**Buiild an etl pipline with azure synapse with dataflow running on it**

Read me: page 1 to 5 – steps to do. Page 6 to 8 Conclusion

Create Synapse workspace



Go to develop tab and create dataflow

Add source



Give filter condition for data transformation and sort for data transformation

Coding challenge on DEVOPS Thursday, December 19, 2024

Validate and publish



Create pipeline drag and drop dataflow from move and transform and the created dataflow
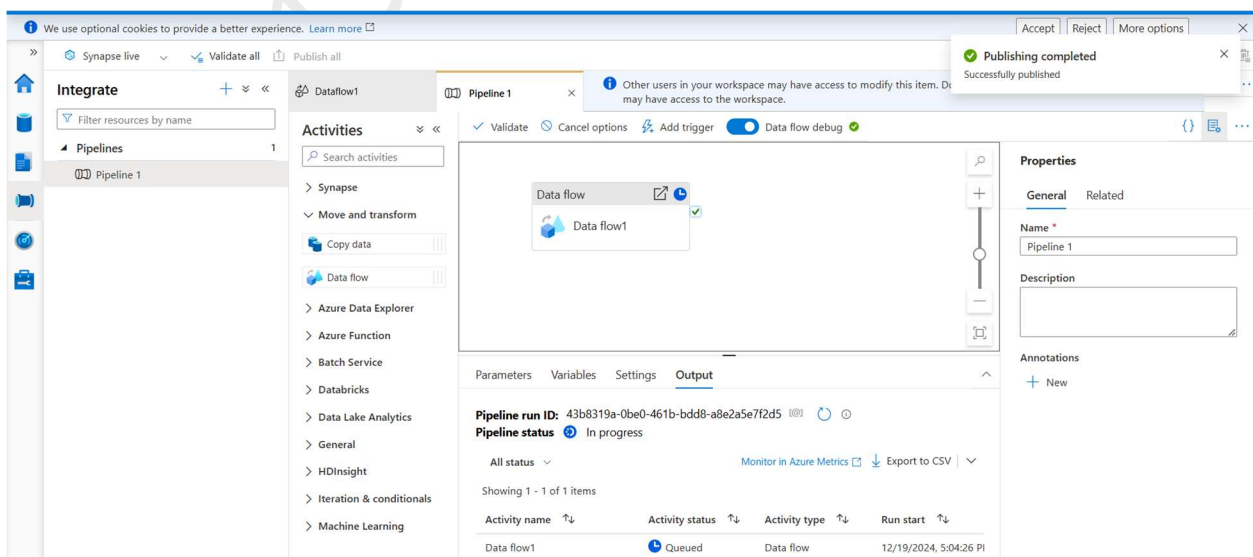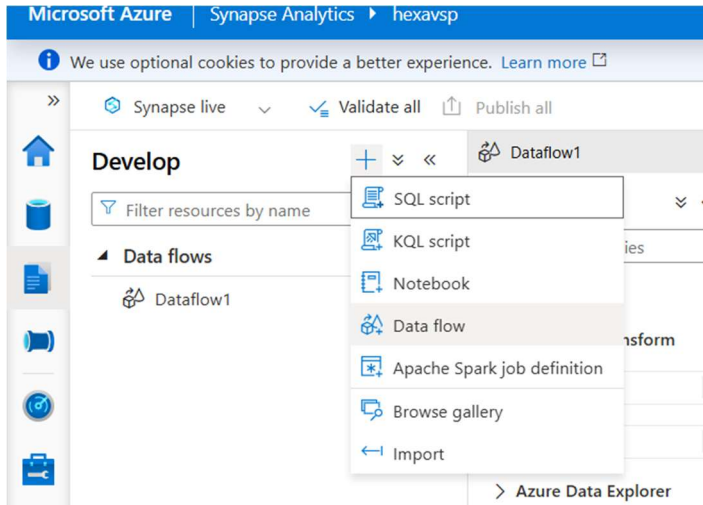


Click on debug to run the pipeline

Coding challenge on DEVOPS Thursday, December 19, 2024
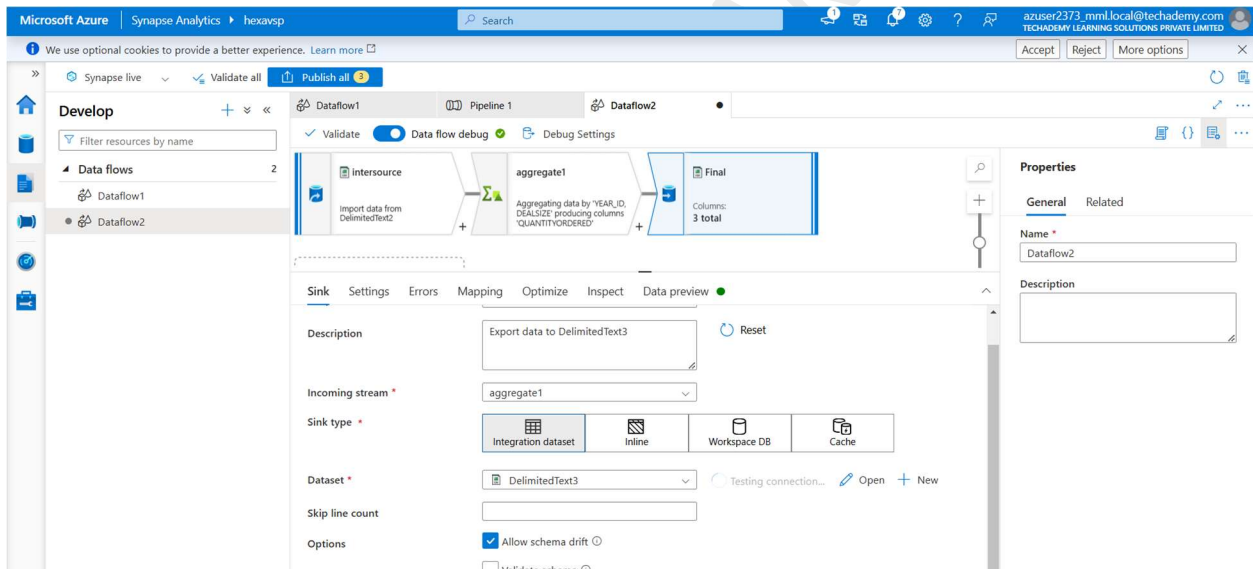
Create new data flow to see the final results



Aggregate it by year and deal size and calculate the quantity ordered.



Add this dataflow in pipeline and run the pipeline to see the results.

Pipeline started



Sort and filter dataflow succeeded now aggregation dataflow is running.



Both ran successfully

Coding challenge on DEVOPS Thursday, December 19, 2024

## CONCLUSION

Input data in datalake



Intermediate result data sorted and filtered in datalake



Result – Aggregated values

Coding challenge on DEVOPS Thursday, December 19, 2024

## DATA ENGINEERING DEVOPS

### DATAFLOW 1



### DATAFLOW2



### PIPELINE

Coding challenge on DEVOPS Thursday, December 19, 2024

## DATA ENGINEERING DEVOPS

MONITORING PIPELINE RUNS

Coding challenge on DEVOPS Thursday, December 19, 2024