

Benchmarking Natural Language to Visualization Models



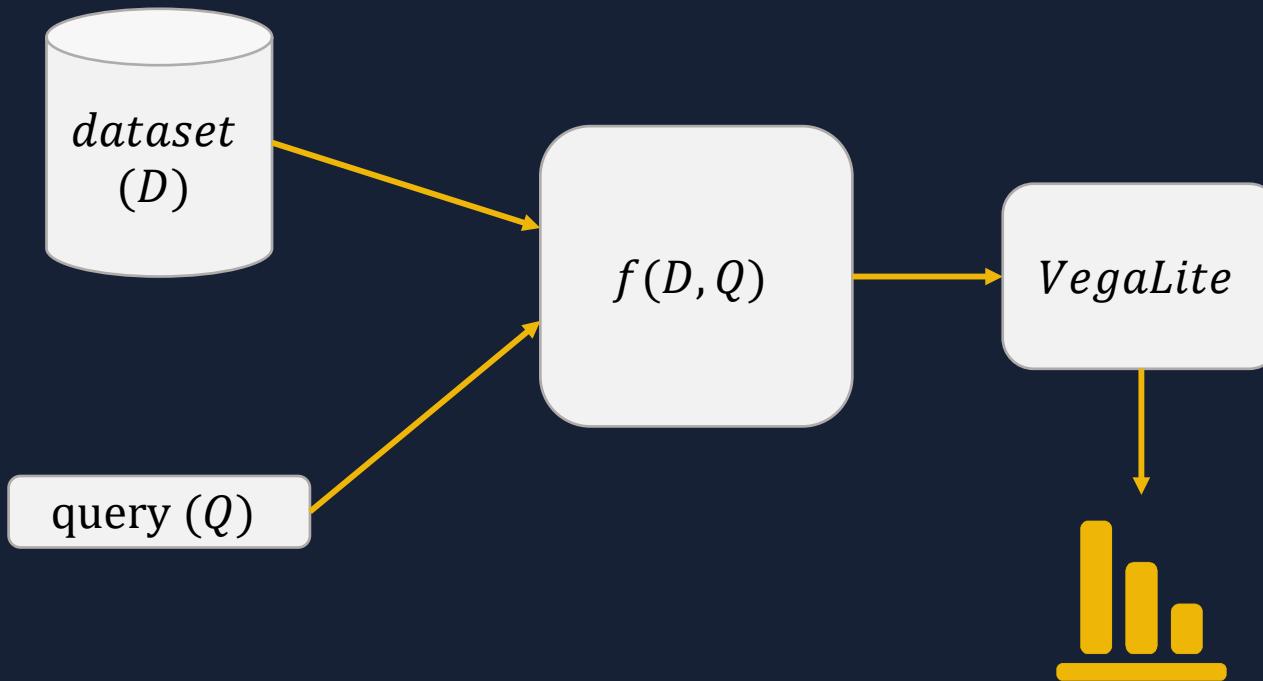
6.S079 Final Project
Enrique Casillas

Motivation

- Data visualization is important
- Turning a thought (natural language) into a visualization is the ultimate goal of visualization platforms



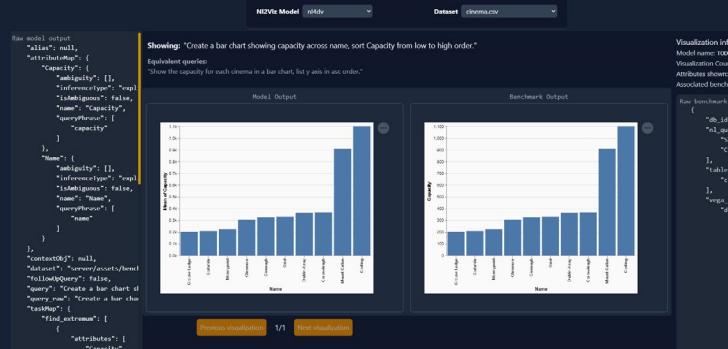
NL2Viz Models



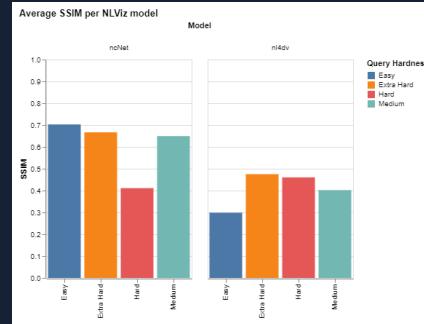
Project Overview

- Primary goal — benchmark NL2Viz models by using common metrics

Web-based tool – **n12viz**



Evaluation of two models, **ncNet** and **n14dv** using the **nvBench** NL2Viz benchmark



Web-based tool – **nl2viz**



nl2Viz Goals

- Visually **compare** a model's output to a benchmark in a couple of clicks
- Get a feel for what the benchmark looks like and **expects**
- Help make your own **conclusions** as to which model is better
- Provide **any** query and receive the model's output

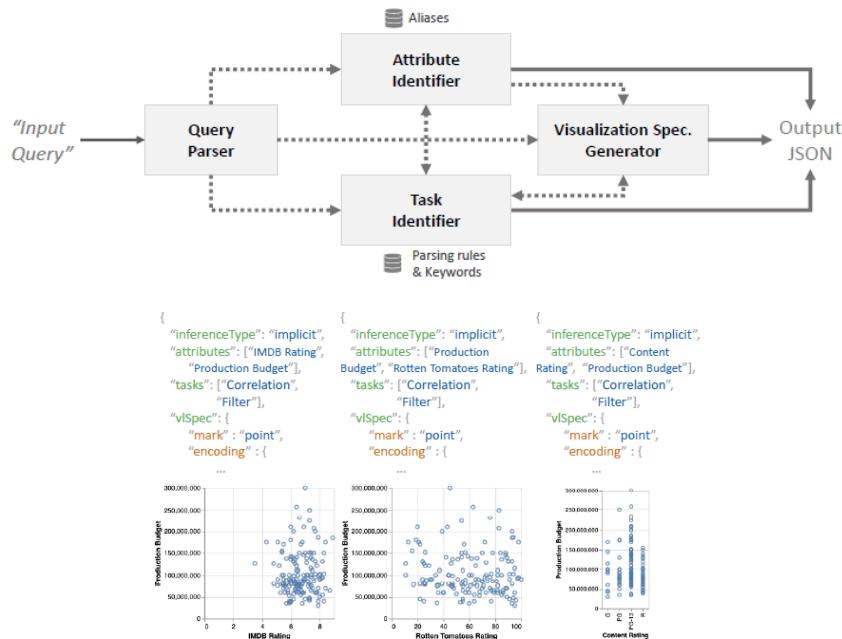
Demo



Model Evaluations



nl4dv



Create a bar chart showing the top 5 states with the most confirmed cases until 2021-03-08

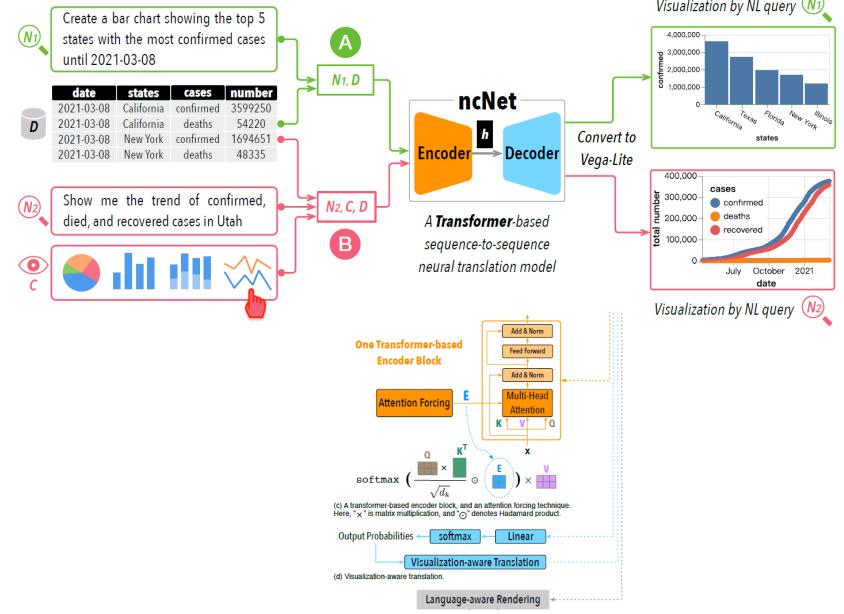
date	states	cases	number
2021-03-08	California	confirmed	54220
2021-03-08	New York	confirmed	1694651
2021-03-08	New York	deaths	48335

Show me the trend of confirmed, died, and recovered cases in Utah



Visualizations generated by NL query (N1)

ncNet



Existing benchmark — nvBench

{ }

```
"vis_query": {  
    "vis_part": "Visualize SCATTER",  
    "data_part": {  
        "sql_part": "SELECT FacID , count(*) FROM Faculty AS T1  
JOIN Student AS T2 ON T1.FacID = T2.advisor GROUP BY T1.FacID",  
        "binning": ""  
    },  
    "VQL": "Visualize SCATTER SELECT FacID , count(*) FROM Faculty  
AS T1 JOIN Student AS T2 ON T1.FacID = T2.advisor GROUP BY  
T1.FacID"  
},  
"chart": "Scatter",  
"hardness": "Medium",  
"db_id": "activity_1",  
"vis_obj": {. . . },  
"nl_queries": [. . . ],
```

How do we compare visualizations?

Common Methods



Crowdsource



Ask Experts



Don't evaluate

Alternative Methods



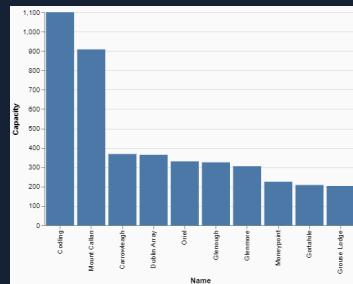
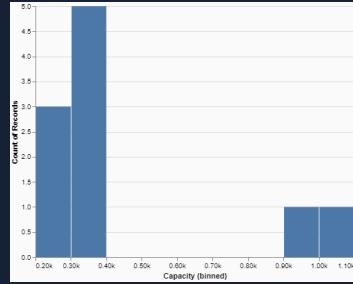
Compare specs



Image similarity

Structural Similarity Index Measure (SSIM)

n14dv



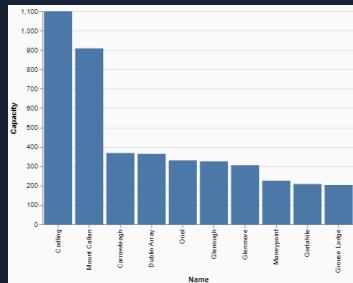
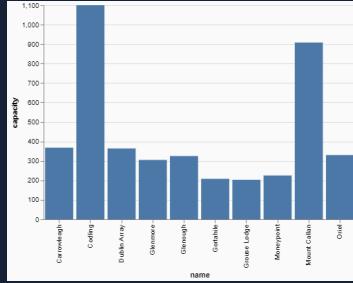
benchmark

$$SSIM(V_1, V_2)$$

0.7164527617582078

Structural Similarity Index Measure (SSIM)

ncNet

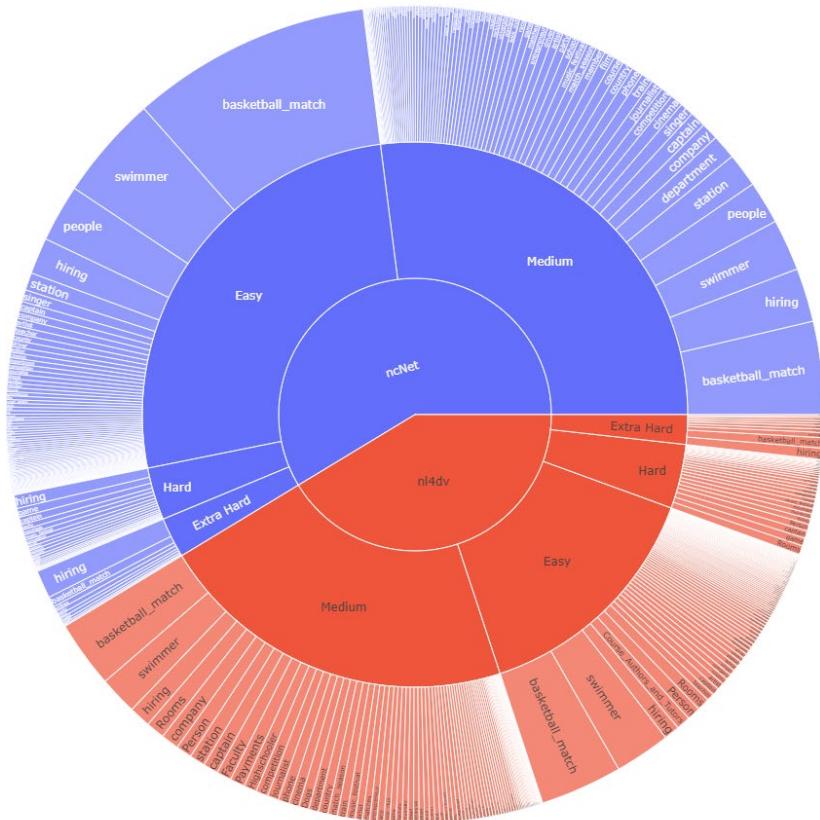


benchmark

$$SSIM(V_1, V_2)$$

0.898257440028004

Results

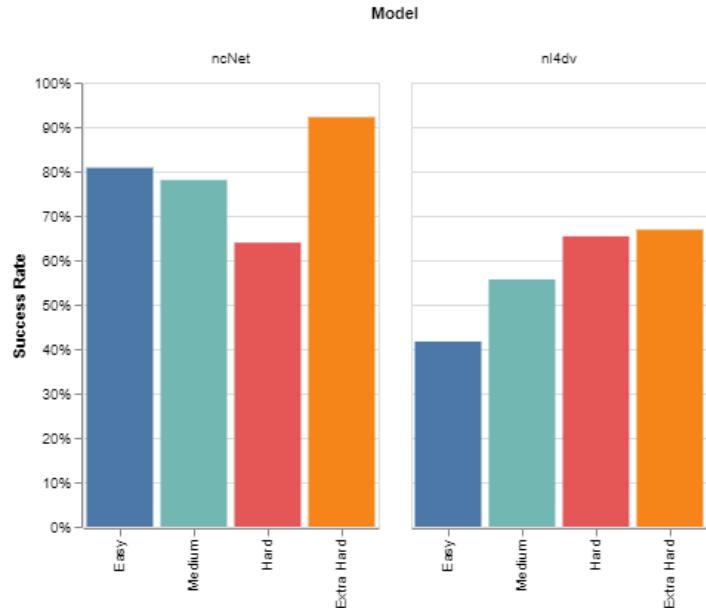


18,490 queries across 2 NL2Viz models

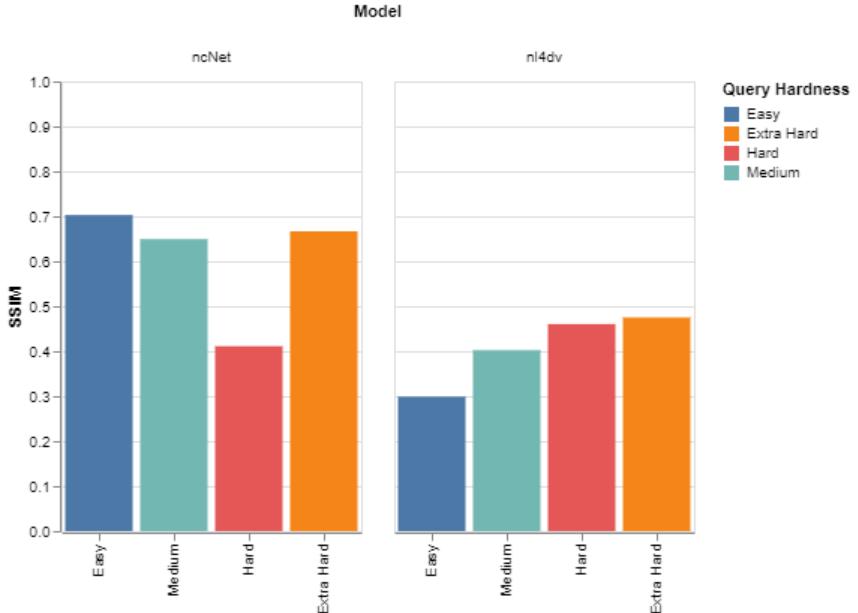
4 “difficulty” categories

141 unique datasets

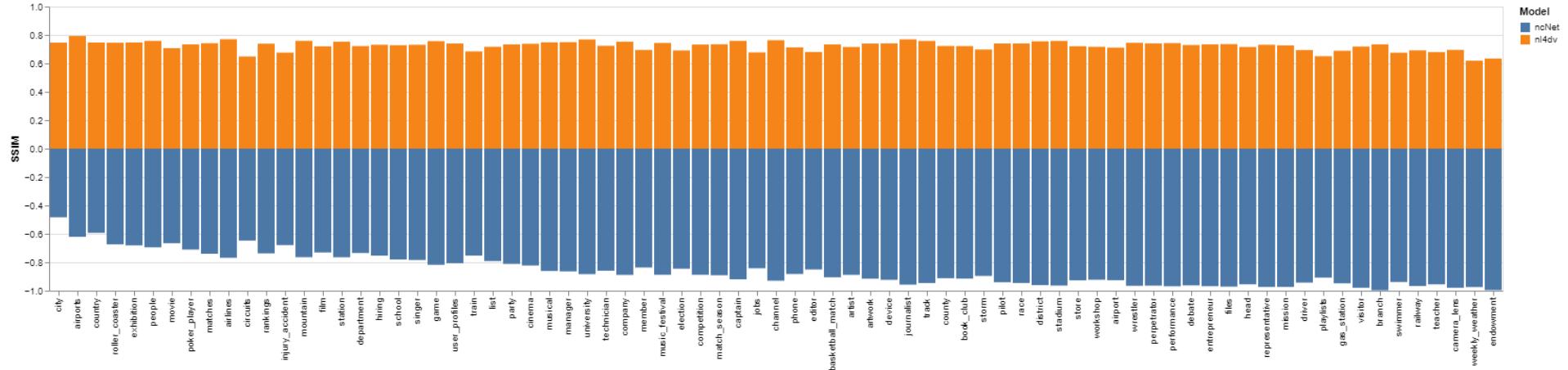
Success Rate by Model



Average SSIM per NLViz model



Average SSIM comparison with benchmark visualizations [0, 1]



Next Steps

- Find/combine more metrics for model evaluation, or come up with a new one
- Evaluate more NL2Viz models using existing pipeline
- Online tool – generate SSIM metrics on-demand
- Online tool – built-in specification editor

Thank you!

Image References

- <https://logos-world.net/wp-content/uploads/2021/10/Tableau-Emblem.png>
- <https://user-images.githubusercontent.com/315810/92254613-279c8000-ee9f-11ea-9b73-5622a7d95f3f.png>
- https://upload.wikimedia.org/wikipedia/commons/thumb/5/58/Vega-Lite_Logo.svg/2560px-Vega-Lite_Logo.svg.png