

## 1 Null

I want the numbering here to match with issues in Github, and there was one issue already so I need a placeholder.

## 2 Dixit, Grossman and Helpman analogy

### Editor Point 3:

You state that your specification of legislature preferences can be seen as a special case of the Dixit-Grossman-Helpman (DGH) model. It would be helpful if you could substantiate this claim. More specifically, let us consider a simple version of the DGH specification in your setting. Suppose you modify your model only in two ways. First, suppose the legislature preferences are given by  $W+g(e)$ , where  $g(e)$  is an increasing and concave function of contributions (while the lobby preferences remain the same). And second, suppose the lobby offers a contribution schedule  $e(t)$  before the legislature chooses the tariff. The question is: would this model deliver the same results as yours (at least qualitatively)? Note that in this setting you can focus on a simple all-or-nothing contribution schedule, of the kind "if you give me a 5% tariff I give you \$100, otherwise you get nothing." Intuitively the promised contribution will just compensate the legislator for the loss associated with the requested tariff, so the analysis might not be hard. If this DGH version of your model yields similar qualitative results, pointing out this "isomorphism" would help you in several ways. First, it would provide "foundations" for your assumed legislature preferences, in terms of a model (DGH) that people are familiar with. Second, this would help address my question 2 above: in the DGH model we can think of  $e$  as money, and I think it would be reasonable to stick with your current definition of aggregate welfare (even though in principle one might question this definition of aggregate welfare when utility is not transferrable). And third, you could examine whether your results rely on the presence of diminishing marginal utility from contributions: what would happen if  $g$  is linear (as in the basic Protection for Sale model) rather than strictly concave?

DGH97 paper:

- Truthful contribution schedule is a device for solving equilibrium in their model. It's not a best response function, but lobby has to be best responding in equilibrium.

- Proposition 3:  $G(a^0, e^T(a^0, u^0)) = \max_a G(a, 0)$
- $e^T(a^0, u^0)$  is essentially  $\phi(a^0, u^0)$ , which is defined implicitly in eqn3 (p. 760) as

$$U[a, \phi(a^0, u^0)] = u^0$$

- if  $G = W + g(e)$  and  $e = 0$  and  $g(0) = 0$ , then  $a^* = \tau^{\text{opt}}$
- Then rewrite as

$$G(\tau^0, e^T(\tau^0, u^0)) = \max_a G(\tau^{\text{opt}}, 0)$$

- \* Right hand side provides a number
- \* Ignoring arguments of  $e^T$  function, LHS traces out  $(\tau, e)$  pairs that satisfy the equation given  $g(e)$ .
- \* For lobby to be best responding, it MUST pick the pair that maximizes  $\pi(\tau) - e$ . *This* concern must be what sets  $u^0$ .
- Corollary to Prop 1 / Prop 3: Gov't gets utility equal to outside option. Is this true when there is just one lobby?
- Combining the two previous facts, gov't getting outside option will set  $u^0$  (eqm utility) and anchor contribution schedule (just have to be careful of zero contributions)

What editor proposes:

- Lobby offers contribution schedule (can be very simple: just one  $(e, \tau)$  pair,  $e = 0$  for everything else)
- Government maximizes  $W + g(e)$ 
  - Note that  $CS_X + \gamma(e) \cdot PS_X + CS_Y + PS_Y + TR = W + (\gamma(e) - 1) PS_X$
  - Also note that there may be a problem of interpretation: DGH fundamentally makes unitary actor indifferent, instead of my set-up where  $e = 0$  and  $e = \bar{e}$  lead to two different decision makers
- Sectioning from my paper:
  - 3.1 Same (execs)
  - 3.2 Trade war
    - Government evaluates welfare (unilateral because  $\tau^*$  doesn't change) at the TIOLI offer of the lobby and at the value  $(\tau^{\text{opt}})$  that satisfies  $\frac{\partial W}{\partial \tau} = 0$  [nothing the leg does here will change the contribution]. Chooses which one maximizes welfare.

- \* DGH Proposition 1: principal (lobby) has to provide at least agent's (gov'ts) outside option; as long as this constraint is satisfied, lobby can propose  $\tau$  and payment that maximizes his own utility
- \* DGH Proposition 3 (p. 760-61): simplifies so we don't need to look for contribution functions, only eqm values
- \* Assume lobby has all the bargaining power as in DGH97. Then lobby calculates  $(\tau, e)$  schedule from

$$W(\tau) + g(e) = W(\tau^{opt}) + g(0)$$

Assume  $g(0) = 0$ . Then

$$\begin{aligned} g(e) &= W(\tau^{opt}) - W(\tau) \\ e &= g^{-1} [W(\tau^{opt}) - W(\tau)] \end{aligned}$$

3.3 Nothing seems to change

3.4 Nothing seems to change

4 I think the break phase is essentially the same, but the way it is calculated is different

- Lobby first determines  $\bar{e}$ , then decides whether it's worthwhile paying  $\bar{e}$
- Need to re-write Equation 12, which defines  $\bar{e}$ 
  - \* What *is*  $\bar{e}$ ? It's what the lobby has to pay to get the legislature to break the trade agreement

$$\begin{aligned} W_{ML}(\gamma(e), \tau^a) + \frac{\delta_{ML} - \delta_{ML}^{T+1}}{1 - \delta_{ML}} W_{ML}(\gamma(e), \tau^a) \geq \\ W_{ML}(\gamma(e), \tau^R(e), \tau^{*a}) + \frac{\delta_{ML} - \delta_{ML}^{T+1}}{1 - \delta_{ML}} W_{ML}(\gamma(e), \tau^{tw}) \quad (1) \end{aligned}$$

- \* In base model, the  $\tau$  for the current period is chosen in response to current period  $e$ ; this decision is made differently (potentially) than the decision to break. Is there a similar separation between the decisions in the DGH version of the model?
- \* When lobby has all the bargaining power, there's a schedule  $(\bar{e}, \tau^b)$  that makes leg indifferent between breaking the trade agreement or abiding by it ( $\tau^b$  no longer a direct fcn of  $\bar{e}$ )
- \* Note: It's not necessarily the minimum  $e$ . Yes,  $\bar{e}$  must satisfy (something like) Equation (7), but lobby may make higher profits choosing an  $e$  higher than the

minimum one in the base model.

$$W(\boldsymbol{\tau}^a) + g(e_a) + \frac{\delta_{ML} - \delta_{ML}^{T+1}}{1 - \delta_{ML}} [W(\boldsymbol{\tau}^a) + g(e_a)] =$$

$$W(\tau^b, \tau^{*a}) + g(\bar{e}) + \frac{\delta_{ML} - \delta_{ML}^{T+1}}{1 - \delta_{ML}} [W(\boldsymbol{\tau}^{tw}) + g(\bar{e})] \quad (2)$$

The lobby's payment / requested tariff pair maximize the lobby's future income stream by paying legislature no more than necessary to break

- If legislature doesn't break, it'll choose  $(\tau^a, e_a)$ , which is set by trade agreement
- Lobby isn't going to set indifference condition to a WORSE outcome (i.e.  $(\tau^{opt}, 0)$ )

Lemma 1:

$$\frac{d\bar{e}}{d\tau^a} = -\frac{\frac{\partial \Omega}{\partial \tau^a}}{\frac{\partial \Omega}{\partial \bar{e}}} = \frac{\left[1 + \frac{\delta_{ML} - \delta_{ML}^{T+1}}{1 - \delta_{ML}}\right] \frac{\partial}{\partial \tau^a} [W(\boldsymbol{\tau}^a) + g(e_a)]}{\left[1 + \frac{\delta_{ML} - \delta_{ML}^{T+1}}{1 - \delta_{ML}}\right] \frac{\partial g(e)}{\partial \bar{e}}} \quad (3)$$

Lemma 2:

$$\frac{d\bar{e}}{d\tau^{*a}} = -\frac{\frac{\partial \Omega}{\partial \tau^{*a}}}{\frac{\partial \Omega}{\partial \bar{e}}} = \frac{\frac{\delta_{ML} - \delta_{ML}^{T+1}}{1 - \delta_{ML}} \frac{\partial}{\partial \tau^{*a}} W(\boldsymbol{\tau}^a)}{\left[1 + \frac{\delta_{ML} - \delta_{ML}^{T+1}}{1 - \delta_{ML}}\right] \frac{\partial g(e)}{\partial \bar{e}}}$$

6 Optimal dispute settlement: depends on how government IC changes

$$\left(1 - \frac{d\pi}{d\bar{e}}\right) \frac{-\frac{\delta_{ML}^{T+1} \ln \delta_{ML}}{1 - \delta_{ML}} [W(\boldsymbol{\tau}^a) + g(e_a) - W(\boldsymbol{\tau}^{tw}) - g(\bar{e})]}{\left[1 + \frac{\delta_{ML} - \delta_{ML}^{T+1}}{1 - \delta_{ML}}\right] \frac{\partial g(e)}{\partial \bar{e}}} + \frac{\delta_L^{T+1} \ln \delta_L}{1 - \delta_L} [\pi(\tau^{tw}) - e_{tw} - \pi(\tau^a) + e_a] \quad (4)$$

Numerator of  $\frac{\partial \bar{e}}{\partial T} \rightarrow 0$  and denominator  $\rightarrow \frac{\delta}{1 - \delta}$ , just as in basic formulation

- \* Need numerator of  $\frac{\partial \bar{e}}{\partial T}$  to be positive, i.e.  $W(\boldsymbol{\tau}^a) + g(e_a) > W(\boldsymbol{\tau}^{tw}) + g(\bar{e})$
- \* If not, trade war is not a punishment for the legislature
- \* Analogy to main formulation: IF there's an equilibrium with  $\tau^a < \tau^{tw}$ ,  $W_{ML}(\gamma(\bar{e}), \boldsymbol{\tau}^a) - W_{ML}(\gamma(\bar{e}), \boldsymbol{\tau}^{tw})$  must be positive.
- \* If there's no  $\bar{e}$  that satisfies inequality (10), then no  $\tau^a < \tau^{tw}$  is self-enforcing.
- \* The last paragraph on pg. 19 restricts attention to this case, **may need to clarify it**; bottom of pg 21 to middle of pg. 22 discusses existence of eqm with  $\tau^a < \tau^{tw}$  in detail

- \* For the purposes of Sections 5 and 6, we must have  $W_{\text{ML}}(\gamma(\bar{e}), \tau^a) - W_{\text{ML}}(\gamma(\bar{e}), \tau^{tw}) > 0$ . Otherwise there would be no interior solution and we wouldn't be interested in these results.

Is it harder to believe  $W(\tau^a) + g(e_a) > W(\tau^{tw}) + g(\bar{e})$ ? No, I don't think so.

- \* Here,  $\bar{e}, \tau^b$  pair is chosen to make it true: they're chosen to make Expression (2) hold with equality by prefix maximizing behavior of lobby. In order to get lobby to defect, by this equation, the above inequality must hold.
- \* To be clear, to get an eqm with  $\tau^a < \tau^{tw}$ , there needs to be a  $\tau^a < \tau^{tw}$  where the lobby loses money at every  $(\tau^b, \bar{e})$  pair that is determined by Expression (2). We simplify this problem by finding the best pair for the lobby and making sure that one is not profitable.
- \* There is no 'extra' relationship with the determination of  $\tau^b$  that could get in the way as in the base formulation
- \* Bottom line: no harder to believe this inequality holds than the DGH version; perhaps even a bit easier

## Comparisons

- Note that in GH94, claim is that IN EQUILIBRIUM, government behaves as if it maximizes a weighted sum of the groups' utilities ( $1 + a$  for those represented by lobbies;  $a$  for those not represented; pg. 10 of PDF). But this is in equilibrium: there's a lot going on out of equilibrium that leads us there!
  - I haven't proved this to myself, but that must also depend on lobbies having all the bargaining power
- In DGH,  $G(\tau^0, e^0) = \max_{\tau}(\tau, 0)$ 
  - This is not generally true in my model
  - I think the slippage is in bargaining power. DGH paradigm assumes lobby essentially make TIOLI offer.
  - $\gamma(e)$  formulation in essence distributes bargaining power more generally
    - \* Is it okay to characterize it this way?
    - \* If so, does bargaining power vary with effort? Seems mixed up with diminishing returns to effort.
  - $W + \Phi(e)$  of Limao and Tovar gives decreasing returns; then explicitly models bargaining power through Nash bargain instead of menu auction (as far as I can tell—I can't find it clearly specified).

- In DGH, essentially lobby chooses its favorite  $(\tau, e)$  pair from among all the possible ones that make the Gov't indifferent.
  - Is it possible that this is equivalent to  $\frac{\partial \pi}{\partial e} = 1$ ?
  - seems like the pair that satisfies that equation might not be available
  - When I calculate the various  $(\tau, e)$  pairs to 'compare' government welfare levels, this comes from the government welfare function: each  $e$  induces an optimal  $\tau$  and then I get an optimal value function
    - \* Lobby doesn't care about these varying welfare levels. Chooses  $(\tau, e)$  pair that maximizes  $\pi - e$

In my model, I argue that the break tariff  $(\tau^b(\bar{e}))$  should be chosen according to the lobbying effort applied in the break stage  $(\bar{e})$  because this determines the identity of the median legislator. The median legislator in the current period gets to pick the tariff in the current period.

- In DGH, the lobby presents a unitary gov't with a schedule of  $(\tau, e)$  pairs.
  - The editor has suggested that I can choose a very simple take-it-or-leave-it schedule with one just pair or  $e = 0$ , but in order for my calculus-based results to hold up, I think I need a differentiable schedule.
  - At any rate, no one in the government gets to 'pick'  $\tau$  on its own; a particular  $\tau$  is demanded in return for a particular  $e$ . They are a bundle constructed by the lobby.
  - In my model, who the decision-maker IS is determined by the lobby's chosen  $e$ .
- Is it possible to have a DGH-style indifference condition in the spirit of my model, i.e.  $e$  determines the median legislator?
  - In a normal repeated game, players choose current period action to maximize the whole stream of payoffs; that's not what happens in my model. The current player makes his decision with a view to influencing who will be the median legislator in the future.
  - Current player finds  $\tau$  that maximizes current payoffs, then checks to see whether it will also maximize future payoffs. If it doesn't, he switches to trade agreement tariff.
    - \* That is, the maximization program for current player is complex: has to find optimal  $\tau$  if he's going to break, then compare total payoff stream under  $\tau^b(\bar{e})$  to total payoff under  $\tau^a$ .

Examples

- Use numeric example from BS2005 in R (DGH.r)
- Isomorphism I've already calculated: if  $G = W + e$  then  $\gamma(e) = 1 + \frac{e}{\pi_x}$
- Lobby's optimization:

$$W(\tau, e(\tau)) = W(0, 0)$$

$$W(\tau) + e = W(0)$$

$$e = W(0) - W(\tau)$$

- This forms a contribution schedule
- Lobby chooses the pair  $(\tau, e)$  that maximizes  $\pi(\tau) - e$
- Make sure these are unilateral changes in  $\tau$  in the program
- Then compare to my old way:  $\frac{\partial \pi}{\partial \tau} \frac{\partial \tau}{\partial \gamma} \frac{\partial \gamma}{\partial e} = 1$  ( $e$  changes with  $\tau$  directly depending on shape of  $\gamma$ , then government maximizes w.r.t.  $\tau$ )
  - NEXT: go into R program and sort it out

## 2.1 Subtract effort in gov't welfare function

### Editor Point 2:

There is a question in my mind about the definition of welfare and the objective functions of the executive and the legislature. As I understand it, the lobbying effort involves a resource cost (this is not a cash transfer), and this cost in principle should be reflected in the expression for welfare. Viewed from a different perspective: you assume that the executive and the legislature care about import-competing profits, and thus about the lobby in a broad sense, but don't care about the resource cost of effort that the lobby incurs, and this seems hard to justify. I am sorry to be raising this issue at this stage and not in the previous round, but I became aware of it recently.

### 3 Interpretation: SOP vs. time inconsistency

#### Referee 1 Point 1:

I believe the interpretation of the model as representing two branches of the government is arbitrary. It is more realistic to interpret the model as representing the preferences of a government with time-inconsistent preferences, such that ex ante the government's objective is to maximize social welfare, but at the time of implementing the agreement it might be influenced by lobby group activities.

I understand that the author has tried to modify the interpretation of the model to address this concern, but in my view the modification was insufficient.

#### Editor Point 5:

Regarding Referee 1's idea about time-inconsistent preferences, I guess I am ok with your separation-of-powers interpretation (subject to the caveats I expressed in my first-round letter), but I also think it would be useful to discuss a possible alternative interpretation along the lines of Referee 1's idea. I am thinking of the Maggi and Rodriguez-Clare setting, where there is a unitary government that cares about welfare and contributions, but at the stage of signing the agreement ("ex ante" stage) the lobby has no influence on the government, because the lobby cares only about the short run (due to the fact that specific capital is mobile in the long run). The way I think about this interpretation is slightly different from Referee 1, in that there is no time inconsistency of preferences: the government has the same preferences across time, but ex-ante the lobby is not active. I am not sure whether a Maggi and Rodriguez-Clare specification of this kind would yield similar results as yours. This might be an interesting question to discuss, and if the answer is yes, it would be useful to point out that the model admits also this alternative interpretation.

- Note this also connects to analysis section
- To generalize Giovanni's 'government has the same preferences across time, but ex-ante the lobby is not active' idea, ex-post, the lobby can be active at different levels and this does not change the 'preferences,' which are really the process by which the median legislator is chosen.
  - Lobby effort is 0 at ex-ante stage because capital is perfectly mobile in the long-run so it's not worthwhile to expend resources to influence the future. In essence, the lobby only cares about the short run.



- **This interpretation works as long as the non-separation-of-powers government is NOT unitary.**
- The lobby doesn't change the preferences of any legislator; it changes who gets to make the decision
  - Perhaps I need to explain better how the legislature works in the paper somewhere
  - **Start by identifying where I talk about this, where I can clarify, if there's a better place to do it**
- Show all the alternatives, say which are special cases, which are richer (hopefully mine!)
  1. MRC-style 'tribute' from two paragraphs above, two versions
    - (a) Non-unitary government, matches mine perfectly
    - (b) Unitary government, which version proposed by referee 1 point 3. Essentially a unitary legislature, or even unitary government w/time inconsistent preferences. Here, main construction stays, but interior optimum in  $T$  goes away
  2. DGH-style, which replicates my results qualitatively as long as we go with the non-unitary interpretation (appendix results)

## 4 Continuation payoffs and changes in the preferences

### Editor Comment 1:

I agree with Referee 2 that there are still problems with the analysis. Correcting these problems is a necessary (though not sufficient) condition for me to move forward with this paper. You will need to convince us beyond the reasonable doubt that the analysis is correct.

### Editor Point 8:

The way you write the key program in (9)-(11) is confusing because you have an “e” floating around, and it is not clear where it should be evaluated. Unless I am missing something, in some places this should be  $\bar{e}(\tau^a)$  and in other places it should be  $e^a$ . It would also be helpful to write  $e^a$  as a function of  $\tau^a$ . Since the choice variable in the program is  $\tau^a$ , you should make clear what is a function of  $\tau^a$  and what is not.

### Referee 1 Point 3:

It is assumed that the current legislator evaluates future welfare (i.e., the continuation payoffs) based on the expected preferences of the future government, which will be induced by lobbying efforts in the future. Alternatively, it could be assumed that the current legislator evaluates future profits based on its current preferences.

The latter assumption might be more consistent with the premise of the model, which is essentially a decision-making model with time-inconsistent preferences. Moreover, I think the results related to self-enforceability of the agreement will continue to hold if the author adopts the latter assumption.

Some discussion of this point could be illuminating.

### Referee 2 Point 3:

I am confused about one basic aspect of the incentive constraint for the legislature. Consider the RHS of (7) on page 16. If the lobby chooses  $e$  and the legislature decides to select its best response given  $e$ ,  $\tau^R(e)$ , and thereby breaks the agreement, then a T-period punishment is launched in the next period. The notation in (7) (and likewise in (10)) suggests that the lobby continues choosing the same  $e$ , thus generating the same  $\gamma(e)$ , during the punishment phase. I don't see why this would be the case. Wouldn't the lobby instead choose the effort level  $e_{tw}$  that is determined in the first-order condition given by (6)? And indeed, if we look at the lobby incentive constraint, as given by the RHS of (8) (and likewise in (11)), we see that that constraint does assume that  $e_{tw}$  is used by the lobby during the punishment phase. I can't tell exactly what is going on here. There could be an oversight, or I could be misunderstanding the notation. At a minimum, some clarification is needed.

This consideration also leads to a further concern/question. If my point above is correct and the lobby should be modeled in (7) as choosing  $e_{tw}$  in the punishment phase, and if  $e_{tw} > e$ , then would the legislature ever deviate (even for  $e$  in the non-triggering range as currently defined) in order to trigger a trade war and thereby enjoy the higher  $e_{tw}$  and thus the higher gamma that the trade war elicits? Recall that  $W_{ML}$  is increasing in  $e$  as an independent argument. Is this potential incentive captured?

**Referee 2 Point v:**

Page 20: Related to comment 3 above, I don't follow why in (12) that  $e$  can't change from  $e$ -bar as we move into the trade war.

My response:

Under the alternative suggested by Referee 2, Equation (7) would be modified from

$$W_{ML}(\gamma(e), \tau^a) + \frac{\delta_{ML} - \delta_{ML}^{T+1}}{1 - \delta_{ML}} W_{ML}(\gamma(e), \tau^a) \geq W_{ML}(\gamma(e), \tau^R(e), \tau^{*a}) + \frac{\delta_{ML} - \delta_{ML}^{T+1}}{1 - \delta_{ML}} W_{ML}(\gamma(e), \tau^{tw}). \quad (5)$$

to

$$W_{ML}(\gamma(e), \tau^a) + \frac{\delta_{ML} - \delta_{ML}^{T+1}}{1 - \delta_{ML}} W_{ML}(\gamma(e^a), \tau^a) \geq W_{ML}(\gamma(e), \tau^R(e), \tau^{*a}) + \frac{\delta_{ML} - \delta_{ML}^{T+1}}{1 - \delta_{ML}} W_{ML}(\gamma(e^{tw}), \tau^{tw}). \quad (6)$$

That is, the current-period median legislator would evaluate the current-period incentive constraint using a mixture of his own weight on import-competing profits and those of the legislators who are median in the future in the trade agreement and trade-war scenarios.

- What if you don't like the median legislator interpretation and want to think of it as a unitary decision-maker who faces different amounts of lobbying effort?
  - Think about what would happen in a simpler repeated game...

## 5 More comparison to MRC

### Referee 1 Point 2b:

However, I find the discussion immediately following the Result 1 quite interesting. The author shows that at the equilibrium, the applied tariff is equal to the negotiated binding. That is, the lobby group exert sufficient effort to induce the government to apply a tariff as high as the negotiated cap. Moreover, it is interesting that in the equilibrium the self-enforcing constraint for the legislature is not binding.

The former result is reminiscent of Maggi and Rodriguez-Clare's finding in their commitment model that the optimal form of trade agreement is a tariff cap, rather than a particular tariff rate. The justification is similar in both papers: by allowing the government to set a tariff below the cap, the lobby groups are induced to pay for the increase in the applied tariff, which reduces the incentives of the lobby groups to push for more protection.

I think more insights could be generated by comparing this paper and those of Maggi and Rodriguez-Clare. Therefore, I recommend the author to provide more discussion of how her paper is related to Maggi and Rodriguez-Clare.

### Editor Point 4:

Referee 1 mentions that your results are reminiscent of Maggi and Rodriguez-Clare's finding that weak bindings are preferable to strong bindings, though for slightly different reasons (there are no issues of self-enforcement in Maggi and Rodriguez-Clare). You do not explicitly compare a weak-binding agreement with a strong-binding agreement (unless I missed it), but if you could show that the former is preferable, this would be an interesting result worthy of emphasis.

Notes:

- Lobby and gov't bargain efficiently over tariff and contribution. Lobby has all the bargaining power, gov't is 'just' willing to improve tariff  $t$
- They can just add utility of both actors and maximize because utility is transferable
- To do bargaining as they do, need generalized NBS since non-transferable utility:

$$U = (G(t, c) - G^0)^\gamma (V(t, c) - V^0)^{1-\gamma}$$

– This from Drazen and Limao

–  $\gamma = 0 \Rightarrow$  lobby has all bargaining power:  $\max V(t, c) - V^0$  s.t.  $G(t, c) = G^0$

- Inherent in this story is the baseline where the lobby is not active (“ $G^0, V^0$ ”), which is hard to square with my formulation

How does weak binding agreement compare to strong-binding agreement?

- Lobby pays less (nothing) for trade agreement tariff
- Median legislator when trade agreement is in place is  $\gamma(0)$  instead of  $\gamma(e_a)$ , so not ‘efficient’ in this sense.
- Nothing changes in trade war
- Median legislator’s repeated-game incentive constraint is unchanged, so Lemmas 1 and 2 are unchanged
- Lobby’s incentive constraint is still satisfied at  $\tau^a = \tau^n$ , so soln always exists:

$$e \geq \pi(\tau^b(e)) - \pi(\tau^a) + \frac{\delta_L - \delta_L^{T+1}}{1 - \delta_L} [\pi(\tau^{tw}) - e_{tw} - \pi(\tau^a)] \quad (7)$$

At  $\tau^a = \tau^{tw}$ , we have

$$e \geq \pi(\tau^b(e)) - \pi(\tau^{tw}) + \frac{\delta_L - \delta_L^{T+1}}{1 - \delta_L} [\pi(\tau^{tw}) - e_{tw} - \pi(\tau^{tw})]$$

$$e \geq \pi(\tau^b(e)) - \pi(\tau^{tw}) - \frac{\delta_L - \delta_L^{T+1}}{1 - \delta_L} e_{tw}$$

$$\pi(\tau^{tw}) + \frac{\delta_L - \delta_L^{T+1}}{1 - \delta_L} e_{tw} \geq \pi(\tau^b(e)) - e$$

Because net profits are maximized at  $\tau^{tw}$ , this must hold for all  $e$ .

- Lobby’s constraint under strong binding is easier to satisfy than under tariff cap:

$$\begin{aligned} e &\geq \pi(\tau^b(e)) - \pi(\tau^a) + e_a + \frac{\delta_L - \delta_L^{T+1}}{1 - \delta_L} [\pi(\tau^{tw}) - e_{tw} - \pi(\tau^a) + e_a] \\ &\geq \pi(\tau^b(e)) - \pi(\tau^a) + \frac{\delta_L - \delta_L^{T+1}}{1 - \delta_L} [\pi(\tau^{tw}) - e_{tw} - \pi(\tau^a)] \end{aligned}$$

- The question: for a given  $T, \delta_L, \delta_{ML}$ , is  $\tau^a$  lower under strong binding or tariff cap?
  - In order to have an interior solution, need  $\bar{e}(\tau^a)$  to not decrease too quickly after its peak before lobby’s constraint can be satisfied (if lobby’s constraint were satisfied before reaching the peak of  $\bar{e}(\tau^a)$ , that would be great, but it doesn’t affect any of the arguments below because it’s the same in both the weak and strong binding cases).

- Recall  $\bar{e}$  is concave in  $\tau^a$  and need  $\bar{e} \geq e^{tw}$ . These facts don't change. The  $\bar{e}(\tau^a)$  schedule does not change.
- As  $\tau^a$  increases, net profits under agreement increase. They go up more under strong binding.  $\bar{e}(\tau^a)$  is decreasing at the same rate regardless of whether lobby is paying  $e_a$ .
- Strong binding version in which  $e_a$  is not being paid has a better chance of the constraint being satisfied / if the constraint is satisfied in both cases, will be satisfied at lower  $\tau^a$ .
- So strong binding is better for welfare-maximizing executive
- Strong binding is worse for the median legislator during the trade agreement phase: there's a mismatch between policy and  $\gamma(0)$  as opposed to perfectly matching the state under tariff cap.
- Strong binding is worse for median legislator under break  $\bar{e}$ ; he has no chance of receiving any lobbying effort during trade agreement phase. But this is all ex-post. Doesn't affect leg IC.
- So can see tariff cap, once again, as mechanism for ensuring that rents are distributed ex-post.

## 6 Punishment length vs. Dispute length

### Referee 2 Point 1:

My first comment concerns the punishment length,  $T$ . The paper highlights an interesting trade-off between the advantages of a high  $T$  (which helps to deter cheating by the legislature) and a low  $T$  (which helps to deter disruptive lobby effort), where the value for  $T$  in turn interacts with how low the cooperative bound tariff level can be pushed. But an ongoing concern for me is the interpretation of the endogenous determination of punishment length,  $T$ . I don't have a good sense of how to interpret the real-world determination of  $T$ . The author acknowledges that the  $T$ -period punishment never occurs along the equilibrium path but argues that the determination of  $T$  can be related to the choice of design for a dispute settlement system. To me, though, I think that there is some tension involved when using a model where all dispute behavior is off the equilibrium path as a means of interpreting an existing dispute settlement system or evaluating a proposed new system. Disputes actually happen, and dispute settlement systems are designed with that in mind.

As the author notes, dispute activity would occur along the equilibrium path if the model were modified to include shocks. Presumably, such a modification would reinforce the value of punishment phases of limited duration. But I worry that the model with equilibrium-path punishments could bring into play new considerations not currently featured in the author's model, since the frequency and duration of the punishment would then directly affect payoffs. And if shocks were public, then arguably the enrichment of the model to consider escape clause rules (contingent trade policies) becomes more compelling.

### Editor Point 6:

Regarding Referee 2's concerns, I agree on a basic point: it's not clear how much we can learn about the design of a dispute settlement system from a model that has no disputes (nor punishments) on the equilibrium path. My suggestion would be to take a more modest approach in pitching your analysis: rather than speaking to the design of dispute settlement procedures, you are making some more limited points about the optimal severity of punishments. You are in good company, by the way: a number of authors (including myself) have written papers on the optimal severity of punishments in trade agreements, and often using models with no punishments on the equilibrium path.

## 6.1 Compare to Park (2011)

### Referee 2 Point 2:

It is interesting in this regard to compare the author's model with that of Park (2011), who also describes trade-offs that lead to the determination of a finite value for  $T$ . In his model, however, the information structure is such that punishment phases are triggered after a government observes an extreme private signal, corresponding to a private shock, and then selects a public tariff that signals the beginning of a punishment phase. (The critical value for  $T$  balances the benefit of selecting a higher public tariff against the cost of triggering a punishment phase.) In his model, punishment phases occur along the equilibrium path, and alternative dispute systems might be pursued with the goal of reducing the frequency or severity of such punishments, if possible, as Park discusses. These kinds of considerations are hard to contemplate in the author's set up, though, since the disputes are off the equilibrium path.

Park (2011) comments on the literature:

- Bagwell and Staiger (2005) and more recently Bagwell (2008) analyse the issue of implementing trade agreements when each government is privately informed about its own domestic political pressure for protection. Their analysis differs from this paper because it focuses on identifying the structure of trade agreements that can induce the truthful revelation of private political pressure.
- Earlier models developed with respect to this issue, such as Dixit (1987), Bagwell and Staiger (1990), and Riezman (1991), suggest that the WTO may serve the role of helping countries coordinate on more efficient equilibria among the multiple equilibria that typically arise in a repeated game set-up. To model a more explicit role of the WTO, Kovenoch and Thursby (1993) assume that the DSP of the WTO has an informational superiority over trading countries in distinguishing between true violations and mistaken perceptions, which in turn enhances a reputation mechanism that supports cooperation. In a multilateral trading environment, Maggi (1999) shows that the WTO may facilitate cooperation-enhancing third-country sanctions by disseminating information about deviations. While these models introduce more specific roles for the WTO to play in coordinating a cooperative equilibrium, the literature has not resolved the question of why the WTO is necessary for coordination because these previous studies offer no theory of why countries could not coordinate a cooperative equilibrium in a non-WTO environment.

This paper represents the emergence of the WTO as a change in the observation structure of a repeated game. The presence of the WTO changes the nature of punishment-triggering signals from private into public. In the absence of the WTO, the private nature of signals of potential violations limits the flexibility of punishment phases that countries can employ because these phases must provide countries with the incentive for truthful revelation of



private signals in triggering punishments. The WTO can publicize its opinions on violations, which relaxes such a constraint in designing an optimal punishment scheme, enabling a better cooperative equilibrium even in the absence of any informational superiority of the WTO. This result contrasts with the analysis of Ludema (2001) who emphasizes that the DSP of the WTO may require trade agreements to be renegotiation-proof by promoting communication among countries prior to starting punishments. According to his analysis, such communication negatively affects cooperation by forcing countries to rely on weaker punishments.

- concealed trade barriers introduced similar to Riezman (1991); I assume that each country cannot directly observe the other country's local market price of its export. For example, a mixture of a consumption tax and a production subsidy can replicate the effect of a tariff
- Maggi and Staiger (2008) analyse the possible role that the DSP of the WTO plays in completing an incomplete contract and characterize the optimal choice of contractual incompleteness and the DSP design. In a related study, Maggi and Staiger (2009) characterize optimal remedies for breaches of trade agreements in the presence of uncertain political pressure for protection, for which the DSP may generate noisy signals. Beshkar (2008) analyses how the rulings of the DSP can affect renegotiation of trade agreements in the context of designing a direct revelation bargaining mechanism. However, they do not introduce imperfect private signals of potential deviations into their models, and so such signals play no role in their analyses of the DSP of the WTO.
- Hungerford (1991) develops a model in which the WTO plays a negative role in enforcing trade agreements because the model assumes that the DSP of the WTO involves uninformative and costly investigation.
- Forward cite search turns up Bajona and Ederington (2012 working paper), Anesi and Facchini (2016 working paper from Midwest)

## 7 Result 1

### **Referee 1 Point 2a:**

As part of Result 1, the author states that “The equilibrium trade agreement is never subject to dispute.” But this is true by construction of the equilibrium. In other words, the author finds an equilibrium that is self-enforcing and, thus, no dispute arises. Therefore, I believe that this statement is not sufficiently interesting or insightful to be part of Result 1.

## 8 Shorten section 8

### Editor Point 7:

Regarding the new section 8 on alternative punishments, I appreciate the additional work you have done, but given that the paper is quite long, I would suggest reducing this section to a couple of paragraphs that summarize the main points and provide the basic intuition. You can place the analysis in an online appendix, or keep it in the working paper version, if you like.

- *I've shortened the section and opted to keep the full analysis in the working paper only.*

## 9 Citation for optimal punishment length of infinity

### Editor Point 9:

I am not sure I understand the way you cite Klimenko-Ramey-Watson (KRW) in section 6. The point that in a standard model (without lobbying and without punishments on the equilibrium path) more severe punishments are always better is a completely standard result that was made in many papers before KRW.

## 10 Referee 2 Minor Comments

- i). Pages 8 and 9: Is  $\pi_x$  in (2) the same as  $\pi$  in (1)?
- ii). Page 12: The notation  $W_E$  and  $W^*_E$  is used in (3) before being formally defined.  $W_E$  is subsequently defined in (4).
- iii). Page 16: The author sometimes refers to "continuation values" as corresponding to current payoffs plus discounted future payoffs, but I think the more conventional use of the term refers only to the latter (i.e., to the  $V$  terms in the incentive constraint).
- iv). Page 19: The author says that, in the absence of a trade agreement, the lobby has no incentive to be active. I don't know what this means - wouldn't the lobby still choose  $e_{tw}$  in that case?
- vi). Page 21: I was a little confused by the statements of Lemmas 1 and 2. Are we assuming here that parameters are in a region such that an equilibrium exists under which  $\bar{e}$  exceeds  $e_{tw}$ ? Existence is treated on the next page, but I wasn't sure what was being assumed on this page and in these lemmas.
- vii). Pages 34-35: I was confused by the discussion of the credibility of the alternative punishment scheme. I probably missed something, but here is where I got confused. In a punishment period, if the lobby deviated with a very high  $e$ , and thereby generated a very high  $\gamma$  value, would the legislature be able to commit not to respond with a higher tariff? I.e., is the legislature's supposed lack of response to lobbying credible?

## Result 3

There is no interior solution in  $T$  under Referee 2's suggestion that gov't welfare in future periods be evaluated at the  $e$  that is realized at that time. I need to be able to explain Results 2 and 3 better as I defend the modeling choice.

- As  $T$  increases, both what the lobby has to pay (this also includes increase to  $\tau^b$ ) and the lobby's profit from the punishment phase increase.
- The increase in  $e$  as  $T$  increases comes from the interaction of a direct effect and an indirect effect.
  - When  $T$  increases, gov't feels punishment for longer so gov't is less willing to break the agreement. That means  $e$  has to be increased to make the gov't indifferent again (note that when  $e$  increases, the gov't becomes more willing to break the agreement because it places more weight on the lobby's profits).
- In my formulation, the direct effect decreases to 0 as  $T \uparrow$ ; the indirect effect (denominator) increases to  $\frac{\delta}{1-\delta}$ . Thus  $\frac{\partial \bar{e}}{\partial T}$  decreases to 0 as  $T \rightarrow \infty$ .
  - The lobby's benefit also decreases to 0 as  $T \rightarrow \infty$ , but at the same rate as the numerator of  $\frac{\delta}{1-\delta}$  so  $\frac{\delta}{1-\delta}$  goes to zero faster than the lobby's benefit.
- In Referee 2's formulation, the essential difference is that the denominator of  $\frac{\delta}{1-\delta}$  is not a function of  $T$  because the current government's  $e$  does not show up in the expression for the punishment period (the current government evaluates the punishment period through the eyes of the government who will be in power at that time).
  - The terms involving  $T$  cancel out of both the cost and benefit side. Thus the constraint varies linearly in  $T$  in this formulation. The change is either a net positive (cost to lobby  $\uparrow\uparrow$  while benefit only  $\uparrow$ ) so want  $T \rightarrow \infty$ ; or the change is a net negative (cost to lobby  $\uparrow$  while benefit  $\uparrow\uparrow$ ) so want  $T \rightarrow 0$ .

Useful equations

$$\frac{\delta_{\text{ML}} - \delta_{\text{ML}}^{T+1}}{1 - \delta_{\text{ML}}} [W_{\text{ML}}(\gamma(e), \tau^a) - W_{\text{ML}}(\gamma(e), \tau^{tw})] \geq$$

$$W_{\text{ML}}(\gamma(e), \tau^b(e), \tau^{*a}) - W_{\text{ML}}(\gamma(e), \tau^a)$$

$$\left(1 - \frac{d\pi}{d\bar{e}}\right) \frac{-\frac{\delta_{ML}^{T+1} \ln \delta_{ML}}{1-\delta_{ML}} [W_{ML}(\gamma(\bar{e}), \boldsymbol{\tau}^a) - W_{ML}(\gamma(\bar{e}), \boldsymbol{\tau}^{tw})]}{\frac{\partial \gamma}{\partial e} [\pi(\tau^b(\bar{e})) - \pi(\tau^a)] + \frac{\delta_{ML} - \delta_{ML}^{T+1}}{1-\delta_{ML}} \frac{\partial \gamma}{\partial e} [\pi(\tau^{tw}) - \pi(\tau^a)]} + \frac{\delta_L^{T+1} \ln \delta_L}{1-\delta_L} [\pi(\tau^{tw}) - e_{tw} - \pi(\tau^a) + e_a] \quad (8)$$

$$\frac{\partial \bar{e}}{\partial T} = \frac{-\frac{\delta_{ML}^{T+1} \ln \delta_{ML}}{1-\delta_{ML}} [W_{ML}(\gamma(\bar{e}), \boldsymbol{\tau}^a) - W_{ML}(\gamma(\bar{e}), \boldsymbol{\tau}^{tw})]}{\frac{\partial \gamma}{\partial e} [\pi(\tau^b(\bar{e})) - \pi(\tau^a)] + \frac{\delta_{ML} - \delta_{ML}^{T+1}}{1-\delta_{ML}} \frac{\partial \gamma}{\partial e} [\pi(\tau^{tw}) - \pi(\tau^a)]}$$

$$\frac{\partial \Omega}{\partial e} \frac{\partial \bar{e}}{\partial T} + \frac{\partial \Omega}{\partial T} = 0$$

$$- \left[ \frac{\partial \gamma}{\partial e} [\pi(\tau^b(\bar{e})) - \pi(\tau^a)] + \frac{\delta_{ML} - \delta_{ML}^{T+1}}{1-\delta_{ML}} \frac{\partial \gamma}{\partial e} [\pi(\tau^{tw}) - \pi(\tau^a)] \right] \frac{\partial \bar{e}}{\partial T} + \left[ -\frac{\delta_{ML}^{T+1} \ln \delta_{ML}}{1-\delta_{ML}} [W_{ML}(\gamma(\bar{e}), \boldsymbol{\tau}^a) - W_{ML}(\gamma(\bar{e}), \boldsymbol{\tau}^{tw})] \right] = 0 \quad (9)$$

- First term: When  $e$  changes, it increases the weight on the lobby's profits relative to everything else. This loosens the constraint.
  - It gets loosened more the larger is  $T$  because the lobby's profits get bigger as  $T$  increases.
- Second term: the per-period punishment is felt for more periods, but the effect is decreasing b/c  $\delta < 1$
- When  $T$  gets very large, the leg constraint is not being loosened much by an extra period of punishment [DIRECT effect]
  - BUT when  $e$  increases, you care about the benefit to the lobby in every period, not just the incremental one [INDIRECT effect].
  - So the tightening of the legislative constraint through  $e$  / the indirect effect is larger than the loosening through the direct effect as  $T$  grows large and so  $\frac{\partial \bar{e}}{\partial T}$  must shrink.
  - At some  $T$ , the increase in  $e$  can no longer outweigh the increase in the benefit to the lobby