



# Power BI on Databricks Best Practices Cheat Sheet

This document summarizes the most relevant best practices when using Power BI

PHASE	DATA PREPARATION	SQL SERVING	POWER BI INTEGRATION	POWER BI REPORT DESIGN
Read if you are	Preparing the dataset that will be served for the report	Setting up and configuring the Databricks SQL Warehouse to which Power BI connects	Configuring the integration between Power BI and Databricks	Designing semantic models and reports in Power BI
Your objective	<ul style="list-style-type: none"> <li>Implement an efficient data model</li> <li>Optimize storage layer for query performance</li> </ul>	<ul style="list-style-type: none"> <li>Meet performance SLA and scalability</li> <li>Keep costs under control</li> </ul>	<ul style="list-style-type: none"> <li>Optimize semantic model for performance</li> <li>Leverage unified governance</li> </ul>	<ul style="list-style-type: none"> <li>Create fast and efficient reports and dashboards</li> <li>Inherit the optimizations performed in all the other layers</li> </ul>
Best practices: What you should do	<ul style="list-style-type: none"> <li>Adopt <b>medallion architecture</b> on <b>Delta Lake</b> (<a href="#">link</a>) and <b>serve the Gold layer only</b> (<a href="#">link</a>) → leverage benefits of Delta and aggregated data for better report performance.</li> <li>When designing your data model, opt if possible for a <b>star schema</b> (<a href="#">link</a>, <a href="#">link</a>) → Power BI performs better.</li> <li>Leverage <b>SQL views</b> and <b>persisted tables</b> with the required granularity (e.g., date/time) → improve performance by pre-aggregating data for the most common and resource-intensive queries.</li> <li>Declare <b>primary and foreign keys</b> (<a href="#">link</a>), and use <b>RELY</b> (<a href="#">link</a>) where possible → Power BI can leverage this information to automatically create table relationships, and the Databricks SQL engine can optimize queries using PK constraints.</li> <li><b>Avoid using wide data types and high-cardinality</b> columns → reduce Power BI semantic model size and improve query performance.</li> <li>Use <b>auto-generated columns</b> when you need to generate a derived value from other existing columns (<a href="#">link</a>) → persisted columns minimize the need for calculating values at query time.</li> <li>Use <b>Liquid Clustering</b> (<a href="#">link</a>); alternatively, use Z-ordering (<a href="#">link</a>) and consider partitioning larger tables &gt;1TB (<a href="#">link</a>) → improve query performance by efficient file and data skipping.</li> <li><b>Optimize your data layout</b> by using predictive optimization (<a href="#">link</a>) or running <b>VACUUM</b> (<a href="#">link</a>) and <b>OPTIMIZE</b> (<a href="#">link</a>) → improve performance by deleting old files and optimizing the physical data layout.</li> <li>Periodically <b>compute statistics</b> (<a href="#">link</a>) or use <b>automatic statistics</b> (currently in preview, <a href="#">link</a>) → improve performance by choosing a more efficient join strategy.</li> <li>Evaluate using <b>materialized views</b> (<a href="#">link</a>) to calculate results incrementally using the latest data in the source tables → improve performance by leveraging precomputed results.</li> </ul>	<ul style="list-style-type: none"> <li>Always use a <b>SQL warehouse</b> (and not all-purpose clusters) (<a href="#">link</a>) → SQL warehouses are optimized for BI workloads.</li> <li>Use SQL Serverless warehouse (<a href="#">link</a>) → optimal price/performance and delivers instant and elastic compute. Leverage Serverless query result cache even when SQL Warehouse is scaled or restarted.</li> <li>Enable SQL warehouse <b>Auto stop</b> (scale-to-zero) only if the SLA permits; the SQL warehouse will take just 5-10s to start up when the first query arrives → save costs when the warehouse is not used.</li> <li>Use <b>higher cluster size for larger datasets</b> (<a href="#">link</a>) → the larger the cluster (M, L, XL, etc.), the faster complex queries run. If having only simple, short-running queries, don't increase the size (might be slower due to data shuffling).</li> <li>Use <b>SQL Warehouse scaling</b> (min/max cluster count) (<a href="#">link</a>) → <b>handle more concurrent users / queries</b>. SQL Warehouse scales out to handle the increased workload. When hitting limits, queries get queued, not rejected.</li> <li>If expecting many concurrent queries, <b>increase the minimum number of clusters</b> → prevent queuing queries waiting for scaling out.</li> <li>Use the <b>same SQL Warehouse</b> whenever the same dataset is queried → will leverage the various layers of caching available, increasing performance.</li> <li>Use <b>separate SQL warehouses</b> for different workloads and/or business units → right-size SQL warehouses to achieve better performance and reasonable costs.</li> <li>If unsure about sizing, start with a <b>Medium SQL Warehouse</b> scaling between 1 and 10; <b>monitor query response time and scaling</b> → then adjust the sizing based on the observed results.</li> <li><b>Leverage</b> the Query History (<a href="#">link</a>) and SQL warehouse events (<a href="#">link</a>) <b>system tables</b> to programmatically monitor SQL Warehouses and query performance → allows for identifying performance bottlenecks and issues, implementing more detailed analysis and setting up alerts.</li> </ul>	<ul style="list-style-type: none"> <li>Ensure that Power BI and Databricks are hosted as closely as possible, ideally in <b>the same region</b> → this would minimize network latency and may help avoid cross-region traffic costs.</li> <li>Use the most appropriate Power BI storage mode: <b>Direct Query</b> for Fact tables, <b>Dual</b> for Dimension tables (NOT Import) (<a href="#">link</a>) → let Power BI generate more efficient SQL queries.</li> <li>Evaluate where and how to use <b>composite models</b> (<a href="#">link</a>) → allows a mixed usage of DirectQuery, Dual, Import mode tables, as well as Aggregation and Hybrid tables.</li> <li>Use <b>hybrid tables</b> whenever you need aggregated historical data augmented with detailed real-time data for the same table (<a href="#">link</a>) → efficient and quick in-memory queries combined with the latest data changes directly from the source.</li> <li>In Import mode, use <b>table partitioning</b> (<a href="#">link</a>). Alternatively use <b>incremental refresh</b> (<a href="#">link</a>) → allows importing data faster and managing larger datasets.</li> <li>Check for <b>query parallelization</b> configuration settings (<a href="#">link</a>) → improves query parallelization and maximizes utilization of SQL warehouse to improve overall performance.</li> <li>Connect Power BI to Databricks using <b>single sign-on (SSO)</b> (<a href="#">link</a>) → allows leveraging security and governance controls implemented in Databricks <b>Unity Catalog</b> (<a href="#">link</a>) and allows audit data access.</li> <li>If you need to connect to different Databricks environments, use <b>Power BI parameters</b> (<a href="#">link</a>) → allows flexibility when connecting to different Databricks workspaces or different Databricks SQL warehouses.</li> <li>Use <b>gateway clusters</b> (<a href="#">link</a>) to connect to IP ACL or private link-secured Databricks workspaces → avoid single points of failure and load balance traffic across gateways in a cluster.</li> <li>Use <b>Publish to Power BI Service</b> (<a href="#">link</a>) → enables seamless catalog integration and data model sync, allowing you to publish datasets directly to Power BI Service without leaving the Databricks UI.</li> <li>Use <b>Automatic Publishing to Power BI</b> (<a href="#">link</a>) → publish datasets from Unity Catalog to Power BI directly from data pipelines.</li> </ul>	<ul style="list-style-type: none"> <li><b>Limit the number of visuals</b> on each report page → limit the number of queries that will be executed.</li> <li><b>Limit the number of rows and columns</b> in semantic models and report visuals → avoid large data transfers.</li> <li>Leverage <b>user-defined aggregations</b> (<a href="#">link</a>) → improve query performance over large DirectQuery semantic models by caching pre-aggregated data.</li> <li>Use <b>automatic aggregations</b> (<a href="#">link</a>) → continuously optimize DirectQuery semantic models by building aggregations based on Query History for maximum report performance.</li> <li>If referential integrity has been validated in the upstream ingestion use <b>"Assume Referential Integrity"</b> when defining table relations (<a href="#">link</a>) → more efficient join strategies in SQL queries.</li> <li><b>Avoid "many-to-many"</b> relationships where possible → decrease complexity and improve Power BI model efficiency.</li> <li>Configure <b>"Is nullable"</b> for table columns where applicable → Power BI generates simpler and more efficient SQL queries.</li> <li><b>"Move left" transformations</b> whenever possible (e.g., prefer SQL views over PowerQuery transformations and DAX formulas) → leverage the power of Databricks SQL engine for more efficient report execution.</li> <li>If using DAX, review your code to use <b>efficient DAX calculations</b> (<a href="#">link</a>) → inefficient calculations can lead to deteriorated performance.</li> <li>Leverage <b>query reduction settings</b> by adding Apply/Clear All Slicers button (<a href="#">link</a>) → prevent a new query from being sent to the data source every time the user interacts with the report's filters.</li> <li><b>Avoid using DAX calculated columns and calculated tables</b> in semantic models → will perform better if defined directly in your Gold tables. Measures that can be precomputed as columns also perform best if done in the Gold layer.</li> </ul>
Troubleshooting: Why is my report slow?	<ol style="list-style-type: none"> <li>Ensure the <b>data layout is regularly optimized</b> with VACUUM and OPTIMIZE.</li> <li>Evaluate generating <b>aggregated views</b> for the most common and resource-intensive queries whenever applicable.</li> <li>In the SQL Warehouse/Monitoring page (<a href="#">link</a>), review the Query History looking for the queries with the longest duration:             <ol style="list-style-type: none"> <li>Ensure the query leverages the way the storage layer has been optimized or adapt it if needed.</li> <li>Open the <b>Query Profile</b>, enable verbose mode, review the tasks that took the most time and/or memory looking for:                 <ol style="list-style-type: none"> <li>"Cloud storage request duration" to check if the cloud storage has been slow in responding.</li> <li>"Files read," "Size of the smallest file read" ensuring not too many small files are read. If so ensure OPTIMIZE is run regularly.</li> <li>"Number of output rows" and "Rows skipped" ensure filters are applied sooner than later.</li> </ol> </li> </ol> </li> </ol>	<ol style="list-style-type: none"> <li>Monitor the SQL Warehouse performance in the SQL Warehouse/Monitoring page (<a href="#">link</a>) looking for:             <ol style="list-style-type: none"> <li><b>Running queries vs. Queued queries.</b> If too many Queued queries: 1) increase the number of clusters if already reaching the maximum allowed, 2) evaluate if worth increasing the cluster size if there are long-running queries.</li> <li>Number of <b>active clusters</b>: 1) evaluate to set a higher minimum number of clusters if queries are queued while scaling out, 2) check if the cluster scaled in to zero when the queries arrived.</li> </ol> </li> <li>Monitor <b>Query History</b> in the SQL Warehouse/Monitoring page (<a href="#">link</a>), reviewing query details for:             <ol style="list-style-type: none"> <li>"Scheduling" (i.e., time spent in queue) vs. "Optimizing Query" (i.e., time spent mainly identifying data to be skipped) vs. "Executing" (i.e., time spent executing the query).</li> <li>"Result fetching by client" (i.e., time spent for the client to download the result set).</li> <li>"Rows returned" to ensure the query does not return too many rows.</li> <li>"Bytes read from cache" to evaluate the disk cache efficiency.</li> <li>"Bytes spilled to disk" to ensure no data is spilled to disk; if so, increase the cluster size.</li> </ol> </li> </ol>	<ol style="list-style-type: none"> <li>Evaluate which Power BI storage mode is used. Prefer <b>DirectQuery for Fact tables</b> and <b>Dual for Dimensions</b> → ideal compromise between performance and freshness of the data.</li> <li>For very <b>small, static and performance-sensitive</b> (&lt;2s) reports, evaluate the usage of <b>Import</b> mode → would provide the best report performance.</li> <li>Especially for DirectQuery, check how many queries Power BI can send <b>in parallel</b> to Databricks. Ensure the Databricks <b>SQL warehouse is sized accordingly</b> to handle the required level of parallelism → avoids queries to be queued, resulting in a slow report.</li> <li>For performance fine-tuning, evaluate the following properties of Power BI semantic models:             <ol style="list-style-type: none"> <li>"Maximum connections per data source"</li> <li>"Maximum number of simultaneous evaluations"</li> <li>"Maximum number of concurrent jobs"</li> <li>"MaxParallelismPerQuery"</li> </ol> </li> <li>Monitor the queries in the SQL Warehouse/Monitoring page, looking at "Started at" and at which <b>time the queries arrived</b> to validate the effective parallelism in Power BI.</li> </ol>	<ol style="list-style-type: none"> <li>Use <b>Power BI Performance Analyzer</b> to examine report element performance (<a href="#">link</a>) → identify the visual that takes the most to load and where the bottleneck is (DAX Query, Visual Display, DirectQuery, etc.).</li> <li>Ensure there are <b>not too many visuals</b> in the same report → many visuals could generate lots of queries that can be queued by Power BI or Databricks (e.g., the query itself runs fast but spends time in the queue).</li> <li>For the most common and resource-intensive queries, <b>ask for creating SQL views or persisted tables</b> in the Gold layer, which provides <b>pre-aggregated data</b> (often valid for Date dimension) → results in overall better performance.</li> <li>Ensure there are no SQL queries returning large result sets (1,000s of records), which is often an indication of inefficient DAX formulas (e.g., TOPN function) → evaluate the complexity of DAX formulas and optimize where possible.</li> </ol>
What you should read	<ul style="list-style-type: none"> <li>Power Up With Power BI and Lakehouse in Azure Databricks: Part 3 — Tuning Azure Databricks SQL (<a href="#">link</a>)</li> <li>Star Schema Data Modeling Best Practices on Databricks SQL (<a href="#">link</a>)</li> <li>One Big Table vs. Dimensional Modeling on Databricks SQL (<a href="#">link</a>)</li> </ul>	<ul style="list-style-type: none"> <li>SQL Warehouse Sizing, Scaling and Queuing Behavior (<a href="#">link</a>)</li> <li>DBSQL Warehouse Advisor (<a href="#">link</a>)</li> <li>Tune Query Performance in Databricks SQL With the Query Profile (<a href="#">link</a>)</li> </ul>	<ul style="list-style-type: none"> <li>Power Up Your BI With Microsoft Power BI and Lakehouse in Azure Databricks: Part 1 — Essentials (<a href="#">link</a>)</li> <li>Power BI — Databricks SQL QuickStart Samples (<a href="#">link</a>)</li> <li>Power BI on DBSQL Design Patterns (<a href="#">link</a>)</li> <li>Boosting Power BI Performance With Azure Databricks Through Automatic Aggregations (<a href="#">link</a>)</li> </ul>	<ul style="list-style-type: none"> <li>Power Up Your BI With Microsoft Power BI and Lakehouse in Azure Databricks: Part 2 — Tuning Power BI (<a href="#">link</a>)</li> <li>Optimization Guide for Power BI (<a href="#">link</a>)</li> <li>Monitor Report Performance in Power BI (<a href="#">link</a>)</li> <li>Troubleshoot Report Performance in Power BI (<a href="#">link</a>)</li> </ul>