

Chapter 3

The Theory of Butter-for-Bombs Agreements: How Potential Power Coerces Concessions

This chapter develops a formal model of costly power shifts. It shows that butter-for-bombs settlements can be sustainable in the long term, even if the rising state can freely renege. Depending on the parameters, the interaction ends in one of three ways. First, if the extent of the power shift is too great, the declining state can credibly threaten preventive war, which in turn makes the rising state's threat to build weapons incredible. Likewise, if the rising state's cost of building is too high, the declining state knows the rising state will never build. In either case, the declining state can offer the rising state no concessions and still induce acceptance. The outcome mirrors a world in which the rising state had no ability to shift power.

Second, if the threat to build is credible but investment costs remain relatively large, the declining state optimally offers immediate concessions to the rising state. The rising state accepts those concessions in the present and continuously in the future. Although the rising state could build and force the declining state to give yet more concessions, those additional concessions do not cover the cost of building. Thus, the rising state extracts concessions using *unrealized* power and maintains the status quo because of the attractiveness of future offers. This in turn allays the fears of the declining state.

Finally, if the cost of shifting power is low, the declining state cannot cheaply buy off the rising state. As a result, the declining state chooses to shortchange the rising state initially, forcing the rising state to shift power. Afterward, the declining state makes great concessions. The declining state could still induce the rising state not to build here, but it simply profits more from stealing as much as it can upfront. Put differently, the declining state's opportunism—not the rising state's opportunism—leads to the shift in power.

The results of the model indicate that, in the context of a bargaining game, the demand for proliferation is rare. In some cases, the declining state's threat of preventive war deters the rising state from building. In other cases, the rising state finds weapons more costly than useful. In between, the declining state can buy off some of the remaining states. Proliferation only occurs in the model when the investment cost is low.¹

This chapter has five additional sections. The next section formally defines the model, describes some key features of the interaction, and derives its solution; in equilibrium, declining states and rising states reach peaceful, stable agreements if the cost to shift power falls within a certain range. The following section explores the robustness of the model, demonstrating that butter-for-bombs bargaining persists under many alternative model specifications. The next two sections interpret the results and broadly illustrate the implications of the model on arms investment, negotiated agreements, and preventive war. A brief conclusion follows.

3.1 Modeling Butter-for-Bombs Agreements

This section introduces the central bargaining model of the book. First, it describes the strategic interaction. Next, it highlights the key features of the model that depart from previous formal work on shifting power. With that, it then derives the game's equilibria and shows that the declining state sometimes offers immediate concessions to convince the rising state not to build, even when conditions appear ripe for proliferation. Lastly, a numerical example illustrates equilibrium game play.

¹However, as Chapter 5 shows, states have incentive to create equilibrium institutions to artificially raise the cost of building. Thus, an inefficiency puzzle remains, which the remainder of the book will address.

3.1.1 Actions and Transitions

Consider an infinite horizon game between two actors, D (the declining state) and R (the rising state), as illustrated in Figure 3.1.² The states bargain over a good standardized to value 1. There are four states of the world: pre-shift bargaining, post-shift bargaining, pre-shift war, and post-shift war. The last two are absorbing.

The game begins in the first period in the pre-shift state, before R develops the weapons technology. D makes a temporary offer $x_t \in [0, 1]$ to R, where t denotes the period. R accepts, rejects, or builds in response. Rejecting results in game ending war; R receives $p_R \in [0, 1)$ in expectation while D receives $1 - p_R$. These payoffs persist through all future periods in this absorbing state, but the states pay respective costs $c_D, c_R > 0$ in each future period regardless.³

If R accepts, the period ends. R receives x_t for the period while D receives $1 - x_t$. The game then returns to this same pre-shift bargaining state, where D makes another temporary offer x_{t+1} .

If R builds, it pays a cost $k > 0$ to begin constructing the new weapons.⁴ D sees this and decides whether to initiate a preventive war or advance to the post-shift state of the world.⁵ Preventive war ends the game and results in the same terminal payoffs as though R had rejected D's offer x_t . If D advances, the period ends, and R receives x_t for the period while D receives $1 - x_t$.

If R successfully builds, the game transitions into the post-shift state, and R's outside option of war improves in all future periods. Similar to before, D makes an offer y_{t+1} to R in such a post-shift period. If R accepts, the period ends, R receives y_{t+1} for the period, D receives $1 - y_{t+1}$ for the period, and the game repeats the post-shift bargaining period, where D makes another offer y_{t+2} . If R rejects, a game-ending war results. Here, R takes $p'_R \in (p_R, 1]$ in expectation while D receives $1 - p'_R$. That is, R expects to receive more from war with the weapons than without, whether because those weapons shift the balance of power or limits D's war aims in the manner Chapter 2

²These labels are a convention from the literature. In the basic model, the rising state rarely rises and the declining state rarely declines. Proliferation decisions occur more frequently in the extensions explored in later chapters.

³The results are the same if costs are only paid in the period of fighting. Moreover, the proof is identical except that we must substitute c_i with c'_i , where $c'_i = \frac{c_i}{1-\delta}$.

⁴Since the bargaining good is fixed at value 1, k implicitly reflects R's resolve as well.

⁵Chapter 7 relaxes this assumption so that D has no direct knowledge whether R built.

described. These payoffs again persist through time in this absorbing state, but the sides still pay their respective costs c_D, c_R .⁶

The states share a common discount factor $\delta \in (0, 1)$. Thus, the states discount period t 's share of the good and costs paid by δ^{t-1} . The discount factor reflects two underlying parameters. First, as is standard, greater values place greater weight on future payoffs. Second, and common to models of shifting power, δ also represents the time it takes R to successfully develop its new weapon. Ineffective programs correspond to lower values, as more time must pass before the states renegotiate their terms of settlement.

3.1.2 Key Features

Before solving for the game's equilibria, four important features of the model are worth highlighting. First, following the second wave of shifting power research (Jackson and Morelli 2009; Chadeaux 2011; Fearon 2011; Debs and Monteiro 2013), the power shift is costly and endogenous. These are minimalist and necessary criteria. The vast majority of major power shifts result from endogenous choices made by rising states (Debs and Monteiro 2013, 4-5). Moreover, keeping power shifts exogenous prohibits the states from bargaining over weapons, since strength appears by assumption. As the robustness checks section will later show, disallowing bargaining leaves both sides worse off than if they could bargain over the weapons program.

Second, the model allows the interaction to continue forever. If the rising state were to lose the ability to proliferate at any point, it would have to build in the periods previous to force the declining state to offer concessions. As such, the rising state maintains the ability to proliferate in every pre-shift period. Later chapters will address what happens if the rising state might be unable to proliferate at a future date through endogenous actions.⁷

Third, the model only permits one-sided armament. This would appear to stack the deck against nonproliferation agreements. If arms races are a form of a repeated prisoner's dilemma, then an explanation for arms treaties seems to exist already—neither proliferates because the other side will proliferate in response, in a manner similar to grim trigger strategies or tit-for-tat (Axelrod

⁶Similar results would obtain if the costs of war changed in the post-shift state.

⁷It is trivial to show that rising states acquire nuclear weapons if proliferation is a now-or-never opportunity, but it is odd to *assume* that a rising state would suddenly lose the ability to proliferate, especially since such an outcome leads to a commitment problem and inefficiency.

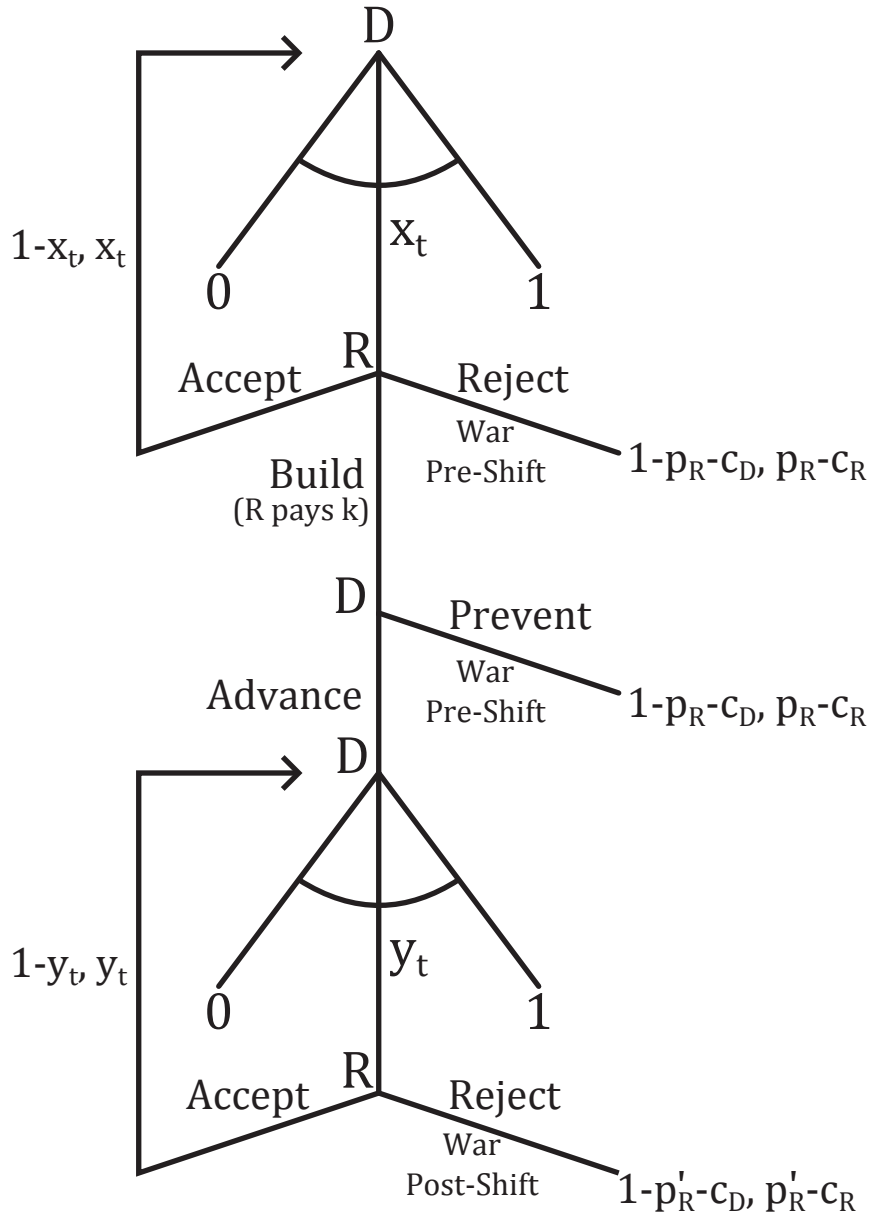


Figure 3.1: The model. All payoffs listed are for the period, though the war outcomes lock in their respective payoffs every period for the rest of time.

1984). Similarly, and in contrast to Kydd (2000), the declining state cannot adjust its military spending to compete with the rising state's armament decision. Thus, any form of nonproliferation agreements must result from a different mechanism.⁸

Finally, the model puts the declining state in a strategically vulnerable position—it must offer a division of the stakes to the rising state before the rising state chooses whether to build, and the declining state cannot retract that offer should the rising state proliferate.⁹ A major policy concern with Iran is that Tehran could take the concessions the United States offers it, renege on any quid-pro-agreement not to build, and proliferate anyway. That being the case, the Washington ought not to give any concessions, since any hypothetical bribe would not alter Iran's endgame behavior. Ordering the moves in this manner directly addresses the policy concern.

It is worth stressing that that structuring bargaining in this way means that the declining state does not directly negotiate over the rising state's weapons program. Given concerns regarding anarchy, making quid-pro-quo offers as in Chadeaux's (2011) model raises the question why rising states simply do not renege after receiving the concessions or why declining states do not renege after the rising state temporarily suspends its weapons program. Yet, interestingly, removing quid-pro-quo bargaining in favor of indirect bargaining does not stack the deck against cooperation.

In that vein, Debs and Monteiro (2013) analyze a similar model in which a rising state chooses whether to develop weapons programs in secret. However, their focus is on the how the secret nature of some weapons programs leads to preventive war. Thus, the structure of their game precludes analysis of negotiating over weapons.¹⁰ This model instead brings negotiations to the

⁸In fact, Chapter 5 shows that the mechanism described in the model sabotages tit-for-tat or grim trigger strategies in two-sided proliferation games precisely because of the attractiveness of butter-for-bombs deals.

⁹The robustness section of this chapter shows that butter-for-bombs settlements are substantially easier to reach when no such vulnerability exists.

¹⁰Specifically, in the finite version of their model, the the rising state chooses whether to build weapons at the beginning of each period. Cooperation is inherently impossible in such a setup because the rising state's decision dictates the period's terms of bargaining. If the rising state chooses not to build, then the declining state needs to only offer the rising state an amount to avert war in the pre-shift state of the world. Therefore, the rising state cannot threaten to proliferate if it receives poor offers, which in turn leads to the rising state receiving bare-bones concessions if it opts not to build at the start. As a result, the rising state must build at the beginning to receive substantial concessions at

forefront. In turn, the rising state acts more strategically, selecting to build only if its offers are insufficient. Knowing this, the declining state has incentive to offer large amounts upfront to dissuade the rising state from building. Both sides finish better off as a result.

3.1.3 Equilibrium

Since this is a dynamic game with an infinite number of periods, this section searches for stationary Markov perfect equilibria (MPE). A stationary MPE is a set of strategies that are sequentially rational, depend only on the state of the world, and specifically are not a function of calendar time.

Before stating the main results, the following lemma will prove useful:

Lemma 3.1. *In every stationary SPE, in every post-shift period, D offers $y_t = p'_R - c_R$, and R accepts.*

The appendix provides a complete proof of Lemma 3.1. However, the intuition is a straightforward application of Fearon's seminal bargaining game. Since war creates deadweight loss to the system, D can always offer enough to satisfy R, and the optimal acceptable offer is preferable to war for D as well. Thus, D offers just enough to induce R to accept, and D keeps all of the surplus. In particular, R earns $p'_R - c_R$ and D earns $1 - p'_R + c_R$ for the rest of time, and peace prevails in the post-shift state of the world.

Overall, Lemma 3.1 shows that R has great incentive to proliferate—nuclear weapons mean greater coercive power, forcing D to offer more concessions to maintain the peace. Consequently, it is not remarkable that declining states want rising states to commit to nonproliferation agreements. What is surprising—especially to those who focus on difficult cases like North Korea, Iraq, and Iran—is that rising states can credibly abide to such deals. The following theorem summarizes the finding:

Theorem 3.1. *For all parameters, R is willing to accept peaceful, nonproliferation settlements.*

Two factors contribute to R's credible commitment to not build: (1) R's satisfaction with future compensation and (2) R's desire to avoid paying for costly weapons. The proof is simple and illuminating. Suppose D offers $x_t \geq \max\{p_R - c_R, p'_R - c_R - \frac{k(1-\delta)}{\delta}\}$ in every pre-shift period. Showing that

any point.

R accepts is sufficient to prove the theorem. R has two alternatives: reject and build. Rejecting yields $p_R - c_R$. However, since D offers at least $p_R - c_R$ in this conjecture, rejecting is not a profitable deviation.

If R builds, D's optimal response is either to prevent or advance to the post-shift state. In the first case, R earns $p_R - c_R$. For the same reason as before, this is not a profitable deviation. In the second case, R receives x_t for the period, earns $p'_R - c_R$ in all subsequent periods (from Lemma 3.1, and pays the cost k . Recall that the value for accepting all offers is at least $p'_R - c_R - \frac{k(1-\delta)}{\delta}$. Thus, accepting is at least as good as building if:

$$p'_R - c_R - \frac{k(1-\delta)}{\delta} \geq (1-\delta)x_t + \delta(p'_R - c_R) - (1-\delta)k$$

$$x_t \geq p'_R - c_R - \frac{k(1-\delta)}{\delta}$$

This holds. Since $\max\{p_R - c_R, p'_R - c_R - \frac{k(1-\delta)}{\delta}\} < 1$ D can always make this offer. Therefore, R is willing to accept some peaceful, nonproliferation settlement.

Before moving on to the solution of the game, there are two notes. First, this is a general result. It is *not* dependent on the structure of the bargaining protocol in pre-shift periods, as it simply states that R prefers these outcomes to investment outcomes. Second, the theorem is *not* an equilibrium claim. Rather, it is a “possibility” theorem—it proves that deals are possible provided that D is willing to make the necessary concessions. Interestingly, this flips the apparent credibility problem of nuclear negotiations. R is always willing to negotiate. The question is whether D is willing to make serious offers. The propositions below address the issue by solving for the game's equilibria.

Proposition 3.1. *If $p'_R - p_R > \frac{c_D + c_R}{\delta}$, D offers $x_t = p_R - c_R$ in the unique stationary MPE. R accepts these offers and never builds.*

Note that the left side of the inequality represents the extent of the power shift and the right side represents the inefficiency of war. When the shift is sufficiently greater than war's inefficiency, the power shift is “too hot.” If R were to build, D would respond with preventive war. As a result, the credible threat of a fight makes R's threat to build incredible. In turn, D can treat the bargaining problem as though R cannot build. Consequently, D offers

$x_t = p_R - c_R$ (the amount R would receive in a static bargaining game), R accepts, and the states avoid war.¹¹

The appendix contains a complete proof. Intuitively, the critical value of $p'_R - p_R$ comes from finding the value for which D prefers preventive war if it offers $x_t = p_R - c_R$ and R attempts to build:

$$1 - p_R - c_D > (1 - \delta)(1 - p_R + c_R) + \delta(1 - p'_R + c_R)$$

$$p'_R - p_R > \frac{c_D + c_R}{\delta}$$

This is the critical value of $p'_R - p_R$ presented in Proposition 3.1.

Note that when D can deter R with the “stick” or preventive war, R receives no “carrots” in the form of butter-for-bombs agreements. This is because credible threats are free for D whereas concessions are costly. As such, D has no reason to give R some of the surplus when it can take all of the benefits for itself. But also note that if the costs of war are sufficiently high—that is, as preventive war becomes sufficiently ineffective—D can never use the preventive threat to deter R.¹² Thus, for some parameters, D must negotiate with R even if the extent of the power shift is extreme.

Proposition 3.2. *If $p'_R - p_R < \frac{k(1-\delta)}{\delta}$, D offers $x_t = p_R - c_R$ in the unique stationary MPE. R accepts these offers and never builds.*

Note the right side of the inequality reflects the time-adjusted cost of building. When the magnitude of the shift is too small relative to that cost, the power shift is “too cold” for the rising state to invest in weapons. D observes that R does not have a credible threat to build and therefore offers the same concessions it would offer if power were static. As a result, though for different reasons, the observable outcome for these parameters are the same as the outcome for Proposition 3.1’s parameters.

The full proof appears in this chapter’s appendix. The critical value of $p'_R - p_R$ comes from finding the value for which R will not build in response to $x_t = p_R - c_R$:

$$p_R - c_R > (1 - \delta)(p_R - c_R) + \delta(p'_R - c_R) - k(1 - \delta)$$

¹¹The fact that preventive war does not occur here should be unsurprising since the game has complete information and power shift is observable and endogenous (Chadefaux 2011).

¹²See Reiter 2006 for a pessimistic outlook on preventive war.

$$p'_R - p_R < \frac{k(1-\delta)}{\delta}$$

This is the critical value of $p'_R - p_R$ presented in Proposition 3.2.

The remaining cases maintain that $p'_R - p_R$ does not fall into the previously discussed cases and assume that $k \in (\frac{\delta(p'_R - p_R - c_D - c_R)}{1-\delta}, \frac{\delta p'_R - p_R}{1-\delta} + c_R)$. The minimum value constraint for k implies D earns more from engaging in a butter-for-bombs settlement than it does from earning its war payoff in the pre-shift stage. The maximum value constraint ensures that R never prefers rejecting to a successful power shift, even if D offers R nothing during the pre-shift periods. For the purposes of this chapter, these cases are theoretically uninteresting and offer no further insight to the analysis.

Proposition 3.3. *If $k > \delta(p'_R - c_R)$, D offers $x_t = p'_R - c_R - \frac{k(1-\delta)}{\delta}$ in all pre-shift periods in the unique efficient stationary MPE; R accepts and never builds.*

Moving outside the parameters of Proposition 3.1 and Proposition 3.2 leaves the world “just right” for a power shift. Nevertheless, if the corresponding investment remains relatively costly, D prefers making immediate concessions. If R were to build in response, it would receive additional concessions in the post-shift state of the world. However, those additional concessions do not cover the costs of building, which in turn convinces R to accept the original offer. In effect, D manipulates R’s opportunity cost for building to the point that investment is no longer profitable.

The appendix contains proof for Proposition 3.3. The equilibrium offer size $p'_R - c_R - \frac{k(1-\delta)}{\delta}$ equals R’s continuation value for building, which is enough to induce R to accept. Note that D receives the remainder, or $1 - p'_R + c_R + \frac{k(1-\delta)}{\delta}$. For D to prefer taking that amount over the long-term to taking everything up front and suffering the consequences of proliferation later, it must be that the investment cost is relatively large, or:

$$1 - p'_R + c_R + \frac{k(1-\delta)}{\delta} > (1-\delta) + \delta(p'_R - c_R)$$

$$k > \delta(p_R - c_R)$$

This is the critical value for k in Proposition 3.3.

Note that Proposition 3.3 states that this is the unique *efficient* stationary MPE. A second MPE exists as well. In this equilibrium, D offers $x_t = 0$ in every period and R builds regardless of the offer. D has no profitable

deviation; R builds regardless of the offer size, so D ought to steal everything it can upfront. R has no profitable deviation either; the lack of *quid-pro-quo* bargaining means that R keeps the offer x_t regardless of whether it builds. Thus, R is willing to build no matter what D offers.

Such “no deal” equilibria are common in games that require mutual compliance to achieve a cooperative outcome. However, no one wins in this equilibrium—*both* are better off in the efficient equilibrium where bargaining succeeds. It is thus natural to focus on Proposition 3.3’s equilibrium.

Before moving on, a couple comparative statics from Proposition 3.3 recur throughout this book, so it is worth understanding them. First, as the extent of the power shift $(p'_R - p_R)$ increases, D’s bribe increases. This might seem counterintuitive—a large power shift means that R is comparatively weak at the beginning. Yet R can use this exact weakness to its advantage because its outside option (investing in weapons) is correspondingly desirable. As such, D must give larger concessions to induce R to accept.¹³

Second, D’s offer decreases as k increases. That is, R receives better butter-for-bombs offers the smaller its investment cost is. Although engaging in butter-for-bombs deals means that R does not build, D knows it can keep more for itself and still induce compliance if the cost to invest is comparably more expensive.

That leaves the final proposition:

Proposition 3.4. *If $k < \delta(p'_R - c_R)$, D offers $x_t = 0$ in all pre-shift periods in the unique stationary MPE; R builds and D does not prevent.*

The proof and intuition follow from Proposition 3.3. Note that the minimalist butter-for-bombs offer $p'_R - c_R - \frac{k(1-\delta)}{\delta}$ increases as k decreases. Thus, the remainder for D decreases as k decreases. So if k is sufficiently small, D takes everything upfront, R proliferates, and D makes great concessions later.

More formally, in any period, R could build and receive a continuation value of $\delta(p_R - c_R) - (1 - \delta)k$. Thus, for R to accept an offer, D must give R at least $p_R - c_R - \frac{k(1-\delta)}{\delta}$ on average in all future periods. The remainder for D equals $1 - \delta p'_R + \delta c_R$. However, D could eschew bargaining, demand the whole pie for the period, and (at worst) induce R to build. D earns

¹³This will be a recurring theme in the later chapters’ case studies: a state feels vulnerable, threatens to proliferate, and extracts concessions. The weakness becomes a bargaining strength.

$1 - \delta p'_R + \delta c_R$ by doing so. This is worth more than attempting bargaining if $k < \delta(p'_R - c_R)$, which is true for the parameter space. \square

There are two perspectives on what causes bargaining to break down here. D deserves part of the responsibility. The proof of Theorem 3.1 shows that R can always credibly commit to accepting any $x_t \geq p'_R - c_R - \frac{k(1-\delta)}{\delta}$. As such, D can always bargain away the problem if it wants to. The issue here is that D prefers taking as much as it can upfront and suffering the consequences later on. Thus, *proliferation benefits D* here, not R.

Given that D is unwilling to buy off R, a commitment problem exacerbates the dilemma. Note that D earns $1 - \delta(p'_R + c_R)$ and R earns $\delta(p'_R - c_R) - k(1 - \delta)$ for this outcome. Consequently, both sides would be better off if R could credibly commit to accepting any x_t greater than $\delta(p'_R - c_R) - k(1 - \delta)$ but less than $\delta(p'_R + c_R)$. However, R's best response to any such offer is to build. Anticipating this, D minimizes its initial offer and accepts the inefficient outcome. Interestingly, and as Chapter 5 discusses at length, the inefficiency benefits no one—both parties would be better off if k were higher and the states reached the equilibrium outcome Proposition 3.3 describes.

3.1.4 Numerical Example

To illustrate the logic of the butter-for-bombs equilibrium, consider the following specific environment. Let $p_R = .2$, $p'_R = .5$, $c_D = .3$, $c_R = .1$, $\delta = .9$, and $k = 1$. These values fit the parameters for Proposition 3.3.

If the game ever reaches the post-shift state, Lemma 3.1 states that D offers $x_t = p'_R - c_R = .5 - .1 = .4$ in every post-shift period. R accepts those offers and earns .4; D earns the remainder, or $1 - x_t = .6$.

In the pre-shift stage, if D offers $x_t = 0$, R earns 0 for accepting and $p_R - c_R = .2 - .1 = .1$ for rejecting. If R builds, D does not prevent, as it earns $1 - p_R - c_D = 1 - .2 - .3 = .5$ for preventing and $1 - \delta + \delta(1 - p'_R + c_R) = 1 - .9 + .9(1 - .5 + .1) = .64 = \frac{32}{50}$ for advancing to the next period. In turn, R earns $0 + \delta(p'_R - c_R) - k(1 - \delta) = .9(.5 - .1) - (1 - .9) = .26$ for building, which is more than it receives for rejecting or accepting.

Alternatively, if D offers $x_t = p'_R - c_R - \frac{k(1-\delta)}{\delta} = .5 - .1 - \frac{1(1-.9)}{.9} = \frac{13}{45}$, R earns $\frac{13}{45}$ for accepting. In contrast, R receives only .1 for rejecting and $\frac{13}{45}$ at most for building, so accepting is optimal. D earns $\frac{32}{45}$ for this outcome, whereas it receives only $\frac{32}{50}$ if it offers $x_t = 0$. Therefore, $x_t = \frac{13}{45}$ is the optimal offer for D.

Figure 3.2 illustrates the bargaining dynamics of this specific example

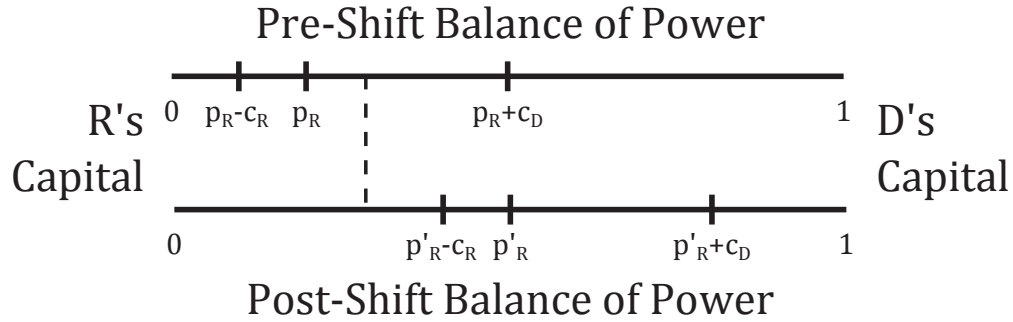


Figure 3.2: The pre-shift and post-shift balances of power. The figure is drawn to scale for the numerical example. In equilibrium, R receives $p'_R - c_R - \frac{k(1-\delta)}{\delta}$ (the dashed line) in every period, and D receives the remainder.

drawn to scale, conceptualized as R and D negotiating over a strip of territory between their respective capitals. The top half represents the pre-shift balance of power. If the rising state rejects the offer or D prevents, the states fight to an expected outcome of p_R , but pay their respective costs c_R and c_D . The bottom half represents the post-shift balance of power. If the states fight here, the average outcome swings in R's favor to p'_R .

The dashed line at $p'_R - c_R - \frac{k(1-\delta)}{\delta}$ is the equilibrium outcome. R must receive more than the minimum acceptable amount ($p_R - c_R$) in a static bargaining game, otherwise it can profitably shift power and enjoy great concessions in the future ($p'_R - c_R$). However, D need not offer $p'_R - c_R$ to induce R not to build. Indeed, the equilibrium outcome sees R receive less in every period than it would in the future if R actually shifted power. D effectively leverages the cost of building against R; the difference between $p'_R - c_R$ and the equilibrium outcome over time equals the discounted cost that R pays to build. The efficiency of the equilibrium outcome ensures D's satisfaction; because R never pays the inefficient cost k , D can extract it out of the negotiated settlement.

3.2 Robustness

As with any stylized model, it is worth asking whether the results are a function of the particular modeling choices or indicative of a broader underlying mechanism. The previous section showed the existence of a bargaining range.

Any settlement within that range leaves both parties better off than a world with investment, and the rising state cannot profitably renege under those terms. Thus, the butter-for-bombs result is not sensitive to particular bargaining protocols. However, the next question is the reasonableness of the structural assumptions. This subsection addresses a few possible issues.

In general, the robustness checks show that butter-for-bombs deals persist under a variety of alternative specifications. Meanwhile, the proliferation outcome Proposition 3.4 describes disappears under a variety of conditions. As such, and as Chapter 5 will pick up on later, the “low cost” explanation for proliferation is unsatisfying, leaving the remainder of the book to provide mechanisms with superior explanatory power.

Punishment for Reneging. Consider off the equilibrium path play in the butter-for-bombs parameter range. If R builds despite D’s generous offer, it receives no punishment. Instead, R keeps everything D offered it for the period and then receives additional concessions once R obtains nuclear weapons. In effect, R is only rewarded for proliferating and faces absolutely no punishment for defying D. This was a deliberate modeling choice, as policymakers worry that immediate concessions leave declining states in this strategically vulnerable position. And even in this worst case scenario, butter-for-bombs agreements worked.

Still, one may wonder what happens if D can punish R in response to R building. For example, D might offer R seats or voting shares in international organizations to induce R not to proliferate; D could easily revoke this if R were to renege. Rather than assuming that R keeps its entire share of the offer if it builds, suppose instead that R only keeps $\alpha \in [0, 1)$ of the offer; equivalently, D recoups $1 - \alpha$ of the concessions.¹⁴ Intuitively, if D attempts to buy R’s compliance but fails, D can cut the remainder of the payment.

Following the proof for Proposition 3.3, R is willing to accept x_t if:

$$x_t \geq \alpha(1 - \delta)(x_t) + \delta(p'_R - c_R) - (1 - \delta)k$$

$$x_t \geq \frac{\delta(p'_R - c_R) - (1 - \delta)k}{1 - \alpha(1 - \delta)}$$

Note that this amount is strictly less than the amount D had to pay previously. This is unsurprising—if D can recoup a portion of its bribe, R

¹⁴When $\alpha = 1$, the interaction is the original model.

finds investment less profitable and is therefore willing to accept a wider range of offers. Thus, butter-for-bombs settlements still exist and become easier to agree on under such conditions.

Moreover, note that for sufficiently small values of α , the proliferation parameters of Proposition 3.4 disappear. To see this, recall that D prefers making the butter-for-bombs offer if the remainder over time is greater than taking everything upfront and suffering the consequences of proliferation later. However, when $\alpha = 0$, D merely has to offer $\delta(p'_R - c_R) - (1 - \delta)k$ to induce R's compliance. Thus, nonproliferation agreements exist if:

$$1 - \delta(p'_R - c_R) - (1 - \delta)k > 1 - \delta + \delta(1 - p'_R + c_R)$$

$$c_R + k > 0$$

Both c_R and k are strictly positive, so the inequality trivially holds. Since D can recoup more and more of the bribe, nonproliferation agreements are easier to sustain. D thus opts exclusively for butter-for-bombs settlements here.

In addition, such *quid-pro-quo* styled offers also eliminate the inefficient equilibrium in Proposition 3.3's parameter space. Such a "no deal" equilibrium only existed because R received the offer x_t regardless of its build decision. In turn, D could not raise its offer from $x_t = 0$ and induce R to accept. Punishment for renegeing breaks this in-period indifference and thus allows D to raise its offer for the period. Only the efficient equilibrium remains.

International institutions commonly adopt measures that tilt the scales in favor of the efficient equilibrium. Specifically, states design institutions to reduce transaction costs (Keohane 1984). This effectively increases the pace at which actions can occur (Stone, Slantchev, and London 2008). Here, that means more opportunities for D to catch R in violation of the agreement and retract an offer. And, indeed, a primary task for the International Atomic Energy Agency (IAEA) is to monitor compliance to nonproliferation agreements. The IAEA can report violations back to leading states, who can then cut inducements.

Bargained resolutions are more likely when the time to proliferation is great. This is for two reasons. First, the delay gives D more time to discover R's violation of the agreement and withdraw concessions.¹⁵ Second, once

¹⁵Implicit in this argument is that weapons programs are not readily observable. See

discovered, D recoups more of the offer for longer. Both factors encourage D to bargain since the greater punishment discourages R from breaking the deal.

Subgame Perfect Equilibria. Stationary Markov perfect equilibrium is a subset of subgame perfect equilibrium (SPE). Solutions to infinite horizon games of this type often fail to characterize SPE due to the level of complexity added when strategies can be a function of a game's history and not just the state of the world. Fortunately, the proofs for Lemma 3.1, Proposition 3.1, Proposition 3.2, and Proposition 3.4 do not use the stationary Markov assumption; thus, the strategies listed are the unique subgame perfect equilibria.

On the other hand, this game has similar complexity and has an infinite number of SPE for Proposition 3.3's parameter space. The logic is essentially an application of the folk theorem. In every equilibrium, from period t forward, R must receive $\delta(p'_R - c_R) - (1 - \delta)k$, the value for receiving nothing in the current period and successfully building. Meanwhile, also from t forward, D must receive $1 - \delta(p'_R - c_R)$, the value for stealing everything in the current period and buying R off in the post-shift periods. Note that a surplus of $(1 - \delta)k$ exists. An equilibrium exists for every division of that surplus.

The logic is as follows. Recall that an inefficient stationary MPE exists for Proposition 3.3's parameter space. Its payoffs match the minimum values both players must receive in any equilibrium. Because MPE is a subset of SPE, the strategies are also an SPE. Consider any schedule of offers $x_t, x_{t+1}, x_{t+2}, \dots$ such that $\sum_{t=n}^{\infty} \delta^{t-1} x_t \geq \delta^n (p'_R - c_R) - (1 - \delta)k$ and $\sum_{t=n}^{\infty} \delta^{t-1} (1 - x_t) \geq 1 - \delta^n (p'_R - c_R)$ for all n .¹⁶ Such offers are supported in SPE if any deviation triggers a switch to the inefficient stationary MPE's equilibrium strategies. Neither D nor R could profitably deviate, since the schedule of offers at any time will ultimately yield a greater payoff than taking a short-term gain in the current period but receiving no surplus in the future periods.

It is worth noting that all of these equilibria disappear if D can recoup a

below for the corresponding discussion of imperfect information and Chapter 7 for a greater exposition on the problem.

¹⁶In words, the payoffs for continued acceptance from any period forward are greater than the minimum payoffs both players must receive in every equilibrium.

sufficiently large share of the bargaining good if R reneges. The logic is identical to the previous robustness check that eliminated the inefficient Markov perfect equilibrium: the potential loss from building breaks R's indifference between building and not building in any period. This allows D to induce R to accept sufficiently large offers, which prevents the mutual punishment trigger strategies from working. Without the threat of reverting to the Pareto inferior equilibrium, the folk theorem result breaks down.

Prior Investment in Nuclear Research. Suppose R has the option to pre-invest (or sink) a certain amount of the cost of nuclear weapons before reaching the bargaining phase, as is standard in research on in the proliferation research (Debs and Monteiro 2014). The original model reveals the outcome of this modification. Depending on how much R pre-invests, two outcomes are possible. First, if the pre-investment is small (and the remaining investment is sufficiently costly), the states will reach a butter-for-bombs agreement. Note that, compared to the original model, R performs better during the bargaining phase when it has already invested in nuclear weapons; with the remaining cost lower, D must give R a greater portion of the pie to successfully negotiate a butter-for-bombs settlement. However, R fares no better than before, as the improved payoff R receives in the bargaining phase equals the upfront cost R pays. On the other hand, D fares substantially worse if R pre-invests because it must compensate R for the pre-investment. This compensation ultimately comes out of D's share of the butter-for-bombs bargain, leaving D in worse shape.

Second, suppose the pre-investment is large. Then D keeps the entire good for itself in the first period, R builds, and then D makes great concessions thereafter. Butter-for-bombs fails. But this only leads to more inefficiency, since R finishes constructing the nuclear weapon in this case. Again, R does not profit from pre-investment, and D is strictly worse off.¹⁷

While international relations does not have a complete theory over how states choose to bargain¹⁸, it would nevertheless be strange to reach either of these outcomes. Inefficiency is understandable when at least one party benefits, incomplete information leads to miscalculations, or commitment problems prohibit mutually preferable alternatives. But here, inefficiency

¹⁷More precisely, *both* are strictly worse off. Chapter 5 explores this point in greater depth.

¹⁸For progress along these lines, see Stone 2011 and Leventoglu and Tarar 2005.

occurs for no discernable reason.

To elaborate, imagine that the states were supposed to interact in the traditional way, where R first chooses whether to build or not. If “anarchy is what states make of it” (Wendt 1992), why wouldn’t R ask D to restructure the interaction such that D first makes an offer to R? D would surely oblige, as it can obtain the surplus. Since negotiations lead to Pareto improvement, R has no reason to decline the restructuring. D therefore has strong incentives to proactively engage R in negotiations, ensure R’s investment cost remains as high as possible, and reach a butter-for-bombs deal. In turn, if the negotiations model presented here is the “wrong” model, a deeper question remains: why can’t states restructure an environment when it would lead to Pareto improvement? Good answers do not seem forthcoming. This is a strong indication that bargaining over weapons is the “right” way to model arms races. The extensive form game presented in this chapter permits such bargaining; previous attempts to explain proliferation do not.

Non-Binary Power Shifts. In the original model, R discontinuously jumps to power level p'_R if it constructs nuclear weapons. A more nuanced model might allow for R to choose its power level from a continuous range. For example, R could endogenously choose $k \in [0, \bar{k}]$, and $p'_R(k)$ is an increasing function that maps expense into power. Put simply, the more R invests, the greater military output it receives, and the more the balance of power shifts in R’s favor.

Butter-for-bombs bargaining survives in such a framework. Consider the maximum of $p'_R(k) - \frac{k(1-\delta)}{\delta}$, the function that measures output of investment minus the cost of the investment. Let k'' be its arg max and $p'_R(k'') = p''_R$. Suppose p''_R and k'' fit the parameters of Proposition 3.3’s butter-for-bombs outcome. Consider the offer $x_t = p''_R - c_R - \frac{k''(1-\delta)}{\delta}$. R cannot profitably deviate to building if offered this amount. Investing any less than k'' yields less realized power for the post-shift periods and thus a smaller payoff than accepting a stream of offers sized $p''_R - c_R - \frac{k''(1-\delta)}{\delta}$ for all of time. Meanwhile, investing any more than k'' yields additional concessions but is not worth the cost because k'' maximizes the tradeoff between building power and paying investment costs. Thus, the model’s binary power shift can be thought of as R deciding whether to pursue its best possible power shift.¹⁹

¹⁹Of course, $p'_R(k) - \frac{k(1-\delta)}{\delta}$ could be less than p_R for all k , in which case Proposition 3.2 applies—all power shifts are incredible and so D treats the situation like a static bargaining

Prestige. In the course of proliferating, many statesmen cite international “prestige” as a benefit to having nuclear weapons. Some researchers have shown concern regarding prestige as well.²⁰ While there are many reasons to be skeptical of this argument²¹, advocates might worry that the prestige negates the cost of proliferating k . Accordingly, k may drop to the critical value for which Proposition 3.4 predicts proliferation. In the extreme, k may even be negative.

Fortunately for the nonproliferation regime, this is a misinterpretation of the parameters. The cost parameter k only affects R’s payoffs directly. However, prestige is zero sum. If *all* states had nuclear weapons, for example, then nuclear weapons would not be prestigious. As such, if nuclear weapons truly provide prestige, each additional state that proliferates drains prestige from the status quo nuclear powers.

While k does not have a zero sum interpretation, recall in contrast that p_R and p'_R refer to a zero sum bargaining good. In a world where nuclear weapons provide prestige, the bargaining good instead represents the bargaining good *and* international prestige. Thus, prestige merely inflates the value of p'_R .²² It does not render nonproliferation agreements impossible.²³ Put differently, if nuclear weapons shift prestige from status quo nuclear states to new proliferators, the status quo states ought to find a way to buy off potential proliferators and reap the benefits of the saved investment costs.²⁴

game.

²⁰See Greenwood et. al. 1977 (50), Meyer 1984 (50-55), Quester 1991 (217), and O’Neill 2006. Gilpin (1981, 215) goes so far as to say “the possession of nuclear weapons largely determines a nation’s rank in the hierarchy of international prestige.”

²¹For example, the nonproliferation regime has succeeded in making nuclear weapons a taboo source of military power (Tannenwald 1999). It is also difficult to disentangle actual beliefs about prestige from bargaining posturing. See Thayer 1995 (468-474), Lavoy 1993 (197-199), and Mueller 2010 (105-108) for counterarguments. Generally, prestige may provide a good narrative for decisions to proliferate but lacks real causal power.

²²In that regard, it is clear why non-nuclear states claim that prestige exists while recognized nuclear weapons states claim the opposite. If it exists, declining states must concede more benefits to deter proliferation; if it does not, rising states receive no additional benefits.

²³States might face a problem if prestige is indivisible, but even then they could negotiate side payments (in the form of the continuous bargaining good) to avoid inefficient outcomes.

²⁴Indeed, the United States offered India “prestige” compensation to induce New Delhi to end its nuclear weapons program (Thayer 1993, 198).

Negative Externalities. Nuclear weapons might have consequences beyond the coercive bargaining relationship between R and D. One common concern is that a nuclear state’s safeguards could fail, allowing a rogue group or terrorist organization to obtain a nuclear weapon. In this manner, proliferation can impose negative externalities on both parties.

Incorporating these concerns into the model shows that nonproliferation agreements become easier to reach if externalities exist. First, recall that k originally represented R’s cost to proliferate. More generally, however, it could be R’s cost *plus* its externality since that value is ultimately subtracted out of its payoff. Consequently, externalities merely shift the parameters to the right on Figure 3.5. Some outcomes that previously led to proliferation under Proposition 3.4 now lead to butter-for-bombs agreements; some outcomes that previously led to butter-for-bombs agreements now lead to nonproliferation under Proposition 3.2’s “too cold” parameters.

Negative externalities for D alters the dynamics in two ways but requires an additional variable. Let $e > 0$ be the (time discounted) externality D pays if R successfully proliferates. First, incorporating the externality into inequality used to derive Proposition 3.1 shows that D finds more power shifts “too hot” to permit. Previously, D intervened if $p'_R - p_R > \frac{c_D + c_R}{\delta}$. Now if D offers R’s reservation value $p_R - c_R$, D is willing to intervene if:

$$1 - p_R - c_D > (1 - \delta) + \delta(1 - p'_R + c_R) - e$$

$$p'_R - p_R > \frac{c_D + c_R - e}{\delta}$$

Thus, D effectively calculates the externality as a *negative* cost for war. In turn, D can credibly deter R from building under a broader set of parameters when proliferation creates negative externalities.

Second, consider the parameters in which the power shift is neither “too hot” nor “too cold.” Before, D was willing to strike an agreement if $k > \delta(p'_R - c_R)$. With the externality, D bargains with R if:

$$1 - p'_R + c_R + \frac{k(1 - \delta)}{\delta} > (1 - \delta) + \delta(1 - p'_R + c_R) - e$$

$$k > \delta \left(p'_R - c_R - \frac{e}{1 - \delta} \right)$$

This inequality is easier to fulfill *ceteris paribus*. Intuitively, D finds taking the entire pie upfront less desirable if doing so triggers proliferation

and the negative externalities. As a result, nonproliferation agreements are easier to reach in a world with negative externalities.

From R's welfare perspective, negative externalities have a mixed effect. When the externality shifts a butter-for-bombs outcome to a "too hot" outcome, R sees its payoff drop since it can no longer credibly leverage the threat of investment to draw concessions. On the other hand, when the externality shifts a proliferation outcome from Proposition 3.4 to a butter-for-bombs agreement, R benefits; D's sudden desire to exterminate nuclear weapons leads to concessions that were not forthcoming previously.

Nondeterministic Proliferation. Successful proliferation is the result of many different factors. As a state begins production, scientists may be unsure whether the investment will produce any results at all. Consequently, one potential alteration to the model is to make proliferation nondeterministic.²⁵ That is, if R invests and D does not prevent, the game shifts to the post-shift state with some probability and returns to the pre-shift state (with the investment wasted) with complementary probability.

The results incentivize nonproliferation agreements. This is largely because probabilistic proliferation effectively inflates k due to uncertainty undermining the attractiveness of the investment. As such, the "too cold" parameters of Proposition 3.2 expand. Meanwhile, the size of D's optimal butter-for-bombs bribe shrinks. D now finds engaging in a deal relatively more attractive than before, thus pushing the parties out of the inefficient 3.4 outcome and into the safety of Proposition 3.3's butter-for-bombs agreements.

Nondeterministic proliferation also has a surprising effect on the "too hot" parameters from Proposition 3.1. Preventive war is less attractive when R's investment might fail. Intuitively, this is because war unnecessarily wastes costs whatever percentage of the time proliferation would have failed. In turn, D is willing to allow power shifts to transpire in more cases than before.

While this would lead to additional proliferation in models without bargaining, the butter-for-bombs logic triumphs once again. With investment still costly (especially since the program might fail entirely), R is willing to accept inducements and not proliferate. So the states agree on a butter-for-bombs deal. But note that R fares better under these conditions even though

²⁵Beliga and Sjöström (2008) model nuclear production in this manner without bargaining or the shadow of preventive war.

the proliferation process is indeterminant. Before the certainty of the shift made preventive war credible and allowed D to stop proliferation via preventive war. Now, because the uncertainty negates the credible preventive war threat, D must offer concessions to stop R. As such, R extracts some of the surplus despite its weaker outside option.

Sanctions. In practice, rivals use trade sanctions to punish states in the process of proliferating (Reynolds and Wan 2012). The baseline model, in contrast, only allows D to use the threat of preventive war to deter R from building.

Appropriately analyzing how sanctions would affect the results first requires understanding what purpose the sanctions serve. The most obvious possibility is that trade sanctions are a costly signal to credibly communicate information (Morrow 1999, 487). While this is not in doubt, such a mechanism inappropriate for this baseline model. Indeed, this chapter is largely an existence proof that demonstrates that credible, mutually preferable alternatives to proliferation exist. To reach that conclusion, the baseline model contains complete information. But that means there is no signal to send. In that light, sanctions as a signaling mechanism fails to overturn the central butter-for-bombs result. After all, sanctions decrease trade efficiency and can be extremely costly for their imposers (Martin 1992). Standard bargaining theory then indicates that some agreement Pareto dominates imposition of sanctions due to deadweight loss (Drezner 2003, 644).

Nevertheless, three complete information mechanisms are worth addressing. First, by crippling the target's economy, the rival may use sanction to remove a regime with unfriendly preferences and hope that a friendlier regime replaces it.²⁶ Sanctions thus serve as a "light" form of preventive war. If the probability of success is sufficiently high, D could credibly threaten sanctions. This could deter R from building in the same manner as Proposition 3.1's logic. If not, R would probabilistically obtain nuclear weapons (since sometimes domestic opposition would oust the current regime). But this means that deadweight loss enters the system both through the cost of weapons and the loss of trade. In turn, D could offer R a deal to match its payoff for attempting proliferation and extract the surplus through a butter-for-bombs deal.

²⁶Or, at minimum, sanctions might sufficiently destabilize the regime so that it concedes the issue out of fear of a domestic uprising (Drezner 1999, 15).

Second, sanctions could shock R's budget constraint. States must allocate resources between domestic production and international coercion. If domestic spending is inelastic, imposing sanctions and shrinking R's overall budget could leave it with insufficient funds to develop a nuclear weapon. This mechanism is comparable to Proposition 3.2's "too cold" mechanism, in which the investment cost k was too large relative to the extent of the power shift. Sanctions merely add to k 's value here.

Lastly, even if sanctions fail to inflate k sufficiently, putting pressure on R's budget constraint could reduce funding to the proliferation program and thereby increase the time it takes to develop a bomb. Proliferation ought not result here—offering R's value for investing and reaching a butter-for-bombs settlement is Pareto dominant. Moreover, sanctions increase the likelihood that the parameters fall in Proposition 3.3's butter-for-bombs outcome rather than Proposition 3.4's proliferation outcome. Recall that Proposition 3.3 prevails if $k > \delta(p'_R - c_R)$. Since sanctions delay successful development and lower values of δ measure longer times to proliferation, the right side of the inequality decreases. This makes it easier for k to be sufficiently large and for butter-for-bombs agreements to win out.²⁷

Bargaining over Objects that Influence Future Bargaining Power.

In the model, the status quo division of the bargaining good does not effect military power in future periods. This caps the minimally acceptable butter-for-bombs at $p'_R - c_R - \frac{k(1-\delta)}{\delta}$, which in turn assures that R will not demand more than that in the period after the first butter-for-bombs agreement. However, if the division of the good affects power, one concern might be that D will refuse to offer concessions upfront—D might be afraid a butter-for-bombs deal will lead to a never-ending stream of increasingly large concessions.

Fortunately, this concern is not an issue. Fearon (1996) provides the intuition. He considers a model in which control over the bargaining good also determines military strength. In equilibrium, the states mitigate the shift in power by moderating each period's transition; the side receiving additional concessions accepts smaller offers at first, knowing the extra control of the good will allow it to coerce yet more concessions out of its rival in the future. As long as the good is infinitely divisible, war never occurs.

²⁷This assumes that D's cost of implementing sanctions is worth the delay. If not, D is no worse than in the baseline model in which sanctions did not exist.

The analogous result would apply in the context of butter-for-bombs. D can offer R a smaller bribe upfront. R is willing to accept because it knows it can extort more later on, while D is satisfied because it receives a larger share earlier. The result remains efficient: R never proliferates, and both sides share the surplus.²⁸

Non-Common Discount Factors. The model gives a common discount factor δ to both parties. Some states may be more patient than others (Horowitz, McDermott, and Stam 2005; Haggard and Kaufman 1995), however, leading to questions whether different valuations of the future might prevent agreement in the present.

Proposition 3.1 and 3.2 have clean translations. D’s discount factor alone determines the cutpoint for the “too hot” range; R’s patience does not affect whether D has a credible threat to intervene. Thus, as D’s patience increases, the the size of Proposition 3.1’s parameter space increases. Similarly, R’s discount factor alone determines the cut point for the “too cold” range; D’s patience does not affect whether R has a credible threat to intervene through a sufficiently attractive investment. As such, as R’s patience increases, the size of Proposition 3.2’s parameter space decreases, since R is willing to take on the investment under a wider range of circumstances.

In contrast, the cutpoint separating Propositions 3.3 and 3.4 depend on both side’s discount factors. This is because the minimally acceptable butter-for-bombs bribe depends on R’s discount factor, but D’s preference between making that optimal bribe versus taking everything upfront depends on its own discount factor. For those reasons, the butter-for-bombs parameter increases relative to the proliferation parameter space if D’s patience increases and R’s decreases. In any case, relaxing the common discount factor assumption does not lead to proliferation for any reason not already covered in the basic model.

²⁸Chapter 4 provides empirical support. Egypt ended its nuclear weapons program after Israel returned the Sinai peninsula to Cairo. Control over the Sinai gave Egypt a tactical advantage it did not have previously. More generally, almost all empirical examples have some flavor of bargaining over objects that influence bargaining power, as declining states make cash payments to appease rising states. The influx of money could potentially be used to construct greater conventional forces in future periods. Chapter 5 also discusses how concessions in the form of nuclear assistance counter-intuitively constrain future nuclear choices (Hymans 2012, 158-202).

Imperfect Monitoring. Debs and Monteiro (2013) focus on monitoring problems in the shadow of proliferation and preventive war. Peace can fail, as D sometimes launches preventive war to deter costly arms construction even if D has no direct evidence that R chose to proliferate. However, their model prohibits the states from bargaining over the weapons to potentially resolve the issue efficiently. Thus, one concern may be whether butter-for-bombs agreements are sustainable without perfect information.

Chapter 7 tackles this extension directly. However, for now, the simple answer is that the butter-for-bombs agreements presented here are resistant to imperfect monitoring. Why? D is not willing to offer concessions if R builds after receiving them. But R chooses not to build here because doing so is simply not profitable. Note that this has nothing to do with D's informational awareness. Thus, having imperfect information does not create an impediment to butter-for-bombs agreements.²⁹

3.3 Interpretation

The butter-for-bombs equilibrium highlights the importance of *potential* power in regard to the stability of settlements. A sizeable literature in international relations debates whether systems with states of relatively equal power are more stable than systems where one state has a preponderance of power.³⁰ The rationalist literature critiques these theories by noting that the difference between relative power and relative benefits underlies incentives for war (Powell 1996; Reed et. al. 2008). In that regard, in the static bargaining model illustrated in Figure 3.2, any settlement on the interval $[p_R - c_R, p_R + c_D]$ is satisfactory to both parties.

Incorporating the possibility of shifting power changes the bargaining range, however. D still must receive no less than its reservation value for pre-shift war, or $1 - p_R - c_D$. Previously, war was R's only outside option, which paid $p_R - c_R$. Here, building is a better outside option. Thus, R must now receive at least $p'_R - c_R - \frac{k(1-\delta)}{\delta}$ to not want to alter the status quo.

²⁹In fact, Chapter 8 shows that butter-for-bombs deals expand to other parameter spaces, as D must pay a premium in the parameter space of Proposition 3.1—with imperfect information, D cannot leverage the stick of preventive war to coerce R not to build and consequently must provide concessions instead.

³⁰Although the literature goes far beyond these two works, see Morgenthau (1960) for the balance of power argument and Blainey (1988) for the preponderance of power argument.

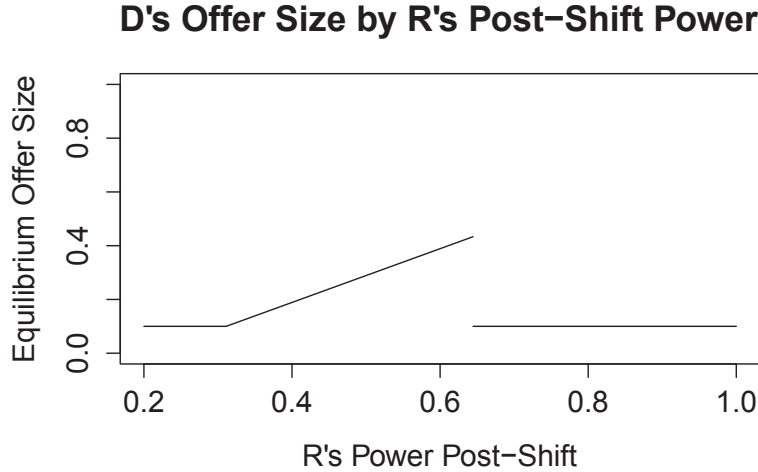


Figure 3.3: R's equilibrium offer size as a function of p'_R , with the same parameters as Figure 3.2. When the shift is too small or too large, the rising state cannot credibly threaten to build and thus receives no concessions. In the middle range, the rising state's potential power coerces concessions, and its payoff is increasing in the extent of the potential shift.

Therefore, the range of stable settlements in which the states do not fight wars and power does not shift is the set $[p'_R - c_R - \frac{k(1-\delta)}{\delta}, 1 - p_R - c_D]$.

Figure 3.3 illustrates D's equilibrium offers in the pre-shift state as a function of p'_R , with the parameters held fixed as in the earlier numerical example. When $p'_R \in (\frac{1}{5}, \frac{14}{45})$, R cannot successfully recoup its building cost. D therefore treats the bargaining problem as though power were static and offers R its pre-shift reservation value for war, which R accepts. When $p'_R \in (\frac{14}{45}, \frac{29}{45})$, R can credibly threaten to build. D utilizes the butter-for-bombs bargaining tactic, which induces R to accept the immediate concessions and not build. Finally, when $p'_R \in (\frac{29}{45}, 1)$, R cannot credibly threaten to build if it receives its pre-shift reservation value for war, as D responds to building with preventive war. Consequently, D stands firm and still induces R to accept.

Further, Figure 3.3 illustrates R's non-monotonic preferences over future power. If the power shift is very small, the ability to build does not affect R's payoff at all. In the middle range, R can successfully threaten to shift

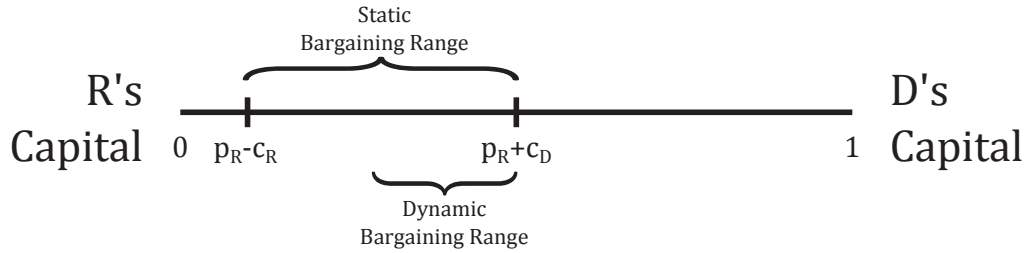


Figure 3.4: The set of Pareto settlements in a static bargaining game versus the dynamic bargaining game presented here.

power, which in turn causes D to make concessions. Moreover, these concessions are increasing in the extent of the power shift. However, the power shift eventually becomes too great, and R cannot successfully build without inducing D to intervene. Thus, R’s payoff drops precipitously, as though R does not have the ability to shift power.

Finally, Figure 3.4 shows the set of stable outcomes for situations of static and dynamic power. If the rising state cannot build additional weapons, then any settlement on the interval $[p_R - c_R, p_R + c_D]$ Pareto dominates war. If the rising state has access to weapons, then the range of settlements that Pareto dominate power shifts and war is $[p'_R - c_R - \frac{k(1-\delta)}{\delta}, p_R + c_D]$. Note that this is a subset of the Pareto dominant set in the static world.

This causes problems for an outside observer trying to understand which game the states are playing. If the observer recorded an outcome on the interval $[p_R - c_R, p'_R - c_R - \frac{k(1-\delta)}{\delta}]$, then she knows the states are in a static environment, otherwise the rising state could increase its payoff by shifting power. However, if the outcome is on the interval $[p'_R - c_R - \frac{k(1-\delta)}{\delta}, p_R + c_D]$, the observer cannot differentiate between a world in which the rising state cannot shift power and a world in which the rising state simply does not want to. Time will not resolve this problem either; since the rising state’s potential power is sufficient to extract the concessions, it never builds the weapons and never demonstrates the dynamic nature of power in the interaction.

3.4 Implications of Butter-for-Bombs Agreements

During his time in office, U.S. President John F. Kennedy feared a world of perhaps twenty-five nuclear states (Reiss 2004, 4). And by 1964, five states (the United States, the Soviet Union, the United Kingdom, France, and China) held nuclear arsenals, perhaps signalling the dawn of a global nuclear age. Yet, since the Nuclear Non-Proliferation Treaty's (NPT) creation in 1968, 190 countries have signed the treaty, and only North Korea has ever withdrawn. Meanwhile, Israel, South Africa, India, and Pakistan are the only other countries to have tested a nuclear bomb.³¹ So, at least thus far, the world has not reached the nuclear tipping point that Kennedy feared.

Yet functional nuclear weapons provide inherent security and allow states to coerce additional concessions out of their rivals during times of crisis (Beardsley and Asal 2009; Kroenig 2013). In light of this, why haven't more states followed in North Korea's footsteps by withdrawing from the NPT and joining the nuclear club?

The model provides a causal explanation: nuclear weapons are simply not in high demand in the context of a bargaining game. Bargaining is constant-sum; if nuclear weapons provide indirect benefits to their possessors, then they must also indirectly hurt their possessors' antagonists. Consequently, those fearing proliferation have incentive to offer attractive deals to shut down the nuclear contagion. Meanwhile, the potential proliferator has incentive to listen. After all, nuclear weapons are far from free. Those states would happily accept most of what they hope to gain from proliferating without investing in an actual nuclear test.

Figure 3.5 shows why demand is so low. When nuclear weapons cause too great of a power shift relative to the declining state's costs of intervention, the rising state declines to proliferate so as to avoid preventive war. Here, the declining state need not offer any carrots to induce compliance, as its stick is a sufficient threat to deter the rising state. Moreover, deterrence gives the declining state its best possible outcome, as it does not have to resort to costly war. In other words, declining states need not use carrots when sticks are credible. Consequently, a state seeking a nuclear arsenal must first shore up its conventional deterrent, otherwise proliferation is not a strategically

³¹Of these, South Africa dismantled its weapons at the end of Apartheid.

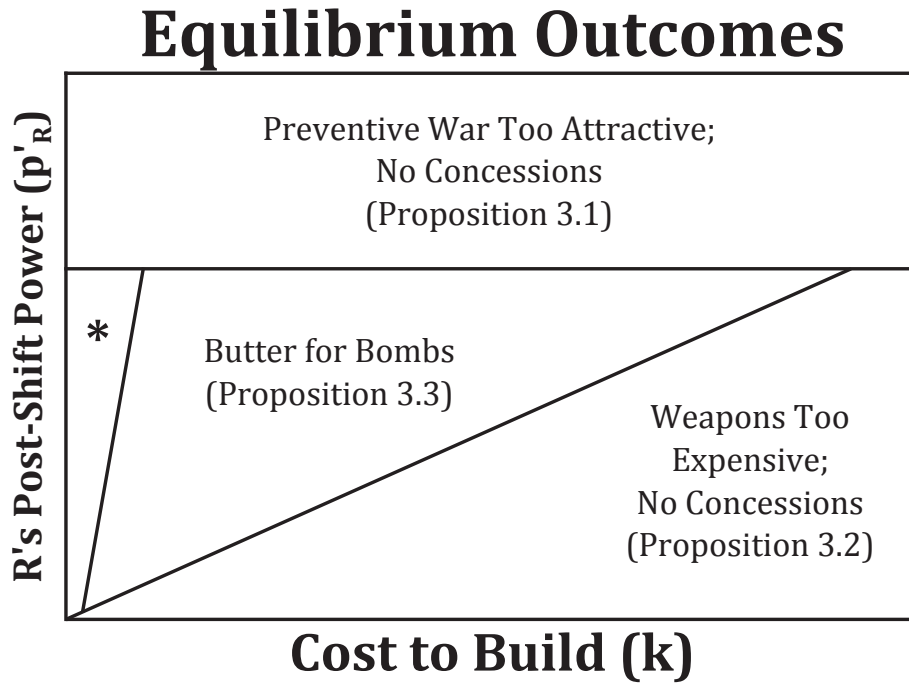


Figure 3.5: Equilibrium outcomes as a function of the cost to build and the rising state's level of future power. Investment only occurs in the region containing the asterisk, as described in Proposition 3.4.

viable option.³²

But even if the potential proliferator can defend itself from an invasion, it still might not want to seek nuclear weapons. After all, bombs are an investment in the future. Such an investment is only sensible if it yields sufficient returns. Thus, states will not proliferate if the financial cost of nuclear weapons is too great. Moldova or Rwanda might view nuclear weapons as attractive in theory. However, the cost to proliferate would bankrupt those countries before they could achieve nuclear capacity. Similarly, states need to have a contentious security issue for proliferation to make sense. Iceland and Ireland have the technical know-how and financial resources to build a bomb, but it is unclear what sort of benefits said bomb could bestow.

Proliferation remains unrealized even as the attractiveness of the invest-

³²Iran has correspondingly placed many of its nuclear facilities underground. This location limits the damage from a possible aerial strike, reducing the Israeli or American ability to effectively intervene.

ment increases. At this point, the potential rising state is conditionally willing to shift power. But it is in the declining state's best interest to bribe the would-be nuclear state and avoid facing the consequences of a much stronger rival. The states ultimately resolve the crisis without proliferation, as the immediate concessions ensure that building a bomb will not lead to a better outcome.

The model also reveals that bargaining over nuclear weapons does not require rising state to commit to the incredible. In negotiating over nuclear weapons programs, many commentators (Bolton 2010, Krauthammer 2009, and Fly and Kristol 2010) have warned that potential proliferators cannot be bought off. That being the case, declining states should hold their ground, as bribes have no effect on tomorrow's power politics.

Yet the model shows that such a strategy creates a self-fulfilling prophecy. Standing firm in the present causes rising states to redouble their efforts. Forestalling negotiations therefore creates the exact nuclear problem declining states wish to avoid. Resolving conflict requires parties to have the correct incentives. Here, with no other bargaining frictions present, it is remarkably easy to obtain a rising state's compliance.

As a final note, these findings instruct us to take a holistic approach to understanding nuclear proliferation. Quantitative studies frequently attempt to understand proliferation behavior by analyzing "supply side" components of nuclear weapons (Meyer 1985; Jo and Gartzke 2007); states with limited nuclear capacities are unlikely to develop nuclear bombs. While the model (via Proposition 3.2) confirms the value of supply side explanations, nuclear capacity is not the sole explanatory variable. Indeed, Figure 3.5 shows that supply side arguments explain outcomes in the bottom right portion of the parameter space only; preventive war and bargaining determine the remaining outcomes.

Ignoring these other factors leads to strange interpretations of the data. Sagan (2011, 229-230) notes that, according to the Jo and Gartzke (2007) estimates, Trinidad and Tobago "had a higher degree of nuclear weapons latency in 2001 than is North Korea, which was only five years away from detonating its first nuclear weapon." But latent capacity does not become active capacity without the will of the state. Trinidad and Tobago has no significant coercive bargaining relationship and maintains an active military force in the thousands. Meanwhile, North Korea has technically been at war since the 1950s and has more than a million active duty soldiers. Thus, latency measures require context. Bargaining relationships explain why states

ultimately choose to develop nuclear capacity.

3.5 Conclusion

This chapter formally investigated the credibility of butter-for-bombs settlements. Although international relations scholars traditionally emphasize how fully realized power extracts concessions, the model demonstrated that *potential* power is sufficient. Declining states have incentive to proactively bargain with rising states, so as to ensure that non-proliferation remains the status quo. Rising states have incentive to welcome the offers, as they can obtain most of their goals without paying costs to develop a weapons program. Credible non-proliferation agreements result.

The model makes a significant contribution to the understanding of costly weapons production. At present, explanations for non-armament are limited to the threat of preventive war and inefficient investments; current models do not explain how carrots convince states to forgo weapons programs. The butter-for-bombs model fills the gap, showing how states can manipulate their rivals' opportunity cost and thereby avoid nuclear proliferation.

While the model reveals the absence of commitment problems and the existence of bargaining space, it fails to provide any intuition as to exactly how states reach butter-for-bombs deals. Consequently, the next chapter provides case studies to corroborate the usefulness of the model. Later chapters then add bargaining frictions—shifting resolve, imperfect information, and incomplete information—to study whether butter-for-bombs agreements remain credible in these contexts.

3.6 Appendix

This appendix covers the proofs of the lemmas and propositions from this chapter. For parsimony, it shows that each equilibrium is unique assuming that R accepts an offer when indifferent between taking an efficient and inefficient action. However, all equilibria are unique with that assumption relaxed for the standard reasons that the equilibrium of an ultimatum game is unique.

3.6.1 Proof of Lemma 3.1

First, in every equilibrium for every history of the game, R's continuation value is at least $p_R - c_R$. This is because R can reject in any period and secure that amount.

Second, R must accept $y_t > p'_R - c_R$ in every equilibrium for every history of the game. Recall R earns $p'_R - c_R$ if it rejects in any period. In contrast, if R receives an offer of $y_t > p'_R - c_R$, accepting generates a payoff of $(1 - \delta)y_t + \delta V_R$, where V_R is R's continuation value. The previous paragraph ensures that $V_R \geq p'_R - c_R$. Using $V_R = p'_R - c_R$ as a lower bound, accepting is strictly better than rejecting if:

$$(1 - \delta)y_t + \delta(p'_R - c_R) > p'_R - c_R$$

$$y_t > p'_R - c_R$$

This holds. So R must accept $y_t > p'_R - c_R$.

Third, D never offers $y_t > p'_R - c_R$ in every equilibrium for every history of the game. Using the one-shot deviation principle, D could instead offer the midpoint between that y_t and $p'_R - c_R$. This amount is still strictly greater than $p'_R - c_R$, so R still accepts. In turn, D receives strictly more for the period and the identical amount in all future periods, so this is a profitable deviation.

Fourth, R rejects $y_t < p'_R - c_R$ in every equilibrium for every history of the game. The first and third step of this proof imply that R's continuation value is no greater than $p'_R - c_R$. Thus, if R accepts $y_t < p'_R - c_R$, it receives strictly less than $p'_R - c_R$ for the game. Consequently, rejecting and earning $p'_R - c_R$ is a profitable deviation.

Fifth, since R's continuation value must be no greater than $p'_R - c_R$ and no less than $p'_R - c_R$, it must be exactly equal to $p'_R - c_R$. Given the previous results on R's accept or reject decision, the only way this is possible is if D offers $y_t = p'_R - c_R$ in every period. R cannot profitably deviate since it receives $p'_R - c_R$ by fighting in any period, which is identical to its payoff for accepting. D cannot profitably deviate because demanding more results in rejection (paying $1 - p'_R - c_D$) while demanding less is a needless concession. \square

3.6.2 Proof of Proposition 3.1

First, in every equilibrium for every history of the game, R's continuation value V_R for any pre-shift period must be at least $p_R - c_R$. The proof is identical to the analogous claim in the proof for Lemma 3.1, swapping y_t for x_t and p_R for p'_R .

Second R must accept $x_t > p_R - c_R$ in every equilibrium for every history of the game. R cannot reject in such circumstances due to the analogous proof in Lemma 3.1. R's only other alternative is to build. However, D prevents if:

$$1 - p_R - c_D > (1 - x_t)(1 - \delta) + \delta(1 - p'_R + c_R)$$

Note that because $x_t \geq p_R - c_R$ in this case, $(1 - p_R + c_R)(1 - \delta) + \delta(1 - p'_R + c_R) \geq (1 - x_t)(1 - \delta) + \delta(1 - p'_R + c_R)$. Therefore, to show that preventing is optimal for D, consider instead the following inequality:

$$1 - p_R - c_D > (1 - p_R + c_R)(1 - \delta) + \delta(1 - p'_R + c_R)$$

$$p'_R - p_R > \frac{c_D + c_R}{\delta}$$

Because R earns $p_R - c_R - (1 - \delta)k$ if D prevents, R must accept $x_t > p_R - c_R$.

Third, in every equilibrium for every history of the game, D never offers $x_t > p_R - c_R$. The proof is identical to the analogous claim in the proof for Lemma 3.1.

Fourth, D never offers $x_t < p_R - c_R$ in every equilibrium for every history of the game. If it did, one of three things could happen in response. First, R could reject. D earns $1 - p_R - c_D$ for this outcome. D could make a one-shot profitable deviation to $x_t = p_R$ in period t . Per above, R accepts. D receives $1 - p_R$ for the period and must earn at least $1 - p_R$ for the rest of time for the same reason, which is greater than $1 - p_R - c_D$. Alternatively, R could build. This is only optimal for R if D does not prevent, as R earns $p_R - c_R - (1 - \delta)k$ in that case, which is less than what it earns for rejecting. So R must earn at least $p_R - c_R$ for this outcome. In turn, D earns *no more* than $1 - p_R + c_R - (1 - \delta)k$, after factoring out R's cost to build. But D could make a one-shot profitable deviation to $p_R - c_R + \frac{(1 - \delta)k}{2}$. Per above, R must accept. This gives D the remainder for the period, and D must receive at least that much in every equilibrium in remaining periods. This generates a

greater payoff for D than offering an amount less than $p_R - c_R$ and inducing R to build. Third, R could accept. But since the rest of this paragraph and the first and third claims ensure that R's continuation value must be less than or equal to $p_R - c_R$, R could profitably deviate to rejecting. In turn, D would have a profitable deviation to offering $x_t = p_R$ for the same reasons as before.

Fifth, since $V_R \leq p_R - c_R$ and $V_R \geq p_R - c_R$, V_R must be exactly equal to $p_R - c_R$ in every equilibrium for every history of the game. Given the above equilibrium constraints, the only way this can happen is if D offers $x_t = p_R - c_R$ in every period and R accepts. D has no profitable deviation since offering more is a needless concession while offering less results in war or a power shift that forces D to give up even more concessions. \square

3.6.3 Proof of Proposition 3.2

First, in every equilibrium for every history of the game, R's continuation value V_R must be greater than $p_R - c_R$ for all pre-shift periods. The proof is the same as the first part of the proof for Proposition 3.1.

Second, R must accept $x_t > p_R - c_R$ in every equilibrium for every history of the game. R has two alternatives: war and building. War generates a payoff of $p_R - c_R$ forever, while $V_R \geq p_R - c_R$ ensures that accepting $x_t > p_R - c_R$ will give a greater amount than rejecting in period t and at least as much in all future periods. Alternatively, R could build. In R's best case scenario, D does not prevent. Using Lemma 3.1, R earns $p'_R - c_R$ in all future periods. Even so, R strictly prefers accepting if:

$$(1 - \delta)x_T + \delta V_R > (1 - \delta)x_t + \delta(p_R - c_R) - (1 - \delta)k$$

Using $V_R = p_R - c_R$ as a lower bound, this holds if:

$$(1 - \delta)x_T + \delta(p_R - c_R) > (1 - \delta)x_t + \delta(p_R - c_R) - (1 - \delta)k$$

$$p'_R - p_R < \frac{(1 - \delta)k}{\delta}$$

This is the cutpoint given in Proposition 3.2.

Third, in every equilibrium for every history of the game, D never offers $x_t > p_R - c_R$. The proof is the same as the third part of the proof for Proposition 3.1.

Fourth, D never offers $x_t < p_R - c_R$ in every equilibrium for every history of the game. The second claim ensures that R will not respond by building. That, combined with the fact that the first and third claims ensure that $V_R \leq p_R - c_R$, imply that R must reject. But D could make a one-shot profitable deviation to $x_t = p_R$. R will accept. That gives D at least as much for the period and at least as much in all future periods. This is greater than earning its war payoff of $1 - p_R - c_D$.

The fifth and final step is identical to the fifth step from the proof for Proposition 3.1. \square

3.6.4 Proof of Proposition 3.3

All stationary MPE can be characterized by an equilibrium value x_t^* offered in each period. Note that if builds and D does not prevent, R earns $(1 - \delta)x_t + \delta(p'_R - c_R) - (1 - \delta)k$. Let V_R be R's continuation value for accepting. Under such conditions, R accepts if:

$$(1 - \delta)x_t + \delta V_R \geq (1 - \delta)x_t + \delta(p'_R - c_R) - (1 - \delta)k$$

$$V_R \geq p'_R - c_R - \frac{k(1 - \delta)}{\delta}$$

In particular, R's decision does *not* depend on the offer x_t ; this is a direct consequence of the bargaining environment without *quid-pro-quo* offers. What matters is the continuation value. Note that R can always accept in any equilibrium. So if $x_t^* \geq p'_R - c_R - \frac{k(1 - \delta)}{\delta}$, R's continuation value must be at least $p'_R - c_R - \frac{k(1 - \delta)}{\delta}$. So R can accept. If $x_t^* < p'_R - c_R - \frac{k(1 - \delta)}{\delta}$, R's only other recourse is to launch preventive war. But that pays $p_R - c_R$, which is less than R's value for receiving nothing in the current period and building. So R must build.

Now consider what are D's optimal offer sizes. Offering $x_t > p'_R - c_R - \frac{k(1 - \delta)}{\delta}$ cannot be optimal. R would not build under such circumstances. But D could make a one-shot deviation to offering the midpoint between that x_t and $p'_R - c_R - \frac{k(1 - \delta)}{\delta}$. The continuation value remains $p'_R - c_R - \frac{k(1 - \delta)}{\delta}$, so R still accepts. However, D receives more of the bargaining good for the period and the same amount in the future, which is a profitable deviation.

Similarly, $x_t \in (0, p'_R - c_R - \frac{k(1 - \delta)}{\delta})$ cannot be optimal either. R must build under such circumstances. But given that R is building, D could make a one-

shot deviation to $x_t = 0$. R still builds. D receives more for the period and the same amount for the remainder of time, which is a profitable deviation.

So x_t^* must equal 0 or $p'_R - c_R - \frac{k(1-\delta)}{\delta}$. In fact, both are supported in equilibrium. Suppose $x_t = p'_R - c_R - \frac{k(1-\delta)}{\delta}$ and R accepts if and only if $x_t \geq p'_R - c_R - \frac{k(1-\delta)}{\delta}$. D cannot make a one-shot profitable deviation in any period; offering more is an unnecessary concession while offering less triggers R to build, which denies D the surplus. R cannot profitably deviate its continuation value (not the offer) completely determines R's optimal strategy. Thus, this is an equilibrium. Likewise, suppose $x_t = 0$ and R always builds regardless of the offer. D cannot affect R's build decision and thus maximizes its payoff by minimizing the offers. R cannot profitably deviate again because its continuation value completely determines its best response.

The final thing to verify is that D would not prevent if R built under these circumstances. Given the parameter space, the only way this could occur is if $x_t > p'_R - c_R - \frac{k(1-\delta)}{\delta}$. But then R would not build anyway, and D could make a one-shot deviation to offering the midpoint between that x_t and $p'_R - c_R - \frac{k(1-\delta)}{\delta}$. \square