

Dataflow for Real-Time ML & Gen AI

Overview

Machine learning (ML) and artificial intelligence (AI) empower businesses to respond to evolving market conditions and tailor their offerings to users and customers. However, decision cycles involving AI and ML can span days or even weeks, particularly when dealing with larger models (model retraining, large inference batch pipelines, etc). This solution guide introduces an architecture designed for real-time predictions, guaranteeing low latency outcomes with both custom and third-party models. Leveraging the capabilities of graphics processing units (GPUs), the proposed architecture effectively reduces prediction times to seconds.

Business Impact

The main benefits derive from the greatly reduced time-to-insight/action, which enables very fast feedback and adaptation loops.

- **Real-time predictions:** Incorporate predictions based on real-time data that enable differentiated product experiences
- **Instantaneous personalization:** Deliver personalized interactions to improve customer satisfaction without having to wait for days for those recommendations to be computed
- **Intelligent agents:** Create bots that can shepherd users through commonly submitted support requests. Use generative AI to expand the range of scenarios that can be supported.
- **Decision agility:** Detect changes and emerging trends faster to capitalize on new business opportunities and neutral emergent risks.
- **Improved customer support:** Predict customer's churn risk in real-time and extend offers & targeted measures to reduce the risk of churn.
- **Predictive maintenance:** Analyze data from sensors or system components that can predict potential failures before they occur, thereby saving cost & time of having to maintain & replace components.

This business impact can be realized across various industries:

- **Retail & E-Commerce:** Real-time personalized recommendations can improve customer satisfaction and increase retention rates. Dynamic pricing responding in real-time to fluctuating demand & inventory levels can maximize revenue per user.
- **Finance:** Real-time ML models can analyze transaction patterns, user behavior, and other data points to quickly identify and flag fraudulent activities, protecting both

businesses and consumers. Agents powered by generative AI can create personalized financial reports, investment recommendations, and retirement planning scenarios based on individual user data and goals.

- **Manufacturing:** Real-time ML models can predict equipment failures based on sensor data from machines and components, which saves customers money and time. Integrating real-time inventory data with ML-powered demand forecasts optimize your supply chain and ensure manufacturing efficiency.
- **Media & entertainment:** Integrate generative AI capabilities into product experience to help creators enhance user-generated content. Collect real-time data on user behavior & viewing patterns to make personalized content recommendations.

Customer Stories

- [Spotify](#) applies heavyweight ML models to process uploaded podcasts and offer relevant snippets for previews. With Dataflow, they have reduced processing time from 2 hours down to 2 minutes.
- [Go-Jek](#) uses Dataflow for its entire machine learning lifecycle (feature creation, standardization, and inference) to forecast demand and dynamically price services for its on-demand gig platform
- [HSBC](#) built a new scenario-risk modeling tool with Dataflow that sped up computations by 16x and empowers traders to better manage their portfolios on an intraday basis

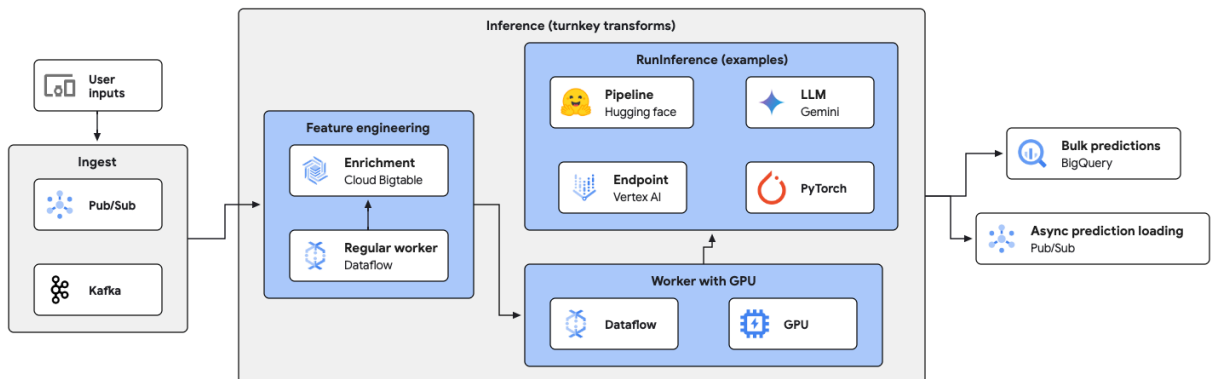
Technical Benefits

- **Developer ease of use with turnkey transforms:** Author complex ML pipelines using utility transforms that can reduce lines code by orders of magnitude:
 - [MLTransform](#) helps you prepare your data for training machine learning models without writing complex code or managing underlying libraries. ML Transforms can generate embeddings that can push data into vector databases to run inference.
 - [RunInference](#) lets you efficiently use ML models in your pipelines, and contains a number of different optimizations underneath the hood that make this an essential part of any streaming AI pipelines
- **Advanced stream processing:** Customers can implement advanced streaming architectures using the open-source [Apache Beam SDK](#), which provides a rich set of capabilities including state & timer APIs, transformations, side inputs, enrichment, and a broad list of I/O connectors.
- **Notebooks integration:** Develop your streaming AI pipeline in a notebook environment, which allows for interactive development and sampling unbounded data sources.

- **Cost efficiency:** Run pipelines without wasting precious resources & cost overruns:
 - **GPU Support:** Accelerate your processing with GPUs, which can return results faster for your most computationally demanding pipelines
 - **Right-fitting:** Deploy pipelines on heterogeneous worker pools. Rightfitting allows you to allocate additional resources to individual stages in your pipeline, which prevents wasteful utilization for stages that don't require the same compute.
- **Open-source compatibility:** Dataflow has support for [running inference with Gemma](#) as well as a strong integration with Tensorflow Extended. Customers should feel comfortable that these pipelines can be ported to any other execution engine with Apache Beam support.

Architecture

The architecture assumes that user inputs are ingested into Pub/Sub by some application or website. The same application/website subscribes to the output topic to receive the predictions when they are ready. This application/website would load the predictions (e.g. recommendations) asynchronously for users, showing them as they become available in the output topic.



The components of the architecture are the following:

- An incoming topic in Pub/Sub or Kafka
- A Dataflow job, possibly leveraging GPUs if using a local model
- Optionally: Bigtable instance or Vertex AI Feature store, with features values
- Optionally: model deployed in a Vertex AI Endpoint (the alternative would be a local model)
- A destination dataset for the analytics output
- A destination topic for real time activation of the predictions.

Design considerations

Streaming inference is complex because of potential delays applying a ML/AI model to a given input of data. For high performance non-blocking processing, you need to decouple the prediction loop from the main application loop. Microservices architectures may be very complex in these situations, and will also require specialized hardware (GPUs) deployed in the endpoints.

A streaming inference pipeline decouples the prediction loop from the main application loop. The input and output of pipelines are message queues that act as buffers between both loops.

The Dataflow pipeline consists of several components, which are shown in the diagram as different boxes:

- The first step is recovering features from a feature store. In this case, since it is streaming, we assume that the features are stored in Bigtable, for very fast lookups. [The Vertex AI Feature store is also supported "out of the box". using the Enrichment turnkey transform available in Beam.](#)
- Once we have recovered all the features, we are ready to send the data to a model.
 - Here we are using Dataflow ML with a GPU for fast inference. Thanks to the RunInference turnkey transform, this inference process could be with a model that is loaded locally in the worker. For instance, a HuggingFace pipeline (simply specified by its id and automatically downloaded to the worker), a model with Tensorflow, PyTorch, etc.
 - In some situations, we don't necessarily need to provision a GPU with the worker. The model could also be deployed as an endpoint in Vertex AI, or we could even leverage Gemini (or any other LLM) through its API, and we would only need a CPU. The web API and the RunInference transforms make it very easy to consume those models, or even any generic API not supported by RunInference, in an optimal way, without having to write additional code.

Once the predictions are ready, there are two alternatives represented in the architecture diagram. These are complementary, not necessarily exclusive:

- The predictions and the features can be written to BigQuery for analytics purposes.
 - This is similar to doing bulk prediction jobs, except that in this case the predictions get constantly updated and refreshed in BigQuery. We could even do upserts for a given user id/entity id if we want, to always have the latest value of a prediction and simplify the queries to always consume fresh data in BigQuery.
- Probably, the most interesting application is reacting in real time to those predictions.

- For that purpose, we also suggest publishing the prediction in Pub/Sub, and connecting Pub/Sub with other systems. These could be simply lazy-loading the predictions in a website (e.g. images generated by a input query using a diffusion model), or any kind of activation/reaction to the prediction. In other architectures in the following slides we are going to see how this pattern is useful for some specific cases (ad bidding, marketing activation, etc)

In summary, you should take into consideration the following points:

Step	Description
Real time input and output (Extract)	The architecture assumes Cloud Pub/Sub or Kafka for receiving data from transactions in real time. Setting up the real time export from those systems into Pub/Sub is out of the scope of this design.
Enrichment (Transform)	<p>We are assuming that the upstream data sources will only populate an "id" (or more than one), and that any other feature required for the inference process is stored in Bigtable, Vertex AI Feature Store, or any other fast-lookup database. In the above architecture, we assume it is Bigtable, since it is supported by the turnkey "Enrichment" transform of Apache Beam. The same would work for Vertex AI Feature store.</p> <p>For any other database, the design may require the development of custom code for lookups. We recommend using the Web APIs I/O connector in those cases.</p> <p>Combining the Enrichment and Web APIs I/O may require using cross language pipelines with the Dataflow runner v2 (see instructions for Python and for Java).</p>
Inference (Transform)	Once the data is hydrated and all the necessary features are extracted, we can apply ML/AI models to our data. The RunInference turnkey transform

	<p>works both with external models deployed and with local models leveraging a GPU.</p> <p>The external models could be deployed in a Vertex AI endpoint, as a HuggingFace pipeline, as a GenAI API such as Gemini or Gemma, and many other options.</p>
Output (Load)	<p>Once the prediction is obtained, the corresponding prediction is published into a Pub/Sub or Kafka topic for low-latency reactions/activation.</p> <p>However, it is interesting to publish in streaming to BigQuery, so it is immediately available for combining with any BigQuery data for reporting, analytics or business intelligence purposes.</p>

Planning your pipelines

We will cover key considerations when planning your Dataflow pipelines in this section.

Service level objectives & indicators (SLOs and SLIs)

An important measure of performance is how well your pipeline meets your business requirements. Service level objectives (SLOs) provide tangible definitions of performance that you can compare against acceptable thresholds. For example, you might define the following example SLOs for your system:

- **Data freshness:** generate 90% of product recommendations from user website activity that occurred no later than 3 minutes ago.
- **Data correctness:** within a calendar month, less than 0.5% of customer invoices contain errors.
- **Data isolation/load balancing:** within a business day, process all high-priority payments within 10 minutes of lodgement, and complete standard-priority payments by the next business day.

You can use service level indicators (SLIs) to measure SLO compliance. SLIs are quantifiable metrics that indicate how well your system is meeting a given SLO. For example, you can measure the example data-freshness SLO by using the age of the most recently processed user activity as an SLI. If your pipeline generates recommendations from user activity events, and if your SLI reports a 4-minute delay between the event time and the time the event is processed, the recommendations don't consider a user's website activity from earlier than 4 minutes. If a pipeline that processes streaming data exceeds a system latency of 4 minutes, you know that the SLO is not met.

Sources & sinks

To process data, a data pipeline needs to be integrated with other systems. Those systems are referred to as sources and sinks. Data pipelines read data from sources and write data to sinks. In addition to sources and sinks, data pipelines might interact with external systems for data enrichment, filtering, or calling external business logic within a processing step.

For scalability, Dataflow runs the stages of your pipeline in parallel across multiple workers. Factors that are outside your pipeline code and the Dataflow service also impact the scalability of your pipeline. These factors might include the following:

- **Scalability of external systems:** external systems that your pipeline interacts with can constrain performance and can form the upper bound of scalability. For example, an [Apache Kafka](#) topic configured with an insufficient number of partitions for the read throughput that you need can affect your pipeline's performance. To help ensure that the pipeline and its components meet your performance targets, refer to the best practices documentation for the external systems that you're using. You can also simplify infrastructure capacity planning by using Google Cloud services that provide built-in scalability. For more information, see [Using Google Cloud managed sources and sinks](#) on this page.
- **Choice of data formats:** certain data formats might be faster to read than others. For example, using data formats that support parallelizable reads, such as Avro, is usually faster than using CSV files that have embedded newlines in fields, and is faster than using compressed files.
- **Data location and network topology:** the geographic proximity and networking characteristics of data sources and sinks in relation to the data pipeline might impact performance.

Regional considerations

Dataflow is offered as a managed service in [multiple Google Cloud regions](#). When choosing a region to use to run your jobs, consider the following factors:

- The location of data sources and sinks
- Preferences or restrictions on data processing locations
- Dataflow features that are offered only in specific regions
- The region that's used to manage execution of a given job
- The zone that's used for the job's workers

For a given job, the region setting that you use for the job and for the workers can differ. For more information, including when to specify regions and zones, see the [Dataflow regions documentation](#).

By specifying regions to run your Dataflow jobs, you can plan around regional considerations for high availability and disaster recovery. For more information, see [High availability and geographic redundancy](#).

Security

As a fully managed service, Dataflow automatically encrypts data that moves through your data pipeline using Google-managed encryption keys for both in-flight data and at-rest data. Instead of using Google-managed encryption keys, you might prefer to [manage your own encryption keys](#). For that case, Dataflow supports customer-managed encryption keys (CMEK) using the [Cloud Key Management Service \(KMS\)](#). You can also use [Cloud HSM](#), a cloud-hosted hardware security module (HSM) service that allows you to host encryption keys and perform cryptographic operations in a cluster of [FIPS 140-2 Level 3](#) certified HSMs.

When you use CMEK, Dataflow uses your Cloud KMS key to encrypt the data, [except for data-key-based operations such as windowing, grouping, and joining](#). If data keys contain sensitive data, such as personally identifiable information (PII), you must hash or otherwise transform the keys before they enter the Dataflow pipeline.

Networking

Your networking and security requirements might mandate that VM-based workloads such as Dataflow jobs use only private IP addresses. Dataflow lets you specify that workers use private IP addresses for all network communication. If public IPs are disabled, you must enable [Private Google Access](#) on the subnetwork so that Dataflow workers can reach Google APIs and services.

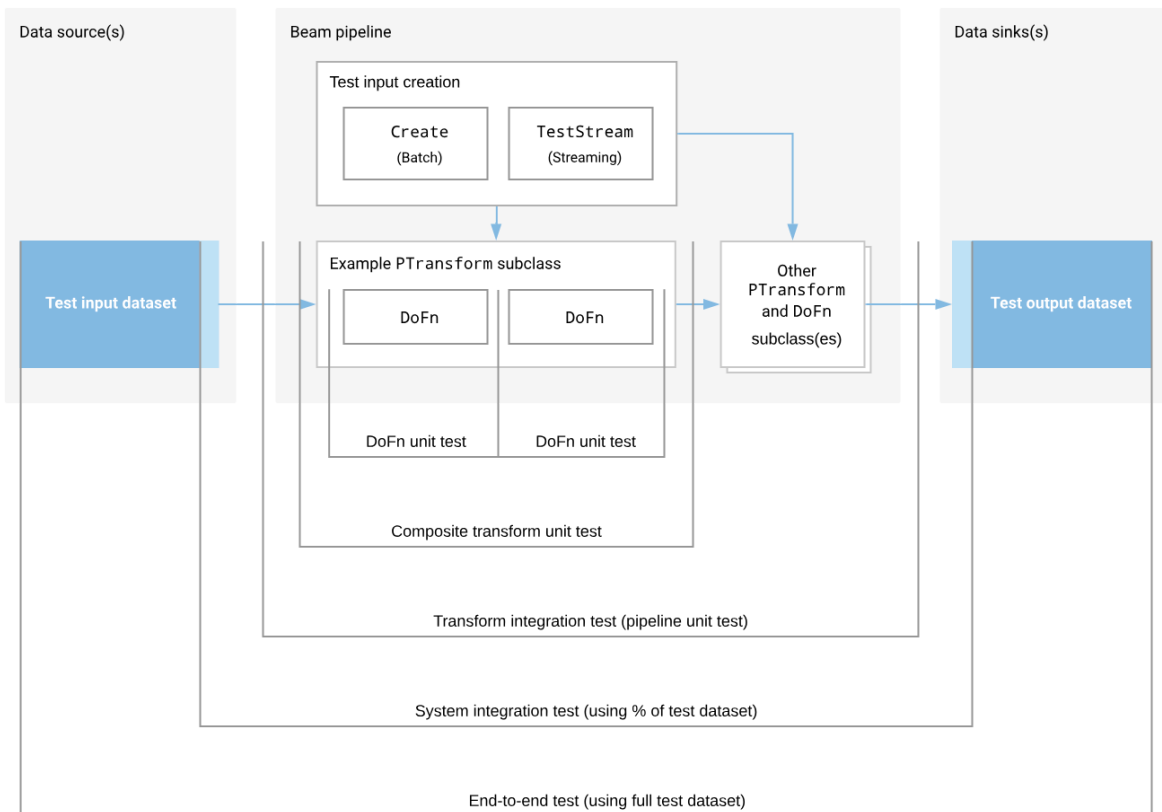
We recommend that you disable public IPs for Dataflow workers, unless your Dataflow jobs require public IPs to access network resources outside of Google Cloud. Disabling public IPs prevents Dataflow workers from accessing resources that are outside the subnetwork or from accessing [peer VPC networks](#). Similarly, network access to VM workers from outside the subnetwork or peer VPC networks is prevented.

For more information about using the `--usePublicIps` pipeline option to specify whether workers should have only private IPs, see [Pipeline options](#).

Developing your pipelines

The way that the code for your pipeline is implemented has a significant influence on how well the pipeline performs in production. In order to develop and deploy battle-tested code, we recommend the following:

- **Pipeline runners:** Use different Apache Beam runners to run pipeline code. The Apache Beam SDK provides a Direct Runner for local development and testing. You can use the Dataflow Runner for ad hoc development testing and for end-to-end pipeline tests.
- **Deployment environments:** Create deployment environments to separate users, data, code, and other resources across different stages of development. Run pipeline locally for development and rapid testing using the Direct Runner. Create a pre-production environment for development phases that need to run in production-like conditions, such as end-to-end testing. The production environment should be a dedicated Google Cloud project, where continuous delivery copies deployment artifacts to the production environment when all end-to-end tests have passed.
- **Leverage open-source:** Apache Beam provides a rich set of pipeline examples in its directories for developers to copy from. The [Dataflow Cookbook](#) provides a library of transformations and common patterns in Java, Scala, Python, and Go. There are also guides for [common use-case patterns](#) found in our documentation. Google-provided [templates are open source](#) under the Apache License 2.0, so you can use them as the basis for new pipelines. The templates are also useful as code examples for reference.
- **Test pipeline code:** Use unit tests, integration tests, and end-to-end tests when applicable. The Apache Beam SDK provides functionality to enable these tests. The Apache Beam SDK provides functionality to enable these tests. Ideally, each type of test targets a different deployment environment. The following diagram illustrates how unit tests, integration tests, and end-to-end tests apply to different parts of your pipeline and data.



Apache Beam best practices

In addition to the guidance listed above, Apache Beam and Dataflow provides features that complement these best practices for improved production readiness.

- Leverage turnkey transforms:** Turnkey transformations provide a utility for developers to accomplish common business logic patterns in the convenience of a transform. These transformations abstract away unnecessary overhead that can take dozens, if not hundreds of lines, to manage.
 - Enrichment:** When you enrich data, you augment the raw data from one source by adding related data from a second source. The additional data can come from a variety of sources, such as [Bigtable](#) or [BigQuery](#). The Apache Beam enrichment transform uses a key-value lookup to connect the additional data to the raw data.
 - RunInference:** Users oftentimes needs to make a call for a prediction from an ML model stored externally to the pipeline. The RunInference API enables you to run models as part of your pipeline in a way that is optimized for machine learning inference. To reduce the number of steps in your pipeline, RunInference supports features like batching.

- **Micro-batch calls to external services:** When you call external services, you can reduce per-call overheads by using the `GroupIntoBatches` transform to create batches of elements of a specified size. Batching sends elements to an external service as one payload instead of individually. In combination with batching, you can limit the maximum number of parallel (concurrent) calls to the external service by choosing appropriate keys to partition the incoming data. The number of partitions determines the maximum parallelization. For example, if every element is given the same key, a downstream transform for calling the external service does not run in parallel.
- **Queue unprocessable data:** Your pipeline might encounter situations where it's not possible to process elements. This situation can occur for different reasons, but a common cause is data issues. Use a pattern called a dead-letter queue (or dead-letter file). Catch exceptions in the `DoFn.ProcessElement` method and log errors as you normally would. Instead of dropping the failed element, use branching outputs to write failed elements into a separate `PCollection` object. These elements are then written to a data sink for later inspection and handling by using a separate transform.

Deploying your pipelines

Deployment

Pipeline development involves different stages and tasks, such as code development, testing, and delivery into production. Deploying a pipeline within a robust continuous integration & continuous delivery framework can ensure safe rollouts that do not introduce regressions.

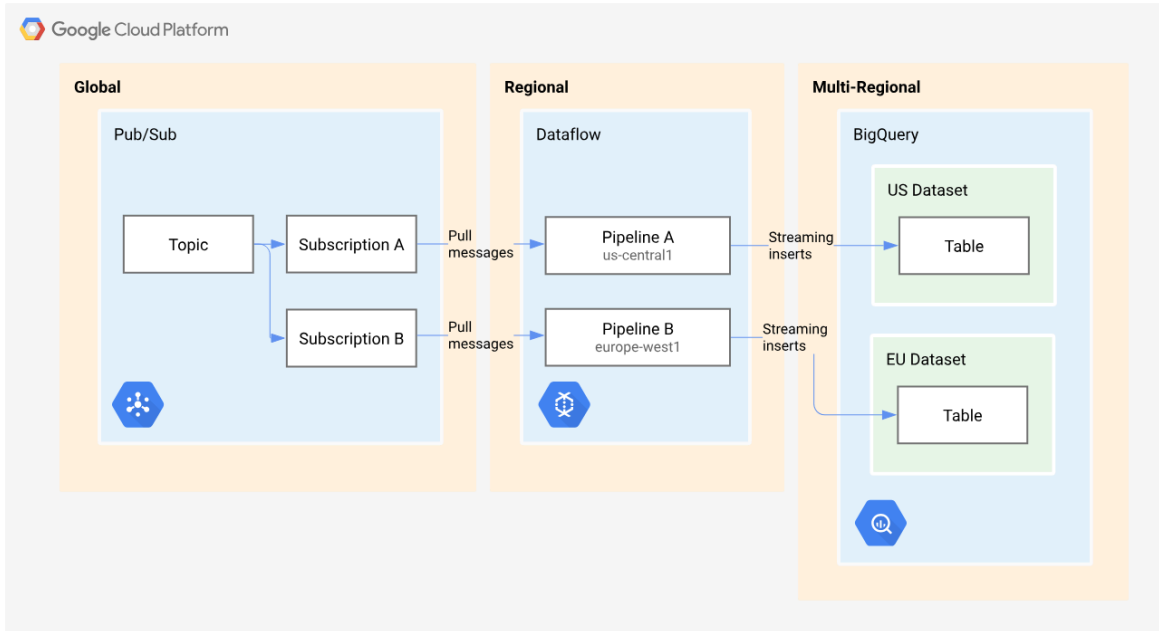
The number and types of deployment artifacts created from a passing build varies depending on how pipelines are launched. Using the Apache Beam Java SDK, you can package your pipeline code [into a self-executing JAR file](#). You can then store the JAR file in a bucket that is hosted in the project for a [deployment environment](#), such as the preproduction or production Google Cloud project. If you use Classic Templates (a type of [templated execution](#)), the deployment artifacts include a [JSON template file](#), the JAR file for your pipeline, and an [optional metadata template](#). You can then deploy the artifacts into different deployment environments using continuous delivery.

You can create a Dataflow job by using the Apache Beam SDK directly from a development environment. This type of job is called a non-templated job. Although this approach is convenient for developers, you might prefer to separate the tasks of developing and running pipelines. To make this separation, you can use [Dataflow templates](#), which allow you to stage and run your pipelines as independent tasks. After a template is staged, other users, including non-developers, can run the jobs from the template using the Google Cloud CLI, the Google Cloud console, or the Dataflow REST API.

Reliability

It is important to consider the failures that might occur in the event of an outage. The following are a few principles to keep in mind:

- **Follow isolation principles:** A general recommendation to improve overall pipeline reliability is to follow the isolation principles behind [regions and zones](#). Ensure that your pipelines don't have critical cross-region dependencies. If you have a pipeline that has critical dependency on services from multiple regions, a failure in any one of those regions can impact your pipeline. To help avoid this issue, deploy to multiple regions for redundancy and backup.
- **Create Dataflow snapshots:** Dataflow offers a snapshot feature that provides a backup of a pipeline's state. You can restore the pipeline snapshot into a new streaming Dataflow pipeline in another zone or region. You can then start the reprocessing of messages in the Pub/Sub or Kafka topics starting at the snapshot timestamp. If you set up regular snapshots of your pipelines, you can minimize Recovery Time Objective (RTO) time.
- **Mitigate regional outages by using high availability or failover:** For streaming jobs, depending on the fault tolerance and budget for your application, you have different options for mitigating failures. For a regional outage, the simplest and most cost-effective option is to wait until the outage ends. However, if your application is latency-sensitive or if data processing must either not be disrupted or should be resumed with minimal delay, there are a couple of architectural options to consider:
 - *High-availability: Latency sensitive with no data loss:* If your application cannot tolerate data loss, run duplicate pipelines in parallel in two different regions, and have the pipelines consume the same data. The same data sources need to be available in both regions. The downstream applications that depend on the output of these pipelines must be able to switch between the output from these two regions. Due to the duplication of resources, this option involves the highest cost compared to other options.



- *Failover: Latency-sensitive with some potential data loss:* If your application can tolerate potential data loss, make the streaming data source available in multiple regions. For example, using Pub/Sub, maintain two independent subscriptions for the same topic, one for each region. If a regional outage occurs, start a replacement pipeline in another region, and have the pipeline consume data from the backup subscription.
- **Don't store data in the broker for long periods of time:** There is no need to do that with Dataflow for increasing robustness. Dataflow ensures exactly once processing and will not pull data more than once from the broker. Dataflow can do live updates to ensure continuity without gaps.

Operating your pipelines

Monitoring

Cloud Monitoring provides powerful logging and diagnostics. Dataflow integration with Monitoring lets you access Dataflow job metrics such as job status, element counts, system lag (for streaming jobs), and user counters from the Monitoring dashboards. You can also use Monitoring alerts to notify you of various conditions, such as long streaming system lag or failed jobs.

A streaming pipeline is a service. Start by defining a SLI and SLO of what you want to achieve with this pipeline, as mentioned when we discussed planning your pipeline. Leverage Cloud

Monitoring to measure the SLO automatically, and determine error budget and other properties without additional tooling needed.

To maintain a robust monitoring posture, the following best practices are recommended:

- **Setup relevant dashboards:** Dataflow publishes [a list of metrics](#) relevant to the performance of your pipeline. These include metrics describing resource utilization, data freshness, system latency, parallelism, and read/write throughput. Creating charts that are directly relevant to your SLOs and SLIs can simplify assessing the performance of your pipeline.
- **Inspect logs & insights:** Worker and job logs can indicate performance issues with your pipeline. Establish a practice of reviewing the diagnostics panel of the Dataflow monitoring UI, and review insights that can improve job performance, reduce cost, or troubleshoot errors.
- **Create alerts:** You can create an alerting policy directly from a metric chart. Creating alerts that align with your SLOs and SLIs can help your organization respond to a regression in your pipeline before a service outage.

Performance

Echoing the guidance regarding defining SLOs and SLIs that are relevant to the key performance indicators that you care about, there are a number of considerations that should be taken into account when optimizing the performance of your pipeline. Delivering better performance for your pipeline (defined by either lower latency or higher efficiency – or both) can generally be done by carefully observing your pipeline in production into your pipeline implementation. The following list are recommended areas to identify for performance improvements (albeit not exhaustive):

- **Region, quotas, and networking:** Confirm that the pipeline is being run in the region where resources are located. Check if networking parameters permit adequate throughput for the pipeline.
- **Excessive data shuffling:** Stages where data is shuffled (i.e. when you are using GroupByKey, CoGroupByKey, or Stateful DoFns) are generally the most computationally expensive parts of a pipeline. These stages consume the most network bandwidth and if you are not using Streaming Engine or Shuffle, will likely add a meaningful amount of processing time to your pipeline. To mitigate against this, be intentional about the data that you are shuffling. Remove as many unused attributes as possible before you invoke a grouping transformations. When possible, use a combiner instead of a GroupByKey.
- **Limited parallelism:** Certain grouping operations might be slow, and you might see “hot key warnings” displayed in the UI. This is usually a result of a skew found in the

distribution of the data, which will cause the overall processing to slow down. If possible, redesign the key partitioning or use combiners. Beam also provides several variations of “withHotKeyFanout” transformations to alleviate the ill effects of hot keys.

- **Batch External API calls:** As mentioned in the best practices section, you can reduce the impacts of making calls to external APIs by microbatching calls. Calling external APIs can sometimes be subjected to quotas that are not known to your pipeline, which can serve to slow down your own pipeline. Implement DoFn’s lifecycle methods (annotated with @Setup, @Startbundle), collect API requests in the @ProcessElement method, and call the API in the @FinishBundle method; or use GroupIntoBatches transform to create batches of certain sizes. It’s important to ensure that calls are idempotent, since bundles can be replayed.
- **IO Specific Recommendations:** Depending on your pipeline, review the best practices guides on IOs. Oftentimes, performance can be significantly constrained by settings configured at the source & sink.

Cost Optimization

Being intentional about upholding performance best practices will generally result in the most cost-efficient pipelines. However, there are important steps to take when monitoring your pipeline costs, as well as configurations that can reduce your overall bill without coming at the expense of your requirements:

- **Monitor Dataflow costs & set alerts:** Enable billing export into BigQuery. This is the most efficient way of observing your Dataflow costs at the project level down to the pipeline and SKU level. Implement a labeling taxonomy using the best practices found [here](#), and add labels to [Dataflow jobs](#). For your critical pipelines, create [monitoring alerts](#) to do real-time notifications and cost control.
- **Number of workers:** Once you have a good sense for the data volume patterns of your pipeline, setting a maximum number of workers can limit any adverse traffic spikes that might cause your pipeline to overprovision workers. On the other hand, if you know that you have a baseline number of workers that you would like to maintain in order to guarantee uptime, you can set a minimum number of workers to ensure your processing doesn’t slow down during an autoscaling event.
- **At-least-once mode:** If your pipeline does not require exactly-once processing (i.e. deduplication can be handled in your sink or destination), you can configure your pipeline to run in at-least-once mode. At-least-once mode can help optimize cost and performance by turning off our exactly-once mode, which adds to pipeline latency & cost due to the reading of metadata required to ensure we process a message no more than once.

