

Mind the Gap: Improving Water Markets with Field-Level Remote Sensing*

Katherine Wright[†] Andrew Ayres[‡] Bryan Leonard[§]

PRELIMINARY DRAFT. PLEASE DO NOT CIRCULATE WITHOUT PERMISSION.

Abstract

Market-based water transfers are increasingly proposed as a climate-adaptation tool in arid regions, but they are hindered by the inability to accurately measure field-scale water use. Instead, existing transfers typically rest on two assumptions: that fallowed land uses no water, and that participation in transfer programs is not systematically related to baseline and predicted water use. We combine satellite evapotranspiration data on 2,500 fields over 20 years with hand-coded contract records to evaluate one of the largest agricultural-to-urban water transfers in the United States—the Palo Verde Irrigation District—Metropolitan Water District Forbearance and Fallowing Program. Using a dynamic difference-in-differences design, we find while fallowing does reduce consumptive water use, the estimated savings are only 53% of reported savings, implying a cumulative shortfall of 624,916 acre-feet over 2005–2021. Roughly half the gap is hydrologic spillover from irrigated neighbors onto fallowed parcels, isolated from atmospheric measurement artifacts via a wind-direction test. The other half is strategic non-additionality—farmers preferentially fallowing fields with lower expected water use—identified by a machine-learning algorithm that predicts water use on fallowed fields.

*For helpful comments we thank Josh Abbott, Anna Boser, Tamma Carleton, Eric Edwards, Robert Heilmayr, Kelsey Jack, Kailin Kroetz, Morgan Levy, Buzz Thompson, and Will Rafey, as well as seminar participants at the Workshop on Advancing Sustainable Water Management at Stanford University, the Bren School of Environmental Science and Management at UC Santa Barbara, the Environmental and Resource Economics Series at UC San Diego, the Land, Environment, Economics and Policy Institute at the University of Exeter, the Agricultural and Resource Economics Department at UC Davis, the Economics Department at the Colorado School of Mines, the Department of Agricultural Economics & Economics at Montana State University, the Department of Applied Economics at Oregon State University, and the UC Riverside Water Dialogues. Talia Greco, Danielle Moon, and Mithran Mathiraj provided excellent research assistance. We have also benefited greatly from data sharing and discussions with staff at Metropolitan Water District.

[†]Hillsdale College kwright1@hillsdale.edu

[‡]University of Nevada, Reno andrew.ayres@unr.edu

[§]University of Wyoming bryan.leonard@uwyo.edu

Introduction

Many arid regions are already experiencing more severe surface water scarcity due to climate change (Rodell et al., 2018; Elliott et al., 2014; Gordon et al., 2024). The Western United States is a prime example: the Colorado River supplies 40 million people and irrigates more than five million acres, and consumptive demands have exceeded the river’s flow in 16 of the years between 2000 and 2020 (Richter et al., 2024). In this setting, ensuring that water is allocated to its highest-valued use has emerged as a central policy challenge. Agriculture accounts for as much as 80% of freshwater use in the basin (Richter et al., 2020; Medellín-Azuara et al., 2024), even though urban willingness-to-pay often far exceeds the market value of water in agriculture (Brewer et al., 2008), and recent evidence suggests that the allocation even within agriculture is inefficient (Arellano-Gonzalez et al., 2021; Rafey, 2026).

Water markets are a leading candidate solution (Brewer et al., 2008; Ayres et al., 2021; Bruno and Jessoe, 2021; Rafey, 2023, 2026; Bruno and Jessoe, 2024), but their expansion has been hindered by transaction costs related to unclear property rights, administrative approvals, conveyance, and third-party objections (Bretsen and Hill, 2008; Leonard et al., 2019; Hagerty, 2023). Many of these costs arise from the difficulty of estimating field-scale consumptive water use to determine how much of a farmer’s diversion right can be transferred without impairing other users. To circumvent these costs, municipalities, NGOs, and the U.S. Bureau of Reclamation increasingly use “alternative transfer mechanisms” (ATMs)—contracts that pay farmers to change on-farm practices (i.e., fallowing or crop-switching), often for limited duration, with the understanding that water savings accrue to the buyer (Dilling et al., 2019). Water transfers of various forms span the Upper and Lower Colorado Basins, the broader Western United States, and other countries including Australia, Spain, and Chile (Grafton et al., 2012; Garrick et al., 2023; Rafey, 2023).

The verification problem that surface water transfers face is not unique. Programs that pay for changes in land use to deliver an environmental outcome—REDD+ forest carbon, agricultural conservation easements, biodiversity offsets, and fallowing contracts—often rely on two crucial assumptions to estimate environmental benefits: that the targeted activity (deforestation, irrigation, development) ceases entirely on enrolled land, and that potential environmental gains are not correlated with enrollment. Both assumptions are testable, but only when reliable measurement of the targeted activity is feasible at the contract scale. Recent advances in satellite measurement now make field-scale verification possible for water use (Volk et al., 2024). Despite this measurement revolution, the magnitude of the verification gap in real-world programs has not been documented. Here, we fill this gap, in a setting where the verification problem is especially tractable: arid, flood-irrigated agriculture without the possibility of groundwater substitution.

We study the Palo Verde Irrigation District–Metropolitan Water District Forbearance and Fallowing Program, one of the largest and longest-running agricultural-to-urban alternative water transfers to date, involving the largest municipal water supplier in the United States. We combine satellite evapotranspiration data from OpenET (Volk et al., 2024) with hand-digitized administrative records of fallowing calls issued by MWD between 2005 and 2021 to estimate the causal

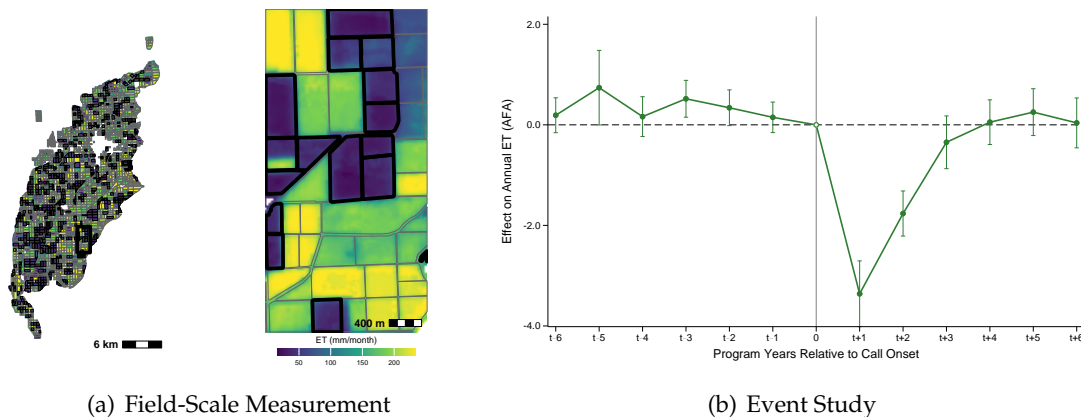
effect of a call on field-level water use. Identification rests on a dynamic difference-in-differences estimator (De Chaisemartin and d’Haultfoeuille, 2020) that accommodates the rotational, on/off nature of fallowing calls and absorbs owner-by-year shocks. We then decompose the gap between our estimated savings and those reported by MWD/USBR into hydrologic spillover from irrigated neighbors (isolated from atmospheric measurement artifact using monthly prevailing wind directions) and strategic non-additionality of selected fields (identified using AR(1) and LASSO predictions of counterfactual ET). The analysis reveals structural biases in common methods currently used by agencies such as the US Bureau of Reclamation to estimate the savings from water transfer programs.

Results

Water Savings from Fallowing Are Substantially Lower than MWD/USBR Estimates

We estimate the causal effect of a fallowing call on field-level water use using a dynamic difference-in-difference framework that compares ET on called fields before versus during a call to ET on un-called fields over the same period. Panel (a) of Figure 1 depicts an illustrative example, with thicker borders indicating called fields, and darker shading indicating less water use. Panel (b) of Figure 1 reports the event-study estimates using the (De Chaisemartin and d’Haultfoeuille, 2020) estimator, which is the only modern DiD framework that accommodates the on/off rotational nature of fallowing assignments. Called fields exhibit similar trends to untreated parcels prior to onset and reduce their water use by an average of 2.53 AFA/year during a fallowing spell before returning to baseline.

Figure 1: Identifying and Measuring the Water-Savings Gap



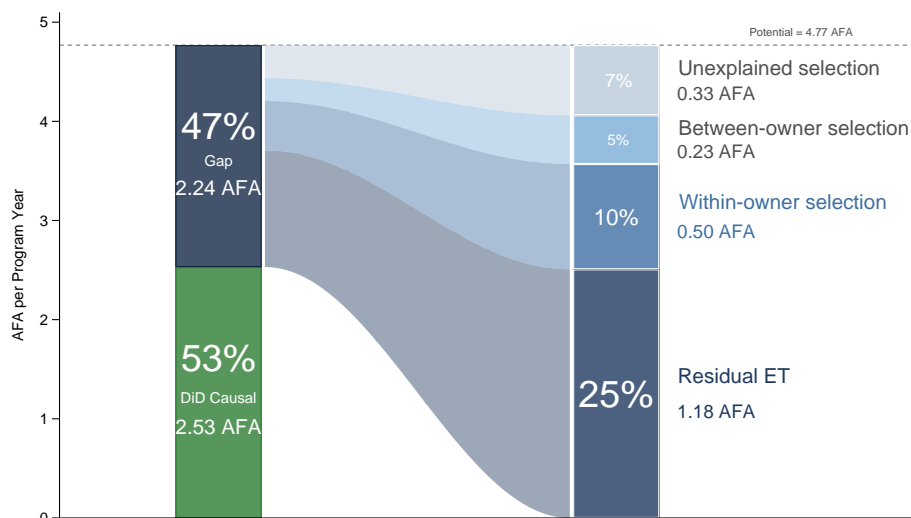
Notes: Panel (a) depicts field-scale water use from the OpenET Ensemble as well as “called” fields in August 2011, a representative example from our data. Panel (b) reports the dynamic event study estimates with six pre-period placebos and six post-onset effects. Connected green line with 95% CI bars, clustered at PLSS section.

We also verify the robustness of our savings estimate. The magnitude varies slightly across the choice of ET model, with estimates for our preferred specification ranging from 1.84 to 3.01 AFA across the seven OpenET models (SI Table S3). The full set of 42 alternative specifications, varying treatment coding, sample restrictions, and control sets, is summarized in Fig. S5. To check whether

farmers offset fallowing on called fields by using more water on their other fields—which would inflate our savings estimates—we examined water-use patterns on *uncalled* fields of called farmers vs. uncalled farmers (Figure S8) and re-estimated the main savings effect under alternative specifications that address this concern (Table S4). Neither shows evidence of significant compensating water use, leaving within-farm reallocation of water as an unlikely driver of our results.

In contrast to the field-scale approach depicted in Figure 1, the MWD/USBR savings methodology multiplies district-wide average water use per acre by the number of contracted fallow acres. This procedure makes three implicit assumptions: that average water use does not differ between called and uncalled fields, that fallowing is fully additional (no fallowed fields would have been fallowed in the absence of the program), and that water use goes to zero on fallowed parcels. Our empirical strategy is designed to relax all three, and allows us to isolate the effect of each assumption on the overall differences between our causal estimates and the MWD/USBR methodology. Figure 2 depicts the overall gap between the causal estimate and the average reported savings in MWD’s verification reports, as well as our decomposition of the differences.

Figure 2: Decomposing the Savings Gap



Notes: Decomposition of the Potential savings rate (4.77 AFA per acre-year, the ET baseline on uncalled ever-enrolled fields). Stage 1 splits Potential into the preferred DiD causal estimate (2.53 AFA, green, 53% of Potential) and a 2.24 AFA residual Gap (dark blue). Stage 2 decomposes the Gap into its four mechanism components: *Residual ET* on called fields (1.18 AFA), non-additionality due to *Between-owner selection* (0.23 AFA), non-additionality due to *Within-owner selection* (0.50 AFA), and *Unexplained* non-additionality (0.33 AFA).

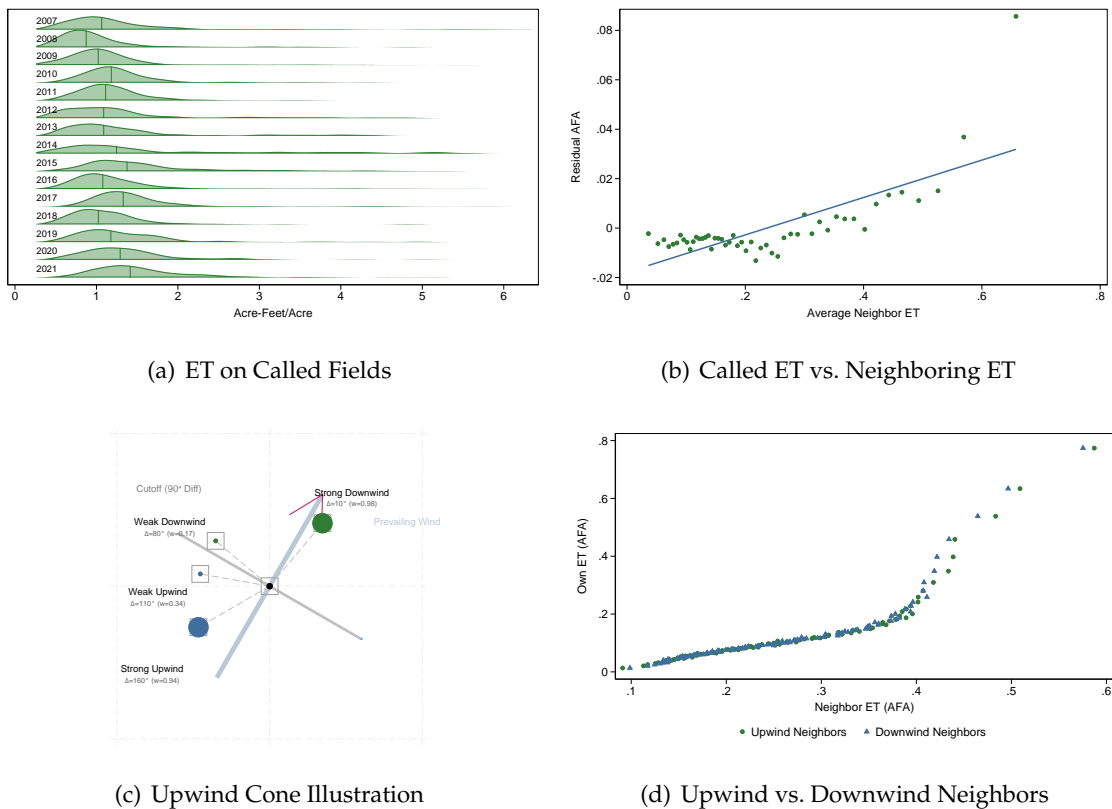
Average annual savings reported by MWD over 2005–2021 are approximated by district-wide average annual water use, at 4.67 AFA. Average water use on *enrolled* fields is slightly higher than the district-wide average, putting potential savings under perfect compliance and additionality at 4.77 AFA (SI Fig. S9). The realized causal savings of 2.53 AFA represent 53% of potential; the remaining 2.24 AFA is the savings gap. Roughly half the gap is persistent water use on called fields (1.18 AFA), and roughly half is non-additionality of the called fields themselves (1.06 AFA). The mechanism evidence below addresses each in turn.

Consumptive Use on Fallowed Fields Driven by Neighbors' Irrigation

A meaningful number of called fields have positive ET, with an overall mean of 1.18 AFA/year across fields and years (panel (a) of Figure 3). This residual accounts for roughly half the savings gap. A similar finding for fallowed land in PVID was previously noted by [Wobus et al. \(2024\)](#), who attributed it to errors in satellite classification of fallowed lands. Such measurement errors cannot explain the persistent use we observe, because our classification of fallowing is based on administrative records of fields that MWD verifies are in fact fallow.

We explore several potential drivers—distance to the Colorado River, soil drainage, cumulative length of fallowing spells, effective precipitation, and water use on neighboring fields (SI Fig. S10 and SI Table S5). Of these, neighbor ET is the only significant predictor: panel (b) of Figure 3 plots within-field residual ET against the field's mean neighbor ET in the same month. The relationship between fallow ET and neighboring ET is statistically significant and large enough to fully explain fallow ET (Materials and Methods).

Figure 3: Water Use on Fallowed Fields



Notes: Panel (a) is a joy plot of annual ET across called field-years, one distribution per year. Panel (b) is a binscatter of within-field residual ET (after partialling out field \times month fixed effects) against the field's average neighbor ET in the same month; green dots = 50 equal-count bins, blue fit line. Panels (c)–(d) illustrate our approach for ruling out atmospheric spillovers: (c) schematic of the prevailing-wind cone classification, where each neighbor is weighted by \cos of its bearing relative to the monthly prevailing direction; (d) scatter of own ET against upwind-weighted (green dots) vs. downwind-weighted (blue triangles) neighbor ET on called field-months — the near-identical slopes bound atmospheric (wind-borne vapor) contamination at $\leq 6\%$ of the total spillover (econometrics and wind-test power calculation in SI).

Next, we assess whether neighbor effects reflect actual spillover of applied water or atmospheric contamination of the satellite ET signature. Evaporative cooling on irrigated parcels could lower surface temperatures on adjacent fallowed parcels, creating a false ET reading. We rule out the latter channel using monthly prevailing wind directions constructed from hourly station data. Atmospheric vapor transport should originate disproportionately from upwind neighbors, and satellite measurement spillovers would reflect this. Panel (d) of Figure 3 plots own ET against cosine-weighted upwind versus downwind neighbor ET (weighting scheme depicted in panel (c)); the slopes are statistically indistinguishable ($p = 0.64$). At 80% power, the design could have detected an upwind–downwind contrast of 0.024 AFA, or 5.9% of the total neighbor-ET slope (SI Table S7). The atmospheric channel is therefore bounded well below the magnitude that would be required for it to be the dominant mechanism.

We also examine whether the residual reflects non-agricultural ET (riparian vegetation, weeds, evaporation). To do so, we adapt the methodology of Boser et al. (2022) to estimate agricultural ET as a sub-component of total ET using machine learning. This approach was originally developed in California’s Central Valley. As depicted in Figure S13 and Table S10, applying the Boser et al. method to this setting produces estimates of non-agricultural ET that are implausibly larger, (and, correspondingly, applied ET estimates that are significantly lower than reported water use in PVID on a monthly and annual basis). When we re-estimate our main specification using these modeled data, the result is roughly 82.6% of the headline DiD effect (SI Table S9), albeit on a slightly different sample. Applying the agricultural share to the full-sample headline yields a smaller implied estimate of -2.090 AFA, and hence a larger gap from MWD/USBR’s reported savings.

Farmer Selection of Fields into Fallowing Calls Reduces Potential Savings

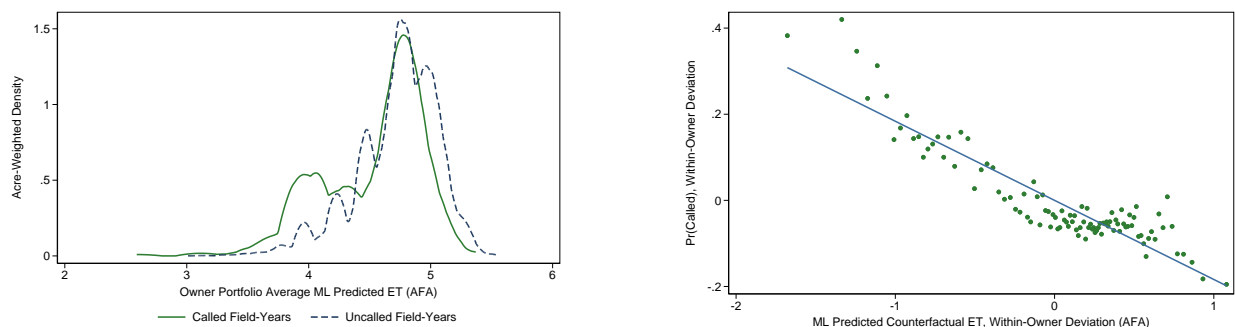
The remaining 1.06 AFA of the savings gap reflects non-additionality—fallowing of fields that would have used less water even without the program. We characterize the magnitude and structure of non-additionality by predicting counterfactual water use on called fields with two different models. The first is an AR(1) model with field-by-year fixed effects. The second is a cross-validated LASSO with 2-year lags, neighbor-ET histories, and year fixed effects. We use these predictions to measure the extent of strategic selection of fields with lower-than-average predicted ET. First, we verify that fields with lower predicted ET are more likely to be selected for a call (Table S8).

Next, we distinguish between-owner from within-owner selection, focusing on the Lasso predictions. In panel (a) of Figure 4, we plot the distributions of acre-weighted portfolio-average predicted counterfactual ET separately at the owner-year level for called acres vs. acres left in production. By construction, this average does not vary within an owner-year, because the within-owner choice over which specific field to fallow is filtered out by the portfolio-averaging. Differences between the two densities are therefore driven entirely by compositional variation in which portfolios contribute to the called pool. The mean shift between the two densities indicates that called acres are drawn disproportionately from owner portfolios with somewhat lower predicted ET than the enrolled benchmark. Panel (b) shows the within-owner channel: residualizing on

owner-by-year fixed effects, fields with lower predicted ET are more likely to be called than other fields associated with the same owner.

Panels (c) and (d) quantify the total observable selection effect, which is calculated as the difference between predicted ET on uncalled vs. called field-years among ever-enrolled fields. The total effect is 0.44 AFA under AR(1) (15% between, 85% within) and 0.73 AFA under LASSO (32% between, 68% within). The LASSO decomposition absorbs roughly 70% of the 1.06 AFA non-additionality component of the savings gap. The remaining 0.33 AFA, which we label *unexplained*, captures selection on factors our prediction models cannot recover from spatial and historical covariates—private farmer information about expected ET, idiosyncratic operational decisions, or measurement noise in the prediction. The within-owner channel is the dominant component of observable selection under both prediction models, indicating that the dominant non-additionality margin is the field-level decision within an enrolled portfolio.

Figure 4: Non-Additionality and Selection into Following



(a) Between-Owner Selection into Following

(b) Within-Owner Selection into Following

Component	AFA	Share
Between-Owner Selection	0.07	15%
Within-Owner Selection	0.37	85%
Total Predicted	0.44	100%

(c) Selection Decomposition, AR(1)

Component	AFA	Share
Between-Owner Selection	0.23	32%
Within-Owner Selection	0.50	68%
Total Predicted	0.73	100%

(d) Selection Decomposition, LASSO

Notes: Panel (a) plots the acre-weighted distribution of owner-year portfolio-average predicted counterfactual ET, separately for called vs. uncalled acres in a given year. Panel (b) is a within-owner-year binscatter of the call indicator (residualized on owner \times year fixed effects) against predicted counterfactual ET (same residualization). Panels (a) and (b) both use the ML counterfactual prediction. Panels (c)–(d) decompose the selection-on-observables portion of the non-additionality gap into between-owner and within-owner components. Panel (c) uses an AR(1) model of own ET, yielding 0.44 AFA total predicted selection (15% between, 85% within). Panel (d) uses a cross-validated LASSO with 12- and 24-month lags, neighbor behavior, and year fixed effects, yielding 0.73 AFA (32% between, 68% within). The LASSO-predicted total absorbs \sim 70% of the 1.06 AFA non-additionality component.

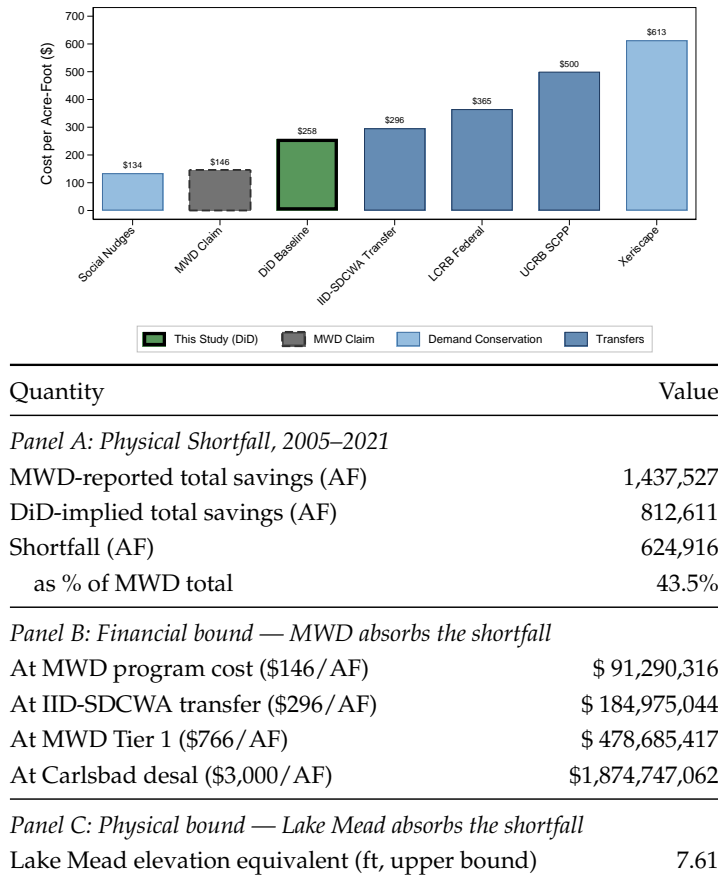
Implications for the Lower Colorado River Basin

To translate per-acre savings into program-aggregate savings, we multiply our annual DiD estimates by the number of fallowed acres MWD reports in each program year. We then calculate the difference between MWD’s aggregate reported savings and the savings implied by our causal

estimate in each year. Figure S12 sums these annual differences over 2005–2021, with the preferred estimate indicated in green and the 41 alternative specifications in gray. The majority of estimates lie between 400,000 and 800,000 AF over the first 17 years of the program; our preferred estimate is 624,916 AF, or 43.5% of MWD’s reported cumulative savings of 1,437,527 AF. PVID as a whole typically consumes a reported 350,000 to 420,000 AF per year, so the cumulative shortfall is roughly 1.5 years of district-wide consumption.

Figure 5 places our findings in the context of water scarcity in the Lower Colorado River Basin. Our findings imply that the cost-per-AF of the program is roughly double what MWD has previously estimated, but this does not change the relative positioning of this particular program compared to other strategies for addressing water scarcity based on administrative estimates and peer-reviewed literature of the cost-per-AF of different approaches (see, e.g., Avila et al. (2025)).

Figure 5: Implications for Water Supply and Cost in the Lower Colorado River Basin



Notes: Upper panel: cost per acre-foot of the PVID following program (green = preferred DiD estimate; neutral gray = MWD’s reported cost) compared to alternative Southern California supply and demand-management options (light blue = demand-side programs; dark blue = agricultural-water transfers). Lower table: translation of the 624,916-AF physical shortfall (2005–2021) into financial and physical bounds on who bears the cost. *Panel A* restates the physical shortfall in acre-feet. *Panel B* values the shortfall at four anchor prices: MWD’s estimated cost per AF in PVID, the IID–SDCWA transfer market, MWD’s Tier 1 untreated marginal supply, and Carlsbad desalination. *Panel C* translates the shortfall into a reservoir-elevation equivalent using observed end-of-calendar-year Lake Mead stages and the USBR area-capacity table to convert each year’s shortfall at its year-specific AF-per-foot (upper bound, assuming 1:1 pass-through).

As we describe in Materials and Methods, tracing the exact incidence of the aggregate $\approx 624,000$ AF savings shortfall is complex, but the lower panel in Figure 5 presents the implications of two boundary cases. If MWD “owns” the shortfall and must find alternative sources of supply to cover the difference, the cost ranges from \$91 million to \$1.8 billion over 2005–2021, depending on the source of supply. The wide dispersion of this estimate reflects broad differences in the costs of alternative supply augmentation strategies, not uncertainty in our estimate of the savings shortfall. On the other extreme, if the estimation error results directly in over-appropriation from the Lower Colorado River, then the program may have lowered the elevation of Lake Mead by as much as 7.6 feet over its first 17 years (Materials and Methods). As we emphasize in the Discussion, these are hypothetical bounds on the implications of the program, and the true incidence of the estimate error likely does not fall at either extreme.

Discussion

Our results indicate that reductions in consumptive agricultural water use due to the PVID-MWD Forbearance and Fallowing Program are substantially lower than estimates associated with the verification methodology employed by MWD and USBR. Crucially, the MWD/USBR method applies district-wide average water use per acre and assumes 100% reductions in consumptive use when estimating savings. Using remote sensing to develop field-scale estimates of ET allows us to improve on this methodology in two key ways.

First, field-scale consumptive use estimates provide a more accurate benchmark for potential savings under the assumption of 100% reductions. Whereas USBR’s approach quantifies consumptive use as a residual difference between measured diversions and mostly unmeasured return flows at the district level, remote sensing of ET directly targets the quantity of interest: evaporation and transpiration that are consumptively used and leave the watershed. Moreover, field-scale estimation facilitates measurement of average water use on the subset of lands enrolled in the program, providing a more accurate measure of potential savings that can be re-scaled to be consistent with USBR’s district wide estimates.

Second, field-scale water use estimates enable researchers to use causal inference techniques to assess the performance of the program itself. In this study, reductions in water use causally related to the program are just 53% of baseline use, as opposed to the assumed 100%. Half of this gap is likely explained by the strategic selection by farmers of which field to fallow in response to calls—called fields have lower predicted water use on average, suggesting that some water savings are not “additional” to the program. The other half is explained by persistent ET on fallowed lands. Much of this ET is driven by spillovers in from neighboring irrigated parcels.

These findings relate to a long-standing debate in water policy and hydrology about how to conceptualize the potential water savings from changes in agricultural practices including both efficiency enhancements and fallowing (Frederiksen and Allen, 2011; Gleick et al., 2011; Frederiksen et al., 2012). Ultimately, the question is whether the consumptive water use that persists on fallowed fields should be deducted from savings estimates or not. Answering that question runs

requires grappling with inherent difficulties associated with scaling changes in field-scale water use to district or basin-level water savings (Lankford et al., 2020). However, regardless of how one conceptualizes this number, our results indicate that half of the water savings gap is explained by behavior, not spillovers. This factor alone would account for over 300,000 AF of water over the life of the program.

The program we study here is designed to free up water for urban customers of MWD. If they have over-estimated true water savings by 624,000 AF, who is absorbing the difference? The answer to that question turns out to be far from straightforward. The MWD/USBR savings estimates contained in this paper are published annually in a Verification Report for the Program, but these estimates do not directly alter MWD's diversion rights to Colorado River Water.¹ Instead, the USBR engages in a complex process of accounting at the end of each year that takes months to complete (Wobus et al., 2024), and MWD is the junior-most user of Colorado River Water in this process. Ultimately, MWD benefits from the fallowing program to the extent that it “moves the needle” sufficiently in USBR's accounting to enable MWD to divert more water. Whether the savings generated in PVID are able to be claimed by MWD depends on how water use changes for every other user of Lower Colorado River Basin water.

Our results underscore how the design and evaluation of various environmental pay-for-practice programs—deforestation and development avoidance, afforestation, crop switching, and fallowing—can be improved with the use of remote sensing to develop contract-scale estimates of both behavioral compliance *and* ultimate environmental benefits. Doing so could help close the gap between expected and realized benefits in these programs. Contract-scale benefit estimation allows managers to more accurately forecast program savings by accounting for baseline differences in land use across enrolled vs. unenrolled parcels, where they were previously blind to these differences. Perhaps more crucially, field-scale estimates enable the use of causal program evaluation to reliably estimate both behavioral compliance and environmental benefits *ex post*, facilitating program tweaks to address non-additionality, compliance gaps (for example through contingent contract terms), and potential spillovers.

Our study has several key limitations. First, our study area likely represents a “best-case scenario” for using total ET to measure consumptive agricultural use, because effective precipitation and natural ET are negligible in PVID. In other settings such as the Central Valley of California, techniques to separate natural ET from applied ET are more essential (Boser et al., 2024). Second, our setting benefits from a lack of groundwater pumping, making it straightforward to relate change in ET to changes in surface water use. In many other settings, it would prove more difficult to relate changes in ET to the supply of irrigation water, which may limit the utility of ET in facilitating transfers of water rights tied to specific sources.

¹These estimates are also used to calculate Intentionally Created Surplus (ICS) credits that MWD can virtually bank in Lake Mead and withdraw later. Hence, the estimation errors revealed by our analysis could inflate MWD's ICS balance and ultimately be associated with excessive drawdown of Lake Mead. However, MWD engages in a litany of programs to produce ICS credits that have hit the annual cap on credit creation almost every year, suggesting that it is relatively unlikely that errors in the savings estimates directly led to changes in the level of Lake Mead.

Materials and Methods

Field-Level Water Use Estimation

We combine two data sets to develop field-scale estimates of water use across PVID. The first is monthly data on evapotranspiration from OpenET (Volk et al., 2024). These data are monthly raster files with measures of total evapotranspiration at a 30-meter resolution. We utilize the OpenET Ensemble mean as well as the six individual ET models that make up the mean. The constituent models include: ALEXI/DisALEXVI v 0.0.32, eeMETRIC v 0.20.26, geeSEBAL v 0.2.2, PT-JPL v0.2.1, SIMS v 0.1.0, and SSEBop v 0.2.6.

We obtained polygons of irrigated fields from the Metropolitan Water District (MWD). These polygons were created by PVID to track “water toll acres” that actually receive deliveries of Colorado River Water. We used zonal statistics to calculate monthly mean ET for each field over 2000 to 2021. Monthly means are depicted in Panel (a) of Figure S1, and annual totals are depicted in Panel (b). Both panels also include MWD and USBR’s estimates of total water in each time period. MWD’s estimates are derived from annual Verification Reports on the fallowing program published by MWD. USBR estimates come from the USBR’s annual “River Accounting” reports.² Because field-level estimates are not possible under these methods, district-wide averages provide the best comparison.

Following Volk et al. (2024)’s comparison of OpenET estimates to ground-truth ET estimates from eddy covariance towers, we compare monthly total ET in PVID from each model to MWD’s estimates of consumptive use. We estimate the following model to evaluate the alignment of each ET estimate with MWD’s estimates:

$$MWD_{mt} = \beta_1 ET_{mt} + \varepsilon_{mt} \quad (1)$$

where MWD_{mt} is MWD’s report estimate of consumptive use in month m in year t , ET_{mt} is total ET in month m in year t (summed across all fields in our sample), and ε_{mt} is an error term. Results are presented in Table S2. Regardless of model, the R-squared exceeds 0.9 and the uncentered slope is near 1. Most slopes cluster around 1.3, suggesting that ET may *understate* consumptive agricultural use in PVID, consistent with Figure S1.

Total ET vs. Agricultural Use

We take two additional measures to probe the relationship between raw total ET and consumptive agricultural use in PVID. First, we assess the extent to which precipitation may confound estimates of consumptive agricultural use. To do so, we download monthly data on total precipitation at a 400-meter resolution from PRISM. We calculate total precipitation and “effective precipitation”—defined as $0.65 \times$ precipitation (Wobus et al., 2024) for each parcel in each month of the sample. Using these data, Figure S2 illustrates that effective precipitation is a small percentage (< 2) of

²<https://www.usbr.gov/lc/region/g4000/wtracct.html>

total ET (panels a and b) and that it varies inversely with the irrigation season (panel c). We also decompose the variation in precipitation into its cross-sectional and time-series components (panel d). Roughly 80% of the variation in precipitation is monthly, with under 20% being parcel-specific. Less than 2% of the overall variation in precipitation is unexplained by parcel and month-of-sample fixed effects, suggesting it is unlikely to affect our results.

The agricultural-ET decomposition follows [Boser et al. \(2022\)](#), who train a pixel-level machine-learning model to predict the agricultural component of evapotranspiration (\widehat{ET}_{ag}) and a total predicted ET (\widehat{ET}) from soil quality and water-holding capacity (gNATSGO), elevation, slope, and aspect (USGS DEM), reference and potential evapotranspiration, and seasonal climate covariates. Their training set consists of pixels confirmed to be fallow in a given year, identified by intersecting the USDA Cropland Data Layer with the California Department of Water Resources crop maps; the non-agricultural residual is then computed as actual ET minus \widehat{ET}_{ag} at each pixel-month.

We adapt the [Boser et al. \(2022\)](#) data-preparation pipeline to our setting and re-run the prediction step over our study period. To avoid contamination of training data with the following program we want to assess, we train the model in nearby Imperial Irrigation District (IID) and average over all available years. We apply the trained model to all PVID pixels for 2001–2023. Pixel-level predictions are aggregated to PVID parcels by spatial intersection and converted from millimeters per month to acre-feet per acre per month to match our master panel.

While this is a potentially useful method for isolating agricultural consumptive use, we caution that this approach was developed in the Central Valley of California and may not be appropriate for our setting for several reasons. First, there is far less scope for naturally occurring ET in Palo Verde due to limited precipitation and limited shallow groundwater. Second, as we demonstrate elsewhere in the paper, there is evidence of significant spillovers of water between fallow and non-fallow fields in this region where flood irrigation is common and soils are sandy. Hence, the ET signature on fallow fields that the [Boser et al.](#) approach attributes to natural factors in the Central Valley may actually be policy-relevant spillovers in our setting.

Figure [S13](#) and Table [S10](#) provide evidence that the [Boser et al. \(2024\)](#) may not perform as well in PVID. Consistent with the limited scope for natural ET in this setting, district-wide consumptive use reported by MWD is much more aligned with Ensemble ET aggregates than with estimates of Agricultural ET produced with the [Boser et al.](#) method. This finding is consistent with the conjecture that the [Boser et al.](#) is picking up and subsequently purging spillovers that are true agricultural water use, rather than natural ET from precipitation or groundwater. Nevertheless, we perform the [Boser et al.](#) calculation for each year of our sample and re-estimate our preferred specification (described below) on the transformed ET estimates.

Fallowing Calls

We obtained our data on which fields were selected for participation following “calls” each year from annual Verification Reports published by MWD. These reports contain information about MWD/USBR’s monthly estimates of total water use and total fallowed acreage for each year of

the program, as well as maps of fields that were selected for fallowing by farmers in response to that year's call. An example map is depicted in Figure S3. PVID created these maps by hand with a printed aerial map of the district and an overhead transparency—neither MWD nor PVID possessed spatial data describing fallowed fields. We were able to obtain verification maps for program years 2007–2021. We therefore exclude 2005 and 2006—the first two years of the program—from our sample.

We geo-referenced each map from all available verification reports and overlaid these images with the irrigated fields shapefile obtained from MWD. We tasked three separate research assistants to independently code “called” fields from each image. We also had the research assistants create indicator variables for instances where the shaded area on the map represented a partial subset of a polygon in the shapefile, or when image quality made it difficult to determine whether a specific polygon was part of call. From these underlying data, we create indicators for each field in each map to express whether one, two, or all three research assistants indicated that a field was i) called, ii) partial, and iii) of dubious quality. E.g., we have three separate definitions of treatment, three separate notions of which fields were only partially fallowed, and three separate notions of image quality.

The final step in creating a month-year panel of fallowing calls is to interpolate from the images in the verification reports to the intervening months. Typically, each annual verification report contains two images—one in July and one in January/December. From these snapshots, we must infer how to assign treatment to parcels in February-July and August-November. This is made possible by the structure of the program: MWD issues a fallowing call near the end of each calendar year that takes effect the following August and runs for two full years. E.g., a fallowing call issued in December 2012 would correspond to required fallowing from August 2013 through July 2015. With this logic, we can construct “fallowing spells” and interpolate forward and backward from a given image. For each image, we determine candidate 24-month windows that it could be a member of, and then assign treatment within the associated window.

This method is imperfect because the two-year spell duration means that there are two potential spells that could cause a given field to appear as called in the monthly maps we observe. However, the relative frequency of the snapshots allows us to discipline potential over-coding by forcing our assignment to agree with previous and subsequent images. We do this for each version of treatment described above. Figure S4 depicts monthly called acreage implied by our approach, in comparison to reported fallowing by MWD in the Verification Reports. The figure depicts three different conventions for aggregating RA coding into treatment based on whether one, two, or three RAs coded a field as treated. The most inclusive rule that assigns treatment if *any* RA selected a field results in an over-estimate of called acreage in many months, while the restrictive version requiring unanimous agreement tends to result in an under-estimate. Across most months, we obtain close alignment with MWD's official records using the majority rule that treats a field as “called” if at least two RAs coded it as such. We therefore adopt this majority rule as our preferred definition of treatment.

Other Data

We use several other data supplemental data sets in our analysis. We estimate distance to the Colorado River as well as the mean of [Schaetzl \(1986\)](#)'s soil drainage index for each field in the sample. We overlay fields with PLSS sections—square-mile grids that comprise the land survey system—to enable us to cluster at a coarser spatial unit than individual fields.

We also obtained assessor parcel data from Riverside County.³ Irrigated field polygons were spatially joined to the Riverside County assessor parcel layer, with each field assigned to the parcel contributing the largest intersection area. We then grouped parcels on their cleaned assessor mailing address (street + city, upper-cased), which serves as a proxy for common ownership across Assessor Parcel Numbers. The resulting crosswalk covers 2,500 irrigated fields linked to 302 unique owners, of whom 175 own more than one field (127 are single-field owners). This is an imperfect proxy for “farms” in at least two ways. First, it is a static snapshot of land ownership in 2025. Given that our panel begins in 2000, parcel ownership has likely changed over time. Second, if farmers lease acreage in or out, land ownership may not correspond to farm operations. Still, these data form a useful proxy for farm ownership when unpacking the extent of selection into following within vs. between landowners.

We also used USDA Cropland Data Layer (CDL) 30-m rasters (2007–2021) to produce per-field pixel count histograms by USDA crop code. Each field-year was assigned a dominant crop from the pixel counts for cotton, alfalfa, hay, and fallow. Wheat-dominant fields and CDL years 2008 and 2013 are dropped because nearly all PVID fields are misclassified in those cases. The resulting panel has 35,217 field-years across 2,709 fields and 13 years (2007–2021 excl. 2008, 2013), with 60.6% alfalfa, 24.9% fallow, 10.2% cotton, 2.0% hay, and 2.4% other. We then subset this panel to the 2,500 fields used for the main analysis. This field×year crop panel is then used to assign each onset field its most recent non-fallow pre-call crop and to define the modal crop for never-called fields. An important limitation of these data is that they are not available until 2007, two years after the program began.

Finally, we obtain hourly wind speeds from the Iowa Environmental Mesonet. Each hourly observation is decomposed into meteorological vector components ($u = -sknt \cdot \sin \theta$, $v = -sknt \cdot \cos \theta$), averaged by calendar month, and recombined into a prevailing direction via $\text{atan2}(-\bar{v}, -\bar{u})$ converted to a compass bearing; a consistency index is computed as the magnitude of the monthly resultant vector divided by the mean scalar speed (0 = chaotic, 1 = unidirectional). The resulting panel has 288 station-months (24 years × 12 months, 2000–2023), with a mean scalar speed of ~7.0 knots, mean consistency of ~0.49, and a mean prevailing direction of ~251° (WSW). This approach is depicted in Figure S11.

Summary statistics for core variables used in the analysis are reported for the enrolled vs. unenrolled group in Table S1.

³<https://gisopendata-countyofriverside.opendata.arcgis.com/datasets/CountyofRiverside::parcels-crest/about>

Difference-in-Difference Estimation

We estimate the causal impact of a fallowing “call” using a difference-in-difference framework that estimates changes in ET on called fields before vs. during a call compared to changes in un-called fields over the same time period. Recent advances have demonstrated that classic two-way fixed effects regressions can result in biased difference-in-difference (DiD) estimates when treatment is staggered (different units are treated at different times) or when treatment effects change over time. A variety of estimators have been proposed to address this problem. We employ the estimator proposed by [De Chaisemartin and d’Haultfoeuille \(2020\)](#), which constructs a weighted average of all valid comparisons of switchers to non-switchers in the data.

We prefer the [De Chaisemartin and d’Haultfoeuille \(2020\)](#) approach to the prevailing alternatives for several reasons. First and foremost, to our knowledge, this is the only modern DiD estimator that is able to deliver causal estimates when treatment can switch on and off, which it does frequently in our setting of rotational fallowing. Second, the double-robust implementation of the [De Chaisemartin and d’Haultfoeuille](#) estimator constructs a propensity-score weighted average of comparisons that is explicitly designed to address potential selection into treatment—a major concern in our setting. The identifying assumption behind the [De Chaisemartin and d’Haultfoeuille](#) is that there are parallel trends in untreated outcomes between the treated and untreated groups. While the potential untreated outcome for the group that ultimately receives treatment cannot be observed, [De Chaisemartin and d’Haultfoeuille](#) argue that placebo tests to rule out differential pre-trends provide evidence in support of the identifying assumption of the estimator. After discussing our implementation of the estimator, we present the results of these tests.

Ultimately, we are interested in comparing *annual* water savings to MWD’s estimates. However, because fallowing calls start in August and end in July two years later, treatment status can vary within a calendar year. We therefore re-aggregate the monthly data into “program years” that run August–July. We implement the dynamic version of the [De Chaisemartin and d’Haultfoeuille \(2020\)](#) estimator, which allows treatment effects to depend on past treatment status, and also accommodates owner-by-year fixed effects.

We apply the estimator to our setting as follows. Let i index PVID fields, t index program years (Aug–Jul), and $o(i)$ denote the assessor-inferred owner of field i . Let Y_{it} be annual ET in acre-feet per acre, and let $D_{it} \in \{0, 1\}$ be the staggered call indicator. For switching fields, define the onset year $F_i = \min\{t : D_{it} = 1\}$ (with $F_i = \infty$ for never-called fields). The [De Chaisemartin and d’Haultfoeuille](#) dynamic estimator reports, for each post-onset horizon $\ell = 1, \dots, 6$, the cohort-weighted average of within-owner switcher-vs.-stayer comparisons,

$$\hat{\delta}_\ell = \sum_g w_{g,\ell} \underbrace{\left\{ \mathbb{E}[Y_{i,g+\ell} - Y_{i,g-1} \mid F_i = g, o(i)] - \mathbb{E}[Y_{i,g+\ell} - Y_{i,g-1} \mid F_i > g + \ell, o(i)] \right\}}_{\text{DID}_{g,\ell} \text{ within owner cell}}, \quad (2)$$

where the outer expectation averages over owners o and the weights $w_{g,\ell}$ are proportional to the number of switchers in cohort g observed at horizon ℓ .

The inclusion of owner-by-year fixed effects ensures that every comparison in (2) is executed within an owner \times program-year cell, so that $\widehat{\delta}_\ell$ differences out owner-specific shocks of arbitrary shape. Hence, the placebo leads $\widehat{\pi}_k$ for $k = 1, \dots, 6$ replace $Y_{i,g+\ell} - Y_{i,g-1}$ in (2) with $Y_{i,g-k} - Y_{i,g-k-1}$, allow us to test for parallel trends by comparing pre-onset outcome changes between switchers and not-yet-treated fields in the same owner cell. Standard errors are clustered at the PLSS section level. Our estimates and subsequent calculations in the main text focus on the average total effect

$$\widehat{\delta}^{\text{ATT}} = \frac{1}{6} \sum_{\ell=1}^6 \widehat{\delta}_\ell, \quad (3)$$

The preferred specification emphasized in the main text utilizes the majority convention for coding treatment (a field is “called” if at least two RAs coded it as such) and a more restrictive sample quality criterion that drops fields that *any* RA flagged as only partially called or of dubious image quality. We perform this drop at the call-spell level, so if a parcel is ever flagged during a set of months where it is coded as called, we omit it from the entirety of that spell. Table S3 reports the results of estimating Equation 2 on all seven ET models, with the average total effect, individual placebo coefficients, and the p-val associated with the null hypothesis that the placebo estimates are jointly equal to zero. Across four of seven models, we fail to reject the null of no pre-trends at the 10% level, and we fail to reject the null at the 5% level across all seven models.

Robustness

We run 42 alternative specifications to illustrate the stability of our headline estimate and present the results in Figure S5. The specifications span seven OpenET models, three RA call-coding thresholds (any, majority, unanimous), and two control sets (none, +precipitation), holding identification related choices about the estimating sample fixed at their preferred values. Specifically, relaxing the data quality flag threshold destroys parallel trends across nearly every ET model (SI Table S11), and the SUTVA-relevant sample-restriction question is examined directly in SI Table S4 alongside a monthly visualization (SI Figure S8). Results are broadly similar across the 42 specifications, with the preferred specification sitting near the middle of the distribution; full per-cell estimates and placebo diagnostics are reported in SI Tables S12–S14.

Although the significance of the placebo tests varies across specifications, many estimates display an apparent decline in water use on treated parcels prior to the onset of the call. With period $t - 1$ differences normalized to zero, this manifests in positive placebo values for periods before $t - 1$. We posit that this apparent relative dip in period $(t - 1)$ is mechanical in nature. Moving from production to complete fallowing entails a major shift in agricultural operations, and unless this occurs on the final day of a given period, there is likely to be some dip in ET in the period just before fallowing is set to begin. E.g, if a farmer is set to begin fallowing on August 1, unless they wait and clear their field literally overnight on July 31, there will be some decline in ET associated with the foregone production in preparation for fallowing.

We provide evidence on this posited mechanical dip in ET prior to the onset of a call in Figure S6. To do so, we subset the data to about-to-be called fields and never-called parcels and then compare average *monthly* water use across each group in the six months leading up the onset of a following spell. Average monthly water use is closely aligned and not statistically different across groups in January through April, with slight differences emerging in May and a major dip observed for the about-to-be-called group in June and July, consistent with harvesting and preparing fields for fallowing in August.

To understand why this dip emerges as early as May rather than solely in July, we repeat the analysis from Figure S6 separately for alfalfa vs. cotton and present the results in Figure S7. The alfalfa dip does not appear until July, whereas the trend for cotton emerges in May. This reflects differences in the timing of planting across these two crops—farmers are able to grow alfalfa and get a cutting in prior to August, but they must forgo cotton entirely for the season in anticipation of an August call.

To assess whether within-farm reallocation of water onto uncalled fields biases our headline estimate—an identification violation in which control fields are partially treated—we conduct two complementary checks. First, in Figure S8, we plot mean monthly ET around August call onset for three mutually exclusive groups of field-months: fields currently in a clean call window; uncalled fields whose owner had another field in a call window that month; and uncalled fields whose owner had no concurrent call. For the two uncalled groups we drop field-months that occur after any earlier call on the same field, so that post-call carryover does not contaminate the contrast. Second, we re-estimated the De Chaisemartin and d’Haultfoeuille (2020) dynamic treatment effect on three alternative samples: (i) dropping never-called fields entirely; (ii) dropping uncalled fields owned by ever-called farmers and removing the non-parametric owner-trend absorption (which otherwise mechanically collapses this specification to (i)); and (iii) removing the owner-trend absorption alone (Table S4). Across both checks we find that, if anything, water use also declines on uncalled fields. The alternative point estimates all fall within 22% of the headline. They are modestly larger in magnitude on average, indicating that the headline may slightly understate per-field savings. However, the central finding that observed conservation falls well below the program’s reported savings is unchanged.

Comparison to MWD Estimates and Decomposition of Differences

The difference between MWD/USBR’s estimates of average water savings and our estimates is driven by several factors. Here, we decompose and characterize each in turn. First, there may be overall level differences in average water use estimates. Even if we employed MWD/USBR’s method of predicting fallowing savings based on district-wide average use, this could cause our estimates to diverge from theirs. Second, MWD/USBR’s method applies a district-wide average, implicitly assuming that average water use is the same across fields enrolled in the program and unenrolled fields. Even if called fields reduce their water use by 100%, compositional differences in enrolled vs. unenrolled fields would lead to difference between MWD/USBR’s estimate and a

causal estimate of water savings.

Third, MWD/USBR’s method assumes that water use declines by 100% when a field is called to fallow. Even though MWD does carefully enforce fallowing calls and include images to this effect in their Verification Reports, water savings on called fields may not be 100% of baseline water use—e.g., there may still be a “savings gap.” This savings gap can be decomposed into i) residual water use on called fields (water use does not go to zero) and ii) “non-additionality,” or the strategic selection of fields that would have had lower water use even in the absence of a call (some fields that were fallowed would have been fallowed in the absence of a call, resulting in zero net savings relative to a counterfactual no-call state).

To be precise, the differences between the MWD/USBR approach and our approach can be written as:

$$\text{Difference} = S_{\text{MWD}} - S_{\text{Actual}} = (\Omega_{\text{Fallow}} + \Omega_{\text{Add}}) - (\Delta_{\text{Meas}} + \Delta_{\text{Select}})$$

$$S_{\text{Actual}} = S_{\text{MWD}} + \underbrace{\Delta_{\text{Meas}} + \Delta_{\text{Select}}}_{\text{Measurement Differences}} - \underbrace{\Omega_{\text{Fallow}} - \Omega_{\text{Add}}}_{\text{Savings Gap}}$$

Where the components are defined as:

- $\Delta_{\text{Meas}} = ET_{\text{avg}} - S_{\text{MWD}}$ (Measurement adjustment)
- $\Delta_{\text{Select}} = ET_{\text{enrolled}} - ET_{\text{avg}}$ (Enrollment selection bonus)
- $\Omega_{\text{Fallow}} = ET_{\text{residual}}$ (Residual water use)
- $\Omega_{\text{Add}} = (ET_{\text{enrolled}} - \Omega_{\text{Fallow}}) - S_{\text{Actual}}$ (Additionality gap)

Each of these components can be calculated directly from the data to better understand why the causal estimate differs from the MWD/USBR estimate by roughly 50%. As a practical matter, Δ_{Meas} and Δ_{Select} are quite small. MWD/USBR’s average annual savings estimate is 4.67 AF per fallowed acre, whereas we estimate district-wide average water use to be 4.63 AFA, implying that $\Delta_{\text{Meas}} = 0.06$. As depicted in Figure S9, average water use on enrolled parcels is slightly but systematically higher than on unenrolled parcels. Average annual water use on enrolled parcels is 4.77 AFA, implying that $\Delta_{\text{Select}} = 0.14$, also quite small.

After accounting for measurement differences and selection into the program, potential water savings from ET-based measurement are 4.77 AFA per year. The causal estimate of 2.53 implies a 2.24 AFA savings gap which can be further decomposed into Ω_{Fallow} and Ω_{Add} . Moreover, Ω_{Add} can be further decomposed into a measurable between-farmer component, a measurable within-farmer component, and an unexplained residual component. The next section describes our approach to quantifying each component of the savings gap.

Breaking Down the Savings Gap

Our approach to unpacking the savings gap is to first focus on directly measurable residual water use on called fields. After documenting non-trivial water use on called fields, we explore several potential explanations. We then turn our attention to the remaining portion of the gap and assess the extent to which modeling predicted water use on about-to-be-called fields can explain the remainder of the gap that we have attributed to non-additionality. The residual from these two exercises is non-additionality driven by factors that our model cannot predict (such as private farmer information). For both exercises, we focus our attention on enrolled parcels only.

Residual ET

Panel (a) of Figure 3 depicts the distribution of average water use across called fields in each year for which we have call data. A significant mass of called fields have non-zero ET, with an overall mean of 1.18 AFA across fields and years. Hence, roughly 54% of the savings gap can be explained by residual ET on called fields. Contrary to MWD/USBR's assumption that water use declines by 100% in response to a call, water use remains roughly 25% of baseline even after a field is called.

Figure S10 presents several variables that could explain residual water use: distance to the Colorado River and Schaetzl (1986)'s drainage index, as well as the cumulative length of a given fallowing spell. None of these variables has significant power in explaining the variation in ET on called fields. We also explore whether spillovers from irrigation on neighboring parcels seep onto fallowed lands. We estimate average ET across each fallowed parcel's neighbors on a monthly basis and provide a binned scatter plot in panel (b) of Figure 3.

Table S5 reports the results of several regressions of residual ET on each potential explanatory variable depicted in Figure S10, as well as water use on neighboring fields (depicted in panel b of Figure 3). There is a strong positive relationship only between neighbor ET and own ET, with a slope of about 0.34, suggesting that surface spillovers and seepage are a major driver of residual ET. Average annual neighbor water use for called fields is 4.43 AFA. Hence, the applying the monthly linear prediction for ET on a fallow field would be $0.34 \times 4.43 = 1.5$, implying that the magnitude of the estimated spillovers could fully explain residual ET on called fields.

One possible concern is that these estimated neighbor effects could reflect measurement error in ET itself, rather than ground-level spillovers. ET models work by using variation in surface temperature to infer evaporation and transpiration associated with evaporative cooling. It may be that there are low-level atmospheric spillovers associated with evaporative cooling, such that the air around a fallow field is cooler when its neighbors are evaporating more, leading to a false ET signature.

We assess the plausibility of this alternative mechanism by exploiting variation in wind direction. The intuition is that evaporative/atmospheric spillovers should be stronger from upwind neighbors than from downwind neighbors. If neighbor effects are not correlated with wind direction, it suggests the effects are indeed ground-level spillovers. We calculate prevailing winds

for each month in the sample by obtaining hourly wind speeds and directions and constructing a vector average for each month. We then construct a field’s average upwind vs. downwind ET by taking the cosine of the difference between the angle between two parcel’s centroids and the angle of the prevailing wind for a given month. The assigned weight is then the maximum of the cosine and zero. This means that a parcel directly downwind of its neighbor would receive a weight of one, whereas a parcel at exactly 90 degrees would obtain a weigh of zero. Panel (c) of Figure 3 illustrates.

Panel (d) of Figure 3 demonstrates that neighbor spillovers are not correlated with wind direction. This can be seen visually in the figure, which plots a fallow parcel’s own ET against its upwind vs. downwind neighbors. Regression results presented below the figure confirm this by regressing fallow ET on total neighbor ET and upwind vs. downwind neighbor ET. Table S6 presents regression results which confirm that upwind and downwind neighbors do not have a differential effect on called parcels’ ET. The p-value on the hypothesis test that upwind and downwind ET matter equally ranges from 0.57 to 0.70. Moreover, the sum of upwind and downwind effects is roughly equal to the total effect.

A non-rejection of the upwind=downwind null is consistent with either the absence of an atmospheric channel or with the test being underpowered, so we quantify what the design could have detected. Taking the cluster-robust standard error of the upwind–downwind contrast from column (3) of SI Table S6, the minimum detectable effect (MDE) at conventional thresholds is $MDE = (z_{1-\alpha/2} + z_{1-\beta}) \widehat{SE}(\hat{\beta}_{up} - \hat{\beta}_{down})$. At 80% power and $\alpha = 0.05$ (two-sided), this yields an MDE of 0.024 AFA, or roughly 5.9% of the total neighbor-ET slope; the observed contrast is 0.004 ($p = 0.64$). The non-rejection therefore bounds any atmospheric contribution to the neighbor spillover well below the magnitude that would be required for it to be the dominant mechanism. The full set of MDEs across (α , power) combinations is reported in SI Table S7.

Non-Additionality

To isolate the mechanisms driving non-additionality in the water fallowing program, we estimate the extent to which farmers strategically select specific fields for enrollment based on heterogeneous water-use expectations. The total observed non-additionality gap—the difference between the total savings gap (2.24) and observed residual ET on called fields (1.18)—is 1.06.

We decompose this 1.06 AFA gap into three distinct factors: (1) observable between-farmer selection, (2) observable within-farmer selection, and (3) private, unobservable transient information. We achieve this by horse-racing a baseline autoregressive model against a cross-validated machine learning approach, and algebraically partitioning the resulting selection effects.

We first estimate a counterfactual ET for all field-months using an AR(1) specification on historically uncalled fields. This model captures baseline field-specific seasonality and mean reversion:

$$ET_{imt} = \alpha_{im} + \beta(ET_{i,t-12} - \alpha_{im}) + \varepsilon_{imt} \quad (4)$$

where α_{im} represents field-by-month fixed effects. This specification allows every field to exhibit a unique seasonal water-use curve (e.g., distinguishing between summer-peaking alfalfa and spring-peaking winter wheat), adjusting for a 12-month lag to avoid seasonal confounding.

To ensure our estimate of structural selection is not downwardly biased by the functional form assumptions of the AR(1) model, we construct an upper bound on observable predictability using a cross-validated LASSO model. This approach includes the field’s seasonal baseline, extended 24-month lags to capture multi-year crop rotations, spatial neighbor histories (average 12-month lagged ET of adjacent fields), and year fixed effects to capture district-wide weather or macroeconomic shocks. The machine learning model is tuned via 5-fold cross-validation to select the optimal penalty parameter, preventing overfitting while maximizing out-of-sample predictive power. The regression coefficients reported in Table S8 confirm that predicted counterfactual ET is a statistically significant predictor of the call decision: fields with lower predicted water use are more likely to be selected for a call, both between and within owners, across both the AR(1) and LASSO models. Overall predictive accuracy of the two models is broadly comparable—the LASSO does not dramatically improve out-of-sample fit relative to the AR(1)—but the additional dimensions it captures (long-horizon cycles, spatial structure, aggregate shocks) are precisely those that correlate with the call decision. The decomposition below makes this concrete.

For both models, we aggregate the predicted counterfactual ET to the program-year level to quantify *total predicted non-additionality* as the difference between average enrolled water use and either AR(1) or LASSO-predicted water use called fields each year. Let \hat{Y}_{ift} denote the predicted counterfactual ET for field i in owner f ’s portfolio in program-year t ; let $C_{ift} \in \{0, 1\}$ indicate whether the field is called for fallowing in that year. The total observable selection effect is the difference between the predicted ET on uncalled and called field-years among ever-enrolled fields:

$$\Delta_{\text{total}} = \mathbb{E}\left[\hat{Y}_{ift} \mid C_{ift} = 0\right] - \mathbb{E}\left[\hat{Y}_{ift} \mid C_{ift} = 1\right]. \quad (5)$$

We can then decompose the total predicted non-additionality exactly into a between- and a within-owner component. Let \bar{Y}_{ft} denote the owner-year portfolio average of \hat{Y} across owner f ’s enrolled fields in year t . Adding and subtracting the owner-year portfolio mean conditional on a call partitions Δ_{total} exactly into two interpretable margins:

$$\Delta_{\text{between}} = \mathbb{E}\left[\hat{Y}_{ift} \mid C = 0\right] - \mathbb{E}\left[\bar{Y}_{ft} \mid C = 1\right], \quad (6)$$

$$\Delta_{\text{within}} = \mathbb{E}\left[\bar{Y}_{ft} \mid C = 1\right] - \mathbb{E}\left[\hat{Y}_{ift} \mid C = 1\right], \quad (7)$$

$$\Delta_{\text{total}} = \Delta_{\text{between}} + \Delta_{\text{within}}. \quad (8)$$

Δ_{between} compares the enrolled-baseline mean of predicted ET (uncalled field-years) to the average owner-portfolio mean across called field-years. Each called field-year contributes equally to that average, so owners enter the called-side expectation in proportion to the number of called fields they have. Because the program calls a uniform fraction of every enrolled owner’s land in a

given year rather than activating some owners and not others, this margin does not identify year-to-year owner-level participation. It captures the compositional effect of differential enrollment intensity—owners who enrolled a larger absolute share of their land contribute more called field-years. Hence, the between-owner component is driven by any systematic relationship between enrollment intensity and portfolio ET.

Δ_{within} is the average gap, within an owner-year, between that owner’s full portfolio mean and the predicted ET of the specific fields the owner called. It captures the micro-level strategic sorting of land within a farm.

The residual non-additionality not absorbed by a given prediction model is

$$\Delta_{\text{unobs}} = \text{Gap} - \Delta_{\text{total}},$$

where $\text{Gap} = 1.06$ AFA is the share of the total savings gap remaining after subtracting residual ET on called fields. Under the LASSO prediction, Δ_{unobs} serves as a bound on the contribution of farmers’ private information (beyond what we condition on) to selection.

Applying the two sets of predictions to equations 5–6 reveals three distinct tiers of adverse selection. The baseline AR(1) model explains 0.44 AFA (panel (c) of Fig. 4), or roughly 42% of the 1.06 AFA non-additionality component of the savings gap. This indicates that farmers systematically rely on basic historical performance when making following decisions, selecting fields with historically lower baselines. This effect is driven primarily by within-owner field sorting (0.37 AFA) rather than between-owner portfolio selection (0.07 AFA).

The LASSO model explains 0.73 AFA (panel (d) of Fig. 4), or roughly 69% of the 1.06 AFA non-additionality component of the savings gap. The improvement is heavily concentrated in the between-owner margin, which more than triples from 0.07 AFA to 0.23 AFA, while the within-owner margin rises more modestly from 0.37 AFA to 0.50 AFA. The growth in the between-margin does not reflect a year-to-year participation decision (every enrolled owner is called for the same fraction of their enrolled acreage), but rather the LASSO’s improved ability to capture owner-level structural variation in ET—through multi-year lags, neighbor histories, and aggregate month-by-year shocks—that the AR(1) cannot. Mechanically, this owner-level predictable variation enters the between-margin because owners enrolled different absolute shares of their land in the program: owners with larger enrolled portfolios contribute more called field-years, so any correlation between enrollment intensity and portfolio ET manifests as between-owner selection.

Subtracting the LASSO prediction from the true causal gap leaves a residual unobservable non-additionality of 0.33 AFA, or roughly 31% of the gap. Because this remainder cannot be predicted even by a flexible machine learning algorithm armed with extensive spatial and historical covariates, it serves as a lower bound—conditional on the covariate space we consider—on the contribution of farmers’ private information to selection. Approximately one-third of the program’s non-additionality therefore appears to be driven by factors not recoverable from our historical or spatial covariates and that program administrators relying on the same data would be unable to anticipate.

Aggregate Implications

The estimates in Section imply a sustained gap between realized water savings from following in PVID and the savings MWD and USBR have reported on the same acres. MWD and PVID both diver water from the Colorado River, which raises questions about where within the Lower Colorado River Basin the difference is being realized. Here, we quantify the overall magnitude of that gap and discuss the implications for water users in the Lower Colorado River Basin under a variety of assumptions about which users within system absorb the difference. All estimates in this section combine the preferred dynamic DiD per-acre savings rate ($\hat{\delta}^{\text{ATT}} = -2.53$ AFA) with MWD’s own contracted-acre-year ledger from the Verification Reports over the program window (program years 2005–2021; 17 years).

Aggregate shortfall. Across this window, MWD’s reports credit the program with **1,437,527 AF** of saved water. Applying our DiD rate to the same acre-year ledger yields **812,611 AF** of physically defensible savings, leaving a cumulative shortfall of

$$\Delta\text{AF} = \underbrace{1,437,527}_{\text{MWD reported}} - \underbrace{812,611}_{\text{DiD implied}} = 624,916 \text{ AF}, \quad (9)$$

or **43.5%** of the program’s claimed yield. On an annual basis the shortfall averages **36,760 AF per program year**, equivalent to roughly **2.2%** of MWD’s ~ 1.7 MAF total annual deliveries. The cumulative gap traced out in Figure S12 of Fig. 1 is exactly this 624,916 AF spike by construction.

Characterizing the implications of this difference in the context of Lower Basin river accounting is not straightforward. While the Verification Reports provide a detailed accounting of MWD/USBR’s initial *estimate* of savings from the program in each year, these numbers do not necessarily translate directly into increased deliveries to MWD. Instead, as the junior-most claimant on the Lower Colorado River, MWD’s annual deliveries from the USBR depend on water use not only in PVID, but in several other irrigation districts in the LCRB. It is entirely possible that any savings generated by the PVID program are either augmented or offset by changes in water use across these other users in the basin who are senior to MWD. This complex accounting takes months to resolve each year (Wobus et al., 2024), and makes it nearly impossible to directly attribute estimation errors in PVID to water deliveries to MWD.

Rather than attempt to fully resolve LCRB accounting, we provide two bounds on the potential impacts of the 625,000 AF gap between MWD’s estimates and realized savings. One interpretation is that this difference is passed through to MWD on a one-to-one basis. E.g., if USBR’s year-end water accounting is perfectly accurate, it should capture actual, rather than estimated water savings. In this idealized scenario, each AF of reduced savings is an AF of reduced deliveries to MWD. Under this assumption, we provide several bounds on the financial implications for MWD of paying for water that they are not ultimately receiving.

On the other extreme, the estimation error may not be passed on at all. In other words, we can quantify the impacts of the estimation error on the overall LCRB supply under the assumption

that MWD’s deliveries increase on a one-to-one basis with the *reported* savings in the Verification Reports in each year. Under this edge case, every AF of mis-estimated savings is an extra AF of water that was diverted from Lake Mead. We emphasize that this provides an extreme upper bound on the aggregate implications of the estimation error, and do not intend to suggest that this is the most accurate understanding of the program’s implications. Likely, the reality lies in between these two extremes.

Financial bound. The first panel of Figure 5 also provides some additional context for these estimates by comparing MWD/USBR’s estimated cost per acre-foot of the program, actual cost per acre-foot implied by our estimates, and cost estimates for alternative supply augmentation or demand management strategies facing municipal providers such as MWD. Estimates of the cost of social nudges come from [Bernedo et al. \(2014\)](#); [Price et al. \(2014\)](#), whereas xeriscaping estimates come from [Brelsford and Abbott \(2021\)](#).

Multiplying the shortfall by four reference per-AF prices brackets what the over-reported savings would have cost MWD if it had to acquire an equivalent volume from each of its alternative supply margins. The table in Figure 5 reports the resulting dollar values; the four anchor prices are MWD’s *own* program cost (program-cost / claimed-AF \approx \$146/AF), the most recent IID–San Diego County Water Authority transfer rate (\$296/AF), MWD’s Tier 1 untreated rate (\$766/AF), and Carlsbad desalination as a marginal manufactured-supply alternative (\$3,000/AF). At MWD’s own cost the over-reporting represents a **\$91 million** mismatch between dollars spent and physical water delivered; at the IID–SDCWA transfer price it is **\$185 million**; at MWD’s Tier 1 marginal supply price it is **\$479 million**; and at the desalination margin — the price MWD’s customers pay for the next manufactured AF — it reaches **\$1.87 billion**. We treat the desalination figure as an upper bound on opportunity cost rather than a realized expenditure.

Physical bound: Lake Mead. The natural physical denominator for foregone basin storage is Lake Mead elevation. Because Lake Mead’s volume-per-foot of stage rises sharply with elevation — from $\sim 38,000$ AF/ft near dead pool to $\sim 148,000$ AF/ft near full pool — a flat conversion overstates the elevation contribution of any shortfall accumulated at the low-stage levels that prevailed during this window. We therefore evaluate the conversion year-by-year. For each program year we attach the observed end-of-calendar-year Hoover Dam stage from USBR Boulder Canyon Operations Office records (mean stage 1,098.4 ft over 2005–2021) and apply the USBR area–capacity table, linearly interpolating between published anchor points to obtain the volume-per-foot at each year’s stage (mean 83,440 AF/ft over the window). Summing the per-year contributions, the cumulative shortfall is equivalent to

$$\Delta h_{\text{Mead}} = \sum_{t=2005}^{2021} \frac{\Delta \text{AF}_t}{(\text{AF/ft})_t} \approx 7.6 \text{ ft of Lake Mead elevation.} \quad (10)$$

We treat (10) as a strict upper bound on basin-storage impacts: it assumes 1:1 pass-through of every shortfall AF into Lake Mead.

References

- Arellano-Gonzalez, J., A. AghaKouchak, M. C. Levy, Y. Qin, J. Burney, S. J. Davis, and F. C. Moore (2021). The adaptive benefits of agricultural water markets in California. *Environmental Research Letters* 16(4), 044036.
- Avila, P., M. Nemati, D. Crespo, A. Dinar, Z. Frankel, and N. Halberg (2025, October). Public spending and water scarcity: An empirical analysis of usbr investments in the Colorado River basin. *JAWRA Journal of the American Water Resources Association* 61(5), e70042.
- Ayres, A., K. Meng, and A. Plantinga (2021, October). Do environmental markets improve on open access? Evidence from California groundwater rights. *Journal of Political Economy* 129(10).
- Bernedo, M., P. J. Ferraro, and M. Price (2014). The persistent impacts of norm-based messaging and their implications for water conservation. *Journal of Consumer Policy* 37(3), 437–452.
- Boser, A., K. Caylor, A. Larsen, M. Pascolini-Campbell, J. Reager, and T. Carleton (2022). Field scale crop water consumption estimates reveal potential water savings in California agriculture. *Nature Communications* 15(2366).
- Boser, A., K. Caylor, A. Larsen, M. Pascolini-Campbell, J. T. Reager, and T. Carleton (2024). Field-scale crop water consumption estimates reveal potential water savings in California agriculture. *15(1)*, 2366.
- Brelsford, C. and J. K. Abbott (2021). How smart are ‘Water Smart Landscapes’? *Journal of Environmental Economics and Management* 106, 102402.
- Bretsen, S. N. and P. J. Hill (2008). Water markets as a tragedy of the anticommons. *Wm. & Mary Envtl. L. & Pol’y Rev.* 33, 723.
- Brewer, J., R. Glennon, A. Ker, and G. Libecap (2008). 2006 presidential address water markets in the west: prices, trading, and contractual forms. *Economic Inquiry* 46(2), 91–112.
- Bruno, E. M. and K. Jessoe (2021). Missing markets: Evidence on agricultural groundwater demand from volumetric pricing. *Journal of Public Economics* 196, 104374.
- Bruno, E. M. and K. Jessoe (2024). Designing water markets for climate change adaptation. *Nature Climate Change* 14(4), 331–339.
- De Chaisemartin, C. and X. d’Haultfoeuille (2020). Two-way fixed effects estimators with heterogeneous treatment effects. *American Economic Review* 110(9), 2964–2996.
- Dilling, L., J. Berggren, J. Henderson, and D. Kenney (2019). Savior of rural landscapes or Solomon’s choice? Colorado’s experiment with alternative transfer methods for water (atms). *Water Security* 6, 100027.

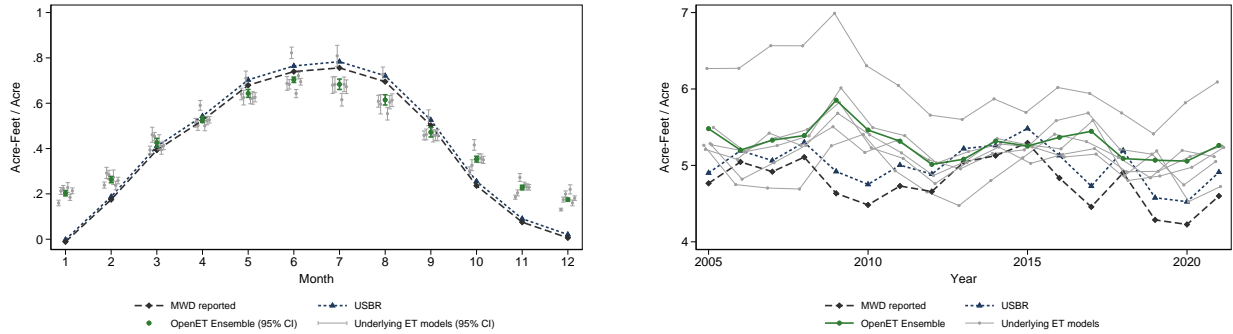
- Elliott, J., D. Deryng, C. Müller, K. Frieler, M. Konzmann, D. Gerten, M. Glotter, M. Flörke, Y. Wada, N. Best, et al. (2014). Constraints and potentials of future irrigation water availability on agricultural production under climate change. *Proceedings of the National Academy of Sciences* 111(9), 3239–3244.
- Frederiksen, H. D. and R. G. Allen (2011). A common basis for analysis, evaluation and comparison of offstream water uses. *Water International* 36(3), 266–282.
- Frederiksen, H. D., R. G. Allen, C. M. Burt, and C. Perry (2012). Responses to gleick et al., which was itself a response to frederiksen and allen. *Water International* 37(2), 183–197.
- Garrick, D. E., S. Balasubramanya, M. Beresford, A. Wutich, G. G. Gilson, I. Jorgensen, N. Brozović, M. Cox, X. Dai, S. Erfurth, R. Rimsàitè, J. Svensson, J. Talbot-Jones, H. Unnikrishnan, C. Wight, S. Villamayor-Tomas, and K. Vazquez Mendoza (2023). A systems perspective on water markets: Barriers, bright spots, and building blocks for the next generation. *Environmental Research Letters* 18(3), 031001.
- Gleick, P. H., J. Christian-Smith, and H. Cooley (2011). Water-use efficiency and productivity: rethinking the basin approach. *Water International* 36(7), 784–798.
- Gordon, B. L., G. F. Boisrame, R. W. Carroll, N. K. Ajami, B. Leonard, C. Albano, N. Mizukami, M. A. Andrade, E. Koebele, M. H. Taylor, et al. (2024). The essential role of local context in shaping risk and risk reduction strategies for snowmelt-dependent irrigated agriculture. *Earth's Future* 12(6), e2024EF004577.
- Grafton, R. Q., G. D. Libecap, E. C. Edwards, R. J. O'Brien, and C. Landry (2012). Comparative assessment of water markets: Insights from the Murray–Darling Basin of Australia and the Western USA. *Water Policy* 14(2), 175–193.
- Hagerty, N. (2023, 23 February). What holds back water markets? Transaction costs and the gains from trade. *Working Paper, Montana State University*.
- Lankford, B., A. Closas, J. Dalton, E. L. Gunn, T. Hess, J. W. Knox, S. Van der Kooij, J. Lautze, D. Molden, S. Orr, et al. (2020). A scale-based framework to understand the promises, pitfalls and paradoxes of irrigation efficiency to meet major water challenges. *Global Environmental Change* 65, 102182.
- Leonard, B., C. Costello, and G. Libecap (2019, Winter). Expanding water markets in the western United States: Barriers and lessons from other natural resource markets. *Review of Environmental Economics and Policy* 13(1), 43–61.
- Medellín-Azuara, J., A. Escriva-Bou, A. C. Gaudin, K. A. Schwabe, and D. A. Sumner (2024). Cultivating climate resilience in california agriculture: Adaptations to an increasingly volatile water future. *Proceedings of the National Academy of Sciences* 121(32), e2310079121.

- Price, J. I., J. M. Chermak, and J. Felardo (2014). Low-flow appliances and household water demand: An evaluation of demand-side management policy in Albuquerque, New Mexico. *Journal of Environmental Management* 133, 37–44.
- Rafey, W. (2023, February). Droughts, deluges, and (river) diversions: Valuing market-based water reallocation. *American Economic Review* 113(2).
- Rafey, W. (2026, 5). Measuring water misallocation in california. (35176).
- Richter, B. D., D. Bartak, P. Caldwell, K. F. Davis, P. Debaere, A. Y. Hoekstra, T. Li, L. Marston, R. McManamay, M. M. Mekonnen, et al. (2020). Water scarcity and fish imperilment driven by beef production. *Nature Sustainability* 3(4), 319–328.
- Richter, B. D., G. Lamsal, L. Marston, S. Dhakal, L. S. Sangha, R. R. Rushforth, D. Wei, B. L. Ruddell, K. F. Davis, A. Hernandez-Cruz, et al. (2024). New water accounting reveals why the colorado river no longer reaches the sea. *Communications Earth & Environment* 5(1), 134.
- Rodell, M., J. S. Famiglietti, D. N. Wiese, J. Reager, H. K. Beaudoin, F. W. Landerer, and M.-H. Lo (2018). Emerging trends in global freshwater availability. *Nature* 557(7707), 651–659.
- Schaetzl, R. J. (1986). Soilscape analysis of contrasting glacial terrains in wisconsin. *Annals of the Association of American Geographers* 76(3), 414–425.
- Volk, J. M., J. L. Huntington, F. S. Melton, R. Allen, M. Anderson, J. B. Fisher, A. Kilic, A. Ruhoff, G. B. Senay, B. Minor, et al. (2024). Assessing the accuracy of openet satellite-based evapotranspiration data to support water resource and land management applications. *Nature Water* 2(2), 193–205.
- Wobus, C., C. Nash, P. Culp, M. Kelly, and K. Kennedy (2024). Simplified agricultural water use accounting in the colorado river basin using openet. *Environmental Research Letters* 20(1), 014020.

Supplementary Material

Figures

Figure S1: District-Wide Consumptive Use: ET Models vs. MWD Reported

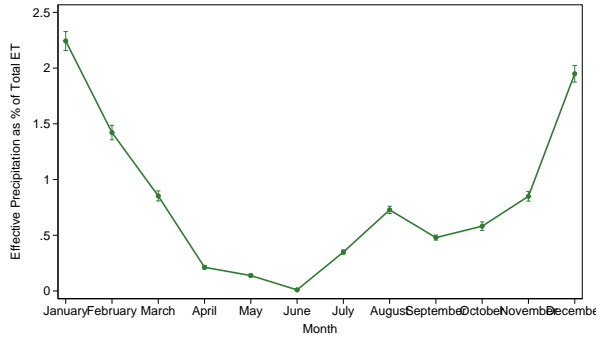


(a) Monthly AF Consumed (District-Wide)

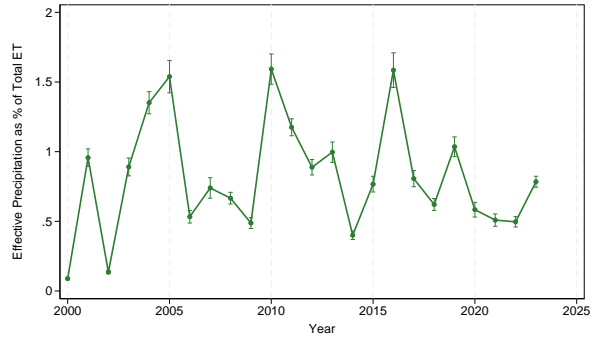
(b) Annual AF Consumed (District-Wide)

Notes: District-wide consumptive-use aggregates compare the seven OpenET models (ENSEMBLE, PTJPL, SIMS, SSEBOP, GEEBPM, EEMETRIC, DISALEXI) against MWD's reported consumptive-use aggregated from the annual USBR Decree Accounting Reports and MWD Verification Reports. Panel (a) plots monthly total AF consumed across all PVID fields; panel (b) plots the same aggregated to calendar year. The ensemble and per-model totals track the MWD series closely in the aggregate, motivating the use of field-level ET as the outcome for the DiD. Companion forces-through-origin regressions are reported in SI Table S2.

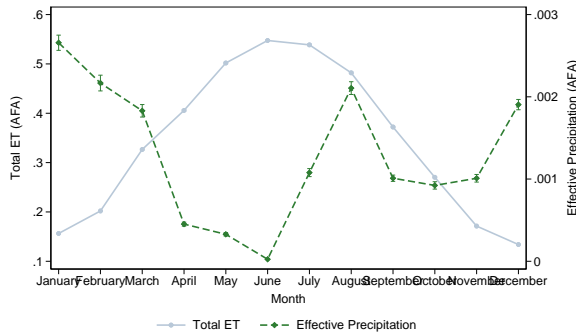
Figure S2: Effective Precipitation on Called Fields



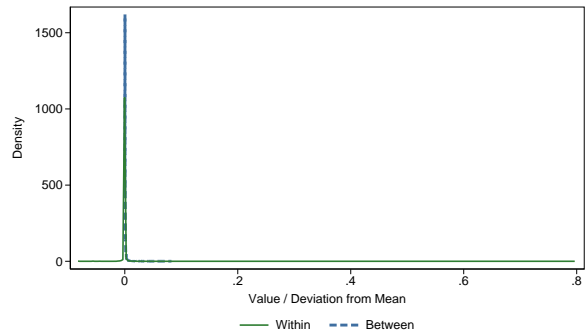
(a) Monthly Precipitation as % of ET



(b) Annual Precipitation as % of ET



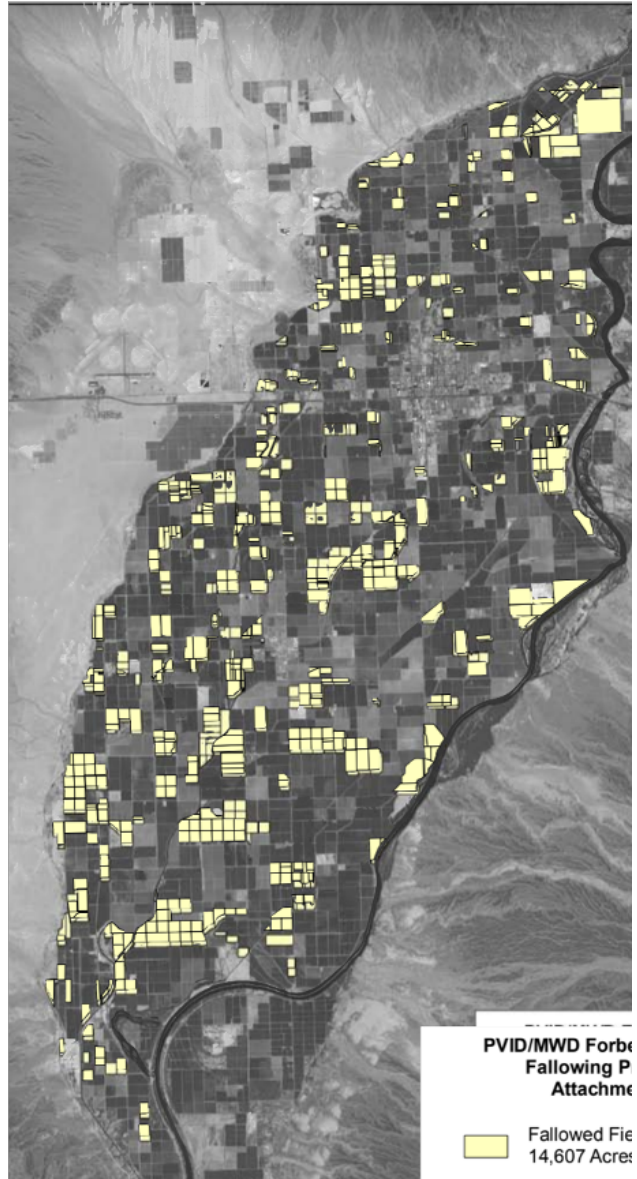
(c) Monthly Precipitation vs. ET



(d) Precipitation Decomposition of Residual ET

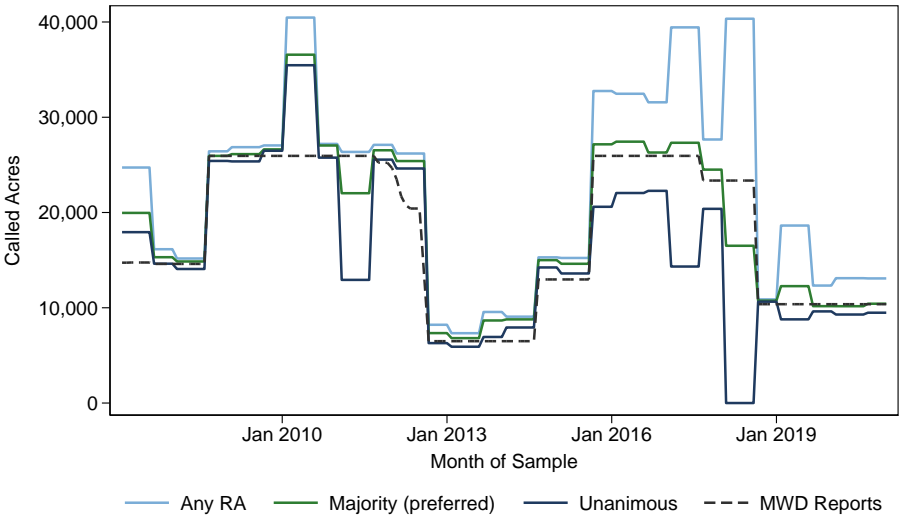
Notes: The PVID service area is hyper-arid (long-run precipitation ~ 9 cm/yr from PRISM), so effective precipitation is a trivial contributor to called-field ET. Panel (a) plots monthly PRISM precipitation as a share of monthly ET on called fields, pooled across parcel-months; panel (b) plots the same aggregated to prog-year. Even in the wettest months, rainfall explains less than 10% of observed ET; annual shares are an order of magnitude smaller. Panel (c) is a bivariate scatter of monthly precipitation against ET on called fields, with the best-fit line; the flat slope confirms precipitation does not covary materially with month-to-month ET variation in this sample. Panel (d) decomposes total residual ET on called fields into the portion explained by concurrent precipitation ($< 5\%$ AFA) and the unexplained remainder, which is the focus of the hydrologic-spillover analysis in the main text (Fig. 3).

Figure S3: Example MWD Verification-Report Map



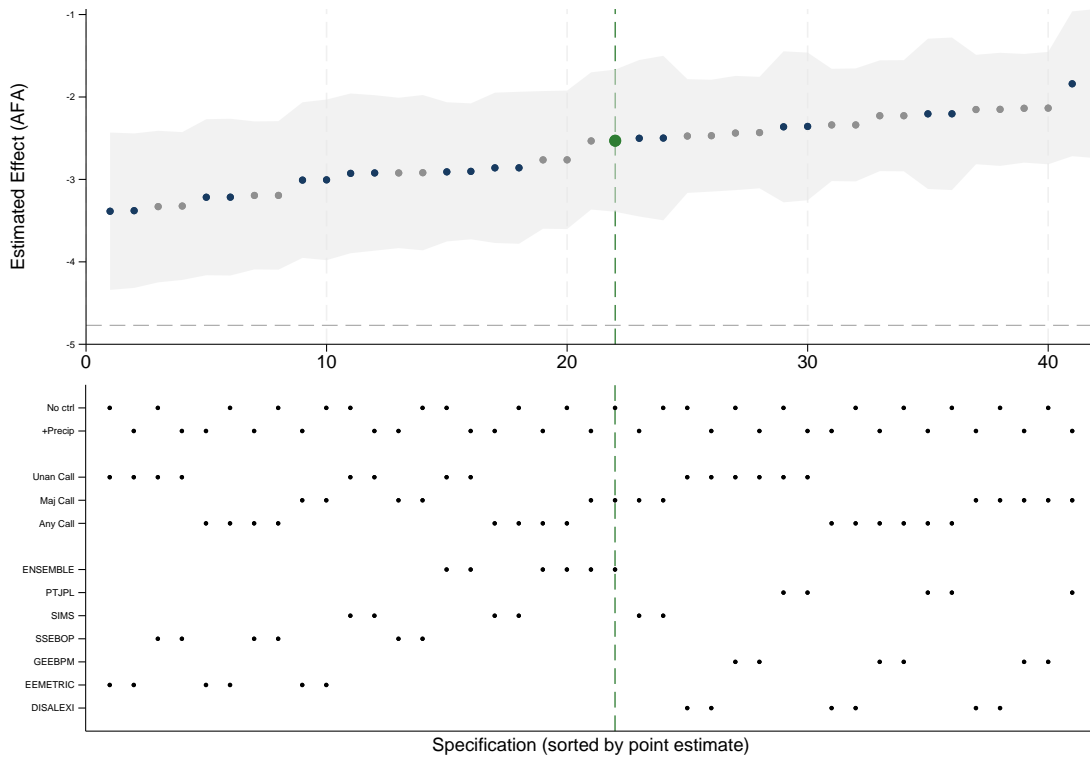
Notes: Representative page from an MWD Verification Report.

Figure S4: Called Acreage by RA Agreement Threshold vs. MWD Reported



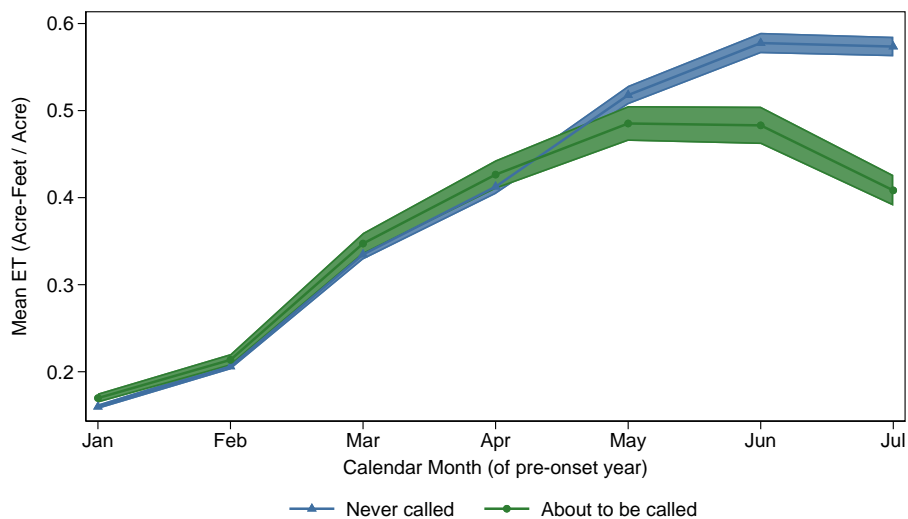
Notes: Monthly called acreage under the three RA agreement thresholds — any (light blue), majority (green; the preferred spec), unanimous (dark blue) — overlaid on the MWD-reported Fallowed Land acre-month series (dark gray dashed). The three RA series bracket the MWD series through the program window, with the majority (≥ 2 RAs) threshold tracking MWD most closely.

Figure S5: Specification Curve



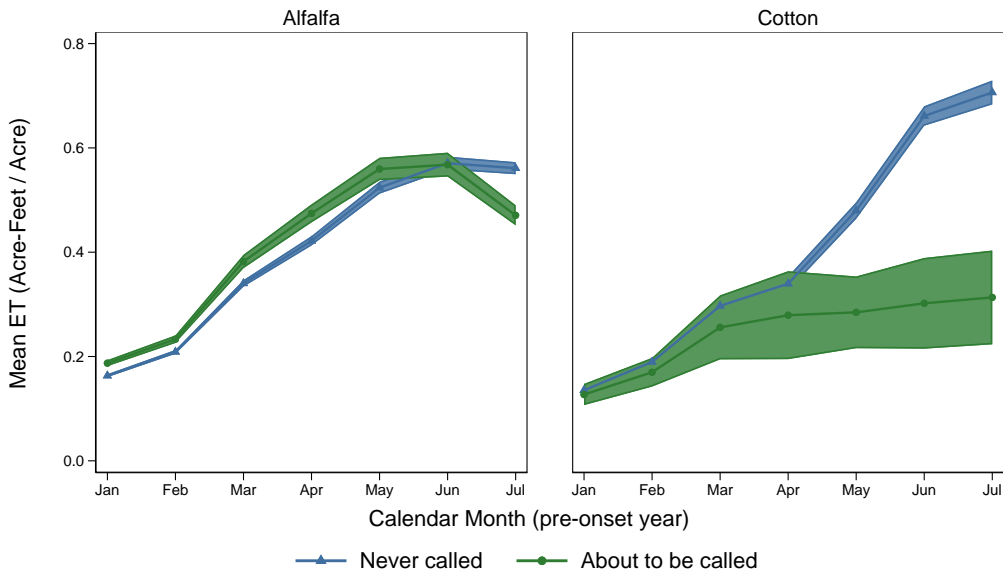
Notes: Coefficient + 95% CI for each of 42 specifications obtained by crossing seven OpenET models, three RA call thresholds (any, majority, unanimous), and two control sets (none / +precipitation). Identification-relevant choices are held fixed at the preferred values: sample = All, spatial-flag threshold ≥ 1 RA (the most conservative; see SI Table S11). Specifications sorted ascending by coefficient. Top-panel dot coloring distinguishes specifications that pass the joint-nullity placebo test at $p > 0.10$ (blue) from those that do not (grey); the preferred specification (ENSEMBLE, call = 2, no controls) is the green dot. The vertical green dashed line marks the preferred-spec rank and aligns with the indicator grid below. Horizontal dashed reference line at -4.77 AFA marks the maximum potential savings, equal to average water use on enrolled parcels.

Figure S6: Pre-Onset Monthly ET on About-to-Be-Called vs. Never-Called Fields



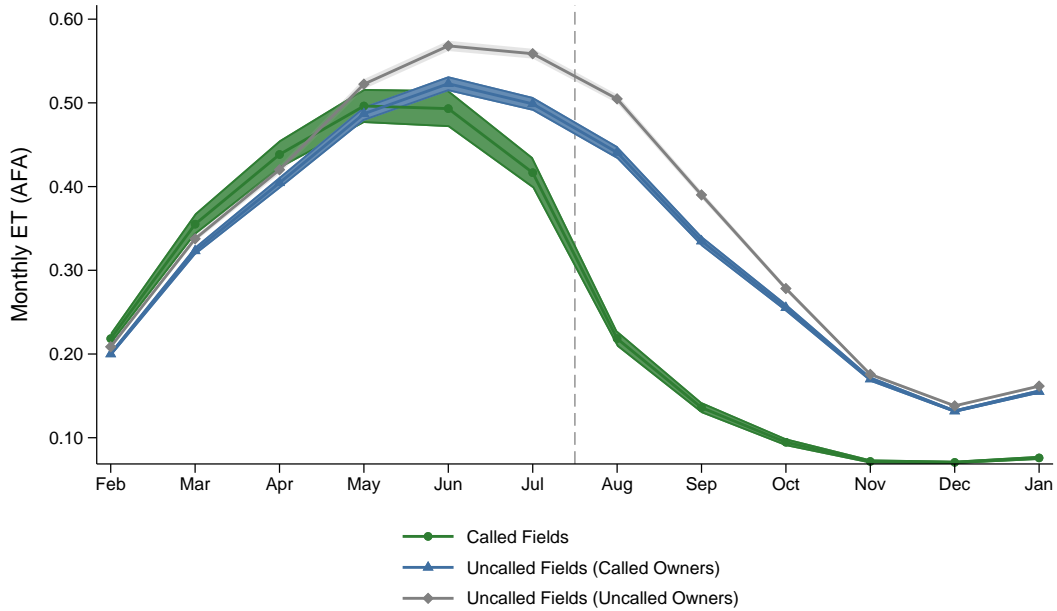
Notes: Monthly ET for fields in the six months *prior* to a call onset (green) versus fields that are never called (blue). 95% CIs on each point. The two series track one another January–May, then the about-to-be-called line drops sharply in June and July. This is the Jun/Jul following-prep phenomenon described in the main text.

Figure S7: Pre-Onset Monthly ET Drop by Crop Type



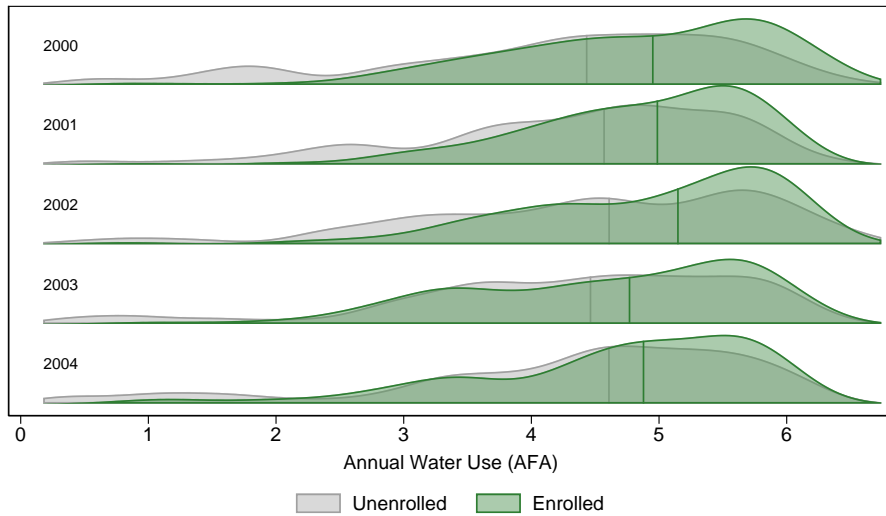
Notes: Monthly difference in ET (about-to-be-called vs. never-called) by CDL crop type. Crop classifications are from the USDA Cropland Data Layer matched to PVID field boundaries at the pre-call cropping decision (the CDL layer for the year before call onset). The June–July dip documented in SI Figure S6 is concentrated in the crops most amenable to mid-season termination — alfalfa/forage and annual row crops — and is absent in permanent crops, which are not plausibly subject to following-prep behavior.

Figure S8: Monthly ET around Call Onset: Called vs. Enrolled-Uncalled vs. Unenrolled



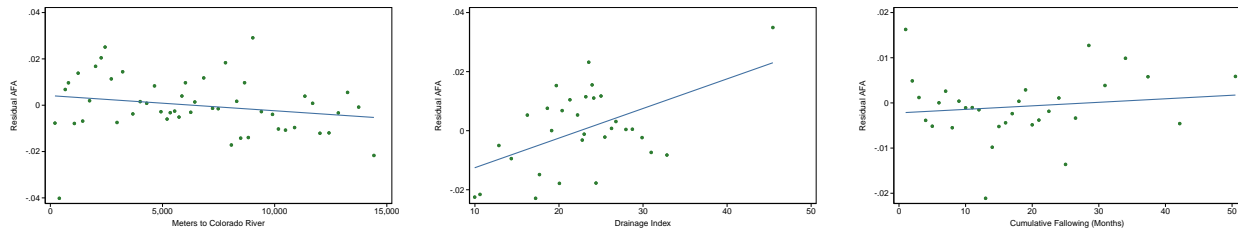
Notes: Monthly per-acre ET on three mutually exclusive groups of fields during a 12-month window around a August call onset. Called (green) are in-window observations of fields whose Aug-onset spell passes the spell-level quality screen. Uncalled fields associated with farmers who have some called fields are depicted in blue. Uncalled fields associated with farmers who have zero called field are depicted in grey.

Figure S9: Enrollment Selection: Pre-Program ET Distributions



Notes: Joy plot of pre-program annual ET distributions on fields that are ever enrolled in the following program (green) versus never enrolled (gray). Each row is one calendar year; the horizontal axis is AFA per field-year.

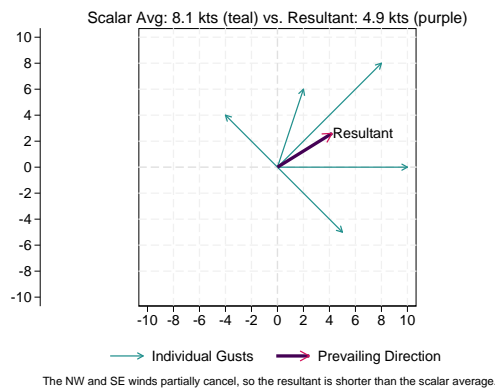
Figure S10: Residual-Water Mechanism Binscatters: Additional Covariates



(a) Distance to Colorado River (b) Drainage Index (c) Cumulative Call Spell Length

Notes: Additional binscatters of within-field residual ET (after partialling out field \times month-of-sample fixed effects) against three field-level drivers that the main-text Fig. 3 does not report. Panel (a) plots residual ET against Euclidean distance from the field’s centroid to the Colorado River mainstem. Panel (b) plots residual ET against a within-field Drainage Index (Schaetzl, 1986), which captures similar hydrologic structure. Panel (c) plots residual ET against cumulative months called for a given following spell.

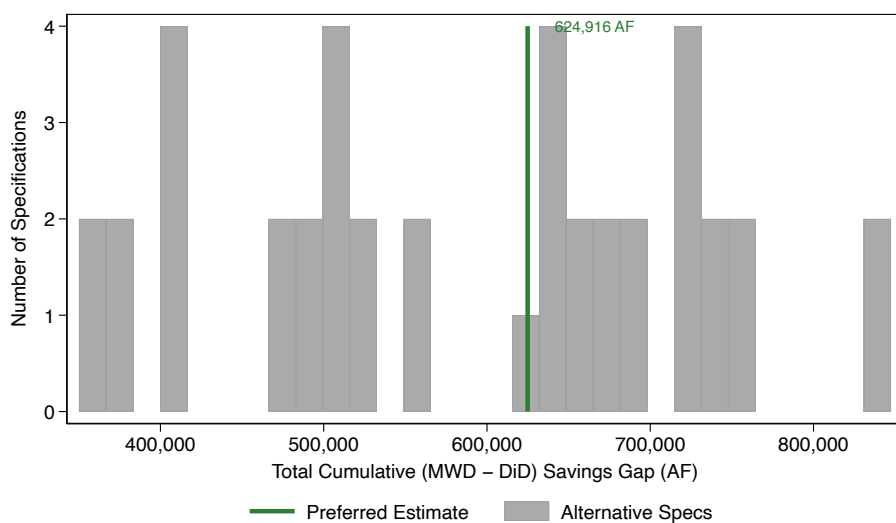
Figure S11: Prevailing Wind Direction at PVID



The NW and SE winds partially cancel, so the resultant is shorter than the scalar average.

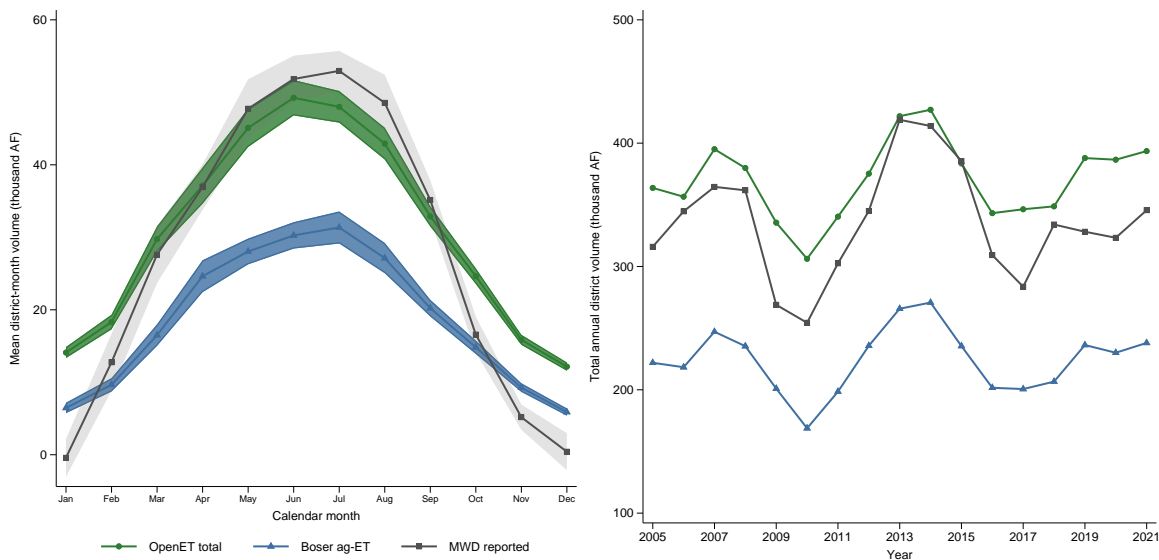
Notes: Wind rose illustrating how we aggregate hourly wind data from Iowa Environmental Mesonet into monthly prevailing wind directions using vector algebra, accounting for both wind direction and wind speed for each unique month-of-sample.

Figure S12: Cumulative Water-Savings Gap



Notes: This figure integrates annual implied differences over the program window: each gray bar is one of the 42 robustness specifications' total cumulative (MWD – DiD) savings gap summed over program years 2005–2021, and the vertical green spike marks the preferred specification's value of 624,916 AF, or 43.5% of MWD's reported total.

Figure S13: Comparison to Boser et al. (2024) Agricultural ET Adjustment



Notes: Left panel: mean district-wide ET (thousand AF) by calendar month with 95% CI ribbons; right panel: total annual district ET by calendar year, 2005–2021. Boser et al. (2024) (blue), OpenET Ensemble (green), MWD reported consumptive use (grey). ET volumes are computed as AFA per field × field acres summed over all PVID fields. MWD reported is from the annual USBR Decree Accounting Reports and MWD Verification Reports. Companion regression slopes in SI Table S10.

Tables

Table S1: Field-Level Balance and Summary Statistics

	Never-Called		Ever-Called		Difference	
	Mean	SD	Mean	SD	Estimate	(SE)
<i>Panel A: Field characteristics (time-invariant)</i>						
Acres	25.8	22.1	35.4	23.3	9.6***	(1.4)
Distance to Colorado River (m)	6462	5217	6018	4015	-444	(613)
Drainage index	22.269	8.628	23.048	5.701	0.779	(0.660)
<i>Panel B: Baseline ET, uncalled field-prog-years (AFA per year)</i>						
Own field ET	4.36	1.37	4.76	1.05	0.40***	(0.09)
Mean neighbor ET	4.24	1.12	4.43	0.86	0.19**	(0.08)
<i>Panel C: Crop shares, 2007–2010 (CDL)</i>						
Alfalfa	0.664	0.339	0.500	0.348	-0.165***	(0.023)
Cotton	0.116	0.236	0.059	0.167	-0.057***	(0.015)
Hay	0.031	0.113	0.021	0.100	-0.010	(0.007)
Fallow	0.152	0.249	0.412	0.345	0.260***	(0.020)
Other	0.036	0.123	0.009	0.054	-0.028***	(0.007)
<i>Fields</i>	726		1,985		2,711	
<i>Panel D: Treatment intensity (ever-called fields, N = 1,985)</i>						
Cumulative fallowing months			mean 42.8, median 36			
Longest spell (program-years)			mean 2.5, median 2			
<i>Panel E: Panel structure (prog-year analysis sample)</i>						
Field × prog-year observations			48,798			
Unique fields			2,711			
Program years			18			
Unique owners			302			

Notes: Comparison of treated (ever-called) and control (never-called) fields across field characteristics, pre-program ET outcomes, and early-sample crop shares. Panel A reports time-invariant covariates (acreage; Euclidean distance to the Colorado River; drainage index). Panel B reports baseline ET by ever-enrolled status. For ever-enrolled fields this equals the paper’s “Potential” reference of 4.77 AFA. Standard errors are clustered at the PLSS section level. Panel C reports field-level shares from the USDA Cropland Data Layer over the earliest available CDL years (2007–2010). Difference columns report the regression coefficient from $y_i = \alpha + \beta \text{Ever-Called}_i + \varepsilon_i$ with cluster-robust standard errors at the PLSS section level; significance: * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$. Panel D summarizes treatment intensity among the ever-called fields (cumulative months in any call; longest contiguous call spell measured in program-years).

Table S2: Regressions of MWD-Reported Consumptive Use on Each OpenET Model

	ENSEMBLE	PTJPL	SIMS	SSEBOP	GEEBPM	EEMETRIC	DISALEXI
Slope	0.992*** (0.0153)	1.021*** (0.0213)	1.001*** (0.0163)	1.007*** (0.0143)	0.998*** (0.0187)	0.865*** (0.0130)	1.037*** (0.0149)
adj. R^2	0.930	0.895	0.924	0.937	0.916	0.932	0.944
RMSE	9161.7	11265.6	9571.0	8712.5	10069.3	9037.4	8190.0
N	204	204	204	204	204	204	204

Notes: Slope is the elasticity of MWD AF on OpenET AFA (no intercept); adjusted R^2 and RMSE are of the same regression. All seven models are close to a 1:1 correspondence with MWD's reported totals at the district-monthly aggregation, with slope estimates ranging from 0.865 (EEMETRIC) to 1.037 (DISALEXI). Companion aggregate time-series in SI Figure S1.

Table S3: Full Regression Results (Preferred Specification)

	ENSEMBLE	PTJPL	SIMS	SSEBOP	GEEBPM	EEMETRIC	DISALEXI
Treatment	-2.530*** (0.440)	-1.837*** (0.462)	-2.498*** (0.509)	-2.918*** (0.481)	-2.135*** (0.348)	-3.005*** (0.497)	-2.150*** (0.350)
Placebo 1 (t-1)	0.149 (0.155)	0.128 (0.166)	0.143 (0.182)	0.146 (0.165)	0.134 (0.127)	0.150 (0.180)	0.178 (0.131)
Placebo 2 (t-2)	0.340* (0.181)	0.311 (0.193)	0.442* (0.227)	0.356* (0.194)	0.304** (0.143)	0.333 (0.204)	0.302** (0.140)
Placebo 3 (t-3)	0.518*** (0.187)	0.442** (0.211)	0.573** (0.229)	0.654*** (0.198)	0.455*** (0.147)	0.356 (0.257)	0.455*** (0.151)
Placebo 4 (t-4)	0.163 (0.202)	0.231 (0.246)	0.088 (0.229)	0.158 (0.222)	0.170 (0.154)	0.327 (0.210)	0.264 (0.162)
Placebo 5 (t-5)	0.739* (0.379)	0.862** (0.436)	0.693* (0.393)	0.820* (0.428)	0.608** (0.295)	0.881** (0.442)	0.663** (0.303)
Placebo 6 (t-6)	0.191 (0.177)	0.235 (0.290)	0.242 (0.248)	0.225 (0.183)	0.135 (0.117)	0.258 (0.209)	0.059 (0.110)
p -value (joint nullity of placebos)	0.132	0.335	0.163	0.055	0.057	0.363	0.070
DiD / MWD reported	54.1%	39.3%	53.5%	62.4%	45.7%	64.3%	46.0%
MWD reported baseline (AFA)				4.67			
Field \times prog-year observations				12,018			
Switcher fields				106			
Clusters (PLSS sections)				195			

Notes: All columns use the same preferred spec (call=2, flag=1, all-sample, no controls, did_multiplegt_dyn with effects(6) placebo(6)), varying only the ET outcome. Standard errors clustered at the PLSS section are in parentheses below each coefficient. Significance: * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$. The row " p -value (joint nullity of placebos)" reports the Wald test of the six pre-period placebos jointly equal to zero; all models reject only weakly or fail to reject, consistent with no pre-trend violation. The preferred ENSEMBLE column yields -2.53 AFA (54.1% of MWD's reported baseline of 4.67 AFA).

Table S4: SUTVA Sensitivity: Preferred Specification vs. Drops Pure Controls

	(1) Preferred (samp = All)	(2) Drop pure control	(3) Purified ctrl, no FE	(4) No owner trend
Treatment (Av_tot_eff)	-2.530*** (0.440)	-3.086*** (0.464)	-3.046*** (0.418)	-2.905*** (0.398)
Placebo 1 (t-1)	0.149 (0.155)	0.050 (0.159)	0.115 (0.131)	0.097 (0.129)
Placebo 2 (t-2)	0.340* (0.181)	0.194 (0.209)	0.249 (0.187)	0.235 (0.178)
Placebo 3 (t-3)	0.518*** (0.187)	0.292 (0.206)	0.476*** (0.175)	0.467*** (0.160)
Placebo 4 (t-4)	0.163 (0.202)	-0.416** (0.206)	0.374** (0.176)	0.320* (0.169)
Placebo 5 (t-5)	0.739* (0.379)	0.023 (0.339)	1.076*** (0.263)	0.887*** (0.261)
Placebo 6 (t-6)	0.191 (0.177)	-0.363 (0.239)	0.707** (0.313)	0.558* (0.307)
<i>p</i> -value (joint nullity of placebos)	0.132	0.048	<0.001	0.006
Field × prog-year obs	12,018	5,447	9,569	12,018
Unique fields	1,170	803	1,033	1,170
Clusters (PLSS sections)	195	177	190	195

Notes: Column (1) reports the preferred specification, which uses the full sample and absorbs non-parametric trends in landowner identity. Column (2) drops never-called fields, so identification comes only from variation in call timing across fields that are ever called. Column (3) implements a retains ever-called fields together with fields whose owner never has any field called (i.e., it drops uncalled fields of called owners); it also omits the owner-trend absorption, which would otherwise mechanically collapse this sample to column (2) by absorbing never-called owners entirely. Column (4) retains the full sample but drops the owner-trend absorption. Six placebo coefficients reported below the treatment estimate; joint-nullity *p*-value reported as “*p*-value (joint nullity of placebos).” Standard errors clustered at PLSS section in parentheses. Significance: **p* < 0.10, ***p* < 0.05, ****p* < 0.01. Companion monthly visualization in SI Figure S8.

Table S5: Predictors of Residual ET on Called Fields

	Cross-sectional gradient		Within-field over time	
Meters to Colorado River	-0.000 (0.000)			
Drainage Index	0.000 (0.000)			
Cumulative Fallowing (Months)			-0.000 (0.000)	
Average Neighbor ET (AFA)			0.344*** (0.033)	
Constant	0.111*** (0.005)	0.099*** (0.009)	0.108*** (0.003)	0.026*** (0.008)
Field FE	No	No	Yes	Yes
Month-of-sample FE	Yes	Yes	Yes	Yes
Observations	17,299	17,299	17,299	17,299
R^2	0.217	0.216	0.491	0.557

Notes: OLS regressions of monthly per-acre ET (AFA, ENSEMBLE) on parcel characteristics hypothesized to explain positive water use on called/fallow fields. The sample is restricted to field-months in which the field is called and not flagged or partially fallowed. Columns (1) and (2) examine cross-sectional gradients in time-invariant covariates and absorb month-of-sample fixed effects only; columns (3) and (4) identify within-field variation over time and absorb both field and month-of-sample fixed effects. *Note:* standard errors clustered at the PLSS section level in parentheses; significance: * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.

Table S6: Regressions of Own ET on Neighbor ET by Wind Orientation (Called Field-Months)

	(1)	(2)	(3)	(4)
All Neighbors' ET Mean		0.404*** (0.0107)		0.283*** (0.0217)
Upwind Neighbors' ET Mean			0.187*** (0.00648)	0.0641*** (0.0100)
Downwind Neighbors' ET Mean			0.183*** (0.00627)	0.0624*** (0.0109)
N	95866	95866	95866	95866
adj. R^2	0.386	0.477	0.471	0.479
p-val (Upwind = Downwind)			0.643	0.845

Notes: All four columns condition on field \times month-of-sample fixed effects; standard errors clustered at the PLSS section. Column (1) includes the field's mean neighbor ET (all neighbors equally weighted); column (2) splits the neighbor mean into upwind and downwind components using each month's prevailing wind; column (3) includes both the all-neighbor mean and the upwind/downwind split as a joint test. Upwind and downwind slopes are statistically indistinguishable ($p=0.64$ in column 2, 0.85 in column 3), consistent with the residual-ET spillover being mediated by sub-surface hydrology rather than atmospheric vapor transport.

Table S7: Minimum Detectable Effects (MDE) for the Upwind-vs.-Downwind Difference

Quantity	Value
Total neighbor effect	0.404 (0.011)
Upwind neighbor effect	0.187 (0.006)
Downwind neighbor effect	0.183 (0.006)
Upwind – Downwind	0.004 (0.009)
p-value (H_0 : equal)	0.64
MDE (power=0.80, α =0.05)	0.024 (5.9% of total)
MDE (power=0.90, α =0.05)	0.028 (6.8% of total)
MDE (power=0.80, α =0.10)	0.021 (5.2% of total)

Notes: Computed from the observed SE on the upwind–downwind contrast in column (3) of SI Table S6. At conventional thresholds (80% power, α =0.05 two-sided), the smallest atmospheric contribution to the neighbor effect the design could have detected is 0.024 AFA, or 5.9% of the total neighbor-ET coefficient. Given the observed upwind–downwind difference of 0.004 (p =0.64), the atmospheric channel is bounded well below the full neighbor spillover.

Table S8: Within-Owner Selection Regressions

	AR(1)	AR(1)+Fld FE	Lagged	Both	ML (LASSO)
Predicted CF (AR1)	-0.054*** (0.012)	-0.247*** (0.020)		0.014* (0.008)	
Lagged ET			-0.067*** (0.006)	-0.071*** (0.005)	
Predicted CF (ML)					-0.114*** (0.013)
._cons	0.314*** (0.058)	1.231*** (0.096)	0.369*** (0.026)	0.320*** (0.047)	0.595*** (0.063)
N	18244	18235	14401	14401	18244
adj. R^2	0.143	0.258	0.249	0.250	0.186

Notes: Dependent variable: field-prog-year call indicator (`called_2`). All regressions condition on owner \times year fixed effects so that coefficients reflect *within-owner* selection among simultaneously eligible fields. Columns vary the counterfactual-ET prediction strategy: (1) AR(1); (2) AR(1) with field fixed effects; (3) simple lagged ET; (4) AR(1) + lagged; (5) cross-validated LASSO (`ml_cf_pred`) including 12- and 24-month lags, neighbor behavior, and month \times year fixed effects. Standard errors clustered at the owner level; significance * p < 0.10, ** p < 0.05, *** p < 0.01. The LASSO coefficient of -0.114 (column 5) is the largest in magnitude and absorbs the selection signal left in the simpler specifications.

Table S9: Results Using Boser et al. (2024) Adjusted ET

	(1) OpenET total (Boser subsample)	(2) Boser ag-only (Boser subsample)	(3) Full-sample OpenET headline	(4) Boser-corrected (headline \times ag share)
DiD coef. (AFA)	-3.288*** (0.159)	-2.716*** (0.141)	-2.530*** (0.440)	-2.090*** (0.363)
% of OpenET total	100%	82.6%	100%	82.6%
Unique fields	1,148	1,148	1,170	1,170
Field \times prog-year obs	11,623	11,623	12,018	12,018

Notes: Columns (1) and (2) re-estimate the preferred dCdH specification on the 1,148-field Boser-covered subsample using OpenET total ET (col. 1) and Boser ag-only ET (col. 2). The ratio of column 2 to column 1 yields an implied Boser ag share of 82.6%. Columns (3) and (4) apply that ag share to the full-sample OpenET headline of -2.530 AFA, yielding a Boser-corrected headline of -2.090 AFA (ag-only causal effect). See SI Figure S13 for the levels comparison and SI Table S10 for the scaling regression. Significance: * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.

Table S10: Comparison of Boser et al. (2024) to MWD-Reported Benchmark

	(1) Boser ag-ET	(2) OpenET total ET
<i>Panel A: MWD reported = $\beta \times$ ET aggregate (no constant, robust SE)</i>		
β	1.603*** (0.025)	0.993*** (0.015)
p -value, $H_0: \beta = 1$	<0.001	0.639
R^2	0.942	0.931
N (year-months)	204	204
<i>Panel B: ET aggregate as % of MWD reported (district-month, 2005–2021)</i>		
Mean (% of MWD)	66.9%	110.4%

Notes: Panel A regresses MWD reported on each ET aggregate without a constant. Column (1): Boser ag-only. Column (2): OpenET Ensemble. Panel B reports mean district-monthly volumes (thousand AF) for the LHS benchmark and each RHS ET aggregate, plus the implied share (RHS / LHS). Significance: * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.

Table S11: Pre-Trend p -Value by Spatial-Flag Threshold and ET Model

Flag threshold	ENSEMBLE	PTJPL	SIMS	SSEBOP	GEEBPM	EEMETRIC	DISALEXI
≥ 1 RA (any; PREFERRED)	0.132	0.335	0.163	0.055	0.057	0.363	0.070
≥ 2 RAs (majority)	0.000	0.000	0.000	0.000	0.000	0.000	0.000
≥ 3 RAs (unanimous)	0.000	0.000	0.000	0.000	0.000	0.000	0.000

Notes: Joint-nullity p -value of the six pre-period placebo coefficients, varying the RA data-flag threshold (rows) and ET model (columns). All other dimensions held fixed at preferred values (call = 2, sample = All, no controls). Bold cells pass the parallel-trends test at $p > 0.10$.

Appendix specification tables: variation across ET model and call coding

Each of the three tables below covers one RA call-coding threshold (the digitization aggregation rule), holding sample, spatial-flag threshold, and estimator fixed at their preferred values. Within each table, columns report the seven OpenET models and panels vary the two control specifications (no controls vs. + precipitation), so each table summarizes fourteen specifications and the bundle as a whole summarizes forty-two.

Table S12: Specification Robustness Across ET Models — Call Coding ≥ 1 RA (Any)

	ENSEMBLE	PTJPL	SIMS	SSEBOP	GEEBPM	EEMETRIC	DISALEXI
<i>Panel A: No controls</i>							
Treatment	-2.763*** (0.429)	-2.204*** (0.472)	-2.859*** (0.470)	-3.194*** (0.460)	-2.227*** (0.344)	-3.215*** (0.485)	-2.339*** (0.350)
<i>p-value (joint nullity of placebos)</i>	0.064	0.187	0.105	0.038	0.019	0.197	0.012
<i>Panel B: + precipitation</i>							
Treatment	-2.763*** (0.426)	-2.204*** (0.465)	-2.859*** (0.465)	-3.194*** (0.458)	-2.227*** (0.343)	-3.216*** (0.483)	-2.339*** (0.347)
<i>p-value (joint nullity of placebos)</i>	0.060	0.184	0.101	0.036	0.018	0.194	0.011
N (switcher-year cells, dCdH)	2,925						

Notes: dCdH average treatment-on-treated coefficients for the post-call ET effect under the ≥ 1 RA (any) call-coding rule — the most permissive coding; treats any RA-flagged month as a call. Sample includes all PVID fields (ever-called and never-called); spatial-flag threshold = ≥ 1 RA (drops any field flagged as a fuzzy image by any RA). Columns are OpenET ensemble and per-model outputs. Panel A omits time-varying controls; Panel B adds monthly field-level precipitation. The italic row in each panel reports the joint-nullity *p*-value of the six pre-period placebo coefficients. Standard errors clustered at the PLSS section in parentheses. Significance: **p*<0.10, ***p*<0.05, ****p*<0.01.

Table S13: Specification Robustness Across ET Models — Call Coding ≥ 2 RAs (Majority; Preferred)

	ENSEMBLE	PTJPL	SIMS	SSEBOP	GEEBPM	EEMETRIC	DISALEXI
<i>Panel A: No controls</i>							
Treatment	-2.530*** (0.440)	-1.837*** (0.462)	-2.498*** (0.509)	-2.918*** (0.481)	-2.135*** (0.348)	-3.005*** (0.497)	-2.150*** (0.350)
<i>p-value (joint nullity of placebos)</i>	0.132	0.335	0.163	0.055	0.057	0.363	0.070
<i>Panel B: + precipitation</i>							
Treatment	-2.533*** (0.425)	-1.840*** (0.448)	-2.501*** (0.484)	-2.921*** (0.465)	-2.137*** (0.336)	-3.008*** (0.482)	-2.152*** (0.339)
<i>p-value (joint nullity of placebos)</i>	0.089	0.213	0.103	0.038	0.041	0.329	0.059
N (switcher-year cells, dCdH)	1,929						

Notes: dCdH average treatment-on-treated coefficients for the post-call ET effect under the ≥ 2 RAs (majority; the preferred coding) — requires at least two of the three RAs to agree. Sample includes all PVID fields (ever-called and never-called); spatial-flag threshold = ≥ 1 RA (drops any field flagged as a fuzzy image by any RA). Columns are OpenET ensemble and per-model outputs. Panel A omits time-varying controls; Panel B adds monthly field-level precipitation. The italic row in each panel reports the joint-nullity *p*-value of the six pre-period placebo coefficients. Standard errors clustered at the PLSS section in parentheses. Significance: * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.

Table S14: Specification Robustness Across ET Models — Call Coding ≥ 3 RAs (Unanimous)

	ENSEMBLE	PTJPL	SIMS	SSEBOP	GEEBPM	EEMETRIC	DISALEXI
<i>Panel A: No controls</i>							
Treatment	-2.908*** (0.431)	-2.362*** (0.467)	-2.927*** (0.494)	-3.330*** (0.469)	-2.437*** (0.353)	-3.386*** (0.487)	-2.474*** (0.352)
<i>p-value (joint nullity of placebos)</i>	0.142	0.301	0.140	0.045	0.060	0.125	0.065
<i>Panel B: + precipitation</i>							
Treatment	-2.902*** (0.421)	-2.357*** (0.457)	-2.921*** (0.481)	-3.323*** (0.457)	-2.432*** (0.345)	-3.379*** (0.478)	-2.469*** (0.346)
<i>p-value (joint nullity of placebos)</i>	0.121	0.278	0.111	0.035	0.051	0.128	0.060
N (switcher-year cells, dCdH)	1,629						

Notes: dCdH average treatment-on-treated coefficients for the post-call ET effect under the ≥ 3 RAs (unanimous) — the most conservative coding; requires all three RAs to agree. Sample includes all PVID fields (ever-called and never-called); spatial-flag threshold = ≥ 1 RA (drops any field flagged as a fuzzy image by any RA). Columns are OpenET ensemble and per-model outputs. Panel A omits time-varying controls; Panel B adds monthly field-level precipitation. The italic row in each panel reports the joint-nullity *p*-value of the six pre-period placebo coefficients. Standard errors clustered at the PLSS section in parentheses. Significance: * $p < 0.10$, ** $p < 0.05$, *** $p < 0.01$.