

# Compstak Analysis Progress: Part 2

William Clinton Co

Department of Economics, University of British Columbia

April 30, 2025

## Abstract

This is a continuation of Part 1. Reading part 1 is unnecessary and this document is structured to read part 2 as a stand alone version. This paper analyzes the CompStak commercial property dataset, focusing on classification consistency and geographic coverage. While ‘Property Type’ and ‘Property Subtype’ are generally populated, subtype fields are missing in 28% of records and show some inconsistency. Property IDs are stable, allowing reliable property counts. By mapping CompStak to DOE building categories, we estimate 35% national coverage. However, representation is uneven, with strong western and large-market bias. These findings highlight classification challenges and non-uniform geographic coverage.

# 1 Initial Analysis

Each field was reviewed individually to identify potential patterns, and all fields appeared relevant for analysis. The broader data set includes several types of industry classifications, such as ‘Property Type’, ‘Property Sub type’, ‘Space Type’, ‘Tenant SIC Code/Description’, and ‘Tenant NAICS Code/Description’. However, there is no clear documentation explaining the differences between these classifications or indicating which should be preferred. In our case, only ‘Property Type’ and ‘Property Sub type’ are available, so our analysis will focus on these variables. It is important to note that these classifications are not consistently available across both lease and sales records, which may introduce challenges to our analysis. Additionally, NAICS and SIC codes are not included in our data set because they are a premium add-on. Finally, ‘Space Type’ is only available for lease data, which may also pose limitations.

We observe “fill rate” (Table 1), wherein a low fill rate would correspond to a large portion of the entries having Nan or missing values. Unique values shows us that the categories used for property type and property sub type are consistent.

Table 1: Data set Initial Analysis

Data set	Type	Fill Rate	Unique Values
Sales	Property type	95%	8
Lease	Property type	99%	8
Sales	Property Subtype	75%	56
Lease	Property Subtype	66%	56

We also investigate whether property IDs are consistent. We expect that a given commercial

property would exhibit stable characteristics, with its location and industry classification remaining unchanged over time. Our analysis confirms this expectation: property IDs are stable and pass state consistency checks. However, it is important to note that property types and sub types are not consistent. See Figure 1 and Table 6

```
Property Id 646606 has multiple Property Types: [nan, 'Land']
Property Id 665103 has multiple Property Types: ['Other', 'Office']
Property Id 1019665 has multiple Property Types: [nan, 'Other']
Property Id 1235355 has multiple Property Types: ['Retail', 'Industrial']
Property Id 1251623 has multiple Property Types: ['Land', 'Office']
Property Id 1375074 has multiple Property Types: ['Land', 'Multi-Family']
Property Id 1470728 has multiple Property Types: ['Retail', 'Office']
Property Id 1702654 has multiple Property Types: ['Retail', 'Multi-Family']
Property Id 1721284 has multiple Property Types: ['Land', 'Industrial']
Property Id 1721935 has multiple Property Types: ['Retail', nan]
Property Id 1725734 has multiple Property Types: ['Office', 'Retail']
Property Id 1768020 has multiple Property Types: ['Land', 'Retail']
Property Id 1783304 has multiple Property Types: ['Land', nan]
Property Id 1818335 has multiple Property Types: ['Retail', 'Multi-Family']
Property Id 1825935 has multiple Property Types: ['Land', 'Multi-Family']
Property Id 1832887 has multiple Property Types: ['Land', nan]
Property Id 1973729 has multiple Property Types: ['Other', 'Retail']
Property Id 1981969 has multiple Property Types: ['Retail', 'Industrial']
Property Id 1992092 has multiple Property Types: ['Land', 'Retail']
Property Id 2017235 has multiple Property Types: ['Retail', 'Mixed-Use']
Property Id 2023360 has multiple Property Types: ['Retail', 'Office']
Property Id 2106390 has multiple Property Types: ['Retail', nan]
Property Id 2126846 has multiple Property Types: [nan, 'Retail']
Property Id 2186407 has multiple Property Types: [nan, 'Other']
Property Id 2295396 has multiple Property Types: [nan, 'Land']
Property Id 2310705 has multiple Property Types: ['Retail', 'Office']
Property Id 2423644 has multiple Property Types: ['Retail', 'Other']
Property Id 2720291 has multiple Property Types: ['Land', 'Retail']
Property Id 3141000 has multiple Property Types: ['Office', 'Other']
Property Id 3587132 has multiple Property Types: ['Retail', 'Mixed-Use']
```

Figure 1: Property ID Unstable Classifications

Despite concerns about inconsistency, this is not a major issue, as the number of unstable observations is relatively small compared to the total number of properties. Specifically, property type inconsistency is observed in only 39 cases, and property sub type inconsis-

tency in 96 cases. This pales in comparison to the 759,623 unique properties in our data set.

What is more concerning is the Nan surrounding each property type and sub type. Wherein the Nan values for property sub types can be as high as 28% of the data. See Table 2

Table 2: Error Associated with Property Types

Error Type	Industry Category		
	Type	Number	Percentage
Inconsistent Category	Property Type	39 / 759623	0.00513%
Inconsistent Category	Property Subtype	96 / 759623	0.01264%
Nan	Property Type	37020 / 759623	4.87%
Nan	Property Subtype	213255 / 759623	28.07%

## 2 Strategy

We will begin by matching the number of buildings. Specifically, we observe unique property IDs in the data set, allowing us to estimate the number of commercial properties in the United States. The total number of commercial properties is a relatively stable metric to study compared to more volatile measures such as valuations or square footage.

Using this approach, we can first compare aggregate numbers, starting with the total number of properties in the United States, and then work our way downward to industry-level and state-level comparisons.

We first assess how much of the national commercial property market is captured in our data set. We assume that the CompStak data set represents only a small subsection of the total U.S. market. However, a key concern is whether these observations are uniformly distributed across regions and industries, or if there are systematic biases. For example, are retail properties in California more likely to be reported than warehouses in Michigan? To begin, we pull external estimates of the total number of commercial properties from [DOE Dataset](#). Although these external figures are not entirely accurate, they provide a useful starting point for benchmarking and analyzing the coverage of our data set.

### 3 Methodology

The key to this analysis is matching the different category structure that exists within the DOE data set and the Compstak data set.

First we discuss the categories in the DOE dataset. There are direct matches. For example “industrial” appears as a category in both data sets. At the same time there are ambiguous data sets such as “speciality”, which can take the form of a casino to a recycling center. In order to work around this we do the matching into two steps. The first step is the direct mapping, where categories with the same name are matched one for one. The second step is the ambiguous mapping, where categories are mapped based on their subcategory. To understand this approach, we note how the data set is structured. There is a main category type along with a sub type. The ambiguous mapping maps based on the subcategory, disregarding the main category. It is also to note that no entries or Nan entries are automatically assigned to “others” category. This potentially overestimates our “others” category, as we see later on in [Figure 2](#). The mapping can be found in [Table 5](#).

Finally, we discuss the categories in the Compstak dataset. Ideally, we would want to match the DOE categories exactly to the Compstak categories. The reason being that DOE has 11 categories and Compstak has 8 categories ( Table 4) . Unfortunately, the Compstak categories also contains ambiguous definitions that are unattributable to the DOE categories, namely “Mixed-Use” and “Land”. We create mappings to address this by using their respective subcategories (Table 3).

As a result we end up with 6 categories (Hotel, Industrial, Multi-Family, Office, Other, Retail).

## 4 Results

We first look at Figure 2 We see that the scales of the categories are mostly consistent in both of the datasets. Furthermore, we see the overestimation bias we have with “other” category. Beyond that our results suggest consistency within the Compstak data set.

Next we look into the percentages of the DOE vs Compstak data set. We expect that the percentage of each category with respect to the total of each data set should be consistent. It is to note we remove “other” category due to overestimation. We see that the pie chart is mostly consistent with the DOE, suggesting valid consistent data. Figure 3

Assuming DOE dataset, contains the total commercial property in the us. We create a variable called coverage rate, which is the number of commercial properties of Compstak divided by the DOE data set. From this we observe a few things, “other” category is once again an outlier. Aside from this, we see that all other categories have decent coverage rates Figure 4 Figure 5 .

We also see decently overall coverage rate of almost 35% (assuming DOE is the total).

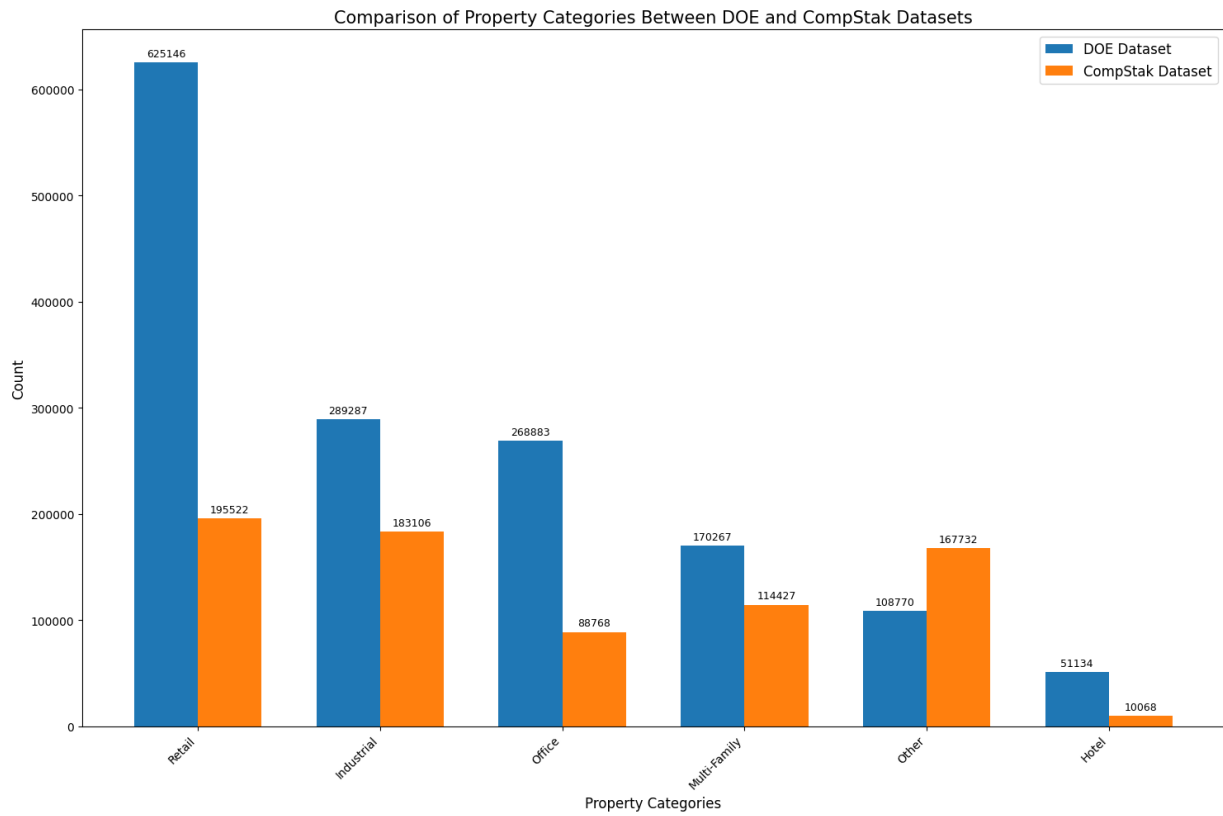


Figure 2

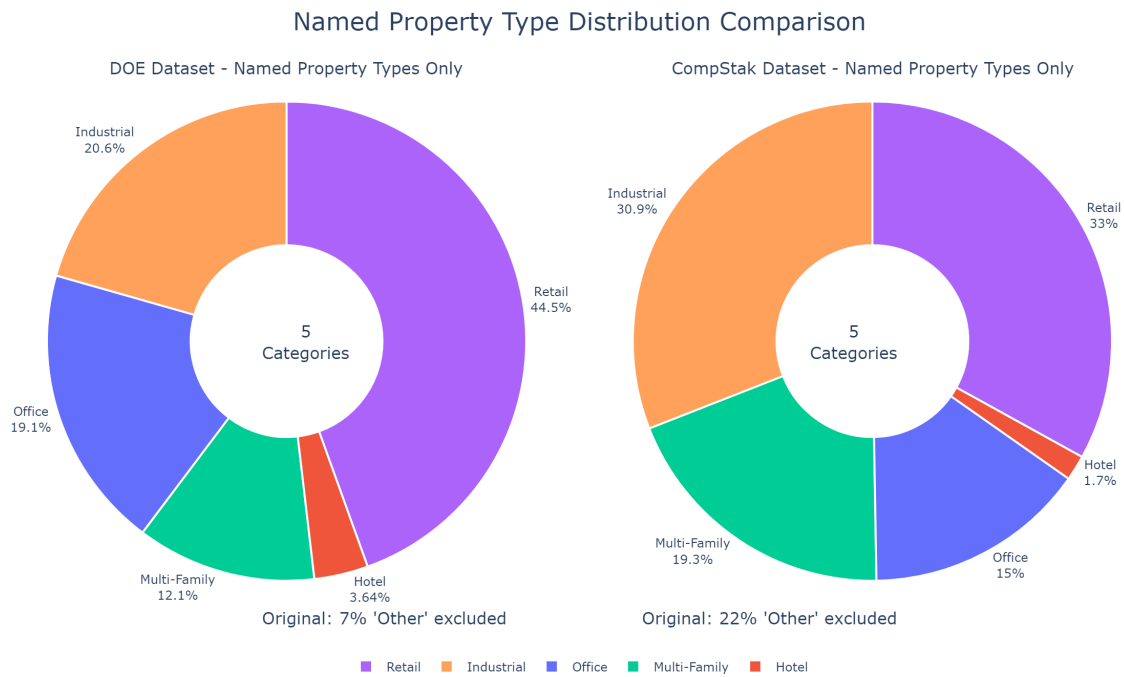


Figure 3

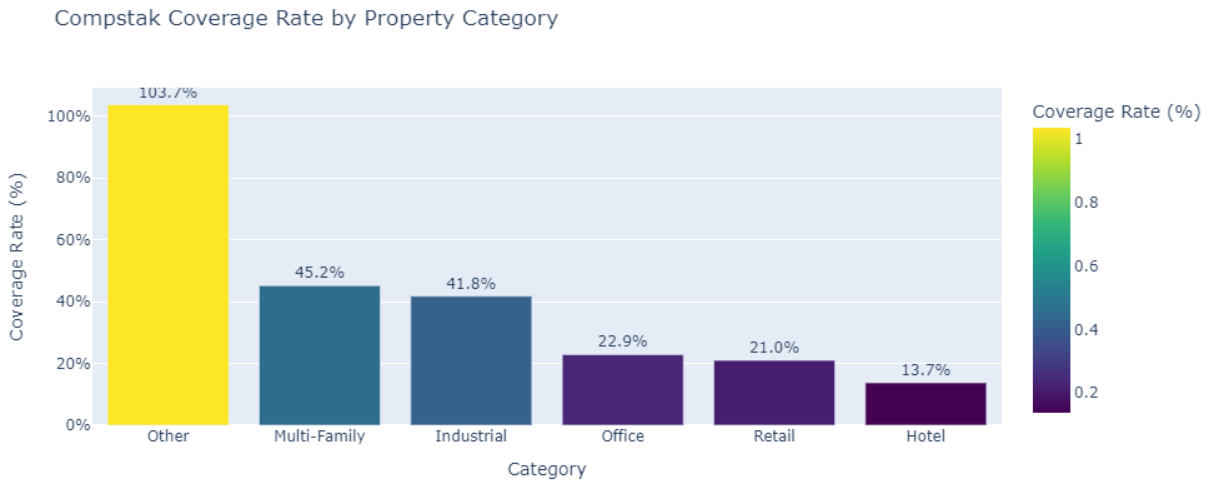


Figure 4

Overall Coverage: DOE Data Covered by Compstak

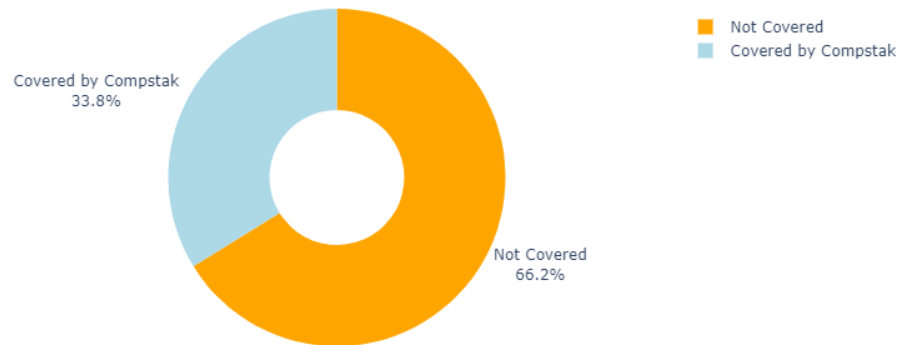


Figure 5



In line with our coverage rate analysis, we now plot by state. We see a big western bias for west coast states. Red indicates a high coverage. Figure 6 . In theory, uniform coverage would imply a consistent number across all states. This suggest that our data is ununiform and biased to western states. CA covers as high as 70%, this is in start difference to the 33% aggregate coverage rate we saw prior.

CompStak Coverage of DOE Data by State (%)

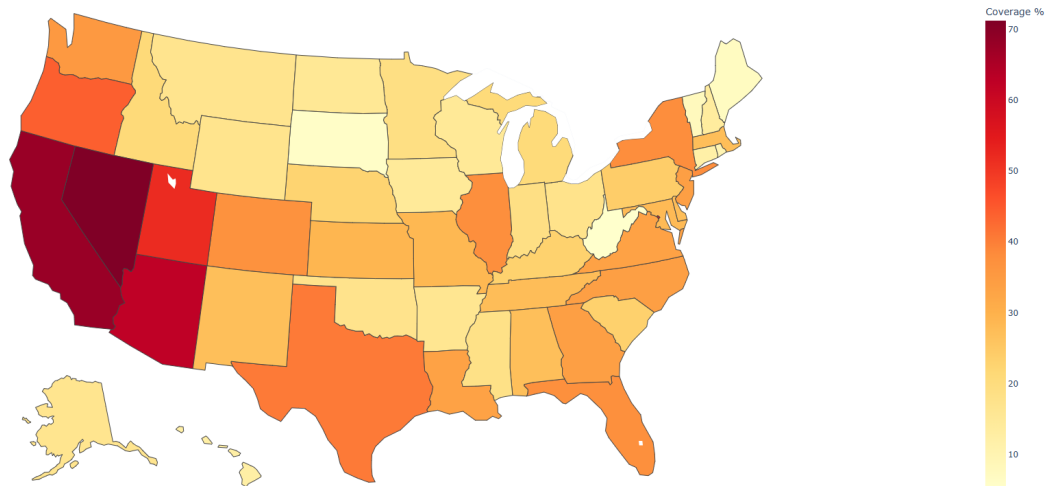


Figure 6

We have now establish a western bias to coverage rates. But now we also investigate weather the number of properties within the state is also a variable to understand coverage rates As shown in Figure 7 , there is evidence to support this theory. States with a greater number of commercial properties tend to have higher coverage rates. This pattern may be explained that there are fixed costs associated with data collection in each location, and that data providers are more likely to specialize in areas with larger commercial markets. As a result, the data set exhibits increasing returns to scale, which introduces bias into our observations.

Next, we use both state and categorization to produce a visualization. Figure 8 . From this

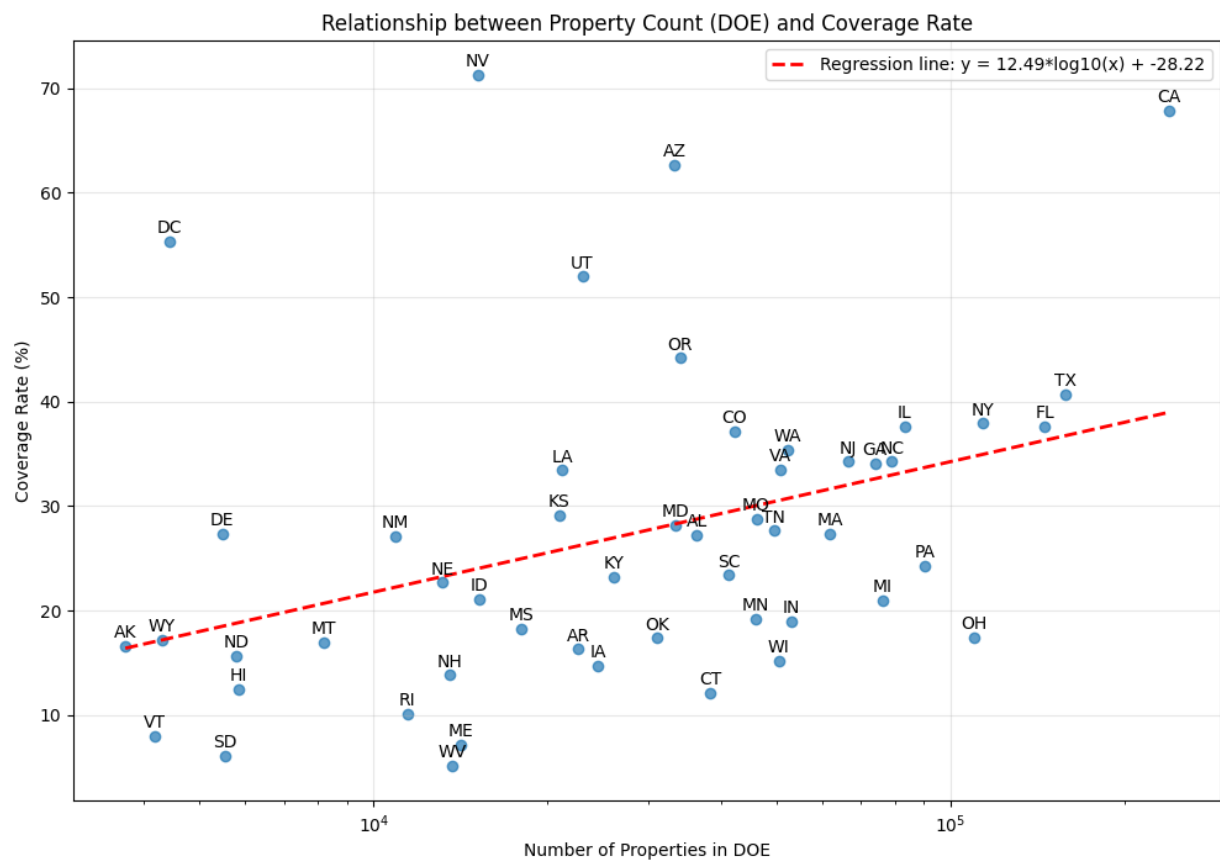


Figure 7

we see that “others” category is introducing bias into our data set. Therefore, we remove “others” category for our next visualization.

Figure 9. The unevenness of the data set is clearly illustrated in this visualization. Notably, the industrial and multifamily categories in Nevada, California, and Arizona show consistently higher coverage rates compared to other regions. Additionally, New York exhibits particularly high coverage for office properties relative to its industrial sector, which is an interesting deviation. These patterns further support the theory of increasing returns to scale in data collection. states with larger commercial property markets tend to have better data set representation.

## 5 Concluding Thoughts

Clustering the categories into smaller subset of types will be difficult to accomplish. The Compstak data set has a high Nan rate and does not suggest a completely consistent categorization methodology. There is no industry dictionary glossary to reference (unlike DOE, [CoStar Glossary](#)). With this said, It is best to stick to higher order property types, which is prone to less errors. It may also be worthwhile to use standardized property codes, which is available in the data dictionary given, but is a “premium add on”. The Compstak data set covers industries reasonably but the geography coverage seems to be western and size biased.

Next steps, will now be working on producing the data set with the following fields Lat/lon , State, ZCTA, Type of business

## 6 Appendix

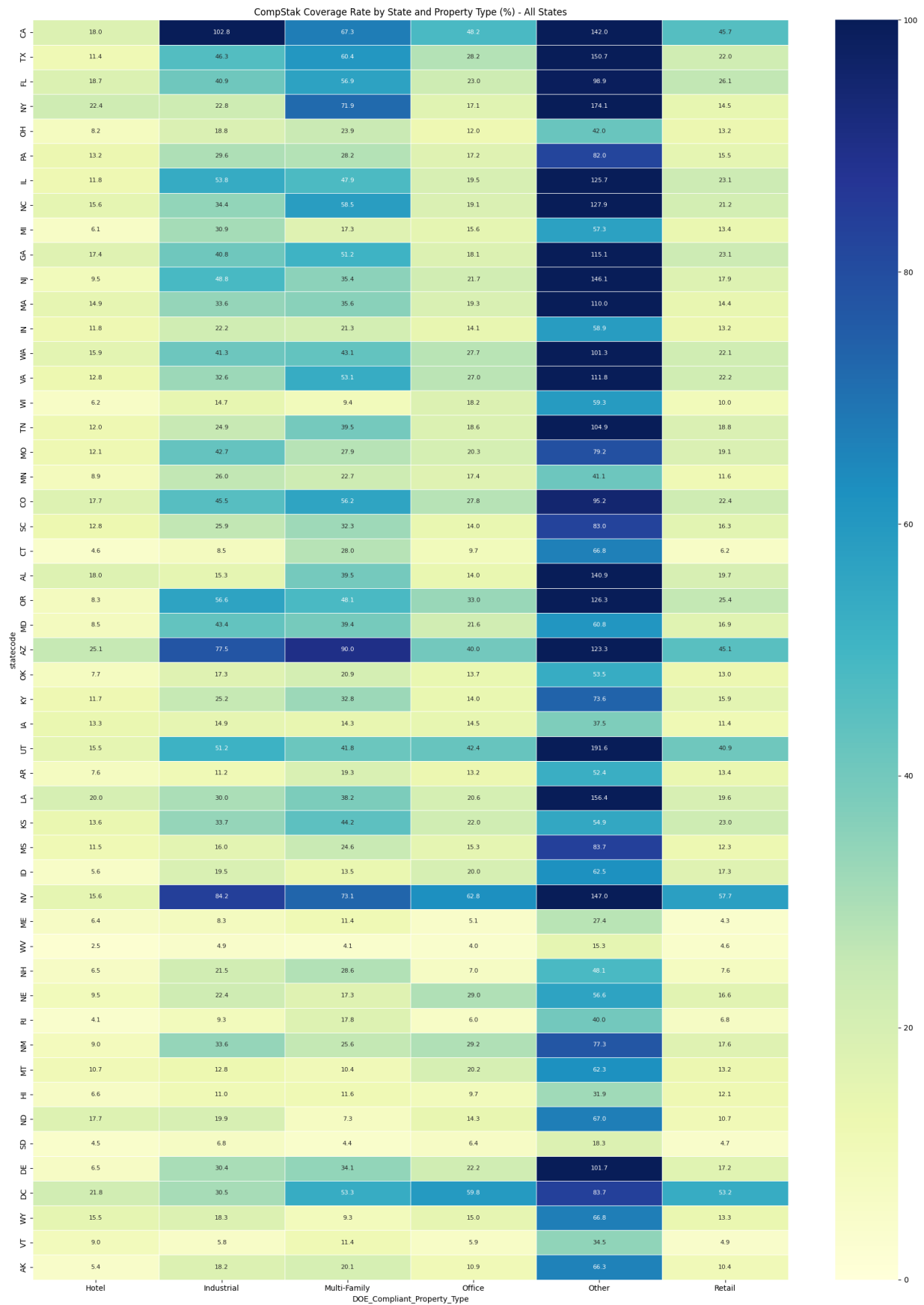


Figure 8  
12

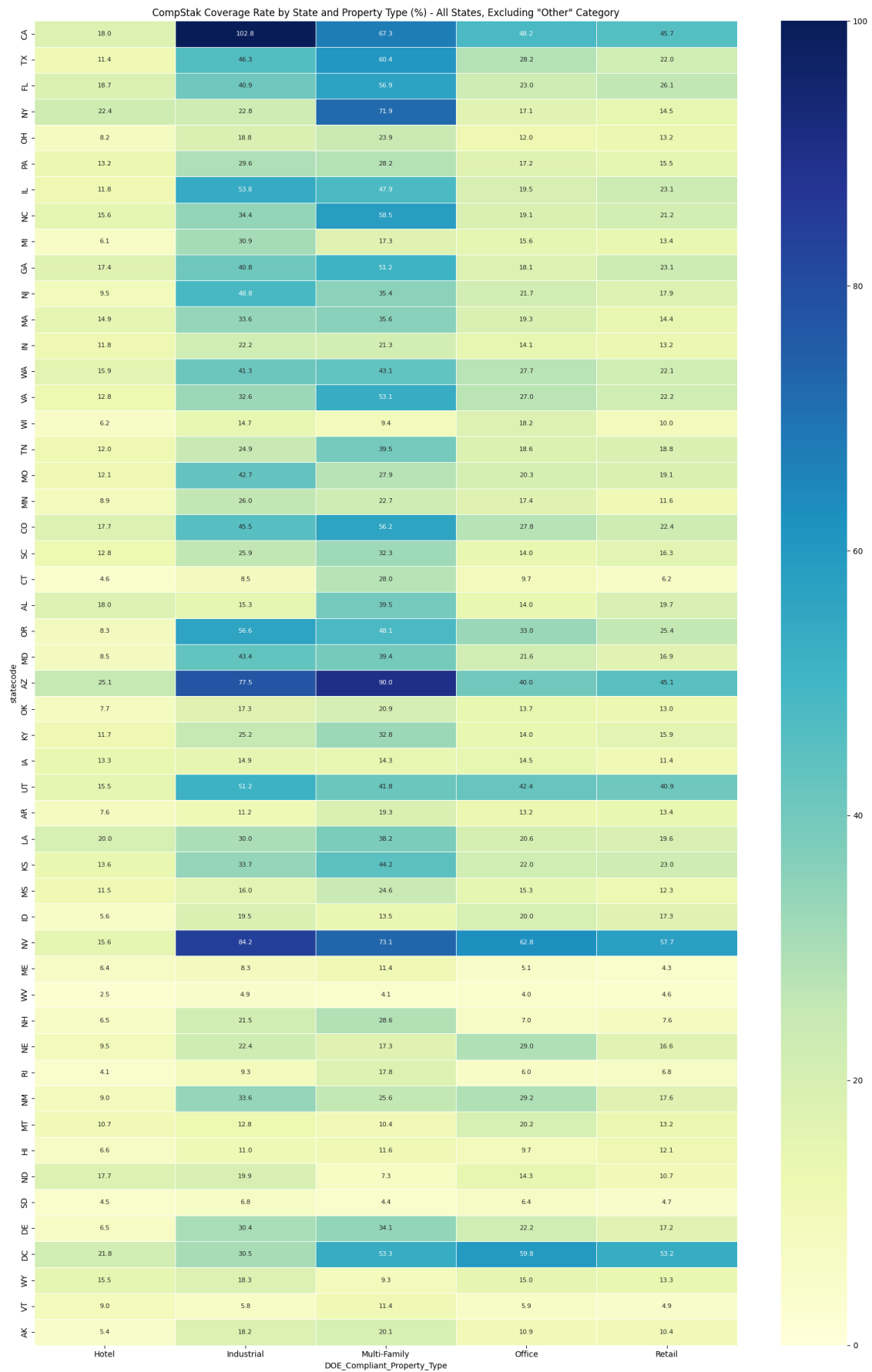


Figure 9

Table 3: Land and Mixed-Use (Compstak) Mapping

Property Type	DOE Category
Street Retail/Storefront	Retail
General Retail	Retail
Super-Regional Center/Mall	Retail
Shopping Centers	Retail
Neighborhood Shopping Center	Retail
Convenience/Strip Center	Retail
Community Shopping Center	Retail
Department Store	Retail
Restaurant/Bar	Retail
Fuel & Service Station	Retail
Freestanding	Retail
Automotive	Retail
Outlet	Retail
Drive Thru	Retail
Warehouse/Distribution	Industrial
Light Industrial	Industrial
Heavy Industrial	Industrial
Refrigerated/Cold Storage	Industrial
Manufacturing	Industrial
Special Industrial	Industrial
Industrial Outdoor Storage	Industrial

Property Type	DOE Category
Flex/R&D	Industrial
Processing	Industrial
Life Science/Lab	Industrial
Business Park	Office
Professional Building	Office
Financial Building	Office
Bank	Office
Creative	Office
Apartments	Multi-Family
Student Housing	Multi-Family
Mobile Home Park	Multi-Family
Condominium	Multi-Family
Senior Housing	Multi-Family
Housing	Multi-Family
Hospitality Related	Hotel
Self-Storage	Other
Vacant Land	Other
Mixed-Use	Other
Medical/Healthcare	Other
Hospital/Healthcare Facility	Other
Day Care Facility	Other
Parking	Other
Funeral/Mortuary	Other

Property Type	DOE Category
Communication/Data Center	Other
Assembly/Meeting Place	Other
Municipality/Public Service	Other
Educational/School	Other
Community/Recreation Center	Other
Sports & Recreation	Other
Transportation	Other
Special Purpose	Other
Live/Work	Other
Under Construction	Other

Table 4: Categories within DOE and Compstak

Compstak Property Types (8 total)	DOE Property Types (11 total)
Hotel	Flex
Industrial	General Retail
Land	Health Care
Mixed-Use	Hospitality
Multi-Family	Industrial
Office	Multi-Family
Other	Office
Retail	Retail
	Specialty



Compstak Property Types (8 total)	DOE Property Types (11 total)
	Sports & Entertainment
	Unknown

Table 5: DOE Mapping

Property Type	Subtype	Standardized Category
Industrial	All	Industrial
Multi-Family	All	Multi-Family
Office	All	Office
Retail	All	Retail
General Retail	All	Retail
Hospitality	All	Hotel
Flex	Light Distribution	Industrial
Flex	Light Manufacturing	Industrial
Flex	R&D	Industrial
Flex	Showroom	Retail
Flex	Telecom Hotel/Data Hosting	Other
Flex	All Others	Industrial
Health Care	Assisted Living	Other
Health Care	Congregate Senior Housing	Other
Health Care	Continuing Care Retirement Community	Other
Health Care	Hospital	Hotel
Health Care	Rehabilitation Center	Other

Property Type	Subtype	Standardized Category
Health Care	Skilled Nursing Facility	Other
Specialty	Airplane Hangar	Other
Specialty	Airport	Other
Specialty	Auto Salvage Facility	Other
Specialty	Car Wash	Retail
Specialty	Cement/Gravel Plant	Industrial
Specialty	Cemetery/Mausoleum	Other
Specialty	Chemical/Oil Refinery	Industrial
Specialty	Contractor Storage Yard	Industrial
Specialty	Correctional Facility	Other
Specialty	Drive-in Movie	Other
Specialty	Landfill	Other
Specialty	Lodge/Meeting Hall	Hotel
Specialty	Lumberyard	Industrial
Specialty	Marina	Other
Specialty	Movie/Radio/TV Studio	Other
Specialty	Parking Garage	Other
Specialty	Parking Lot	Other
Specialty	Police/Fire Station	Other
Specialty	Post Office	Other
Specialty	Public Library	Other
Specialty	Radio/TV Transmission Facilities	Other
Specialty	Railroad Yard	Industrial

Property Type	Subtype	Standardized Category
Specialty	Recycling Center	Industrial
Specialty	Religious Facility	Other
Specialty	Residential Income	Multi-Family
Specialty	Schools	Other
Specialty	Self-Storage	Other
Specialty	Shelter	Other
Specialty	Shipyard	Industrial
Specialty	Sorority/Fraternity House	Other
Specialty	Trailer/Camper Park	Other
Specialty	Utility Sub-Station	Other
Specialty	Water Retention Facility	Other
Specialty	Water Treatment Facility	Other
Specialty	Winery/Vineyard	Other
Specialty	All Others	Other
Sports & Entertainment	Amusement Park	Other
Sports & Entertainment	Baseball Field	Other
Sports & Entertainment	Casino	Hotel
Sports & Entertainment	Golf Course/Driving Range	Other
Sports & Entertainment	Horse Stables	Other
Sports & Entertainment	Race Track	Other
Sports & Entertainment	Skating Rink	Other
Sports & Entertainment	Swimming Pool	Other
Sports & Entertainment	Theater/Concert Hall	Other

Property Type	Subtype	Standardized Category
Sports & Entertainment	All Others	Other
Unknown	All	Other

Table 6: Property Sub type Inconsistency

Property ID	Inconsistent Property Subtypes
21036	General Retail, nan
353844	nan, Vacant Land
354461	nan, Municipality/Public Service
374802	nan, Vacant Land
398906	nan, Vacant Land
415099	General Retail, nan
418514	Apartments, Sports & Recreation
420902	nan, Apartments
422199	Outlet, Vacant Land
433935	nan, Municipality/Public Service
434400	nan, Vacant Land
443800	Vacant Land, Super-Regional Center/Mall
444319	Apartments, nan
445159	General Retail, nan
448126	Apartments, nan
449199	General Retail, nan
466525	General Retail, nan

Property ID	Inconsistent Property Subtypes
476566	Parking, Apartments
485298	nan, Shopping Centers
489216	nan, Parking
490863	nan, Automotive
491210	Apartments, General Retail
491863	Apartments, Vacant Land
496416	nan, Apartments
508849	General Retail, Shopping Centers
520174	nan, Vacant Land
538162	Vacant Land, nan
567596	Super-Regional Center/Mall, Neighborhood Shopping Center
581027	nan, General Retail
581309	Apartments, nan
623357	Vacant Land, nan
624813	Condominium, Apartments
633831	Apartments, nan
646161	General Retail, nan
669817	nan, Apartments
679028	Apartments, nan
699484	nan, Apartments
702675	Apartments, nan
703449	nan, Apartments
731139	Apartments, Convenience/Strip Center

Property ID	Inconsistent Property Subtypes
742440	nan, Apartments
745417	Vacant Land, nan
754857	nan, Apartments
755143	General Retail, nan
757418	Apartments, nan
849172	Flex/R&D, Business Park
1204302	General Retail, nan
1211995	nan, Sports & Recreation
1212507	Special Purpose, nan
1235355	General Retail, Vacant Land
1254651	Self-Storage, nan
1255443	Vacant Land, Condominium
1261872	nan, Mixed-Use
1272015	Freestanding, General Retail
1311302	General Retail, Freestanding
1418012	Warehouse/Distribution, Special Industrial
1421693	Apartments, nan
1431770	nan, Vacant Land
1448113	nan, General Retail
1449588	Apartments, Financial Building
1451345	nan, Apartments
1684081	nan, Light Industrial
1705515	Manufacturing, Light Industrial

Property ID	Inconsistent Property Subtypes
1721935	Day Care Facility, nan
1722765	Super-Regional Center/Mall, Convenience/Strip Center
1724884	Vacant Land, nan
1725903	General Retail, Community Shopping Center
1743582	nan, General Retail
1765432	Apartments, nan
1822722	Vacant Land, nan
1858406	nan, Apartments
1866157	Hospitality Related, Apartments
1922481	nan, Apartments
1929300	Condominium, nan
2049015	General Retail, nan
2050907	Apartments, nan
2054074	nan, Super-Regional Center/Mall
2057688	nan, Vacant Land
2096757	nan, Vacant Land
2106390	Automotive, nan
2144456	Parking, Restaurant/Bar
2266042	Apartments, nan
2288779	Parking, Warehouse/Distribution
2292241	nan, Warehouse/Distribution
2295396	nan, Vacant Land
2330278	Apartments, nan

---

Property ID	Inconsistent Property Subtypes
2331364	Vacant Land, General Retail
2423080	nan, Vacant Land
2425229	nan, Restaurant/Bar
2720291	Vacant Land, Mixed-Use
3417569	nan, Medical/Healthcare
3417659	nan, Mixed-Use
3418319	nan, Mixed-Use
3464690	Warehouse/Distribution, Manufacturing
3575908	Community Shopping Center, Vacant Land
3587132	Apartments, Mixed-Use

---

## References