

Economica Coase Lecture

Reference Points and the Theory of the Firm

By OLIVER HART*

Harvard University

Final version received 11 September 2007.

I argue that it has been hard to make progress on Coase's theory of the firm agenda because of the difficulty of formalizing haggling costs. I propose an approach that tries to move things forward using the idea of grievement costs, and apply it to the question of whether a transaction should be placed inside a firm (in-house production) or in the market place (outsourcing).

INTRODUCTION

Let me begin by saying that it has been a great pleasure to write this paper for two reasons. First, I am a great admirer of Ronald Coase and his work, which has had an enormous influence on me. Second, I take as my starting-point Coase's 1937 paper on the nature of the firm, which Coase began while he was an undergraduate at LSE and published in *Economica* a few years later when he was on the staff.

I will divide the paper into two parts. In the first I will discuss Coase's 1937 paper and what followed; one of my themes will be that it has been quite difficult to make progress on the Coasian agenda. In the second part I will turn to some recent work (with John Moore) on contracts as reference points, and argue that this may be helpful in moving things along. Not surprisingly, this part of the paper is preliminary and speculative.

I. COASE (1937) AND WHAT FOLLOWED

In his 1937 article, Coase first raised the question of why we have firms at all in a modern market economy. If, as economists usually suggest, markets are so good at allocating resources, why do we need firms? Coase recognized that the converse question also has to be answered. Firms cannot always be better at allocating resources, since if they were we would not have markets.

In D. H. Robertson's words, we find 'islands of conscious power in this ocean of unconscious co-operation like lumps of butter coagulating in a pail of buttermilk' (1923, p. 85). The challenge is to explain the relative amounts of butter and buttermilk. Coase's question is breathtakingly simple and original. It is hard to find anyone earlier who had formulated the issue in such stark terms.

In order to convince oneself that Coase's question is still relevant in the early part of the twenty-first century, consider the following:

1. In February 2007 there were 41 companies in the world with a market value of equity greater than US \$100 billion.¹

*This paper is based on the inaugural Coase Lecture, presented at the London School of Economics in February 2007.

TABLE 1
EMPLOYEE- WEIGHTED AVERAGE SIZE OF FIRMS

	1988	2001
France	811	727
Germany	769	725
Italy	474	296
Spain	306	328
United Kingdom	859	935

2. In 2007 Wal-Mart, the largest US employer, had 1.8 million employees.²
3. The value of transactions in US firms is approximately equal to that in US markets.³
4. As Table 1 shows, the employee-weighted average size of firms—that is, the average number of employees in the same firm as a randomly selected employee—is sizeable, ranging from 296 to 935 in different countries in 1988 and 2001.⁴
5. In a sample of 43 countries, two-thirds of the growth in industries over the 1980s came from growth in the size of existing firms.⁵
6. The boundaries of firms keep changing. In 2006 the world-wide value of mergers and acquisitions exceeded US \$4 trillion.⁶

How did Coase answer the question of why firms and markets coexist? Coase argued that each way of organizing economic activity is costly. He suggested that the two most obvious costs of using the market/price mechanism are (i) discovering what the market prices are and (ii) negotiating a contract for each exchange transaction. Economists since Coase have referred to these as ‘haggling’ costs (although I do not believe that Coase uses this term). ‘Argument’ costs might also be appropriate for (ii).

According to Coase, haggling costs are avoided inside the firm because bargaining is replaced by authority: an employer tells an employee what to do and (within limits) the employee obeys. However, authority is also costly. Coase emphasizes two main costs. First, as a firm gets bigger, the entrepreneur or manager running it will find it increasingly difficult to organize the firm’s activities given his or her limited (intellectual) capacity. Second, and related to this, the person in charge will make mistakes.

Coase’s questions about why firms and markets co-exist are brilliant, but his answers are less satisfactory. There are at least three problems with them. First, it has been very difficult to formalize or operationalize haggling costs, as I will explain further below. Second, Coase’s costs of using the firm are unconvincing. Why cannot an overstretched manager hire another manager to help him out? Third, it seems optimistic and unrealistic to suppose that placing a transaction inside a firm eliminates all haggling: there are many examples in practice of serious and costly disagreements inside firms. Coase has made life hard for his followers by never attempting to write down a formal model. Interestingly, as far as I know Coase has also never shown any indication that he thinks that such an activity is in the least bit worth while!

Why has it been difficult for economists since Coase to operationalize haggling costs? The reason for this can be traced to a paper published in 1960 by an economist every bit as important and famous as Ronald Coase, called—Ronald Coase. Coase (1960), writing on the seemingly entirely different topic of externalities, developed a theory with the following implication. Suppose that you can provide a good (or service) that is worth more to me than it costs you to produce—say it is worth 20 to me and costs you 10. Then

it would be silly for us not to trade the good at *some* price between 10 and 20, since we are both made better off. We will probably argue about what the price should be—I would like $p = 10$ and you would like $p = 20$. But presumably we'll settle on something in between; for example, we might split the difference at $p = 15$.

Now if we realize that this is going to happen, why do we not agree on $p = 15$ right away and avoid the haggling costs? This argument—that rational parties will avoid, or bargain around, haggling costs 'in the twinkling of an eye'—is an implication of what has become known as the Coase theorem.⁷

To emphasize the point, Coase (1960) is problematic for followers of Coase (1937) because it suggests that haggling costs may be small or even zero. In other words, markets may not be so costly after all.

How has the literature dealt with this difficulty? Mainly by sidestepping the issue and introducing new features. For example, the transaction cost literature (see Williamson 1971; Klein *et al.* 1978) has emphasized that firms will be efficient relative to markets when parties make large relationship-specific investments and cannot write good (i.e. complete) contracts. The property rights approach (see Grossman and Hart 1986; Hart and Moore 1990) has refined these ideas by arguing that, in the presence of relationship-specific investments and incomplete contracts, firms and markets are both costly; which organizational form is chosen will depend on such things as the relative importance of the parties' relationship-specific investments.⁸

The transaction cost literature does not take a formal approach and has implicitly assumed the existence of haggling (or rent-seeking) costs: it has not confronted the issue of how to model them. The more formal property rights approach has sided with Coase (1960), thereby avoiding haggling costs. In a typical model the parties bargain costlessly *ex post*, and the focus is on *ex ante* investment inefficiencies. I have argued elsewhere (Hart and Moore 2007) that, while such an approach can yield useful insights about optimal asset ownership, it is unlikely to be helpful for studying the internal organization of large firms. Specifically, in a world of Coasian bargaining, it is hard to see why important aspects of organizational form such as authority, hierarchy and delegation matter. Why would the parties not simply bargain about everything all the time, using monetary side-payments?

In my view, in order to make progress on the Coasian agenda, we must move away from Coase (1960) and back in the direction of Coase (1937). We need to bring back haggling costs! This will be the subject of the next part of the paper.

II. PUTTING HAGGLING COSTS BACK INTO THE PICTURE

A recent paper with John Moore (Hart and Moore 2008) tries to incorporate something akin to haggling costs. I will begin by describing it and will then apply it to the theory of the firm.

The best way to introduce this paper is to go back to the example where a seller (henceforth S) can provide a good that costs 10 and is worth 20 to a buyer (henceforth B). To fix ideas, imagine that we are talking about a musical evening that B is arranging at his house, and at which B wants S to sing. The musical evening is worth 20 to B, and S's cost of performing (an effort cost, say) is 10.

For the moment, ignore the fact that B could engage other singers or that S could perform elsewhere on the night in question. Earlier I argued that B and S might agree to trade at a price of 15. That discussion implicitly assumed that, once B and S agreed, trade

would proceed smoothly. But suppose that is not so. In particular, assume that **B** and **S** each have some discretion about the ‘quality’ of performance they provide, i.e. how pleasant they make the experience for the other party. **S** can perform within the letter rather than the spirit of the contract, or can stint on quality; for example, she can be rude to **B**’s guests or refuse to give autographs. **B** can quibble about the details of performance or be slow in paying.

To use the language of Hart and Moore (2008), each party has the discretion to provide ‘perfunctory’ (basic) or ‘consummate’ (exemplary) performance. It is worth emphasizing that this is a significant departure from the standard contracting literature. The literature usually assumes that trade is perfectly enforceable *ex post* (e.g. by a court of law). Here we are assuming that only perfunctory performance can be enforced: consummate performance is always discretionary.

What determines whether a party provides consummate performance? Hart and Moore (2008) appeal (quite loosely) to a number of ideas from the behavioural economics literature. It is assumed that a party is roughly indifferent between providing perfunctory and consummate performance—consummate performance costs only slightly more or may even be slightly more pleasurable—and will provide consummate performance if he is ‘well treated’ but not if he is ‘badly treated’ (negative reciprocity).

Importantly it is supposed that a party feels ‘well treated’ if he gets what he believes he is entitled to; that a contract is a reference point for perceived entitlements; and that in the absence of a reference point entitlements can diverge, possibly wildly.

Let us apply these ideas to our 20/10 example. First, let us put in a time line (Figure 1). The time line captures the idea that **B** and **S** will typically write a contract some months before the musical evening takes place (date 0) rather than the night before (date 1). One reason for this is that each will have more alternatives earlier on. In fact, I am going to assume that there is a competitive market for sellers, i.e. singers, at date 0.

Let us suppose first that, although **B** and **S** sign a contract at date 0, they leave the determination of how much **B** will pay **S** until the night before the concert, date 1. This may seem odd, and indeed I will show that it is a bad idea. If **B** and **S** do not specify price, then p can be anywhere between 10 and 20. What might each party feel entitled to?

As mentioned, Hart and Moore (2008) take the view that entitlements can diverge. **S** can convince herself that she is hugely talented and that her presence is the entire reason the evening will be a success. **S** feels she is entitled to $p = 20$. **B** has a dimmer view of **S**’s ability and contribution and thinks that **S** is worth much less: p should be 10.

Even though **B** and **S** have these different views of what p should be, they are rational enough to come to some agreement; let us say they split the difference at $p = 15$. However, each feels short-changed and aggrieved. **B** thinks that he has paid 5 too much; **S** thinks that she has been paid 5 too little. Neither of them is in the mood to provide consummate performance.

The precise assumption made in Hart and Moore (2008) is that each party feels entitled to the best outcome consistent with the contract and ‘shades’ on consummate performance in proportion to the amount he feels aggrieved. Since **B** is aggrieved by 5, **B** shades to the point where **S**’s payoff falls by 5θ , where θ is the constant of



FIGURE 1. Time line.

proportionality: it might be 0.2, say. And since S is aggrieved by 5, S shades to the point where B's payoff falls by 5θ .

The bottom line is that, if B and S leave the price open until the night before the concert, there will be a total deadweight loss of 10θ owing to shading. This is money down the drain. It reduces the value of B and S's relationship from 10 to $10(1-\theta)$. If $\theta = 0.2$, the relationship is worth 8 instead of 10.

Economists do not like deadweight losses (nor does anyone else). Can anything be done to avoid them here? The answer is yes. Note first that *ex post* Coasian bargaining at date 1 does not do the job. The reason is that shading is not contractible, and therefore a contract not to shade is not enforceable. To put it another way, if B offers to pay S more not to shade, e.g. offers $p = 16$, then, while this will indeed reduce S's shading (from 5θ to 4θ), since S will feel less aggrieved, it will increase B's shading (from 5θ to 6θ) because B will feel more aggrieved! Total deadweight losses remain at 10θ . However, there is a simple solution: to put the price in the contract at date 0. Since I have supposed that there is a competitive market for singers at date 0, B will be able to hire S for $p = 10$. With $p = 10$ specified in the contract, there is nothing for B and S to argue about at date 1. The fact that B and S may disagree about S's talents as a singer does not matter any more. B and S have agreed that B will pay S 10, and neither B nor S will be disappointed or aggrieved when that happens given that the contract is a reference point for entitlements.

In short, a contract that sets $p = 10$ in advance eliminates *ex post* argument and aggrievement, and hence both parties will be willing to provide consummate performance. Deadweight losses will be zero and the first-best will be achieved.⁹

Aggrievement and shading costs are a bit like haggling costs, and so the above model moves us in the direction of Coase (1937). However, we have to add one further ingredient to get to a theory of the firm.

Let us now introduce the realistic notion that not all the details of the musical evening can be anticipated at date 0. To make it simple, imagine that the musical evening can be carried out in two ways, i.e. according to two methods. (We might be talking about the exact songs, who are the other performers, the order of the program, etc.)

In Figure 2, method 1 yields value 20 and costs 10, as above. Method 2 yields value 14 to the buyer and costs the seller 8. Assume that the methods cannot be specified in the date 0 contract, e.g. because they are too complicated to describe in advance. However, the choice between them becomes clear at date 1. Note that with these numbers method 1 is more efficient than method 2 since it generates higher surplus.

Compare two different organizational forms. In the first, B and S fix the price of the good at date 0 (at 10, say) and determine that S will be an independent contractor. In other words, this is a market exchange between two distinct economic 'entities'. I will take this to mean that S has the right to decide on the details of production, i.e. on the choice between methods 1 and 2.¹⁰ I believe that this accords with the common understanding of what it means to be an independent contractor.

What will S do? Given that the price is fixed, S will pick method 2, since it is cheaper. This is inefficient. B will then be aggrieved that S didn't choose method 1—B will feel

	<u>Method 1</u>	<u>Method 2</u>
Value	20	14
Cost	10	8
Surplus	10	6

FIGURE 2. Payoffs from methods.

entitled to this, and will regard S's choice as ungenerous; B is short-changed or aggrieved by 6 (his payoff would be 6 higher under method 1), and he will shade to the point where S's payoff falls by 6θ . Total surplus = $6-6\theta$.

Now consider a second organizational form. B and S agree at date 0 that S is an employee: S will work for B at a fixed wage (10, say). I am going to take this to mean that B has the right to decide on the choice between method 1 and method 2—and indeed, this accords with common usage of the term 'employment'. Given the fixed wage, B will of course choose method 1, since it gives B more value. This is efficient. S will be aggrieved that B did not choose method 2, but S's aggrievement is only 2. Total surplus = $10-2\theta$.

The conclusion is that in this example employment is the better arrangement. Employment is good for two reasons. First, the production method matters more to B than to S and so it is efficient that B chooses it. Second, and related, S's aggrievement will be low because S does not care that much.¹¹

Now change the numbers. Keep method 1 the same but suppose that method 2 yields value 14 and costs 2 (see Figure 3). Method 2 is now more efficient. Under employment, however, the buyer will choose method 1, yielding surplus $10-8\theta$. Independent contracting is superior here because the seller will select method 2, yielding surplus $12-6\theta$.¹² We see that employment is good if the production method matters more to B than to S, while independent contracting is good if the production method matters more to S than to B.

One point worth emphasizing is that in neither of the above examples is the following contract optimal: to leave the choice of price *and* method until date 1, i.e. to rely on unconstrained Coasian bargaining. This would always yield the efficient method, but the aggrievement costs would be high. In Figure 2 the parties would agree on method 1; however, since there are 10 dollars of surplus to argue over, shading costs equal 10θ : net surplus = $10(1-\theta)$, which is less than that obtained under the employment contract. In Figure 3 there are 12 dollars of surplus to argue over and net surplus = $12(1-\theta)$, which is less than that obtained under independent contracting.

The examples in Figures 2 and 3 are obviously 'toy' ones (for instance, there are only two production methods), but I believe that they contain the ingredients of a theory of the choice between doing a transaction 'in the market place', i.e. through independent contracting, and 'inside the firm', i.e. through employment. The theory is in the spirit of Coase (1937), but perhaps a bit more satisfactory in some respects. I have replaced haggling costs by aggrievement costs, but they are not so different, since both have to do with not getting your way. I have stressed that who controls or decides the production method is a key issue in choosing between the two organizational forms. I have also emphasized that aggrievement costs arise under employment as well as under independent contracting. I have not had to suppose, as Coase did, that managers of large firms make mistakes, in order to explain why not everything is done inside the firm.

There are, of course, many things missing from the framework described above. One important one is that it has been supposed that the cost of production is always borne by

	<u>Method 1</u>	<u>Method 2</u>
Value	20	14
Cost	10	2
Surplus	10	12

FIGURE 3. Modified payoffs.

the seller. In practice, it may be possible to transfer some costs to the buyer, contractually, through a cost-sharing arrangement. Note, however, that if B bears most of the costs it may also be efficient for B to control the production method, since B now cares more about this than S. That is, one might expect a positive correlation between B's bearing the production cost and the use of the employment relationship.¹³ Of course, a disadvantage of transferring costs to B is that S will have little incentive to keep them low. Introducing cost sharing and worker/manager effort would enrich the model greatly.

Even with these qualifications, our model has some plausible empirical implications. The model suggests that outsourcing is likely to be efficient when a detailed contract can be written about the nature of the good to be delivered, since in this case B's value will be pretty insensitive to the choice of production method while S's cost may not be (see Figure 3). In contrast, if a detailed contract is hard to write and B's value is very sensitive to the details of production, then in-house production may be better (see Figure 2).

Hart *et al.* (1997), in a paper about the choice between government and private ownership, argue that municipal rubbish collection probably falls into the first category, and fighting wars into the second. The provision of prison services may be somewhere in between. Using a model based on non-contractible relationship-specific investments, Hart *et al.* derive some conclusions about the costs and benefits of privatizing these services. It may be interesting to revisit their analysis, and the issue of outsourcing more generally, using an *ex post* inefficiency model of the type described here.

III. CONCLUSIONS

I have argued that it has been hard to make progress on Coase's 1937 agenda because of the difficulty of formalizing haggling costs. I have sketched an approach based on the idea of aggrievement costs, and have suggested that it can throw light on whether a transaction should be placed inside a firm (in-house production) or in the market place (outsourcing). To emphasize the obvious, the analysis I have presented is preliminary and rudimentary. How useful it will be in the development of a general theory of the firm, time will tell.

ACKNOWLEDGMENTS

I am grateful to Paul Niehaus for research assistance and to the US National Science Foundation through the National Bureau of Economic Research for financial support.

NOTES

1. http://www.forbes.com/lists/2006/18/06f2000_The-Forbes-2000_MktVal.html, accessed 2 February 2007.
2. <http://www.walmartfacts.com/featuredtopics/?id=3>, accessed 2 February 2007.
3. See Lafontaine and Slade (2007).
4. These data come from the Structural Business Statistics database, which is maintained by Eurostat, the statistical department of the European Commission. They are available online at <http://epp.eurostat.ec.europa.eu/> where they were accessed during February 2007. Employee-weighted average firm sizes were estimated using the method of Kumar *et al.* (2002). They are biased downwards, given the coarseness of the firm size bins in which the underlying Eurostat data are reported. They are also quite sensitive to the sizes of, and averaging process for, the largest firms; e.g. a merger or divestiture could have a large impact on the data for a country. This sensitivity may be responsible for some of the variability in Table 1.
5. See Rajan and Zingales (1998).
6. See *The Economist*, Indicators section, 13 January 2007.

7. Of course, the Coase theorem requires assumptions such as symmetric information, the existence of enforceable contracts, etc. For a recent summary, see de Meza (1998).
8. For recent summaries of the transaction cost and property rights literatures, see Gibbons (2005).
9. An obvious question to ask is what changes between dates 0 and 1? Why does a date 0 contract that fixes p avoid aggravement, whereas a date 1 contract that fixes p does not? The *ex ante* market at date 0 provides a crucial role here. It provides an objective measure of what B and S bring to the relationship. Given that there are many sellers willing to supply at $p = 10$, S accepts that she cannot expect to receive more than 10, while B understands that he can't expect to pay less. Thus, neither party is aggravated by $p = 10$. See Hart and Moore (2008) for further discussion.
10. That is, S has the residual rights of control in the sense of Grossman and Hart (1986).
11. Readers may wonder whether there are contracts that do even better than the employment contract, i.e. achieve surplus greater than $10 - 2\theta$. One possibility is to set a very high contract price, e.g. $p = 20$. Suppose that neither party feels entitled to more than 100% of the gains from trade, and each party has the right to quit the relationship. Then method 1 will be chosen both under employment (for the same reason as in the text) and under independent contracting (if S chooses method 2, B will quit) and there will be no aggravement since there is nothing to argue about. In other words, the first-best can be achieved. However, as Hart and Moore (2008) show, this is a knife-edge result, depending on the absence of uncertainty. With a small bit of uncertainty, the employment contract with $p = 10$, say, is optimal; i.e. maximum surplus = $10 - 2\theta$.
12. I have implicitly assumed that under independent contracting the seller can choose the method without violating the original contract. But if the original contract is tight, this may not be true. Switching methods may correspond to a 'change order', in which case both parties may have to approve. This raises interesting new possibilities. For a more general discussion of change orders, see Bajari and Tadelis (2001).
13. Correlations like this are studied in Holmstrom and Milgrom (1994).

REFERENCES

- BAJARI, P. and TADELIS, S. (2001). Incentives versus transaction costs: a theory of procurement contracts. *RAND Journal of Economics*, **32**, 387–407.
- COASE, R. (1937). The nature of the firm. *Economica*, **4**, 386–405.
- (1960). The problem of social cost. *Journal of Law and Economics*, **3**, 1–44.
- DE MEZA, D. (1998). The Coase theorem. *Palgrave Dictionary of Economics and the Law*. Basingstoke: Palgrave Macmillan.
- GIBBONS, R. (2005). Four formal(izable) theories of the firm? *Journal of Economic Behavior and Organization*, **58**, 200–45.
- GROSSMAN, S. and HART, O. (1986). The costs and benefits of ownership: a theory of vertical and lateral integration. *Journal of Political Economy*, **94**, 691–719.
- HART and MOORE, J. (1990). Property rights and the nature of the firm. *Journal of Political Economy*, **98**, 1119–58.
- and ———. (2007). Incomplete contracts and ownership: some new thoughts. *American Economic Review*, **97**, 182–86.
- and ———. (2008). Contracts as reference points. Forthcoming in *Quarterly Journal of Economics*.
- , SHLEIFER, A. and VISHNY, R. (1997). The proper scope of government: theory and an application to prisons. *Quarterly Journal of Economics*, **112**, 1127–61.
- HOLMSTROM, B. and MILGROM, P. (1994). The firm as an incentive system. *American Economic Review*, **84**, 972–91.
- KLEIN, B., CRAWFORD, R. and ALCHIAN, A. (1978). Vertical integration, appropriable rents, and the competitive contracting process. *Journal of Law and Economics*, **21**, 297–326.
- KUMAR, K., RAJAN, R. and ZINGALES, L. (2002). *What determines firm size?* Mimeo, University of Chicago Graduate School of Business.
- LAFONTAINE, F. and SLADE, M. (2007). Vertical integration and firm boundaries: the evidence. *Journal of Economic Literature*, **45**, 629–85.
- RAJAN, R. and ZINGALES, L. (1998). Financial dependence and growth. *American Economic Review*, **88**, 559–86.
- ROBERTSON, D. (1923). *The Control of Industry*. Hitchin, Herts: Nisbet.
- WILLIAMSON, O. (1971). The vertical integration of production: market failure considerations. *American Economic Review*, **61**, 112–23.

Copyright of *Economica* is the property of Blackwell Publishing Limited and its content may not be copied or emailed to multiple sites or posted to a listserv without the copyright holder's express written permission. However, users may print, download, or email articles for individual use.