

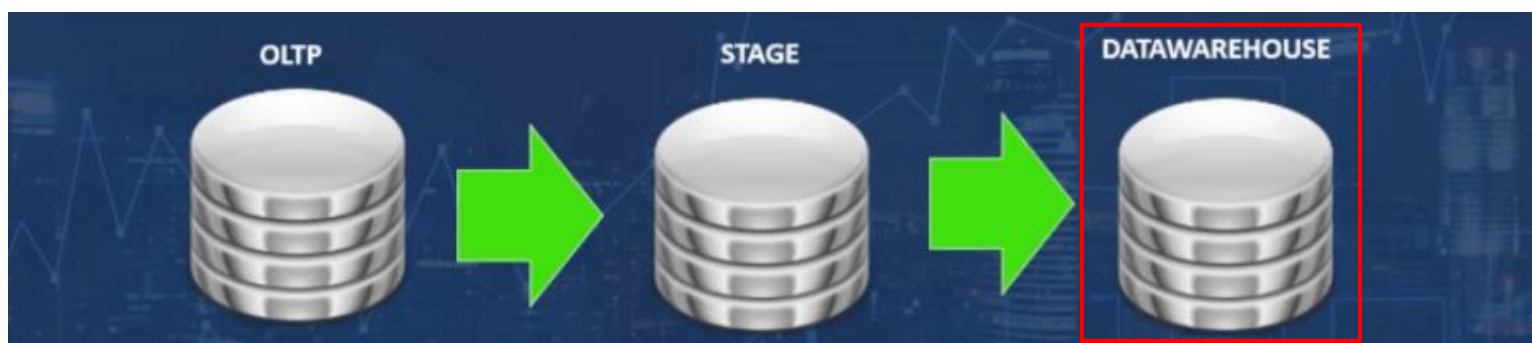
DATA WAREHOUSE – GOIÁS MARKET

O *Data Warehouse*, ou DW, é uma das partes mais importantes do projeto de *Business Intelligence*. Aqui, culmina-se todos os esforços de cargas e transformações de dados – ETL – necessários para armazenamento de dados concentrando-as no formato correto para análises, construção de *dashboards* e *reports* necessários para os analistas, gestores e tomadores de decisões no nível gerencial.

A partir do Data Warehouse, assim como os próximos projetos, irão tratar do ambiente OLAP, destinado à análises e automação de relatórios.

Um rápido paralelo aos conceitos de modelagem de um DW, pode-se citar os *Data Mart*, que são bancos de dados com a mesma função do DW, porém, com um assunto (ou Fato) específico. Como exemplo, *Data Mart* destinados ao Marketing, ou Logística, ou Financeiro, ou Comercial, etc.

Neste projeto, não serão modelados *Data Mart*. Apenas o *Data Warehouse* que irá armazenar todos os dados do projeto.



Portanto, após as etapas de ETL do sistema OLTP para a *STAGE AREA* e, em seguida, para o *Data Warehouse*, os dados já estão prontos para ingestão e análise.

Como será definido mais adiante, é importante ressaltar que o DW também possui a função de guardar o histórico de modificações das entradas do OLTP.

Uma vez que algumas informações podem ser modificadas, tal como endereço, local de residência, telefone, sexo,

OBJETIVO DO PROJETO – DATA WAREHOUSE

O projeto do *Data Warehouse* tem como objetivo o armazenamento dos dados oriundos do sistema transacional do negócio (banco OLTP), a fim de diminuir o tráfego de informações que causam lentidões, e disponibilizar esses dados de forma estruturada para serem analisados e interpretados pelos analistas e tomadores de decisões.

Para este projeto, o objetivo do DW visa responder as seguintes questões de negócio:

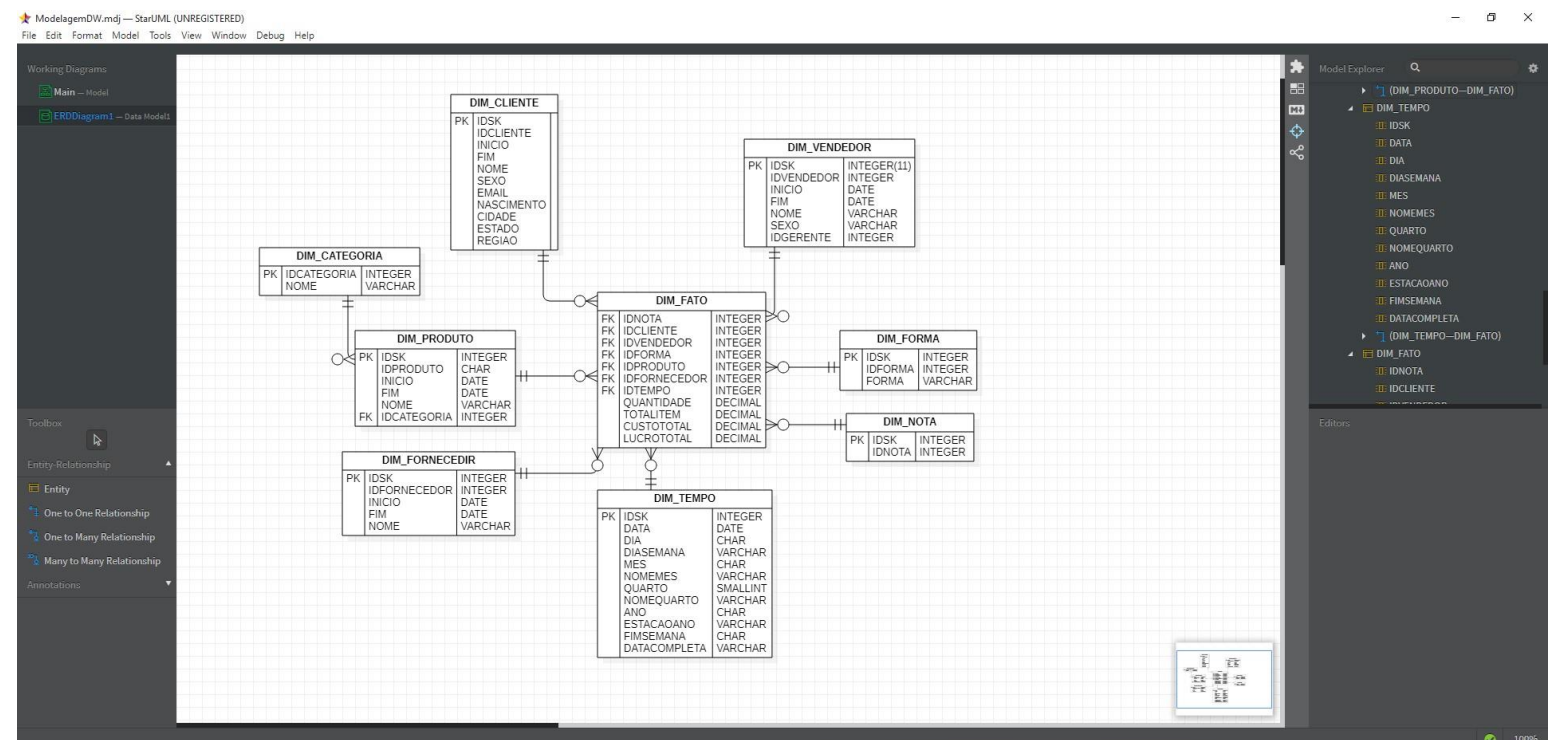
- Quem são os melhores clientes?
- Quem são os melhores vendedores?
- Qual categoria rende mais?
- Qual a minha relação com os fornecedores?
- Qual meu pior e melhor produto?
- Em qual região eu vendo mais?

MODELAGEM LÓGICA

A modelagem do *Data Warehouse* possui grande foco na integração das tabelas. A integridade relacional deve ser garantida a fim de que todos os dados estejam conectados corretamente através das *Primary Key* (PK, ou chave primária) e *Surrogate Key* (SK, ou chave substituta).

No modelo a seguir, pode-se destacar a tabela DIM_FATO, que é a tabela principal do DW. Essa tabela poderia ser facilmente escrita por um Query SQL, trazendo todos os dados na seleção. Entretanto, a tabela FATO consolida a função do *Data Warehouse*, facilitando a consulta de dados.

A tabela FATO, como demonstrado abaixo, recebe os relacionamentos de cardinalidade de 1 para muitos (1 x N). Isto significa que a tabela fato recebe os registros únicos de cada tabela relacionada, bem como o seu histórico de modificação, consolidando as informações de uma PK em um único lugar.



Os dados armazenados no *Data Warehouse* já estão prontos para consultas pelos analistas, sejam elas através de *queries* diretamente do banco de dados, por queries SQL, por exportação para um sistema de análise em ambiente OLAP (como o SQL Server Analysis Services) ou por uma conexão com alguma ferramenta de *DataViz* (como o Microsoft Power BI, Tableau, Qlik Sense e etc).

Nesse ponto, é válido também a conexão com ferramentas de Data Science, como R e Python, que são capazes de acessar diretamente o DW através de bibliotecas específicas para conexão e interface com SQL.

Os pipelines em R e Python serão modelados futuramente neste mesmo projeto, para modelagem de dados, aplicação de modelos de *Machine Learning* e *Data Mining*.



Ao modelar o DW, há um novo atributo nas tabelas em relação ao banco OLTP e a *Stage Area*, são as colunas “INICIO” e “FIM”, do tipo *DATE*. Estas colunas serão responsáveis por manter o registro de alteração de algum item.

O conceito que se aplica em tabelas com capacidade de armazenar histórico por alguma alteração, são definidas como dimensões mutáveis, ou SCD (*Slowly Changing Dimensions*)

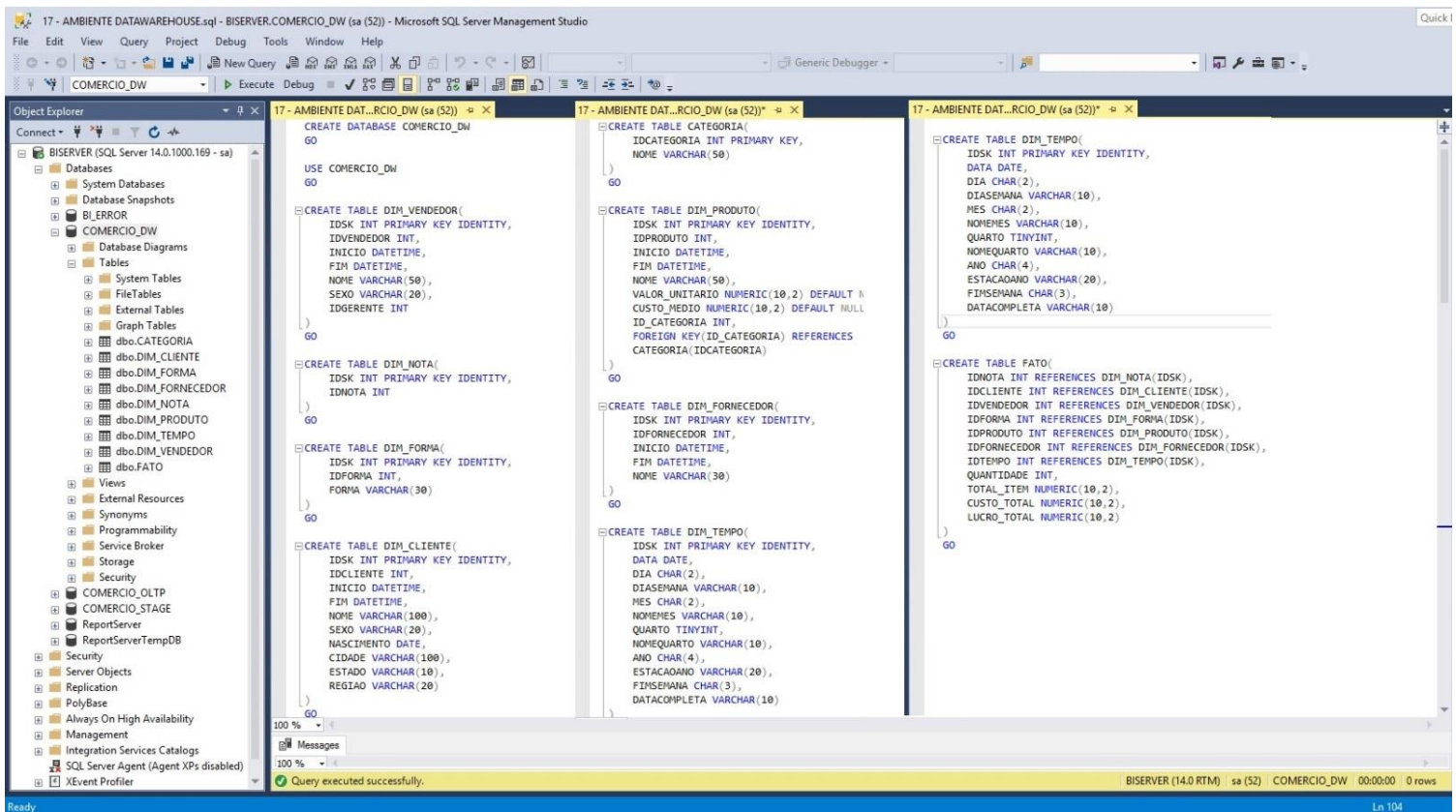
Dessa forma, é possível guardar um registro de alteração de um item ao longo tempo sem perder a rastreabilidade, pois a chave primeira (PK) se manterá igual, e os outros atributos que serão alterados terão uma data de validade, representada na coluna “FIM”, que é a data em que foi feita a alteração e este item será substituído pela sua alteração.

Feito isso, cria-se uma nova entrada no DW com a mesma *Primary Key* (PK), porém com uma *Surrogate Key* (SK) diferente, já com as devidas alterações.

As tabelas que podem sofrer esse tipo de alteração, são: DIM_PRODUTO, DIM_CLIENTE, DIM_FORNECEDOR E DIM_VENDEDOR. Todas essas tabelas podem sofrer alteração em alguma de suas entradas mediante a necessidade e serem salvas no *Data Warehouse* durante o processo de ETL. O *Integration Services* será responsável por esse processo e será exemplificado a seguir.

Com os conceitos já definidos e a modelagem lógica feita, partimos então para a modelagem física do *Data Warehouse*.

Começamos então a modelagem física no banco de dados, pelo SQL Server, com a criação do banco de dados do *Data Warehouse*. Neste projeto, o DW será chamado de COMERCIO_DW.



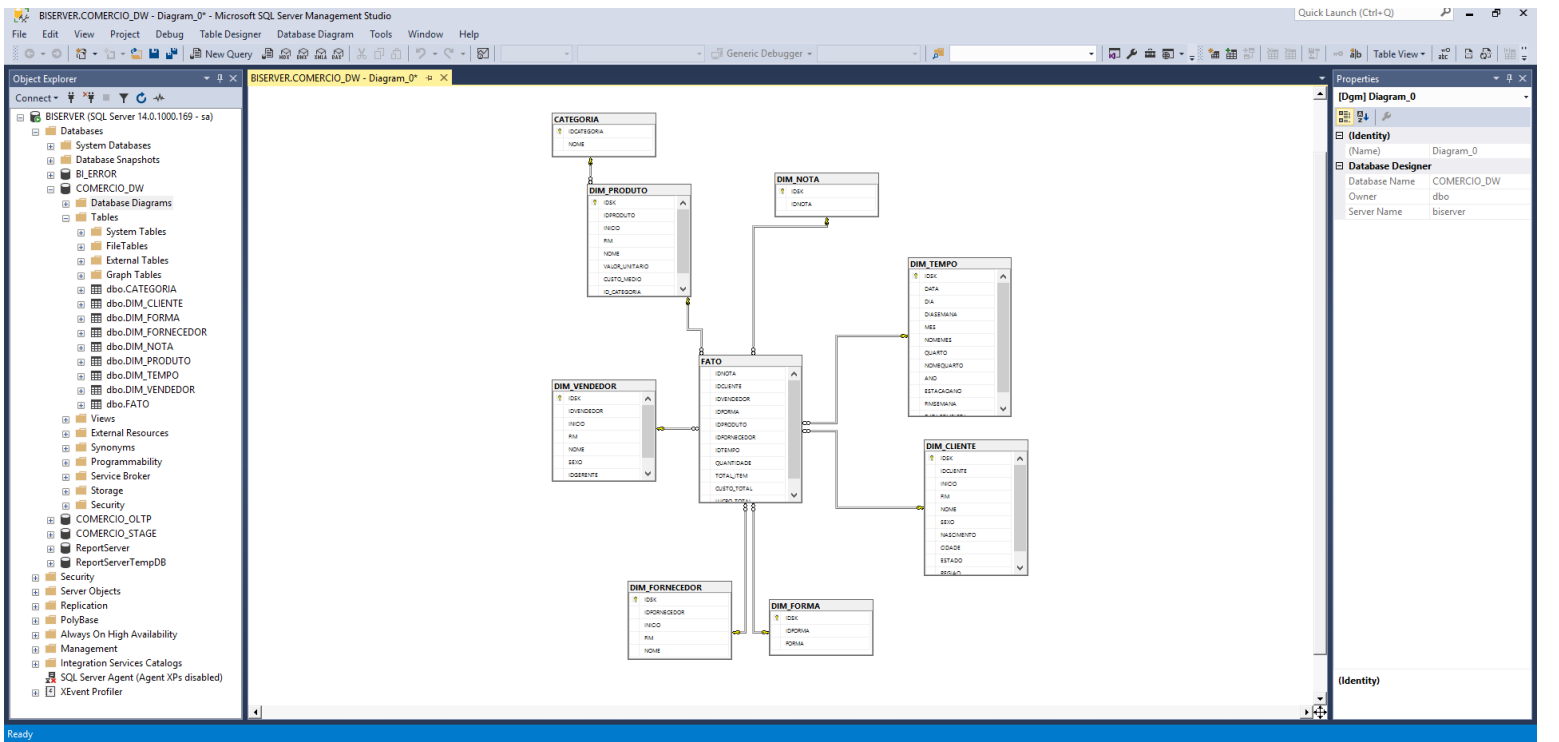
As tabelas foram criadas com o prefixo “DIM_”, indicando as dimensões do DW separadas por assunto.

As desse banco foram criadas em uma ordem específica já com os seus relacionamentos (ou *Constraints*) definidos, interligando as chaves primárias de cada uma. Apesar de que o procedimento de criação de *constraints* possa ser feito a parte, foram aplicados mediante a criação das tabelas para evitar qualquer fuga de relacionamentos.

A tabela DIM_FATO, ilustra bem a modelagem de *constraints* pois é a tabela mais importante do DW, que necessita estar totalmente referenciada com as dimensões que a compõe.

A função SQL utilizada para referenciar as tabelas durante a criação, segue na última *query*:
REFERENCES

Neste momento, a banco já está criado, como segue o diagrama gerado pelo SQL Server, e as tabelas já estão prontas para receber os dados por ETL.

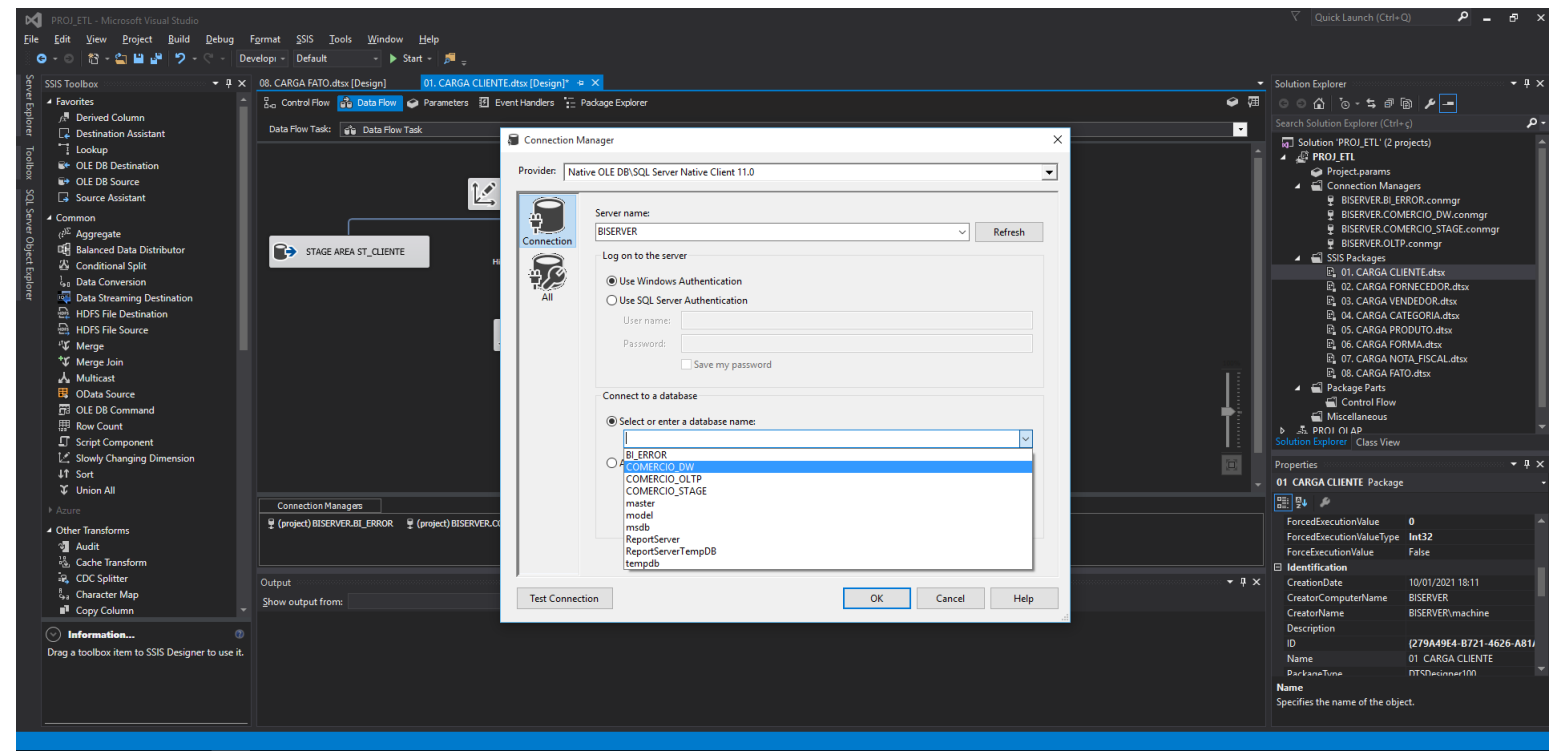


SINCRONIZAÇÃO DE ETL COM INTEGRATION SERVICES

O *Integration Service* (SSIS) é uma versátil ferramenta de ETL, construindo os pacotes de cargas e os gerenciando por projetos e por solução, como segue no menu à direita.

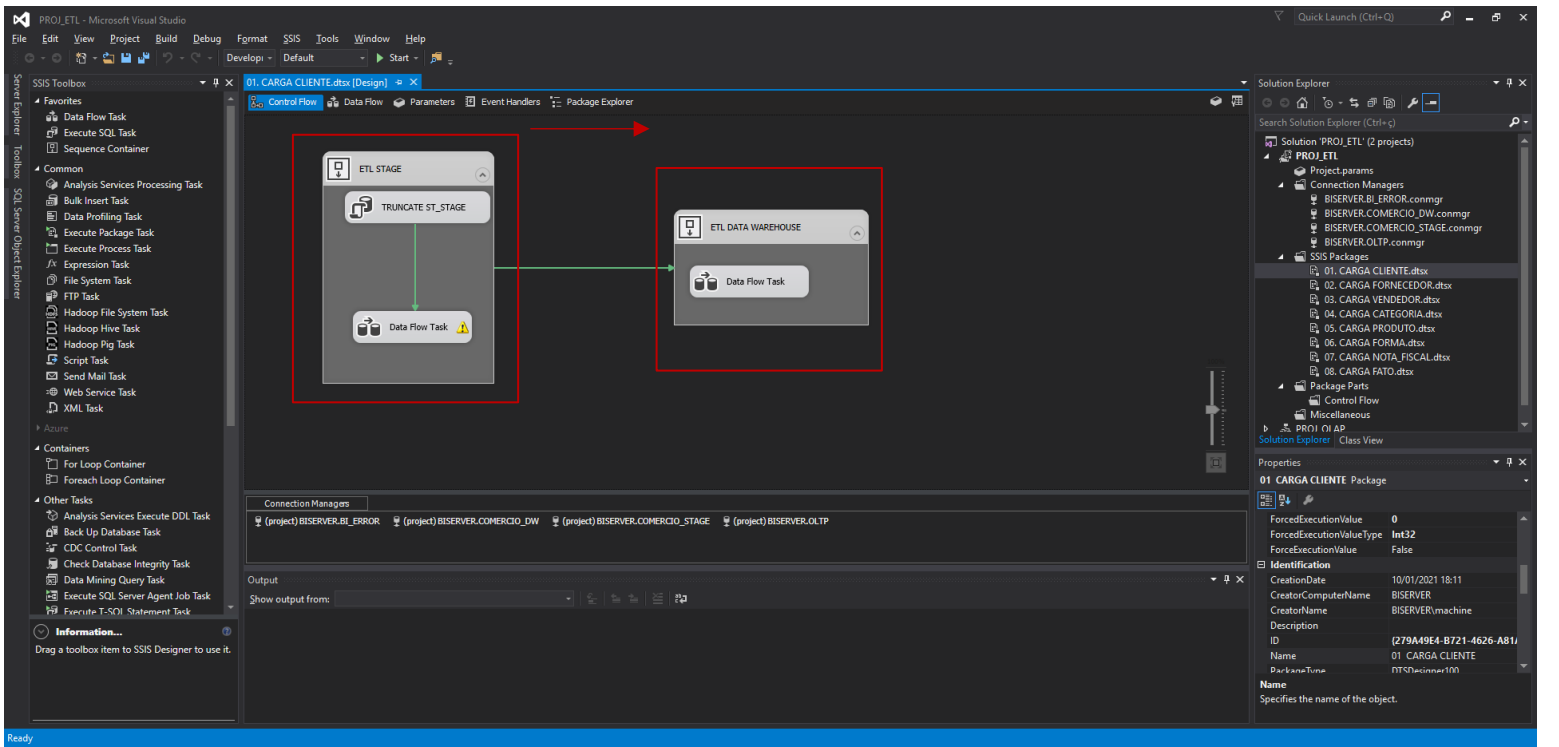
Nesse projeto, os pacotes terão o processo de ETL sincronizados com a carga da Stage Area. Isso significa que o pacote irá fazer separadamente e consequentemente as cargas dos bancos.

As conexões entre os bancos foram definidas previamente e foram compartilhadas como conexões de toda essa solução BI. Neste momento, é criada a conexão com o banco COMERCIO_DW, que é o *Data Warehouse*.



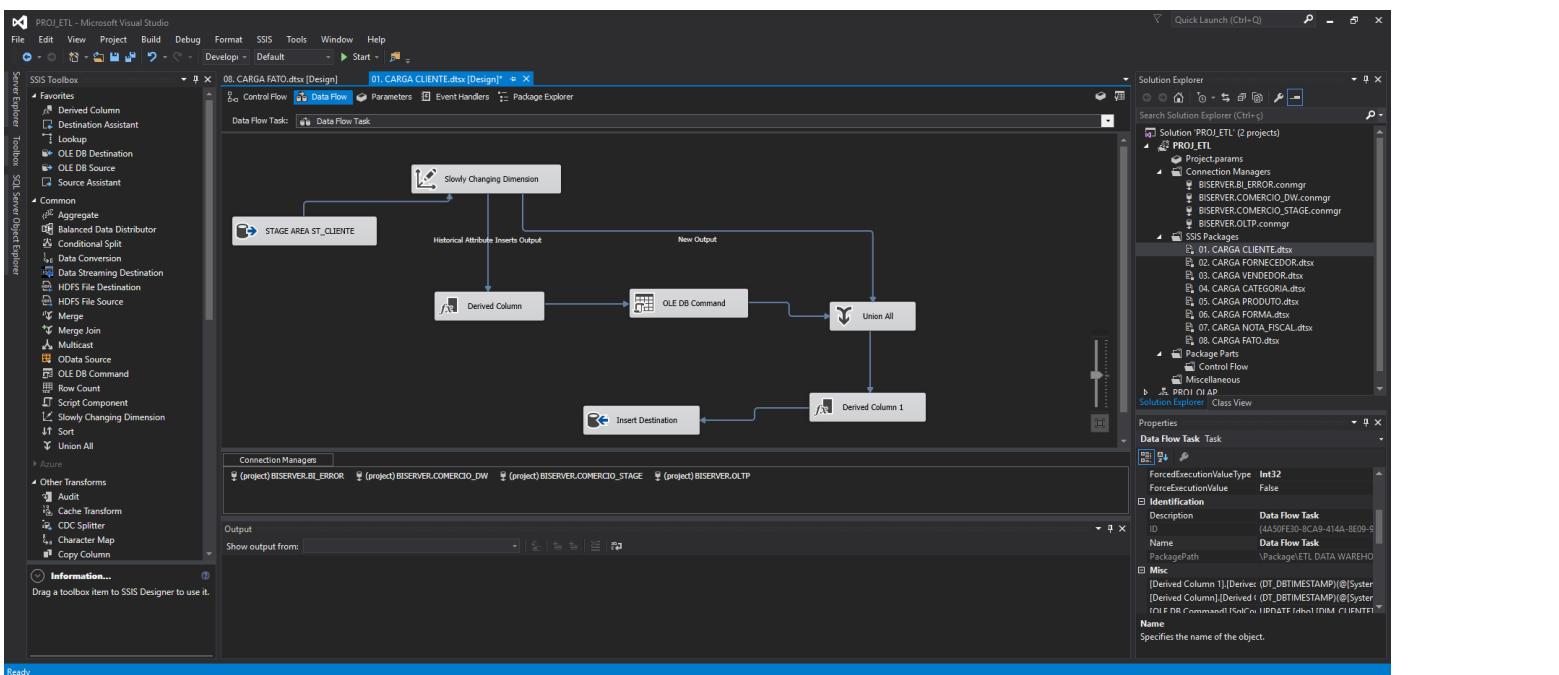
Como mostra a *screenshot* abaixo, os processos foram separados por uma ferramenta chamada ***Sequence Container***. Essa ferramenta é capaz de separar os processos de ETL da área de Stage, que é feita no primeiro momento, para então, depois de finalizado, fazer a carga do *Data Warehouse*.

Este processo faz todo sentido uma vez que a *Stage Area* serve apenas para carga e transformações necessárias dos dados. Então, ao finalizar a carga da *Stage* pelo pacote ‘ETL_STAGE’, inicia-se o container seguinte, ‘ETL DATA WAREHOUSE’, com a carga do DW.

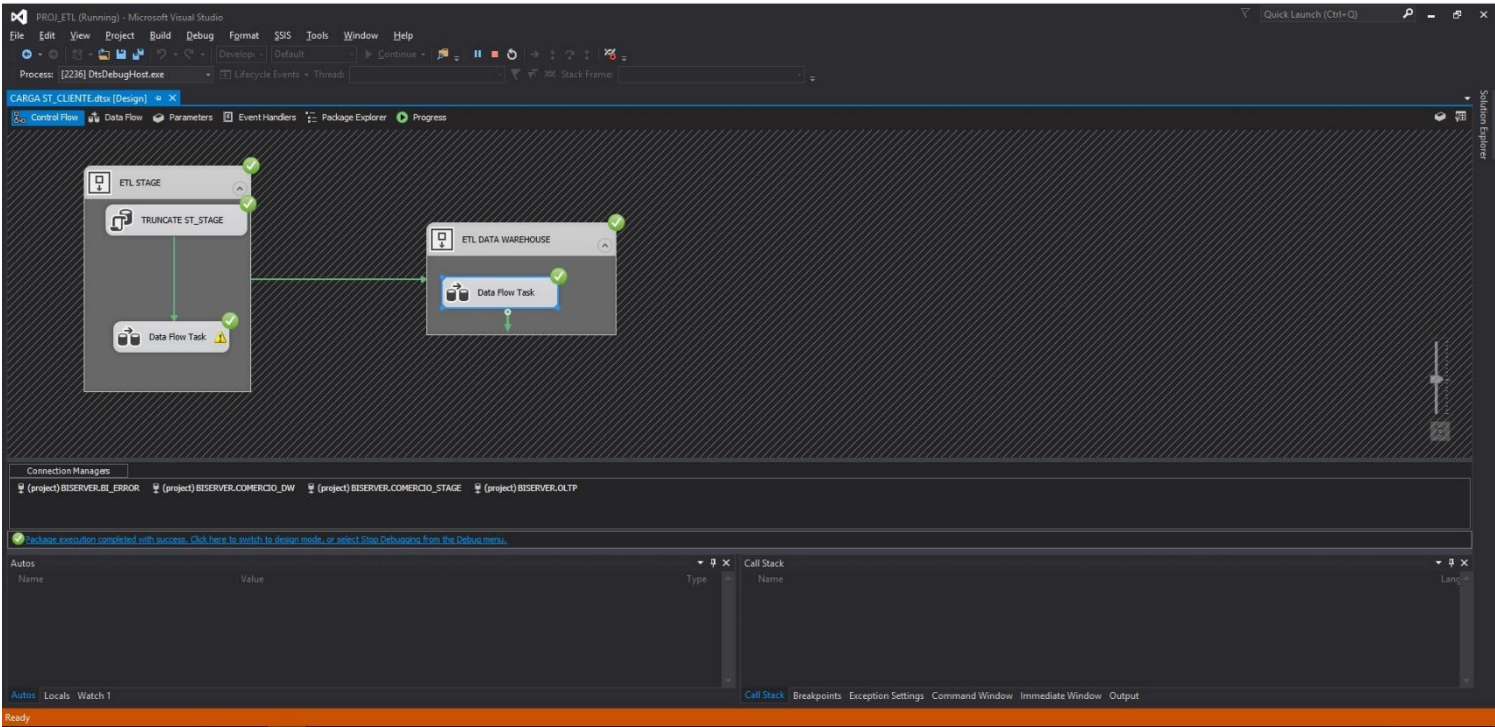


O processo interno de ETL do pacote do *Data Warehouse*, ou ***Data Flow***, possui um processo de carga um pouco diferente das demais, uma vez que esta carga se baseia em SCD (Slowly Changing Dimensions ou dimensões mutáveis) para garantir o armazenamento do histórico.

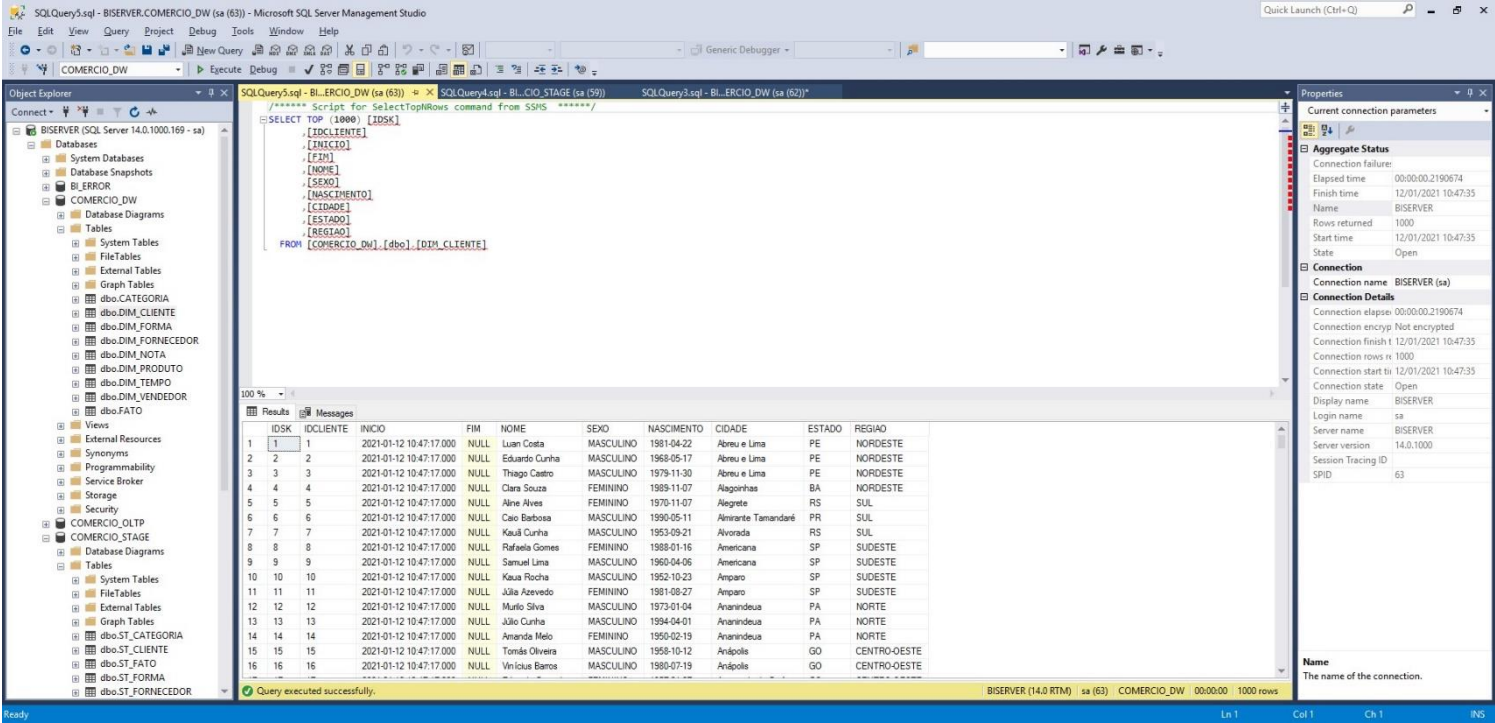
É adicionado uma ferramenta ao pacote chamada ‘Slowly Changing Dimensions’, encontrada no menu de transformações comuns do SSIS. O assistente da ferramenta define as opções para a tabela e cria automaticamente o *Data Flow*, já conectado ao source de dados (‘STAGE AREA ST_CLIENTE’).



Ao finalizar o processo de ETL, podemos ver a confirmação de êxito do processo, como na *screenshot* abaixo:



Ao finalizar o pacote do SSIS, podemos consultar no SQL Server a carga dos dados, realizando um SELECT simples na tabela DIM_CLIENTE



Nesse ponto, podemos ver o *Data Warehouse* já disponível para queries e conexão com as ferramentas de BI. Claro, após a carga de todas as outras tabelas.

O processo de ETL demonstrado aqui se repete igualmente para todas as outras dimensões. Portanto, não será demonstrado aqui todas as tabelas, uma vez que o processo é igual.

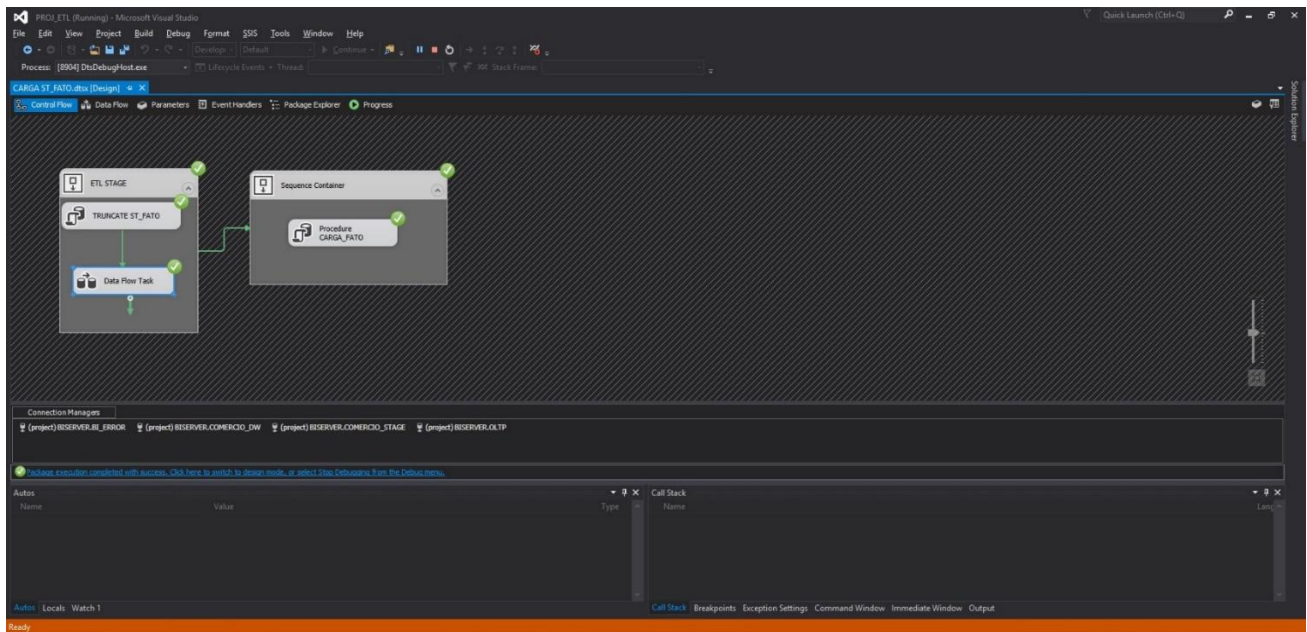
CARGA DA TABELA DIM_FATO

Por final, carregamos a tabela fato, que é a consolidação do Data Warehouse.

A tabela DIM_FATO é carregada utilizando um script em SQL que é capaz de filtrar os dados já existentes, comparando as entradas já existentes no DW com os novos resultados. Dessa forma, o ETL carrega apenas os novos dados inseridos no sistema OLTP.

Por motivos de privacidade, o script que foi fornecido pelo instrutor Felipe Mafra, não será descrito aqui nesse projeto. Entretanto, este pode ser obtido ao adquirir um de seus cursos na plataforma Udemy. O mesmo motivo é válido também para a DIM_TEMPO, que foi criada separadamente do SQL Server por um script de inteligência de data.

Os scripts, de forma genial, podem ser replicados em outros projetos de forma inteiramente funcional.



Após a conclusão do processo de carga da DIM_FATO, podemos verificar os dados através de uma query de consulta.

IDNOTA	IDCLIENTE	IDVENDEDOR	IDFORMA	IDPRODUTO	IDFORNECEDOR	IDTEMPO	QUANTIDADE	TOTAL_ITEM	CUSTO_TOTAL	LUCRO_TOTAL
1	19	778	6	9	58	74684	1	1150.00	900.00	250.00
2	20	251	21	9	111	74706	2	2600.00	1400.00	1200.00
3	20	251	21	9	65	74706	3	300.00	180.00	120.00
4	22	859	6	7	65	74674	3	300.00	180.00	120.00
5	24	43	19	4	121	74735	2	9600.00	5998.00	3602.00
6	26	450	1	7	134	74782	1	90.00	38.00	52.00
7	28	504	2	25	124	74534	2	3000.00	1600.00	1400.00
8	28	504	2	25	3	74534	3	144.00	72.00	72.00
9	30	958	1	11	117	74664	1	3300.00	2700.00	600.00
10	32	453	5	23	95	74797	3	150.00	75.00	75.00
11	32	453	5	23	194	74797	4	348.00	288.00	60.00
12	34	34	3	21	85	74806	2	5000.00	2600.00	2400.00
13	34	34	3	21	70	74806	1	890.00	500.00	390.00
14	36	171	24	9	13	74539	1	189.00	60.00	129.00
15	36	171	24	9	66	74539	3	1680.00	1155.00	525.00
16	38	944	21	6	122	74735	3	10500.00	8400.00	2100.00
17	40	956	18	9	25	74726	1	88.00	46.00	52.00
18	42	472	11	4	103	74565	3	1650.00	600.00	1050.00
19	44	976	3	15	71	74525	1	500.00	360.00	140.00
20	46	471	4	7	39	74572	1	145.00	78.00	67.00
21	48	283	13	13	116	74598	3	5610.00	3900.00	1710.00
22	52	13	11	13	3	74479	2	96.00	48.00	48.00
23	52	1000	21	26	165	74628	3	3300.00	2550.00	750.00
24	54	155	18	3	43	74493	2	290.00	134.00	156.00
25	56	181	1	4	158	74777	3	327.00	243.00	84.00

CONCLUSÃO DATA WAREHOUSE

O projeto de modelagem do *Data Warehouse* finalmente conclui-se com êxito em todos os processo de ETL e definição de relacionamentos.

O DW está pronto para ser consultados por queries em SQL e/ou receber conexões diretamente de ferramentas de *Data Visualization* para construção de dashboards.

Entretanto, nesse ponto, inicia-se os projetos do sistema OLAP, utilizando o *Analysis Services* (SSAS) e o *Reporting Services* (SSRS) para criar os cubos analíticos e automação dos relatórios.