

WEB CRAWLERS USING CARLIST.MY

Group B - Data Drillers

THEVAN RAJU A/L JEGANATH (A22EC0286)
MULYANI BINTI SARIPUDDIN (A22EC0223)
ALIATUL IZZAH BINTI JASMAN (A22EC0136)
MUHAMMAD ANAS BIN MOHD PIKRI (A21SC0464)

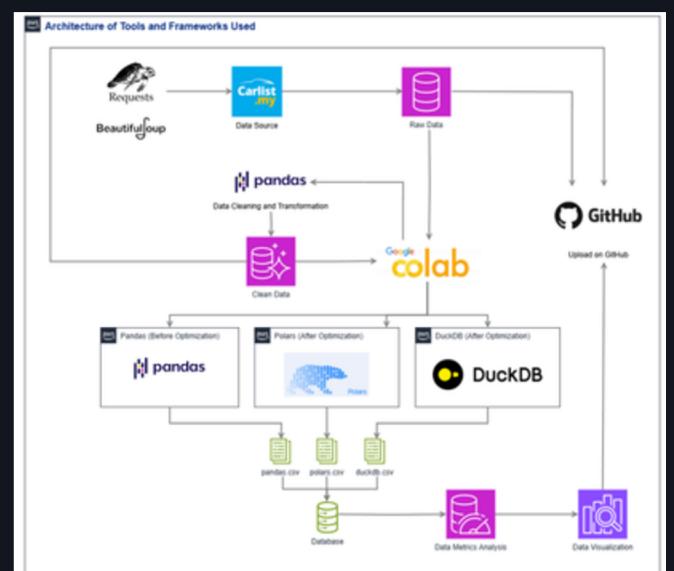
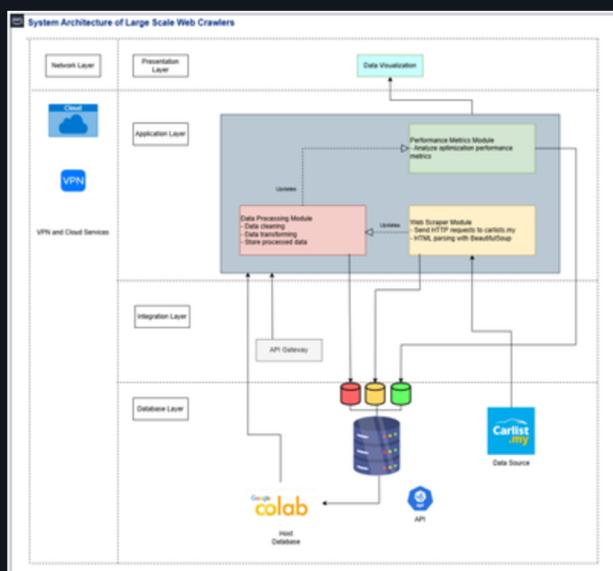


• • •

SYSTEM DESIGN AND ARCHITECTURE

SYSTEM ARCHITECTURE

TOOLS AND FRAMEWORKS



TOOLS AND FRAMEWORKS USED ...



Beautifulsoup



DuckDB

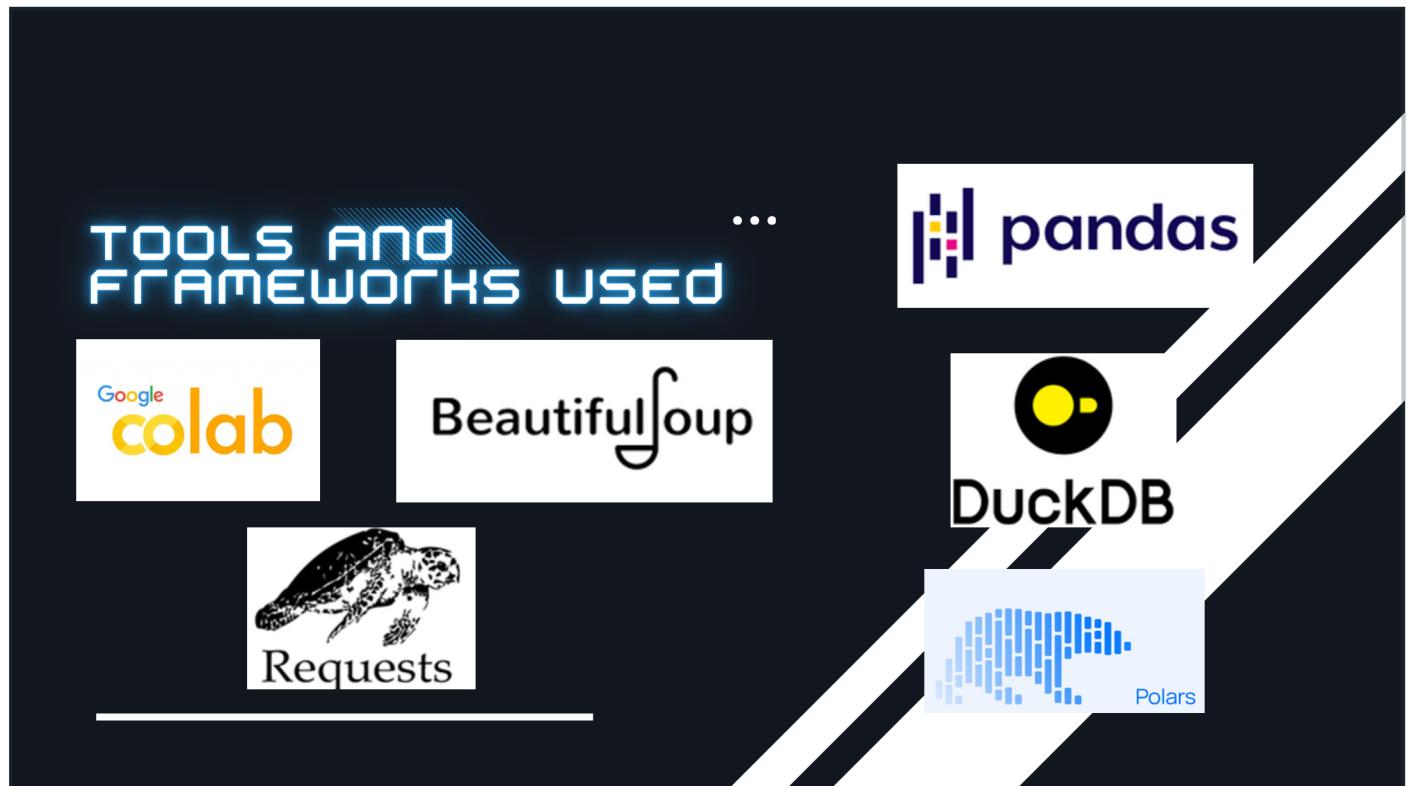


Table 2: Attributes collected from web scraping	
Data	Description
car_name	The car's title or name on the listing is defensible
price	Asking price for the car
location	City/ locality where the car is sold
region	Malaysian state or territory
brand	Car manufacturer
model	Car model type
year	Manufacturing year
mileage	Driven distance specifically in kilometres (km)
fuel_type	Fuel type, such as petrol or diesel
color	The body colour of the car
body_type	Car type, such as hatchback, sedan, etc
seating_capacity	The seating number in the car

DATA COLLECTION & ENTITIES

condition	Vehicle condition, such as used, new, etc
image	Vehicle image link to identify the car
description	Seller's description of the car
url	Full list link.

RAW DATA

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R																							
1	Car Name	Price (MYR)	Currency	Location	Region	Brand	Model	Year	Mileage	Fuel Type	Color	Body Type	Seating Capacity	Condition	Image	Description	URL																								
2	Toyota Vellfire	420000	MYR	Subang Jaya	Selangor	Toyota	Vellfire	2023	2500	Petrol - Ur Black	MPV		7	RefurbishedCc	https://im_2023_TOYC	<a href="https://www.carlist.i</td></tr> <tr> <td>3</td><td>Mazda CX-5</td><td>36900</td><td>MYR</td><td>Seri Kembangan</td><td>Selangor</td><td>Mazda</td><td>CX-5</td><td>2016</td><td>82500</td><td>Diesel - Grey</td><td>SUV</td><td></td><td>5</td><td>UsedCondition</td><td>https://im_2016_MAZI	<a href="https://www.carlist.i</td></tr> <tr> <td>4</td><td>MINI 3 Door</td><td>152000</td><td>MYR</td><td>TTDI</td><td>Kuala Lumpur</td><td>MINI</td><td>3 Door</td><td>2020</td><td>17500</td><td>Petrol - Ur Black</td><td>Hatchback</td><td></td><td>4</td><td>RefurbishedCc</td><td>https://im_0728Y2	<a href="https://www.carlist.i</td></tr> <tr> <td>5</td><td>Nissan X-Trail</td><td>59800</td><td>MYR</td><td>Johor Bahru</td><td>Johor</td><td>Nissan</td><td>X-Trail</td><td>2020</td><td>52500</td><td>Hybrid - White</td><td>SUV</td><td></td><td>5</td><td>UsedCondition</td><td>https://im_WhatsApp	<a href="https://www.carlist.i</td></tr> <tr> <td>6</td><td>Ford Focus</td><td>13900</td><td>MYR</td><td>Seri Kembangan</td><td>Selangor</td><td>Ford</td><td>Focus</td><td>2013</td><td>132500</td><td>Petrol - Ur White</td><td>Sedan</td><td></td><td>5</td><td>UsedCondition</td><td>https://im_2013_FOCE	<a href="https://www.carlist.i</td></tr> <tr> <td>7</td><td>Lexus RX300</td><td>261000</td><td>MYR</td><td>Subang Jaya</td><td>Selangor</td><td>Lexus</td><td>RX300</td><td>2021</td><td>17500</td><td>Petrol - Ur White</td><td>SUV</td><td></td><td>5</td><td>RefurbishedCc</td><td>https://im_2021_LEXU	<a href="https://www.carlist.i</td></tr> <tr> <td>8</td><td>Audi A8</td><td>46900</td><td>MYR</td><td>Seri Kembangan</td><td>Selangor</td><td>Audi</td><td>A8</td><td>2013</td><td>97500</td><td>Petrol - Ur Black</td><td>Sedan</td><td></td><td>4</td><td>UsedCondition</td><td>https://im_2013_AUDI	<a href="https://www.carlist.i</td></tr> <tr> <td>9</td><td>Land Rover Defender</td><td>419000</td><td>MYR</td><td>Subang Jaya</td><td>Selangor</td><td>Land Rove</td><td>Defender</td><td>2023</td><td>7500</td><td>Petrol - Ur Bronze</td><td>SUV</td><td></td><td>5</td><td>RefurbishedCc</td><td>https://im_2023_LANCI	<a href="https://www.carlist.i</td></tr> <tr> <td>10</td><td>Citroen Grand C4 SpaceTourer</td><td>60900</td><td>MYR</td><td>Seri Kembangan</td><td>Selangor</td><td>Citroen</td><td>Grand C4 S</td><td>2020</td><td>77500</td><td>Petrol - Ur White</td><td>MPV</td><td></td><td>7</td><td>UsedCondition</td><td>https://im_2020_CITR	<a href="https://www.carlist.i</td></tr> <tr> <td>11</td><td>Audi Q7</td><td>49900</td><td>MYR</td><td>Seri Kembangan</td><td>Selangor</td><td>Audi</td><td>Q7</td><td>2014</td><td>112500</td><td>Petrol - Ur White</td><td>SUV</td><td></td><td>7</td><td>UsedCondition</td><td>https://im_2014_AUDI	<a href="https://www.carlist.i</td></tr> <tr> <td>12</td><td>Porsche Cayenne</td><td>69900</td><td>MYR</td><td>Seri Kembangan</td><td>Selangor</td><td>Porsche</td><td>Cayenne</td><td>2011</td><td>127500</td><td>Diesel - Silver</td><td>SUV</td><td></td><td>5</td><td>UsedCondition</td><td>https://im_2011_PORS	<a href="https://www.carlist.i</td></tr> <tr> <td>13</td><td>Land Rover Defender</td><td>379000</td><td>MYR</td><td>Subang Jaya</td><td>Selangor</td><td>Land Rove</td><td>Defender</td><td>2022</td><td>17500</td><td>Petrol - Ur Blue</td><td>SUV</td><td></td><td>7</td><td>RefurbishedCc</td><td>https://im_LAND_ROV	<a href="https://www.carlist.i</td></tr> <tr> <td>14</td><td>Mazda Biante</td><td>60900</td><td>MYR</td><td>Seri Kembangan</td><td>Selangor</td><td>Mazda</td><td>Biante</td><td>2017</td><td>62500</td><td>Petrol - Ur Silver</td><td>MPV</td><td></td><td>8</td><td>UsedCondition</td><td>https://im_2017_MA2I	<a href="https://www.carlist.i</td></tr> <tr> <td>15</td><td>Mercedes-Benz C200</td><td>188000</td><td>MYR</td><td>Ampang</td><td>Selangor</td><td>Mercedes-</td><td>C200</td><td>2020</td><td>12500</td><td>Petrol - Ur Black</td><td>Sedan</td><td></td><td>5</td><td>RefurbishedCc</td><td>https://im_MODEL_M	<a href="https://www.carlist.i</td></tr> <tr> <td>16</td><td>Toyota Harrier</td><td>203000</td><td>MYR</td><td>Old Klang Road</td><td>Kuala Lumpur</td><td>Toyota</td><td>Harrier</td><td>2023</td><td>12500</td><td>Petrol - Ur White</td><td>SUV</td><td></td><td>5</td><td>RefurbishedCc</td><td>https://im_2023_Toyo	<a href="https://www.carlist.i</td></tr> <tr> <td>17</td><td>Porsche Panamera</td><td>518000</td><td>MYR</td><td>Old Klang Road</td><td>Kuala Lumpur</td><td>Porsche</td><td>Panamera</td><td>2019</td><td>32500</td><td>Petrol - Ur Silver</td><td>Hatchback</td><td></td><td>4</td><td>RefurbishedCc</td><td>https://im_2019_Pors	<a href="https://www.carlist.i</td></tr> <tr> <td>18</td><td>Toyota Alphard</td><td>351000</td><td>MYR</td><td>Johor Bahru</td><td>Johor</td><td>Toyota</td><td>Alphard</td><td>2024</td><td>2500</td><td>Petrol - Ur White</td><td>MPV</td><td></td><td>6</td><td>RefurbishedCc</td><td>https://im_GENUIN	<a href="https://www.carlist.i</td></tr> <tr> <td>19</td><td>Lexus NX250</td><td>283000</td><td>MYR</td><td>Old Klang Road</td><td>Kuala Lumpur</td><td>Lexus</td><td>NX250</td><td>2023</td><td>12500</td><td>Petrol - Ur White</td><td>SUV</td><td></td><td>5</td><td>RefurbishedCc</td><td>https://im_2023_Lexu	<a href="https://www.carlist.i</td></tr> <tr> <td>20</td><td>Mercedes-Benz G350</td><td>578000</td><td>MYR</td><td>Old Klang Road</td><td>Kuala Lumpur</td><td>Mercedes-</td><td>G350</td><td>2021</td><td>37500</td><td>Diesel - Black</td><td>SUV</td><td></td><td>5</td><td>RefurbishedCc</td><td>https://im_2021_Merc	<a href="https://www.carlist.i</td></tr> <tr> <td>21</td><td>Lexus RX300</td><td>261000</td><td>MYR</td><td>Subang Jaya</td><td>Selangor</td><td>Lexus</td><td>RX300</td><td>2021</td><td>27500</td><td>Petrol - Ur Black</td><td>SUV</td><td></td><td>5</td><td>RefurbishedCc</td><td>https://im_2021_RXE	<a href="https://www.carlist.i</td></tr> <tr> <td>22</td><td>Toyota Alphard</td><td>323000</td><td>MYR</td><td>Old Klang Road</td><td>Kuala Lumpur</td><td>Toyota</td><td>Alphard</td><td>2023</td><td>27500</td><td>Petrol - Ur White</td><td>MPV</td><td></td><td>7</td><td>RefurbishedCc</td><td>https://im_2023_Toyo	<a href="https://www.carlist.i</td></tr> <tr> <td>23</td><td>MINI Clubman</td><td>213000</td><td>MYR</td><td>Old Klang Road</td><td>Kuala Lumpur</td><td>MINI</td><td>Clubman</td><td>2022</td><td>17500</td><td>Petrol - Ur White</td><td>Wagon</td><td></td><td>5</td><td>RefurbishedCc</td><td>https://im_2022_MINI	<a href="https://www.carlist.i</td></tr> <tr> <td>24</td><td>Toyota Vellfire</td><td>395000</td><td>MYR</td><td>Old Klang Road</td><td>Kuala Lumpur</td><td>Toyota</td><td>Vellfire</td><td>2023</td><td>22500</td><td>Petrol - Ur White</td><td>MPV</td><td></td><td>7</td><td>RefurbishedCc</td><td>https://im_2023_Toyo	<a href="https://www.carlist.i</td></tr> <tr> <td>25</td><td>Lexus IS300</td><td>237000</td><td>MYR</td><td>Old Klang Road</td><td>Kuala Lumpur</td><td>Lexus</td><td>IS300</td><td>2021</td><td>22500</td><td>Petrol - Ur White</td><td>Sedan</td><td></td><td>5</td><td>RefurbishedCc</td><td>https://im_2021_Lexu	<a href="https://www.carlist.i</td></tr> <tr> <td>26</td><td>Volvo XC60</td><td>128900</td><td>MYR</td><td>Seri Kembangan</td><td>Selangor</td><td>Volvo</td><td>XC60</td><td>2020</td><td>102500</td><td>Petrol - Ur White</td><td>SUV</td><td></td><td>5</td><td>UsedCondition</td><td>https://im_2020_VOLV	<a href="https://www.carlist.i</td></tr> <tr> <td>27</td><td>Mitsubishi Triton</td><td>19800</td><td>MYR</td><td>Kajang</td><td>Selangor</td><td>Mitsubishi</td><td>Triton</td><td>2010</td><td>152500</td><td>Diesel - White</td><td>Pickup Truck</td><td></td><td>5</td><td>UsedCondition</td><td>https://im_WELCON	<a 154="" 56="" 953="" 970"="" data-label="Page-Footer" href="https://www.carlist.i</td></tr> </tbody> </table> </div> <div data-bbox="> <p>HPDP Slide</p>

DATA CLEANING



CLEANED DATA

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R
1	car name	price (myr)	location	region	brand	model	year	mileage	fuel type	color	body type	seating capacity	condition	image	descripitior url			
2	toyota vellfire	420000	subang jaya	selangor	toyota	vellfire	2023	2500	petrol - un	black	mpv	7	refurbishe: https://im12023toyo	https://www.carlist.my/cars-for				
3	mazda cx-5	36900	seri kembangan	selangor	mazda	cx-5	2016	82500	diesel	grey	suv	5	used	https://im12016mazc	https://www.carlist.my/cars-for			
4	mini 3 door	152000	ttdi	kuala lumpur	mini	3 door	2020	17500	petrol - un	black	hatchback	4	refurbishe: https://im10720Y2	https://www.carlist.my/cars-for				
5	nissan x-trail	59800	johor bahru	johor	nissan	x-trail	2020	52500	hybrid	white	suv	5	used	https://im12020whatsapp	https://www.carlist.my/cars-for			
6	ford focus	13900	seri kembangan	selangor	ford	focus	2013	132500	petrol - un	white	sedan	5	used	https://im12013ford	https://www.carlist.my/cars-for			
7	lexus rx300	261000	subang jaya	selangor	lexus	rx300	2021	17500	petrol - un	white	suv	5	refurbishe: https://im12021lexus	https://www.carlist.my/cars-for				
8	audi a8	46900	seri kembangan	selangor	audi	a8	2013	97500	petrol - un	black	sedan	4	used	https://im12013audi	https://www.carlist.my/cars-for			
9	land rover defender	419000	subang jaya	selangor	land rover	defender	2023	7500	petrol - un	bronze	suv	5	refurbishe: https://im12023land	https://www.carlist.my/cars-for				
10	citroen grand c4 spacetourer	60900	seri kembangan	selangor	citroen	grand c4 sp	2020	77500	petrol - un	white	mpv	7	used	https://im12020citro	https://www.carlist.my/cars-for			
11	audi q7	49900	seri kembangan	selangor	audi	q7	2014	112500	petrol - un	white	suv	7	used	https://im12014audi	https://www.carlist.my/cars-for			
12	porsche cayenne	69900	seri kembangan	selangor	porsche	cayenne	2011	127500	diesel	silver	suv	5	used	https://im12011orsi	https://www.carlist.my/cars-for			
13	land rover defender	379000	subang jaya	selangor	land rover	defender	2022	17500	petrol - un	blue	suv	7	refurbishe: https://im1land	https://www.carlist.my/cars-for				
14	mazda biante	60900	seri kembangan	selangor	mazda	biante	2017	62500	petrol - un	silver	mpv	8	used	https://im12017mazc	https://www.carlist.my/cars-for			
15	mercedes-benz c200	188000	ampang	selangor	mercedes-	c200	2020	12500	petrol - un	black	sedan	5	refurbishe: https://im1model	https://www.carlist.my/cars-for				
16	toyota harrier	203000	old klang road	kuala lumpur	toyota	harrier	2023	12500	petrol - un	white	suv	5	refurbishe: https://im12023toyo	https://www.carlist.my/cars-for				
17	porsche panamera	518000	old klang road	kuala lumpur	porsche	panamera	2019	32500	petrol - un	silver	hatchback	4	refurbishe: https://im12019pors	https://www.carlist.my/cars-for				
18	toyota alphard	351000	johor bahru	johor	toyota	alphard	2024	2500	petrol - un	white	mpv	6	refurbishe: https://im1genuine	https://www.carlist.my/cars-for				
19	lexus rx250	283000	old klang road	kuala lumpur	lexus	rx250	2023	12500	petrol - un	white	suv	5	refurbishe: https://im12023lexus	https://www.carlist.my/cars-for				
20	mercedes-benz g350	578000	old klang road	kuala lumpur	mercedes-	g350	2021	37500	diesel	black	suv	5	refurbishe: https://im12021merc	https://www.carlist.my/cars-for				
21	lexus rx300	261000	subang jaya	selangor	lexus	rx300	2021	27500	petrol - un	black	suv	5	refurbishe: https://im12021lexus	https://www.carlist.my/cars-for				
22	toyota alphard	323000	old klang road	kuala lumpur	toyota	alphard	2023	27500	petrol - un	white	mpv	7	refurbishe: https://im12023toyo	https://www.carlist.my/cars-for				
23	mini clubman	213000	old klang road	kuala lumpur	mini	clubman	2022	17500	petrol - un	white	wagon	5	refurbishe: https://im12022mini	https://www.carlist.my/cars-for				
24	toyota vellfire	395000	old klang road	kuala lumpur	toyota	vellfire	2023	22500	petrol - un	white	mpv	7	refurbishe: https://im12023toyo	https://www.carlist.my/cars-for				
25	lexus is300	237000	old klang road	kuala lumpur	lexus	is300	2021	22500	petrol - un	white	sedan	5	refurbishe: https://im12021lexus	https://www.carlist.my/cars-for				
26	volvo xc60	128900	seri kembangan	selangor	volvo	xc60	2020	102500	petrol - un	white	suv	5	used	https://im12020volvc	https://www.carlist.my/cars-for			
27	mitsubishi triton	19800	kajang	selangor	mitsubishi	triton	2010	152500	diesel	white	pickup truck	5	used	https://im1welcom	https://www.carlist.my/cars-for			
28	isuzu d-max	19800	kajang	selangor	isuzu	d-max	2011	142500	diesel	black	pickup truck	5	used	https://im1welcom	https://www.carlist.my/cars-for			

DATA OPTIMIZATION

POLARS

POLARS IS AN OPEN-SOURCE LIBRARY FOR DATA MANIPULATION, KNOWN FOR BEING ONE OF THE FASTEST DATA PROCESSING SOLUTIONS ON A SINGLE MACHINE. IT FEATURES A WELL-STRUCTURED, TYPED API THAT IS BOTH EXPRESSIVE AND EASY TO USE.



PANDAS
PANDAS IS A POWERFUL, OPEN-SOURCE DATA ANALYSIS AND MANIPULATION LIBRARY FOR PYTHON. IT PROVIDES EASY-TO-USE DATA STRUCTURES AND TOOLS FOR WORKING WITH STRUCTURED DATA.

DUCKDB

DUCKDB SERVES AS A POWERFUL ENGINE FOR ANALYZING AND TRANSFORMING DATA RESIDING IN VARIOUS FORMATS LIKE CSV, PARQUET, JSON, OR DATABASES SUCH AS POSTGRESQL AND SQLITE. IT CAN BE USED TRANSIENTLY FOR DATA MANIPULATION OR TO CREATE PERSISTENT TABLES FOR ANALYTICAL QUERIES.



...

COMPARISON

TOTAL PROCESSING TIME (SECONDS)

Optimization Stage	Run 1	Run 2	Run 3	Average
Pandas Optimization	5.35	3.49	3.44	4.09
DuckDB Optimization	0.45	0.89	0.43	0.59
Polars Optimization	0.43	0.26	0.61	0.43

CPU USAGE (%)

Optimization Stage	Run 1	Run 2	Run 3	Average
Pandas Optimization	14.67	16.67	3.6	11.65
DuckDB Optimization	19.8	50	46.2	38.67
Polars Optimization	83.1	36.2	14.8	44.7

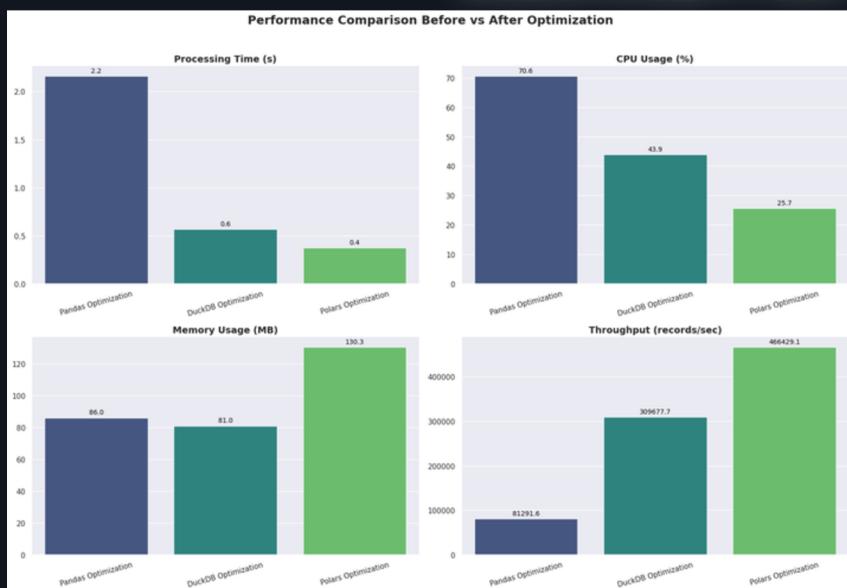
MEMORY USAGE (MB)

Optimization Stage	Run 1	Run 2	Run 3	Average
Pandas Optimization	102.34	119.98	125.8	116.04
DuckDB Optimization	65.79	98.83	81.64	82.09
Polars Optimization	28.63	3.81	129.63	54.02

THROUGHPUT

Optimization Stage	Run 1	Run 2	Run 3	Average
Pandas Optimization	32827.06	50356.60	51003.49	44729.0
DuckDB Optimization	390337.6	196502.90	406446.43	331095.6
Polars Optimization	412467.7	676460.87	289688.39	459539.0

VISUALIZATION



• • •

CHALLENGES AND LIMITATIONS



LONG DATA SCRAPING
DURATION DURING
EARLY PHASE

FINDING SUITABLE
OPTIMIZATION
LIBRARY

SUDDEN CHANGES AND
RESTRICTIONS BY
CARLISTS.MY

...

THANK YOU

...