

Hidden Environmental Footprints of AI Training: Evidence from Google’s Gemini 1.0

Junho Choi

Columbia University, PhD in Sustainable Development



Key Takeaway

Context: During Google's Gemini 1.0 Training (Council Bluffs, Iowa; 2023)

Method: Compare **electrically proximate vs. distant fossil-fuel plants** using a **difference-in-differences** design

Finding: **Natural gas** plants ~322 tons (0.2SD) in CO₂/day (Jan-Apr '23)
Coal plants ~4.81 kilotons (0.5SD) in CO₂/day (May-Sep '23) (*p*<0.05), and similar for **local emissions** (SO₂ and NO_x)

Scale: With SCC of \$51/CO₂, extra emissions are worth about **113% of Google’s \$35 million pledge for carbon removal credits**

Implication: Without aligning training with **renewable** availability, scaling AI models risks amplifying **local and global emissions burdens**

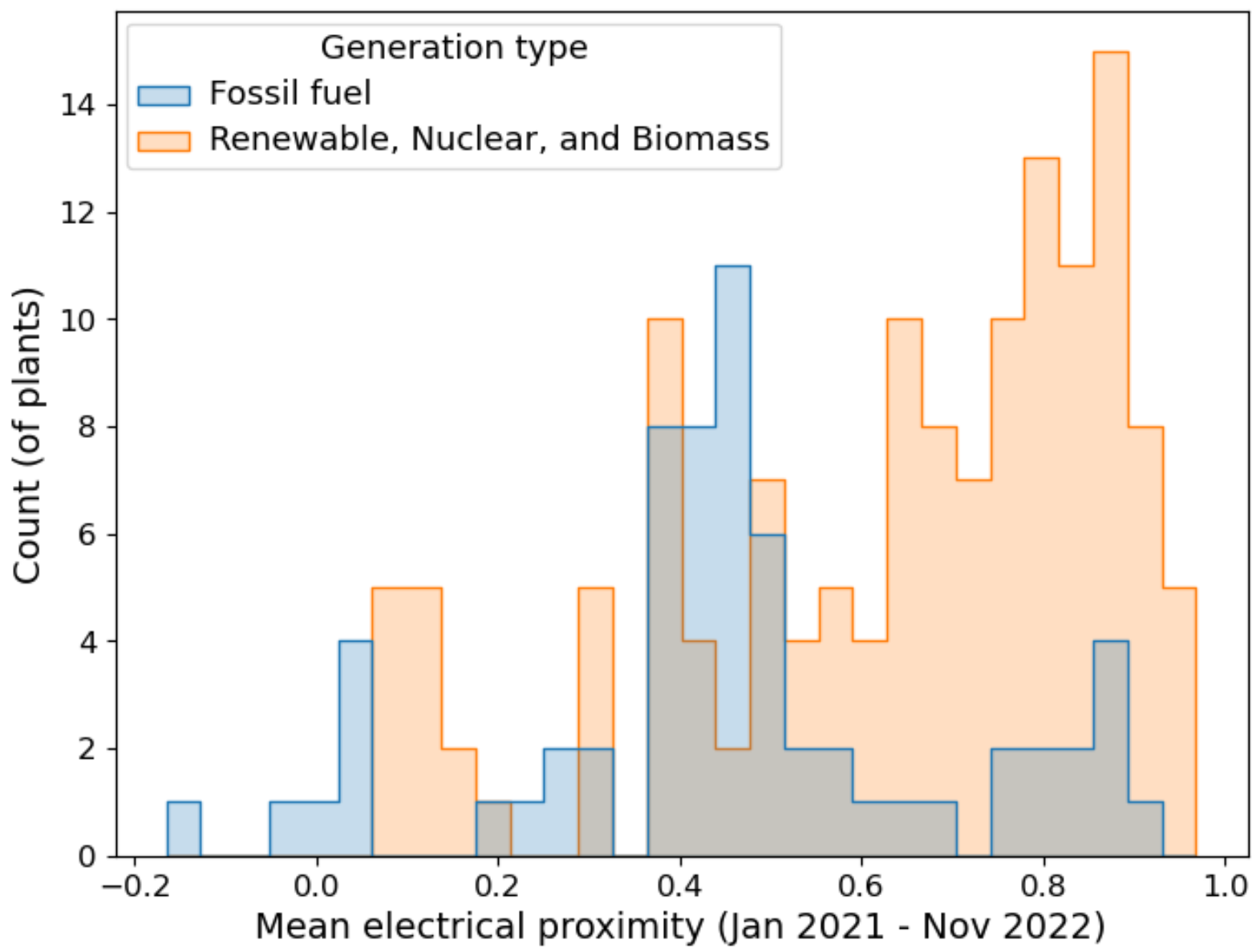
Data

CEMS	Hourly power plant emissions (CO ₂ , SO ₂ , NO _x)
MISO	Locational marginal prices and congestion component
ERA5 + ECWMF	Weather variables
IA DNR	Monthly operating reports of water supply
Homebase	Labor outcomes (to be added)
FERC-715	For resistance-based electrical proximity (under CEII)

Methodology

A. Calculate pre-training electrical proximity

- Identify loading zone (LZ) for Council Bluffs DC
- Identify power plants associated with gennodes
- Using **congestion component of pre-training period**, calculate **corr(LZ, gennode)** → “**electrical proximity**”



B. Define treated vs. control units (plants)

Treated

Fossil-fuel plants most likely to respond to DC demand
= electrically proximate
(top 10%)

Control

Comparable FF plants less likely to respond to DC
= electrically distant
(bottom 60%*)

*For sensitivity checks, other definitions (from bottom 10% to 60%) are also tested.

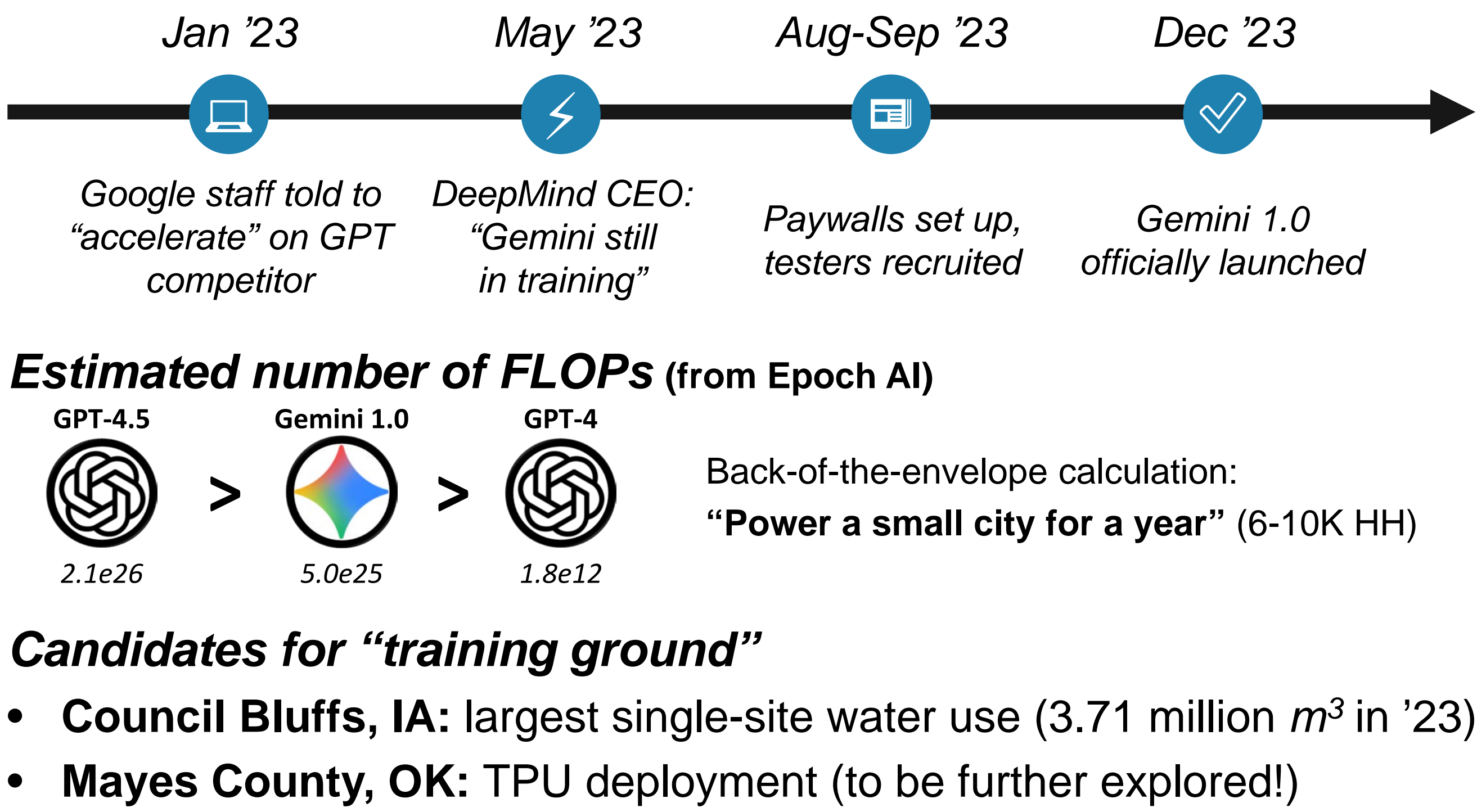
C. Difference-in-differences (DiD)

Key idea: comparing **treated vs. control plants** (difference) during and outside of **Gemini 1.0 training period** (in differences)

$$Y_{it} = \sum_{k=1}^5 \beta_k T_i \times C_{k,t} + \rho_i + \iota_{ym(t)} + X_{it}\gamma + \epsilon_{it}$$

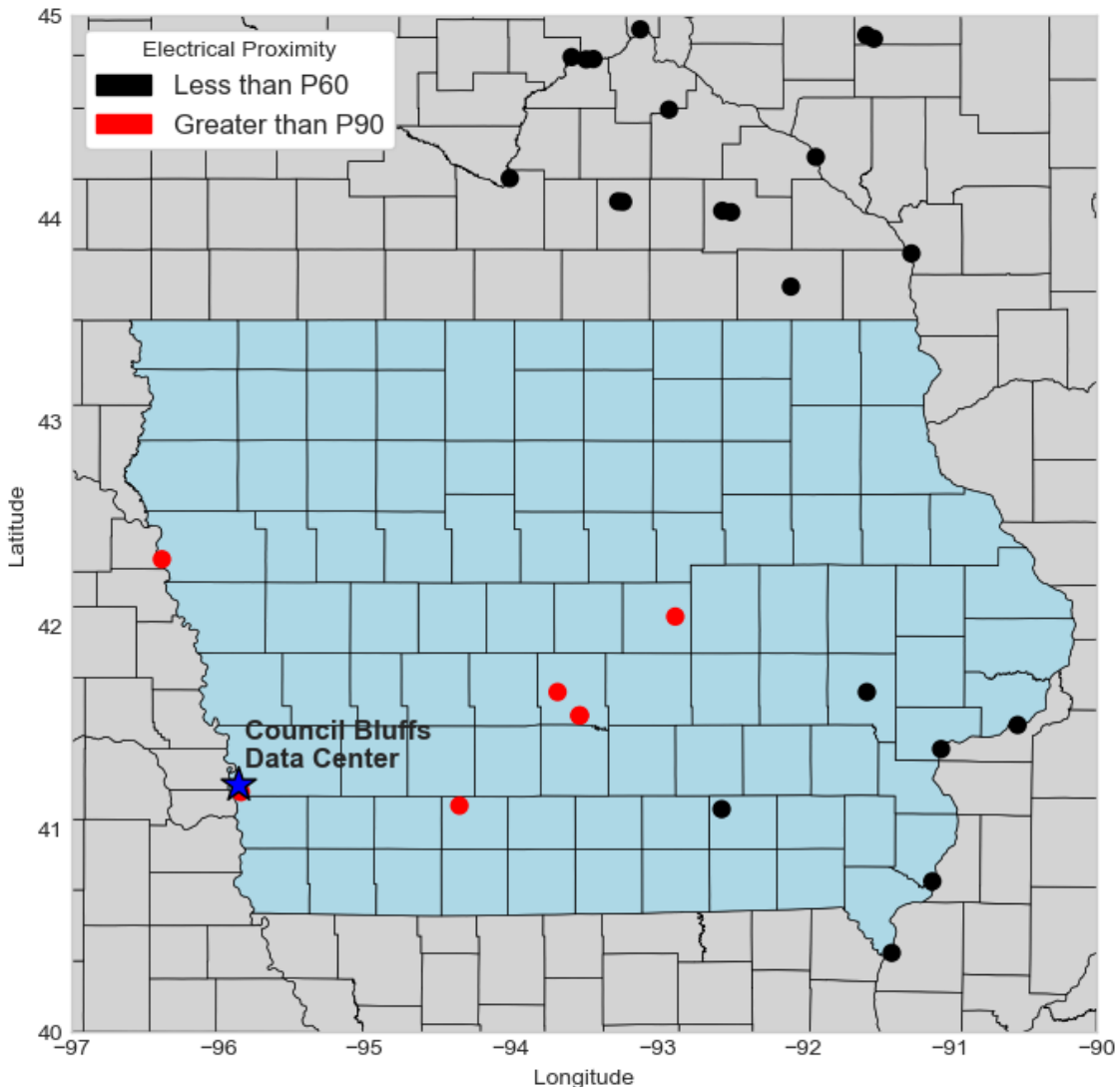
- T_i*: binary, =1 if treated (electrically proximate)
- C_{k,t}*: binary, =1 if during training
 - k* = 1: early, suspected training (Jan-Apr '23)
 - k* = 2: confirmed heavy training (May-Jul '23)
 - k* = 3: testing and fine-tuning (Aug-Sep '23)
 - k* = 4, 5: latter periods for Gemini 1.5 training (auxiliary tests)

Background

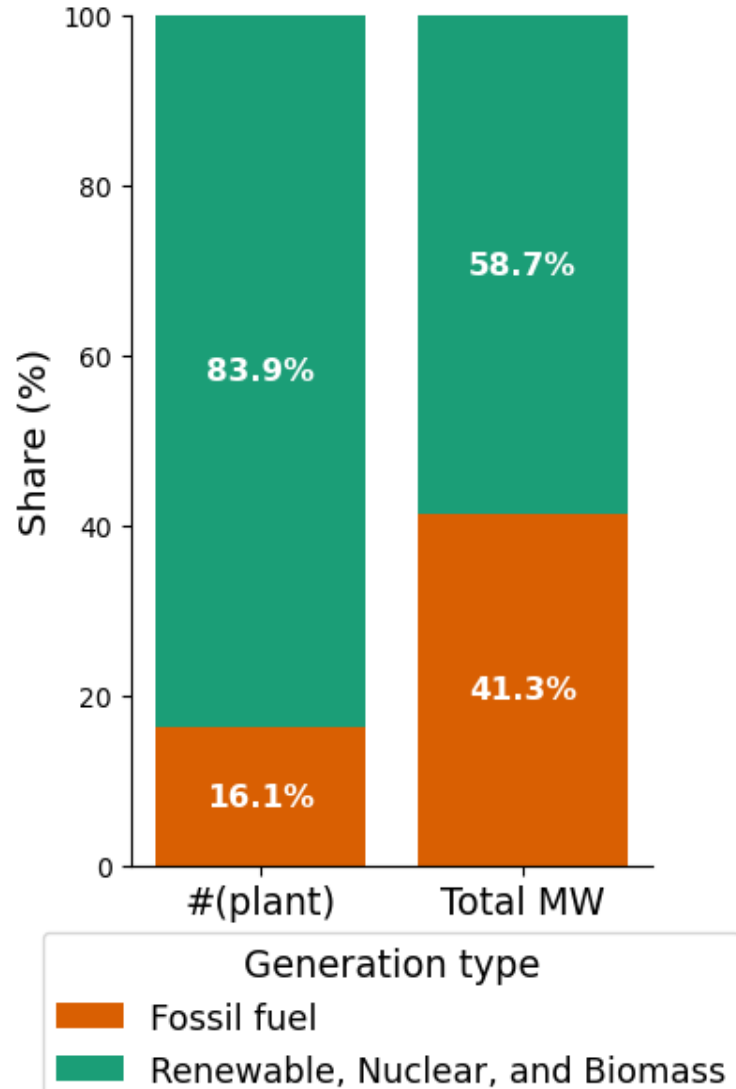


Council Bluffs, IA

Electrically proximate and far fossil fuel power plants



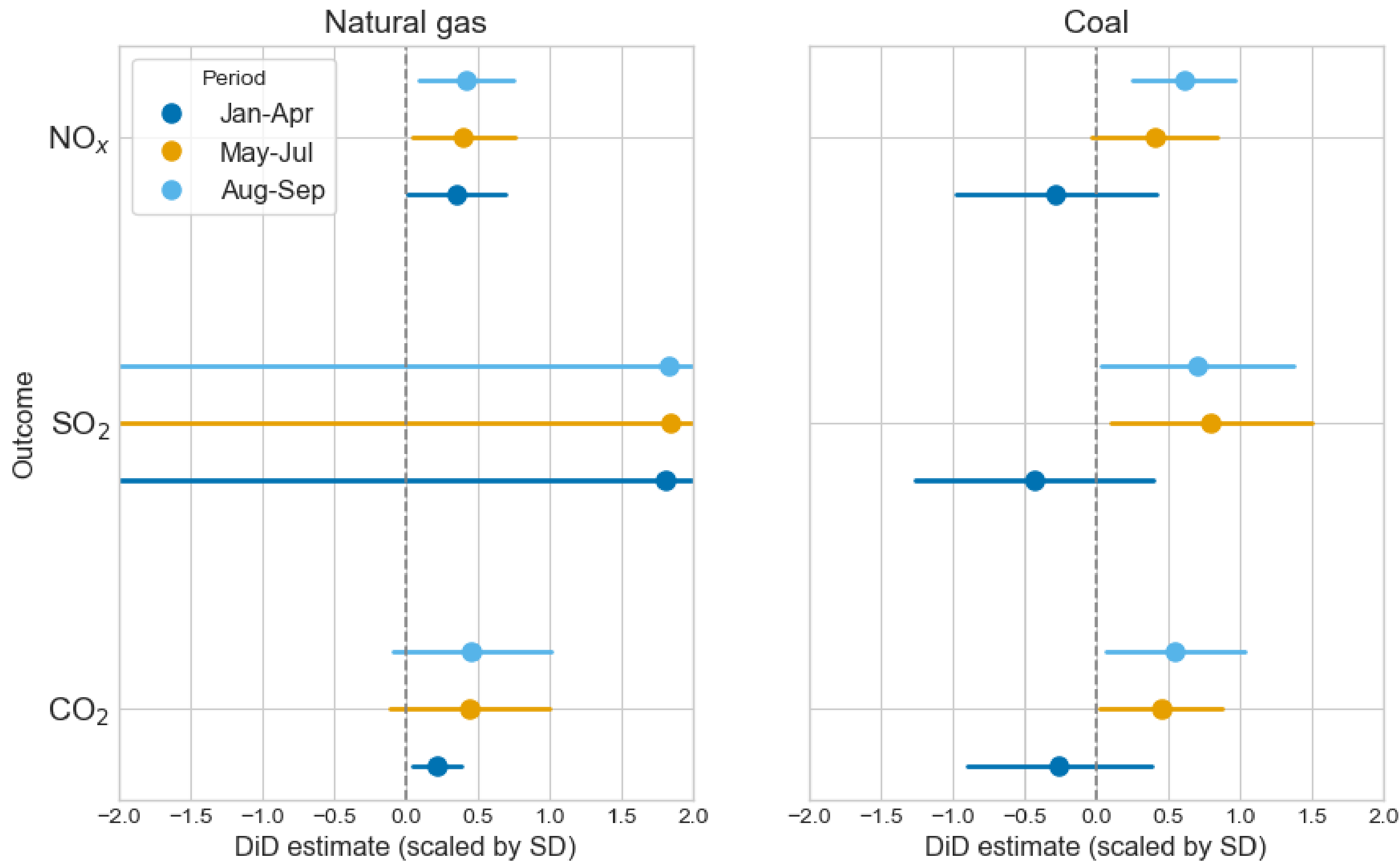
Presence of FF generation among proximate plants



Electrically close to renewable sources, but also to **fossil-fuel plants** with **substantially larger per-plant generating capacity**

Results

DiD effects (with 95% CI) by fuel, emission, and training period



Note: SE clustered two-way at the plant and year-by-month levels.

Emissions summary statistics for reference

	Natural gas			Coal		
	CO ₂ (kT)	SO ₂ (ton)	NO _x (ton)	CO ₂ (kT)	SO ₂ (ton)	NO _x (ton)
Mean	0.78	1.00	0.24	10.49	4.98	7.20
SD	1.49	0.16	0.60	9.62	7.15	7.81

Data center demand spikes trigger a systemic fallback to fossil fuels.

Discussion and Next Steps

- CO₂**: ~774 kilotons rise over the training period (SCC of \$51: USD 39 million; SCC of \$190: USD 147 million)
- Quantifying local effects (SO₂ and NO_x):** Homebase and AERMOD
- Preliminary evidence on water usage:** deriving compute intensity?
- Expanding beyond Iowa and Gemini 1.0:** Oklahoma and newer models