DATA+AI
SUMMIT 2022

# Databricks Observability

**Enterprise Observability With Overwatch**

ORGANIZED BY databricks

**Daniel Tomes & Mohan Baabu**
Principal Architect, Databricks

1

# Who We Are

**Daniel**

- Principal Architect

- Big Data Since 2010 – Big Oil

- Big Science Since 2014 – Big Oil
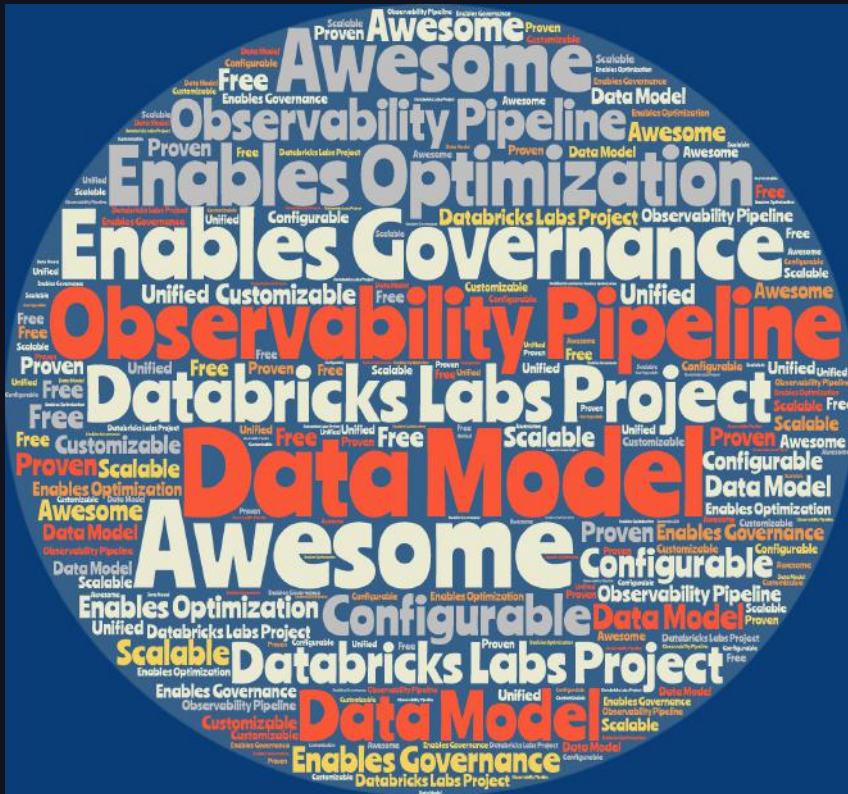
- Cloudera

- Databricks since early 2017

**Mohan**

- Data Analyst

- Software Developer since 2019 – Esales Technologies

- Data Analyst since early 2021 – Kushagramati Analytics

- Solution Consultant since 2021 - Databricks

# Agenda

- Intros
- What is Overwatch
- Motivation
- What Can It Do (Demo)
- Architecture & Implementation
- Roadmap
- Challenges / Solutions
- The Future
- Getting Started
- Surprise

**DATA+AI**
SUMMIT 2022

# What Is Overwatch

## Overwatch_Gold 0.6.0

Daniel Tomes | December 13, 2021

### Databricks Platform

**job**
- organization_id
- job_id
- action
- unixTimeMS
- timestamp
- date
- job_name
- job_type
- timeout_seconds
- schedule
- notebook_path
- new_settings
- cluster
- acl_permission_set
- grants
- target_user_id
- session_id
- request_id
- user_agent
- response
- source_ip_address
- created_by
- created_ts
- last_edited_by
- last_edited_ts
- deleted_by
- deleted_ts

**jobRun**
- organization_id
- run_id
- run_name
- job_runtime
- job_id
- id_in_job
- job_cluster_type
- job_task_type
- job_terminal_state
- job_trigger_type
- cluster_id
- notebook_params
- libraries
- children
- workflow_context
- task_detail
- request_detail
- time_detail

**accountMod**
- organization_id
- mod_unixTimeMS
- mod_date
- action
- endpoint
- modified_by
- user_name
- user_id
- group_name
- group_id
- from_ip_address
- user_agent
- request_id
- response

**notebook**
- organization_id
- notebook_id
- notebook_name
- notebook_path
- cluster_id
- action
- unixTimeMS
- timestamp
- date
- old_name
- old_path
- new_name
- new_path
- parent_path
- user_email
- request_id
- response

**accountLogin**
- organization_id
- login_unixTimeMS
- login_date
- login_type
- login_user
- user_email
- from_ip_address
- user_agent
- request_id
- response

**jobRunCostPotentialFact**
- organization_id
- run_id
- job_id
- id_in_job
- job_runtime
- run_terminal_state
- run_trigger_type
- run_task_type
- cluster_id
- cluster_name
- cluster_type
- custom_tags
- driver_node_type_id
- node_type_id
- dbu_rate
- running_days
- avg_cluster_share
- avg_overlapping_runs
- max_overlapping_runs
- run_cluster_states
- worker_potential_core_H
- driver_compute_cost
- driver_dbu_cost
- worker_compute_cost
- worker_dbu_cost
- total_driver_cost
- total_worker_cost
- total_compute_cost
- total_dbu_cost
- total_cost
- spark_task_runtimeMS
- spark_task_runtime_H
- job_run_cluster_util

**clusterStateFact**
- organization_id
- cluster_id
- cluster_name
- custom_tags
- unixTimeMS_state_start
- unixTimeMS_state_end
- timestamp_state_start
- timestamp_state_end
- state
- driver_node_type_id
- node_type_id
- current_num_workers
- target_num_workers
- uptime_since_restart_S
- uptime_in_state_S
- uptime_in_state_H
- state_dates
- days_in_state
- cloud_billable
- databricks_billable
- isAutomated
- dbu_rate
- worker_potential_core_H
- core_hours
- driver_compute_cost
- driver_dbu_cost
- worker_compute_cost
- worker_dbu_cost
- total_driver_cost
- total_worker_cost
- total_compute_cost
- total_dbu_cost
- total_cost
- driverSpecs
- workerSpecs

**dbuCostDetails**
- organization_id
- sku
- contract_price
- is_active
- activeFrom
- activeUntil

**cluster**
- organization_id
- cluster_id
- cluster_name
- action
- unixTimeMS
- timestamp
- date
- driver_node_type
- node_type
- num_workers
- autoscale
- auto_termination_minutes
- enable_elastic_disk
- is_automated
- cluster_type
- security_profile
- cluster_log_conf
- init_scripts
- custom_tags
- cluster_source
- spark_env_vars
- spark_conf
- acl_path_prefix
- instance_pool_id
- instance_pool_name
- driver_instance_pool_id
- driver_instance_pool_name
- spark_version
- idempotency_token
- deleted_by
- created_by
- last_edited_by

**instancePool**
- organization_id
- instance_pool_id
- instance_pool_name
- actionName
- timestamp
- node_type_id
- idle_instance_autotermin...
- min_idle_instances
- max_capacity
- preloaded_spark_versions

**instanceDetails**
- organization_id
- API_Name
- instance
- vCPUs
- memory_gb
- linux_vm_price_hour
- on_demand_cost_hourly
- linux_reserved_cost_hourly
- Hourly_DBUs
- isActive
- activeFrom
- activeUntil

### SparkUI

**sparkExecutor**
- organization_id
- spark_context_id
- cluster_id
- executor_id
- executor_info
- removed_reason
- executor_alivetime
- unixTimeMS
- timestamp
- date
- event_log_start
- event_log_end
- Pipeline_SnapTS
- Overwatch_RunID

**SparkStream**
- organization_id
- spark_context_id
- cluster_id
- stream_id
- stream_name
- stream_run_id
- stream_batch_id
- stream_timestamp
- streamSegment
- streaming_metrics
- execution_ids

**SparkStage**
- organization_id
- spark_context_id
- cluster_id
- stage_id
- stage_attempt_id
- unixTimeMS
- timestamp
- date
- stage_runtime
- stage_info
- event_log_start
- event_log_end
- Pipeline_SnapTS
- Overwatch_RunID

**SparkTask**
- organization_id
- spark_context_id
- cluster_id
- task_id
- task_attempt_id
- stage_id
- stage_attempt_id
- executor_id
- host
- unixTimeMS
- timestamp
- date
- task_runtime
- task_metrics
- task_info
- task_type
- task_end_reason
- event_log_start
- event_log_end
- Pipeline_SnapTS
- Overwatch_RunID

**sparkExecution**
- organization_id
- spark_context_id
- cluster_id
- execution_id
- description
- details
- unixTimeMS
- timestamp
- date
- sql_execution_runtime
- event_log_start
- event_log_end
- Pipeline_SnapTS
- Overwatch_RunID

**SparkJob**
- organization_id
- spark_context_id
- cluster_id
- job_id
- job_group_id
- execution_id
- stage_ids
- notebook_id
- notebook_path
- user_email
- db_job_id
- db_id_in_job
- db_job_type
- unixTimeMS
- timestamp
- date
- job_runtime
- job_result
- event_log_start
- event_log_end
- Pipeline_SnapTS
- Overwatch_RunID

**Key**
- Key To DB Platform
- Value Struct
- Partition Column
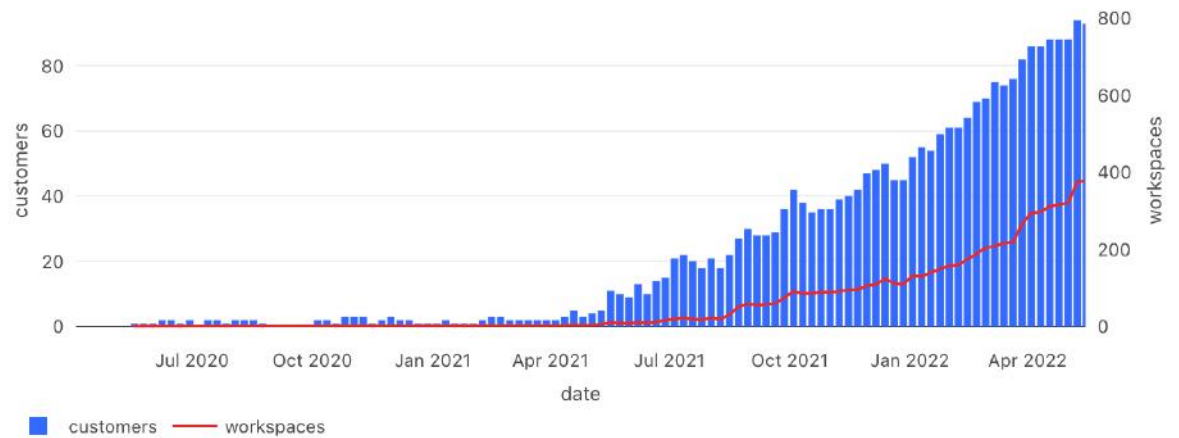- ZOrder Column
- TO BE Integrated

# Motivation

- Customers struggled to understand spend / how to implement governance

- Databricks' footprint continues to grow

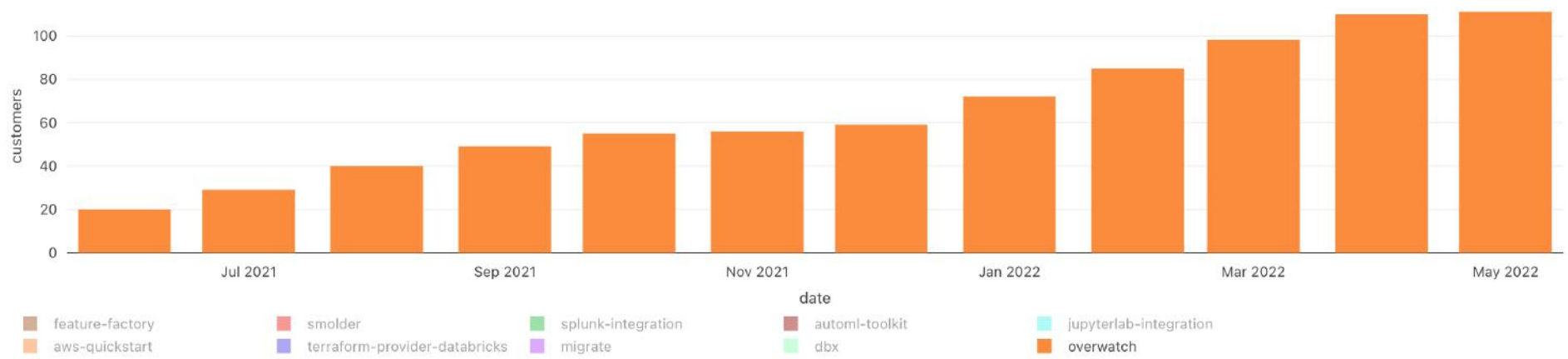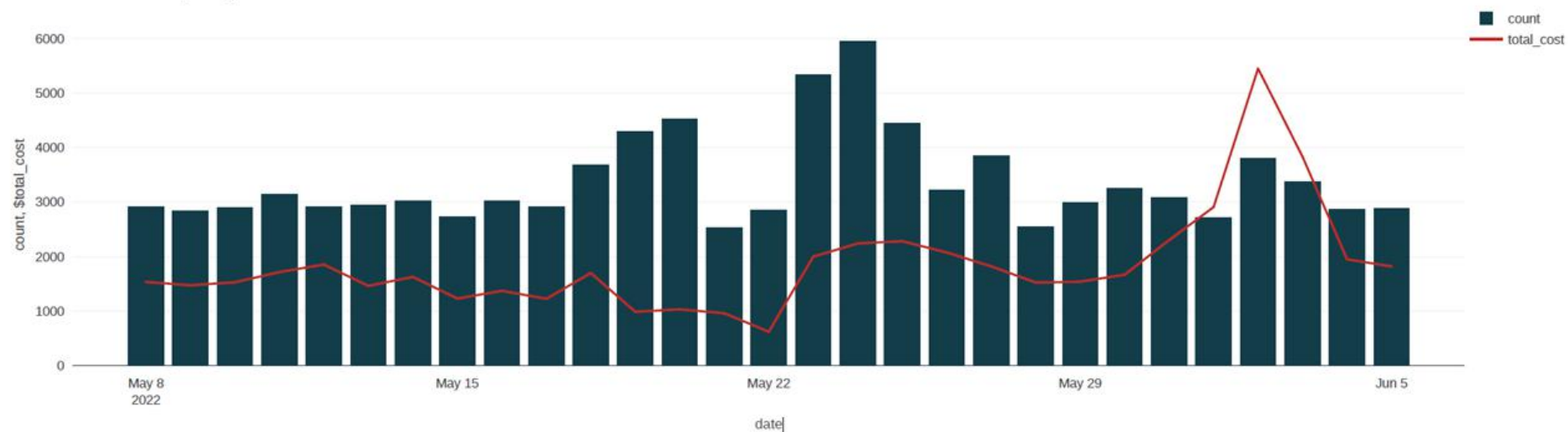- Optimization

- Unification

- Measure governance

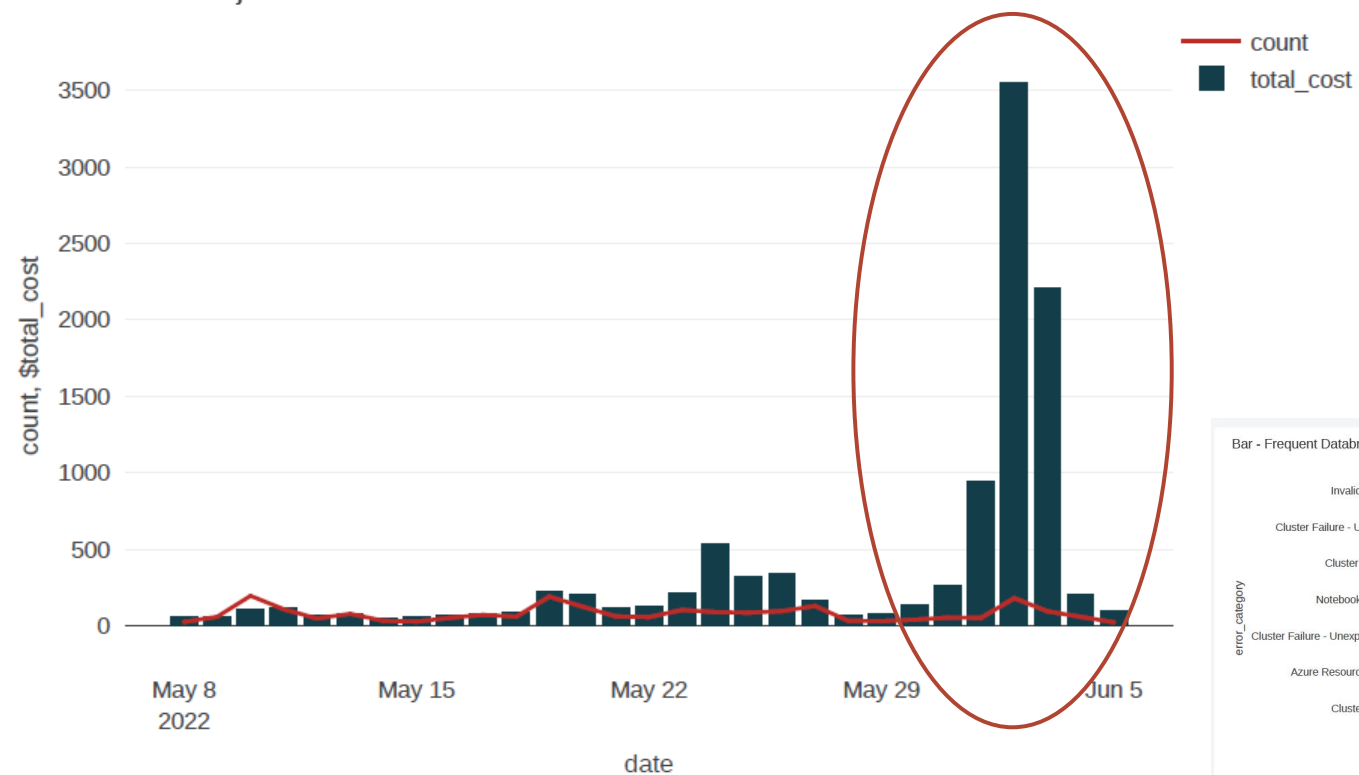# What Can This Puppy Do?

# Jobs Efficiency

- Maximize Workload Throughput

- Minimize Cost / Failures / SLA Breaches / Anomalies
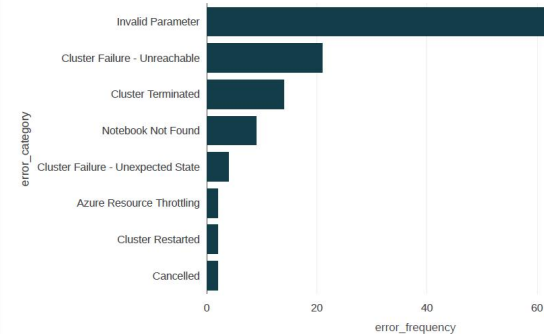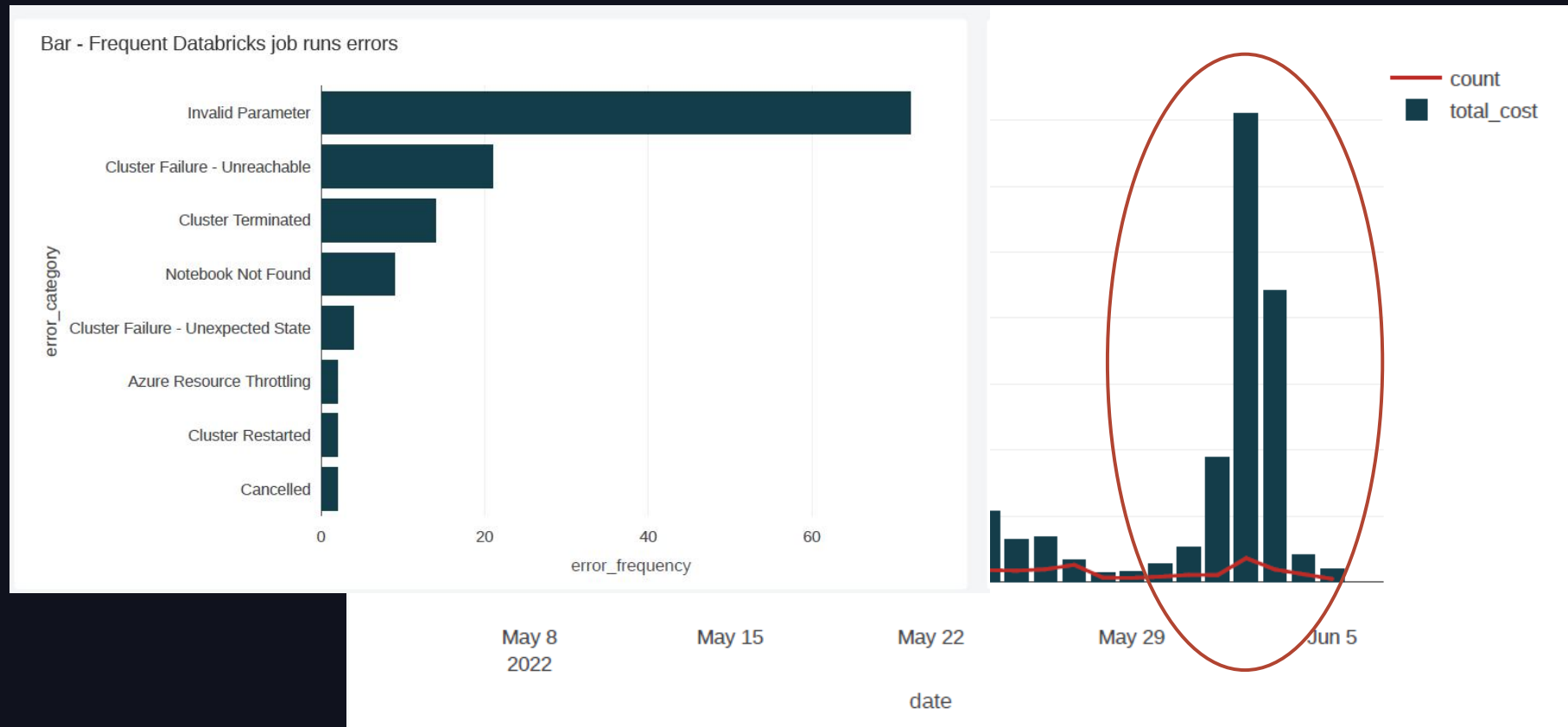


Combo - Job runs total per day

# Jobs Efficiency



Combo - Cost of job runs failures
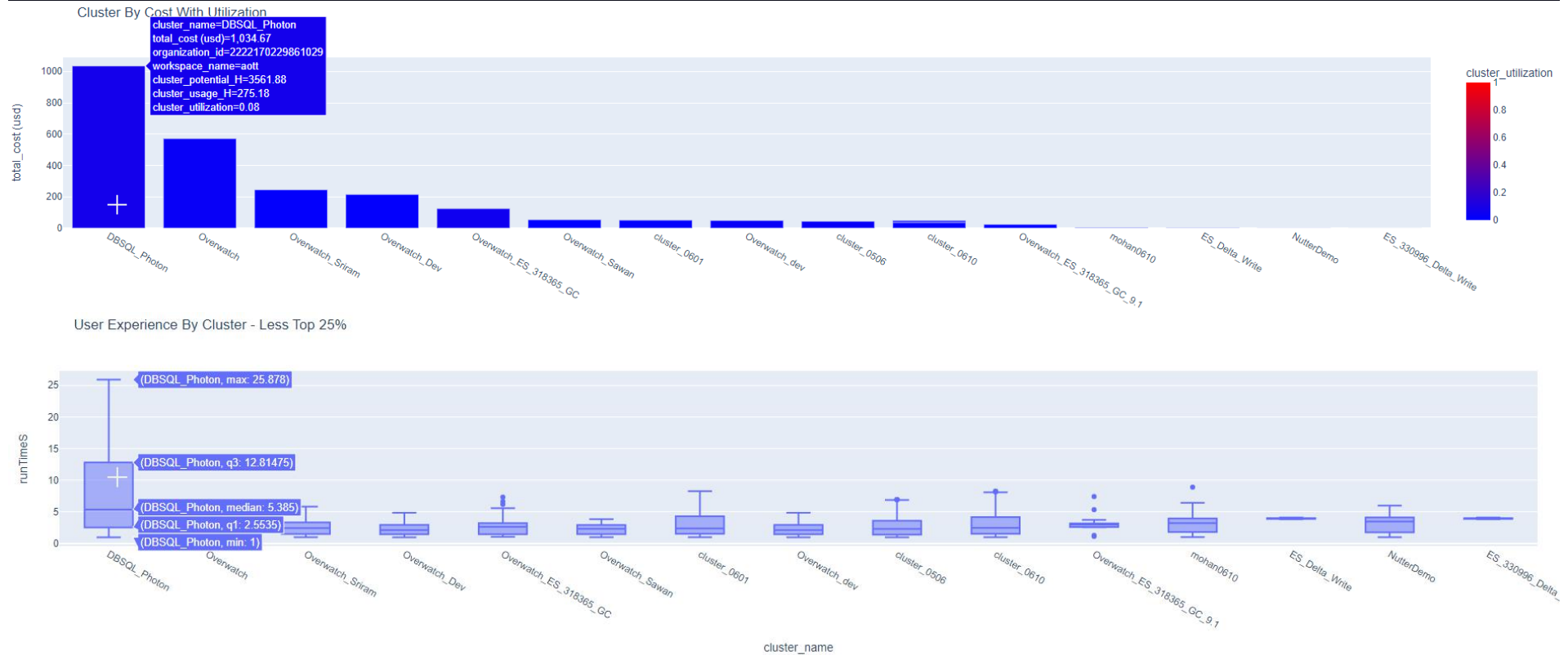


Bar - Frequent Databricks job runs errors

# Jobs Efficiency

# Interactive Efficiency

- Balancing Act

- Maximize utilization WHILE maximizing User Experience
  - User Experience == fast queries

- Optimized through workload classification and governance



https://images.all-free-download.com/images/graphiclarge/demand_and_income_concept_balance_icon_6829040.jpg

# Interactive Efficiency

# Interactive Efficiency



RunTime over Cores (Minutes)

# Autoscaling Efficiency


Bar - Number of autoscaling events of interactive clusters

FREQUENT

EFFICIENT?


Line - cost of autoscaling events of interactive clusters

# Autoscaling Efficiency



PERIOD P = 15m

Scale Down

Scale Down

Opportunity

Opportunity

Review Result
Update Query
Resubmit

Review Result
Update Query
Resubmit

Scale Up

Scale Up

$Q1_{Start}$

$Q1_{End}$

$Q2_{Start}$

$Q2_{End}$

**Jitter**:
Scale up Request issue within n seconds of scale down OR
1+ consecutive scale down events surrounded by scale up events within some threshold period P

**Example Above**:
Scale up request ⭐ within 5m (30s) of scale down event ⭐

Scale Down Event ⭐ surrounded by scale up requests ⭐ & ⭐ within a determined threshold period of 15 min

# Autoscaling Efficiency



Interactive Jitter Impact

cluster_id=1105-021745-f
jitterCosts=4,052.88
organization_id=
workspace_name=
cluster_name=        by Ca        Cluster
jitterEvents=2084

- **Solution:** Tune autoscaling for high–jitter clusters

DATA+AI
SUMMIT 2022

# Other Sweet Use Cases

- Maximize Spot Market
  - By node family, type, region (AZ) and spot availability
  - Stagger job runtimes by market availability

- Optimize local disk configuration
  - High shuffle – fast / substantial local

- Anti-Pattern Identification (by user/group/etc)
  - Micro-Targeted Training

- Policy Efficacy & Value
  - Quantifies value of FTEs in governance teams

- Streams Optimization & Stabilization
  - Optimize cost while still meeting SLAs

# Architecture

# Implementation

- New Job
- Configure Cluster
- Configure Overwatch
- Run


Small Cluster


Schedule and Forget


Custom Config to Meet Your Needs

# Roadmap

- Stronger Tests

- Add latest SKUs to get complete cost coverage

- Low-Latency Bronze Layer
  - Facilitates higher-frequency alerts without overhead of entire pipeline

- Additional Gold Model Entities
  - DBSQL
  - JDBC
  - Additional Fact Tables
  - Delta Live Tables / Delta Sharing / Unity Catalog

- Databricks Curated Bronze Data

- Roadmap Maintained On Git
  - Submit your feature requests

https://1.bp.blogspot.com/-BcfQpJiH48M/XaWwi4Kxt_I/AAAAAAAAA9I/5W33CJEWJWoZHjEX9ZRvCQf2jJ2W5_rcACLcBGAsYHQ/s1600/SeekPng.com_roadmap-png_4940541.png

# Challenges

- ## Keeping Pace
  - Keeping up is challenging

- ## Bronze Layer Data
  - Using data that wasn't built for this
  - Complex ETLs



http://blog.markallencoaching.com/wp-content/uploads/2020/08/Challenge-Yourself.jpg

- ## Data Quality
  - Customers use and abuse Databricks in all the ways
  - Maintaining high-quality in all scenarios

# Solutions

- ## Keeping Pace
  - Integrate with the product directly

- ## Dirty Bronze
  - Data contracts from product
  - Publish directly to Unity Catalog (System Tables)

- ## Data Quality
  - Strong Integration Tests
  - Engineered datasets
  - Cleaner Bronze

# Getting Started



## OVERWATCH

- Search...
- Home
- Getting Started
- Environment Setup

### Project Overview

Overwatch was built to enable Databricks' customers, employees, and partners to quickly / easily understand operations within Databricks deployments. As enterprise adoption increases there's an ever-growing need for strong governance. Overwatch means to enable users to quickly answer questions and then drill down to make effective operational changes. Common examples of operational activities Overwatch assists with are:

- Cost tracking by various dimensions

https://www.dailydot.com/wp-content/uploads/2018/06/overwatch-characters-ranked-1024x512.jpg

Account Team

| | project | go-live-prep (#101) | 13 months ago |
| | src | 061 cluster view columns | 15 days ago |
| | terraform | first attempt on improving terraform templates (#267) | 6 months ago |
| | .gitignore | limits added to clusterEvents max History | 12 months ago |
| | Building.md | Removing previously added documentation | 14 months ago |
| | CONTRIBUTING.md | required updates from legal (#98) | 13 months ago |
| | LICENSE | Simplify and optimize 2 (#54) | 15 months ago |
| | NOTICE | required updates from legal (#98) | 13 months ago |
| | README.md | Update README.md | 8 months ago |
| | build.sbt | rev version to 061 | 15 days ago |

monitoring  databricks

- Readme
- View license
- 70 stars
- 17 watching
- 24 forks

Releases

- v0604  Latest
  on Feb 28

+ 15 releases

GIT

# SURPRISE

# Instructions: Read me!

**Getting started with our slide template**

When using this template, create your new slides at the very top of the slide order, above this slide. Explore the advice and example slides below to find useful layouts and graphics to pull into your design. **When your slide deck is complete, delete this slide and every slide below it.**

**DATA+AI**
SUMMIT 2022

# Presentation best practices

## Less is more

**Clarity over density**

Don't try to cram everything onto a limited number of slides. More slides with less text per slide is easier to digest.

**Make it scannable**

Use text hierarchy to create order and keep your content scannable. No walls of text! Try to keep headlines short.

**Get creative**

There are great baseline slides in this template, but it may not have everything you need. Don't be afraid to craft your own layouts! Just pay attention to the font and grid guidelines, and take advantage of starter shapes.

DATA+AI
SUMMIT 2022

# Font Guidance

## Font selection

All text in our slide decks should use one of two available event brand fonts: **DM Sans** or **DM Mono**.

If you do not see these fonts in your font selection menu, they can be added by selecting "More fonts" and searching for "dm." Click on DM Sans and DM Mono, then hit OK.

**1**

| Arial ▾ | — | 20 | + | **B** | *I* |

A₊  More fonts

**2**

Fonts

dm    🔍    Scripts: All Scripts ▾   Show: All fonts ▾   Sort: Popularity ▾

✓  DM Sans

DM Serif Display

DM Serif Text

✓  DM Mono

Goldman

OK    Cancel

# Font Guidance (Cont.)

## Font sizing

Using consistent type sizing is a good way to help your slides feel uniform. When selecting type sizes, try to stick to multiple of 8, with the exceptions of 12 and 20 as in-betweens.

64  DATA+AI Summ
56  DATA+AI Summ
40  DATA+AI Summit
32  DATA+AI Summit
24  DATA+AI Summit
20  DATA+AI Summit
16  DATA+AI Summit
12  DATA+AI Summit

**DATA+AI**
SUMMIT 2022

# Grid Guidance

**Keep it orderly**

Your presentation template has a 12 column grid to help you organize the elements on your slides. When laying out objects, consider using the grid to help.

Toggle the grid visibility by navigating to *View > Guides > Show Guides*.

# Color Guidance

## Keep it on brand

When customizing charts or adding other visual elements, do your best to stay within our defined event color palette. This will ensure that all your content looks great together and doesn't clash with the slide template design.

**Always use black text when placing content over a colored background.** The only exception is when using a black background. Any color text is acceptable on black.

| 10121E | 00B6E0 | 85DDB5 | F16047 |
|--------|--------|--------|--------|
| EDEEF1 | 8FDDEF | AFE9CF | F3A89B |

# Example Slides

**DATA+AI**
SUMMIT 2022

# Choose Your Title Slide

Eighteen colorful title slide options with varying shapes

ORGANIZED BY databricks

**Add your Name**
Add your title, company

# DATA+AI
## SUMMIT 2022

# Choose Your Title Slide

Eighteen colorful title slide options with varying shapes

ORGANIZED BY ◈ databricks

Add your Name
Add your title, company

**DATA+AI**
SUMMIT 2022

# Choose Your Title Slide

Eighteen colorful title slide options with varying shapes

ORGANIZED BY ⬨ databricks

**Add your Name**
Add your title, company

# Basic Content Slide

**Your all-purpose zone**

Use this slide as a starting point for crafting your own layouts, or for simple text slides.

**DATA+AI**
SUMMIT 2022

# Activate Dark Mode

Mix in black slides to add contrast and variety

Or make your whole presentation dark!

DATA+AI
SUMMIT 2022

37

# Insert your charts or images

## Take advantage of the content panels

**Insert Image by URL**

If you want to insert a gif or other image from the web, simply navigate to *Insert > Image > by URL.*

Crop and resize your image to fit within content panels, if you're feeling fancy.

"With just a few adjustments to text size and alignment, you can use the basic content slide for other types of content such as quotes."

**Andrew Pons**
Slide Designer

| | Column A | Column B | Column C | Column D | Column E | Column F |
|---|---|---|---|---|---|---|
| Row A | You can create simple tables to help organize information. | | | | | |
| Row B | | | | | | |
| Row C | | | | | | |
| Row D | | | | | | |
| Row E | | | | | | |
| Row F | | | | | | |
| Row G | | | | | | |
| Row H | | | | | | |

# Timeline Style One

## Your subtitle here

Timeline Item

Timeline Item

Timeline Item

Timeline Item

Timeline Item

Timeline Item

Timeline Item

# Timeline Style Two

## Your subtitle here

| Q1 | Q2 | Q3 | Q4 |
|---|---|---|---|

Your gantt chart item

Your gantt chart item

Your gantt chart item

Your gantt chart item

Your gantt chart item

Your gantt chart item

Your gantt chart item

Your gantt chart item

# Single Column

## Content Tile

**Multi-purpose**

Use this panel for content, images, diagrams, or whatever else you want to include. You can use the line tool to divide this panel into multiple sections if you want.

# Two Column

## Content Tile

**Multi-purpose**

Use these slides for comparing two topics or just for splitting your content into multiple pieces.

**Multi-purpose**

Use these slides for comparing two topics or just for splitting your content into multiple pieces.

# Three Column

**Column 1**

**Column 2**

**Column 3**

# Four Column

**Column 1**

**Column 2**

**Column 3**

**Column 4**

# Half Panel

## Right aligned

**Open Content**

This space is great for supporting text that compliments whatever content is inside the panel.

**Panel Content**

This space can be for text content, images, diagrams, or whatever you need

# Half Panel

**Left aligned**

**Panel Content**

This space can be for text content, images, diagrams, or whatever you need

**Open Content**

This space is great for supporting text that compliments whatever content is inside the panel.

# ⅔ Panel

**Right aligned**

## Open Content

This space is great for supporting text that compliments whatever content is inside the panel.

### Panel Content

This space can be for text content, images, diagrams, or whatever you need

# ⅔ Panel

Left aligned

**Panel Content**

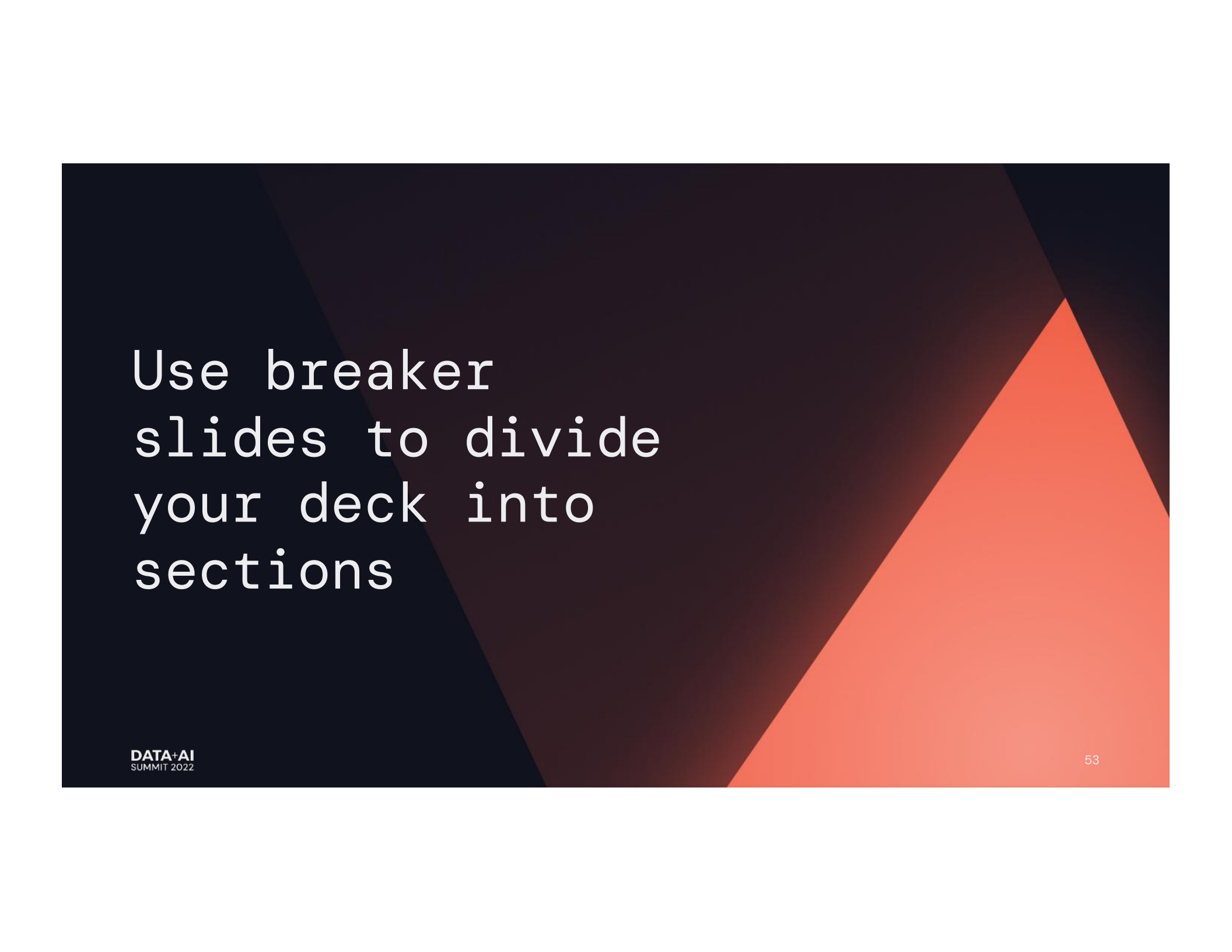This space can be for text content, images, diagrams, or whatever you need

**Open Content**

This space is great for supporting text that compliments whatever content is inside the panel.

DATA+AI
SUMMIT 2022

# Code Display

Paste snippets

```
1  // Open connection to SQL Server database
2  SQLServerConnection Conn;
3  Conn = new SQLServerConnection("host=nc-star;port=4100;User ID=test01;
4  Password=test01;Database Name=Test");
5  try
6  {
7  Conn.Open();
8  Console.WriteLine ("Connection successful!");
9  }
10
```

# Use breaker slides to divide your deck into sections

# Use breaker slides to divide your deck into sections

# Starter Shapes

Copy and paste these wherever you need them

Floating panel for text or graphics

Medium pill label

Medium pill label

SMALL PILL LABEL

SMALL PILL LABEL

# Logos

## Partners and cloud platforms

**DATA+AI**
SUMMIT 2022

# Logos

Open source projects

# DATA+AI
## SUMMIT 2022

# Thank you

**Your Name**
You Title