



Forward reliability markets: Less risk, less market power, more efficiency

Peter Cramton^{a,1}, Steven Stoft^{b,*}

^a Economics Department, University of Maryland, College Park, MD, USA

^b 2910 Elmwood Court, Berkeley, CA 94705, USA

Received 10 January 2008; accepted 12 January 2008

Abstract

A forward reliability market is presented. The market coordinates new entry through the forward procurement of reliability options—physical capacity bundled with a financial option to supply energy above a strike price. The market assures adequate generating resources and prices capacity from the bids of competitive new entry in an annual auction. Efficient performance incentives are maintained from a load-following obligation to supply energy above the strike price. The capacity payment fully hedges load from high spot prices, and reduces supplier risk as well. Market power is reduced in the spot market, since suppliers enter the spot market with a nearly balanced position in times of scarcity. Market power in the reliability market is addressed by not allowing existing supply to impact the capacity price. The approach, which has been adopted in New England and Colombia, is readily adapted to either a thermal system or a hydro system.

© 2008 Published by Elsevier Ltd.

1. Understanding the generation adequacy problem

The reliability of a power system depends on how it is operated in the short term (security) and medium term (firmness), but if there is not enough generating capacity, it will not be possible to serve all load and achieve security and firmness.² In this way, adequate generation is the most fundamental reliability issue, and it is also the one most distant from the spot market because it is the most long-term aspect of reliability. This paper focuses only on the long-term issue of generation adequacy. However, the importance of the prescribed reliability-options approach to adequacy is that it facilitates the solution to the three worst problems of contemporary electricity markets: investment risk, market power and inefficient pricing.³

Contrary to conventional wisdom, the level of generation adequacy is not much of a problem for two reasons. First, marginal generating capacity is relatively inexpensive when compared with other costs of delivered energy. Second, as with most optima, the derivative of net benefit with respect to capacity is zero at the optimal capacity level. For example, an extra 10% of capacity increases capacity costs by much less than 10% because peaking capacity is by far the cheapest kind of capacity, and adding peak capacity does not increase fuel costs, transmission costs or administrative costs. As a consequence, increasing total capacity by 10% will cost consumers only, perhaps, 2% extra. But there is some benefit to the resulting extra reliability. So the loss of net benefit is less than 2%. A good regulatory approach is unlikely to overshoot by more than 10% on average, and the best market-based approach will not be perfect. Hence, the net benefit of improved adequacy from any market-based approach is necessarily quite small—likely less than 1% of total retail cost. The true cost of the adequacy problem has been the distortion of market designs by misguided attempts to solve it. These designs cause risks, inefficiencies and regulatory responses that are far more costly than any likely mistake in the provision of adequacy. Of course the adequacy problem needs a solution,

* Corresponding author. Tel.: +1 510 644 9410.

E-mail addresses: p.cramton@gmail.com (P. Cramton), steven@stoft.com (S. Stoft).

¹ Tel.: +1 207 42 7211.

² For a more complete explanation of reliability see the paper in this issue by Batlle and Pérez-Arriaga.

³ These problems are actually solved by the capacity-market/reliability option design, provided that the rest of the market design is reasonable and the market does not have structural problems such as high supplier concentration.

but that solution should not exacerbate already difficult aspects of electricity markets.

The misconception most responsible for the current state of affairs is the notion that a cleverly designed “energy-only” market can induce optimal adequacy, or something close to it, even while the market has insufficient demand elasticity. Interestingly, when the notion of reliability markets was first developing back in the late 1990s, the importance of adequate demand response in an “energy-only” market was fully recognized.

In this approach there is no price cap that limits the market price. It is assumed that the elasticity of demand to prices is enough to prevent any occurrences of market’s failure to supply because of lack of generation adequacy.

Pérez-Arriaga (1999)

Unfortunately, when the idea of reliability options spread to the US, these assumptions were sometimes replaced with the idea that options priced only on the basis of their financial cost would solve the reliability problem in spite of a lack of sufficient demand elasticity.

In an ideal market, with sufficient demand elasticity, the market always clears. This means there can be no adequacy problem because involuntary load shedding occurs only when the market fails to clear and demand exceeds supply. In a market that always clears, energy prices do not and cannot determine the level of reliability for, with respect to adequacy, the market is automatically 100% reliable. Instead, energy prices do what every economics text says they do, they determine the efficient (not reliable) level of capacity. Any less capacity would cause expensive *voluntary* load reductions, and any more capacity would mean that too much had been spent on capacity. The crucial point is that no energy-only market, even with ideal demand elasticity, can solve the adequacy problem. No energy market, on its own, can ever answer the question: “What level of capacity provides optimal reliability?” That would be the same as a market discovering a positive and optimal amount of time that it should fail to clear. No economic theory claims any market can do this.

Of course if reliability is sold as a product, so that the customer can pay more to gain reliability, the market can answer the reliability question. But that is not an energy-only market, and it requires technology that is not yet in place. At present, you cannot buy more reliability than your neighbor, because you are on the same physical circuit and neighbors will always be blacked out together.

Rather than attempting a reliability market, it may be better to install real-time meters and use real-time pricing to increase elasticity to the point where the market becomes perfectly reliable with regard to adequacy. In principle this can be achieved with prices well below the Value of Lost Load (VoLL). As long as price stays below that level, such a market will be more efficient than an inelastic market with even a perfectly optimal level of adequacy. Even so, one should not hope for dramatic efficiency gains because peaking capacity is cheap relative to the total cost of power.

As just noted, those who believe an energy market can solve the adequacy problem have simply misunderstood the economic theory of optimal investment. In a competitive market, optimal investment has nothing to do with reliability. To understand why ~~is~~ sometimes appears as if there is a connection, a more detailed look at energy markets is helpful.

Energy markets fall into one of two categories: Case 1, they always clear and have no adequacy problem; or Case 2, they can fail to clear and do have an adequacy problem. To classify a market correctly it is necessary to consider how it would perform without market power or regulation.

Case 1 occurs when there is enough demand elasticity so that the supply and demand curves always intersect as they do in any normal market. With the market clearing at all times, there will be a long-run capacity equilibrium, but it will not be efficient if the price sometimes exceeds the VoLL. VoLL varies with time, but at any given time it is the average price that all consumers in the blackout region would pay to avoid the blackout.⁴ Suppose this value is \$2000/MWh. It is quite possible for the demand curve to be downward sloping at \$50,000/MWh or any other value. In this case the market could clear at \$50,000 even though VoLL was only \$2000. Paying \$50,000 for power that is only worth \$2000 is clearly inefficient, and because it pays suppliers too much, it causes excess entry, and the capacity level will end up inefficiently high. Economics does not predict, as many imagine, that simply because a market clears, it is efficient. Efficiency depends on all consumers expressing their demand in the market. Next consider the sub-case in which there is just enough demand elasticity to keep spot prices in Case 1 below VoLL, both risk and market power are almost certain to reach problematic levels. Under such conditions, the regulated, reliability market described in this paper is still almost certain to outperform an unregulated energy-only market. Only when a market’s demand side is functioning quite normally and spot prices stay well below VoLL, is a pure energy market at all likely to outperform a good capacity market, and even then, better investor coordination may favor the capacity market.

Case 2, in which reliability is not guaranteed to be 100%, is the case considered in this paper. In most markets, few observers are willing to guarantee that involuntary load shedding is out of the question if the market is left to its own devices. Even so, some markets appear to be reliable on their own. Does this mean competitive energy prices are coming close to solving the reliability problem? Not at all. There are three possible explanations. First, suppliers may be exercising enough market power to attract new entrants. Monopoly power can easily produce more than enough reliability, but the resulting level has nothing to do with efficiency. Second, reliability can be the result of unrecognized regulatory interventions, such as tampering with the demand for operating reserves.

⁴ Although VoLL is not measurable that does not mean it does not exist or does not matter. Although inter-temporal preference make VoLL difficult to define in real-world situations, the above analysis is rigorous in simplified models, and in complex models the inefficiencies would be no less. The complexities of the real world are extremely unlikely to produce an efficient outcome in situations where simple models predict inefficiency.

Third, it is possible, though unlikely, that LSEs are expressing their own demand for reliability in a market that, in effect, is selling reliability to LSEs. For such an LSE-reliability market to work, the system operator would need to monitor LSE purchases of power, and the real-time performance of their suppliers, and cut-off LSEs that have purchased insufficient or defaulting supply. We do not believe that system operators are currently basing load shedding on contract positions, real-time purchases by LSEs, and the real-time performance of their suppliers. Even if this were the case, reliability would not be determined by consumers' value of lost load, but by the LSE's value of lost load, which would reflect the likelihood and severity of regulatory punishment—not a market signal from consumers.

To summarize Case 2, the case in which there is an adequacy problem, the reliability level is determined by regulators or market power and not by unregulated competitive energy prices. In fact, industry consultants have argued for years that electricity markets need market power to allow suppliers to cover fixed costs. This approach is common. But so are regulatory adjustments to guide the reliability level. When dispatchers (the engineering arm of the regulator) are worried about adequacy, they find more reasons to allow prices to rise and are more generous in their demand for operating reserves. When they see the system as overly reliable and prices as too high, they find ways to suppress the spot price. If the reliability problem grows more severe, regulatory intervention becomes more open, sometimes resulting in the purchase of capacity. When market power becomes more severe, regulatory intervention is also likely but it works to lower prices and decrease reliability. This is natural and indeed beneficial, but it is not what economics means by a competitive market.

The common "proof" that a "pure-market" design is doing a "pretty good job" of letting the market determine reliability goes as follows. The lights are staying on, and there is no capacity market, therefore "the market" must be choosing "about the right level" of capacity. The flaw in the proof is that it overlooks both market power and regulators. In fact, "pure-market" approaches are simply inadvertent hoaxes. But the hope for an energy-only, pure-market solution, is more deeply flawed than this mechanistic analysis indicates. The concept of an energy-only market solving the reliability problem without selling a reliability product is logically impossible. It suggests that "the market" can do something "fairly well" when logic shows that it cannot do it at all. Here is the argument that "energy-only" competitively determined reliability is simply impossible.

If nothing is known of consumers' utility for reliability, no limits can be put on the level of optimal reliability. Current markets have no access to information concerning how consumers value reliability. In other words, consumers take no market actions that are even partially based on reliability considerations. This is obviously true for consumers who do not have real-time meters and who cannot be individually interrupted. These customers would be foolish to cut back when prices are high, because they cannot receive credit for their efforts, and they would be foolish to pay more for reliable

service since they cannot physically be given more reliability than their neighbor who does not pay for more reliability. Many large consumers do have real-time meters and some can be interrupted. But, system operators are not prepared to black them out based on the performance of their contractual arrangements for power and the performance of their suppliers. If they are never blacked out for contracting for too little power, they will not pay for reliability. So the market receives no signal of how they value reliability. They may well reduce their demand if the spot price rises to say \$1000/MWh, but this tells nothing about the costs they would suffer if they were blacked out. Their price response is a step toward Case 1, in which reliability is perfect, but it provides no information about their VoLL.

In summary, customers currently reveal nothing about the value they place on reliability. Without this information, it makes no sense to think that a market can approximately determine how much reliability they would be willing to pay for. Currently, reliability is determined entirely by market power and regulators. Neither is ideal but at least regulators may attempt to solve the problem. But even unadorned market power is better than one particular "pure-market" design that has become popular in Australia and in the Western US. That design calls for increasing market risk for consumers so that generators can make more money selling them hedges. Although this can increase generator profits by assisting their exercise of market power in forward markets, it has nothing to do with determining optimal reliability, and is certainly not suggested by any economic theory. This approach is an attempt to increase the market power in forward markets as a way of producing the right level of investment and reliability. One cannot get much further from the idea of competitive market efficiency.

2. The basic reliability-options solution

As just explained, if a market has an inherent reliability problem, then it must be solved by regulation or market power. The goal then is to design a regulatory approach that mimics a competitive market as closely as possible. By setting the spot market price cap at VoLL and using options to suppress risk and market power while preserving marginal incentives, regulation can be confined to the single task for which it is needed: determining the adequate level of capacity. Of course reducing market power moves the market closer to the competitive ideal, but what is often overlooked is that reducing risk does too. While, present energy markets are risky because of price spikes, these are caused by the extremely low demand elasticity. If the market did not have demand-side flaws, demand elasticity would be far higher, and risk far lower. Reducing risk moves the market toward this competitive ideal.

We will now describe a simplified design that illustrates the principles involved. Related approaches are described in Bidwell (2005), Chao and Wilson (2004), Cramton and Stoft (2006, 2007), Oren (2005), and Vazquez et al. (2002). However, not all of these include all essential features. The design

is assumed to be implemented in a single isolated market. Later we will discuss the problem of trade. At first, the energy market will be assumed to contain only the real-time market, but later, settlements that include forward markets will be described.

Step one of the design is to set the price cap on the spot market to the best estimate of VoLL. Of course this is a poor estimate, but any value in the €3000–€30,000 range should provide good dispatch incentives. Since, after step two, VoLL will play no role in determining investment, the level is not critical. The difference between energy supplied under a €3000 cap and a €30,000 cap is minuscule, because, with the market power well controlled, these prices rarely are reached and by the time the price reaches €3000, there are not many generators that can provide power but are choosing not too.

Step two of the design is to provide load with a complete hedge in the form of a reliability option with a strike price of, say, €300. Load purchases the target quantity of physical capacity together with a load-following call option at the strike price. In this way, load is 100% hedged from energy prices in excess of €300. If the energy price in an hour is €1000, each supplier has a financial obligation to serve its share of load. Since deviations are priced at the €1000 spot price, the supplier continues to be motivated by the spot price, although both load and suppliers are hedged from price volatility above the €300 price.

The key point is that all generation still faces the spot price even though it is hedged. Suppose a supplier owns 100 MW of capacity. If it provides 80 MW of power for the hour in question and has a 90-MW obligation, it is paid €80,000 because the spot price is €1000, but it must pay $(90 \text{ MW}) \times €(1000 - 300)$ as an option payment. If it provides 90 MW of power, it is paid €90,000 and is obliged to make the same €63,000 option payment. If it produces 100 MW, it is paid €100,000, and again makes the same hedge payment. For every MW it increases or decreases its production, its net revenue increases or decreases by €1000.

Note that when the spot price is €300 or above, it makes sense for virtually every generator to be producing, since marginal cost typically is less than €300. As long as the suppliers produce their share of load, they will earn the strike price for all of their output. In other words, a generator with average performance is fully hedged against spot prices above €300 by its physical generator.⁵ This approach greatly reduces the risk from weather related price fluctuations. It does not, however, reduce performance risk. Although all risk is costly and hence undesirable, performance risk cannot be eliminated

⁵ A generator with a marginal cost of €310 will lose €10/MWh when it runs, but will still make money from the auction payment for its capacity. It will be fully hedged, and will still be motivated to run by the full spot price. If its marginal cost is extremely high, it may find it more profitable to sell a smaller amount of capacity, so it can exceed its share of production and earn the full spot price. These difficulties should be small, and there should be little generation in this category.

without eliminating the performance incentive. Generators will argue for a long list of exceptions due to “acts of God,” but virtually all of these are to some extent under the generators’ control, and since performance risk is quite small, it is best to simply ignore these arguments. In the full-price market of step one, no one would dare argue for exceptions, even though their excuses would be equally valid (or equally invalid) in a pure-market setting. In a pure energy market, no matter why a generator is out of service when the price is €10,000/MWh, no one would consider paying it for power it did not produce due to an “act of God.” This obvious fact has proven extremely difficult for suppliers to remember during capacity-market design negotiations.

To summarize our progress with the design, the first step (price cap equals VoLL) assured that the market conforms to a classic competitive design, although it still includes the currently unavoidable demand-side flaws (which lead to the absence of a robust demand response), and the compensating intervention of a VoLL price cap. As is well known, with this cap, and ignoring risk and market power, the market will provide both optimal dispatch incentives and optimal investment incentives. Of course the regulator is in control of investment incentives, and can induce any level of reliability. But if the regulator sets the cap to VoLL (which cannot be determined by the market), investment will be optimal.

Step two (the complete hedge) preserves the dispatch incentives perfectly on the generation side, but destroys the investment incentive, because although generators face the spot price on the margin, their revenues are limited exactly as if there were a price cap at €300/MWh. Why is this progress? Two advantages are obvious: a dramatic reduction in risk and in market power. The third advantage is that VoLL has become much less important. It no longer affects investment. If a higher VoLL is used, the price cap will be higher and the incentive to perform on peak will be greater, but because of the hedge, this will not increase the earnings of generators. This means changing VoLL does not change investment or adequacy. VoLL is only linked to real-time performance. In this role, VoLL makes little difference for the following reason. If VoLL and the price cap are €2000, all functioning generators will produce as much as they can, and if VoLL and the cap are €20,000, they will do almost exactly the same thing. Perhaps, at the higher value, some load with real-time pricing will be reduced, or some generator will squeeze out one more MW, but little will change. Because the estimation of VoLL is always controversial, this is an advantage.

Step three, the final step of the design, introduces the capacity auction. This sets the payments to generators for providing reliability options just high enough to induce optimal investment and adequate capacity. An annual auction is used to purchase new capacity up to the level required for reliability. These auctions determine the price of reliability options that is just sufficient to induce the required new investment. For example, a strike price of €300/MWh might result in the average annual loss of €40,000 of revenue per MW of capacity relative to the spot market of step one. This is often termed the “missing money.” In this case, new entrants will

bid the price of reliability options down to €40,000/MW-year. If the cost of constructing new capacity increases or decreases, due to environmental restrictions or new technology, new entrants will bid just enough higher or lower to maintain a normal rate of return.

The result is that the regulator controls the level of capacity, but the market controls the price of capacity and the type and quality of capacity built. Hence the regulatory intervention has been strictly limited to the determination of the one factor about which the market has little information—the adequate level of capacity.

Although the auction design requires care to address the potential exercise of market power, the following simple procedure would work quite well. Every September an auction is held for reliability options (ROs), which take effect on January 1, just over three years in the future. Existing generators may choose either to enter the auction with a zero bid, or not to sell ROs. New projects are allowed to bid without restriction. The regulator bids a demand curve that intersects the target adequacy level at the most recent RO price and slopes down to the right by 5% in price for each 1% increase in capacity. It slopes up from the same point by 20% for each 1% decrease in capacity. The auction is held using a descending clock procedure (see Ausubel and Cramton, 2004). All accepted bids are paid the clearing price, but existing generation receives one-year contracts while new generation may choose any contract length up to seven years. Once a new generator's initial contract expires it becomes an "existing" generator. If no new generation is purchased in a given year, all existing generators that bid have their contracts extended for one year.

There is one further rule, which assures full hedging and limits market power in both the spot and capacity markets. Any generator that does not sell ROs for its full capacity receives only the spot price capped at the strike price. In other words, such generators, in effect, provide the hedge without compensation. In fact a few extremely unreliable generators may opt not to sell capacity, and others will decide not to because they are selling their capacity into another system. These are good reasons and cause no trouble. Neither is only withholding capacity to exercise market power. However, withholding to exercise market power is discouraged.

Some designs omit this rule. Objections to it are based on ideological grounds. We find these unpersuasive and prefer to rely on economic analysis. As explained by Hogan and Harvey (2000) during the California crisis, suppliers will not give up their market power for free by entering into long-term contracts. They realize that their market power is valuable, and will extract approximately the value of this market power before relinquishing it by selling ROs. This rule prevents the exercise of market power in the capacity market, as well as the spot market. Only permanent retirements can reduce the quantity of existing supply.

One puzzle is why the spot price would ever exceed the option's strike price if all load and generation is fully hedged. Remember that suppliers only hedge a specific quantity of power—their share of load. Power produced beyond this level is unhedged for the supplier. Not only can some suppliers

produce more than their share, but some always will, simply because it is impossible for a hundred or more suppliers to all supply exactly their shares. In any case, the spot price can be high, depending on market rules, because the system operator bids a high price or because suppliers bid a high price for supply that may be beyond their market share. In fact, except for the reduction in market power, prices will be just as high as in an unhedged market.

3. How reliability options work in practice

There have been concerns that the mandatory RO system requires a centralized day-ahead market such as found in the Northeastern ISOs of the US (PJM, ISO-NE, and NYISO). Fortunately these concerns are unfounded. There is also concern that the RO mechanism is incompatible with a high level of bilateral contracting. In fact there always has been a high level of bilateral contracting in New England and other markets where this design has been adopted. How ROs integrate with bilateral markets and trade with other markets are explained below.

Assume a 100-MW generator has sold a reliability option with a strike price of €300, and the Transmission System Operator is handling the reliability market. There are also long-term bilateral markets, a day-ahead market run by APX, and a balancing market run by the TSO.

Currently, European markets can be thought of as two settlement systems. A generator sells total energy of Q_{Forward} in forward markets, either bilateral or centralized like APX, and Q_{Forward} is scheduled with the TSO. The generator then delivers an amount Q_{RT} in real time. The forward quantity is settled at privately determined prices regardless of what is delivered in real time, while the deviation from the forward sale, $Q_{\text{RT}} - Q_{\text{Forward}}$, is paid the balancing market price, P_{Balance} . Of course the deviation can be negative, in which case the generator pays the TSO. So without reliability options or when the balancing price is below the strike price, the settlement works as follows:

$$\text{Generator Revenue} = P_{\text{Forward}} \times Q_{\text{Forward}} + P_{\text{Balance}} \times (Q_{\text{RT}} - Q_{\text{Forward}})$$

With reliability options, and the balancing price above the strike price, the RO is settled second in what is essentially a three-settlement system. Each supplier is responsible for a share of the real-time load that is proportional to the quantity of reliability options it has sold. For example, if it has sold 100 MW of options out of a total of 10,000 MW of reliability options, it is responsible for 1% of the load in every hour.⁶ Call

⁶ Why base the hedge on load share? New England has about 30 GW of capacity, but sometimes, due to cold weather, many generators cannot run, and the price has spiked with as little as 20 GW of load. If reliability options covered the full 30 GW, then load would be paid for 30 GW times the \$1000 spot price less the \$300 strike price. Hence load would profit by \$7 million dollars per hour during such an incident. This needlessly upsets generators, and causes them to worry that extra capacity will be purchased so load can profit more in this way. Basing the reliability option on load share solves this problem by putting the generators in a nearly balanced position in every hour.

the supplier's share Q_{Share} . The three-settlement system works as follows:

$$\text{Generator Revenue} = P_{\text{Forward}} \times Q_{\text{Forward}} + P_{\text{Strike}} \times (Q_{\text{Share}} - Q_{\text{Forward}}) + P_{\text{Balance}}(Q_{\text{RT}} - Q_{\text{Share}})$$

If the generator supplies exactly its load share, so that $Q_{\text{RT}} = Q_{\text{Share}}$, then it is fully hedged against the balancing price. However, if it deviates either up or down from its share, it is paid or must pay the balancing price. Hence its incentive to perform has not changed, and the balancing market continues to play its traditional role. As can be seen, there is not much interaction between the forward transactions and the reliability option. In particular there is no reason generators cannot sell all of their power in the forward markets and sell reliability options for all of their capacity in the capacity market. Three things have changed as a result of the hedge built into the reliability option: the average generator earns at most the strike price in the balancing market, both load and generation are less at risk, and the forward contract needs only to cover prices below the strike price.

Notice that the sum over all generators of Q_{RT} , the total power delivered, equals the total load, which equals the sum over all generators of Q_{Share} . This means the sum over all generators of $P_{\text{Balance}}(Q_{\text{RT}} - Q_{\text{Share}})$ is zero. This is exactly true, because generation shares are determined after the fact. The summed terms are the incentive payments for under and over performance by generators, relative to Q_{Share} . So these payments do not affect load, but are simply payments from poorly performing generators to the better performers. As a consequence of this fact, generators do not have to be concerned that the TSO will under-procure capacity because, on average, this cannot change the total performance payments to generators which always sum to zero. Under-procurement of capacity will only increase the number of hours when there is a shortage and generators are paid the strike price.

Another concern is that European markets are open to trade between countries. If the price is high in Germany, but Dutch generators are under reliability options, they might leave the Dutch market where, on average, they can earn only the strike price, and sell as much as they can into the German market. To examine this possibility, the settlement must include a term for exports. Since the concern is with the effect of reliability options, the real-time price must be above the strike price, so it is safe to assume that a generator will produce the most it is capable of, Q_{Max} . Next note that only the power delivered to the domestic balancing market, Q_{Domestic} , receives the balancing price. The export quantity is then, $Q_{\text{Max}} - Q_{\text{Domestic}}$, and the settlement works as follows.

$$\text{Generator Revenue} = P_{\text{Forward}} \times Q_{\text{Forward}} + P_{\text{Strike}}(Q_{\text{Share}} - Q_{\text{Forward}}) + P_{\text{Balance}}(Q_{\text{Domestic}} - Q_{\text{Share}}) + P_{\text{Export}}(Q_{\text{Max}} - Q_{\text{Domestic}})$$

Because the derivative of Generator Revenue with respect to Q_{Domestic} is $(P_{\text{Balance}} - P_{\text{Export}})$, the incentive to export is

exactly the same as without reliability options. Hence there is no justifiable concern with a disruption of the balance of electricity trade, or a collapse of the domestic market.

Another practical concern is the assignment to load serving entities (LSEs) of responsibility for the cost of ROs. Because the options are procured by the TSO, the LSEs are not burdened with purchasing ROs and need not make any long-term commitments by purchasing them. This is a great advantage because it means that LSEs are at little risk from consumers moving from one LSE to another. The cost assignment is simply adjusted each year and based on the LSE's coincident peak load during the year. This determination is best made after the fact. For example, option cost responsibility for 2010 should be based on the peak loads during 2010. To reduce randomness in loads on any given day, a weighted average of the three highest peak-load days could be used.

Another practical concern is that load is not exposed to the full spot price because of the hedge. Although most load is still not on real-time meters, these are becoming more prevalent. Moreover, LSEs can implement various programs to encourage conservation during times of peak load, and it would be worthwhile to properly motivate them to do so. This can be accomplished as follows. First compute the "peak energy costs" of each LSE. This is simply the integral of $\{\text{its load times } \text{Max}\{0, (P_{\text{Balance}} - P_{\text{Strike}})\}\}$ over the year. Then, since each LSE is assigned a reliability share, L_{Share} , based on its coincident peak load, this share can be used to compute its share of the total of all peak energy costs. Each LSE then pays a penalty equal to the amount by which its actual peak energy cost exceeds its share of total peak energy costs.

$$\text{Peak-Load Penalty} = \text{peak energy cost} - L_{\text{share}} \times (\text{Sum of all peak energy costs})$$

The sum of the penalties is zero, and the derivative of a load's penalty with respect to its own peak energy cost is $(1 - L_{\text{Share}})$, which is near one as long as LSEs are small. This means that purchasing an MW of power, when the balancing price is above the strike price, costs each LSE an amount $P_{\text{Strike}} + (1 - L_{\text{Share}})(P_{\text{Balance}} - P_{\text{Strike}})$, which is very nearly P_{Balance} . In other words the penalty makes the LSEs face the balancing price on the margin even though they pay no more on average. There will be some risk to loads from this performance penalty but it is small and it is only what is inevitable if loads are to face the real-time price on the margin. They are still completely hedged against price spikes caused by weather, nuclear outages, or other events out of their control. This same technique can be used by LSEs to pass real-time price signals through to their loads equipped with real-time meters.

One final concern is that reliability options may impose burdensome new information requirements. However, the

TSOs are already aware of the quantities transacted in the forward markets because these must be scheduled. They are also aware of actual production and of which generators are exporting how much power. Because the TSO will conduct the capacity auctions, it will know who owns the reliability options. It also knows the daily loads of the LSEs, and the balancing market prices. This is all the information needed to implement this reliability option design. In particular there is no need to collect more information about bilateral transactions.

4. Reliability options in a hydro-dominated system

If supply is mainly from hydro-electric generation, the limiting factor is not likely to be capacity (the ability to provide energy in peak hours) but rather, firm energy (the ability to provide energy in dry periods).⁷ As a consequence the TSO will need to purchase firm energy options. Just as capacity is the physical basis for reliability options, so firm energy options have a physical basis that involves a longer-term supply of energy. Firm energy is defined as the amount of energy a generator can deliver per month during an exceptionally dry period (a worst-case benchmark). A typical thermal unit is certified at its nameplate capacity times its average availability, such as 92% of nameplate. In contrast, a hydro resource's firm energy may be well below its nameplate, say 35%, due to a limited water reservoir and low inflows during dry periods. A new unit's firm energy contribution is how much less energy the system would have without the unit in the worst-case benchmark.

The second difference in a firm energy market is the way load share is defined when the real-time price rises above the strike price and the option comes into play. The total quantity of the call option follows load, and is divided between thermal and hydro resources as follows. The load obligation is first divided among the thermal resources in proportion to their firm energy ratings. Once load exceeds the total firm energy ratings of thermal capacity, the excess load is divided among the hydro generators in proportion to their firm energy ratings. Economically, this division of load between thermal and hydro generators makes sense, since it is consistent with efficient dispatch of the units. During scarcity hours the hydro opportunity cost sets the price and thermal resources run at capacity if available. One might think that this approach is biased against hydro resources, since hydro is asked to do more load following than the thermal resources. Indeed, there would be a bias if the ability to follow load was scarce; however, since there is a surplus of capacity, load following is not costly and the capability is efficiently priced at zero.

Given these definitions, settlement proceeds exactly as it does in the reliability-options market just described.

In a market in which it is unclear whether capacity or firm energy is the scarce resource to achieve generation adequacy, the approach can be enhanced to accommodate both products.

⁷ This applies to stored water. Run of river can be handled like wind in a capacity market.

Each resource is rated for both its capacity and its firm energy, and offers the two products as a package. There is a target for each product. The descending clock auction then has two prices, one for capacity and one for firm energy. The prices descended until a supply and demand balance has been achieved in both products. In a hydro-dominated system, this will imply a zero price for capacity and a positive price for firm energy. In a thermal-dominated system, it will imply a positive price for capacity and a zero price for firm energy.

5. Implementation in New England and Colombia

Reliability markets based on this reliability options approach have been adopted in New England's thermal-dominated market (Cramton, 2006) and Colombia's hydro-dominated market (Cramton and Stoft, 2007). The markets are currently in a transition period in which the capacity and firm energy prices are set administratively. Both markets will have their first auctions in early 2008. Both markets differ somewhat from the specific approach described above.

New England has locational pricing, so it was important to implement the capacity market on a zonal basis, so that adequacy would be achieved throughout New England. Another important difference is the way that market power is addressed in the capacity market is more complex than what we have suggested above.

Colombia's firm energy market is quite similar to what we describe above. The strike price of the reliability option is only about €100/MWh. This is still above the marginal cost of nearly all generation, which means that suppliers have a physical hedge to protect against prices above the strike price. Another change is a longer planning period of four years, rather than three to accommodate longer-lead time hydro projects. Indeed, to accommodate large hydro products, suppliers can sell firm energy up to seven years ahead at the four-year-ahead price, subject to some restrictions. The large hydro project is a price taker in the auction, since it is not selling firm energy four years ahead, but up to seven years ahead. At the conclusion of the auction, the investor specifies the fraction of the firm energy from the project it desires to lock-in at the four-year ahead auction price. Only a fraction of the predicted need for new generation can be sold in this way.

6. Risk analysis based on market data

In both New England and Colombia, we have examined the supplier risk associated with the reliability options approach. We have examined company risk using historical market data as well as detailed market simulations carefully calibrated to the specific market setting.

The findings in Colombia (Cramton and Stoft, 2007), where we did more detailed simulations, are as follows.

- Lumpy investment means that few new units are added each year. Indeed, in 27% of the years no new entry occurs. (This is an overestimate to the extent that the size

of proposed projects reflect the actual need in the year, as one might expect.)

- The mandatory hedge is remarkably successful in reducing risk. In the benchmark case, where we assume demand has constant elasticity of -0.05 for prices above the strike price (a 20% increase in price produces a 1% decline in demand), the hedge reduces aggregate profit risk by a factor of 7. More importantly, the hedge reduces company risk by a factor of 4.5 in the benchmark case. Even when we assume a high level of demand response so that prices remain low during scarcity periods and there is less profit risk to start with, the hedge reduces company risk by 55%.
- A higher strike price increases risk. Increasing the strike price shifts the profit distribution toward the no hedge case (a strike price of infinity). This results in a large increase in energy rent risk and a small decrease in hedge payment risk. The overall impact is a large increase in profit risk.

Taken together, the simulation results demonstrate the risk reducing benefits of the reliability market. Provided there is competitive new entry in response to load growth, the reliability market should work well at coordinating investment in new supply, while minimizing supplier and consumer risks.

7. Conclusion

The forward reliability market approach described here is the product of a systematic development based on clear economic principles. Indeed, once it is understood it seems almost inevitable. The first step is the classic VoLL pricing system which is known to be both short- and long-run optimal if VoLL is known and risk and market power are assumed to be costlessly suppressed by unspecified means. The second step is to suppress risk and market power, the two evils of VoLL pricing, by introducing reliability options. These do not interfere with real-time price signals, so on the second step, the market retains the dispatch optimality of classic VoLL pricing, but without the problems.

The second step suppresses a significant amount of generation revenue and thereby destroys the investment incentive, so the third step restores it by introducing a capacity market, which induces an adequate level of investment by procuring the appropriate quantity of reliability options. All generation, new and existing will want to sell reliability options for their full capacity because these options fetch a high price relative to the financial cost of the option. Participation in this market is guaranteed by the rule that non-participating generators receive the spot price capped at the option strike price.

The auction for procuring reliability options takes place three years in advance of the effective date, so that there is time for new entry to back the options. To suppress market power in the reliability-option auction, only new capacity bids are allowed to set the price. The auction is a descending clock auction. By inducing investment with an auction instead of high prices, not only does the regulator have better control of the average reliability level, but far better investment coordination is assured.

Investment coordination prevents the boom-bust cycles which increase both investor risk and reliability risk for consumers. Without an auction, as the market tightens it offers an increasingly large prize for the next entrant. However, entry is a secretive process, and so simultaneous entry is possible. Aware of this, investors are torn between holding off until the prize is large enough to support some simultaneous entry and entering quickly to ward off competition. The optimal strategy is mixed and the outcome chaotic. With an auction, the TSO coordinates entry without reducing competition. There are multiple simultaneous bids, but the TSO selects only as many as needed. This stability benefits both consumers and investors.

While there are many concerns about the use of reliability options, with the designs specified here, none prove warranted. There is no difficulty in deciding which private contracts are acceptable substitutes for reliability options, because no substitutes are accepted. One hundred percent coverage by reliability options does not interfere with 100% coverage with bilateral contracts. Reliability options provide price coverage above the strike price; bilateral contracts provide price coverage below the strike price. Although reliability options limit the average real-time price to the strike price, the marginal price for both load and investors remains the balancing market prices. This preserves incentives and prevents any increase in exports relative to the present system even on days when other countries have high prices. The benefits of this design are significant. The design minimizes risk and market power, while coordinating efficient entry.

Uncited reference

Cramton et al. (2006).

References

- Pérez-Arriaga, Ignacio J., 1999. Reliability in the new market structure. In: IEEE PES 1999 Summer Meeting, Plenary session.
- Ausubel, Lawrence M., Cramton, Peter, April–May 2004. Auctioning many divisible goods. *Journal of the European Economic Association* 2, 480–493.
- Bidwell, Miles, June 2005. Reliability options. *Electricity Journal*.
- Chao, Hung-po, Wilson, Robert, 2004. Resource Adequacy and Market Power Mitigation via Option Contracts. EPRI, Palo Alto, CA.
- Cramton, Peter, 2006. New England's forward capacity auction, Working paper, University of Maryland.
- Cramton, Peter, Stoft, Steven, March 2006. The convergence of market designs for adequate generating capacity, White Paper for the California Electricity Oversight Board.
- Cramton, Peter, Stoft, Steven, 2007. Colombia Firm Energy Market. In: Proceedings of the Hawaii International Conference on System Sciences.
- Cramton, Peter, Stoft Steven, West Jeffrey, 2006. Simulation of the Colombian firm energy market, Working paper, University of Maryland.
- Hogan, William W., Harvey, Scott M., 2000. California electricity prices and forward market hedging, Working paper, Harvard University.
- Oren, Shmuel S., November 2005. Generation adequacy via call option obligations: safe passage to the promised land. *Electricity Journal*.
- Vazquez, Carlos, River, Michel, Perez Arriaga, Ignacio, 2002. A market approach to long-term security of supply. *IEEE Transactions on Power Systems* 17 (2), 349–357.