



Nobel Prize Data Warehouse



Introduction

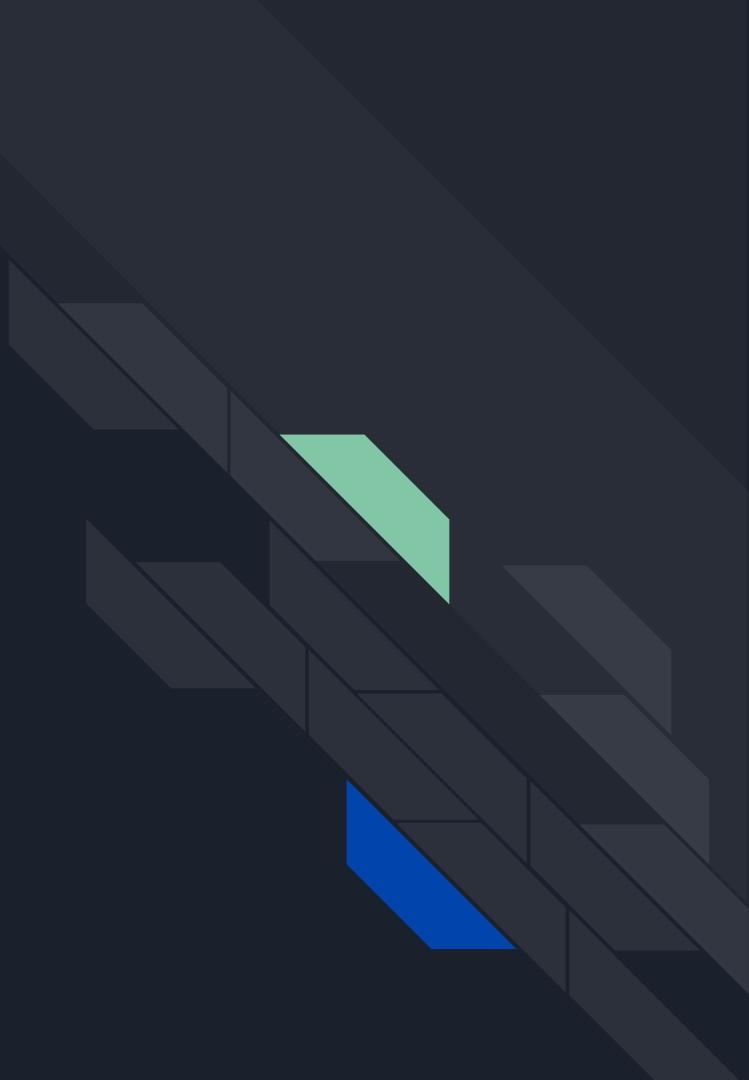
For the final project, our group chose to create a data warehouse to determine if a person's origin country correlates to their likelihood of winning a Nobel Prize. We would like to use our data warehouse to determine if the development of a country affects the amount of people born there who have won Nobel Prizes. For example, is it appropriate to theorize that a person from a first world country is more likely to win a Nobel Prize, compared to someone from a developing country?



Data Source

The data that we used comes from the [Harvard Dataverse](#), which has published a dataset that contains three files; information about Nobel laureates, the prize they have won, and a list of the countries. We use this data to create dimensions such as laureate, overall motivation, prize category and institution.

Preparation and Planning



BEAM Matrix

EVENT	Importance	Estimate	LAUREATES	INSTITUTION	PRIZE	COUNTRY	OVERALL MOTIVATION	INSTITUTION TYPE	PRIZE CATEGORY	PRIZE MOTIVATION	stakeholder group		
				100	100	30	80	20	20				
				25	10	10	10	5	5				
LAUREATE SPLIT PRIZE MONEY Laureate splits Prize Money	1	5		✓	✓	✓			✓				
INSTITUTION MEMBER ENROLLMENT AND VALUE	2	3			✓		✓			✓			
SPLIT PRIZE MONEY	3	3			✓					✓			
Event Count			75	2	2	0	0	0	2	0	0	0	0
				0	0	0	0	0	0	0	0	0	0

The BEAM event matrix which shows the various dimensions and events.

BEAM Dimensions

LAUREATES HV														
Laureate ID	First Name	Last Name	Date of Birth	Date of Death	Birth Country	Birth City	Gender	Year Prize Received	Prize Category	Prize Amount	Motivation	Share	Institution Name	Institution Country
BK, NN			FV	FV	FV	FV					NN	NN		
I	(LAUREATE.CSV)	(LAUREATE.CSV)	(LAUREATE.CSV)	(LAUREATE.CSV)	(LAUREATE.CSV)	(LAUREATE.CSV)	(LAUREATE.CSV)	(LAUREATE.CSV)	(LAUREATE.CSV)	(LAUREATE.CSV)	(LAUREATE.CSV)	(LAUREATE.CSV)	(LAUREATE.CSV)	
C100	C100	C100	D	D	C100	C100	C10	FV, NN	FV, NN	NN	NN	NN		
							D	(LAUREATE.CSV)	(LAUREATE.CSV)	(LAUREATE.CSV)	(LAUREATE.CSV)	(LAUREATE.CSV)	(LAUREATE.CSV)	(LAUREATE.CSV)
1	Wilhelm Conrad	Rontgen	1845-03-27	2/10/1923	Prussia (now Germany)	Lennep (now Remscheid)	male	1901	physics	1000000	"in recognition"	1	Munich University	Germany
2	Hendrik Antoon	Lorentz	1853-07-18	2/4/1928	the Netherlands	Arnhem	male	1902	physics	500000	"for discovery"	2	Leiden University	the Netherlands
3	Pieter	Zeeman	1865-05-25	10/9/1943	the Netherlands	Zonnemaire	male	1902	physics	500000	"for recognition"	2	Amsterdam University	the Netherlands
4	Antoine Henri	Becquerel	1852-12-15	8/25/1908	France	Paris	male	1903	physics	500000	"in recognition"	2	Polytechnique	France
5	Pierre	Curie	1859-05-15	4/19/1906	France	Paris	male	1903	physics	250000	"for discovery"	4	(Municipal School of Industrial Physics and Chemistry)	France
6	Marie	Curie Skłodowska	1867-11-07	7/4/1934	Russian Empire (now Poland)	Warsaw	female	1903	physics	250000	"for recognition"	4	Frankfurt-on-the-Main University	Germany
6	Marie	Curie, Skłodowska	1867-11-07	7/4/1934	Russian Empire (now Poland)	Warsaw	female	1911	chemistry	1000000	"in recognition"	1	Sorbonne University	France
8	Lord Rayleigh	John William Strutt	1842-11-12	6/30/1919	United Kingdom	Langford Grove, Maldon, Essex	male	1904	physics	1000000	"for discovery"	1	Royal Institution of Great Britain	United Kingdom
9	United Nations	NULL	0000-00-00	0000-00-00	NULL	NULL	NULL	1905	peace	1000000	"for recognition"	1	NULL	NULL

A screenshot of the Laureates dimension from the BEAM Modelstormer.

BEAM Dimensions

INSTITUTION				
Institution Key	Institution ID	Institution Name	Institution City	Institution Country
SK, NN	BK, NN	NN, ND		
I	C6	{Institution.CSV}	{Institution.CSV}	{Institution.CSV}
1	IN001	Kiel University	Kiel	Germany
2	IN002	University of Cambridge	Cambridge	United Kingdom
3	IN003	University of Chicago	Chicago, IL	USA
4	IN004	Sorbonne University	Paris	France
5	IN005	Marconi Wireless Telegraph Co. Ltd.	London	United Kingdom
6	IN006	Strasbourg University	Strasbourg	Alsace (then Germany, now France)
7	IN007	Amsterdam University	Amsterdam	the Netherlands
8	IN008	Harvard	Cambridge, MA	United States
9	IN009	Swedish Gas-Accumulator Co.	Stockholm	Sweden

A screenshot of the Institution dimension from the BEAM Modelstormer.

BEAM Dimensions

INSTITUTION TYPE			
Institution Type Key	Institution Type ID	Institution Type	Institution Type Description
SK, NN	BK, NN	{InstitutionType.CSV}	{InstitutionType.CSV}
I		C30	C200
1	INT001	Educational	involves schools and universities either being private or public
2	INT002	Religious	involves religious groups i.e. Christianity, Judaism, Islam, Hinduism, Buddhism, and Sikhism
3	INT003	Medicinal	involves hospitals and other health care institutions either being private or public
4	INT004	Government	involves governments that are either monarchy, oligarchy, dictatorship, and democracy
5	INT005	Economic	involves foundations focused on collecting economic data with the idea of providing a good or service that is important to the economy
6	INT006	Research	involves establishments that specialize in advancing scholarly activity through research and experimental development

A screenshot of the Institution Type dimension from the BEAM Modelstormer.

BEAM Dimensions

COUNTRY			
Country Key	Country ID	Country Name	Country Abbreviation
SK, NN	BK, NN {Country.CSV}	NN {Country.CSV}	{Country.CSV}
I	C5	C100	C3
1	C100	Alsace, then Germany	DE
2	C101	Alsace	DE
3	C102	Germany	DE
4	C103	Argentina	AR
5	C104	Australia	AU
6	C105	Austria	AT
7	C106	Belgium	BE
8	C107	Burma	MM
9	C108	Old Republic	NULL

A screenshot of the Country dimension from the BEAM Modelstormer.

BEAM Dimensions

PRIZE			
Prize Key	Prize ID	Prize Year	Prize Category
SK, NN	BK, NN {Prize.CSV}	NN {Prize.CSV}	NN {Prize.CSV}
I	C5	C4	C20
1	P001	2016	physics
2	P002	2016	physics
3	P003	2016	physics
4	P004	2016	chemistry
5	P005	2016	chemistry
6	P006	2016	chemistry
7	P007	2016	medicine
8	P008	2016	literature
9	P009	2016	peace

PRIZE CATEGORY			
Category Key	Category ID	Category	
SK, NN	BK, NN {Category.CSV}	ND {Category.CSV}	
I	C5	C20	
1	NPC1	physics	
2	NPC2	chemistry	
3	NPC3	medicine	
4	NPC4	peace	
5	NPC5	literature	
6	NPC6	economics	

A screenshot of the Prize and Prize Category from the BEAM Modelstormer.

BEAM Dimensions

MOTIVATION

Motivation Key	Motivation ID	Motivation
SK, NN I	BK, NN (Motivation.CSV) C5	NN (Motivation.CSV) C300
1	M001	in recognition of his work in thermochemistry
2	M002	for their discoveries concerning the molecular structure of nucleic acids and its significance for information transfer in living material
3	M003	for their discoveries concerning the structural and functional organization of the cell
4	M004	for the discovery of the quantized Hall effect
5	M005	for their pathbreaking contribution to the theory of international trade and international capital movements
6	M006	for having created new poetic expressions within the great American song tradition
7	M007	for his discovery of the organizer effect in embryonic development
8	M008	for his development of nuclear magnetic resonance spectroscopy for determining the three-dimensional structure of biological macromolecules in solution
9	M009	for their discoveries relating to the hormones of the adrenal cortex, their structure and biological effects

OVERALL MOTIVATION

Overall Motivation Key	Overall Motivation ID	Overall Motivation
SK, NN I	BK, NN (OverallMotivation.CSV) C5	NN (OverallMotivation.CSV) C100
1	OM001	for pioneering experimental contributions to lepton physics
2	OM002	for contributions to the developments of methods within DNA-based chemistry
3	OM003	for basic work on information and communication technology
4	OM004	for the development of methods for identification and structure analyses of biological macromolecules
5	OM005	for discoveries concerning channels in cell membranes
6	OM006	for pioneering contributions to the development of neutron scattering techniques for studies of condensed matter

BEAM Facts

LAUREATE PRIZE							
LAUREATE KEY	originates COUNTRY KEY	belongs INSTITUTION KEY	receives PRIZE KEY	in CATEGORY KEY	SHARE	earned PRIZE AWARD	
[who]	[where]	[what]	[what]	[what]	[how]	[\$]	
53	87	88	124	2	1	1000000	
88	81	86	11	3	2	500000	
30	17	47	7	2	2	500000	
5	70	89	13	1	2	500000	
59	21	75	134	2	4	250000	
26	47	14	37	1	4	250000	
84	22	80	76	2	1	1000000	
91	87	4	188	6	1	1000000	
73	37	77	196	6	1	1000000	

A screenshot of the Laureate Shared Prize fact from the BEAM Modelstormer.

BEAM Facts

INSTITUTION VALUE & MEMBERS			has Institution Value in Millions	with Institution Member Count
Institution Key	Institution Type Key	Country Key	[\$]	
[who/what]	[what]	[where]	[\$]	
88	1	87	16.6	424
86	2	81	723.4	7744
47	3	17	40.3	399333
89	4	70	51.7	234
75	5	21	3.2	32
14	6	47	1.9	6533
80	1	22	500	66250
4	2	87	221.2	29309
77	3	37	3.3	437

A screenshot of the Institution Value and Members fact from the BEAM Modelstormer.

BEAM Facts

INSTITUTION VALUE & MEMBERS			has Institution Value in Millions	with Institution Member Count
Institution Key	Institution Type Key	Country Key	[\$]	
[who/what]	[what]	[where]	[\$]	
88	1	87	16.6	424
86	2	81	723.4	7744
47	3	17	40.3	399333
89	4	70	51.7	234
75	5	21	3.2	32
14	6	47	1.9	6533
80	1	22	500	66250
4	2	87	221.2	29309
77	3	37	3.3	437

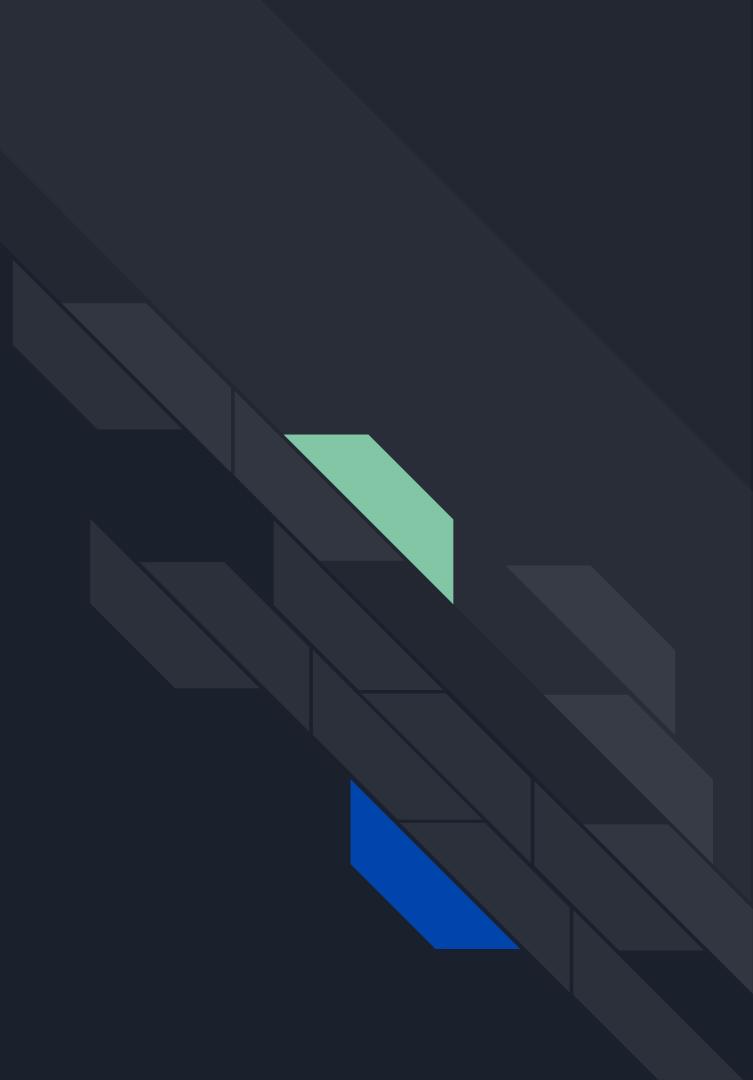
A screenshot of the Institution Value and Members fact from the BEAM Modelstormer.

BEAM Facts

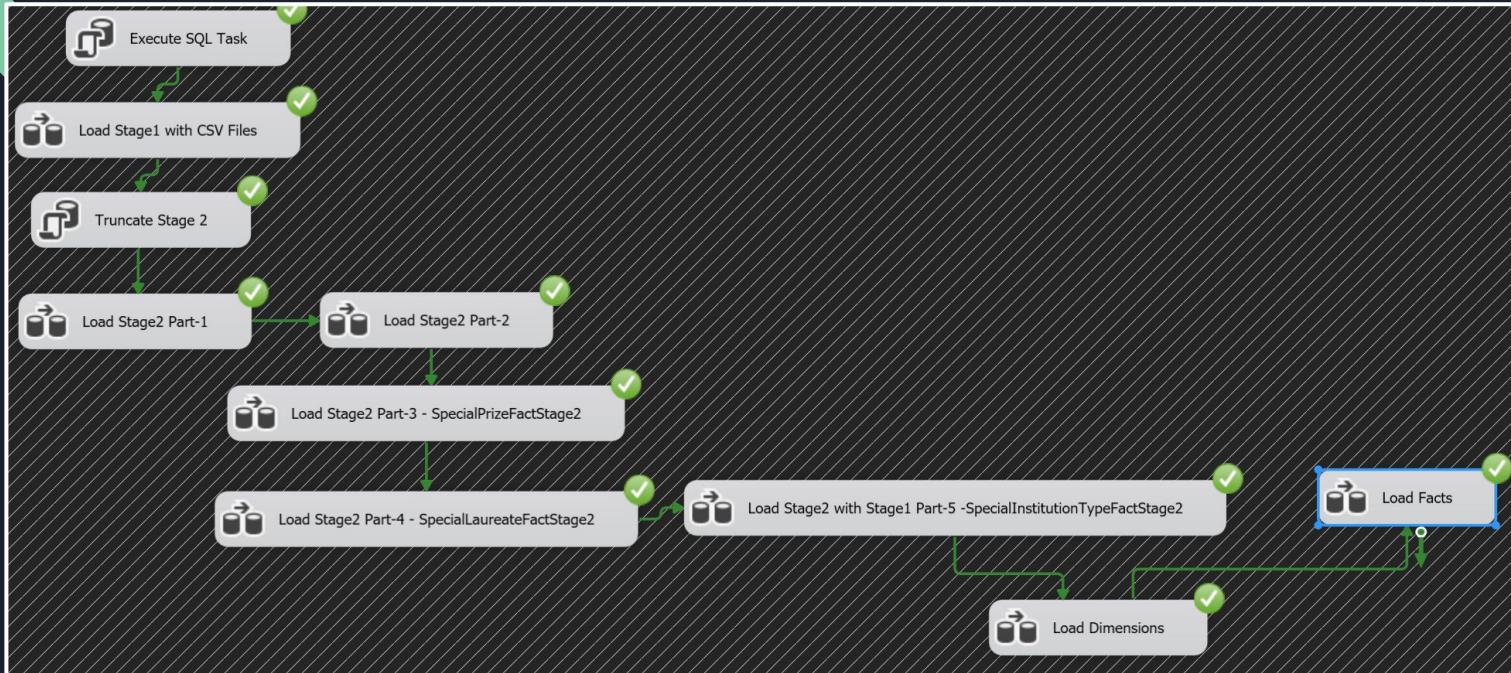
DISTRIBUTED MONEY					
Laureate Key	receives Prize Key	in Category Key	has to Share	results Prize Distribution per Winner	
[who]	[what]		[how]	[\$]	
53	124	2	1	500000	
88	11	3	2	250000	
30	7	2	2	250000	
5	13	1	2	333333.3333	
59	134	2	4	333333.3333	
26	37	1	4	333333.3333	
84	76	2	1	1000000	
91	188	6	1	1000000	
73	196	6	1	1000000	

A screenshot of the Distributed Money fact from the BEAM Modelstormer.

ETL Process

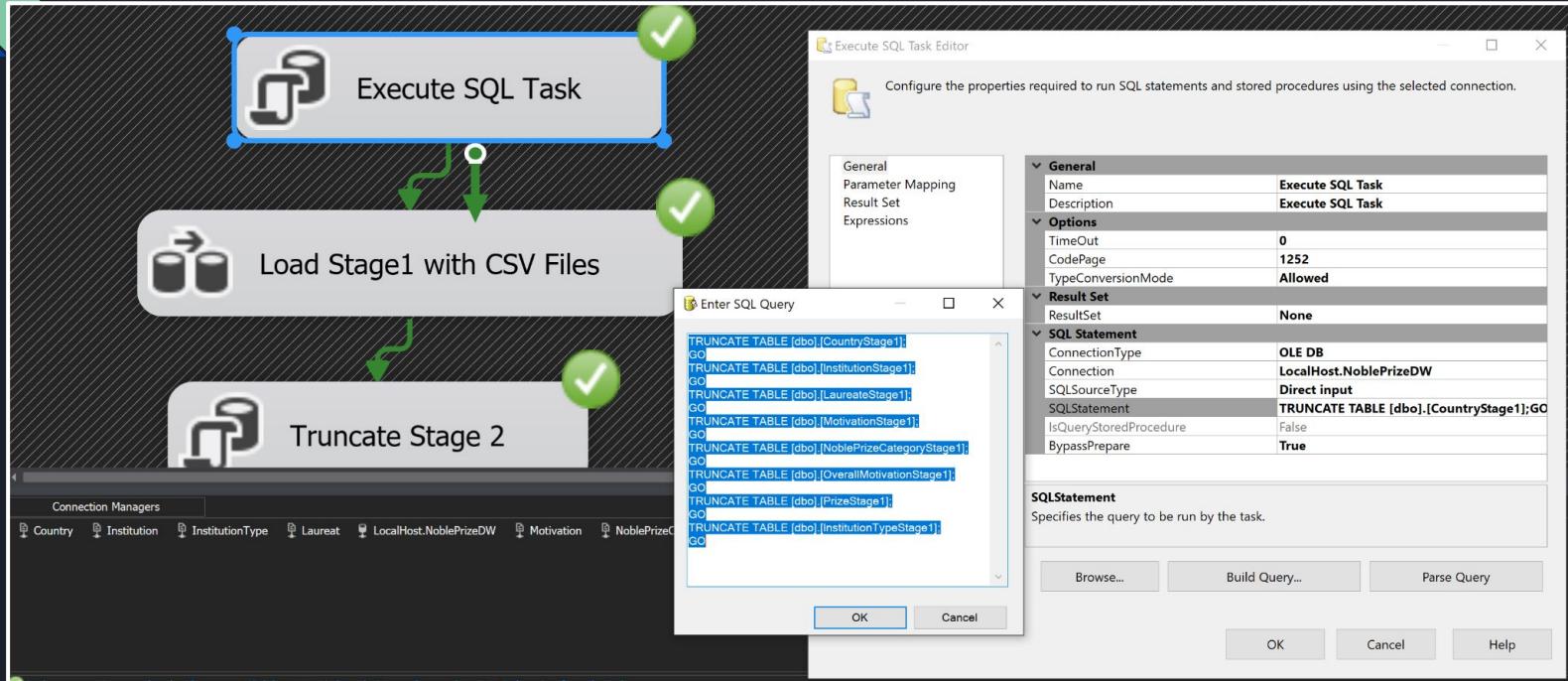


ETL Process: Whole Diagram



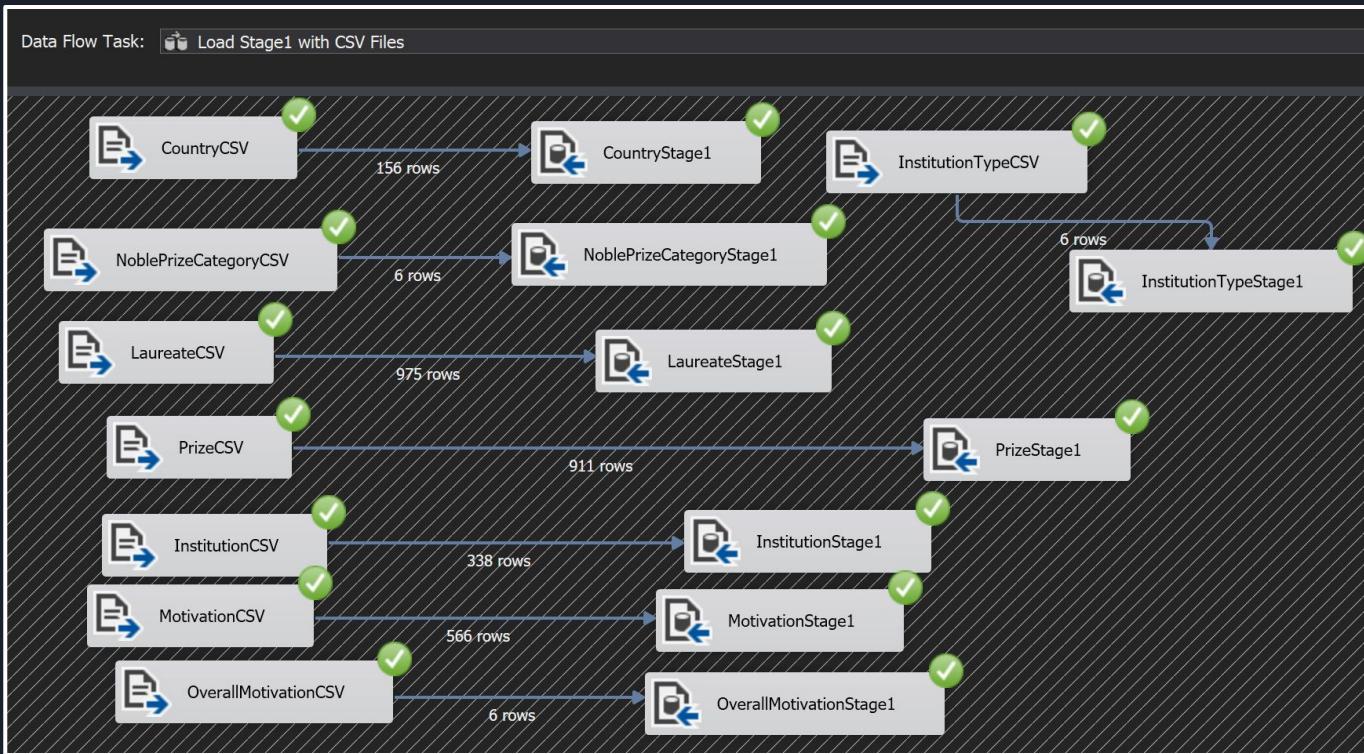
Shows the entire ETL Process' Steps from “Execute SQL Task” to “Load Facts”

ETL Process: Execute SQL Task



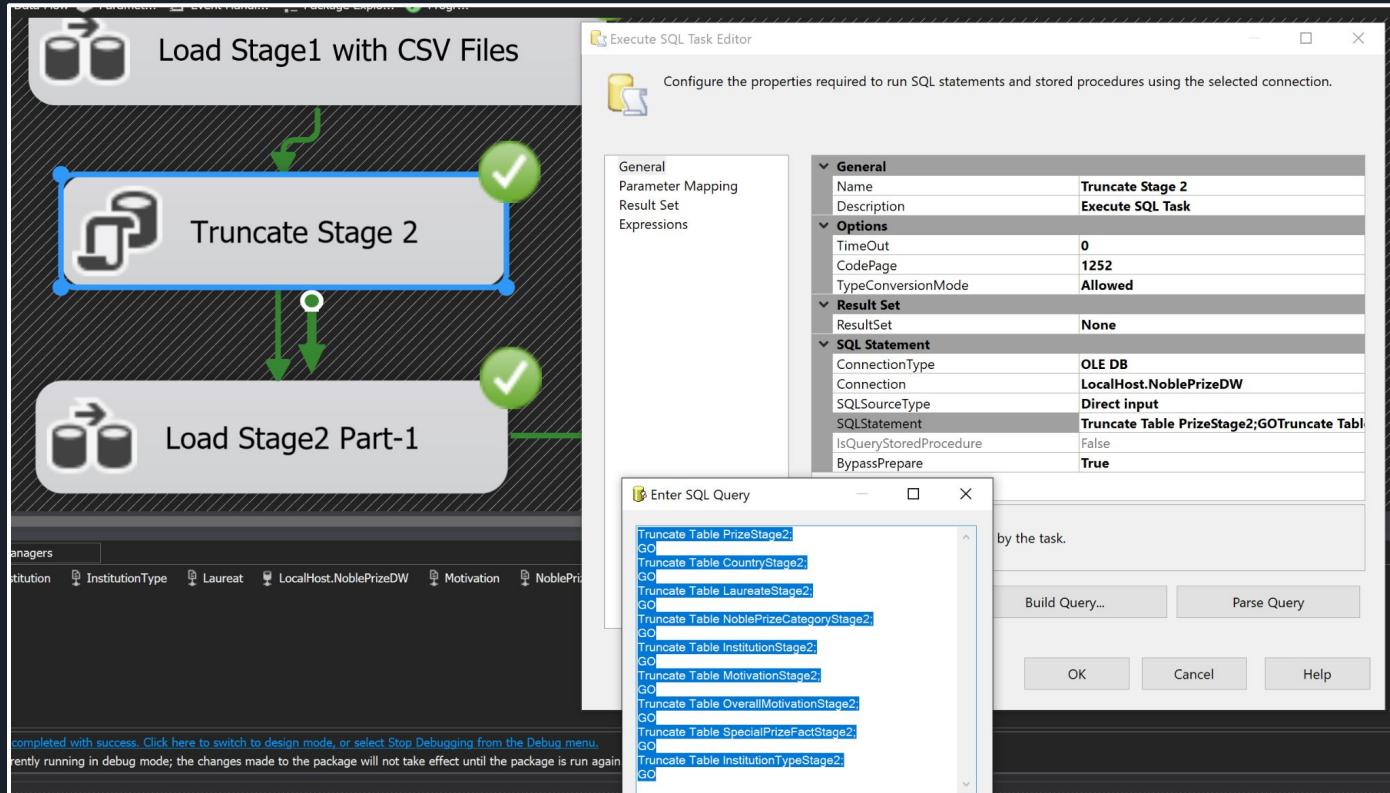
Extracts the data from the 8 CSV files to Load Stage1 Data Flow Task through Truncation

ETL Process: Load Stage1 with CSV Files



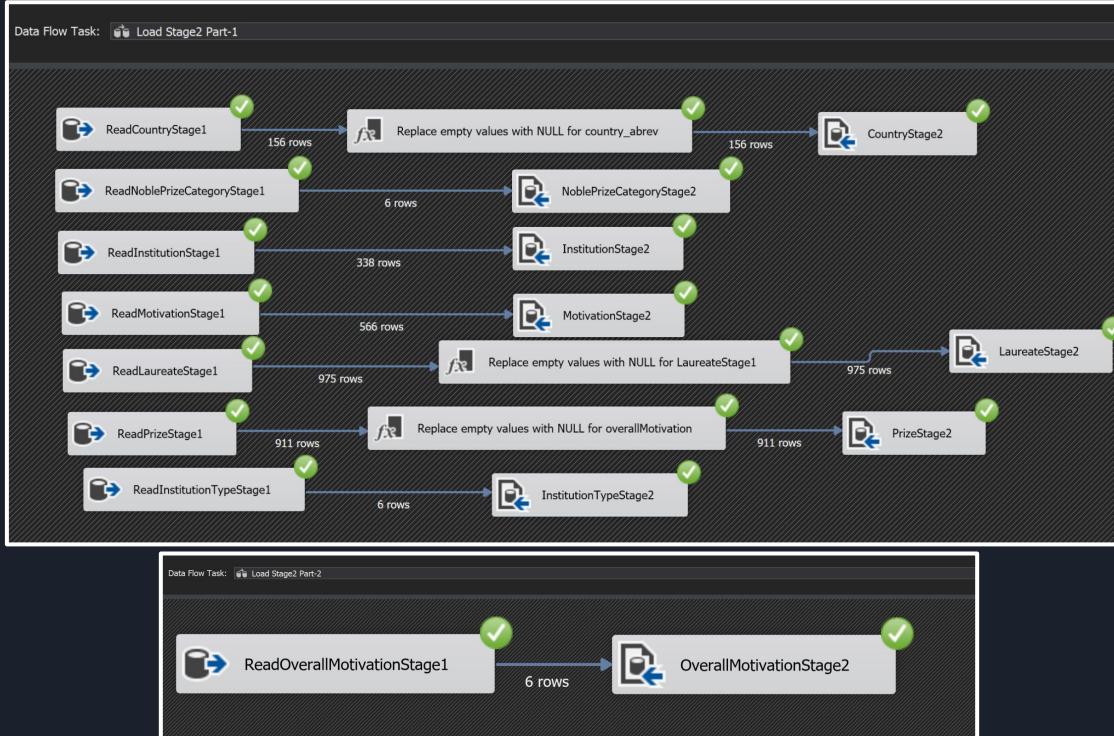
Extracts the data from the 8 CSV files to Stage1 SQL Server Destinations

ETL Process: Truncate Stage 2



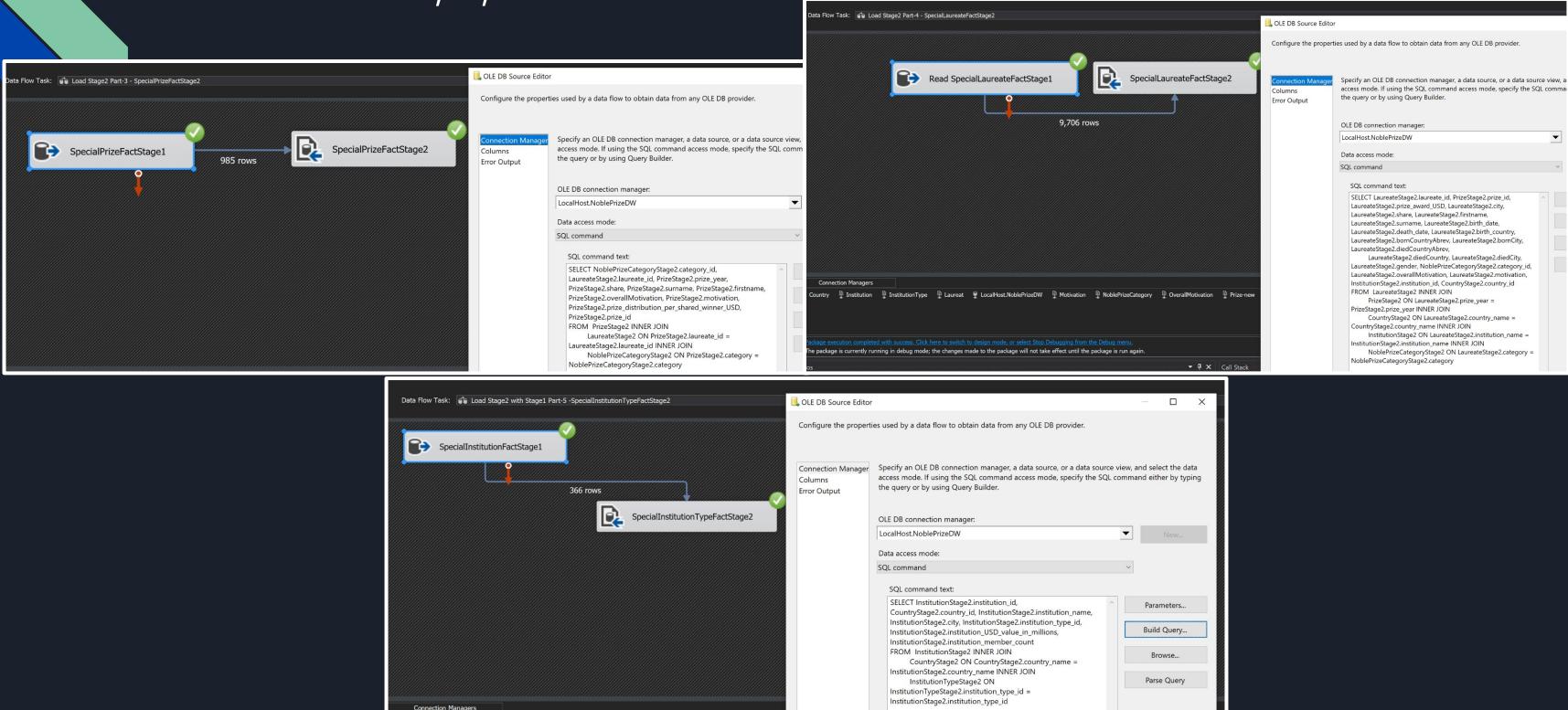
Extracts the data from the Stage1 to Load Stage 2 Data Flow Task through Truncation

ETL Process: Load Stage2 with Stage1 Part 1 and 2



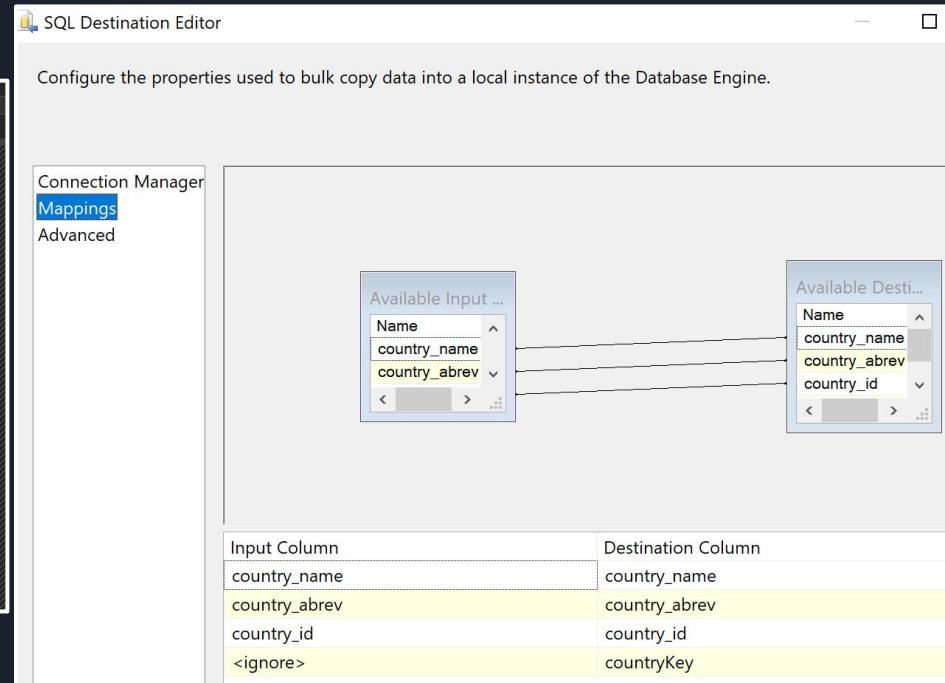
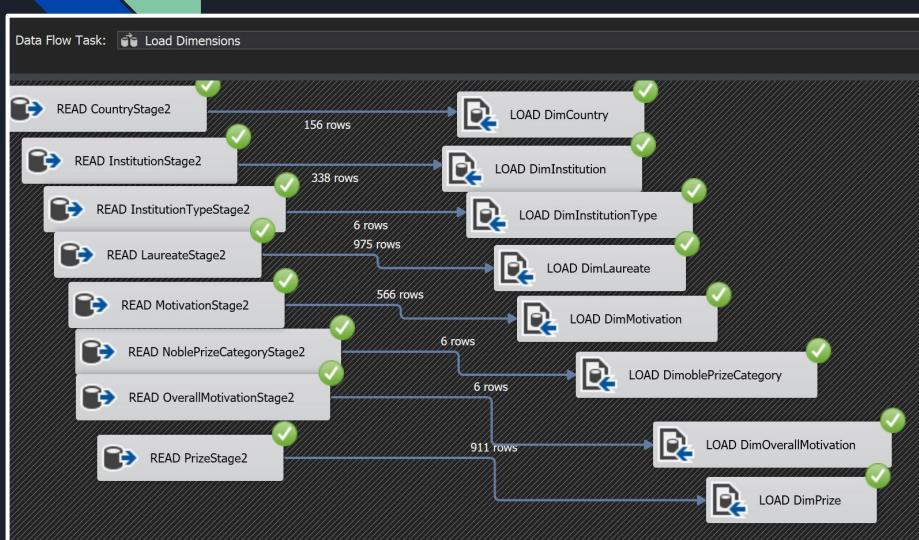
Extracts the data from the Stage1 OLE DB Sources to Stage 2 SQL Server Destinations with some Stages Replacing the Empty Values with NULL

ETL Process: Load Stage2 with Stage1 Part 3,4, and 5



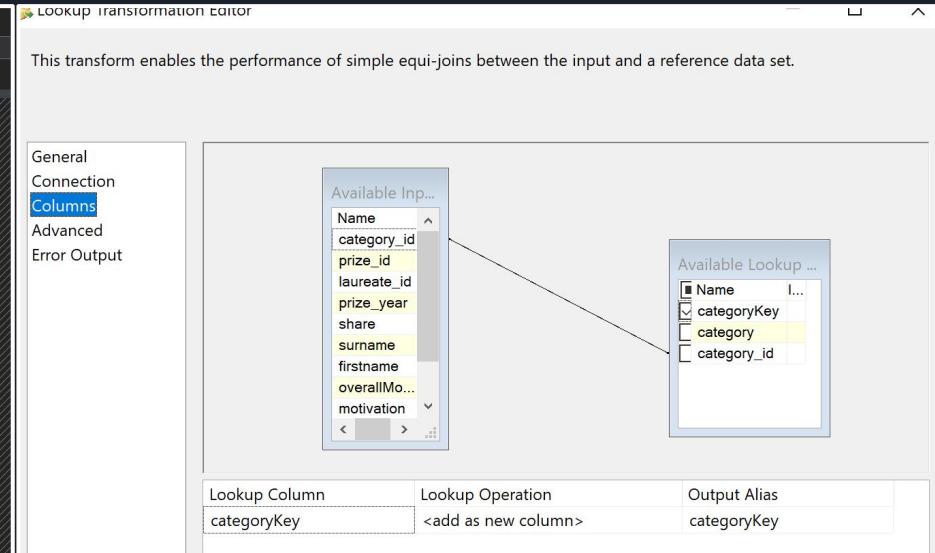
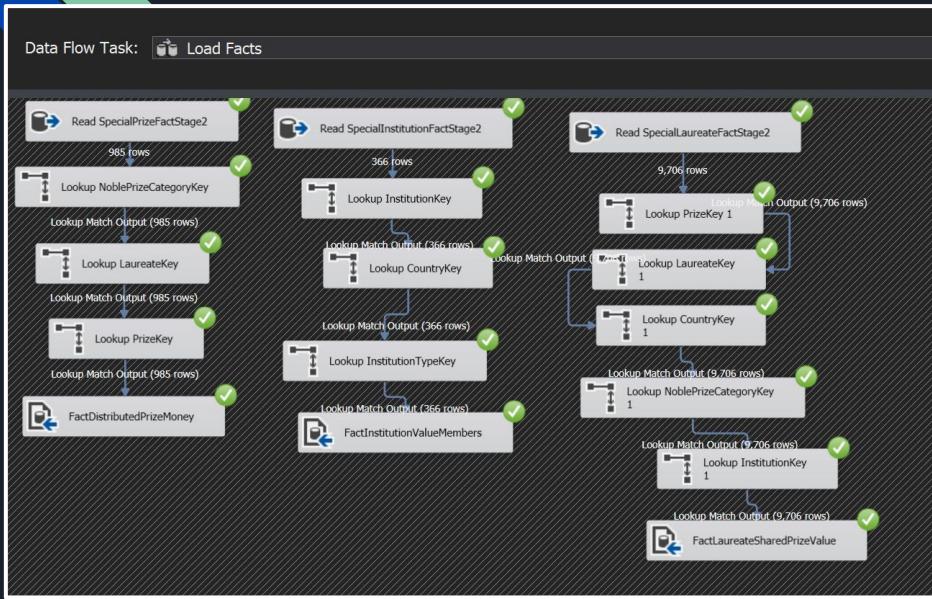
Extracts the data from the Stage1 OLE DB Sources to Stage 2 SQL Server Destinations with Stage 1 made through a SQL Commands

ETL Process: Load Dimensions



Extracts the data from the Stage2 OLE DB Sources to Load Dimension SQL Server Destinations with Loading Dimension SQL Server Destinations Creating Surrogate Keys that Replaces the Natural Keys of each Stage

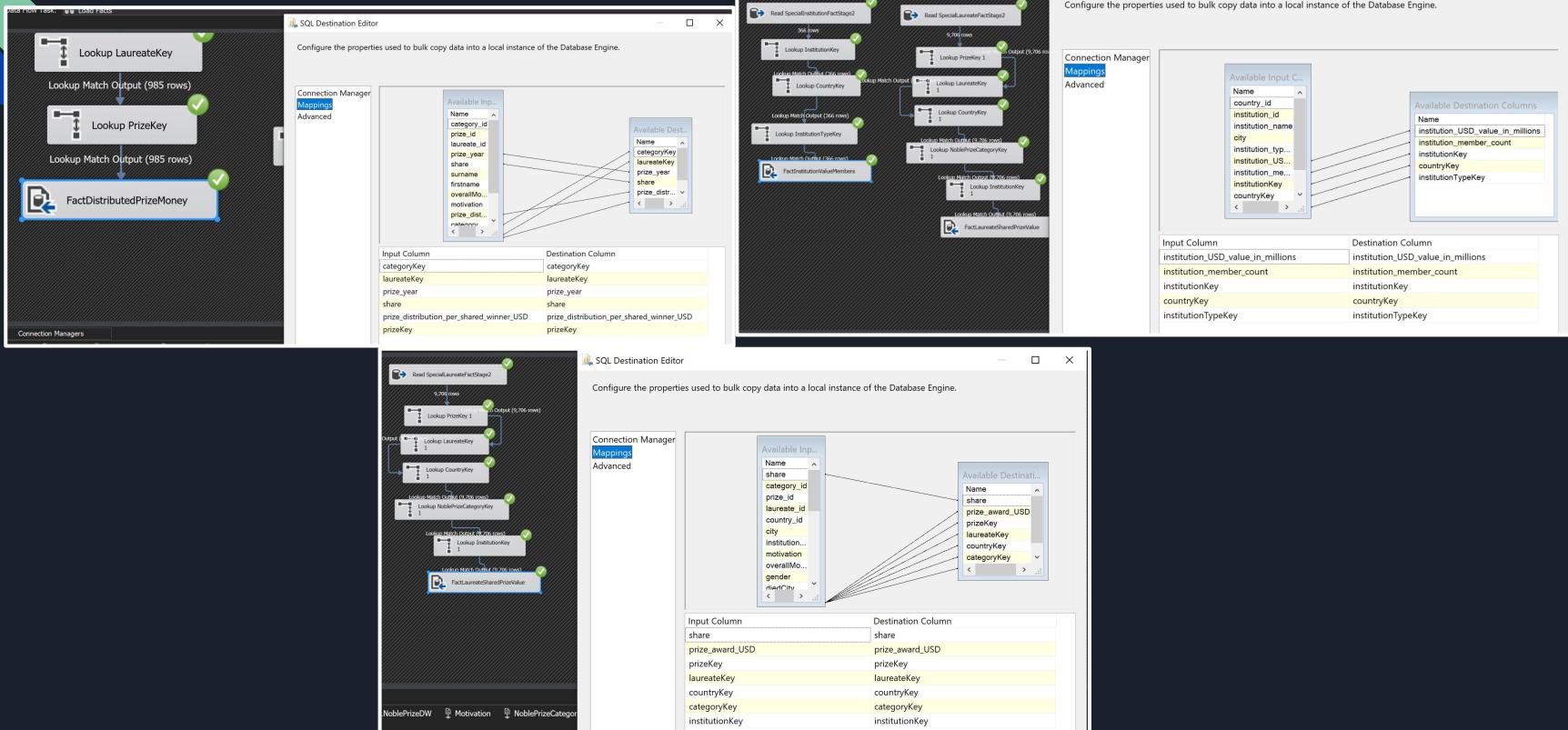
ETL Process: Load Facts - Part One



Extracts Data from Load Fact Stage 2 OLE DB Sources to Fact SQL Server Destination with Lookup Surrogate Key Pipelines for each Fact

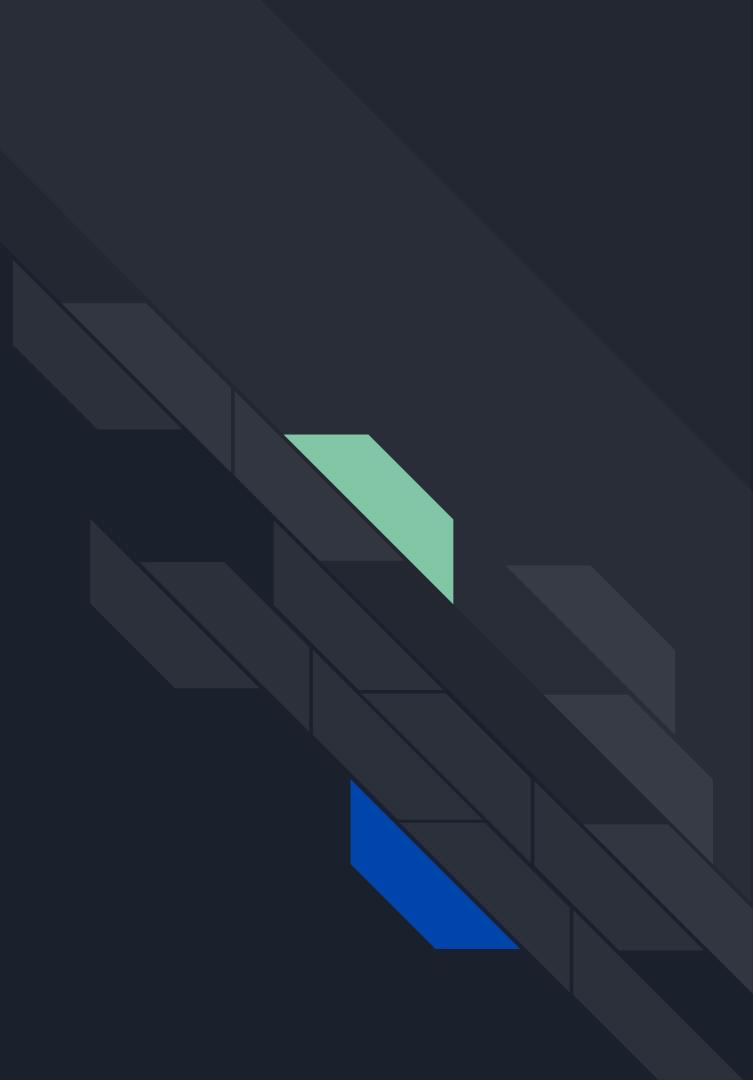
Lookup Surrogate Key Pipeline of Dimension Category as Example

ETL Process: Load Facts - Part Two



Extracts the data from the Stage2 OLE DB Sources to Fact SQL Server Destinations with Lookup Surrogate Key Pipelines

Data Warehouse Mart Views



Nobel Prize Data Warehouse with 3 Marts

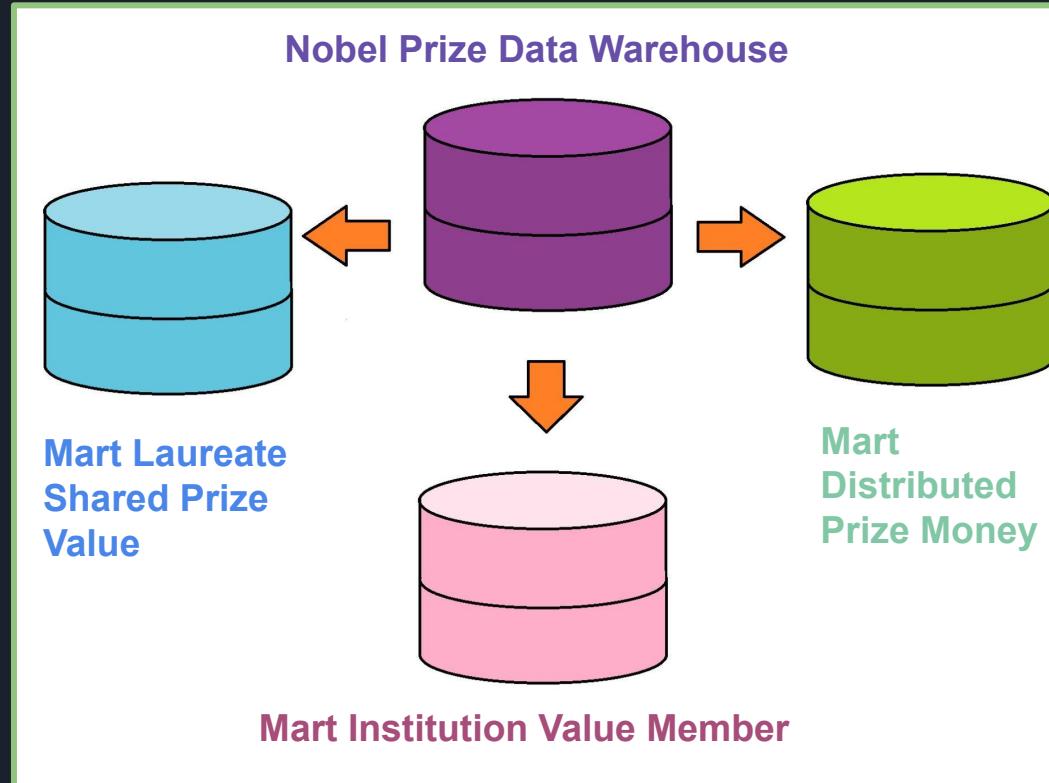
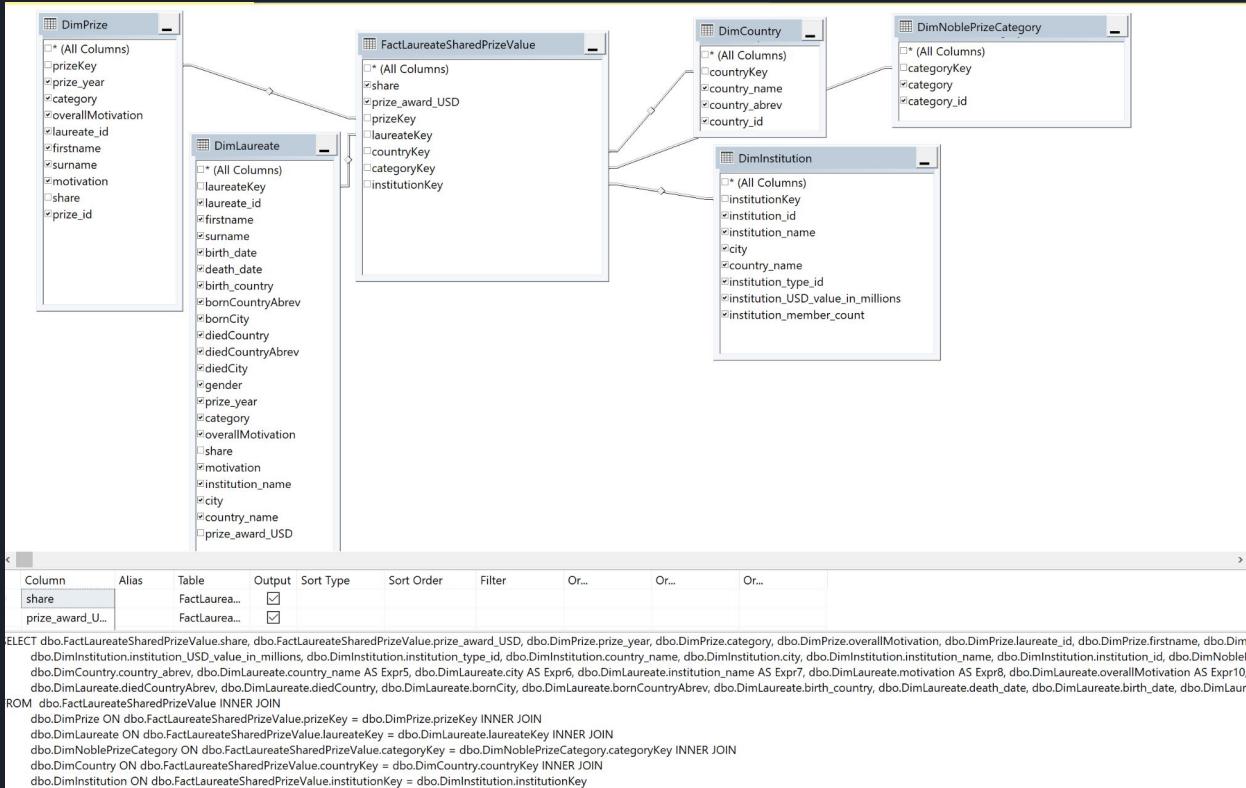


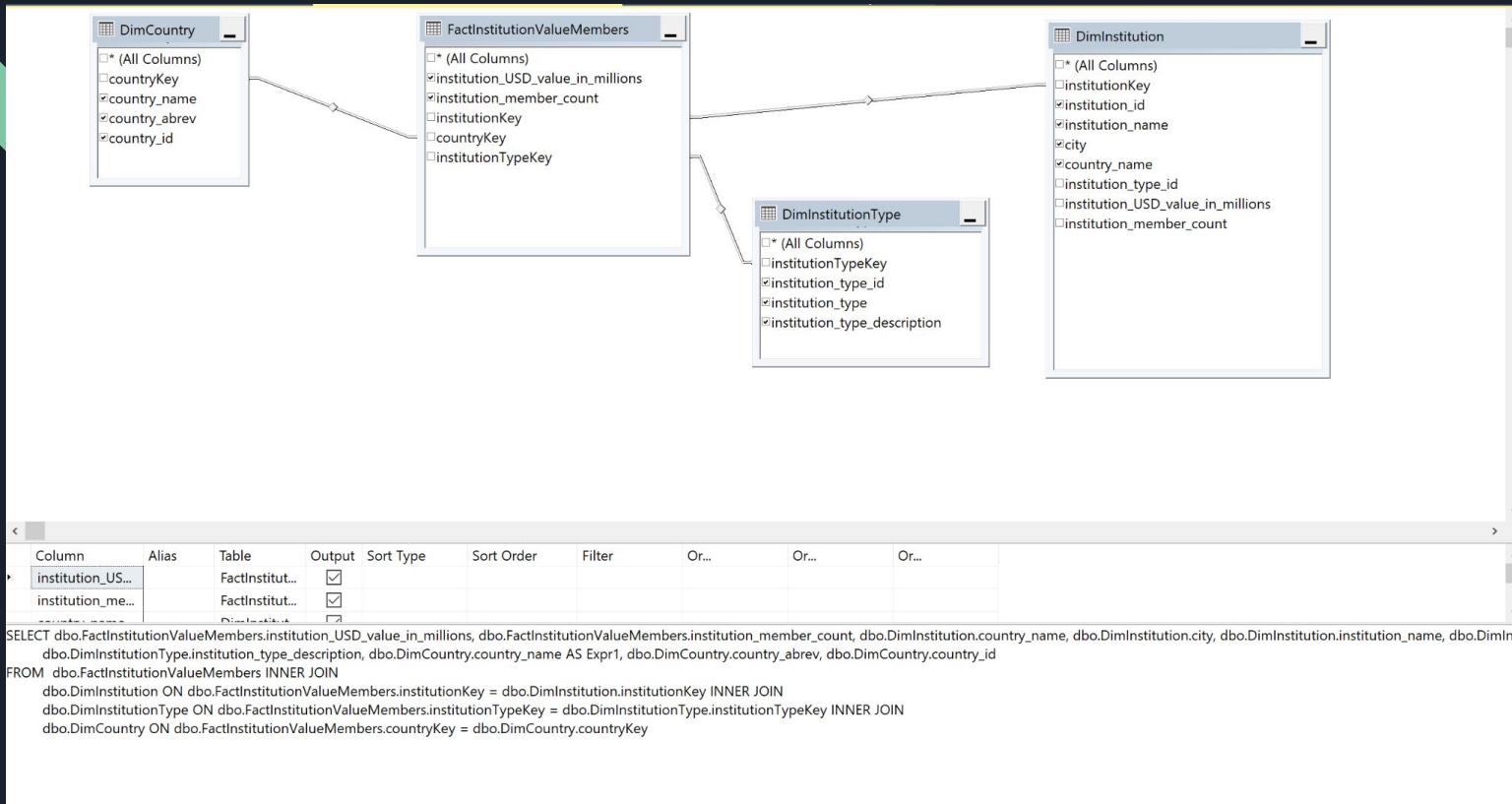
Diagram of the Nobel Prize Data Warehouse with the Laureate Shared Prize Value Mart, Institution Value Member Mart, and Distributed Prize Money Mart

Mart Laureate Shared Prize Value



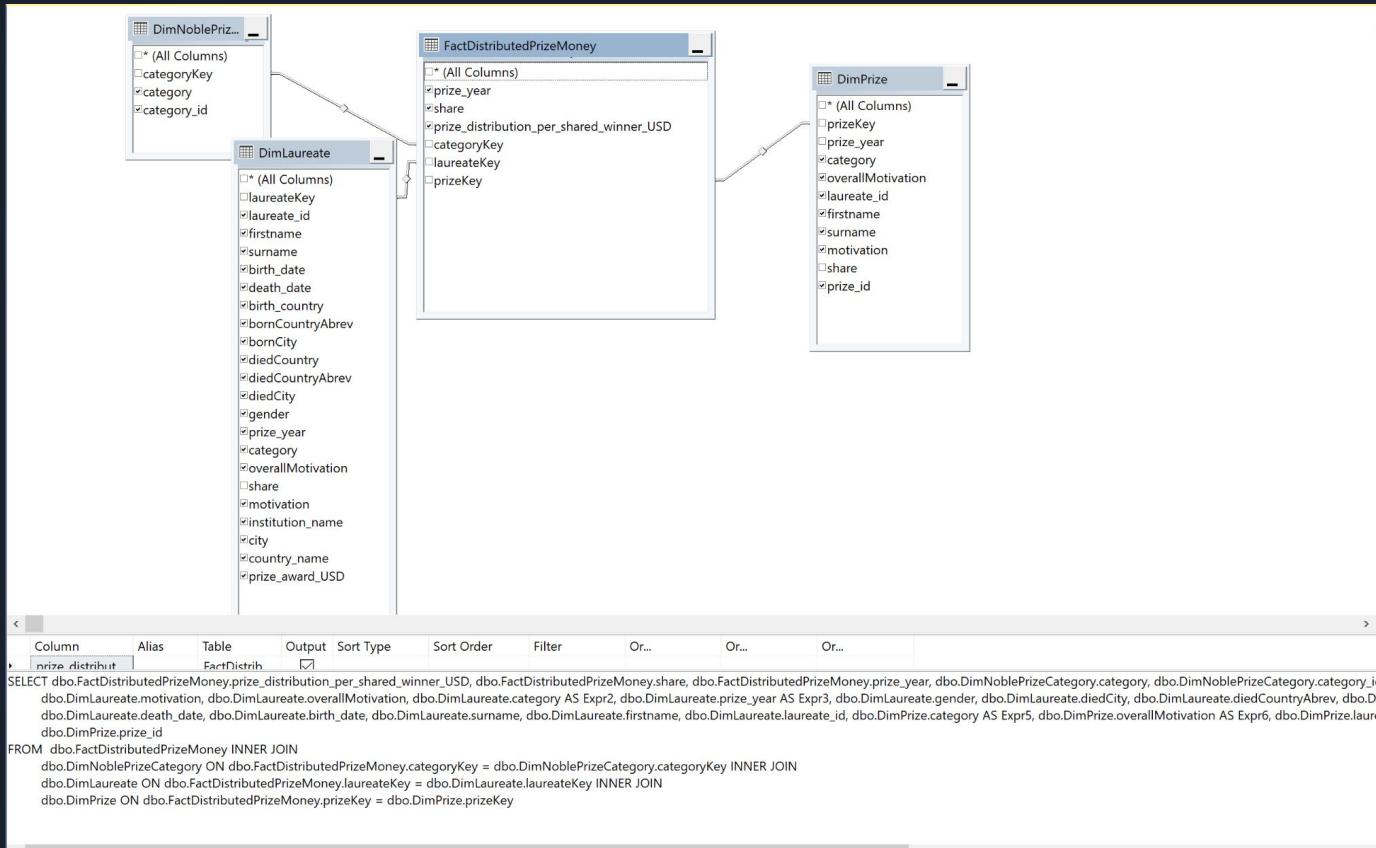
Nobel Prize Data Warehouse View of the Laureate Shared Prize Value Mart with the Laureate Shared Prize Value Fact Table and the Prize, Laureate, Country, Institution, and Noble Prize Category Dimension Tables

Mart Institution Value Member



Nobel Prize Data Warehouse View of the Institution Value Members Mart with the Institution Value Members Fact Table and the Country, Institution Type, and Institution Dimension Tables

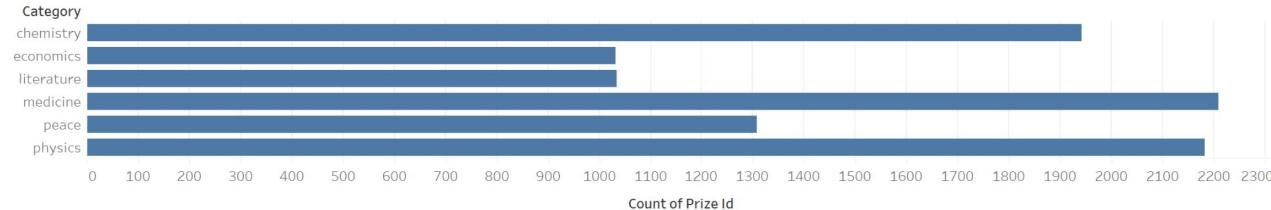
Mart Distributed Prize Money



Nobel Prize Data Warehouse View of the Distribution Prize Money Mart with the Distributed Prize Money Fact Table and the Laureate, Nobel Prize Categories, and Prize Dimension Tables

Tableau Data Visualizations

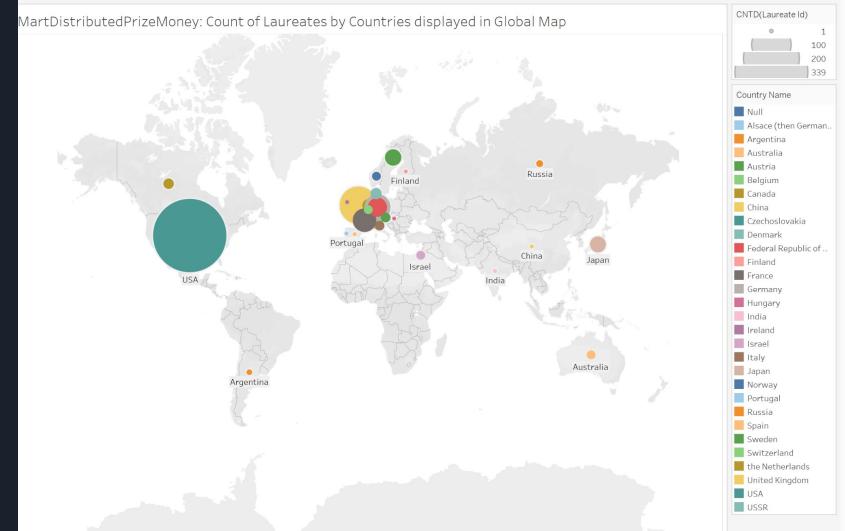
MartLaureateSharedPrizeValue: Count of Prizes Won by Nobel Prize Category



MartInstitutionValueMembers: Tree Map of the Sum of Institution USD Value in Millions by Institutions



MartDistributedPrizeMoney: Count of Laureates by Countries displayed in Global Map



SQL Queries

```

SELECT [prize_year], [1] as SHARED_ONE, [2] as SHARED_TWO, [3] as SHARED_THREE, [4] as SHARED_FOUR
FROM
(
    SELECT [prize_year]. [share], [prize_distribution_per_shared_winner_USD]
    FROM [NoblePrizeDW].[dbo].[MartDistributedPrizeMoney]
) ps
PIVOT
( AVG([prize_distribution_per_shared_winner_USD])
    FOR [share] IN ([1],[2],[3],[4])
) AS pvt

```

```

=SELECT InstitutionStage2.institution_id, CountryStage2.country_id, InstitutionStage2.institution_name,
InstitutionStage2.city, InstitutionStage2.institution_type_id, InstitutionStage2.institution_usd_value_in_millions,
InstitutionStage2.institution_member_count
FROM [NoblePrizeDW].[dbo].InstitutionStage2 INNER JOIN
[NoblePrizeDW].[dbo].CountryStage2 ON CountryStage2.country_name = InstitutionStage2.country_name INNER JOIN
[NoblePrizeDW].[dbo].InstitutionTypeStage2 ON
InstitutionTypeStage2.institution_type_id = InstitutionStage2.institution_type_id

```

```

:SELECT [NoblePrizeDW].[dbo].[NoblePrizeCategoryStage2].[category_id],
[NoblePrizeDW].[dbo].[LaureateStage2].[laureate_id], [PrizeStage2].[prize_id], [PrizeStage2].[prize_year],
[PrizeStage2].[prize_name], [PrizeStage2].[prize_award_usd], [PrizeStage2].[share],
[PrizeStage2].[prize_distribution_per_shared_winner_usd], [PrizeStage2].[motivation],
[PrizeStage2].[prize_stage], [PrizeStage2].[category]
FROM [NoblePrizeDW].[dbo].[NoblePrizeCategoryStage2] ON
[NoblePrizeDW].[dbo].[NoblePrizeCategoryStage2].[category_id] = [NoblePrizeDW].[dbo].[LaureateStage2].[laureate_id]
ON [NoblePrizeDW].[dbo].[PrizeStage2].[category] = [NoblePrizeDW].[dbo].[NoblePrizeCategoryStage2].[category]

```

```

% - 4 Results Messages
category_id laureate_id prize_year share category
NobelPrizeDW].[dbo].[NoblePrizeCategoryStage2].[category_id], [NoblePrizeDW].[dbo].[LaureateStage2].[laureate_id]
[NobelPrizeDW].[dbo].[LaureateStage2].[laureate_id], [PrizeStage2].[prize_id], [PrizeStage2].[prize_year],
[PrizeStage2].[prize_name], [PrizeStage2].[prize_award_usd], [PrizeStage2].[share],
[PrizeStage2].[prize_distribution_per_shared_winner_usd], [PrizeStage2].[motivation],
[PrizeStage2].[prize_stage], [PrizeStage2].[category]
FROM [NoblePrizeDW].[dbo].[NoblePrizeCategoryStage2] ON
[NoblePrizeDW].[dbo].[NoblePrizeCategoryStage2].[category_id] = [NoblePrizeDW].[dbo].[LaureateStage2].[laureate_id]
ON [NoblePrizeDW].[dbo].[PrizeStage2].[category] = [NoblePrizeDW].[dbo].[NoblePrizeCategoryStage2].[category]

```

```

% - 4 Results Messages
category_id laureate_id prize_award_usd [PrizeStage2].[prize_id], [PrizeStage2].[prize_year] INNER JOIN
[NobelPrizeDW].[dbo].[LaureateStage2].[laureate_id], [PrizeStage2].[prize_id], [PrizeStage2].[prize_year], [PrizeStage2].[share], [NoblePrizeDW].[dbo].[LaureateStage2].[city],
[NobelPrizeDW].[dbo].[LaureateStage2].[firstname], [NobelPrizeDW].[dbo].[LaureateStage2].[surname], [NobelPrizeDW].[dbo].[LaureateStage2].[date_of_birth], [NobelPrizeDW].[dbo].[LaureateStage2].[date_of_death],
[NobelPrizeDW].[dbo].[LaureateStage2].[birth_country], [NobelPrizeDW].[dbo].[LaureateStage2].[death_country], [NobelPrizeDW].[dbo].[LaureateStage2].[birth_date],
[NobelPrizeDW].[dbo].[LaureateStage2].[death_date], [NobelPrizeDW].[dbo].[LaureateStage2].[category_id], [NobelPrizeDW].[dbo].[LaureateStage2].[category], [NobelPrizeDW].[dbo].[LaureateStage2].[motivation],
[NobelPrizeDW].[dbo].[LaureateStage2].[category], [NobelPrizeDW].[dbo].[LaureateStage2].[category], [NobelPrizeDW].[dbo].[LaureateStage2].[category]
FROM [NobelPrizeDW].[dbo].[NoblePrizeCategoryStage2] ON
[NobelPrizeDW].[dbo].[NoblePrizeCategoryStage2].[category_id] = [NobelPrizeDW].[dbo].[LaureateStage2].[laureate_id]
INNER 2018

```

```

% - 4 Results Messages
category_id laureate_id prize_award_usd [PrizeStage2].[prize_id], [PrizeStage2].[prize_year] INNER JOIN
[NobelPrizeDW].[dbo].[LaureateStage2].[laureate_id], [PrizeStage2].[prize_id], [PrizeStage2].[prize_year], [PrizeStage2].[share], [NoblePrizeDW].[dbo].[LaureateStage2].[city],
[NobelPrizeDW].[dbo].[LaureateStage2].[firstname], [NobelPrizeDW].[dbo].[LaureateStage2].[surname], [NobelPrizeDW].[dbo].[LaureateStage2].[date_of_birth], [NobelPrizeDW].[dbo].[LaureateStage2].[date_of_death],
[NobelPrizeDW].[dbo].[LaureateStage2].[birth_country], [NobelPrizeDW].[dbo].[LaureateStage2].[death_country], [NobelPrizeDW].[dbo].[LaureateStage2].[birth_date],
[NobelPrizeDW].[dbo].[LaureateStage2].[death_date], [NobelPrizeDW].[dbo].[LaureateStage2].[category_id], [NobelPrizeDW].[dbo].[LaureateStage2].[category], [NobelPrizeDW].[dbo].[LaureateStage2].[motivation],
[NobelPrizeDW].[dbo].[LaureateStage2].[category], [NobelPrizeDW].[dbo].[LaureateStage2].[category], [NobelPrizeDW].[dbo].[LaureateStage2].[category]
FROM [NobelPrizeDW].[dbo].[NoblePrizeCategoryStage2] ON
[NobelPrizeDW].[dbo].[NoblePrizeCategoryStage2].[category_id] = [NobelPrizeDW].[dbo].[LaureateStage2].[laureate_id]
INNER 2018

```



Final Comments and Difficulties Faced

- Cleaning the Data
 - Foreign names
 - Countries that no longer existed
- ETL process
 - Stage 2
- Tableau and Data Visualization
 - Unfamiliar software
 - Applying effective data visualization



References

Kuzmenko, Maryna. (2016). *Nobel Prize - Dataset with Information about Prizes, Laureates and Countries*. Harvard Dataverse, V1, UNF:6:McdDh+IdUTGgZDs5XVOQUA== [fileUNF].
<https://doi.org/10.7910/DVN/AGAFAQ>

The Nobel Foundation. (2022, April 13). *The Nobel Prize*. NobelPrize.Org.
<https://www.nobelprize.org/>