the deterrence hypothesis. However, on average participants spent similar amounts of money on rockets (see figure 3a, Base: AU\$7.22; Hidden: AU\$7.60; $p = 0.954$, rank-sum test). The distribution of rocket expenses is also similar between the two treatments (see figure 4, $p = 0.138$, Kolmogorov-Smirnov test). These results imply that participants in the Hidden treatment were more likely to escalate *unilaterally* than those in the Base treatment, although escalation in the Hidden treatment has proven to be almost always unnecessary, as deactivations are rare (only three cases, see previous paragraph). The rocket-buying decision in general, therefore, seems to be mainly driven by the fear of being deactivated and the insurance motive.

Nevertheless, this aggregate result masks the fact that rocket escalation was terminated much more often and much earlier in the Base treatment than the Hidden treatment. For example, in the Base treatment, almost 50% of deactivations (6 out 13) occurred in the second round. This could have limited the scope for observing escalation if the deactivated participants would have followed through with the deterrence strategy to achieve a MAD situation, and therefore caused an underestimate of the effect of deterrence. When focusing on groups where deactivation did not occur, we found that participants spent significantly more money on rockets in the Base treatment than the Hidden treatment (see figure 3b, Base: AU\$10.38; Hidden: AU\$7.96; $p =$

0.004). This, too, is consistent with the hypothesis that visibility of one's strengths (or weaknesses) strengthens the motive to arm.¹³

Result 1: Deactivation occurred more often in the Base treatment than the Hidden treatment. While overall escalation levels were similar in both treatments, the escalation level was higher in the deactivation-free groups in the Base treatment than the Hidden treatment.

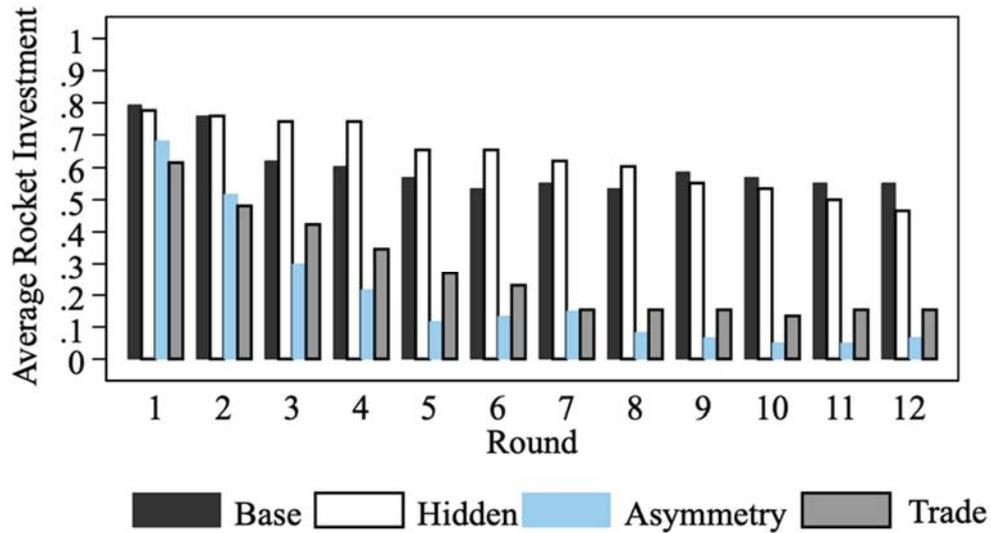


Figure 3a: Average rocket investment over rounds

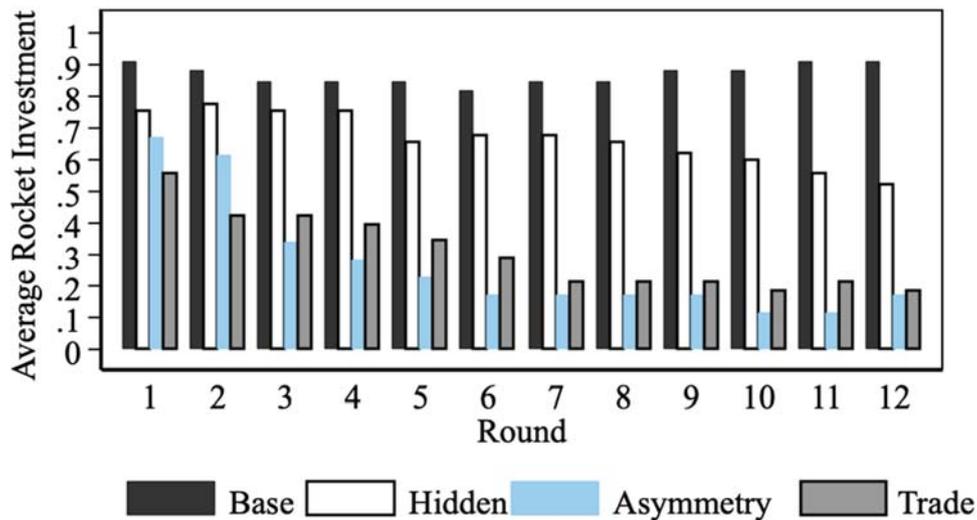


Figure 3b: Average rocket investment in deactivation-free groups over rounds

Notes: These figures show the average percentage of rocket investment in each round. 3a includes all groups, while 3b excludes groups where deactivation occurred. Note that the subgame perfect equilibrium prediction of the rocket investment is zero for all treatments.

¹³ We treat this evidence as suggestive because we implicitly assume that participants in deactivation-free groups of the Hidden treatment understood the deterrence strategy and would perform similarly to those in deactivation-free groups of the Base treatment had they participated in the Base treatment themselves. The main purpose of this exercise is to show that the effect of deterrence could have been underestimated if there happened to be more participants who did not act according to the logic of deterrence in the Base treatment than the Hidden treatment. This could be either because they did not understand this logic or because they hoped to build trust with the adversary by not escalating.

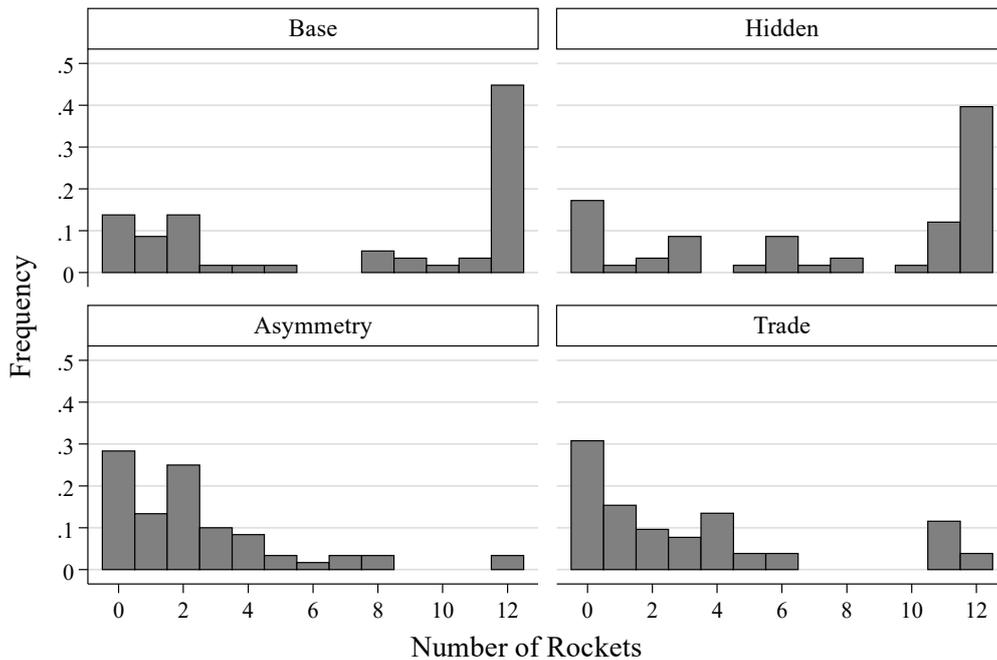


Figure 4: Distribution of rocket investment

Notes: This figure shows the distribution of players' total number of rocket purchases across all rounds in each treatment.

3.2 Balance of Power

We next study the balance of power by comparing the Base treatment and the Asymmetry treatment, in which one side is endowed with an initial half-rocket advantage. A balance of power is therefore impossible to achieve, and the tension of deactivation is present in every round. Over the course of the game, deactivation occurred in 21 out of 30 groups (70%) in the Asymmetry treatment, a frequency significantly higher than the Base treatment (see figure 2; $p = 0.050$, Z test). Nine advantaged players (those who had the half-rocket endowment) deactivated their adversary in the first round. Eight advantaged players did so after some delay.¹⁴ In the other four groups, the power balance was tilted toward the initially disadvantaged players because their opponents neither deactivated nor escalated in the first round and were eventually deactivated by the initially disadvantaged players. These results support the balance of power hypothesis and are similar to the cases where rocket asymmetry arose endogenously through investment decisions in the Base treatment.

Turning to investment behavior, participants spent on average AU\$2.43 on rockets, which is significantly lower than in the Base treatment (see figure 3a; $p < 0.001$, rank-sum test). Only one group escalated neck-and-neck till the end of the game without deactivation. Even among

¹⁴ Provocation can contribute to such deactivations. Specifically, among these 8 groups, 5 out of the 8 disadvantaged players escalated at least once before being deactivated by the advantaged player.

the nine peaceful groups, these participants only spent on average AU\$3.17 on rockets (see figure 3b). These results also support the balance of power hypothesis.¹⁵

It is interesting that some groups managed to sustain peace with low levels of escalation. As alluded to in section 2.3, peace might be reached if the advantaged players could keep at least one-and-half rockets ahead of their adversaries and the disadvantaged players did not provoke their opponents by buying rockets. Indeed, of the nine deactivation-free groups, seven groups kept at least one-and-half rockets apart. Among these groups, three disadvantaged players did not buy any rockets. The other four disadvantaged players bought at least one rocket; in fact, two of them overtook the advantaged player in rockets and kept buying more to maintain the buffer of at least one-and-half rockets. Thus, while not precisely described by our hypothesized strategy, it seems that the condition for peace in the Asymmetry treatment is a safe disparity in rocket investment levels.¹⁶ In contrast, of the 21 groups where deactivations occurred, 15 took place when the two players were only half a rocket apart.

Result 2: Deactivations occurred more often in the Asymmetry treatment than the Base treatment. The escalation level was lower in the Asymmetry treatment than the Base treatment. Most deactivation-free groups in the Asymmetry treatment managed to sustain peace by preserving a margin of at least one and half rockets.

3.3 Trade Opportunities

Last, we look at whether trade opportunity alleviated the tension of deactivation in the Asymmetry treatment. Over the 13 rounds, deactivation occurred in 7 out of 26 groups (26.9%) in the Trade treatment, a frequency significantly lower than the Asymmetry treatment (see figure 2; $p = 0.001$, Z test). To elaborate, two advantaged players deactivated their adversary in the first round. Four advantaged players did so after some delay; in all these cases, the disadvantaged players also escalated at least once. In one group, the disadvantaged player overtook and deactivated the advantaged player. As hypothesized, trade opportunities reduce deactivation, but do not eliminate it, as the fear of being deactivated still appears to haunt some participants.

With regard to investment behavior, participants spent on average AU\$3.27 on rockets, a similar amount to that spent in the Asymmetry treatment (see figure 3a; $p = 0.572$, rank-sum test). Among the 19 deactivation-free groups, these participants spent on average AU\$3.63 on rockets (see figure 3b). The distributions of rocket expenditure in the Asymmetry and Trade treatments (for all groups) are similar (see figure 4, $p = 0.132$, Kolmogorov-Smirnov test).

¹⁵ Since there were more deactivations in the Asymmetry treatment than the Base treatment, the effect may not be causal: there might be a selection effect of deactivation-free groups and we could not know what the investment decisions of the deactivated would have been if they were still active. Nevertheless, the substantial drop in the average investment level in the Asymmetry treatment is unlikely to be fully explained by a selection effect.

¹⁶ One could imagine another treatment featuring even more initial power asymmetry, e.g., the initial token endowment for one side being one and a half. This treatment could more directly test this hypothesis for sustaining peace.

Thus, trade opportunities do not increase escalation, implying that participants did not jostle to gain an upper hand in bargaining or in ability to deactivate.

How did these 19 groups manage to sustain peace? 11 groups maintained a safe distance of at least one-and-a-half rockets between each side. Of these, in four groups the advantaged player only bought one rocket in the first round, while the disadvantaged player did not buy any rockets. In another four groups, both sides escalated at least once; in the remaining three groups, the disadvantaged player overtook the advantaged player. The other eight deactivation-free groups kept only half a rocket apart in almost every round. Among these, in four groups neither side bought any rockets, while in the other four both sides escalated almost every round. In sum, the ways that groups sustained peace are much more varied when trade opportunities are introduced. This suggests that rocket investment is less important to the deactivation decision when trade opportunities are available.

Table 2: The determinants of bargaining demand

	<i>Dependent Variable: Bargaining Demand</i>			
	(1)	(2)	(3)	(4)
Constant	3.449*** (0.066)	3.454*** (0.034)	3.418*** (0.068)	3.471*** (0.019)
Lottery Earning	-0.005 (0.007)	0.001 (0.005)	-0.006 (0.007)	0.000 (0.005)
Advantaged Player	0.067 (0.050)	0.052 (0.045)		
Own Token			0.035 (0.021)	0.034* (0.021)
Other's Token			-0.033 (0.021)	-0.030 (0.020)
Controls	Yes	No	Yes	No
Observations	486	486	486	486

Notes: 1) *, **, *** stand for $p < 0.10$, $p < 0.05$ and $p < 0.01$, respectively. 2) Standard errors clustered at the group level are in parentheses. 3) We controlled for gender, age, field of study and round dummies in regressions in columns (1) and (3).

Last, does the trade dynamics have anything to do with arming and deactivation decisions? Thus, we examine bargaining behavior in more detail. It is possible that bargaining failure might frustrate both players and trigger more armaments and deactivations. Or conversely, a player with more armaments (and thus the deactivation power) might try to influence the bargaining process to their advantage. Among the 19 deactivation-free groups, only two had ever failed in bargaining. In fact, out of a total of 228 instances, bargaining failed only four times, and three of these failures were caused by one group. Further, the bargaining outcome is overwhelmingly egalitarian (i.e. both sides ultimately demanded AU\$3.50; 89.3%) and not

correlated with the lottery payoff or the initial endowment (see columns (1) and (2) in Table 2 for random effects regression evidence).¹⁷ We also did not observe much influence of armaments on bargaining outcomes: they are largely not correlated with either one's own armament level or the opponent's armament level, though column (4) in Table 2 suggests marginal evidence that players with more armaments obtained slightly more payoff.

On the other hand, among the 7 groups where deactivation occurred, two had failed in bargaining before the round when deactivation occurred. Of the total of 17 instances of bargaining in these groups, three were failures. While the evidence is only suggestive, it seems to imply that bargaining failures increased the chance of deactivation. Overall, participants almost always reached an agreement in bargaining, which largely prevented the side with more rockets from deactivating the other side, although it did not completely eliminate escalation.

Result 3: Deactivations occurred less often in the Trade treatment than the Asymmetry treatment. Escalation levels were similar in both treatments. There is suggestive evidence that bargaining failures might lead to deactivation.

3.4. Spite or fear?

The Hobbesian trap describes a situation in which two actors may attack each other although neither has an actual desire to harm. This, however, does not mean that attacking is necessarily irrational. Consider Schelling's burglar example. It may turn out that this particular burglar had no intention to kill the homeowner, but this is unknown to the latter. It is sufficient that there is a positive probability for the burglar to be malevolent to rationalize a first strike. This is particularly true in our game since maintaining the ability to defend oneself, i.e. maintaining a balance of power or one's own superiority requires expenses in weapons. In this section we explore the roles of spite (including malevolence due to the desire to save future expenses on rockets) and fear in driving deactivation decisions in our experiment.¹⁸

One clue about spite (or rather lack thereof) would be the number of subjects who never attack, even though in some round they know they are stronger than the opponent. (A spiteful subject who becomes stronger might enjoy waiting and attacking in a later round, but if he never attacks then he is unlikely to be spiteful.) In the Base treatment, asymmetries in the stock of rockets arose in 17 out of 29 pairs. Of those, nine stronger players attacked immediately, four did so later, and only in four groups peace prevailed. Thus, less than a quarter (23.5%) can clearly be identified as not spiteful. In the Hidden treatment, asymmetries arose in 16 out of 29 groups. All remained peaceful, though the interpretation of this observation is difficult as players were not aware of their superiority. In the Asymmetry and Trade treatments, one player by design always had the opportunity to attack. 9 out of 30 groups (Asymmetry) and 19 out of 26 groups

¹⁷ Figures B5 and B6 in Appendix B plot both sides' bargaining demands over time for each group in the first and last round respectively. They show that a group very often quickly reaches a consensus on equal split of the bargaining pie and that the time taken to reach such a consensus tends to be shorter in later rounds.

¹⁸ A previous study (Abbink and de Haan 2014) on the first-strike game comprised a treatment in which spite was the only motive for deactivation. The authors found zero deactivations. This may serve as a prior for our analysis, though the setting was considerably different from ours, as there was no armament to create or prevent deactivations.

(Trade) remained peaceful. Thus, the fraction of non-spiteful players are 30% and 73.1%, respectively. Note that the cost of acting spitefully are much higher in the Trade treatment, since attacking also means losing future bargaining surplus.

The timing of attacks could also provide a clue on whether strikes are motivated by spite or fear. Very late attacks are most likely driven by spite, since the saving on future investments becomes irrelevant, and also the damage done to the recipient is highest. Figure 5 shows the number of attacks in the 13 rounds of the experiment. Overall, very late attacks are rare. Only in the Base treatment, three players attacked in the last round, when spite is the only explanation for an attack. Most attacks occurred early in the game, when material gains (saving of future investments) can explain deactivations.

Overall, deactivation seems to be driven by both fear and spite (particularly when one can save on future investments), though the extreme spiteful type who aims to minimize the opponent's payoff is rarely observed. Cleanly disentangling fear from spite is difficult in the current setting: If investments were costless, arms races would be trivial (though this situation might be a unnatural diagnostic treatment). Alternatively, if no one feared, the desire to save future expenses would be pointless. In the Hidden treatment, spiteful players would have little reason to invest since they can never be sure that attacks are safe. Therefore, the desire to save future expenses should also be limited. However, the observation that the investment rate is around two-thirds seems to contradict that spite is the only motive. While the verdict is still out there, we believe fear must hold some truth in creating arms races.

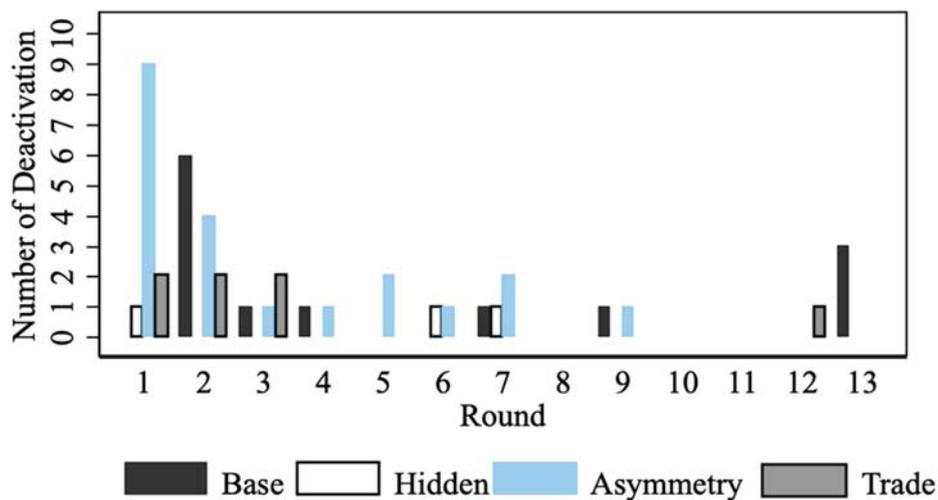


Figure 5. Number of deactivations by round

4. Related literature

This study provides a unified and tractable experimental framework for studying arms races, with specific focus on the role of deterrence, balance of power, and trade in mediating arms races. There exists a vast empirical literature on the relationship between each of these three factors and war. For example, evidence from historical cases shows that the likelihood of

deterrence successes in terms of preventing outbreaks of war depends on a number of factors such as military capability, diplomatic policy and past record of military actions (Huth (1988); see Huth (1999) for a survey). In particular, studies using Militarized Interstate Dispute (MID) data find that symmetric nuclear weapon possession between two states appeared to reduce the odds of war compared to asymmetric possession (Bueno de Mesquita and Riker (1982); Rauchhaus (2009); although see Bell and Miller (2015) for countervailing evidence; see Geller (2017) for a survey). These findings reflect the importance of balance of power in deterrence success, though its effect on escalation and war is still much debated (Geller 1990; Wagner 1994; Bueno de Mesquita, Morrow, and Zorick 1997). Empirical data on the effect of trade on war also have reached mixed conclusions. Some studies find a negative causal link between the frequency of war and countries' bilateral trade (Oneal, Russett, and Berbaum 2003; Hegre, Oneal, and Russett 2010) or multilateral trade among networks of military alliances (Jackson and Nei 2015). However, other studies find a positive relationship (Barbieri 1996; 2002; Martin, Mayer, and Thoenig 2008; Li and Reuveny 2011).

There are several *theoretical* papers studying arms races, to which our games are loosely connected. One strand of this literature focuses on either arming decisions or war decisions. For example, in Baliga and Sjöström's (2004) model, players decide whether or not to arm (with the cost of arming as private information) and the payoff resembles a stag hunt game. The authors show a spiraling effect that leads to the unique Bayesian-Nash equilibrium of both players choosing to arm as long as each party assigns a small but positive probability to their opponent arming as his dominant strategy. Chassang and Miquel (2010) study a setting where players, with exogenous armament levels, decide whether to strike (defect) in a repeated Prisoner's Dilemma game. They show that under complete information, arms have deterrent effects, but under strategic uncertainty, arms can destabilize coordination as players may second-guess each other's predatory incentives and launch a preemptive attack. The other strand of arms race literature studies both arming and war decisions. For example, Jackson and Morelli (2009) develop a repeated two-stage game in which players first simultaneously invest in arms and then choose whether to go to war after observing each other's armament level. They characterize a Markov perfect equilibrium in which each player's armament level and the equilibrium of peace versus war depend on the costs of war. Meirowitz and Sartori (2008) study a model with a similar time line, although in their model armament levels are not observable at the time of the war decisions. Strategic uncertainty about military strength risks war; however, players may have incentives to keep their strength unknown and unpredictable.¹⁹ In both models, deterrence is achieved through different kinds of randomization in arming decisions. Interestingly, the main difference between these two models echoes that between our Base and Hidden treatments, namely whether one's military capacity is known to the opponent. However, in our game, this information is exogenous and the standard theory does not produce different predictions regarding arming decisions.

Besides Jackson and Morelli (2009) and Meirowitz and Sartori (2008), other models in this vein include Brito and Intriligator (1985), Morrow (1989), Powell (1993), Kydd (1997, 2000),

¹⁹ See also Baliga and Sjöström (2008) on strategic uncertainty in conflicts.

and Slantchev (2005). In some of these models (Brito and Intriligator 1985; Kydd 2000), deterrence is not understood in the *strict sense* of MAD, in which arms races are seen as a political means of improving a state's bargaining position. More generally, in all these models, there is always some gain from winning a war such that the defensive side, which possesses valuable resources, has incentives to deter the offensive side from attacking (as is also true in all the experimental competition games discussed earlier). This, too, is not compatible with the strict definition of MAD, which does not presume conflicts over resources. Borrowing the phrase used by Kydd (2000), we are studying deterrence from "the dove perspective."²⁰ Yet it emerges that resource waste and conflict propensity are not negligible even in this seemingly unthreatening environment.

A competing theory to MAD in explaining arms races is the spiral model, which holds that the fear of being unexpectedly attacked leads countries to increase arms stockpiles (Jervis 1976). This increase in arms in turn generates more fear on the other side, driving further arms buildup. The spiral model suggests that psychological bias is responsible for preemptive attacks, with actors mistakenly assuming that the other side's arms are a sign of aggression while the purpose of one's own arms as self-defense is transparent. However, unlike the logic of MAD, the spiral model places little weight on credible retaliation and therefore largely neglects the possibility of peace even under great mutual fear.

5. Conclusion

We introduce an experimental paradigm to study factors modulating arms races and conflict in a highly stylized setting. We find that the possibility of encountering a highly armed opponent is a strong deterrent to conflict, even (and especially) if the opponent's strength is unknown. Open arms races, in which exact power distribution is visible to all, often lead to conflict when opportunity arises. In both cases, we observe high levels of investment in arms. This indicates that investments in arms are made mainly for the purpose of insurance rather than deterrence. Interestingly, the finding that secrecy helps peace making is seemingly at odds with Fearon's (1995) account that lack of information about an opponent's capability might explain costly decisions for war. One reason, Fearon argued, is that secrecy prevents rational players to locate a mutually preferable settlement since players may have incentives to misrepresent their private information. The ability to communicate private information is indeed a missing element in the current study. A situation where players can strategically choose either secrecy or openness about their weapons and communicate with each other may shed more light on this discrepancy.

Our results also show the importance of a balance of power. When power asymmetry arises during an arms race (in the Base treatment), conflict becomes very likely (13 out of 17 groups, 76%). In a treatment with perpetual power asymmetry, conflict is also very frequent (70%).

²⁰ In economics literature, deterrence is often studied in market entry games (Rapoport et al. 1998; Camerer and Lovo 1999; Duffy and Hopkins 2005; Morgan, Orzen, and Sefton 2012). Our game is distinct from this literature in that deterrence is symmetrical and can lead to escalation.

Escalating arms races, on the other hand, become less frequent even in peaceful groups when there is an imbalance of power.

Trade, as represented by a bargaining opportunity in our experiment, is very effective in reducing destructive conflict as well as the occurrence of arms races. Bargaining almost never fails in our experiment, and hence creates profitable exchanges that build up enough trust for the majority of players to find it unnecessary to build up arsenals as insurance.

It is worth noting that all our results concerning arms races and deactivation contradict the subgame perfect equilibrium, which predicts neither. Failures of backward induction in predicting behavior are also famously found in extensive form games such as centipede games (McKelvey and Palfrey, 1992). This has led game theorists to realize the importance of *common knowledge of rationality*, as well as the difficulty of achieving this state of rationality in practice (Aumann 1992; 1995). The pervasive arms races observed in the Base treatment are a testament to this difficulty. In reality, if this game were to play out between state leaders backed by their “super-rational” think tanks, the situation might be sufficiently close to common knowledge of rationality for arms-free peace to obtain. Nonetheless, as Aumann (1995) stresses, even the smallest departure from common knowledge of rationality may induce rational players to deviate significantly from SPE play. It therefore remains a nontrivial task to eradicate arms races and deactivation even among rational players.

In our game, we have used the notion of fear to capture a subject’s psychological reaction to the possible failure of common knowledge of rationality. To sketch a mechanism that fear could cause escalation and deactivation in the Base treatment, we consider two player types: a strategic type and a non-strategic type. The strategic type perfectly understands the logic of the game and thus will neither escalate nor deactivate if he believes the other side is of the same type. The non-strategic type can have plural motives: he might be one of the spiteful types who strive to get ahead of others in earnings or one of the uncertainty averse types who do not like any degree of insecurity about his payoff and thus buy rockets to not fall behind others. Therefore, the non-strategic type with either motive will escalate and probably deactivate the other side when possible, irrespective of his belief of the other’s type. Now if the strategic type has the slightest belief that he is paired with a non-strategic type (or the belief that the other side believes she is paired with a non-strategic type, *ad infinitum*), then it will only be rational for him to also escalate. Hence, even when both sides are the strategic type, the logic of fear could easily set off an arms race. We will leave rigorous theoretical modeling for future research.

Beyond our present findings, our experimental framework is flexible enough to capture more nuanced views about arms races, nuclear wars and other types of conflict that researchers have in mind. For example, we can incorporate endogenous economic growth (Romer 1986) by allowing the supplemental gains from lotteries or trade to be a convex function of the current wealth level. Alternatively, to investigate endogenous technological growth of weapons build-up due to greater economic resources, we may set the cost of building new rockets as a function of the current wealth level. In addition, it might be interesting to study arms races involving conventional weapons by allowing the destructive power of war to increase smoothly with the

amount of accumulated weapons. Furthermore, we may add a “de-escalate” button to allow parties to destroy their existing rockets to see how this affects peace and conflict. Finally, since we have observed fewer deactivations in the Hidden treatment, we may investigate whether players would strategically choose to keep their arms secret or reveal them to deter others’ aggression. The voluntary decision of revealing arsenal represents an intermediate case between the idealized extremes of openness and secrecy and is probably closer to the real world. Will a player initially commit to a policy of secrecy and revert to openness once he believes he is strong enough to deter an attack? We leave this question for future research. In summary, we believe our game provides a fertile platform for testing various hypotheses with regards to arms races and conflict.

Of course, our study has limitations. The most obvious one is the small scale compared to a nuclear arms race in the real world. Deactivating a fellow subject is a cruel act by the standards of a laboratory session, but it is still not the same as annihilating entire countries. Nevertheless, we believe that the study of these situations in the laboratory can give us valuable insights. In this area, we hope that we will never have a sufficient number of real-world observations for an econometric analysis (and once we do, there is probably no one left to do the analysis). Nonetheless, we need to devise institutions that minimize the risk of arms races or, worse, the outbreak of conflict. Understanding the factors affecting the behavior of our institutions is paramount in this endeavor, and the laboratory, with its controlled, replicable and damage-free environment can be a helpful tool.

References

- Abbink, Klaus, and Jordi Brandts. 2008. “24. Pricing in Bertrand Competition with Increasing Marginal Costs.” *Games and Economic Behavior* 63 (1): 1–31.
- Abbink, Klaus, and Thomas de Haan. 2014. “Trust on the Brink of Armageddon: The First-Strike Game.” *European Economic Review* 67 (April): 190–96.
- Abbink, Klaus, and Abdolkarim Sadrieh. 2009. “The Pleasure of Being Nasty.” *Economics Letters* 105 (3): 306–8.
- Aumann, Robert J. 1992. “Irrationality in Game Theory.” In *Economic Analysis of Markets and Games: Essays in Honor of Frank Hahn*, edited by Partha Dasgupta, Douglas Gale, Oliver Hart, and Eric Maskin. Cambridge, MA: MIT Press.
- . 1995. “Backward Induction and Common Knowledge of Rationality.” *Games and Economic Behavior* 8: 6–19.
- Baliga, Sandeep, and Tomas Sjöström. 2004. “Arms Races and Negotiations.” *Review of Economic Studies* 71 (2): 351–69.
- . 2008. “Strategic Ambiguity and Arms Proliferation.” *Journal of Political Economy* 116 (6): 1023–57.
- Barbieri, Katherine. 1996. “Economic Interdependence: A Path to Peace or a Source of Interstate Conflict?” *Journal of Peace Research* 33 (1): 29–49.
- . 2002. *The Liberal Illusion: Does Trade Promote Peace?* Ann Arbor: University of Michigan Press.

- Battalio, Raymond, Larry Samuelson, and John Van Huyck. 2001. "Optimization Incentives and Coordination Failure in Laboratory Stag Hunt Games." *Econometrica* 69 (3): 749–64.
- Bell, Mark S., and Nicholas L. Miller. 2015. "Questioning the Effect of Nuclear Weapons on Conflict." *Journal of Conflict Resolution* 59 (1): 74–92.
- Bolle, Friedel, Jonathan H.W. Tan, and Daniel John Zizzo. 2014. "Vendettas." *American Economic Journal: Microeconomics* 6 (2): 93–130.
- Boone, Jan, María Jose Larraín Aylwin, Wieland Müller, and Amrita Ray Chaudhuri. 2012. "Bertrand Competition with Asymmetric Costs: Experimental Evidence." *Economics Letters* 117 (1): 134–37.
- Brito, Dagobert L., and Michael D. Intriligator. 1985. "Conflict, War and Redistribution." *American Political Science Review* 79 (4): 943–57.
- Bueno de Mesquita, Bruce, James D. Morrow, and Ethan R. Zorick. 1997. "Capabilities, Perception, and Escalation." *American Political Science Review* 91 (1): 15–27.
- Bueno de Mesquita, Bruce, and William H. Riker. 1982. "An Assessment of the Merits of Selective Nuclear Proliferation." *Journal of Conflict Resolution* 26 (2): 283–306.
- Camerer, Colin, and Dan Lovallo. 1999. "Overconfidence and Excess Entry: An Experimental Approach." *American Economic Review* 89 (1): 306–18.
- Chassang, Sylvain, and Gerard Padró I. Miquel. 2010. "Conflict and Deterrence under Strategic Risk." *Quarterly Journal of Economics* 125 (4): 1821–58.
- Chuah, Swee-Hoon, Robert Hoffmann, and Jeremy Lerner. 2014. "Elicitation Effects in a Multi-Stage Bargaining Experiment." *Experimental Economics* 17 (2): 335–45.
- Cracau, Daniel, and Abdolkarim Sadrieh. 2016. "Coexistence of Small and Dominant Firms in Bertrand Competition: Judo Economics in the Lab." *Journal of Institutional and Theoretical Economics* 172 (4): 665–93.
- Duffy, John, and Ed Hopkins. 2005. "Learning, Information, and Sorting in Market Entry Games: Theory and Evidence." *Games and Economic Behavior* 51 (1): 31–62.
- Dufwenberg, Martin, and Uri Gneezy. 2000. "Price Competition and Market Concentration: An Experimental Study." *International Journal of Industrial Organization* 18 (1): 7–22.
- Fearon, James D. 1995. "Rationalist Explanations for War." *International Organization* 49 (3): 379–414.
- Fischbacher, Urs. 2007. "Z-Tree: Zurich Toolbox for Ready-Made Economic Experiments." *Experimental Economics* 10 (2): 171–78.
- Fitzgerald, F. 2000. *Way Out There in the Blue: Reagan, Star Wars, and the End of the Cold War*. New York, NY: Simon & Schuster.
- Garfinkel, Michelle R., and Stergios Skaperdas. 2000. "Conflict without Misperceptions or Incomplete Information." *Journal of Conflict Resolution* 44 (6): 793–807.
- Geller, Daniel S. 1990. "Nuclear Weapons, Deterrence, and Crisis Escalation." *Journal of Conflict Resolution* 34 (2): 291–310.
- . 2017. "Nuclear Weapons and International Conflict: Theories and Empirical Evidence." In *Oxford Research Encyclopedia of Politics*. Oxford University Press.
- Hegre, Håvard, John R. Oneal, and Bruce Russett. 2010. "Trade Does Promote Peace: New

- Simultaneous Estimates of the Reciprocal Effects of Trade and Conflict.” *Journal of Peace Research* 47 (6): 763–74.
- Hoffman, David. 1999. “I Had A Funny Feeling in My Gut.” *Washington Post Foreign Service*, 1999.
- Holt, Charles A., and Susan K. Laury. 2002. “Risk Aversion and Incentive Effects.” *American Economic Review* 92 (5): 1644–55.
- Huth, Paul K. 1988. “Extended Deterrence and the Outbreak of War.” *American Political Science Review* 82 (2): 423.
- . 1999. “DETERRENCE AND INTERNATIONAL CONFLICT: Empirical Findings and Theoretical Debates.” *Annual Review of Political Science* 2 (1): 25–48.
- Jackson, Matthew O., and Massimo Morelli. 2009. “Strategic Militarization, Deterrence and Wars.” *Quarterly Journal of Political Science* 4 (4): 279–313.
- Jackson, Matthew O., and Stephen Nei. 2015. “Networks of Military Alliances, Wars, and International Trade.” *Proceedings of the National Academy of Sciences* 112 (50): 15277–84.
- Jervis, Robert. 1976. *Perception and Misperception in International Politics*. Princeton, NJ: Princeton University Press.
- Kant, Immanuel. 1795. *Perpetual Peace: A Philosophical Sketch*. *The Grotius Society Publications, No. 7*. Helen O’Brien, Trans. London: Sweet and Maxwell.
- Kydd, Andrew. 1997. “Game Theory and the Spiral Model.” *World Politics* 49 (03): 371–400.
- . 2000. “Arms Races and Arms Control: Modeling the Hawk Perspective.” *American Journal of Political Science* 44 (2): 222–38.
- Lacomba, Juan A., Francisco Lagos, Ernesto Reuben, and Frans van Winden. 2014. “On the Escalation and De-Escalation of Conflict.” *Games and Economic Behavior* 86 (July): 40–57.
- Levitt, Steven D., John A. List, and Sally E. Sadoff. 2011. “Checkmate: Exploring Backward Induction among Chess Players.” *American Economic Review* 101 (2): 975–90.
- Li, Quan, and Rafael Reuveny. 2011. “Does Trade Prevent or Promote Interstate Conflict Initiation?” *Journal of Peace Research* 48 (4): 437–53.
- Mansfield, Edward. 1994. *Power, Trade, and War*. Princeton, NJ: Princeton University Press.
- Martin, Philippe, Thierry Mayer, and Mathias Thoenig. 2008. “Make Trade Not War?” *Review of Economic Studies* 75 (3): 865–900.
- McKelvey, Richard D., and Thomas R. Palfrey. 1992. “An Experimental Study of the Centipede Game.” *Econometrica* 4 (60): 803–36.
- Meirowitz, Adam, and Anne E. Sartori. 2008. “Strategic Uncertainty as a Cause of War.” *Quarterly Journal of Political Science* 3 (4): 327–52.
- Morgan, John, Henrik Orzen, and Martin Sefton. 2012. “Endogenous Entry in Contests.” *Economic Theory* 51 (2): 435–63.
- Morrow, James D. 1989. “A Twist of Truth: A Reexamination of the Effects of Arms Races on the Occurrence of War.” *Journal of Conflict Resolution* 33 (3): 500–529.

- Nagel, Rosemarie, and Fang Fang Tang. 1998. "Experimental Results on the Centipede Game in Normal Form: An Investigation on Learning." *Journal of Mathematical Psychology* 42 (3): 356–84.
- Norris, Robert S., and Hans M. Kristensen. 2010. "Global Nuclear Weapons Inventories, 1945–2010." *Bulletin of the Atomic Scientists* 66 (4): 77–83.
- Oneal, John R., Bruce Russett, and Michael L. Berbaum. 2003. "Causes of Peace: Democracy, Interdependence, and International Organizations, 1885-1992." *International Studies Quarterly* 47 (3): 371–93.
- Palacios-Huerta, Ignacio, and Oscar Volij. 2009. "Field Centipedes." *American Economic Review* 99 (4): 1619–35.
- Pinker, Steven. 2011. *The Better Angels of Our Nature: Why Violence Has Declined*. New York, NY: Viking Books.
- Powell, Robert. 1993. "Guns, Butter, and Anarchy." *American Political Science Review* 87 (1): 115–32.
- Rankin, Frederick W., John B. Van Huyck, and Raymond C. Battalio. 2000. "Strategic Similarity and Emergent Conventions: Evidence from Similar Stag Hunt Games." *Games and Economic Behavior* 32 (2): 315–37.
- Rapoport, Amnon, Darryl A. Seale, Ido Erev, and James A. Sundali. 1998. "Equilibrium Play in Large Group Market Entry Games." *Management Science* 44 (1): 119–41.
- Rauchhaus, Robert. 2009. "Evaluating the Nuclear Peace Hypothesis: A Quantitative Approach." *Journal of Conflict Resolution* 53 (2): 258–77.
- Romer, Paul M. 1986. "Increasing Returns and Long-Run Growth." *Journal of Political Economy* 94 (5): 1002–37.
- Schelling, Thomas C. 1960. *The Strategy of Conflict*. Cambridge, MA: Harvard University Press.
- Siegel, Ron. 2009. "All-Pay Contests." *Econometrica* 77 (1): 71–92.
- Slantchev, Branislav L. 2005. "Military Coercion in Interstate Crises." *American Political Science Review* 99 (4): 533–47.
- Wagner, R. Harrison. 1994. "Peace, War, and the Balance of Power." *American Political Science Review* 88 (3): 593–607.
- Walsh, Walter. 1910. "Europe's Optical Illusion." *The Advocate of Peace* 72 (7): 166–69.
- World Bank. 2018. "Military Expenditure (% of GDP)." 2018. <https://data.worldbank.org/indicator/MS.MIL.XPND.GD.ZS?locations=US>.

Appendix A. Experiment Instructions

[Note: this appendix presents the experimental instructions for all treatment. Below, we first present the instructions for the Base, Hidden and Asymmetry treatments, along with their main screenshots. Capitalized texts are the additional words for the Hidden Treatment and italicized texts are the additional words for the Asymmetry Treatment.]

Welcome to this experiment on decision-making. You have earned \$5 for showing up on time. Please read the following instructions carefully. During the experiment, you will be asked to make a number of decisions. Your decisions and the decisions of other participants will determine your cash earnings. Your anonymity is ensured for the decisions you make. The experiment will consist of 12 rounds. No communication during the experiment is allowed. Final payment will be rounded to the nearest 10 cents. If you have a question, please raise your hand.

Matching: In each of the 12 rounds all participants will be matched in pairs. The pairs will be the same for all rounds. So you will be matched with another person and you will stay matched to this person throughout the whole experiment. You will NOT be rematched. [*In each pair, the computer will randomly select one person to be person X and the other person Y.*]

Lotteries: In each round, you can earn money by choosing between lotteries: you will see on your screen two lotteries displayed, lottery A and lottery B. You can choose either lottery A or B. After you made the choice, you can spin the wheel and the chosen lottery will be played out. You will see on your screen the lottery outcome and the amount earned will be added to your total earnings. You have a start fund of \$5.

Tokens: Apart from the lottery decisions you have to make in each round, both you and the person you are matched with will have the option to invest \$1 to buy a token *in each round*. That is, in each round, you will decide whether or not to buy a \$1 token. [*Person X starts with 0.5 token and person Y starts with 0 token in round 1.*] The money spend in tokens is non-refundable and tokens are accumulated from round to round. Note that the token you have bought in one round will only be effective for the next round. For example, if you already have 2 tokens at the beginning of a round and decide to buy one more token in this round, it means you only have 2 tokens for this round, and you will have 3 tokens for the next round.

Deactivation: In each round, you will see your own and the other person's token balance, that is, the numbers of tokens bought up till this round. [IN EACH ROUND, YOU WILL ONLY SEE THE NUMBERS OF TOKENS YOU BOUGHT UP TILL THIS ROUND, AND YOU WILL NOT KNOW OTHER PERSON'S TOKEN.] Both of you can decide whether to press the "deactivate" or "not deactivate" button. Pressing the "deactivate" button costs \$1.5. You only have one chance to press the "deactivate" button.

Depending on the current token balance, pressing the "deactivate" button has different consequences for the other person and yourself.

If you have more tokens than the other person and you press “deactivate”, the other person will be deactivated, meaning:

1. All earnings so far (including the start fund) of the other person will be set to \$0.
2. All payoffs in the future lotteries for the other person will be divided by 10.
3. The other person will not be able to deactivate you in future rounds.

[Note: the following condition about equal number of tokens is irrelevant and therefore absent in the Asymmetry treatment.]

If you and the other person have equal number of tokens and you press “deactivate”, both the other person and you will be deactivated, meaning:

1. All earnings so far (including the start fund) for the other person and you will be set to \$0.
2. All payoffs in the future lotteries for the other person and you will be divided by 10.

If you have fewer tokens than the other person and you press “deactivate”, only you will be deactivated, meaning:

1. All earnings so far (including the start fund) for you will be set to \$0.
2. All payoffs in the future lotteries for you will be divided by 10.
3. You will not be able to deactivate the other person in future rounds.

In summary, once a person is deactivated either directly by the other person or indirectly by his/her own action, it will wipe out his/her total earnings from all previous rounds and decrease the potential earnings from future rounds by 90%. Furthermore, a deactivated person will not be able to deactivate the other person.

There are two further rules:

- a) If you and the other person press the “deactivate” button in the same round, the rules described above apply to both you and the other person. But no one can be deactivated twice.
- b) If no one has been deactivated after the last (12th) round, you will have the option to press the “deactivate” button one last time though there will be no lottery decision to make any more.

The 12 lottery choices

Choice if you have not been deactivated			Choice if you have been deactivated		
No.	Lottery A	Lottery B	No.	Lottery A	Lottery B
1	15% of \$2.70, 85% of \$2.30	15% of \$3.40, 85% of \$1.70	1	15% of \$0.27, 85% of \$0.23	15% of \$0.34, 85% of \$0.17
2	30% of \$2.70, 70% of \$2.30	30% of \$3.40, 70% of \$1.70	2	30% of \$0.27, 70% of \$0.23	30% of \$0.34, 70% of \$0.17
3	45% of \$2.70, 55% of \$2.30	45% of \$3.40, 55% of \$1.70	3	45% of \$0.27, 55% of \$0.23	45% of \$0.34, 55% of \$0.17
4	60% of \$2.70, 40% of \$2.30	60% of \$3.40, 40% of \$1.70	4	60% of \$0.27, 40% of \$0.23	60% of \$0.34, 40% of \$0.17
5	75% of \$2.70, 25% of \$2.30	75% of \$3.40, 25% of \$1.70	5	75% of \$0.27, 25% of \$0.23	75% of \$0.34, 25% of \$0.17
6	90% of \$2.70, 10% of \$2.30	90% of \$3.40, 10% of \$1.70	6	90% of \$0.27, 10% of \$0.23	90% of \$0.34, 10% of \$0.17
7	15% of \$2.60, 85% of \$2.40	15% of \$3.30, 85% of \$1.80	7	15% of \$0.26, 85% of \$0.24	15% of \$0.33, 85% of \$0.18
8	30% of \$2.60, 70% of \$2.40	30% of \$3.30, 70% of \$1.80	8	30% of \$0.26, 70% of \$0.24	30% of \$0.33, 70% of \$0.18
9	45% of \$2.60, 55% of \$2.40	45% of \$3.30, 55% of \$1.80	9	45% of \$0.26, 55% of \$0.24	45% of \$0.33, 55% of \$0.18
10	60% of \$2.60, 40% of \$2.40	60% of \$3.30, 40% of \$1.80	10	60% of \$0.26, 40% of \$0.24	60% of \$0.33, 40% of \$0.18
11	75% of \$2.60, 25% of \$2.40	75% of \$3.30, 25% of \$1.80	11	75% of \$0.26, 25% of \$0.24	75% of \$0.33, 25% of \$0.18
12	90% of \$2.60, 10% of \$2.40	90% of \$3.30, 10% of \$1.80	12	90% of \$0.26, 10% of \$0.24	90% of \$0.33, 10% of \$0.18

You make these 12 lottery choices in random order and you will not know the other person's earnings. If you have not been deactivated, then you choose lotteries from the left-hand side of the table. If you have been deactivated, then you choose the corresponding lotteries from the right-hand side of the table, where all payoffs are divided by 10. Likewise, if you deactivate the other person, then the other person from then on chooses lotteries from the right-hand side of the table, where all payoffs are divided by 10.

This completes the instruction. Before we begin the experiment, to make sure that every participant understands the instructions, please answer several review questions on your screen.

Screenshot for the Hidden Treatment

The screenshot displays the experiment interface. At the top, it shows 'Round 3' and 'Remaining Time (sec): 76'. The main area is divided into several sections:

- Lottery A:** 45% \$2.70, 55% \$2.30
- Lottery B:** 45% \$3.40, 55% \$1.70
- SCOREBOARD:**
 - Start Fund: 5.00
 - Total Lottery Earnings So Far: 5.10
 - Total Investment So Far: 2.00
 - Deactivation Expense So Far: 0.00
 - Total Earnings So Far: 8.10
 - Show-up fee: 5.00
- Brief Instruction:**
 - Choose a lottery between A and B by clicking on the lottery itself. Then spin a wheel to determine the outcome of the lottery.
 - Choose whether to buy a \$1 token. The token you buy this round will be effective for the next round.
 - You can choose to deactivate the other person. However, how the deactivation will affect the other person and yourself depends on the numbers of tokens you and the other person have so far bought respectively (see rules and consequences of deactivation in the instructions). Pressing the "deactivate" button costs \$1.5 no matter what.
 - There will be 12 rounds. After the last round, each player may have to make the deactivation decision one last time.

At the bottom, there is a 'Your Tokens' bar chart showing 2 tokens (represented by red blocks) and buttons for 'BUY A TOKEN', 'DO NOT BUY', 'DEACTIVATE', and 'NOT DEACTIVATE'. A 'CONFIRM' button is also visible between the lottery options.

Screenshot for the Asymmetry Treatment

The screenshot displays the experimental interface for the Asymmetry Treatment. At the top, it shows 'Round 2' and 'Remaining Time (sec): 85'. The main area is divided into several sections:

- Lottery A:** 90% \$2.70, 10% \$2.30
- Lottery B:** 90% \$3.40, 10% \$1.70
- SCOREBOARD:**
 - Start Fund: 5.00
 - Total Lottery Earnings So Far: 1.70
 - Total Investment So Far: 1.00
 - Deactivation Expense So Far: 0.00
 - Total Earnings So Far: 5.70
 - Show-up fee: 5.00
- Token Counts:** Two vertical bar charts show 'Your Tokens' (red) and 'Other's Tokens' (black). The y-axis ranges from 1 to 12. 'Your Tokens' is currently at 2, and 'Other's Tokens' is at 1.
- Buttons:** 'BUY A TOKEN', 'DO NOT BUY', 'DEACTIVATE', and 'NOT DEACTIVATE' are visible.
- Brief Instruction:**
 - 1) Choose a lottery between A and B by clicking on the lottery itself. Then spin a wheel to determine the outcome of the lottery.
 - 2) Choose whether to buy a \$1 token. The token you buy this round will be effective for the next round.
 - 3) You can choose to deactivate the other person. However, how the deactivation will affect the other person and yourself depends on the numbers of tokens you and the other person have so far bought respectively (see rules and consequences of deactivation in the instructions). Pressing the "deactivate" button costs \$1.5 no matter what.
 - 4) There will be 12 rounds. After the last round, each player may have to make the deactivation decision one last time.

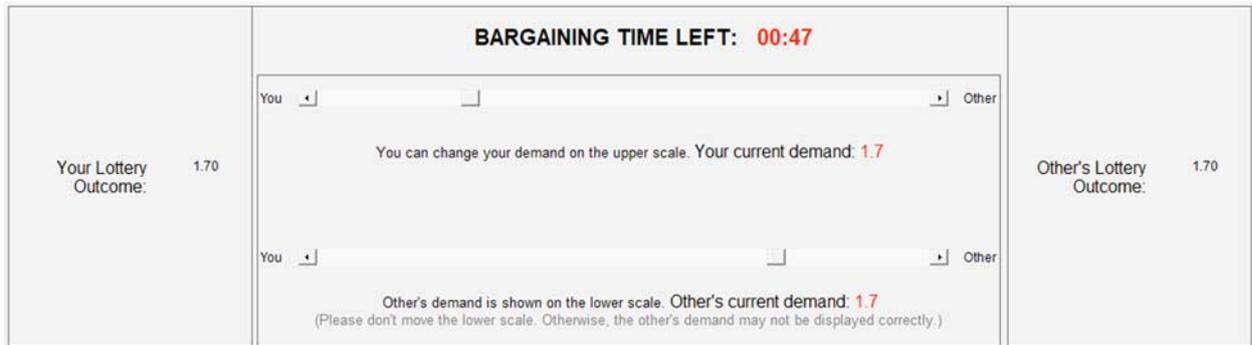
[Note: below we present the experimental instructions for the Trade treatment.]

Welcome to this experiment on decision-making. You have earned \$5 for showing up on time. Please read the following instructions carefully. During the experiment, you will be asked to make a number of decisions. Your decisions and the decisions of other participants will determine your cash earnings. Your anonymity is ensured for the decisions you make. The experiment will consist of 12 rounds. No communication during the experiment is allowed. Final payment will be rounded to the nearest 10 cents. If you have a question, please raise your hand.

Matching: In each of the 12 rounds all participants will be matched in pairs. The pairs will be the same for all rounds. So you will be matched with another person and you will stay matched to this person throughout the whole experiment. You will NOT be rematched. In each pair, the computer will randomly select one person to be Person X and the other Person Y.

Lotteries: In each round, you can earn money by choosing between lotteries: you will see on your screen two lotteries displayed, lottery A and lottery B. You can choose either lottery A or B. After you made the choice, you can spin the wheel and the chosen lottery will be played out.

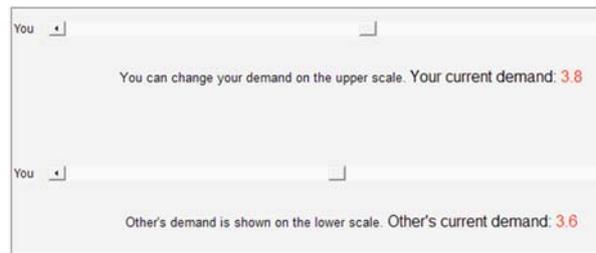
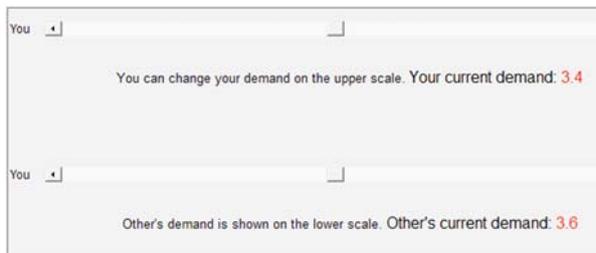
Bargaining: After both you and your pair have made the lottery decisions, you will be told both your and your pair's lottery earnings. Then, the pair will bargain over a fixed amount of money of \$7 by clicking on a scale from \$0 to \$7 (see figure below). In 60 seconds, you can choose your demand on the upper scale and you can see your pair's demand in real time on the lower scale. On both scales, the amounts from the cursor to the left end represents your payoff and the amounts from the cursor to the right end represents your pair's payoff. (Please don't move the lower scale, which displays the other's current demand. It is not possible for you to change the other's demand by moving this scale.)



After the 60 seconds, if the sum of your and your pair’s demands is no greater than \$7 (i.e. a deal is made, see, e.g., figure A1), then each person’s payoff equals to the amounts he/she demands in the bargaining. For example, in figure 2, in this round, you receive \$3.4 and your pair receives \$3.6. If the sum is greater than \$7 (i.e. no deal is made, see, e.g., figure A2), then each person’s payoff equals to his/her own lottery earnings. The payoff earned will be added to each person’s total earnings.

Figure A1

Figure A2



Tokens: Apart from the lottery decisions you have to make in each round, both you and the person you are matched with will have the option to invest \$1 to buy a token *in each round*. That is, in each round, you will decide whether or not to buy a \$1 token. Person X starts with 0.5 token and person Y starts with 0 token in round 1. The money spend in tokens is non-refundable and tokens are accumulated from round to round. Note that the token you have bought in one round will only be effective for the next round. For example, if you already have 2 tokens at the beginning of a round and decide to buy one more token in this round, it means you only have 2 tokens for this round, and you will have 3 tokens for the next round.

Deactivation: In each round, you will see your own and the other person’s token balance, that is, the numbers of tokens bought up till this round. Both of you can decide whether to press the “deactivate” or “not deactivate” button. Pressing the “deactivate” button costs \$1.5. You only have one chance to press the “deactivate” button.

Depending on the current token balance, pressing the “deactivate” button has different consequences for the other person and yourself.

If you have more tokens than the other person and you press “deactivate”, the other person will be deactivated, meaning:

1. All earnings so far (including the start fund) of the other person will be set to \$0.
2. All payoffs in the future lotteries for the other person will be divided by 10.

- The other person will not be able to deactivate you in future rounds.

If you have fewer tokens than the other person and you press “deactivate”, only you will be deactivated, meaning:

- All earnings so far (including the start fund) for you will be set to \$0.
- All payoffs in the future lotteries for you will be divided by 10.
- You will not be able to deactivate the other person in future rounds.

In summary, once a person is deactivated either directly by the other person or indirectly by his/her own action, it will wipe out his/her total earnings from all previous rounds and decrease the potential earnings from future rounds by 90%. Furthermore, a deactivated person will not be able to deactivate the other person.

There are two further rules:

- If you and the other person press the “deactivate” button in the same round, the rules described above apply to both you and the other person. But no one can be deactivated twice.
- If no one has been deactivated after the last (12th) round, you will have the option to press the “deactivate” button one last time though there will be no lottery decision to make any more.

The 12 lottery choices

Choice if you have not been deactivated			Choice if you have been deactivated		
No.	Lottery A	Lottery B	No.	Lottery A	Lottery B
1	15% of \$2.70, 85% of \$2.30	15% of \$3.40, 85% of \$1.70	1	15% of \$0.27, 85% of \$0.23	15% of \$0.34, 85% of \$0.17
2	30% of \$2.70, 70% of \$2.30	30% of \$3.40, 70% of \$1.70	2	30% of \$0.27, 70% of \$0.23	30% of \$0.34, 70% of \$0.17
3	45% of \$2.70, 55% of \$2.30	45% of \$3.40, 55% of \$1.70	3	45% of \$0.27, 55% of \$0.23	45% of \$0.34, 55% of \$0.17
4	60% of \$2.70, 40% of \$2.30	60% of \$3.40, 40% of \$1.70	4	60% of \$0.27, 40% of \$0.23	60% of \$0.34, 40% of \$0.17
5	75% of \$2.70, 25% of \$2.30	75% of \$3.40, 25% of \$1.70	5	75% of \$0.27, 25% of \$0.23	75% of \$0.34, 25% of \$0.17
6	90% of \$2.70, 10% of \$2.30	90% of \$3.40, 10% of \$1.70	6	90% of \$0.27, 10% of \$0.23	90% of \$0.34, 10% of \$0.17
7	15% of \$2.60, 85% of \$2.40	15% of \$3.30, 85% of \$1.80	7	15% of \$0.26, 85% of \$0.24	15% of \$0.33, 85% of \$0.18
8	30% of \$2.60, 70% of \$2.40	30% of \$3.30, 70% of \$1.80	8	30% of \$0.26, 70% of \$0.24	30% of \$0.33, 70% of \$0.18
9	45% of \$2.60, 55% of \$2.40	45% of \$3.30, 55% of \$1.80	9	45% of \$0.26, 55% of \$0.24	45% of \$0.33, 55% of \$0.18
10	60% of \$2.60, 40% of \$2.40	60% of \$3.30, 40% of \$1.80	10	60% of \$0.26, 40% of \$0.24	60% of \$0.33, 40% of \$0.18
11	75% of \$2.60, 25% of \$2.40	75% of \$3.30, 25% of \$1.80	11	75% of \$0.26, 25% of \$0.24	75% of \$0.33, 25% of \$0.18
12	90% of \$2.60, 10% of \$2.40	90% of \$3.30, 10% of \$1.80	12	90% of \$0.26, 10% of \$0.24	90% of \$0.33, 10% of \$0.18

You make these 12 lottery choices in random order and you will not know the other person’s earnings. If you have not been deactivated, then you choose lotteries from the left-hand side of the table. If you have been deactivated, then you choose the corresponding lotteries from the right-hand side of the table, where all payoffs are divided by 10. Likewise, if you deactivate the other person, then the other person from then on chooses lotteries from the right-hand side of the table, where all payoffs are divided by 10.

This completes the instruction. Before we begin the experiment, to make sure that every participant understands the instructions, please answer several review questions on your screen.

Appendix B. Additional Figures

[Note: this appendix shows additional figures. Figures B1-B4 show escalation in tokens and timing of deactivation per group per treatment, respectively. Whenever a “Deactivated” symbol falls upon a player’s escalation path, it means that this player has been deactivated from that round onward. Figures B5-B6 depict bargaining processes for each active group in the first and last rounds of the Trade treatment. Most groups agree on equally splitting the bargaining pie at the end of the bargaining.]

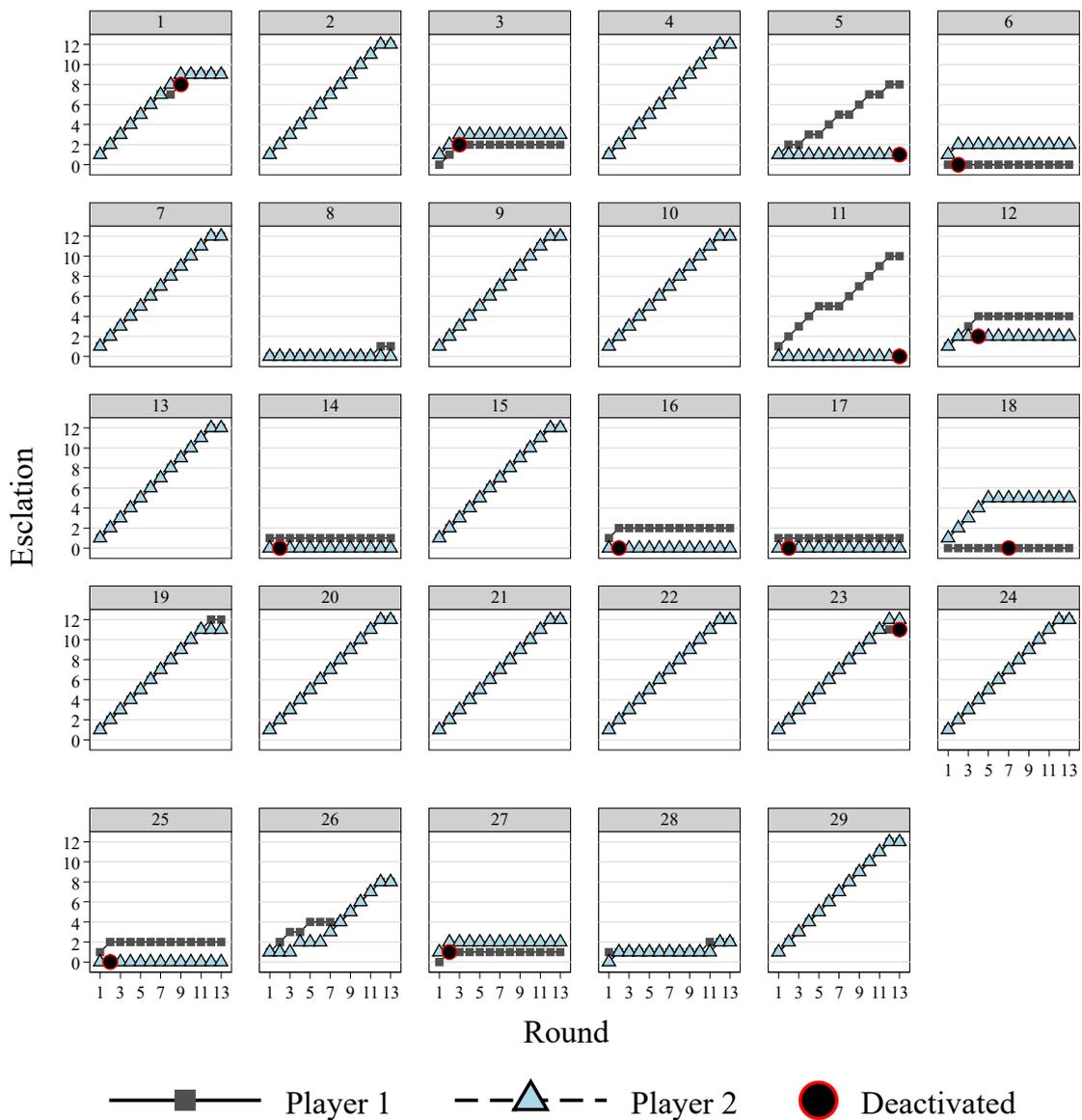


Figure B1: Group level results in the Base treatment

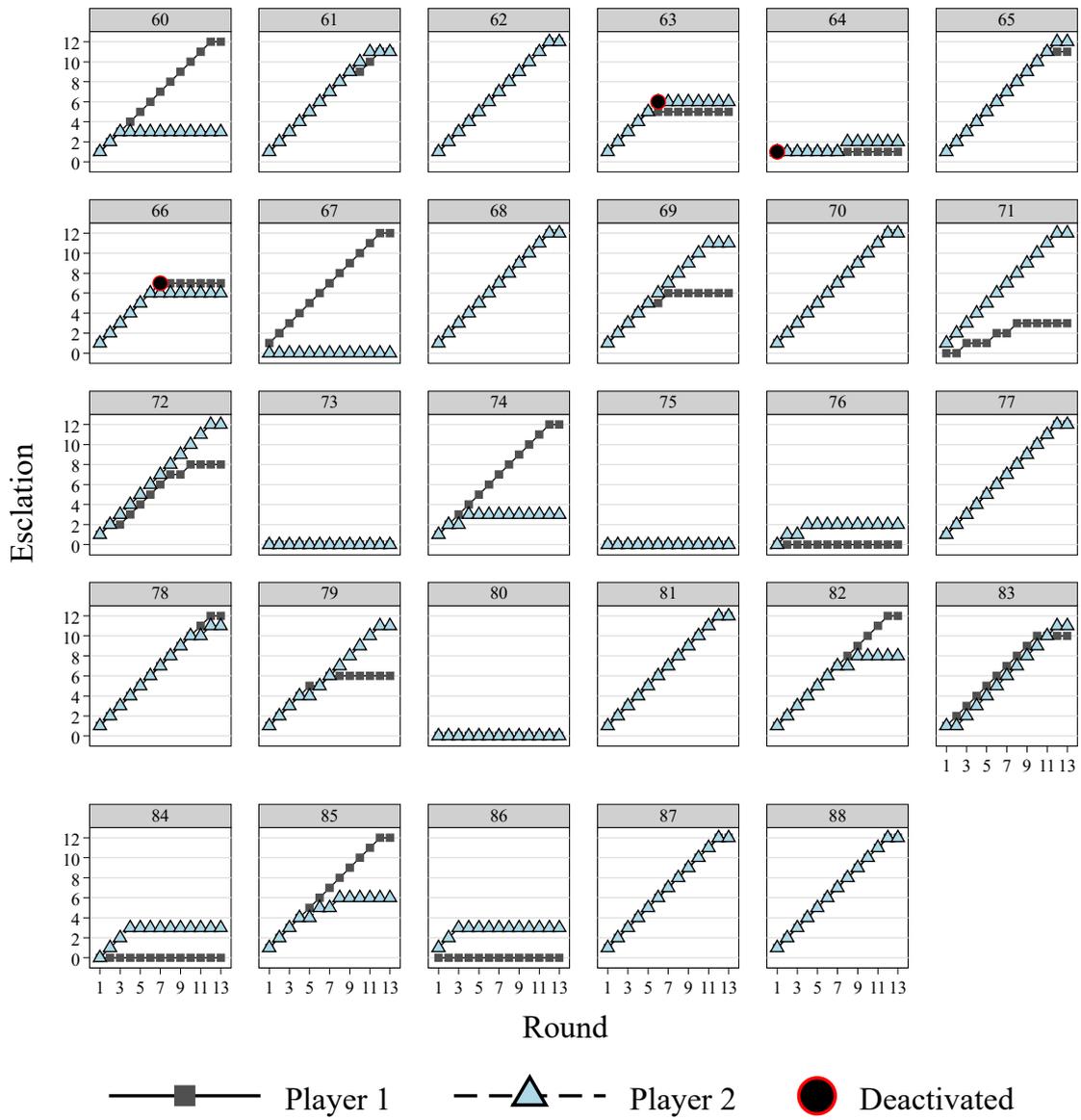


Figure B2: Group level results in the Hidden treatment

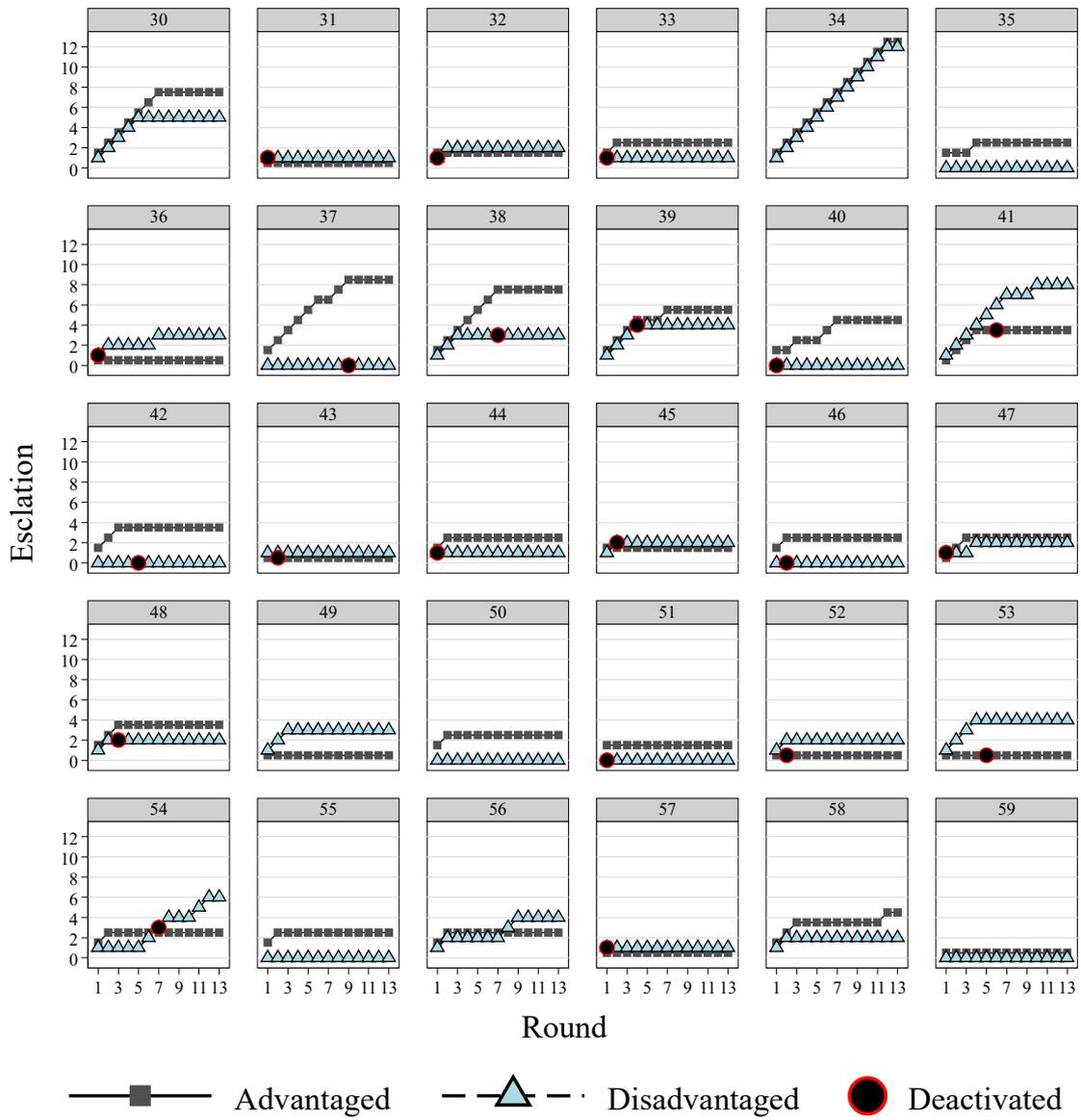


Figure B3: Group level results in the Asymmetry treatment

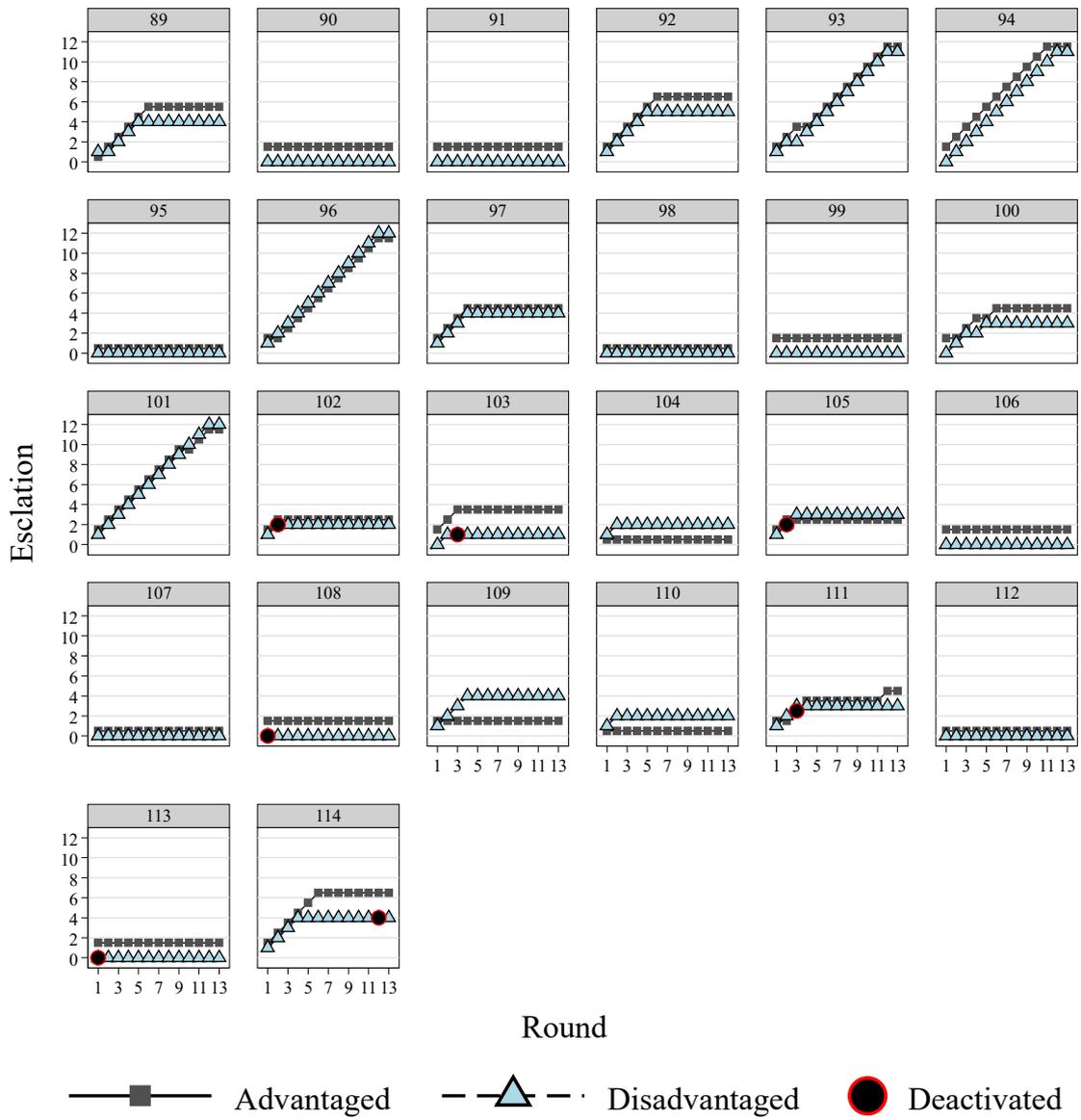


Figure B4: Group level results in the Trade treatment

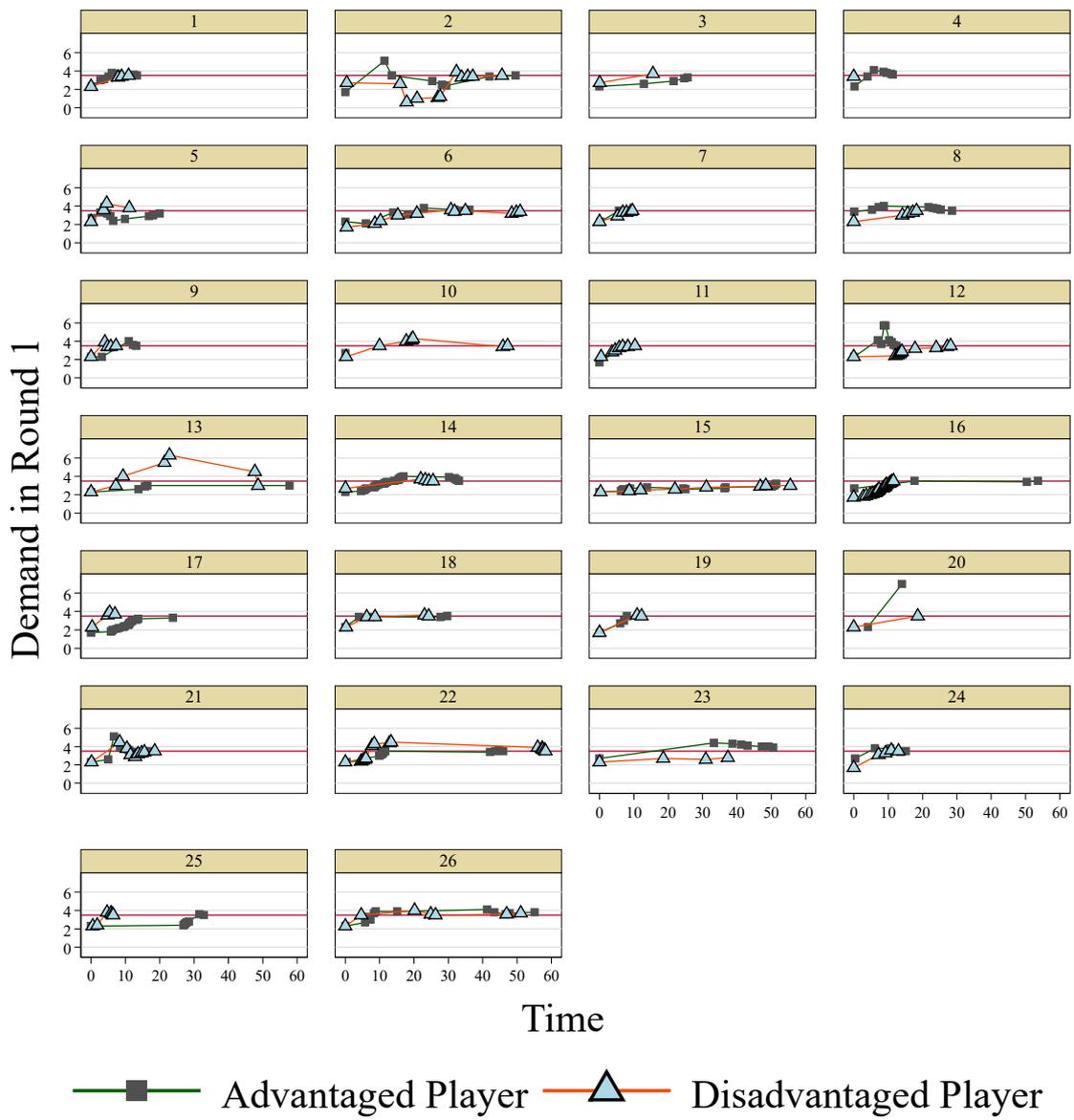


Figure B5: Group level bargaining demands over time in round 1

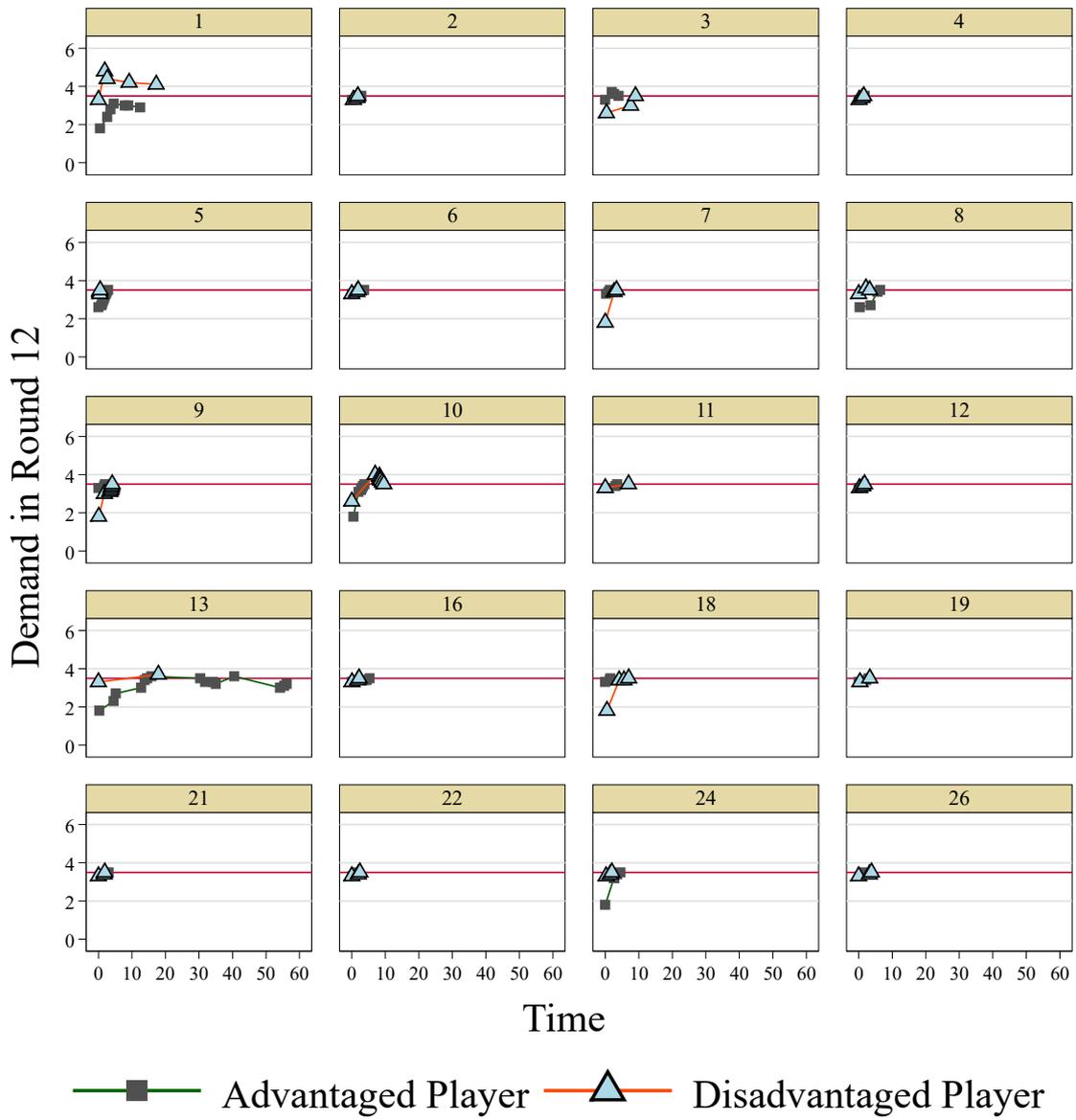


Figure B6: Group level bargaining demands over time in round 12

Online Appendix. Proofs of Unique Subgame Perfect (or Sequential) Equilibrium

A.1. Subgame perfect equilibrium in the Base, Asymmetry and Trade treatments

Preliminaries. All games have two players denoted as player $i, j \in \{1, 2\}$. All games have 13 rounds. In the first round, both players i simultaneously choose actions a_i^1 from choice sets $A_i(h^0) = \{R\&D; R\&ND; NR\&D; NR\&ND\}$ ($R = \text{buy a rocket}; NR = \text{buy no rocket}; D = \text{deactivate}; ND = \text{don't deactivate}$). (For clarity, we will ignore the lottery decisions, and the bargaining decisions in the Trade treatment, since they do not affect the equilibrium analysis as long as lottery and bargaining payoffs are treated as some exogenous endowment). At the end of the first round, both players observe the action profile $a^1 \equiv (a_1^1, a_2^1)$. At the beginning of the second round, players know history h^1 , which can be identified with a^1 . In our games, the actions player i has available in round 2, that is, the availability of the deactivation decision, depend on what has happened previously. Let $A_i(h^1)$ denote the set of possible actions when the history is h^1 . In particular, if deactivation happened in the first round $A_i(h^1)$ shrinks to $\{R; NR\}$, and otherwise $A_i(h^1)$ remains the same as $A_i(h^0)$. Iteratively, we define h^k , the history at the end of round k , to be the sequence of actions in the previous rounds, $h^k = (a^1, a^2, \dots, a^k)$, and we let $A_i(h^k)$ denote player i 's feasible actions in round $k + 1$, for $k = 0, \dots, 12$. In particular, for $k = 12$, in the last round the only possible action available is the deactivation decision, that is, $A_i(h^{12}) = \{R; NR\}$ if deactivation has not occurred in previous rounds and otherwise $A_i(h^{12}) = \emptyset$.

A pure strategy for player i is a contingent plan of how to play in each round k for possible history h^{k-1} . Let H^{k-1} be the set of all round- $(k-1)$ histories, and let $A_i(H^{k-1}) \equiv \bigcup_{h^{k-1} \in H^{k-1}} A_i(h^{k-1})$. A pure strategy for player i is a sequence of maps $\{s_i^k\}_{k=1}^{13}$, where each s_i^k maps H^{k-1} to the set of player i 's feasible actions $A_i(H^{k-1})$, $s_i^k(h^{k-1}) \in A_i(h^{k-1})$ for all h^{k-1} . Since the terminal histories represent an entire sequence of play, we can represent each player i 's payoff as a function $u_i: H^{13} \rightarrow \mathbb{R}$. With some abuse of notations, the payoff function at the terminal node is written as follows:

$$u_i = \begin{cases} \sum_{k=r+1}^{12} (0.1 * \text{lottery payoff}^k - 1_R^k) - 1.5 * \sum_{k=1}^{13} 1_D^k & \text{if deactivated in round } r \\ \sum_{k=1}^{12} (\text{lottery payoff}^k - 1_R^k) - 1.5 * \sum_{k=1}^{13} 1_D^k & \text{if never deactivated} \end{cases},$$

where 1_R^k is an indicator function of whether the player buys a rocket in round k ; 1_D^k is an indicator function of whether the player presses the deactivate button in round k . Note that according to the rules of our games, a player will be deprived of this deactivation action as long as one of the two players has pressed the button in an earlier round, no matter which player has been deactivated. In that case, we also let $1_D^k = 0$ starting from the round when the deactivation decision becomes unavailable. Finally, *lottery payoff* ^{k} is a positive amount of money which may vary from round to round and always larger than \$1. (In the Trade treatment, on top of lottery payoffs, players may earn extra money from the bargaining. For simplicity, we only refer to *lottery payoff* ^{k} .)

A subgame at the beginning of round k is defined as $G(h^{k-1})$. To define the payoff function in this subgame, let the final history be $h^{13} = (h^{k-1}, a^k, a^{k+1}, \dots, a^{13})$ and, therefore, the payoffs will be $u_i(h^{13})$. Strategies in $G(h^{k-1})$ are $\{s_i^k | h^{k-1}\}_{k=1}^{13}, s_i^r(h^{r-1}) \in A_i(h^{r-1})$ for all h^{r-1} consistent with h^{k-1} . Any strategy profile s of the whole game induces a strategy profile $s|h^{k-1}$ on any $G(h^{k-1})$ with the restriction of player i 's strategy to be $s_i|h^{k-1}$.

Definition 1: A strategy profile s is a subgame perfect equilibrium if, for every h^k , $s_i|h^k$ is a Nash equilibrium of the subgame $G(h^k)$.

Proposition 1: There is a unique subgame perfect equilibrium in which neither player buys any rocket nor deactivate in any round in the Base, Asymmetry and Trade treatments.

Proof: With Definition 1, we can apply backward induction reasoning to find all subgame perfect equilibria. Note that players must pay \$1.5 to press the “deactivate” button, even if they deactivate themselves and their previous earnings from lotteries minus any rocket investment are set to 0. And the outcome of deactivation does not depend on who pressed the button but on the relative number of rockets. Thus, deactivation is a strictly dominated action in the last round. Let's replace the last-round strategies by the dominant strategies, $a^{13} = (a_1^{13}, a_2^{13}) = (ND, ND)$. (Off this path where deactivation happened in a previous round, there is no action to make.) In the penultimate round, each player's dominant strategy is again ND no matter which side has more rockets. Given a^{13} , buying a rocket is also a strictly dominated action since trying to gain an advantage in rocket number is meaningless for the last round. Thus, the penultimate-round strategies can be replaced by $a^{12} = (a_1^{12}, a_2^{12}) = (NR\&ND, NR\&ND)$. (Off this path where deactivation happened in a previous round, the Nash equilibrium for this subgame is (NR, NR) .) Then consider the 11th round where a player may have incentive to deactivate in order to save on future expenses on rockets. To see that this incentive does not exist, note that a subgame at this round must be at one of the three situations: a player has more, equal or fewer rockets than the other player. However, in none of these situations does any player have incentive to deactivate as a preemption given that the penultimate-round optimal strategy profile specifies no deactivation. We consider one earlier round and apply the same reasoning, and so on, we find $a^1 = a^2 = \dots = a^{12}$. Thus, the strategy profile based on the sequence of dominant strategies $\{a^k\}_{k=1}^{13}$ in each subgame constitutes a subgame perfect equilibrium. Since this is the only subgame perfect equilibrium found by backward induction, uniqueness follows. The above reasoning for finding the unique subgame perfect equilibrium holds equally well for the Base, Asymmetry and Trade treatments. Q.E.D.

A.2. Sequential equilibrium in the Hidden treatment

Defining sequential equilibrium requires a few more preliminary steps. A system of beliefs μ specifies beliefs at each information set h : $\mu(x)$ denotes the probability a player assigns to node x conditional on reaching information set $h(x)$. Specifically, since in the Hidden treatment the stock level of rocket is unknown to each other, a player's information set is different from histories defined in the complete information treatments. Let h_i^k denote player i 's information set at round k . $h_i^k \equiv (a_i^k, y^k)$ where $y^k \equiv (y_1, y_2, \dots, y_k)$ and $\forall t < k, y_t \in$

$\{D, ND\}$. $y_t = D$ means that deactivation occurs in or before round t , and otherwise $y_t = ND$. $y_k = D$ iff $\exists t < k, \exists i \in \{1, 2\}$ such that $a_i^t \in \{R\&D, NR\&D\}$.

Let $u_{i(h)}(s|h, \mu)$ be the expected payoff of player i at information set h with the player's beliefs given by $\mu(h)$ and strategies s . An assessment (s, μ) is *sequentially rational* if each player believes that the other player will adhere to the equilibrium strategy profile s at every information set whether reached or not reached in equilibrium. Formally, for any information set h and alternative strategy $s'_{i(h)}$,

$$u_{i(h)}(s|h, \mu) \geq u_{i(h)}(s'_{i(h)}, s_{-i(h)}|h, \mu).$$

To discipline beliefs at an information set off the equilibrium path, an assessment is consistent such that it is a limit point of some totally perturbed assessment (totally mixed strategies and associated beliefs pinned down by Bayes' rule). Formally, let Σ denote the set of all completely mixed strategies with profile σ such that $\sigma_i(a_i|h) > 0$ for all h and $a_i \in A(h)$. Let Ψ denote the set of all assessments (σ, μ) such that $\sigma \in \Sigma$ and μ is derived from σ by Bayes' rule. Players' beliefs are as if there were a small probability of a "tremble" at each information set with these trembles statistically independent of each other. An assessment (s, μ) is *consistent* if

$$(s, \mu) = \lim_{n \rightarrow \infty} (\sigma^n, \mu^n) \text{ for some sequence } (\sigma^n, \mu^n) \text{ in } \Psi.$$

Definition 2: A sequential equilibrium is an assessment (s, μ) that satisfies sequential rationality and consistency.

Proposition 2: There is a unique sequential equilibrium (in the sense of strategies but not necessarily the belief system) in which neither player buys any rocket nor deactivate in any round in the Hidden treatments.

Proof: First, we show that an assessment, in which no player buys any rockets or deactivates in any round and the players' belief system assigns probability 1 to this strategy in each round, is a sequential equilibrium. This assessment is sequentially rational since the deviation of a player to pressing the deactivate button in any round k cannot be profitable, given his belief that the other player will not deactivate in any round. As a result, any deviation to buying a rocket in some round also cannot be profitable. Consistency is also satisfied since the players' belief system μ simply corresponds to the pure strategy profile s in this essentially simultaneous move game within each round. This follows by using a sequence of mixed strategies σ^n in which each player plays the strategy of $NR\&ND$ in every round (and ND in the last round and NR if deactivation happened in a previous round) with probability $\frac{n-1}{n}$ and plays all other strategies with combined probability $\frac{1}{n}$, for each $n \in \mathbb{N}$, and a corresponding sequence of μ^n by Bayes rule.

Next, to show the uniqueness (in the sense of strategies but not necessarily the belief system), consider an assessment in which player i 's strategy involves pressing the deactivate button in

round k and player j 's strategy involves pressing the button in round l , $l > k$. Player j would be notified of the consequence when player i pressed the deactivate button in round k and player j would actually not have the deactivation decision to make. Sequential rationality, nevertheless, requires a player's action to be optimal even at an information set not reached in equilibrium. Thus, player j 's strategy which involves pressing the deactivate button in round l is not optimal (nor feasible): player j could have been better off by pressing the deactivate button before round k . Second, an assessment in which both players deactivate in the same round is not sequentially rational. A player could have saved the deactivation cost since his action would not affect the consequence of the other player's pressing the deactivate button. Third, any assessment in which only one player plans to deactivate in some round cannot be sequentially rational since this player faces no risk of being deactivated and thus has no incentive to deactivate the other player. Finally, any assessment in which no player presses the deactivate button but some player buys rocket(s) cannot be sequentially rational because rockets are essentially useless if no player has a plan to press the deactivate button. That leaves only the case where no player will deactivate or buy any rocket and hence the uniqueness holds. Q.E.D.