

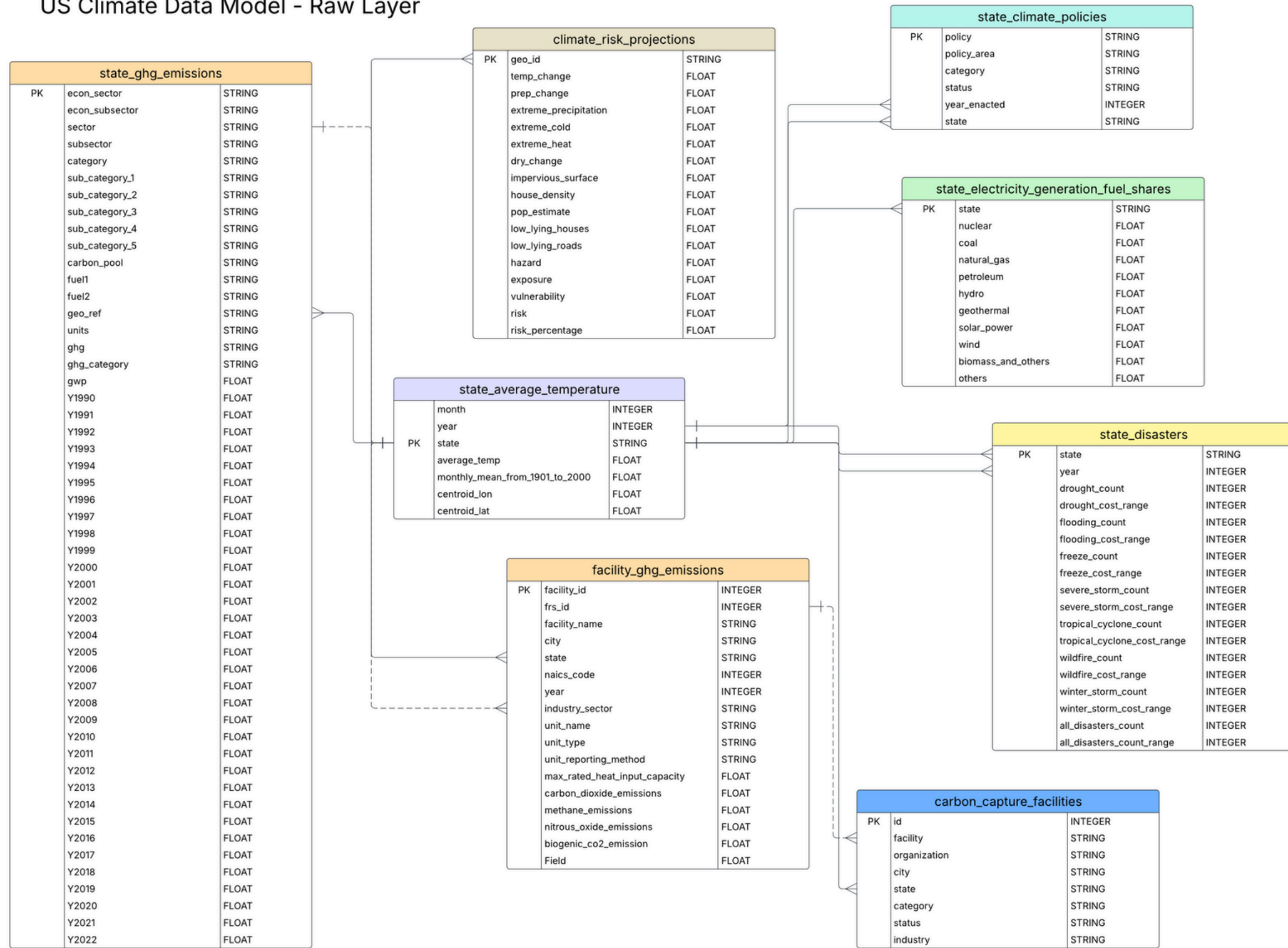


U.S. CLIMATE DATA WAREHOUSE

Erica Zhao, Kiara Foght

RAW DATA

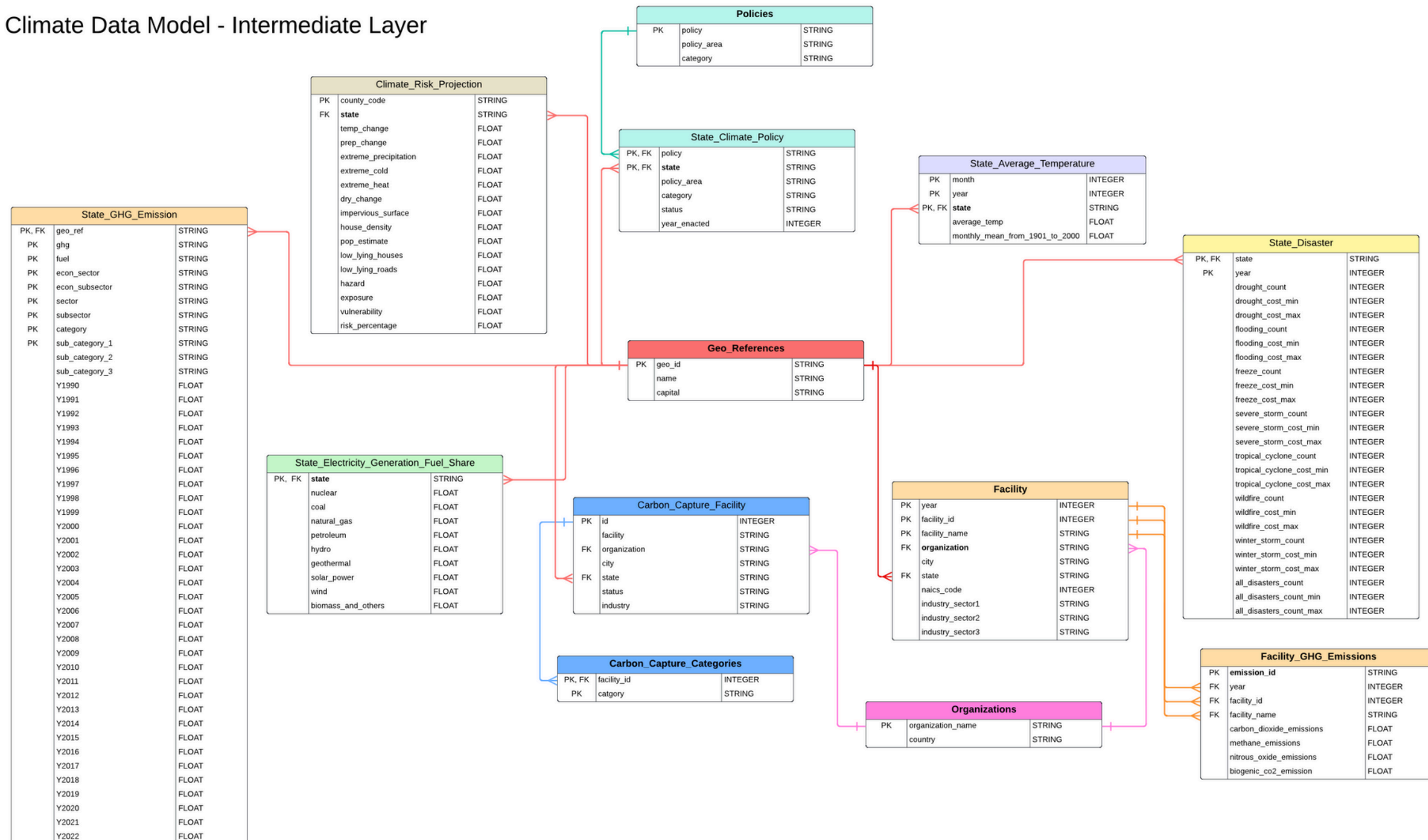
US Climate Data Model - Raw Layer



- EPA
- Kaggle
- Climate XChange
- NOAA
- NEI
- BBI International
- NASA

INT ERD

US Climate Data Model - Intermediate Layer



ORGANIZATION TABLE

Use LLM to Identify Organizations

We used an LLM to identify organizations by querying facilities without an assigned organization and extracting names based on facility details. CCF and GHG facilities were processed separately, and results were stored in a temporary table.

Normalize Organization Names

Next, we normalized organization names by extracting distinct names, using the LLM to standardize variations (e.g., "3M" vs. "3M Company"), and appending country information where possible. The standardized names were stored in a mapping table.

Create the Final Organization Table

We then created the final Organization Table by merging CCF and GHG organizations, removing duplicates, and assigning primary keys for data integrity.

Validate and Finalize the Organization Table

Finally, we validated and finalized the table by ensuring name uniqueness, creating the final BigQuery table, and cleaning up temporary tables to optimize storage.

ORGANIZATION TABLE

```
prompt_ghg = """Given a facility from the GHG emissions dataset:
facility_id, facility_name, city, state, naics_code, industry_sector1,2,3.

Identify the organization that owns or operates this facility, or return null if unknown.
Return EXACTLY one JSON line:
{
  "facility_id": <string>,
  "organization_name": <string or null>
}
No extra text or explanation.
"""
```

01

Prompt to Find
Organization

02

Create Check-
Point Table

```
sql_orgs = """
SELECT DISTINCT organization_name
FROM us_climate_int.tmp_ccf_facilities_llm_org_checkpoint
WHERE organization_name IS NOT NULL

UNION DISTINCT

SELECT DISTINCT organization
FROM us_climate_stg.carbon_capture_facilities
WHERE organization IS NOT NULL
"""
```

03

Prompt to
Normalize
Organization

```
batch_size = 50
min_batch_size = 5
sleep_time = 5
max_retries = 5

results_ghg = []
i = 0

while i < len(df_ghg):
    try:
        batch_df = df_ghg.iloc[i:i + batch_size]
        batch_results = [find_org_for_ghg(row) for _, row in batch_df.iterrows()]

        df_batch = pd.DataFrame(batch_results)
        pandas_gbq.to_gbq(
            df_batch,
            "us_climate_int.tmp_ghg_facilities_llm_org_checkpoint",
            project_id=project_id,
            if_exists="append"
        )

        results_ghg.extend(batch_results)
        i += batch_size
        print(f"Processed {i}/{len(df_ghg)} records.")

        time.sleep(sleep_time)

    except GoogleAPIError as e:
        print(f"Quota error encountered: {e}. Retrying with backoff...")

        for retry in range(1, max_retries + 1):
            time.sleep(sleep_time * retry)
            print(f"Retrying (attempt {retry}/{max_retries})...")
            try:
                batch_df = df_ghg.iloc[i:i + batch_size]
                batch_results = [find_org_for_ghg(row) for _, row in batch_df.iterrows()]

                df_batch = pd.DataFrame(batch_results)
                pandas_gbq.to_gbq(
                    df_batch,
                    "us_climate_int.tmp_ghg_facilities_llm_org_checkpoint",
                    project_id=project_id,
                    if_exists="append"
                )

                results_ghg.extend(batch_results)
                i += batch_size
                print(f"Processed {i}/{len(df_ghg)} records after retry.")

                break

            except GoogleAPIError:
                if retry == max_retries:
                    print("Max retries reached. Reducing batch size.")
                    batch_size = max(batch_size // 2, min_batch_size)
                    if batch_size == min_batch_size:
                        print("Minimum batch size reached. Exiting.")
                        break
```

UNIVERSAL IDENTIFIER

STATE

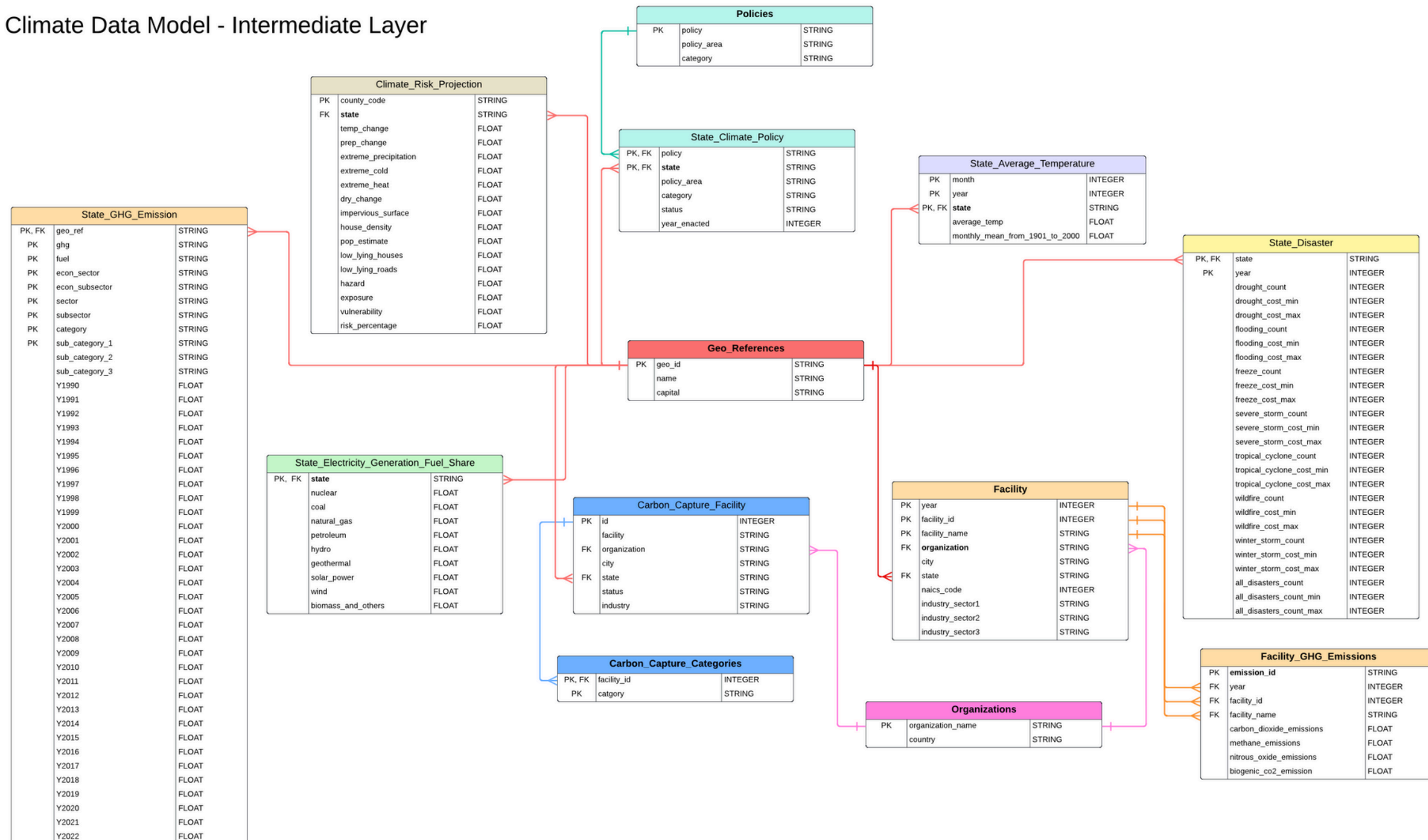
From the raw data, only saw a way to create relationships between tables through state field.

- Climate Risk Projections table only had county code
- Mix of state abbreviations and spelled out
- Had data for U.S. territories and some Canadian provinces

Created Geo References table using LLM to match state/territories/provinces abbreviation to full name and vice versa. Also used it to find the capitals.

INT ERD

US Climate Data Model - Intermediate Layer



**THANK
YOU**

