

## Goals

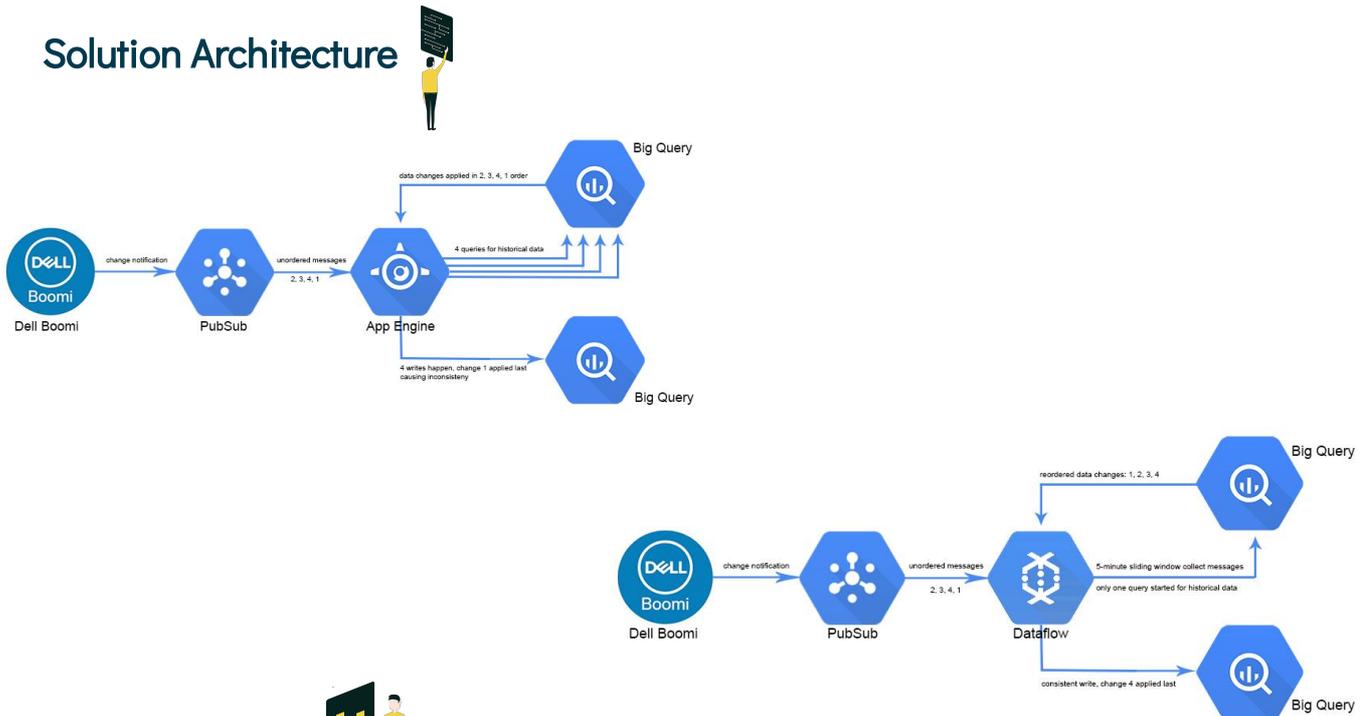


Dagrofa is Denmark's largest food wholesaler with 20 % market share. Dagrofa has 14000 employees and it is behind two successful department store chains, more than 500 grocery stores and 450 self-employed merchants.

To operate successfully on a large scale Dagrofa needs to make informed, data-driven decisions and for that constant market analysis. The company puts great value on big data analytics that allows them to adjust their strategies based on actionable insights. In the previous phases of the project a data warehouse was built for Dagrofa. The solution already supported analysis of historical and current state data and it allowed a wide variety of stream and batch jobs. But the existing system was also a basis for improvements to be even more aligned with Dagrofa's goals.

<b>Challenge:</b> <b>Out of order message processing</b>	<b>Solution:</b> <b>efficient, consistent message processing with time windows</b>
<p>In the previous solution a <b>Dell Boomi data preparation</b> service pushed a large amount of messages to servlets. The data was processed, deserialized and validated on <b>Google App Engine</b>. The target Big Query table was queried to reevaluate history, then came postprocessing and finally the results were <b>written to BQ</b>. The solution had multiple weaknesses.</p> <ul style="list-style-type: none"> <li>• Large spikes: GAE handled badly when thousands of messages arrived within minutes</li> <li>• Out of order data: some changes came in quick succession</li> <li>• Concurrency: messages were lost when more arrived at the same time</li> <li>• High cost: every single message started BQ queries</li> </ul> <p>Dagrofa wanted to get consistent and correct data despite parallelism, while lowering BQ costs by skipping the unnecessary queries.</p>	<p>In response to the problems of the customer the Aluz team made a plan to replace the existing system with a new, better one within 6 months by the following steps.</p> <ul style="list-style-type: none"> <li>• Collect issues from the client caused by the insufficient previous solution, specify requirements and estimate cost reduction</li> <li>• Replace GAE processing with a <b>Dataflow</b> pipeline</li> <li>• Batch messages before processing, reorder them within a <b>5-minute sliding window</b></li> <li>• Cut expenses: reduce query cost by querying once per window, and by reducing the duplications in the resulting data</li> <li>• Optimize post processing and BQ writes</li> <li>• Start to use new solution for 25 different type of entities</li> </ul>

## Solution Architecture



## Business value



With the newly implemented solution Dagrofa's data became more reliable. Right now if the relative data latency is under 6 hours, then unordered executions cause no inconsistencies. Dataflow also allows to avoid problems caused by splitting messages between multiple paths or parallel executions. The results remain consistent even when a message and a user modify the same data at the same time - previously this resulted in one of the changes being lost.

The increased correctness of the most important data types benefited the client's decision making by more accurate comparisons, trend analysis and as a stable, valid source for consumers of the data, it serves as the basis for automated analytical solutions.

One of the most valuable data types for Dagrofa describes their items in various stores and how they are handled. The store item type was in the highlight of attention during the entire development process. The item data often arrives in large spikes containing ten thousands of messages, while the changes most frequently affect only the item range column. In this case it was very important to apply the messages in the correct order and to keep the item history consistent because complex post processing follows the inserts that prioritizes the store items.

Another major factor was the significant cost reduction. By lowering the number of queries and duplicate data while also optimizing the processing the solution cut costs by 35%. Beside the obvious savings the maintenance of the system became easier. The unified process will make debugging faster, and it is also easier to introduce further changes and improvements to the system.