

# M2 release notes & upcoming releases

## M2 TLDR

### Overall

1. **MLflow Integration:** Most `rag_studio` SDKs are now part of `mlflow`: logging, evaluation, chain parameterization. MLflow Tracing integration is improved.
2. **Developer workflow improvements:** Simplified chain configuration that accepts Dictionaries and YAML.

### Evaluation & Databricks LLM Judges

3. **Evaluation is faster & easier to use:** Faster evaluation with progress bar. Accepts single Pandas Dataframe as input and provides single output Dataframe with results. Metrics & data logged to MLflow vs. Delta Tables.
4. **Improved LLM Judges & Metrics docs:** Clarity in metric names and how metrics evaluate chains
5. **New metrics:** Added new metrics such as token count and latency
6. **LLM Judge Tuning:** Few-shot examples can enhance LLM Judge agreement with human ratings
7. **Customer-defined LLM Judges:** Custom judges can be defined to assess use-case specific criteria.

### Logging & monitoring

8. **Human Feedback and Monitoring:** New Delta Tables that automatically ETL feedback collected via Review App and Feedback API

### Chain compatibility

9. **Chains support External Models & Provisioned Throughput:** Chains can use any Model Serving model, including OpenAI and Provisioned Throughput models with automatic credential provisioning.
10. **Chains support Third-Party APIs:** Chains can utilize third-party APIs with secure credential management through Databricks Secrets.

### Review App

11. **Review App SSO:** Review App is accessible to any SSO users even if they aren't a Databricks user

### Data Pipelines

12. **Flexible data pipeline:** New data processing pipeline template supports easily testing out different parsing, chunking, embedding strategies - several off the shelf implementations are included.

## Details and upcoming releases

Theme	M2 release	Upcoming releases
Integration with MLflow	<p>RAG Studio's integration with MLflow is significantly improved for both chain developing / logging and chain evaluation.</p> <p>Evaluation:</p> <ul style="list-style-type: none"><li>• <code>rag_eval.evaluate(...)</code> is now called via <code>mlflow.evaluate(..., model_type="databricks-rag")</code></li><li>• Evaluation metrics &amp; evaluation data are logged to MLflow Runs (versus a Delta Table)</li></ul>	Customer defined LLM judges configuration will be moved from YAML into the native MLflow SDK

	<p>Chain development:</p> <ul style="list-style-type: none"> <li>• <code>rag.log_model(...)</code> is merged into <code>mlflow.langchain.log_model(...)</code> SDK</li> <li>• Chain configuration via YAML is now part of MLflow e.g., <code>rag.RagConfig(...)</code> moved to <code>mlflow.models.ModelConfig(...)</code></li> </ul>	
One-click RAG Data Pipeline accelerator	<p>Comes with various pre-built parsing / chunking / embedding strategies or build your own.</p> <p>100% of the code is available to you, so this accelerator can be fully customized</p>	Improvements based on your feedback
Faster and easier to use evaluation capabilities	<p>Beyond tighter integration with MLflow, Evaluation now:</p> <ul style="list-style-type: none"> <li>• Evaluation runs faster</li> <li>• Evaluation status displayed via a progress bar w/ ETA</li> <li>• <b>Breaking change</b> Improved schema for evaluation input and output</li> <li>• Accepts a single Pandas Dataframe as input vs. requiring 2 separate Delta Tables</li> <li>• Evaluation results are available as a single output table vs. split across assessments &amp; metrics, making review of evaluation easier <ul style="list-style-type: none"> <li>◦ Returned as Pandas Dataframes vs. saved to Delta Tables</li> </ul> </li> <li>• Evaluation Set schema for retrieval ground truth is more flexible and does not require chunk text</li> <li>• <b>Breaking change</b> the config YML syntax has changed to accommodate the new functionality of custom judges and few-shot examples</li> </ul> <p><a href="#">Documentation</a></p>	Native MLflow-SDK approach for defining custom LLM judges and metrics.
Improved integration with MLflow Tracing	<p>MLflow Tracing provides a way to instrument and visualize the steps in your chain.</p> <p>In this release, MLflow Tracing is more tightly integrated with RAG Studio:</p> <ul style="list-style-type: none"> <li>• When running your chain during local development, traces are logged to MLflow and visualized in a Notebook-based UI</li> <li>• The data from the trace can be accessed via the MLflow SDK</li> </ul>	<p><b>Breaking change</b> Step-by-step chain traces are stored in the same MLflow tracing format used during development. This replaces the current trace schema.</p> <p>Traces collected online can be viewed using the same MLflow UI as in development.</p>

Support for External Models & Provisioned Throughput	<p>Your chain can now use any External Model e.g., OpenAI, etc and any Provisioned Throughput model. Credentials for these models are automatically provisioned when calling <code>deploy_model()</code>.</p> <p>Previously, only FM API Pay-Per-Token models were available.</p>	Support
Support for 3rd party APIs	<p>Your chain can now use 3rd party API services e.g., non-Databricks vector databases, etc with secrets stored in Databricks Secrets Manager via environment variables.</p> <p>Note: Doing so requires manually updating the Model Serving endpoint after calling <code>deploy_model()</code>.</p>	Native support for using Databricks Secrets via <code>deploy_model()</code> .
Fast chain deployment		If using a Databricks defined set of Python packages in your chain, chains now deploy to REST APIs and Review App in 2 - 3 minutes vs. 15 - 20 minutes.
Chain token streaming support		<p>If your LangChain chain supports streaming, your deployed REST API and Review App will also support streaming.</p> <p>Support for token streaming in PyFunc based chains</p>
Human feedback collection & online chain monitoring	<p>Chat requests to a deployed chain's REST API and chat requests via the Review App are available as Delta Tables.</p> <p>These tables are ETL'd versions of the raw Inference Table that contains the same data, but has a raw schema and duplicate rows.</p> <p>2 tables are generated:</p> <ul style="list-style-type: none"> <li>• <code>request_log</code>: Keyed to request_id, contains one row for every trace. A trace represents one turn of conversation in a chat and is linked to other traces in that conversation via conversation_id. Contains: <ul style="list-style-type: none"> <li>○ Request</li> <li>○ Response</li> <li>○ Step-by-step chain trace</li> </ul> </li> <li>• <code>assessment_log</code>: Keyed to request_id, contains one row for every piece of human feedback collected via the Review App or the Feedback API, including</li> </ul>	<p><b>Breaking change</b> Improvements to the assessment log schema to enable the data to be more easily used as input to tune the LLM Judge models and create Evaluation Sets.</p>

	<p>feedback for live chats &amp; offline review of traces. Contains:</p> <ul style="list-style-type: none"> <li>○ Thumbs up / down + why <ul style="list-style-type: none"> <li>■ Free text &amp; select-many choices</li> </ul> </li> <li>○ Corrected bot responses</li> </ul> <p><i>Note: Due to how the ETL job is scheduled, rows in these tables can take up to ~2 hours to appear. Raw logs appear in the Inference Table in ~15 minutes.</i></p> <p><a href="#">Documentation</a></p>	
Review App	<p>Access to the Review App can be granted to any user in your company's SSO even if they do not have access to your Databricks Workspace.</p> <p>To use this feature, your account admin must enable <a href="#">SCIM</a>.</p> <p><a href="#">Documentation</a></p>	<p>Support for collecting 👍👎 feedback on individual chunks, which can be used to generate Evaluation Sets and tune the Chunk Relevance LLM judge.</p> <p>Improvements to the UX of Review App that make reviewing simpler / faster for stakeholders.</p>
LLM Judges & Metrics documentation	<p>Improved the clarity of the metric &amp; LLM Judge names.</p> <p>Added more clear documentation on what each metric &amp; LLM judge evaluates and the data used.</p> <p><a href="#">Documentation</a></p>	<p>How to guide &amp; sample code explaining techniques to try based on the results of your evaluation.</p>
New Metrics	<p>Added token count and latency metrics.</p>	<p>Fully custom metrics e.g., use a user-defined function as a metric</p>
Increasing the quality of Databricks provided LLM Judges	<p>Improved support for providing few-shot examples to the LLM Judges to tune their agreement with human raters.</p> <p><code>evaluate()</code> now accepts a Dataframe of examples that follows the same schema as the judges output, making it easier to modify incorrect judge scores / tag high quality judge scores and provide them as examples to improve the judge quality.</p> <p><a href="#">Documentation</a></p>	<p>Improvements in the quality of the LLM Judges e.g., LLM judges are more likely to agree with human rated assessments.</p>
Customer defined LLM judges	<p>Added support for customer defined LLM judges. Each judge is defined by:</p> <ul style="list-style-type: none"> <li>● An instruction (e.g., "asses if the response follows my company tone of voice, which is xyz")</li> <li>● A set of trace fields for the judge to review (e.g.,</li> </ul>	<p>Customer defined LLM judges configuration will be moved from YAML into the native MLflow SDK</p>

	<p>request, response, ground truth, etc)</p> <ul style="list-style-type: none"> <li>Optionally, 1+ examples of “good” and “bad” responses to tune the judge’s quality</li> </ul> <p><a href="#">Documentation</a></p>	
Improvements to the chain developer workflow	<p>Various improvements:</p> <ul style="list-style-type: none"> <li>Chain configuration can be passed a Dictionary vs. only as a YAML file</li> <li>Single wheel that is available on PyPi</li> </ul>	<p><code>dbutils</code> use in the <code>chain.py</code> is automatically removed vs. the user needing to comment this code out before logging the model</p> <p><code>delete_deployment()</code> API to turn off a Chain deployment</p>
Improved support for other types of Chains	Support for custom PIP dependencies with a LangChain based chain	<p>Support for logging PyFunc chains:</p> <ul style="list-style-type: none"> <li>Rather than requiring PyFunc chains to be pickled/serialized (which often fails) - PyFunc-based Chains are captured as Code + Config</li> <li>Support for non-Dataframe based PyFunc signature e.g., <code>predict(param1, param2)</code> vs. <code>predict(model_input : dataframe)</code>.</li> </ul> <p>PyFunc chains support YAML/Dictionary based configuration.</p> <p><b>Breaking change</b> More flexible schema support for chains - you can now use custom parameters in your chain’s input / output signature. Previously, you could only use one pre-defined signature for your chain.</p>
Improved compatibility w/ security postures		<p>Support for PrivateLink customers to use Model Serving with Vector Search securely</p> <p>Support for Azure Storage Firewall customers to use Inference Tables</p>

