

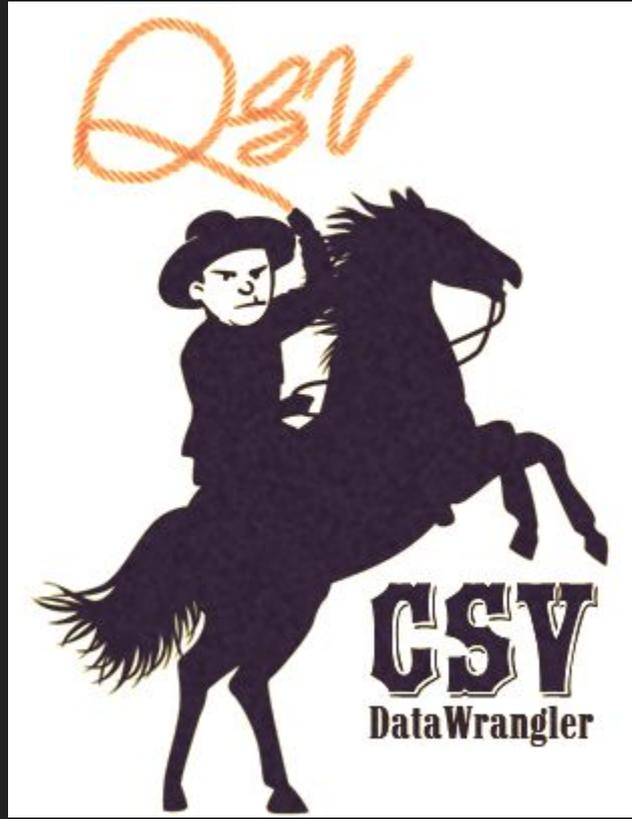
# qsv

A blazing-fast, multi-platform, command-line,  
Data-wrangling toolkit

Joel Natividad  
csv.conf.v8  
May 2024  
v2.0.0









<https://github.com/jqnatividad/qsv/releases/tag/0.128.0>

A little history....

# Born of Open Data

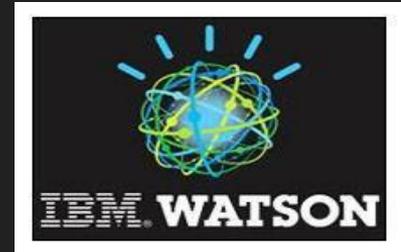


1<sup>st</sup> US  
Professional  
Services  
Partner

2013



Reinvent 311  
Winner  
Jan 2014



Finalist  
IBM Watson Mobile Challenge  
May 2014

2016  
Apr

# OpenGov Acquires Open Data Leader Ontodia



*Combination of Ontodia's Open Data and Performance Management Capabilities with OpenGov's leading Financial Intelligence and Transparency Platform Ushers in Next Wave of Open and Efficient Government*

**REDWOOD CITY, Calif. –April 13, 2016** – Today OpenGov, the world leader in government financial intelligence, planning, and transparency, expands its platform with the acquisition of Ontodia. Ontodia is the leading provider of Open Data and performance management solutions using CKAN, the premier open-source data-portal for governments around the world. CKAN powers Data.gov, the home of U.S. Government's Open Data initiative.



**~100 installations  
across the US and  
there was one  
recurring problem...**

# Data Quality

TECHNOLOGY

# For Big-Data Scientists, 'Janitor Work' Is Key Hurdle to Insights

By STEVE LOHR AUG. 17, 2014

- EMAIL
- FACEBOOK
- TWITTER
- SAVE
- MORE



Technology revolutions come in measured, sometimes foot-dragging steps. The lab science and marketing enthusiasm tend to underestimate the bottlenecks to progress that must be overcome with hard work and practical engineering.

The field known as “big data” offers a contemporary case study. The catchphrase stands for the modern abundance of digital data from many sources — the web, sensors, smartphones and corporate databases — that can be mined with surprising discoveries and insights. It is a smarter, data-driven economy in every field. That is why it is the economy’s hot new job.

Yet far too much hand-wringing by data scientists call “data munging” and “data



red. Data scientists, according to interviews a  
from 50 percent to 80 percent of their time  
lane labor of collecting and preparing unruly

# “Data Wrangling” Challenges

- Brittle data pipelines
- Larger & larger datasets
- “Regular” tools cannot scale (i.e. Excel)
- Specialized tools
  - Platform specific
  - Expensive
- Specialized “data science” skills
- Slow
  - Ramp-up time
  - Preparation time
  - Execution time



# datHere launching 2020

The band gets back  
together to take on  
Data Quality...

datHere - Data Engineering Infrastructure X +

dathere.com

Info@dathere.com +1 732 707 1866

Home Services Solutions Blog About Contact

## DATA INFRASTRUCTURE ENGINEERING

*Standards-based, best-of-breed, open source solutions to make your Data Useful, Usable & Used*

Learn More

— SERVICES —

Reduce your  
Data Pipeline Debt

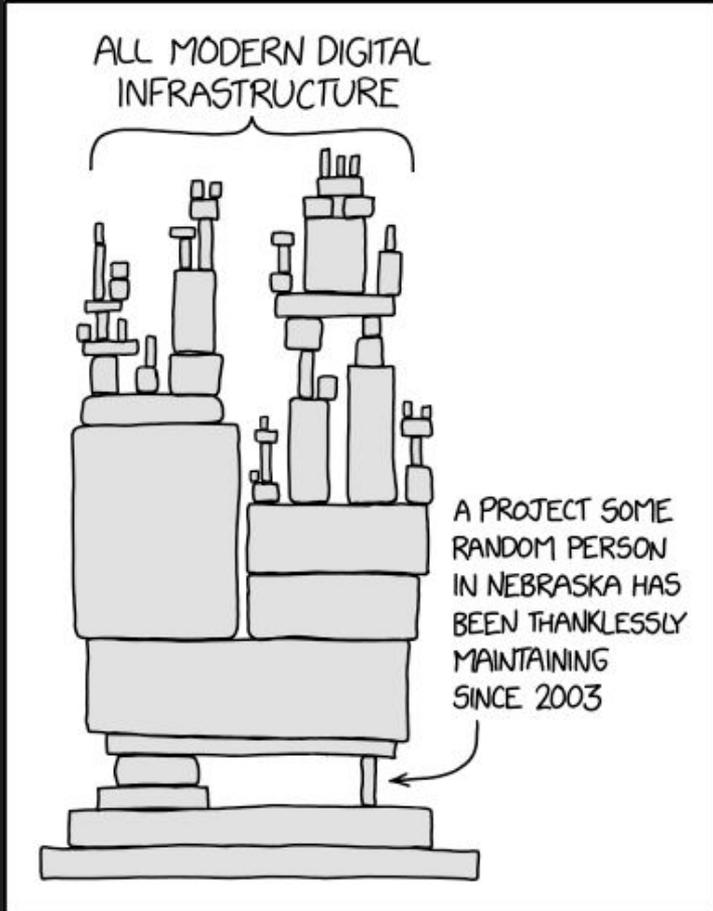


*Standards-based, best-of-breed, open source solutions  
to make your Data Useful, Usable & Used*

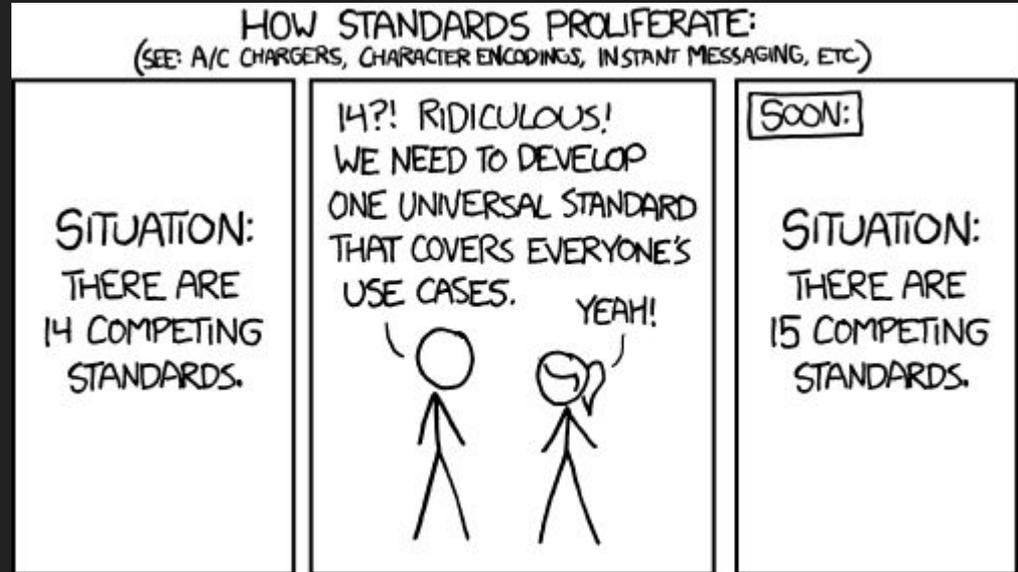
We deploy & co-create Data Infrastructure:

- Open Data Portals
- Internal Data Exchange
- Data Libraries
- Data Pipelines
- Water Data Practice

# Open Source



# Open Standards



But then

2020

happened...





### When Satan Met 2020

4,684,366 views • Dec 3, 2020

196K

4.2K

SHARE

SAVE



# We needed a “Data Wrangler”

- Works with a universal data format
- Cross-platform
- Fast, blazing Fast!
- Open Source
- Easy to Learn
- Easy to Use for initial investigations
- But powerful enough to integrate into mission-critical data pipelines



CSV



**CSV**  
DataWrangler

# Origins

It all started with a **failed pilot** with a Hedge Fund to build an Internal Data Portal

- Brand new startup during COVID
- Data Portals anybody?
- An Internal Data Catalog Pilot, populated with latest metadata
- Traditional metadata ingestion pipeline (csvkit) was too slow
- Forked xsv to start qsv...

# qsv "Data Wrangler" Goals

- Works with a universal data format
- Cross-platform
- Open Source
- Easy to Learn
- Easy to Use for initial investigations
- But powerful enough to integrate into mission-critical data pipelines

*CSV, EXCEL, JSON, JSONL,  
POSTGRESQL, SQLITE, PARQUET,  
DATA PACKAGE, AVRO +  
RECOGNIZES 130 FILE FORMATS*

*LINUX, MACOS + WINDOWS*

*FAST! BLAZING "SPEEDY GONZALES" FAST!!!*

# How fast is Blazing, “Speedy Gonzales” fast?

For a 1 million row sample of NYC’s 311 data (41 columns, 520 mb):

- 11 “streaming” summary statistics in **0.204 secs**
- 21 more statistics & infer dates(19 formats recognized) in **1.97 secs**
- Frequency table in **1.129 secs**
- Count rows in **0.05 secs**
- Validate against RFC 4180 CSV standard in **0.5 secs**
- Validate against a JSON Schema in **1.266 secs**
- Run a simple SQL query in **0.05 secs**, a SQL aggregation in **0.082 secs** & a very inefficient SQL aggregation in **0.144 secs**
- Reverse geocode WGS84 coordinate against Geonames in **3.59 secs**
- And more...

<https://qsv.dathere.com/benchmarks>

field	type	is_ascii	sum	max	max	range	min_length	max_length	mean	stddev	variance	nullcount	max_precision	sparsity	mad	lower_outer_fence	lower_inner_fence
Unique Key	Integer		32687965858032	11465364	48478173	37012809	8	8	32687965.85	9013895.3358	81250309125279.6000	0		0	7577800.5	-19639208.5	2803282.25
Created Date	DateTime			2010-01-01 T00:00:00+00:00	2020-12-23 T01:25:51+00:00	4009.05962			2015-11-10 T18:05:22.615+00:00	1155.01606	1334062.09198	0		0	965.58623	1997-01-08T17:56:34	2005-02-08T08:58:19
Closed Date	DateTime			1900-01-01 T00:00:00+00:00	2100-01-01 T00:00:00+00:00	73049			2015-11-14 T10:16:16.743+00:00	1314.70016	1728436.50813	28619		0.0286	955.59374	1997-04-12T11:33:24	2005-04-09T10:53:21
Agency	String	FALSE		3-1-1	TLC		3	42				0		0			
Agency Name	String	FALSE		3-1-1	Valuation Policy		3	82				0		0			
Complaint Type	String	TRUE		../WEB-INF/web.xml;x=	ZTESTINT		3	41				0		0			
Descriptor	String	TRUE		1 Missed Collection	unknown odor/taste in drinking water (QA6)		0	80				3001		0.003			
Location Type	String	TRUE		1-, 2- and 3- Family Home	Wooded Area		0	36				239131		0.2391			
Incident Zip	String	TRUE		*	XXXXX		0	10				54978		0.055			
Incident Address	String	TRUE		**	west 155 street and edgecomb e avenue		0	55				174700		0.1747			
Street Name	String	TRUE		*	wyckoff avenue		0	55				174720		0.1747			
Cross Street 1	String	TRUE		1 AVE	mermaid		0	32				320401		0.3204			
Cross Street 2	String	TRUE		1 AVE	surf		0	35				323644		0.3236			
Intersection Street 1	String	TRUE		1 AVE	flatlands AVE		0	35				767422		0.7674			
Intersection Street 2	String	TRUE		1 AVE	glenwood RD		0	33				767709		0.7677			
Address Type	String	TRUE		ADDRESS	PLACENA ME		0	12				125802		0.1258			
City	String	TRUE		*	YORKTOWN HEIGHTS		0	22				61963		0.062			
Landmark	String	TRUE		1 AVENUE	ZULETTE AVENUE		0	32				912779		0.9128			

COMPREHENSIVE SUMMARY STATS IN 1.97 SECONDS!

```

SELECT
  A.Agency,
  A.Borough,
  COUNT(*) AS total_incidents,
  SUM(
    CASE
      WHEN A."Complaint Type" LIKE 'Noise%' THEN 1
      ELSE 0
    END
  ) AS noise_related_incidents,
  SUM(
    CASE
      WHEN A.Status = 'Closed' THEN 1
      ELSE 0
    END
  ) AS closed_incidents,
  SUM(
    CASE
      WHEN A.Status != 'Closed' THEN 1
      ELSE 0
    END
  ) AS open_incidents,

  SUM(
    CASE
      WHEN POSITION('Water' IN A."Complaint Type") > 0 THEN 1
      ELSE 0
    END
  ) AS water_related_incidents,
  MAX(LENGTH (A."Complaint Type")) AS max_complaint_type_length,
  SUM(
    CASE
      WHEN UPPER(A."Complaint Type") = UPPER(A."Complaint Type")
      THEN LENGTH (A."Complaint Type")
      ELSE 0
    END
  ) AS sum_complaint_type_lengths,
  COUNT(DISTINCT A."Complaint Type") AS distinct_complaint_types
FROM
  read_csv ('NYC_311_SR_2010-2020-sample-1M.csv') A
GROUP BY
  A.Agency,
  A.Borough
ORDER BY
  total_incidents DESC;

```

**ANSWERED IN 0.144 SECONDS!**



+



polars

# How is it so Fast?

by standing on the  
Shoulders of Giants & the  
Ecosystem

- Rust
- Multi-threaded, Multi-I/O
- Performance architecture
  - Indexed access
  - Various caching techniques
  - Performance oriented memory allocator
- Built on a solid foundation (xsv)
- Polars Dataframes Engine
- Vibrant Rust & Polars Ecosystems

# Why the Obsessive Need for Speed?

What does it unlock?

- Big Data is getting Bigger
- Embedding into other Systems
- Quicker Data Investigations
- Enables new Data Workflows
  - Preemptive metadata inferencing
  - Compile Extended Data Dictionaries
  - Interactive Data-Wrangling
  - Leverage AI

# Datapusher+

Embedded use case

- Next-gen Data Ingestion extension for CKAN
- Guaranteed Data Type inferences
- Data Validation
  - Dedupe
  - PII screening
  - As context for AI - “describeGPT”
  - Extended Data Dictionary
  - Pre-calculate metadata (spatial extent, date range for time-series data, etc.)
  - Pre-populate DCAT 3 recommended metadata fields

<https://ckan.org/events/ckan-datapusher-plus-automagical-metadata>

*Standards-based, best-of-breed, open source solutions to make your Data Useful, Usable & Used*

Data that is  
Useful,  
Usable &  
Used



*We have a solution for this with DP+ & qsv*

*But what about actually **Using the Data**  
to gain **Actionable Insight**,  
to drive **Evidence-based Decisions**?*

# qsv pro

Cross-Platform Desktop  
Data-Wrangling & Query tool  
for the Rest of Us

<https://qsvpro.dathere.com>

- OpenRefine + Excel + qsv + CKAN + recipes + High Value Curated Data = qsv pro
- Familiar spreadsheet interface
- No need to know complex Command Line Interface (CLI) commands
- FAST! Blazing Fast!
- Interactive Data Wrangling
- Recipes! (desktop ETL)
- Integration with datHere's upcoming cloud-based services
  - High Value Data Feeds
  - Data Enrichment
  - Data Normalization
  - Geocoding
- Natural Language Interface

## Cross-platform Desktop Data Wrangling & Query tool for the Rest of Us

- For a Data Analyst Audience
- You don't need to be a Developer
- Use ready-made Recipes for common tasks (e.g. Scan for PII, geocode, deduplicate records, etc.)
- Create/modify/combine Recipes using either Lua or Python
- Share your Recipes on the datHere Recipe Catalog
- Pre-process security-sensitive data on your desktop without uploading it first
- Enrich your data with datHere's ever-expanding corpus of High Value Data like the Census, Bureau of Labor Statistics, etc.
- Use the "Answering People Interface" on your data or of other CKAN portals
- Upload to your CKAN or to datHere's Data Catalog to share your data with the world!

qsv pro (preview)

qsvpro v0.3.1 (preview) Workflow Configurator

# Workflow

Drag and drop a file to start.

## Recipes

Reusable scripts to modify your CSV file.

Type a command or search...

All Recipes

Sort in lexicographical order

Remove duplicate rows

Remove rows with Personally Identifiable Information (PII)

0 recipe(s) applied.

## Action Logs

History of actions performed based on your CSV file.

[7:43:31 AM] Analysis completed in: 59C  
[7:43:31 AM] [Analysis] Ran qsv sniff for  
[7:43:31 AM] [Analysis] Ran qsv sortche  
[7:43:31 AM] [Analysis] Computed frequ  
[7:43:31 AM] [Analysis] Computed advai  
[7:43:31 AM] [Analysis] Computed basic  
[7:43:31 AM] [Analysis] Indexed nyc311-  
[7:43:31 AM] Rendered table for nyc311-  
[7:43:31 AM] [Pre-Analysis] Computed c  
[7:43:31 AM] [Pre-Analysis] Computed r  
[7:43:31 AM] Screening completed in: 11  
[7:43:31 AM] [Screening] No empty ba

## Data Table

Preview your file as you transform it.

Choose File

Accepts: csv, tsv, tab, xlsx, xls, ods, xism, xisb

Analyzed 50k rows, compiling stats and frequency tables instantly!

Directly upload to any CKAN running v2.9 and above!

### nyc311-50k.csv

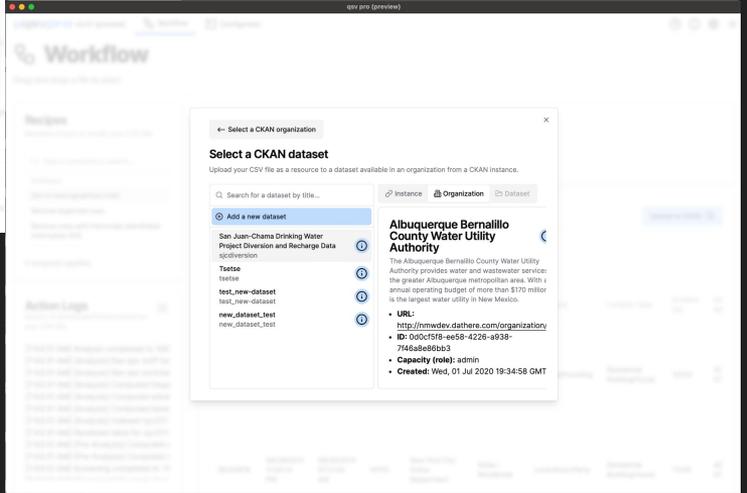
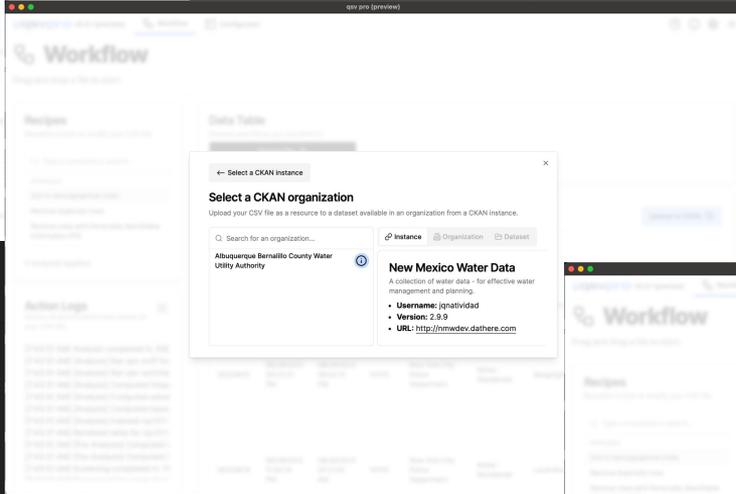
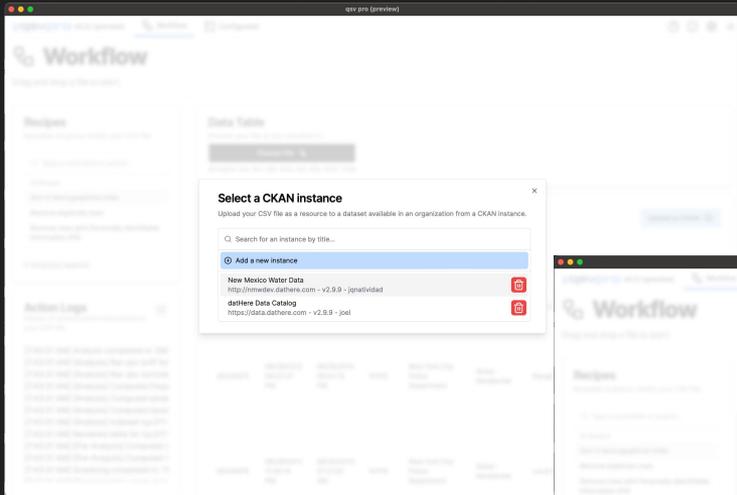
/Users/joelnatividad/Downloads/qsv\_test/nyc311-50k.csv

Upload to CKAN

Table Stats Frequency Metadata

Rows per page 10 Page 1 of 5000

Unique Key	Created Date	Closed Date	Agency	Agency Name	Complaint Type	Descriptor	Location Type	Incident Zip	Inc Ad
26220675	08/29/2013 09:23:37 PM	08/29/2013 09:42:15 PM	NYPD	New York City Police Department	Noise - Residential	Banging/Pounding	Residential Building/House	10024	10 ST
26220679	08/29/2013 11:05:14 PM	08/30/2013 07:21:02 AM	NYPD	New York City Police Department	Noise - Residential	Loud Music/Party	Residential Building/House	11226	40 ST



Successfully ran SQL query!  
Estimated elapsed time: 1139ms

Ran SQL query in 1139ms!

Upload to CKAN

Choose File

### Recipes

Reusable scripts to modify your CSV file.

Search for a recipe...

- All Recipes
- Sort in lexicographical order
- Remove duplicate rows
- Remove rows with Personally Identifiable Information (PII)

0 recipe(s) applied.

### Action Logs

History of actions performed based on your CSV file. Durations and timestamps are estimates.

- [12:30:18 PM] Completed SQL query run ir
- [12:27:35 PM] Completed SQL query run ir
- [12:25:00 PM] Estimated total processing
- [12:25:00 PM] Analysis completed in: 4717
- [12:25:00 PM] [Analysis] Ran qsv sniff for
- [12:24:59 PM] [Analysis] Ran qsv sortchec
- [12:24:58 PM] [Analysis] Computed limitec
- [12:24:48 PM] [Analysis] Computed freque
- [12:24:35 PM] [Analysis] Computed advan
- [12:24:17 PM] [Analysis] Computed basic
- [12:24:12 PM] Rendered table for NYC\_311
- [12:24:10 PM] [Des. Analysis] Computed

18 action(s) performed.

### NYC\_311\_SR\_2010-2020-sam

D:\Work\datHere\Projects\sample\NYC\_311\_SR\_2010-2020-sample-1M.csv

Table | Stats | Frequency | Metadata | SQL Query

### Run a SQL Query

Run a Polars SQL query on your CSV file. Refer to your CSV file as a table named `_t_1`.

Enter your natural language query. It should be converted to a SQL query below.

What are the most common complaint types by borough?

Ask

Enter your SQL query.

```
SELECT
  Borough,
  "Complaint Type",
  COUNT(*) AS Complaint_Count
FROM
  _t_1
GROUP BY
  Borough,
  "Complaint Type"
ORDER BY
  Borough,
  Complaint_Count DESC;
```

... an LLM we prompt to create a SQL query based on the Natural Language query & the context we provided

Recent query's estimated elapsed time: 1139ms

Run SQL query | Reset SQL query | Save output to file | Decrease code size | Increase code size

Rows per page: 10 | Page 1 of 129

Reproducible, hallucination-free answers

Borders | Wrap Rows

Borough	Complaint Type	Complaint_Count
BRONX	Noise - Residential	24284
BRONX	HEAT/HOT WATER	18584
BRONX	Street Light Condition	8354
BRONX	HEATING	7000

## DMS Applications



Open Data  
Portal



Internal  
Data Exchange



Enterprise  
Data Catalog



Regional  
Data Center



Data  
Library



Water  
Data Hubs

## datHere DMS Framework

datHere DMS Distribution

### High Value Data Sources



## Data Management System (DMS) Platform

Metadata



Analytics



Data Wrangling



Data Enrichment

Curated High Value  
Reference Data  
Geocoded Data  
Political Contexts  
Census Contexts



# DMS Framework

more than a Data Portal, a  
**Data Management System Framework**  
you can build on

- Built around CKAN
- Certified CKAN Extensions
- Bundled with other Best-of-Breed open source tooling
- Integrated Data Enrichment
- Build DMS applications like
  - Water Data Hubs
  - Open Data Portals
  - Internal Data Exchange
  - Data Library
  - Enterprise Data Catalog
  - and more...

# Pathways to Open Source Ecosystems (POSE)

- NSF initiative that *”aims to harness the power of open-source development for the creation of new technology solutions to problems of national and societal importance”*.
- In 2023, University of Pittsburgh and datHere conducted Phase 1 study on how to scale up Civic Data Ecosystem around CKAN and other open source **Data Infrastructure** initiatives.
- Summer 2024, we anticipate we’ll announce new initiatives to implement our Phase 1 scale up proposal
- More info at [civicdataecosystems.org](https://civicdataecosystems.org)



*Standards-based, best-of-breed, open source solutions to make your Data Useful, Usable & Used*

<https://datHere.com>

V1.0.0 (2024-05-29) - original slide deck presented at csv,conf,v8

V2.0.0 (2024-06-01) - fine-tuned/corrected

- Corrected number of implementations (50 to 100 - old number from ~2017)
- Replaced old stats screenshot with nicer looking version using latest qsv. Added elapsed time.
- Added “interactive data-wrangling” as reason for need for speed