

# The Paradox of Misaligned Profiling: Theory and Experimental Evidence

Charles A. Holt, University of Virginia,  
Andrew Kydd, University of Wisconsin  
Laura Razzolini, Virginia Commonwealth University  
Roman Sheremeta, Chapman University

June 5, 2013

## Abstract

This paper implements an experimental test of a game-theoretic model of equilibrium profiling. Attackers choose a demographic “type” from which to recruit, and defenders choose which demographic types to search. Some types are more reliable than others in the sense of having a higher probability of carrying out a successful attack if they get past the security checkpoint. In a Nash equilibrium, defenders tend to profile by searching the more reliable attacker types more frequently, whereas the attackers tend to send less reliable types. Data from laboratory experiments with financially motivated human subjects are consistent with the qualitative patterns predicted by theory. However, we also find several interesting behavioral deviations from the theory.

*JEL Classifications:* C72, C91, J16

*Keywords:* terrorism, profiling, game theory, laboratory experiment

---

Corresponding author: Laura Razzolini, E-mail: [lrazzolini@vcu.edu](mailto:lrazzolini@vcu.edu)

We thank Rachel Croson, Catherine Eckel, seminar participants at the University of Virginia, the University of Texas at Dallas, SUNY-Buffalo and participants at the North American Economic Science Association Conference in Tucson. We also wish to thank Michael Patashnik for research assistance. Research support provided by the University of Texas at Dallas, the Virginia Commonwealth Presidential Research Incentive Program and NSF grant NSF/NSCC-0904695 to Razzolini is gratefully acknowledged. Holt’s work on the project was funded in part by NSF/NSCC grants 0904795 and 0904798 and the UCS CREATE National Center for Risk and Economic Analysis of Terrorism Events.

## 1. Introduction

An important problem facing security personnel is to identify terrorists within large groups of mostly innocent people. This problem arises at checkpoints of all kinds, such as roadblocks, permanent checkpoints between different regions or countries and airport security counters. In such settings, large numbers of people pass through a screening process designed to detect and detain terrorists. Typically, the volume of traffic in comparison with the number of potential terrorists is so large that it is not economically or politically sensible to screen everyone with the intensity required to detect a terrorist. The security personnel, therefore, face a difficult task, whom to take out of the line and subject to greater scrutiny, when there are many innocent people and few terrorists. For instance, this issue was faced at U.S. military checkpoints that attempted to secure the Green Zone in Baghdad during the U.S. occupation. This issue is also faced daily at busy airports, especially during periods of high security alert.

The problem is compounded by the fact that a terrorist group may be strategic in the sense of being able to respond quickly to any targeted screening program. If the government screens one category that is closely associated with the terrorist group (for instance, young men from a certain province in the Iraq case), then the terrorist group faces a strong incentive to begin recruiting and sending people outside that category, for example, women. The government would then rationally respond by searching women, at least to some extent. The equilibrium outcome of this interaction is unclear. Some argue for a completely random search process, but that would give the terrorists an incentive to send only their core supporters, since they would be no more likely to be searched than anyone else, and these core supporters *would* presumably be more likely to carry out a successful attack if not searched. There is a need for careful analysis to determine what kind of search strategy is rational, implementable and efficient in an environment with threats that evolve in response to securities measures.

Profiling occurs when a certain characteristic or signal, such as race or ethnicity, is used to decide who to subject to a more intrusive investigation. Racial profiling, for instance, is based on the belief that certain crimes are committed disproportionately by the members of a particular race. The scholarly debate in the social science literature is mostly focused on whether profiling occurs and whether it is effective, rather than on normative and constitutional issues. Profiling is not *per se* illegal in most countries, and legal scholars have discussed the pros and cons of various types of profiling and the circumstances under which it might be justified (Barnes and

Gross, 2002; Ellmann, 2003).<sup>1</sup> Economists have studied racial profiling issues from a different perspective: the question is whether profiling is a rational use of limited enforcement assets. A recent literature pioneered by Knowles, Persico and Todd (2001) has characterized the Nash equilibrium of a simultaneous move game between the police and a specific group of the population, such as motorists/drivers, who may commit a crime. The police objective is to minimize crime when deciding which vehicles to search, while motorists choose whether to carry contraband or not. The equilibrium involves unequal investigation rates across different demographic groups, even if police officers are unbiased, as long as the members of one group incur higher costs of carrying contraband than those of another group. In this sense, profiling can result from a type of rational experienced-based or “statistical” discrimination. Antonovics and Knights (2009) point out that if statistical discrimination alone is used to explain differences in the rates at which vehicles of drivers of different races are searched, then these search decisions should be independent of police officers’ *own* race. They test this prediction using data from the Boston Police Department and find that officers are more likely to conduct a search if the race of the officer differs from the race of the driver.

The debate on profiling has significantly changed after the terrorist attacks of September 11th, 2001. Screening procedures have included different versions of the Computer Assisted Passenger Prescreening System (CAPPS) and the Secure Flight Passenger Screening Program, a computerized tool to select passengers for screening, and more recently the full body image scanning. Passengers with elevated ratings according to these mechanisms are selected for additional searches and for baggage inspection, while some other passengers are still searched at random. If previously the discussion was about whether demographic profiling was happening, after 9/11 researchers and politicians have focused on the conditions under which such profiling is acceptable, either constitutionally or as a policy matter. All screening and profiling mechanisms have encountered criticisms and often legal actions. Full body imaging is opposed because it produces a nude image of the passengers. Computerized searching mechanisms are

---

<sup>1</sup> In *Protecting Liberty in an Age of Terror*, Heymann and Kayyem (2005) discuss the tension between avoiding past abuses of profiling and the need to confront high-stakes threats. They suggest more reliance on nationality-based profiling, as opposed to pure racial profiling. Barak-Erez (2007) offers an effects-based consideration: if profiling really is necessary, then it should be used more often in situations in which it is less likely to have a long-lasting effect on the lives of those being profiled. Harcourt (2007) provides a thoughtful discussion of the unintended consequences of profiling, such as the tendency for “false positive” searches to induce more terrorism.

often accused of inducing racial or religious profiling and discrimination. In fact, there is a debate among experts whether profiling strategies are more effective than pure random searches.

Several years ago, on November 22, 2010, National Public Radio held an Oxford-style debate at New York University with the two teams arguing the motion “Should U.S. Airports Use Racial and Religious Profiling?” Advocates of the motion supported the use of profiling specifically concentrating on young fundamentalist Muslim males from the Middle East, as the majority of recent terrorist attacks have been associated with this type of individual. Opponents argued that profiling just invites terrorist groups to recruit agents who do not fit the profile. Bin Laden himself, in his hand-written journals, “exhorted followers to explore ways to recruit non-Muslims (...) – particularly African Americans and Latinos” (the *New York Times*, 5/12/2011).

This paper contributes to this debate by providing theoretical analysis and experimental validation to guide policy makers to improve the effectiveness of targeted and/or profiled screening. Our research investigates the conditions under which profiling is a rational and efficient counterterrorism policy. Although our work is related to the previous economic literature, we take a somewhat different approach by assuming that the terrorist group rationally chooses individuals with certain characteristics to carry out its attacks, rather than viewing terrorism as the result of decentralized individual choices.<sup>2</sup> In the model, the terrorist group (attacker) decides which demographic “type” to send through a security checkpoint. The security officials (defender) decide which type to subject to an extensive search. Some types (for instance, young males with military and ideological training) are more “reliable” than others (women and children) in the sense of having a higher probability of mounting a successful attack once they pass undetected through a security checkpoint. If the attackers were not selective, when sending a mix of types, then defenders would use limited resources to search the most reliable types who would cause the greatest damage if they passed security. Consequently, attackers would respond by sending less reliable types more often. In turn, defenders would respond by defending less reliable types more often. In equilibrium, attackers and defenders should not have any additional incentive to change their strategies. We show that, in a mixed

---

<sup>2</sup> Our profiling game is related to the hide-and-peek game introduced by Rubinstein, Tversky and Heller (1996) and further examined by Crawford and Iriberry (2007). Our game is also related to the multi-player parasite game introduced by Avrahami, Güth and Kareev (2005). Similarly to the hide-and-peek and the parasite game, the defender desires to match by searching a category selected by the attacker, who in turn desires to select a category that does not match. However, a significant departure of our profiling game is that even when the selected categories by defender and attacker match, the outcome of the game is still probabilistic and it depends on the reliability of a selected type. Therefore, our game can be considered as a generalization of the hide-and-peek and the parasite game.

strategy Nash equilibrium, there is a tendency to use low-reliability attack strategies and high-reliability defense strategies. Thus, attack and defense strategic patterns are seemingly “misaligned,” even though both players are rational and there are no surprises in terms of observed behavior.

We use an experiment to assess the extent to which individual decisions are consistent with theoretical predictions of misaligned profiling. The experiments are motivated, in part, by the somewhat counter-intuitive nature of equilibrium patterns of the randomized strategies. In particular, the theory produces a paradox of misaligned profiling: in equilibrium the high reliability categories are searched more intensively, even though they are used less intensively by the terrorist organization. Field experiments with “professional” terrorists and security officials to test these predictions would be expensive and controversial, if possible at all, and the results would surely be confidential. Instead, we rely on laboratory experiments, which provide the ability to replicate and control the environment, even though the laboratory environment is admittedly highly simplified. The results of the experiment reveal behavioral patterns that are consistent with predicted patterns. However, we also find several interesting behavioral deviations from the theory.

## **2. Theoretical Model and Predictions**

To address the problem of profiling consider a simple two-player model of screening at security checkpoints. The first player represents a government agency that is attempting to discover terrorists, e.g., military officers at a checkpoint, Transportation Security Agency officials at an airport security counter, or other agencies dealing with homeland defense, such as the Coast Guard, the Border Patrol, the Customs Service, or the Immigration and Custom Enforcement department of the Department of Homeland Security. Their objective is to identify any terrorists attempting to penetrate their checkpoint and, thereby, block an attack. The second player represents a terrorist group, which is assumed to be centrally directed, strategically rational, and motivated by a desire to penetrate the defenses and commit an attack. Experts agree that there is usually a strategy behind terrorists’ actions. Whatever form it takes, terrorism is typically not random, or blind; it is a deliberate use of violence against civilians for political or religious reasons. Therefore, following the spirit of most game theoretic literature on terrorism,

we model terrorists as rational actors (for an excellent survey, see Sandler and Arce, 2007).<sup>3</sup> In what follows, we will refer interchangeably to the terrorist as the “attacker” and to the government agency as the “defender.”

The main motivation for our research can be illustrated with a simple two-category example derived from Kydd (2011). Specifically, assume that the attacker can send two types of individuals (type 1 and type 2) passing through the checkpoint, each characterized by a “reliability” parameter  $r_i$  – the probability of mounting a successful attack after making through security. Also, without loss of generality, assume that type 1 is more likely to succeed if undetected than type 2, i.e.,  $r_1 > r_2$ . Both types could be identified and determined by different criteria, such as age, gender, and country or region of origin, religion, or any other observable personal characteristic. For instance, the *New York Times* on June 27, 2011 reported that in a remote area in central Afghanistan “insurgents tricked an 8-year-old girl ... into carrying a bomb wrapped in cloth that they detonated remotely when she was close to the police vehicle.” This is an example of an attack using a person less likely to be searched from a less reliable type or category.<sup>4</sup>

A successful attack by a person from any category will result in a gain of  $G$  for the terrorist and a loss of  $L$  for the government security agency. We assume that the defender selects one category to search, and that the search is fully effective in detecting the attacker.<sup>5</sup> Hence, the attacker of type  $i$  who is searched would fail, and one who is not searched would carry out a successful attack with probability of success  $r_i$ . Let  $d_i$  denote the probability of defending against type  $i$ , and let  $a_i$  denote the probability of attacking with type  $i$ , with  $0 < d_i, a_i \leq 1$ . Then, a person from category  $i$  would succeed only if not searched and able to carry off a successful attack, which occurs with probability  $(1 - d_i)r_i$ . The attacker, therefore, faces a tradeoff between sending highly reliable type 1 who is more likely to be searched and sending less reliable type 2 who is

---

<sup>3</sup> Our models of strategic terrorism can be modified in a straightforward manner to include the possibility of exogenous, decentralized sources of terrorism of the emotional or “home-grown” variety.

<sup>4</sup> To be clear, there is heterogeneity within a category, and the “reliability” of a category represents the average effectiveness of the best people in that category who can be successfully identified, recruited and trained, i.e., the “tail” of the distribution in terms of physical and emotional fitness, willingness to risk extreme injury or death, unwavering loyalty to the cause, and the ability to mask emotions and improvise to defeat unanticipated challenges.

<sup>5</sup> This assumption, used in Kydd (2011) and Basuchoudhary and Razzolini (2006), can be generalized to derive comparative statics predictions for an improvement in search technology (see Holt, Kydd, Razzolini, and Sheremeta, 2011).

less likely to be searched. The defender's probability of a loss from defending against type  $i$  is  $(1 - a_i)$  times the average reliability of the other type.

Before deriving the equilibrium, it is useful to provide some intuition. The equilibrium involves randomization, since a deterministic attack via one type would lead to a sure defense there, and a deterministic defense against one type would lead to a sure attack via another. In equilibrium with randomization, the expected payoffs for all decisions used must be equal, otherwise the player would prefer decisions with higher expected payoffs. The main result reveals a *paradox of misaligned profiling*: in equilibrium the high reliability types are searched *more* intensively, even though they are used *less* intensively by the terrorist organization. This is a paradox in the sense that the equilibrium pattern makes the defense strategy appear to be misguided, and hence, ineffective, which is not the case.

The paradox can be illustrated for the special case of a zero-sum game in which a successful attack results in payoffs of 1 for the attacker and  $-1$  for the defender. To be willing to randomize, the attacker's expected payoff for sending either type must be equal, i.e., the product of the probabilities of not being searched and of succeeding are the same for both categories:  $(1 - d_1)r_1 = (1 - d_2)r_2$ . Since  $1 - d_2 = d_1$ , we obtain a single equation,  $(1 - d_1) r_1 = d_1 r_2$ , which can be solved for defense probability against type 1:

$$\text{Defense probability against type 1: } d_1 = \frac{r_1}{r_1 + r_2}.$$

Since type 1 is more reliable,  $r_1 > r_2$ , it follows that  $d_1 > d_2$ , or the defender searches the more reliable type 1 more often, which is intuitive.<sup>6</sup> Conversely, for any given attack probabilities  $a_1$  and  $a_2$ , a defense against type 1 will result in a loss with probability  $a_2 r_2$  whereas a defense against type 2 will result in a loss with probability  $a_1 r_1$ . Since  $a_2 = 1 - a_1$ , the equality of expected defender payoffs results in an equation,  $(1 - a_1) r_2 = a_1 r_1$ , which can be solved for the attack probability via type 1:

$$\text{Attack probability via type 1: } a_1 = \frac{r_2}{r_1 + r_2}.$$

---

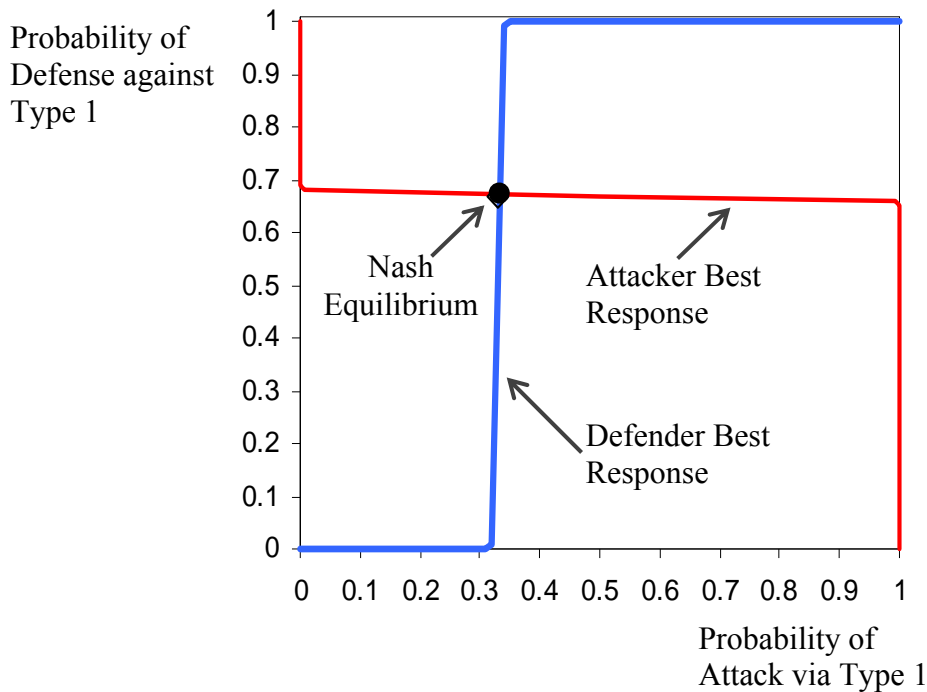
<sup>6</sup> The use of expected payoffs implicitly assumes risk neutrality, but a reformulation in terms of expected utility yields the same monotonicity, since the attacker's expected utility would still be decreasing in  $d_i$  and increasing in  $r_i$ , so equating expected payoffs would still imply that high reliability types are searched more often.

Hence the counter-intuitive result that  $a_1 < a_2$ ; that is, in equilibrium, the attacker uses the *more* reliable type 1 *less* often, since  $r_2 < r_1$ .<sup>7</sup>

### 3. Experimental Design and Procedures

Our objective is to evaluate the extent to which observed behavior is consistent with theoretical predictions. To this end, we design two treatments with  $(r_1 = 0.67, r_2 = 0.33)$  and  $(r_1 = 0.80, r_2 = 0.20)$ . In both treatments  $r_1 + r_2 = 1$ . In the 67/33 treatment, type 1 is twice as reliable as type 2. Theoretical prediction is that the defender searches type 1 with probability  $d_1 = 2/3$  and the attacker uses type 1 with probability  $a_1 = 1/3$ . Conversely, the probability of defense against type 2 is  $d_2 = 1/3$  and the probability of attack via type 2 is  $a_2 = 2/3$ . The best-response functions for this game are shown in Figure 1.<sup>8</sup>

**Figure 1: Best Responses and Equilibrium**



<sup>7</sup> Given the theoretical predictions about the defense and attack probabilities, the expected payoff of the attacker is  $r_1 r_2 / (r_1 + r_2)$  and the expected payoff of the defender is  $1 - r_1 r_2 / (r_1 + r_2)$ .

<sup>8</sup> The best-response functions should actually have sharp corners, since perfectly rational players are assumed to respond sharply to small differences in expected payoffs. The lines in the figure are plotted with a slight amount of curvature to make it easier to visually separate the two best-response lines.



For example, if the probability of an attack via type 1 is low (left side), the probability of a defense against type 1 is low (bottom left part of the figure). Conversely, if the probability of defense against type 1 is high (top), then the attacker will use type 1 with probability 0 (upper left side). The intersection of the best response lines determines the equilibrium, with a 2/3 probability of a defense against type 1, and a 1/3 probability of an attack via type 1.

In the 80/20 treatment, we increase the reliability of type 1 so that type 1 is four times more reliable than type 2, i.e.,  $r_1 = 0.80$  and  $r_2 = 0.20$ . The prediction of the theory is that since type 1 is more reliable, the probability of defense against type 1 should *increase* to  $d_1 = 4/5$ , while the probability of attack via type 1 should *decrease* to  $a_1 = 1/5$ . Correspondingly, the probability of defense against type 2 should decrease to  $d_2 = 1/5$ , while the probability of attack via type 2 should increase to  $a_2 = 4/5$ .

Subjects for the experiment were recruited from student populations at the University of Virginia, with the promise that they will “participate in a research experiment” and will receive a fixed payment of \$6 plus additional cash earnings, which will depend on their own and others’ decisions. When subjects arrived in the lab, they were seated in visually separated cubicles with networked computers. The software kept track of total earnings, and subjects were paid in cash at the end of each session, after they signed receipt forms. A total of 144 subjects participated in the experiment with 72 subjects (36 pairs) participating in one treatment and 72 subjects (36 pairs) participating in the other treatment. Instructions for the experiment are included in the Appendix. The screen displays listed the most reliable category on the left side for half of the subject pairs, and on the right side for the other half. The experiment was run in a series of sessions consisting of 12 to 18 subjects each, with fixed pair matching for 50 rounds. Subjects in each pair were given the role of attacker or defender and stayed in that role in all rounds of the experiment.

In each round, the attacker chose a category corresponding to a type of terrorist agent (type 1 or type 2), and the defender chose a profiling strategy, or type of person to search. If the selected categories matched then the attack failed. If the selected categories did not match, then the attack success probability was determined by the reliability of the attacker’s category choice. A successful attack resulted in a fixed payoff of 1 dollar to the attacker and a loss of 1 dollar for the defender. The payoffs were added to private incomes of \$1 for the defender and \$0.60 for the attacker in each round. These outside incomes were selected to equalize final payoffs and were

private information (defenders did not know attacker incomes, and vice versa). On average subjects earned \$26 and the experiment lasted for about 30 minutes.<sup>9</sup>

#### 4. Results

Table 1 reports for the two treatments the predicted and average observed probabilities of defense and attack for type 1, respectively from the first round, the second half and all 50 rounds of play.

**Table 1: Experimental Data and Predictions**

	Treatment			
	67/33		80/20	
	$d_1$	$a_1$	$d_1$	$a_1$
Predicted probabilities	0.67	0.33	0.80	0.20
1 <sup>st</sup> round average data	0.81	0.36	0.92	0.19
2 <sup>nd</sup> half average data	0.78	0.30	0.88	0.20
All rounds average data	0.79	0.31	0.88	0.22

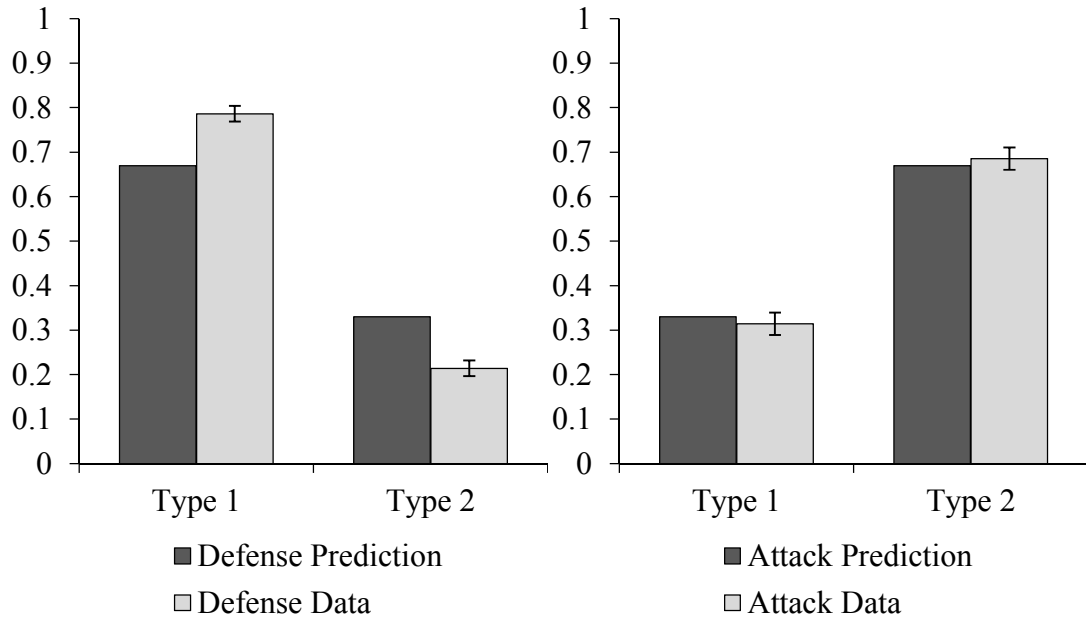
We begin by analyzing the data from the 67/33 treatment. Figure 2 displays the average defense and attack probabilities, while Figure 3 displays the average attack and defense probabilities for each of the 36 fixed pairs. As predicted by the theory, there is strong “misaligned profiling.” Specifically, the defenders search more reliable type 1 with higher probability than less reliable type 2 (0.79 versus 0.21; Wilcoxon signed-rank test, p-value < 0.01, n = 36).<sup>10</sup> Conversely, the attackers employ more reliable type 1 with *lower* probability than less reliable type 2 (0.31 versus 0.69; Wilcoxon signed-rank test, p-value < 0.01, n = 36).

Relative to theoretical point predictions, we find that defenders tend to defend against the more reliable type 1 more than predicted (0.79 versus 0.67; Wilcoxon signed-rank test, p-value < 0.01, n = 36). The behavior of attackers is not significantly different from theoretical predictions for category 1 (0.31 versus 0.33; Wilcoxon signed-rank test, p-value = 0.42, n = 36).

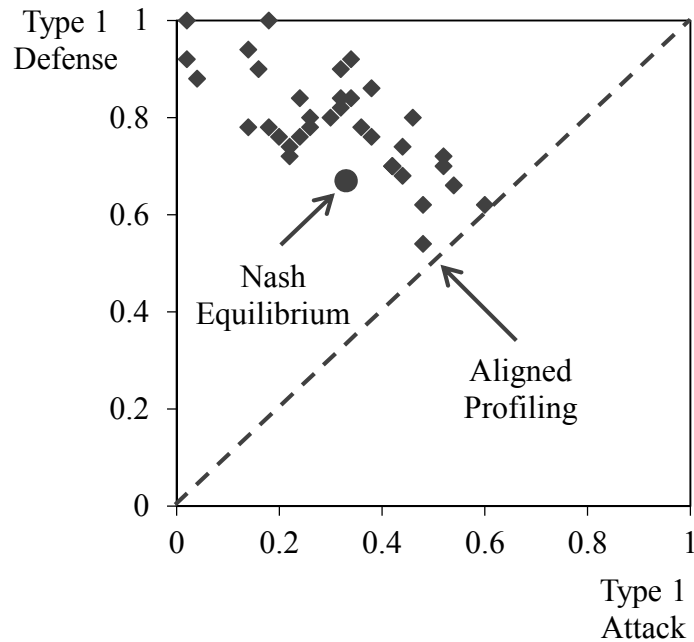
<sup>9</sup> The data for all sessions are available on the web at: [http://www.people.virginia.edu/~cah2k/profiling\\_data.htm](http://www.people.virginia.edu/~cah2k/profiling_data.htm). The 10 sessions are numbered adn24-adn28, adn30, adnp1-adnp4. For each session, the data table provides links to a graph of data averages for each round, a color-coded verbal/numerical summary of each attacker-defender interaction, and a presentation of all data in column form, which was used to create an Excel data file for each session (also included in the table).

<sup>10</sup> In conducting statistical tests, we treat the average over all 50 rounds of the experiment by the same subject as one observation. The results also hold if we analyze only the first round of the experiment and are available upon request.

**Figure 2: Average Data for All Rounds and Theoretical Predictions in the 67/33 Treatment**



**Figure 3: Attack and Defense Probabilities for 36 Fixed Pairs in the 67/33 Treatment**

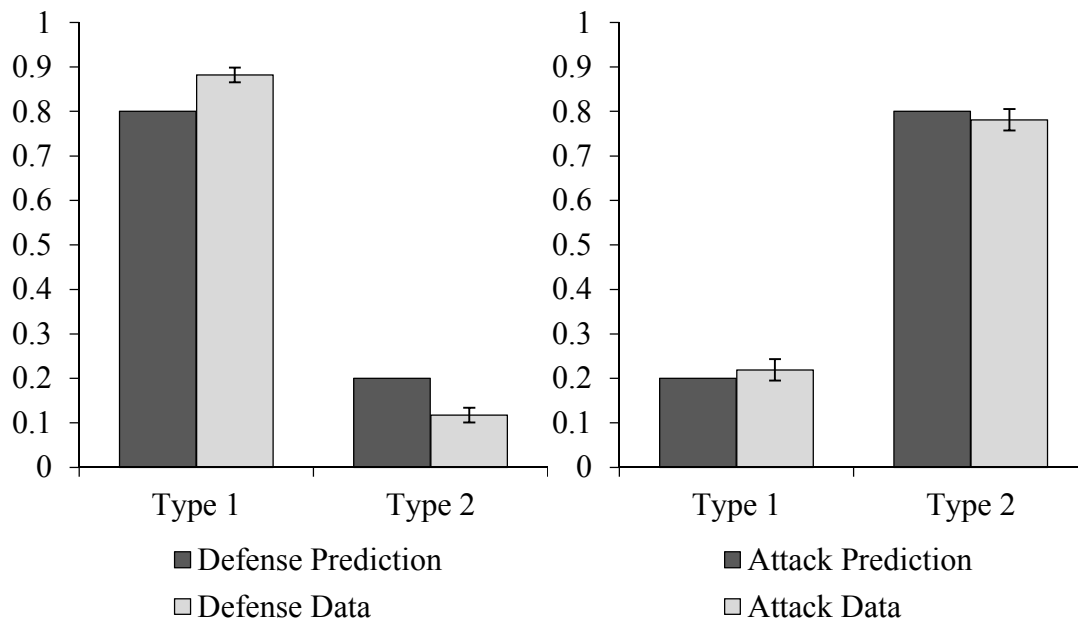


One may argue that subjects in a role of defender search type 1 more often because this option is presented to the left from type 2. In a related literature on multi-battle contests, called Colonel Blotto games, where attackers and defenders allocate resources on multiple battlefields,

it is documented that subjects often exhibit allocation bias towards left battlefields (Chowdhury, Kovenock and Sheremeta, 2013).<sup>11</sup> Nevertheless, it is unlikely that allocation bias can explain our data, since in half of the sessions type 1 option was presented to the left from type 2 and in the other half it was presented to the right. Moreover, the probability that the defender searches type 1 is virtually the same disregarding whether type 1 is located to the left or to the right of type 2 (0.79 versus 0.79; Wilcoxon rank-sum test, p-value = 0.72,  $n_1 = n_2 = 18$ ).

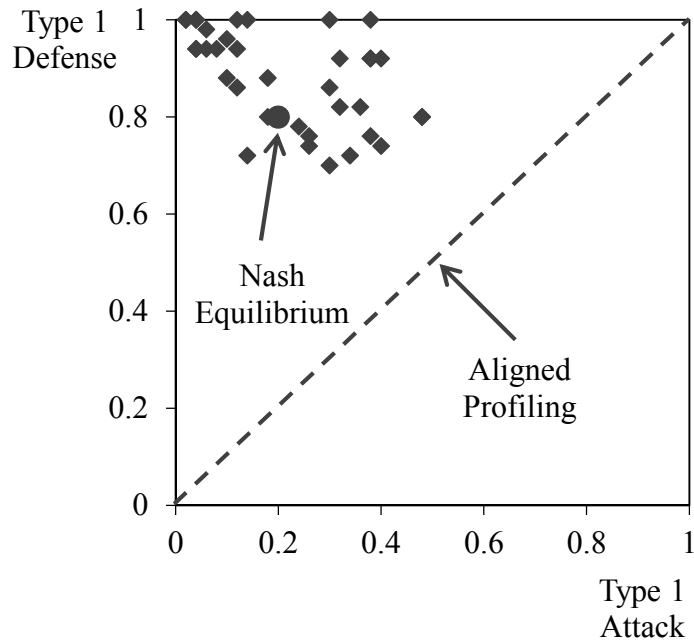
The pattern of data that we observe in the 67/33 treatment is also observed in the 80/20 treatment with  $r_1 = 0.80$  and  $r_2 = 0.20$ . Figures 4 and 5, displaying the average defense and attack probabilities for the 80/20 treatment, show even stronger “misaligned profiling.”

**Figure 4: Average Data for All Rounds and Theoretical Predictions in the 80/20 Treatment**



<sup>11</sup> For a comprehensive review of this literature see Dechenaux, Kovenock and Sheremeta (2012). The two most relevant studies that investigate the behavior of attackers and defenders are Kovenock, Roberson and Sheremeta (2010) and Deck and Sheremeta (2012). Both studies examine behavior in attacker-defender games, where the defender needs to win all targets, while the attacker needs to win only one target to secure the prize.

**Figure 5: Attack and Defense Probabilities for 36 Fixed Pairs in the 80/20 Treatment**



Specifically, the defenders search more reliable type 1 with much higher probability than less reliable type 2 (0.88 versus 0.12; Wilcoxon signed-rank test,  $p$ -value  $< 0.01$ ,  $n = 36$ ). Conversely, the attackers use more reliable type 1 with much *lower* probability than less reliable type 2 (0.22 versus 0.78; Wilcoxon signed-rank test,  $p$ -value  $< 0.01$ ,  $n = 36$ ). As before, the defenders tend to defend against the more reliable attacker type 1 more than predicted (0.88 versus 0.80; Wilcoxon signed-rank test,  $p$ -value  $< 0.01$ ,  $n = 36$ ). This behavior cannot be explained by the allocation bias, since the probability that the defender searches type 1 is very similar in sessions where type 1 is located to the left versus sessions where type 1 is located to the right of type 2 (0.86 versus 0.90; Wilcoxon rank-sum test,  $p$ -value = 0.19,  $n_1 = n_2 = 18$ ).

A comparative static prediction of the theory is that increasing reliability of type 1 should increase the defense probability against this type. This is exactly what we observe in the experiment. When the reliability of type 1 increased from  $r_1 = 0.67$  to  $r_1 = 0.80$ , the defense probability against type 1 increased from 0.79 to 0.88 (Wilcoxon rank-sum test,  $p$ -value  $< 0.01$ ,  $n_1 = n_2 = 36$ ). The prediction for the attacker is the opposite, increasing reliability of type 1 should *decrease* the probability of using type 1. Again, we observe this in the experiment. When the reliability of type 1 increased from  $r_1 = 0.67$  to  $r_1 = 0.80$ , the probability of using type 1 by the attacker decreased from 0.31 to 0.22 (Wilcoxon rank-sum test,  $p$ -value = 0.01,  $n_1 = n_2 = 36$ ).

## 5. Behavioral Explanations

The data averages are close to Nash predictions for the two treatments. However, there is a slight tendency for defenders to defend against the more reliable type more often than predicted (the attackers' behavior, on the other hand, is quite close to the Nash predictions). Although deviations from these predictions are not large in magnitude, it is notable that deviations are largely concentrated among defenders in both treatments. There are several alternatives that can explain this behavior.

First consider whether a non-equilibrium model of initial responses based on "level-k" thinking may explain the observed deviations from Nash equilibrium in our data (Stahl and Wilson, 1994, 1995). Level-k thinking would imply that level-0 attackers randomly choose between types 1 and 2, and level-0 defenders randomly search type 1 and type 2, with equal frequency. In response to level-0 behavior, level-1 attackers should use type 1 and level-1 defenders search type 1. Similarly, in response to level-1 behavior, level-2 attackers should use type 2 and level-2 defenders should search type 1. Given that the majority of subjects are documented to be of level-1 and level-2 (Crawford, Costa-Gomes and Iriberry, 2013), we should expect subjects in a role of defender to search type 1 more often than predicted, since both level-1 and level-2 behavior predicts searching type 1. In our opinion, a level-k analysis is more appropriate for explaining the first round data, where subjects do not have actual experience on which to base expectations, which, in turn, are then formed with levels of introspective thinking. The subjects in our experiments have data on which to base beliefs, at least after round 1. Even in the first round, however, the level-k analysis makes the same predictions for both treatments, since level-0 attackers and defenders would select type 1 in each treatment, and the best response to level-0 is for attackers to attack with type 2 and defenders to defend type 1, in each treatment. The analysis of higher levels also yields identical predictions for each level for the two treatments. Hence, the level-k analysis cannot explain the most salient feature of the data even in the initial round, i.e., the strong treatment effect of increasing the reliability parameter from  $r_1 = 0.67$  to  $r_1 = 0.80$ , i.e., the observed increase in defenses with type 1 and decreases in attacks with type 1. In contrast, this treatment effect of the increase in  $r_1$  is picked up by the Nash predictions.

An alternative approach would be to stay with an equilibrium model but introduce some behavioral "noise" in decision making (McKelvey and Palfrey, 1995). The effect of this noise is to inject a little more curvature into the "better response" functions. It can be shown that this

would result in an intersection of best response functions that is moved up and to the right in Figure 1, which is not in the direction of the deviation from the Nash prediction. A quantal response model would track the main qualitative feature of the data, the strong predicted treatment effect of an increase in  $r_1$ , but the deviations of data from Nash predictions (upward and to the left of the Nash prediction for both treatments in Figure 1) are not explained by smoothing off the best response functions (as would happen with a quantal response equilibrium) since the "S shaped" defender best reply function would curve to the right in the upper half of the figure, moving intersections up and to the right, not up and to the left. This intuition is confirmed by our estimates of a quantal response error parameter (available on request from the authors), which only provided a slight improvement over Nash predictions for these two treatments.<sup>12</sup>

Another possibility for why subjects in a role of defender tend to search type 1 more often than predicted is loss aversion (Kahneman and Tversky, 1979). If type 1 is not searched, there is a high probability of a successful attack if the attacker decides to employ type 1. On the other hand, if type 2 is not searched there is relatively low probability of a successful attack even if the attacker decides to employ type 2. Let  $F$  denote the cost of an attack that fails. Thus the attacker's payoff is 0 if the attack is stopped by the defender at the search point, and it is  $-F$  if the attacker makes it past the search point and fails in its implementation.<sup>13</sup> Our motivation in considering a behavioral explanation for observed behavior is based on some informal discussions with subjects after several of the sessions in which those with attacker roles expressed frustration with selecting an attack type that was not defended but that failed nevertheless. This regret, which we model as having a cost of  $F$ , is more likely to occur when the less reliable type 2 is used.

The analysis of a model with asymmetric failure cost will involve mixed strategies and equality in expected payoff. As before, the payoff for a successful attack is 1. An additional assumption necessary to avoid a boundary solution is that the possibility of success ( $r_i$  times the payoff of 1 for type  $i$ ) is not overwhelmed by the regret cost for either type:  $r_i > F(1 - r_i)$  for  $i = 1, 2$ , where the right side of the inequality is the expected regret cost from a failed attack. To be

---

<sup>12</sup> Of course, this does not mean that a quantal response analysis would not provide useful predictions in a richer model with payoff asymmetries, but the inclusion of noise in the simple model is not useful in explaining the types of behavioral asymmetries that we observe.

<sup>13</sup> It is possible to imagine scenarios in which such an asymmetry could exist, but there is no *ex ante* reason to expect that failure after the defense point would be more damaging to the attacker than a failure at the defense point. Although a failure in implementation after getting through the defenses might leave a lot of incriminating evidence (positive cost), a failure of *resolve* might leave no evidence.

willing to randomize, the attacker’s expected payoff (including regret costs) for sending either type must be equal, i.e.,  $(1 - d_1)r_1 - (1 - d_1)F(1 - r_1) = (1 - d_2)r_2 - (1 - d_1)F(1 - r_2)$ . The first term on each side of the equation is the same as before, the product of the probability of not being searched and of succeeding. The second term on each side represents the expected regret cost. Since  $1 - d_2 = d_1$ , the above equation can be solved to express  $d_1$  as a ratio of regret-adjusted attacker payoffs:

$$\text{Adjusted defense probability against type 1: } d_1 = \frac{r_1 - F(1 - r_1)}{r_1 - F(1 - r_1) + r_2 - F(1 - r_2)}.$$

It is straightforward to use the assumption that  $r_i > F(1 - r_i)$  to show that  $d_1 > r_1 / (r_1 + r_2)$ . The intuition is that the asymmetric cost of failure provides attackers with a motivation to use the more reliable type, which in turn, requires that the defender must search the more reliable type more often in equilibrium (to keep the attacker indifferent).

In Figure 1, the effect of this asymmetric attacker failure cost is to shift the horizontal crossover portion of the attacker best response upward, which increases the value of  $d_1$  at the equilibrium intersection. Notice that this analysis of failed attack cost has no effect on defenders’ expected payoffs, so the defender best response line in the figure would remain unchanged. Finally, it is worth considering a symmetric loss aversion, i.e., a failure cost of  $F$  applied equally to attack failure from any source, whether caused by encountering a defense or by a subsequent reliability issue. It is straightforward to show that this symmetric attacker loss aversion would have no effect on the equilibrium,  $d_1 = r_1 / (r_1 + r_2)$ . The intuitive reason for this invariance is that, in equilibrium, a failure from any source is equally likely for each attacker type since the more reliable type is defended against more often.

## 6. Conclusions

This paper implements an experimental test of a game-theoretic model of equilibrium profiling. Attackers choose a demographic “type” from which to recruit, and defenders choose which demographic types to search. Some types are more reliable than others in the sense of having a higher probability of carrying out a successful attack if they get past the security checkpoint. Using a controlled laboratory experiment with financially motivated human subjects, we find strong support for game-theoretic model of equilibrium profiling. Consistent with theoretical predictions, the defenders search more reliable types with *higher* probability, while



the attackers employ more reliable types with *lower* probability than less reliable types. However, we also find systematic deviations with defenders searching more reliable types more often than predicted. This type of behavior in our experiment can be partially explained by a model based on asymmetric loss aversion, and in particular on attacker aversion to failed attacks that are caused by reliability problems after making it through the defense search point.

There are several important implications of our findings. The Department of Homeland Security (DHS) in October 2010 issued the following statement: “As a precaution, DHS has taken a number of steps to enhance security. [...] Passengers should continue to expect an unpredictable mix of security layers that include explosives trace detection, advanced imaging technology, canine teams and pat downs, among others.” As our theoretical model predicts and experimental results confirm, a security agency such as the DHS or the Transportation Security Agency, must respond optimally to terrorist organizations’ actions and pre-empt any terrorist attack by identifying terrorists within large groups of mostly innocent people. In this context, profiling is rational and the government should actually screen individuals according to their potential to be reliable recruits for the terroristic organization. The security agency should search more often the individuals belonging to the most reliable categories with an apparently unfair profiling practice. The intense search directed toward high reliability individuals should induce the terrorists to send less reliable categories more often. Our findings should be interpreted with caution, however, since our experiment does not shed light on the indirect effects of profiling, e.g., the possible increases in a propensity for terrorist activity among groups being profiled (Harcourt, 2007).

There are many possible avenues for future research. For example, it would be interesting to investigate both theoretically and experimentally a non-constant sum version of the profiling game, where each player can choose to attack (or defend) more categories at once and the cost of attacking (or defending) is increasing in the number of categories attacked (or defended). Another avenue would be to introduce multiple attackers and defenders, where defenders trying to defend against multiple independent (or dependent) attackers. Also, it is important to examine the problem of profiling in the setting of incomplete information, i.e., when defenders and attackers know only the distribution of the reliability of each type. In such case, players would need time and experience to learn how reliable different types are. These are all very interesting questions and we leave them for future research.

## References

- Antonovics, K., and Knight, B. G. (2009). A New Look at Racial Profiling: Evidence from the Boston Police Department. *Review of Economics and Statistics*, 91(1), 163-177.
- Avrahami, J., Güth, W., and Kareev, Y. (2005). Games of Competition in a Stochastic Environment. *Theory and Decision*, 59(4), 255-294.
- Barak-Erez, D. (2007). Terrorism and Profiling: Shifting the Focus from Criteria to Effects. *Columbia Law Review*, 29(1), 1-9.
- Barnes, K., and Gross, S. (2002). Road Work: Racial Profiling and Drug Interdiction on the Highway. *Michigan Law Review*, 101, 653–754.
- Basuchoudhary, A., and Razzolini, L. (2006). Hiding in Plain Sight – Using Signals to Detect Terrorists. *Public Choice*, 128(1), 245-255.
- Chowdhury, S. M., Kovenock, D., and Sheremeta, R. M. (2013). An Experimental Investigation of Colonel Blotto Games. *Economic Theory*, 52, 833-861.
- Crawford, V. P., and Iriberri, N. (2007). Fatal Attraction: Salience, Naivete, and Sophistication in Experimental Hide-and-Seek Games. *American Economic Review*, 97, 1731-1750.
- Crawford, V. P., Costa-Gomes, M. A., and Iriberri, N. (2013). Structural Models of Nonequilibrium Strategic Thinking: Theory, Evidence, and Applications. *Journal of Economic Literature*, 51(1), 5-62.
- Dechenaux, E., Kovenock, D., and Sheremeta, R. M. (2012). A Survey of Experimental Research on Contests, All-Pay Auctions and Tournaments. Chapman University, Working Papers.
- Deck, C., and Sheremeta, R. M. (2012). Fight or Flight? Defending Against Sequential Attacks in the Game of Siege. *Journal of Conflict Resolution*, 56, 1069-1088.
- Ellmann, S. (2003). *Racial Profiling and Terrorism*. New York Law School Review, 46, 675-730.
- Harcourt, B. E. (2007). Muslim Profiles Post-9/11: Is Racial Profiling and Effective Counter Terrorist Measure and Does it Violate the Right to be Free from Discrimination? In B. J. Gould and L. Lazarus, Eds., *Security and Human Rights*. Oxford: Hart Publishing, 73-98.
- Heymann, P. B., and Kayyem, J. N. (2005). *Protecting Liberty in an Age of Terror*. Cambridge, MA: MIT Press.
- Holt, C. A., Kydd, A., Razzolini, L., and Sheremeta, R. M. (2011). Theory and Experiments on Profiling and Terrorism. Draft NSF Proposal.

- Kahneman, D., & Tversky, A. (1979). Prospect Theory: An Analysis of Decision under Risk. *Econometrica*, 47, 263-291.
- Knowles, J., Persico, N., and Todd, P. (2001). Racial Bias in Motor Vehicle Searches: Theory and Evidence. *Journal of Political Economy*, 109(1), 203-229.
- Kovenock, D., Roberson, B., and Sheremeta, R. M. (2010). The Attack and Defense of Weakest-Link Networks. Chapman University, Working Paper.
- Kydd, A. (2011). Terrorism and Profiling. *Terrorism and Political Violence*, 23(3), 458-73.
- McKelvey, R., and Palfrey, T. (1995). Quantal Response Equilibria for Normal Form Games. *Games and Economic Behavior*, 10, 6-38.
- Rubinstein, A., Tversky, A., and Heller, D. (1996). Naive Strategies in Competitive Games. In Albers, W., W. Guth, P. Hammerstein, B. Moldovanu and E. van Damme, Eds., *Understanding Strategic Interaction - Essays in Honor of Reinhard Selten*, Springer-Verlag, 394-402.
- Sandler, T., and Arce, D. G. (2007). Terrorism: a game-theoretic approach. In T. Sandler and K. Hartley, Eds., *Handbook of Defense Economics: Defense in a Globalized World*, vol. 2, Amsterdam: North-Holland.
- Stahl, D., and Wilson, P. (1994). Experimental Evidence on Players' Models of Other Players. *Journal of Economic Behavior and Organization*, 25, 309-327.
- Stahl, D., and Wilson, P. (1995). On Players' Models of Other Players: Theory and Experimental Evidence. *Games and Economic Behavior*, 10, 218-254.

### Appendix - Instructions

- **Rounds and Matching:** The experiment consists of a number of rounds. Note: You will be matched with the same person in all rounds.
- **Interdependence:** Your earnings are determined by the decisions that you and the other person make.
- **Roles:** In each pair of people, one person will be given the role of "attacker" and the other will be given the role of "defender." Your role will be (attacker or defender) in all rounds.
- **Locations:** There are 2 locations that will be designated as: L1 and L2. An attack can be targeted to either of these locations, and a defense can be augmented at either of these locations.
- **Attack Success Probabilities:** An attack will always fail at a site that is defended. An attack at an undefended site may or may not succeed, and the probability of attack success at undefended sites will depend on the site, as explained later.

#### Instructions (page 2)

<b>Location:</b>	<b>L1</b>	<b>L2</b>
<b>Attacker Gain from Successful Attack:</b>	<b>\$1.00</b>	<b>\$1.00</b>
<b>Defender Loss from Successful Attack:</b>	<b>\$1.00</b>	<b>\$1.00</b>
<b>Position Your Asset</b>	<input type="checkbox"/> <b>L1</b>	<input type="checkbox"/> <b>L2</b>

- **Attacker Gains:** If an attack is successful at a location, the attacker earns an amount of money shown in the Attacker Gain row of the table, for that location.
- **Defender Losses:** A successful attack at any location results in a loss to the defender, as shown in the Defender Loss row of the table, for that location.
- **Available Assets:** Each attacker has 1 asset to allocate (1 attack), and each defender has 1 asset to allocate (1 defense).

#### Instructions (page 3)

- **Attack Outcomes:** If a site is defended, an attack at that site will fail. The chances of a successful attack at an undefended site depend on the site, as shown in the table below, which will be reproduced for you when you submit your decision.

<b>Location:</b>	<b>L1</b>	<b>L2</b>
<b>Probability of Attack Success at a Defended Site:</b>	<b>0</b>	<b>0</b>
<b>Probability of Attack Success at an Undefended Site:</b>	<b>0.67</b>	<b>0.33</b>

- **Random Outcome Determination:** Consider an attack on site L2. If this site is defended, the attack will fail. If this site is undefended, the probability of attack success at that site is 0.33. You can think of this process as spinning a Roulette wheel with

stops labeled 1, 2, ... 100 and the outcome is a success if the wheel stops on a number that is less than or equal to 33, so a probability of 0.33 corresponds to 33 chances out of 100 of attack success.

- **View Failed Attacks:** After all decisions are made and confirmed, the defender will always be able to see where an attack occurred, even if it fails.
- **Visibility of Defense Assets:** The attacker will NOT be able to see where a particular defense asset is located prior to making an attack decision.
- **Cause of Failed Attack:** At the end of each round the attacker's results table will indicate whether a site was defended or not. Thus if an attack does fail, the attacker will be able to see whether it failed because the site was defended or because the attack at an undefended site failed due to random causes.
- **Private Incomes:** In addition to the earnings, losses, and costs that result from asset allocations and attack outcomes, each person will receive a fixed income in each round. Attackers and defenders may have different private incomes, which are not public information. As a (Attacker or Defender), your income will be: \$\*.\*\* per round.

### Instructions (summary page)

- There will be one or more rounds in this part of the experiment, and the final round will not be announced in advance.
- You will be matched with the **same** person in all rounds.
- In each group, there will be **1 attacker** and **1 defender**.
- Your role is that of \*\*\*\*\*
- Defenders each have 1 asset to allocate across the 2 sites, and attackers each have 1 asset to allocate to one of the 2 sites.
- If a site is defended, an attack at that site will fail. The chances of a successful attack at an undefended site depend on the site, as shown in the table below.
- A successful attack at a site reduces the earnings for the defender, as indicated by the Defender Loss for that site.
- A successful attack at a site increases the earnings for the attacker, as indicated by the Attacker Gain for that site.
- The defender will always be able to see where an attack occurred in previous rounds, even if it fails.
- The attacker will not be able to see where a defense asset is located (before the attack decision is made).
- In addition, your payoff will be raised by an amount \$\*.\*\* in each round, which is your private income, not observed by the others with a different role.
- **Special Earnings Announcement:** Your cash earnings will be **50%** of your total earnings at the end of the experiment.