# **Extracting Structured Policy Information from Climate Action Plans**

#### Tom Corringham

Scripps Institution of Oceanography University of California San Diego La Jolla, CA, USA tomc@ucsd.edu

## Nupoor Gandhi

Language Technologies Institute Carnegie Mellon University Pittsburgh, PA, USA

#### **Bryan Flores**

Independent Researcher San Diego, CA, USA

#### Emma Strubell

Language Technologies Institute Carnegie Mellon University Pittsburgh, PA, USA

## Sireesh Gururaja

Language Technologies Institute Carnegie Mellon University Pittsburgh, PA, USA

## **Jacob Dunafon**

Independent Researcher San Diego, CA, USA

#### Tristan Romanov

Independent Researcher San Francisco, CA, USA

## **Abstract**

Most of the world's climate action policies are planned and implemented at the local level, through city and regional climate action plans (CAPs). To assess global progress in climate mitigation and adaptation, as in forthcoming assessments such as the 2027 IPCC Special Report on Climate Change and Cities, we need systematic ways to track and analyze these plans. However, CAPs are dispersed across thousands of jurisdictions, vary widely in structure and format, and are often difficult to access. We propose a standard CAP ontology, and a retrieval- and extractionoriented pipeline that leverages recent advances in natural language processing (NLP) and information retrieval (IR) to transform CAPs into a structured, verifiable dataset of climate policies. As a case study, we focus on California, where more than 260 local governments have published one or more CAPs since 2006. We develop an annotated benchmark dataset of 17 San Diego County CAPs with over 1,800 extracted policies and associated attributes. Unlike prior efforts that rely on small annotated corpora or industry-specific disclosures, our system explicitly grounds every extracted element in its underlying PDF, ensuring transparency and reducing hallucination in the produced dataset. Addressing these challenges will enable large-scale comparative analyses of CAPs across jurisdictions worldwide, supporting policymakers, sustainability officers, and hazard managers, and accelerating climate adaptation and mitigation efforts.

# 1 Introduction

Climate action plans (CAPs) represent a growing class of gray policy literature: lengthy, heterogeneous PDFs produced by local governments that encode critical commitments for climate mitigation and adaptation [Gandhi et al., 2024]. Their diversity in format, terminology, and availability makes them difficult to retrieve, parse, and analyze at scale. Sample CAP pages illustrating this heterogeneity are shown in Figure A1, highlighting challenges for automated information extraction. Despite their

policy importance, CAPs remain largely inaccessible to large-scale comparative analysis, creating an opportunity for natural language processing (NLP) to have real-world impact.

We propose a pipeline that combines agentic web retrieval with grounded information extraction (IE). Using headless browsing, our system autonomously navigates local government websites to surface CAPs that are often unindexed by search engines and unavailable to conventional scraping. Once retrieved, CAPs are transformed into structured datasets through IE methods designed for ecological validity: outputs are directly useful to policymakers and sustainability officers, while every extracted element is linked back to its source passage in the PDF to ensure verifiability and reduce hallucination.

From a research perspective, CAPs present unsolved challenges that make them an ideal testbed for advancing NLP and IR. Retrieval must operate over dispersed, poorly indexed sources where precision–recall trade-offs have tangible consequences. This extends work in the direction of automatic dataset collection, as in Ma et al. [2025]. Information extraction is necessarily multimodal, an area of active development, especially for vision-language models (VLMs, e.g. Wang et al. [2024], Liu et al. [2024]). Evaluation must account for cases where multiple annotations are equally valid, especially when policies are represented in tables or expressed in ambiguous spans. Addressing these challenges directs progress in web-based agentic IR and document-focused work in NLP towards gray policy literature domains and produces actionable tools for climate governance.

## **Novel contributions:**

- **CAP Ontology** We develop a nested structured ontology that allows for standardized representations of CAPs across jurisdictions (Figures A2, A3).
- Agentic retrieval of unindexed policy documents: Using headless browsing, our system automatically navigates county and municipal websites to surface CAPs that are generally not indexed by major search engines and inaccessible to basic scraping protocols.
- **Grounded extraction for policy trustworthiness:** Every extracted policy element is explicitly linked to its source passage in the PDF, reducing hallucination and enabling practitioners to verify outputs.
- Ecological validity and practitioner utility: Our dataset design emphasizes outputs that are directly useful to sustainability officers, hazard managers, and policymakers, rather than optimizing only for academic benchmarks.
- Calibrated evaluation thresholds: Conventional metrics (e.g., Cohen's kappa) assume rigid annotation consistency, but CAPs often admit multiple valid extractions. We propose evaluation methods that reflect practitioner-relevant agreement standards.

# 2 Proposed Approach

Our pipeline for extracting structured information from Climate Action Plans (CAPs) has three primary components: (1) web retrieval of dispersed PDF documents, (2) large language model (LLM)-based information extraction with grounding in source passages (Figure A4), and (3) evaluation through human annotation and agreement analysis (Figure A5).

**Web Scraper.** CAPs are often buried deep in municipal or county government websites and often not indexed by search engines. To systematically collect them, we employ Selenium-driven headless browsing to autonomously navigate and capture PDF documents. Metadata such as jurisdiction, publication year, and plan version are stored alongside each document, enabling transparent traceability and jurisdiction-level comparisons. As future work, we plan to integrate agent-based web navigation tools [Ma et al., 2023, Huang et al., 2024].

**Information Extraction Task.** The core of our approach is an LLM-based extraction pipeline that processes CAPs at the page level. Earlier versions of our system used a text-only unimodal approach, where we used the Adobe Acrobat API to convert PDFs into structured HTML and JSON, which were then processed with OpenAI models and other LLMs. While this is a common approach in document-based NLP, we found that the lack of multimodal processing rendered model results highly inconsistent. Our current implementation instead renders page images and passes them to OpenAI's multimodal models, prompting them to extract policy statements and associated attributes (Figure A3).

This variability in extraction accuracy reflects open questions in document-based NLP in how best to process multimodal content [Deng et al., 2024].

Policies are represented using an ontology of nested fields that we developed in consultation with climate policy experts (Figure A2). Each record consists of a root policy and associated attributes, including: policy description, sector (e.g., transportation, buildings, waste), target year, quantitative GHG reduction goal, cost allocations (residential, private-sector, municipal), and co-benefits (e.g., equity, resilience, public health). Extraction is explicitly grounded: every attribute is linked to the original text span, table cell, or figure icon in the PDF. This allows practitioners to verify system outputs and ensures transparency in downstream analyses.

A key technical challenge is the length of documents and the nesting structure of policies within CAPs, which often contain hierarchical "goal–approach–action" relationships spread across documents that often exceed 100 pages in length. To address this, we are experimenting with windowed page contexts (e.g., combining two consecutive pages) and hierarchical labeling schemes. We are also developing strategies for handling CAP-specific iconography (e.g., cost or benefit symbols, Figure A6) that define attributes remotely within the document.

**Evaluation.** To ensure both accuracy and ecological validity, we evaluate our system against a human-annotated benchmark of 17 San Diego CAPs with over 1,800 labeled policies. Labels were created using Label Studio, with multiple annotators tagging the same documents to assess agreement. In addition to standard metrics, we compute inter-annotator reliability (Cohen's  $\kappa$ , Krippendorff's  $\alpha$ ) to quantify the inherent ambiguity of policy texts. Evaluation of the LLM outputs is then performed relative to this annotated benchmark, with attention to both span-level overlap and attribute-level correctness. We also assess trade-offs between precision and recall to reflect the practical needs of policymakers: in some applications, high coverage of policies is preferred even at the cost of additional noise, whereas in others strict precision is critical.

# 3 Expected Outcomes and Impact

Our proposed pipeline will deliver outputs of both immediate practical utility and long-term research value. By making local CAPs and climate action policies accessible, structured, and verifiable, we will enable new forms of comparative analysis and decision support.

**Website.** We will build a public-facing website that allows users to interactively explore climate action policies. Users will be able to query policies by sector, timeframe, cost, or co-benefits, and immediately view the corresponding passages in the original CAPs. This grounding enhances transparency, reduces the risk of hallucination, and builds practitioner trust in AI-derived policy databases. A screenshot of the prototype interactive website is shown in Figure A7.

**Database.** The structured policy database will expand over time from our San Diego benchmark set to encompass all California CAPs and eventually U.S. and international plans. Each policy entry will include attributes such as sector, emissions reduction targets, timeframes, responsible entities, co-benefits, and costs (public and private). Linking structured attributes back to the original text ensures accountability and provides a foundation for reproducible research.

**Downstream Analyses.** The resulting dataset will allow comparative, quantitative analysis of local climate governance at unprecedented scale. For example, our pilot analyses reveal notable trends in California CAPs (Figure A8). For example, over time, CAPs have placed increasing relative emphasis on adaptation alongside mitigation, suggesting a growing recognition that projected warming thresholds require both strategies (Figure A9). Metrics of CAP quality, e.g., fraction of populated attributes per policy by CAP, are positively correlated with community wealth, suggesting disparities in local government capacity and resources for climate planning (Figure A10). These kinds of insights are directly relevant to forthcoming efforts such as California's Fifth Climate Change Assessment and the IPCC Special Report on Climate Change and Cities [IPCC, 2024], and illustrate the potential of NLP to illuminate equity and effectiveness in climate policy.

# 4 Conclusion

By coupling retrieval, grounded information extraction, and calibrated evaluation, our system will transform scattered, heterogeneous local climate action plans into structured, verifiable datasets. The

resulting tools will support policymakers, sustainability officers, hazard managers, and researchers in evaluating and improving local climate action, helping accelerate both mitigation and adaptation worldwide.

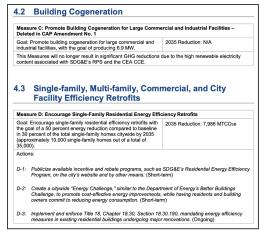
# **Acknowledgments and Disclosure of Funding**

This work was supported by the California Nevada Adaptation Program, a NOAA Climate Adaptation Partnership, and Climate Change AI (CCAI). The authors declare no competing interests.

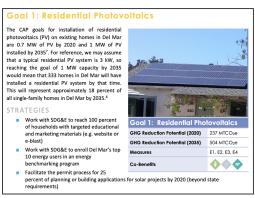
## References

- Naihao Deng, Zhenjie Sun, Ruiqi He, Aman Sikka, Yulong Chen, Lin Ma, Yue Zhang, and Rada Mihalcea. Tables as images? exploring the strengths and limitations of llms on multimodal representations of tabular data. *arXiv preprint arXiv:2402.12424*, 2024.
- Nupoor Gandhi, Tom Corringham, and Emma Strubell. Challenges in end-to-end policy extraction from climate action plans. In Dominik Stammbach, Jingwei Ni, Tobias Schimanski, Kalyan Dutia, Alok Singh, Julia Bingler, Christophe Christiaen, Neetu Kushwaha, Veruska Muccione, Saeid A. Vaghefi, and Markus Leippold, editors, *Proceedings of the 1st Workshop on Natural Language Processing Meets Climate Change (ClimateNLP 2024)*, pages 156–167, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.climatenlp-1.12. URL https://aclanthology.org/2024.climatenlp-1.12/.
- Wenhao Huang, Zhouhong Gu, Chenghao Peng, Jiaqing Liang, Zhixu Li, Yanghua Xiao, Liqian Wen, and Zulong Chen. Autoscraper: A progressive understanding web agent for web scraper generation. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 2371–2389, 2024.
- IPCC. Decision ipcc-lxi-5: Outline for the special report on climate change and cities. https://www.ipcc.ch/site/assets/uploads/2024/07/IPCC-LXI-2-Add.3.pdf, 2024. 61st Session decision document.
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llavanext: Improved reasoning, ocr, and world knowledge, January 2024. URL https://llava-vl.github.io/blog/2024-01-30-llava-next/.
- Kaixin Ma, Hongming Zhang, Hongwei Wang, Xiaoman Pan, Wenhao Yu, and Dong Yu. Laser: Llm agent with state-space exploration for web navigation. *arXiv preprint arXiv:2309.08172*, 2023.
- Tianyi Ma, Yiyue Qian, Zheyuan Zhang, Zehong Wang, Xiaoye Qian, Feifan Bai, Yifan Ding, Xuwei Luo, Shinan Zhang, Keerthiram Murugesan, Chuxu Zhang, and Yanfang Ye. Autodata: A multiagent system for open web data collection, 2025. URL https://arxiv.org/abs/2505.15859.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.

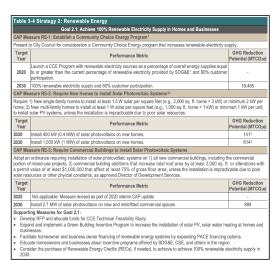
# **Appendix**



City of Carlsbad 2020, p.46



City of Del Mar 2016, p.35



City of Encinitas 2020, p.51



City of San Diego 2022, p.52

Figure A1: Example pages from four San Diego County CAPs (Carlsbad, Del Mar, Encinitas, and the City of San Diego). CAPs vary widely in format, with information often presented in nonstandard tables, using symbols or icons whose definitions appear on other pages or chapters, making automated information extraction especially challenging.

Tag	Sub-Tag	Definition	Example
Mitigation	Quantity	Quantity of GHG reduction	quantitative (60, 500), qualitative (high, medium, low), other
Mitigation	Unit	Unit of GHG reduction	mtCO2e, tons of carbon, other
Mitigation	Type	GHG type	co2, methane, n2o, hfcs, pfcs, sf6, nf3, other
Mitigation	Emitter	Source of GHG emmisions	concrete, industrial plant, tailpipe, other
Mitigation	Sector	Sector of GHG emissions reduction	vehicles, industry, waste, agriculture, land use, buildings, other
Mitigation	Context	Context of emissions	during peak hours, in cold weather
Adaptation	Hazard	Hazard addressed	heat, drought, flood, wildfire, sea-level rise, other
Adaptation	Quantity	Quantity associated with hazard	quantitative (60, 500), qualitative (high, medium, low), other
Adaptation	Metric	Metric associated with hazard	trees planted, acres treated, acre-feet of water saved, other
Policy	Time Frame	Target date	2030, 2040, 2050, already achieved, other
Policy	Implementing Agency	Responsible for implementation	state government, city department of public works, air quality board
Policy	Mechanism	Policy type	requirement, mandate, ordinance, plan, incentive, education program
Policy	Funding	Funding mechanism	tax, bond, grant, fee, general fund
Policy	Cost	Cost of policy	quantitative (\$1 million), qualitative (low, moderate, high), other
Policy	Benefit	Benefit of policy	quantitative (\$5 million), qualitative (low, moderate, high), other
Policy	Co-Benefits	Supplementary benefit	public health, cost reduction, green jobs, other
Policy	Reference	Reference document	Assembly bill 32, CARB scoping plan, IPCC AR6

Figure A2: A simplified representation of our ontology for structured representation of CAP policies. The ontology organizes policy information into top-level tags (e.g., Emissions Reduction, Adaptation, Time Frame, Benefit, Cost, etc.) with sub-tags, definitions, and representative examples. A hierarchical design enables consistent extraction and comparison of heterogeneous policy elements across jurisdictions. By explicitly defining attributes such as hazards, quantities, co-benefits, and emissions reduction metrics, the ontology provides both machine-actionable structure and practitioner-relevant interpretability, bridging the gap between natural language text and standardized policy datasets.

```
Prompt:
"You are an expert in climate policy.
For the following image, extract the nested tree structure of measures,
actions, goals, policies, strategies.

-- Include only spans that appear in the image, without modification.

-- Only Extract spans that appear in the text.
-- Do not extract any other spans that do not appear exactly in the text.
-- Do not modify or combine spans."
Ontology:
from pydantic import BaseModel
from typing import Optional
                                                  class Duration(BaseModel):
                                                      start_date: str
class SpanCluster(BaseModel):
                                                      end_date: str
    label: str
                                                      length_of_time: str
    label_description: str
                                                      elapsed: bool
    spans: list[str]
                                                      in_progress: bool
                                                      still_needs_to_be_initiated: bool
class Clusters(BaseModel):
    clusters: list[SpanCluster]
                                                  class ClimateHazard(BaseModel):
                                                      primary_impact: list[str]
class Mechanism(BaseModel):
                                                      secondary_impact: list[str]
    education outreach: bool
                                                      tertiary_human_impact: list[str]
    mandate: bool
                                                      exposure_subpopulation: list[str]
    incentive: bool
                                                      exposure_geographic_area: list[str]
    disincentive: bool
                                                      severity_area: list[str]
    voluntary_action: bool
                                                      duration: Duration
    tax: bool
    regulation: bool
                                                  class GHGEmission(BaseModel):
    fee: bool
                                                    emission_source: list[str]
                                                    emission_type: list[str]
    objective: list[str]
                                                    quantity: str
    compliance: list[str]
                                                    unit: str
    action_description: list[str]
                                                    duration: Duration
    authority: list[str]
                                                    emitter_stakeholder: list[str]
    target: list[str]
                                                    affected_by_emission_stakeholder: list[str]
```

Figure A3: Illustration of the extraction setup for policy text. The prompt instructs the model to identify nested structures of measures, actions, goals, policies, and strategies directly from the text, enforcing exact-span fidelity. To structure responses, we provide the model with a JSON schema derived from our ontology (excerpt shown). This schema specifies classes (e.g., Mechanism, Duration, ClimateHazard, GHGEmission) and their associated attributes, ensuring that extracted policies are represented consistently and with explicit links to mechanisms, timeframes, hazards, and emissions data. By constraining output to this ontology, the system produces machine-readable representations of policies while retaining their hierarchical and semantic relationships.

```
"doc_id": "cap_measure_e6_example",
  "page": 1,
  "policy": {
    "id": "E-6",
    "title": "Income-Qualified Solar PV Program",
    "description": "Facilitate installation of small-scale on-site solar PV
systems on income-qualified housing by promoting state programs and
collaborating with GRID Alternatives."
  "hierarchy": {
    "existing_efforts": [
      "SASH program: 24 new affordable housing units by People's Self Help
Housing."
      "Collaboration with GRID Alternatives on outreach and eligibility."
    "implementation_actions": [
      "E-6.1: Outreach via GRID Alternatives to developers and homeowners.",
      "E-6.2: Promote solar incentives from the California Solar Initiative
(SASH, MASH)."
  },
  "attributes": {
   "emissions_reduction": {
     "quantity": 87,
      "unit": "MT CO2e",
      "sector": ["Buildings", "Energy"],
      "emitter": "Residential electricity (avoided via rooftop PV)"
    "cost": {
     "public": "Very Low",
      "private": "None"
    "benefit": {
     "public": "None",
      "private" "Medium"
    "co_benefits": []
  },
  "provenance": [
    { "field": "emissions_reduction.quantity", "anchor": "GHG Reduction
Potential: 87 MT CO2e" },
    { "field": "cost.public", "anchor": "City Cost: Very Low" },
    { "field": "benefit.private", "anchor": "Private Savings: Medium" },
    { "field": "policy.title", "anchor": "Measure E-6: Income-Qualified Solar PV
Program" }
  ],
  "status": "model_pred",
  "model meta": {
   "source": "vision+layout pipeline",
    "created_at": "2025-08-20T00:00:00Z"
```

Figure A4: Mock Docling output. Example of structured JSON representation for a climate action plan policy (Measure E-6: Income-Qualified Solar PV Program), showing extracted attributes (emissions reduction, costs, benefits), hierarchical policy structure (existing efforts, implementation actions), and provenance links to the original text.

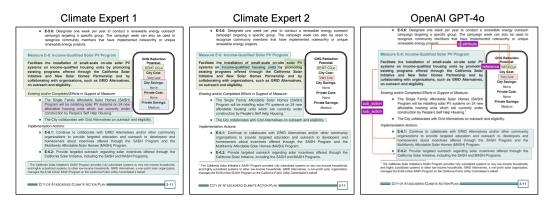


Figure A5: In addition to the 17 San Diego County CAPs, for which all policies and attributes were manually extracted, a subset of pages from other California CAPs were annotated in Label Studio to evaluate model performance. Shown here are three examples: two annotated by climate experts and one by OpenAI GPT-4o. The two human annotators produced notably different label sets, underscoring the need to calibrate performance thresholds rather than relying on fixed values from the literature, which may not be appropriate in this context.



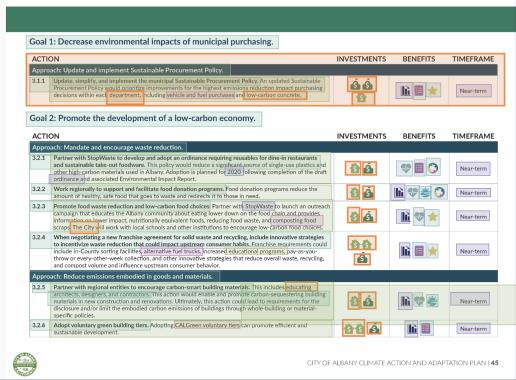


Figure A6: Example from the City of Albany Climate Action and Adaptation Plan illustrating CAP-specific iconography. Policy attributes such as costs, benefits, and co-benefits are encoded with icons (bottom, p. 45), whose meanings are defined remotely elsewhere in the document (top, p. 35). This highlights a key multimodal extraction challenge: text-only methods cannot recover these attributes, whereas multimodal approaches must link icons across pages to correctly capture policy information.

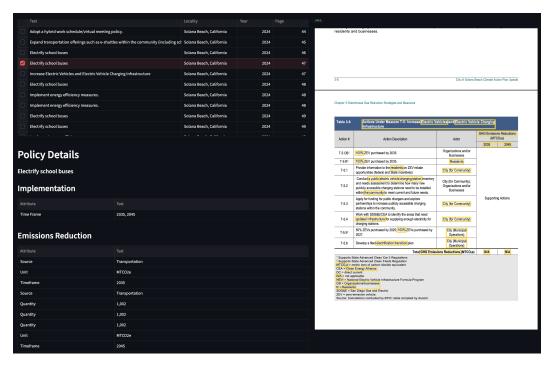


Figure A7: Streamlit prototype website for exploring structured CAPs. The application enables users to browse extracted policies (top left), filter by locality, year, and page, and view structured attributes. In this example, selecting "Electrify school buses" from the Solana Beach CAP reveals implementation timeframes and emissions reduction quantities (bottom left). The right-hand panel displays the source PDF with highlighted spans corresponding to extracted attributes, providing explicit grounding and allowing users to verify outputs directly against the original document.

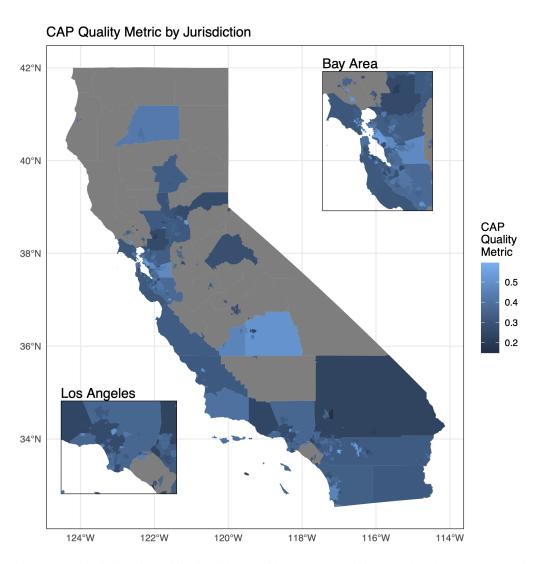


Figure A8: Jurisdictions in California with and without CAPs as of 2022, colored by a CAP "quality metric." Cities (census-designated places) and counties with CAPs are shaded according to the average fraction of policy attributes populated (e.g., sector, implementing agency, cost, emissions reduction, co-benefits), while jurisdictions without CAPs are shown in gray. Insets highlight the Bay Area and Los Angeles regions. The figure shows no clear geographic pattern in CAP quality, suggesting that differences in plan completeness vary jurisdiction by jurisdiction rather than clustering spatially.

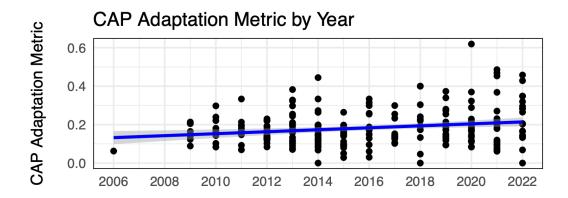


Figure A9: Temporal trend in the share of adaptation policies in California CAPs. Each point represents a CAP, with the adaptation metric defined as the fraction of its policies classified as adaptation rather than mitigation. Classification was performed by prompting GPT-40 to label policy descriptions as "adaptation" (e.g., wildfire defensible space), "mitigation" (e.g., home electric vehicle charging incentives), or "both" (e.g., tree planting for sequestration and heat reduction). The fitted regression line shows a modest but steady increase in adaptation focus over time, suggesting that local governments are devoting increasing attention to climate impacts alongside emissions reduction.

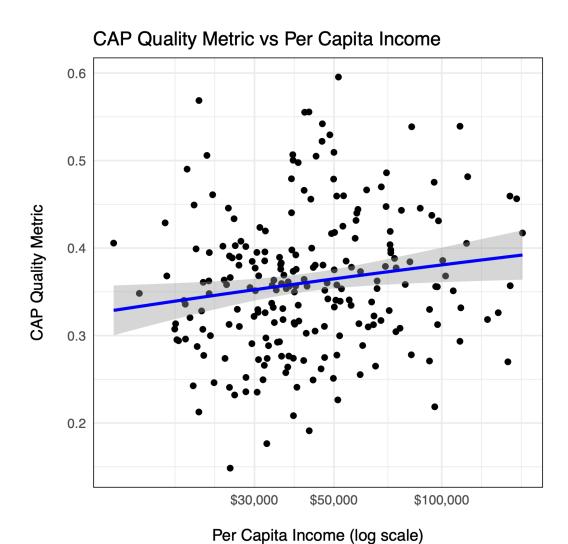


Figure A10: Relationship between CAP quality and community wealth. CAP quality is measured as the average fraction of populated attributes per policy (e.g., inclusion of sector, responsible entity, cost, emissions reduction, and co-benefits). Per capita income is shown on a log scale. The positive association indicates that wealthier jurisdictions tend to produce more complete and detailed climate action plans, suggesting disparities in local government capacity and resources for climate planning.