

# An MDP-Based Winning Approach to Autonomous Power Trading: Formalization and Empirical Analysis

Daniel Urieli and Peter Stone

Dept. of Computer Science  
The University of Texas at Austin  
Austin, TX 78712 USA  
{urieli,pstone}@cs.utexas.edu

## Abstract

With the efforts of moving to sustainable and reliable energy supply, electricity markets are undergoing far-reaching changes. Due to the high-cost of failure in the real-world, it is important to test new market structures in simulation. This is the focus of the Power Trading Agent Competition (Power TAC), which proposes autonomous electricity broker agents as a means for stabilizing the electricity grid. This paper focuses on the question: how should an autonomous electricity broker agent act in competitive electricity markets to maximize its profit. We formalize the complete electricity trading problem as a continuous, high-dimensional Markov Decision Process (MDP), which is computationally intractable to solve. Our formalization provides a guideline for approximating the MDP's solution, and for extending existing solutions. We show that a previously champion broker can be viewed as approximating the solution using a lookahead policy. We present TacTex'15, which improves upon this previous approximation and achieves state-of-the-art performance in competitions and controlled experiments. Using thousands of experiments against 2015 finalist brokers, we analyze TacTex'15's performance and the reasons for its success. We find that lookahead policies can be effective, but their performance can be sensitive to errors in the transition function prediction, specifically demand-prediction.

## 1 Introduction

With the efforts of moving to sustainable and reliable energy supply, electricity markets (aka power markets) are undergoing far-reaching changes: customers are being engaged in power markets, to incentivise flexible demand that adapts to supply conditions (U.S 2003); and wholesale markets are being deregulated and opened to competition (Joskow 2008). In principle, deregulation can increase efficiency. However, in practice the California energy crisis (2001) has demonstrated the high-costs of failure due to flawed deregulation (Stoft 2002; Borenstein 2002), and the importance of testing new market structures in simulation before deploying them. This is the focus of the Power Trading Agent Competition (Power TAC) (Ketter, Peters, and Collins 2013).

In Power TAC, autonomous broker agents compete with each other to make profits in a realistic, detailed smart-grid

simulator with wholesale, retail and balancing power markets, and about 57,000 customers. The stability of electricity grids critically depends on having balanced electricity supply and demand at all times. Broker agents are financially incentivised to maintain supply-demand balance in their portfolio and thus contribute to grid stability. It is likely that autonomous broker agents will be employed in power markets, due to the complexity of the electricity trading domain. The decision-making challenges of such brokers has been under study in the autonomous agents community, but either under limited scope, or limited competitiveness and comparability (Ketter, Peters, and Collins 2013). This paper focuses on the question: how should an autonomous broker agent act in competitive power markets to maximize its profits? we advance the state of the art in the following ways:

- This paper is the first to formalize the *complete* broker's power trading problem. We formalize the problem as a Markov Decision Process (MDP) which, due to its continuous high-dimensional state and action spaces, cannot be solved exactly in practice. Our formalization compactly captures the challenges faced by a broker, and provides a guideline for approximating the solution and for extending existing solutions. While our formalization is based on the Power TAC simulator, we expect it to generalize and be useful in reality, since Power TAC closely models real-world markets.
- We present TacTex'15, which is by many metrics the best Power TAC broker at the current time. TacTex'15 improves upon a previous strategy that can be viewed as an approximate solution to the MDP, using a similar architecture with three strategic improvements. The strategic improvements may seem minor on the surface but result in large performance improvements.
- Using thousands of experiments, we analyze the performance of TacTex'15, and the reasons for its success.

## 2 Power TAC Game Description

Power TAC is an annual competition in which the competitors are autonomous brokers programmed by teams from around the world. The competition includes hundreds of games and takes several days to complete. In a game, the Power TAC simulator runs on a central server, while competing brokers run remotely and communicate with the

server through the internet. Each broker receives partial state information from the server, and responds by communicating the actions it takes. The competition includes different game sizes, ranging from small to large number of competitors. Participants release their broker binaries after the competition, and use them to run controlled experiments.

Power TAC uses a high-fidelity power markets simulator, modeling a smart-grid with more than 57,000 simulated customers (50,000 consumers and 7,000 renewable producers). Power TAC’s customers are *autonomous agents* that optimize the electricity-costs and convenience of their human-owners (Reddy and Veloso 2012). Customers model commercial/residential buildings and solar/wind farms which consume/produce using time-series generators constructed from real-world data, according to weather and calendar factors. The simulation proceeds in 1-hour timeslots for 60 simulated days and takes 2 hours to complete.

In Power TAC, autonomous broker agents compete by acting in three markets: (1) a *wholesale market*, where generation companies sell energy, and brokers procure energy to be delivered in the following 24 hours (or sell surplus) in sequences of 24-double auctions, (2) a *tariff market*, which is a retail market where energy is traded with consumers and distributed renewable energy producers, and (3) a *balancing market*, which ensures that supply and demand are balanced at all times and determines broker imbalance fees.

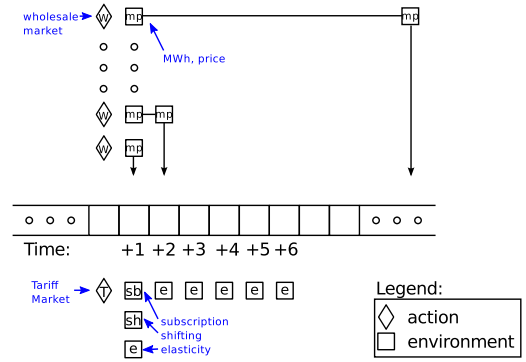
The brokers compete to gain market share and maximize profit by trading electricity. In the tariff market, brokers publish tariff contracts for energy consumption/production. Tariffs may include fixed and varying prices and possibly bonuses/fees. Customers stochastically subscribe to tariffs which maximize their utility (cost-saving and comfort). Customers are equipped with smart-meters, so consumption and production are reported to the broker every hour. Brokers typically balance their portfolio’s net demand by buying in the wholesale market. Full details can be found in the Power TAC Game Specification (Ketter et al. 2015).

### 3 The Broker’s Power Trading Problem

This section formalizes the broker’s power trading problem. Our formalization compactly captures the complex challenges faced by a broker, and provides a guideline for approximating the solution and for extending existing solutions. While our formalization is based on the Power TAC simulator, we expect it to generalize and be useful in reality, since Power TAC closely models real-world markets. We start with an intuitive problem description and continue to our formalization.

Figure 1 illustrates the temporal structure of a broker’s power trading problem. The temporal structure of the tariff and wholesale market actions differ in multiple ways. Tariffs specify energy for *immediate* and *repeated* delivery and are published at *low-frequency* (every one or more days). Wholesale bids typically specify energy for *future, one-time* delivery and are executed at *high-frequency* (every hour).

**Power Trading as an MDP.** Given the internal states of the simulator and competing brokers, the broker’s energy trading problem is a Markov Decision Process (MDP) (Puterman 1994). However, since competitors’ state and parts of



**Figure 1: Temporal structure of the power trading problem.** Time progresses to the right; the notation ‘+i’ stands for ‘i timeslots into the future’. Diamonds stand for broker actions. Squares stand for simulation environment responses. The top part represents the wholesale market: a broker submits limit orders to buy/sell energy for the next 24 hours, then it receives the results of the 24 double-auctions. The bottom part represents the tariff market: a broker may publish one or more tariffs (once every 6 hours), and customers respond by potentially (1) subscribing to new tariffs, (2) shifting consumption to cheaper times, and (3) elastically adapting total consumption based on price.

the simulator state are unobservable, the trading problem is actually a Partially Observable MDP (POMDP). Nevertheless, for computational tractability and modeling clarity, we approximate the trading problem as an MDP, as follows (denoting  $B_0$  as the acting broker):

- **States:**  $S$  is the set of states, where state  $s$  is a tuple  $\langle t, \mathcal{B}, \mathcal{C}, \mathcal{P}, \mathcal{T}, \mathcal{S}_{B_0}, \mathcal{Q}_{B_0}, \mathcal{A}_{B_0}, I_{B_0}, \mathcal{W}, cash_{B_0}, \rho \rangle$  that includes the current time  $t$  (which encapsulates week-day/hour), and the sets: competing broker identities  $\mathcal{B}$ ; identities of consumers  $\mathcal{C}$  and producers  $\mathcal{P}$  (both referred to as *customers*); published tariffs of all brokers  $\mathcal{T} := \cup_{B \in \mathcal{B}} \mathcal{T}_B$ ; customer subscriptions to  $B_0$ ’s tariffs  $\mathcal{S}_{B_0}$ ; current energy consumption/production of  $B_0$ ’s customers  $\mathcal{Q}_{B_0}$ ; recent auction results  $\mathcal{A}_{B_0} := \{ \langle p^c, q^c, \mathcal{O}^c_{B_0}, \mathcal{O}^u, \mathcal{M}_{B_0} \rangle_j \}_{j=t+1}^{t+24}$  including, for each of the following 24 timeslots, the clearing price  $p^c$  and total quantity  $q^c$ ,  $B_0$ ’s cleared orders  $\mathcal{O}^c_{B_0}$ , all brokers’ uncleared orders  $\mathcal{O}^u$ , and  $B_0$ ’s *market-positions*  $\mathcal{M}_{B_0}$  (energy deliveries and charges, updated incrementally from  $\mathcal{O}^c_{B_0}$ );  $B_0$ ’s energy imbalance  $I_{B_0}$ ; current weather and forecast  $\mathcal{W}$ ;  $B_0$ ’s cash balance  $cash_{B_0}$ ; and randomly sampled game-parameters (such as fees and game length)  $\rho$ . **Note:** the underlying state of the game, which includes elements unobserved by the broker, is the tuple  $\langle t, \mathcal{B}^o, \mathcal{G}^o, \mathcal{C}^o, \mathcal{P}^o, \mathcal{T}, \mathcal{S}, \mathcal{Q}, \mathcal{A}, \mathcal{I}, \mathcal{W}, cash, \rho \rangle$  where  $\mathcal{B}^o, \mathcal{G}^o, \mathcal{C}^o, \mathcal{P}^o$  are the identities and internal states of brokers, generation companies, consumers and producers, respectively; and where  $\mathcal{S} := \cup_{B \in \mathcal{B}} \mathcal{S}_B$ ,  $\mathcal{Q} := \cup_{B \in \mathcal{B}} \mathcal{Q}_B$ ,  $\mathcal{A} := \cup_{B \in \mathcal{B}} \mathcal{A}_B$ ,  $\mathcal{I} := \{I_B\}_{B \in \mathcal{B}}$ ,  $cash := \{cash_B\}_{B \in \mathcal{B}}$ .
- **Actions:** A broker’s set of actions  $A := A^\tau \cup A^\omega \cup A^\beta$  is composed of tariff market actions  $A^\tau$ , wholesale market actions  $A^\omega$  and balancing market actions  $A^\beta$ , as follows.

1. *Tariff market actions*  $A^\tau$ : create/modify/revoke tariffs. A tariff is a tuple  $T = \langle type, rates, fees \rangle$  where:
    - $type \in \{\text{consumption, production, ...}\}$  can be general (e.g. production) or specific (e.g. solar-production).
    - $rates$ : a set of rates, each specifying price/kWh and times, and/or usage thresholds where it applies.
    - $fees$ : optional periodic/signup/withdraw payments.
  2. *Wholesale market actions*  $A^\omega$ : submit limit orders of the form  $\langle energyAmount, limitPrice, targetTime \rangle$  to buy/sell energy for one of the next 24 hours.
  3. *Balancing market actions*  $A^\beta$ : submit customers energy curtailment requests (currently unused).
- **Transition Function:** The transition function is partially deterministic and partially stochastic, as follows. The time  $t$  is incremented by 1 hour;  $\mathcal{B}, \mathcal{C}, \mathcal{P}$  remain unchanged;  $\mathcal{T}$  is updated by create/modify/revoke tariff actions, deterministically by  $B_0$ , and stochastically (due to unobservability) by other brokers;  $\mathcal{S}_{B_0}$  is updated stochastically based on customers' decisions;  $\mathcal{Q}_{B_0}$  is determined stochastically based on weather and customers' internal states (shifting and elasticity, see Figure 1);  $\mathcal{A}_{B_0}$  is updated with auction results, stochastically since (i) competitors rely on stochastic information (demand predictions), (ii) competitors' internal states are hidden, and (iii) generation companies bid stochastically;  $I_{B_0}$  is a deterministic function of  $\mathcal{T}_{B_0}, \mathcal{S}_{B_0}, \mathcal{Q}_{B_0}, \mathcal{A}_{B_0}$ ;  $\mathcal{W}$  is stochastic; *cash* is updated deterministically from the recent stochastic reward; and  $\rho$  remains unchanged.
  - **Reward:** Let  $s_t, r_t, a_t$  be the state, reward, and broker-action(s) at time  $t$ . Let  $r^\tau, r^\omega, r^\beta$  be the broker's energy buy/sell payments in the tariff, wholesale, and balancing markets respectively. Let *dist* be the energy distribution fees, and *fees* the tariff-market fees. The reward at time  $t$  can be characterized by the following function.

$$\begin{aligned}
 r_t(s_{t-1}, a_{t-1}, s_t) &:= r^\tau(s_t) + r^\omega(s_t) + r^\beta(s_t) \\
 &+ \text{dist}(s_t) + \text{fees}(s_{t-1}, a_{t-1}, s_t) := \\
 &\underbrace{Q_t^{\text{cons}} p_t^{\text{cons}} - Q_t^{\text{prod}} p_t^{\text{prod}}}_{r^\tau(s_t)} + \underbrace{Q_t^{\text{ask}} p_t^{\text{ask}} - Q_t^{\text{bid}} p_t^{\text{bid}}}_{r^\omega(s_t)} \\
 &\underbrace{\pm \text{bal}(I_{B_0,t})}_{r^\beta(s_t)} - \underbrace{\max(Q_t^{\text{cons}}, Q_t^{\text{prod}}) \times \text{distFee}}_{\text{dist}(s_t)} \\
 &\underbrace{-\text{pub}(a_{t-1}) - \text{rev}(a_{t-1}) \pm \text{psw}(\mathcal{S}_{B_0,t-1}, \mathcal{S}_{B_0,t})}_{\text{fees}(s_{t-1}, a_{t-1}, s_t)} \quad (1)
 \end{aligned}$$

where  $\pm$  denotes components that can be positive or negative;  $Q_t^{\text{cons}}, Q_t^{\text{prod}}$  are the total consumed/produced quantities by  $B_0$ 's customers in the tariff-market (both are sums of entries of  $\mathcal{Q}_{B_0}$ );  $Q_t^{\text{ask}}, Q_t^{\text{bid}}$  are the amounts  $B_0$  sold/procured in the wholesale-market (both are sums of elements of  $\mathcal{M}_{B_0}$  inside  $\mathcal{A}_{B_0}$ );  $p_t^{\text{cons}}, p_t^{\text{prod}}, p_t^{\text{ask}}, p_t^{\text{bid}}$  are the average buying/selling prices (determined by  $\mathcal{T}_{B_0}, \mathcal{S}_{B_0}, \mathcal{Q}_{B_0}$  and  $\mathcal{M}_{B_0}$ );  $\text{bal}(I_{B_0,t})$  is the fee for imbalance  $I_{B_0,t} = Q_t^{\text{cons}} - Q_t^{\text{prod}} + Q_t^{\text{ask}} - Q_t^{\text{bid}}$  (which depends on unobserved other broker imbalances  $I \setminus I_{B_0,t}$ ); *distFee* is a fixed fee per kWh transferred over the grid; *pub, rev*

are tariff publication and revoke fees; *psw* are tariff periodic/signup/withdraw fees/bonuses.

- **Discount Factor:**  $\gamma$  reflects daily interest on cash balance.

## 4 The TacTex'15 Broker Agent

This section characterizes approximate MDP solutions, and describes TacTex'15's approximate solution.

### 4.1 Approximate MDP Solutions

The MDP's solution is an optimal power-trading policy (a mapping from states to actions). There are two problems to solve the MDP exactly: first, the high-dimensional states and actions and the complex reward makes it computationally intractable, and second, some components of the transition and reward functions are unknown to the broker. Therefore, brokers necessarily can only approximate the solution. There are four categories of approximate solutions to large MDPs, of which *lookahead policies* seem suitable for our domain, since they are effective in time-varying settings, where it is unclear how to approximate a value function or find a simple rule that maps states to actions (Powell and Meisel 2015).

Lookahead policies are partial MDP solutions that optimize over simulated trajectories  $s_t, r_t, a_t, s_{t+1}, r_{t+1}, a_{t+1}, \dots$  using generative models that predict *action effects* (next state and reward). Here, the reward is a deterministic function of  $s_{t-1}, a_{t-1}, s_t$  except for the  $\text{bal}(I_{B_0,t})$  component. Therefore a broker needs generative models for  $\text{bal}(I_{B_0,t})$ , for  $\mathcal{T} \setminus \mathcal{T}_{B_0}, \mathcal{S}_{B_0}, \mathcal{Q}_{B_0}$  (to predict  $Q_t^{\text{cons}}, p_t^{\text{cons}}, Q_t^{\text{prod}}, p_t^{\text{prod}}$ ), and for  $\mathcal{A}_{B_0}$  (to predict  $Q_t^{\text{ask}}, p_t^{\text{ask}}, Q_t^{\text{bid}}, p_t^{\text{bid}}$ ).

While these action effects can be predicted independently, actions need to be optimized in conjunction: the  $\text{bal}(I_{B_0,t})$  function is designed such that imbalance fees typically result in negative reward when taking actions of a single type, while positive reward can be achieved by taking actions of multiple types in parallel (to maintain low imbalance). Therefore, any tractable lookahead policy is required to efficiently (i) sample, and (ii) combine the actions to simulate.

The 2013 champion, TacTex'13 (Urieli and Stone 2014), can be viewed as approximating an MDP solution using a lookahead policy. TacTex'13 does not optimize production tariffs, wholesale selling and fees, so  $Q_t^{\text{prod}}, p_t^{\text{prod}}, Q_t^{\text{ask}}, p_t^{\text{ask}}$ , and *psw*() are always zero in Equation 1. TacTex'13's main routine is roughly Algorithm 1. For each tariff in a sample of fixed-rate consumption tariffs (line 1), it uses a *demand-predictor* to predict  $Q_t^{\text{cons}} p_t^{\text{cons}}$  for each  $t$  in the horizon (line 2), assumes  $Q_t^{\text{bid}} = Q_t^{\text{cons}}$  (and therefore marks both as  $Q_t$ ), uses a *cost-predictor* to predict for every  $t$  the price  $p_t^{\text{bid}}$  of buying  $Q_t$  (line 4), predicts a profit (called *utility* or *return*) as the sum of rewards along the horizon (line 5), and executes the utility-maximizing combination of actions (lines 7-8). Therefore, TacTex'13 lookahead efficiently *combines* actions (addressing (ii) from above) by constraining  $Q_t^{\text{bid}} = Q_t^{\text{cons}}$ , instead of examining combinations (therefore  $\text{bal}(I_{B_0,t}) = 0$ ). TacTex'13 efficiently *samples* actions (addressing (i) from above) by sampling fixed-rate tariffs in a limited region in the tariff market, and by treating wholesale actions hierarchically: it (a) treats  $Q_t^{\text{bid}}$  as

an action to be sampled (in  $Q_t^{cons}$  values), and (b) solves a subproblem of finding a cost-minimizing sequential bidding policy  $\pi(Q)$  for procuring quantities  $Q$  on a small MDP isolated from the full MDP.

---

**Algorithm 1** TacTex'13's Lookahead Policy

---

```

1: for trf in sampleCandidateTariffs()  $\cup$  {no-op} do
2:    $\{Q_t, p_t^{cons}\} | t = +1, \dots, +T \leftarrow \text{demandPredictor.predict(trf)}$ 
3:   for t in  $\{+1, \dots, +T\}$  do
4:      $p_t^{bid} \leftarrow \text{costPredictor.predict}(Q_t)$ 
5:     utilities[trf]  $\leftarrow \sum_{t=+1}^{+T} Q_t p_t^{cons} - Q_t p_t^{bid} - \text{dist}(Q_t) - \text{pub}(\text{trf})$ 
6: bestTariff  $\leftarrow \arg \max_{\text{trf}} \text{utilities}[\text{trf}]$ 
7: publish(bestTariff) // tariff market action, possibly no-op
8: procure  $\{Q_t\}_{t=+1, +2, \dots}$  predicted for bestTariff in line 2 // wholesale market

```

---

## 4.2 TacTex'15's Architecture

TacTex'15's architecture is similar to that of TacTex'13 in four main ways. TacTex'15 does not try to benefit from (1) production tariffs (2) wholesale selling, (3) imbalance, and (4) tariff fees, since in preliminary tests (1)-(3) did not seem beneficial and (4) had some simulator implementation issues (see Section 5.1). As a result, TacTex'15 assumes  $Q_t^{prod}$ ,  $p_t^{prod}$ ,  $Q_t^{ask}$ ,  $p_t^{ask}$ ,  $bal()$ , and  $psw()$  in Equation 1 to be zero. Therefore TacTex'15's lookahead policy is quite similar to TacTex'13's (Algorithm 1); we refer the reader to their paper for pseudo-code of the main routines. On the other hand, TacTex'15 introduces three main improvements over TacTex'13: it uses different (1) demand-predictor, (2) cost-predictor (both (1)-(2) are *transition-function* predictors), and (3) wholesale bidding strategy  $\pi$ .

**Demand Predictor.** The demand-predictor predicts customer subscription changes and future demand, which determine  $Q_t^{cons}$ ,  $p_t^{cons}$ . TacTex'13 learned a demand-predictor from data. However, in Power TAC there is no need to do so: these complex stochastic customer behaviors are coded in Power TAC's open-source simulator. Instead, TacTex'15 uses the simulator's customer code as a basis for its demand-predictor. However, this code does not provide a complete demand-predictor: it relies on information hidden from brokers. TacTex'15 seeds this information with expected values: customers of other brokers are assumed to be subscribed to the best tariffs, customer subscriptions are predicted in the limit (expected values after infinite time), and customer demand parameters are set to expected values. In Section 5 we examine the question of how important it is to the broker's overall performance to have an accurate demand-predictor.

**Cost Predictor.** Wholesale costs are determined by procured quantities and brokers' bidding strategies, which may change dynamically. TacTex'15 uses an adaptive cost-predictor  $Q_t^{bid} \mapsto p_t^{bid}$ , described in Algorithm 2. It has two components: a linear regression predictor trained on *boot data* (wholesale transactions sent by the simulator at game start) (line 1), and a real-time correction factor constructed from the last 24 hours' prediction errors (line 2). Since the correction factor is constructed from little data (to ensure responsiveness), we limit it to bias correction. The boot data is larger (336 instances) so we use it to determine the slope. TacTex'13's cost-predictor ignored  $Q_t^{bid}$ , and predicted past average prices based on time. We compare the two predic-

tors in Section 5.

---

**Algorithm 2** cost-predictor( $Q_t^{bid}$ )

---

```

1: reg  $\leftarrow \text{trainLinearRegression}(\{Q_t^{bid}, p_t^{bid}\}_{i \in \text{bootdata}})$ 
2: correctionFactor  $\leftarrow \text{averagePredictionErrorInLast24Hours}()$ 
3: return reg.predict( $Q_t^{bid}$ ) - correctionFactor

```

---

**Wholesale Bidding Strategy.** TacTex'15 hedges between truthful and strategic (i.e. non-truthful) bidding. A truthful bidder sets its limit price to the predicted imbalance fee  $\bar{p}$ . It gets the highest priority among competitors who bid less than  $\bar{p}$  and never pays more than  $\bar{p}$ . However, since the sequential double-auction mechanism is not incentive compatible, truthful bidding is suboptimal in some situations. TacTex'13 used an optimistic strategic (i.e. non-truthful) bidding strategy  $\pi(Q)$  that learns to bid slightly higher than each double-auction's expected clearing price. This strategy is optimal in some situations (e.g. single-buyer or cooperative setups), but can be exploited by competitors who learn to bid slightly higher. Since each of the two strategies is beneficial in different situations, TacTex'15 hedges between them. Let  $\bar{p}$  be the limit price suggested by TacTex'13's strategy, and  $\epsilon$  be the minimum amount that can be traded (0.01 mWh). To bid for a quantity  $Q_t^{bid}$ , TacTex'15 submits the following 25 orders (see MDP wholesale actions, Section 3)  $\langle Q_t^{bid} - 24\epsilon, \bar{p}, t \rangle$ ,  $\{\langle \epsilon, \bar{p} + i \frac{\bar{p} - \bar{p}}{24}, t \rangle\}_{i=0}^{23}$ . This strategy benefits from both worlds: if TacTex'15 sets the price, it will either be the strategic price returned by  $\pi(Q)$ , or the lowest among its higher bids. If another broker sets the price, TacTex'15 will have a higher priority and benefit from the lower price as long as it is not higher than  $\bar{p}$ .

**Future Extensions.** Referring back to the reward specification (Equation 1), our MDP provides a guideline for future extensions of TacTex'15's lookahead policy. In some situations a broker can profit from imbalance. We can relax the assumption that  $Q_t^{cons} = Q_t^{bid}$ , add imbalanced trajectories to our lookahead search, setting  $Q_t^{cons} - Q_t^{bid} = I_{B_0, t}$  for a sample of  $I_{B_0, t}$  values, and predict  $bal(I_{B_0, t})$  using a learned predictor. We can sample production tariffs like consumption tariffs, and treat wholesale sell hierarchically like wholesale buy actions. This addresses requirement (i) from Section 4.1 (sample actions efficiently). However, addressing requirement (ii) (combined actions efficiently) becomes more challenging. In an initial implementation we use an alternating, local improvement based approach which performs well, but more sophisticated methods might be possible. Finally, tariff-revoke actions can be added by simulating lookahead trajectories with each of the active tariffs removed. Initial implementation shows promising results.

## 5 Results

We analyze TacTex'15's performance in competitions (Section 5.1) and controlled experiments (Section 5.2).

### 5.1 Competition Results

The Power TAC 2015 Finals included 11 teams from universities in America, Europe and Asia. 230 games were played continually over a week, in three different sizes: 3-brokers, 9-brokers, and 11-brokers. A day after the finals

ended, 8 of the teams competed in a post-finals, demo-competition with 70 4-broker games. While being unofficial, this competition was run similarly to the finals with one important difference: a simulator-loophole that was exploited during the finals, was fixed. Due to the proximity to the finals, and a parallel workshop, we believe that teams used the same brokers they used in the finals.

Table 1 summarizes the 2015 finals results. While TacTex’15 was officially ranked 2nd, it was the best broker that did not exploit a simulator-loophole: the 1st-ranked broker gained the highest overall score by exploiting a simulator loophole in 3-broker games, which resulted in unrealistic dynamics and an unrealistically high score that biased the final ranking (see dark gray cells in Table 1).<sup>1</sup> Specifically, Maxon15 subscribed customers to inflated tariffs which promised customers large payments if customers *unsubscribed* from them after a period shorter than customers’ minimum response time. Due to the loophole, customers subscribed to these tariffs assuming they could collect the payments, even though they could not.

Table 1: **Power TAC 2015 finals results.** Ranking is determined by the “Total” score, which is a sum of individual z-scores in each game size, displayed in the columns “11-brokers” (10 games played by all brokers), “9-brokers” (45 games played by each broker) and “3-brokers” (45 games played by each broker).

Broker	11-brokers	9-brokers	3-brokers	Total
Maxon15	0.611	0.801	1.990	3.402
TacTex’15	0.897	1.066	0.258	2.221
CUHKTac	0.962	0.859	0.106	1.927
AgentUDE	0.421	0.367	0.809	1.597
Sharp	0.429	0.614	0.521	1.564
COLDPower	0.726	0.397	-0.751	0.371
cwiBroker	-0.002	-0.120	0.465	0.343
Mertacor	0.413	0.142	-1.341	-0.786
NTUTacAgent	-1.017	-1.638	0.453	-2.202
SPOT	-1.052	-0.243	-1.032	-2.327
CrocodileAgent	-2.387	-2.244	-1.479	-6.111

Table 2: **Power TAC 2015 post-finals demo competition results.** 70 games were played in a single game-size (4-brokers). Ranking is determined by z-score.

Broker	4-brokers (profits)	4-brokers (z-score)
TacTex’15	15.0M	1.122
Maxon15	10.7M	0.627
CUHKTac	10.0M	0.537
AgentUDE	9.7M	0.509
cwiBroker2015	7.9M	0.297
Sharp	4.6M	-0.092
COLDPower	-0.8M	-0.724
SPOT	-14.0M	-2.276

After the finals, the loophole was fixed. When replaying 3-broker competition games without the loophole, Maxon15 no longer won by a large gap, but instead lost by a large gap to TacTex’15. When taking into account only 11- and 9-broker games from the finals (where the loophole had no impact), TacTex’15 ended 1st with a total z-score of 0.142 ahead of CUHKTac and 0.551 ahead of Maxon15, finishing slightly behind CUHKTac in 11-broker games (by 0.065) and ahead of CUHKTac in 9-broker games (by 0.207). In the post-finals demo competition with a repaired simulator, TacTex’15 won by a large gap ahead of the others (Table 2),

<sup>1</sup>Maxon was not disqualified: they explained it as an unintended result of automatic parameter tuning right before the finals.

making 50% more profits than the 2nd place (Maxon15). Maxon15 used the same strategy as before, but it was not as effective with the loophole fixed.<sup>2</sup>

## 5.2 Controlled Experiments

While the competition is motivating and its results are illustrative, it cannot isolate specific broker components in a statistically significant way. We therefore subsequently tested TacTex’15 in thousands of games, in two types of controlled experiments: (a) performance tests, and (b) ablation analysis tests, which evaluate the contribution of TacTex’15’s main components to its overall performance.

**Experimental Setup.** Each experiment consisted of running 56 games against a set of opponent brokers, using broker binaries of 2015 finalists. To better evaluate statistical significance, we held most of the random factors in the simulation fixed across experiments (random seeds, weather conditions). To fix weather conditions, we used weather files containing 3 months of real-world weather. To cover year-round weather conditions we used 8 weather files (each file used by 1/8 of the games) with start-dates of January, April, July, October of 2009 and 2010.

**Performance Tests.** A successful broker should perform well in expectation against every set of opponents, under different stochastic conditions (here weather/random seeds). Currently, five 2015 finalists have released their brokers’ binaries. We used these binaries to test TacTex’15’s performance in 2, 3, . . . , 6-broker games. We generated combinations of brokers for each game size, and tested each combination in 56 games, as described above. Figure 2 presents the results. TacTex’15 significantly won against every combination of opponents, typically by a large gap.

**Ablation Analysis.** To understand the reasons for TacTex’15’s success, we tested the contribution of TacTex’15’s main components to its overall performance, in all possible game-sizes (2,...,6). We created three ablated versions of TacTex’15 by disabling each of its main components. For each game size, we selected the “strongest” combination of opponents, against which TacTex’15 had the lowest score. We tested each ablated version against these opponents in a 56-game experiment, holding random seeds and weather conditions fixed to the same values used against TacTex’15. When disabling a component, we used as a baseline the corresponding component used by TacTex’13 (since TacTex’15’s ablated version must have some component in place of a disabled one to run properly). Figure 3 shows the results of our ablation analysis. Disabling the cost-predictor (Abl-cost) did not have significant impact on TacTex’15’s performance (however it can reduce performance, see Figure 4b). Disabling the wholesale-bidding strategy (Abl-bid) significantly hurts TacTex’15’s performance: it reduces TacTex’15’s score in game sizes 2, 4, 5, 6, and it causes TacTex’15 to either lose its lead (in game sizes 2, 3) or have a smaller victory margin (in game sizes 4, 5,

<sup>2</sup>To be fair, one should note that they did not retune their parameters to the repaired simulator. On the other hand, it’s not clear that other parameters would have done particularly better in the absence of the loophole.



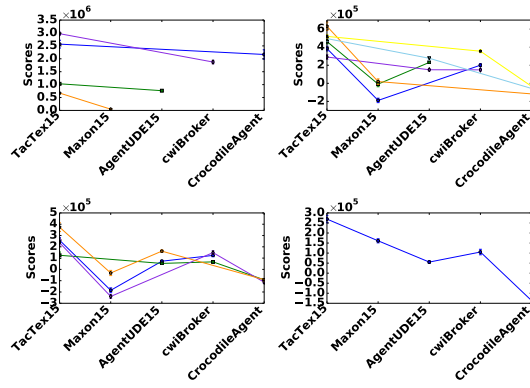


Figure 2: **Performance of TacTex'15 against Power TAC 2015 finalists in controlled experiments of game-sizes of 2-5.** Each line represents the average scores of a combination of brokers playing each other under a variety of conditions (note the small error bars). Results are shown for game-sizes of 2-, 3-, 4-, 5-brokers (top-left, top-right, bottom-left, bottom-right, respectively). Similar results for 6-brokers are omitted. TacTex'15 consistently won against all combinations of brokers, in all game-sizes.

6). Disabling the demand-predictor (Abl-demand) significantly hurts TacTex'15's performance: it drops TacTex'15's score in all game sizes, and causes TacTex'15 to either lose its lead (in game sizes 3, 5, 6) or have a smaller victory margin (in game sizes 2, 4).

Next, we extended our analysis as follows. Figure 4a extends Figure 3's 3-broker demand-predictor ablation to a continuum of ablation levels (0 is TacTex'15's predictor, 1 is TacTex'13's predictor, values in (0,1) are weightings of the two. TacTex'15's score degrades even for small ablation levels. Figure 4b extends the cost-predictor ablation. It shows the average score throughout games where TacTex'15 played Abl-cost and where wholesale costs dropped abruptly in timeslot 1080. TacTex'15's adaptive cost-predictor was quicker to react; TacTex'15 updated its tariffs and won significantly (note the small error margin).

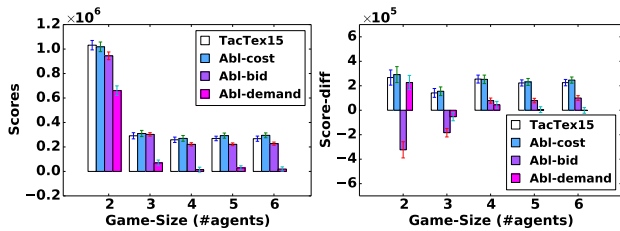
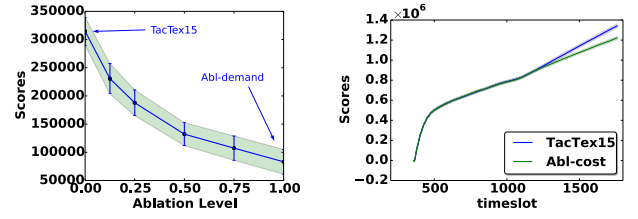


Figure 3: **Ablation analysis for 2-6 broker games.** The performance of TacTex'15 is compared with three of its ablated versions, when playing against the strongest combination of opponents in each game size. Ablated versions are constructed from TacTex'15 by disabling cost predictor (Abl-cost), wholesale-bidding strategy (Abl-bid), and demand-predictor (Abl-demand). The left figure shows the average scores of each version in each game size; the right figure shows the average score-differences of each version from opponents' average score (y-axes' scales are the same).



(a) Demand-predictor

(b) Cost-predictor

Figure 4: **Ablation analysis extensions** (see text for explanation).

## 6 Related Work

This work is the first to formalize the *complete* broker's power trading problem as an MDP, and characterize its approximate solutions. Previous research either used an MDP to model a more abstract trading problem (Reddy and Veloso 2011), or used an MDP to model a subproblem of the complete problem (Peters et al. 2013; Kuate et al. 2013; Kuate, Chli, and Wang 2014; Babic and Podobnik 2014; Urieli and Stone 2014; Ozdemir and Unland 2015). Other approaches to power trading did not directly optimize the predicted utility. AgentUDE14 (Ozdemir and Unland 2015) (1st place, 2014) used an empirically tuned tariff strategy provoking subscription changes and withdraw payments, and Q-learning for wholesale bidding. CwiBroker14 (2nd place, 2014) (Hoogland and Poutre 2015) used tuned heuristics based on domain knowledge. An analysis of the 2014 Power TAC finals can be found at (Babic and Podobnik 2014b). Mertacor13 (Ntagka, Chrysopoulos, and Mitkas 2014) used Particle Swarm Optimization based tariff strategy. CwiBroker13 (Liefers, Hoogland, and Poutre 2014) (2nd place, 2013) used tariff strategy inspired by Tit-for-Tat. Their wholesale strategy introduced the idea of multiple bids, but was based on equilibria in continuous auctions, rather than TacTex'15's hedging between optimistic strategic bidding and truthful bidding. In other trading agent competitions, utility-optimization approaches were used in different market structures (Stone et al. 2003; Pardoe 2011). Other approaches included game theoretic analysis of the economy (Kiekintveld, Vorobeychik, and Wellman 2006) and fuzzy reasoning (He et al. 2005).

## 7 Conclusion

This paper has focused on the question: how should an autonomous electricity broker agent act in competitive electricity markets to maximize its profit. We have formalized the complete electricity trading problem as an MDP, which is computationally intractable to solve exactly. Our formalization provides a guideline for approximating the MDP's solution, and for extending existing approximate solutions. We introduced TacTex'15 which uses a lookahead policy that improves upon that of a previously champion agent, and achieves state-of-the-art performance in competitions and controlled experiments. Using thousands of experiments against 2015 finalist brokers, we analyzed TacTex'15's performance and reasons for its success, finding that while its

lookahead policy is an effective solution in the power trading domain, its performance can be sensitive to errors in the transition function prediction, especially demand-prediction. An important direction for future work is to further close the gap between the current approximate solution to the trading MDP and its fully optimal solution.

## 8 Acknowledgments

This work has taken place in the Learning Agents Research Group (LARG) at UT Austin. LARG research is supported in part by NSF (CNS-1330072, CNS-1305287), ONR (21C184-01), AFRL (FA8750-14-1-0070), and AFOSR (FA9550-14-1-0087).

## References

- Babic, J., and Podobnik, V. 2014a. Adaptive bidding for electricity wholesale markets in a smart grid. In *AAMAS Workshop on Agent-Mediated Electronic Commerce and Trading Agents Design and Analysis (AMEC/TADA 2014)*.
- Babic, J., and Podobnik, V. 2014b. An analysis of power trading agent competition 2014. In Ceppi, S.; David, E.; Podobnik, V.; Robu, V.; Shehory, O.; Stein, S.; and Vetsikas, I. A., eds., *Agent-Mediated Electronic Commerce. Designing Trading Strategies and Mechanisms for Electronic Markets*, volume 187 of *Lecture Notes in Business Information Processing*, 1–15. Springer International Publishing.
- Borenstein, S. 2002. The trouble with electricity markets: Understanding california’s restructuring disaster. *Journal of Economic Perspectives* 16(1):191–211.
- He, M.; Rogers, A.; David, E.; and Jennings, N. R. 2005. Designing and evaluating an adaptive trading agent for supply chain management applications. In Poutre, H. L.; Sadeh, N.; and Sverker, J., eds., *Agent-mediated Electronic Commerce, Designing Trading Agents and Mechanisms: AAMAS 2005 Workshop AMEC 2005, Utrecht, Netherlands, July 25, 2005, and IJCAI 2005 Workshop TADA 2005, Edinburgh, UK, August 1, 2005, Selected and Revised Papers*. Springer. 35–42. Event Dates: August 2005.
- Hoogland, J., and Poutre, H. L. 2015. An effective broker for the power tac 2014. In *AAMAS Workshop on Agent-Mediated Electronic Commerce and Trading Agents Design and Analysis (AMEC/TADA 2015)*.
- Joskow, P. L. 2008. Lessons learned from electricity market liberalization. *The Energy Journal* Volume 29.
- Ketter, W.; Collins, J.; Reddy, P. P.; and Weerdt, M. D. 2015. The 2015 power trading agent competition. *ERIM Report Series Reference No. ERS-2015-001-LIS*.
- Ketter, W.; Peters, M.; and Collins, J. 2013. Autonomous agents in future energy markets: The 2012 Power Trading Agent Competition. In *Proceedings of the Twenty-Seventh AAAI Conference on Artificial Intelligence*. AAAI.
- Kiekintveld, C.; Vorobeychik, Y.; and Wellman, M. 2006. An analysis of the 2004 supply chain management trading agent competition. In Poutre, H.; Sadeh, N.; and Janson, S., eds., *Agent-Mediated Electronic Commerce. Designing Trading Agents and Mechanisms*, volume 3937 of *Lecture Notes in Computer Science*. Springer Berlin Heidelberg. 99–112.
- Kuate, R. T.; He, M.; Chli, M.; and Wang, H. H. 2013. An intelligent broker agent for energy trading: An mdp approach. In *The 23rd International Joint Conference on Artificial Intelligence*.
- Kuate, R. T.; Chli, M.; and Wang, H. H. 2014. Optimising market share and profit margin: Smdp-based tariff pricing under the smart grid paradigm. In *Innovative Smart Grid Technologies Conference Europe (ISGT-Europe), 2014 IEEE PES*, 1–6.
- Liefers, B.; Hoogland, J.; and Poutre, H. L. 2014. A successful broker agent for power tac. In *AAMAS Workshop on Agent-Mediated Electronic Commerce and Trading Agents Design and Analysis (AMEC/TADA 2014)*.
- Ntagka, E.; Chrysopoulos, A.; and Mitkas, P. A. 2014. Designing tariffs in a competitive energy market using particle swarm optimization techniques. In *AAMAS Workshop on Agent-Mediated Electronic Commerce and Trading Agents Design and Analysis (AMEC/TADA 2014)*.
- Ozdemir, S., and Unland, R. 2015. Agentude: The success story of the power tac 2014’s champion. In *AAMAS Workshop on Agent-Mediated Electronic Commerce and Trading Agents Design and Analysis (AMEC/TADA 2015)*.
- Pardoe, D. M. 2011. *Adaptive Trading Agent Strategies Using Market Experience*. Ph.D. Dissertation.
- Peters, M.; Ketter, W.; Saar-Tsechansky, M.; and Collins, J. 2013. A reinforcement learning approach to autonomous decision-making in smart electricity markets. *Machine Learning* 92(1):5–39.
- Powell, W., and Meisel, S. 2015. Tutorial on stochastic optimization in energy – part ii: An energy storage illustration. *Power Systems, IEEE Transactions on* PP(99):1–8.
- Puterman, M. L. 1994. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. New York, NY, USA: John Wiley & Sons, Inc., 1st edition.
- Reddy, P. P., and Veloso, M. M. 2011. Strategy learning for autonomous agents in smart grid markets. In *Proceedings of the Twenty-Second international joint conference on Artificial Intelligence-Volume Volume Two*, 1446–1451. AAAI Press.
- Reddy, P. P., and Veloso, M. M. 2012. Factored Models for Multiscale Decision Making in Smart Grid Customers. In *Proceedings of the Twenty-Sixth AAAI Conference on Artificial Intelligence (AAAI-12)*.
- Stoft, S. 2002. *Index*. Wiley-IEEE Press. 460–468.
- Stone, P.; Schapire, R. E.; Littman, M. L.; Csirik, J. A.; and McAllester, D. 2003. Decision-theoretic bidding based on learned density models in simultaneous, interacting auctions. *Journal of Artificial Intelligence Research* 19:209–242.
- Urieli, D., and Stone, P. 2014. Tactex’13: A champion adaptive power trading agent. In *Proceedings of the Twenty-Eighth Conference on Artificial Intelligence (AAAI 2014)*.
- U.S. Department of Energy. 2003. “Grid 2030” A National Vision For Electricity’s Second 100 Years.