

From Data to Decisions: Building Al Agents That Take Action

Maria Zervou

©2024 Databricks Inc. — All rights reserved

Al is no longer a trend— It's a business

imperative

1. The State of AI in 2024, McKinsey & Company

2. From Potential to Profit: Closing the Al Impact Gap, <u>BCG Al Radar</u>

3. How Generative AI Is Changing The Future Of Work, Oliver Wyman Forum Generative AI Survey

65%

of organizations now use generative AI in at least one business function, nearly doubling from last year

McKinsey & Company¹

75%

of executives rank AI/GenAI as a top three strategic priority

Boston Consulting Group²

>50%

of employees say they use generative Al weekly at work

Oliver Wyman³

GEN AI IS AWESOME! IT'S EVERYWHERE, RIGHT?

RIGHT??

OK... why not?

Challenge: Building and deploying production-quality Gen Al solutions

90%

of enterprises *not* confident going to production



Quantity demanded more than doubles

Mosaic's Law



The amount of learning per dollar will quadruple every year.

4x YoY for 20 years = 10 orders of magnitude to biological parity Why haven't cheap models caused this disruption?

Jevons Paradox



decreased friction With increased officiency comes higher consumption



GenAl disrupts the cost of converting intention \rightarrow action

USER INTERFACE

General Intelligence

Consumer models trained on a broad dataset disconnected from your business data



Data Intelligence

Al that reasons over your enterprise data and is able to solve domain-specific use cases

Difficult to build AI Agents that accurately and securely work on your data



Cannot reason over your data

Agents that lack understanding of enterprise data will not respond accurately Understanding quality and remediation for your custom use case is complex

Difficult to

evaluate accuracy



Inability to fully govern AI apps

Ensuring access controls, guardrails, lineage, and cost controls is disjointed

Agents that reason over your data





Governance across data, models and tools



A unified platform to build AI Agents



Agent Lifecycle

Let's walk through the Agent lifecycle



Agents with Tools



Increasing Capability, but also Complexity

Customer Challenges

How do you **reuse** existing production ready assets **without recreating the wheel?**

Agents with function/tool calling

Some LLMs can be prompted to call functions or tools, to build complex agents.



Step 1: Prepare data, create tools



Why Tools Matter for Al Agents Discussion Points

- 1. What **existing** ML models, APIs, or data sources in your organization could be **exposed as functions** for AI agents?
- 2. How would you determine which processes should remain human-led versus agent-led?
- 3. What level of autonomy would you be comfortable giving to Al agents in your organization?

Agents Iteration

Customer Challenges

- How to enable developers of all skill levels to get started quickly?
- 2. How to **rapidly iterate** *without needing to* touch code, but yet, still have access to the code when I need more control to improve quality?

Step 2: Rapid prototyping w/ LLM judge quality checks

Agent Evaluation	
2 Evaluate prototype	
Agent code	
	Unity Catalog
2 Prototype Agents	Tables
Use Tools	

Step 2: Rapid prototyping w/ LLM judge quality checks



Al Playground lets you

any code

Step 2: Rapid prototyping w/ LLM judge quality checks

Quickly assess initial quality with built-in LLM judges from **Agent Evaluation**



Al Playground lets you

any code

an Agent prototype without

GenAi's biggest challenge is <u>evaluation</u>!

Start with a benchmark!

Customer challenges

- 1. Improve quality with a representative evaluation dataset without SMEs spending months labeling?
- 2. **Collect high-quality feedback** from SMEs who are busy and have limited time?

Step 3: Label evaluation dataset



Build your eval dataset

Eval dataset

	questions		ground_truth
21	> What is the requirement for		> The requirement for enabling a
22	> How can users access the Id	1.1	> Users can access the Ideas Po
23	> How can I disable personal a	т.	> To disable personal access tok
24	> How can I set a maximum lif		> To set a maximum lifetime for
25	> What values do I need to sp		> In the API call for log delivery,
26	> How long should I expect to		> After initial setup or other log
27	> How do I create an IAM role		> To create an IAM role for audit

©2024 Databricks Inc. — All rights reserved

Time for Evaluation

Customer challenges

Iteratively identify & fix quality issues

- What are the **right metrics to evaluate quality?** How do I trust the 1. outputs of these metrics?
- How do I quickly identify the **root cause of quality problems?** 2.
- I need to evaluate many ideas how do I... 3.
 - ...run evaluation quickly so the majority of my time isn't spent waiting? 0
 - ...track these different versions of my agent? 0
 - ...quickly debug my Agent's logic? 0
- How do I easily **compare 2 versions** of my Agent to understand if my fix 4. worked?
- What tools do you wish existed for evaluating Al agents?" 5. Would you use MLflow for AI agent evaluation? What features would make it better?

Agents vs Traditional ML Eval

- 1. Traditional ML evaluation \rightarrow Precision, recall, RMSE, F1-score.
- 2. All agent evaluation needs additional metrics:
 - Function-call success rate
 - 🗸 Response latency
 - 🖌 Correctness of retrieved data
 - V Error handling (does it fail gracefully?)
 - Action impact (did it take the right step?)

MLFlow for evaluation



Build your eval dataset

Eval dataset



©2024 Databricks Inc. — All rights reserved

Grade & compare both answers using MLFlow Evaluation

LLM-as-a-Judge metrics with MLFlow



MLFlow.evaluate() API with GenAl metrics

Here is a question with an answer: <question - answer> You're asked to judge the answer given this expected solution <label>:

- Is the answer provided correct (note 1-5)?
- Is the answer professional (note 1-5)?
- any custom GenAl metrics

ବ୍ଦ

Judge LLM

	Baseline model	New Model
correctness	4.1 (avg)	4.5 (avg)
professionalis m	3.5 (avg)	4.2 (avg)
top 10 less correct answers	[]	[]

Step 4: Iteratively identify & fix quality issues



Customer challenges

Pre-production

- Quickly create a Chat UI for stakeholders to test the agent?
- 2. Track each piece of feedback?
- 3. Understand 100s of pieces of feedback at-scale?

Step 5: Release to pre-production



Time for Monitoring

Customer challenges

Deploy to production & monitor quality

- 1. How do I host my Agent as a production ready, scalable service?
- 2. How do I execute tool code securely and ensure it respects my governance policies?
- 3. How do I enable **telemetry / observability** in development and production?
- 4. How do I monitor my Agent's quality at-scale in production?

Step 6: Release to production



Agent Governance

Governance across data, models and tools



Unity Catalog + Al Agents & Governance Discussion Points

- 1. What **governance challenges do you anticipate** when deploying Al agents that can access enterprise data?
- 2. How do you currently **manage access controls** for data that Al systems might need?
- 3. What **audit and lineage tracking** capabilities would you need for Al agents in production?

Let's Talk Use Cases

Agent Use Case Identification Discussion Points

- 1. Have you or your team built Al agents before? If yes, what use cases?
- 2. What are the **biggest risks or limitations** you see with Al agents?
- 3. If an AI agent made decisions for your business, what would you need to trust it?

How Do I Find the FIRST Use Case?

©2024 Databricks Inc. — All rights reserved

Agent Use Case Identification Discussion Points

- 1. What are some of the biggest pain points your team faces today?
- 2. Where do you see the most repetitive or time-consuming tasks?
- 3. What decisions would you want an AI agent to automate or support?
- 4. Where do you currently lose time or resources due to lack of automation or insight?

Agent Use Case Identification

Best Practices

- 1. **Customer-Centric Approach:** Start by interviewing stakeholders to uncover pain points.
- 2. Mapping Out Processes: Visualize workflows to identify inefficiencies and automation opportunities.
- 3. Data Availability Check: Before pursuing a use case, ensure data is accessible and of sufficient quality.

Ensuring Solutions Address Real Problems

Ensuring Solutions Address Real Problems Best Practices

- 1. User-Centric Design: Involve end-users throughout development, not just at the end.
- 2. Feedback Loops: Build mechanisms for continuous feedback collection and monitoring.
- 3. **Iterative Testing:** Deploy in controlled environments before scaling to production.

Key Takeaways

01	Identify use cases -start with the simplest POC
02	Define success criteria with Agent Evaluation
03	Iterate Fast and Get Feedback

Thank you . . . !