

# Data Warehouse or Data Lake: Which one is right for you?

Traditional Data Warehouse		Data Lake
Structured	<b>Types of data</b>	Structured, semi-structured, unstructured, and raw
Expensive for large volumes of data	<b>Cost effectiveness</b>	Designed for low cost (typically on commodity hardware)
Business users	<b>Intended users</b>	Data scientists and advanced users
Business Intelligence, reporting, dashboards and visualizations	<b>Purpose</b>	Not always predefined, but may be leveraged for machine learning and AI
Expensive and time consuming to change due to schema design implications	<b>Flexibility</b>	More flexible and easier to update due to lack of rigid schema definition, but may introduce challenges due to number of small files or complexity of Hadoop File System (HDFS)
Data is processed and organized into a single schema before storing in the data warehouse	<b>Ingestion process</b>	Raw and unstructured data can be directly ingested into the data lake
Data is already structured and can be consumed by business intelligence tools	<b>Analysis process</b>	Data is organized on demand at the time of analysis. Schema on read rather than write. Datalake is horizontally scalable.
Centralized	<b>Data storage</b>	Centralized
Relational (e.g. SQL)	<b>Schema type</b>	Undefined schema
Data is collected from multiple relational sources	<b>Data sources</b>	Collect data from any source and in any format and provide near real-time ingestion and streaming pipelines
Medium	<b>Data volume</b>	High
Established definitions of data structures and schemas	<b>Data definition</b>	Requires data to be cleaned, enriched and joined with metadata for consumption and or feeding into Data Warehouse