



Building Breakthrough AI Applications with Not Diamond

Not Diamond × Dagster

Feb 11 2025

Table of contents

01

Intros

02

LLM Routing

03

Not Diamond for improved LLM usage

04

The Not Diamond Dagster Integration

05

Integration Usage Examples

06

LLM Best Practices

07

Next Steps

08

Q & A

Introductions

Meet the presenters



Tomás Hernando Kofman

Co-founder & CEO

Not Diamond



Alejandro Companioni

Founding Engineer

Not Diamond



Colton Padden

Developer Advocate

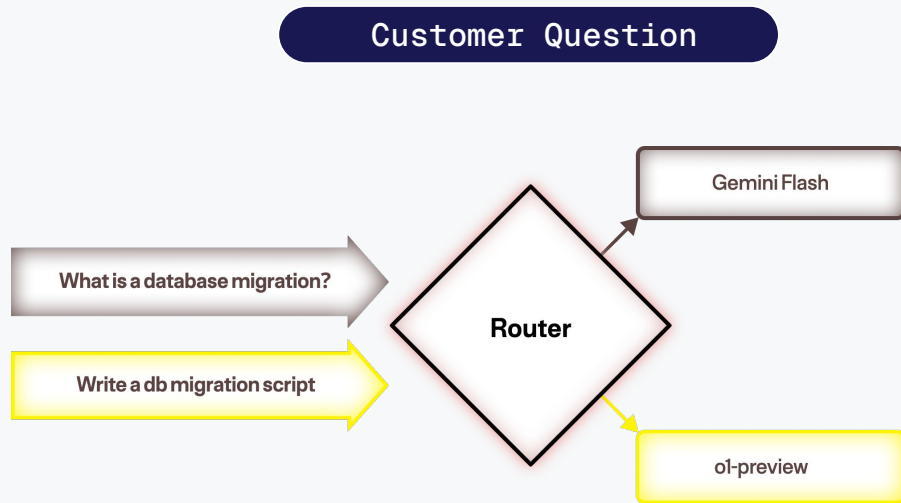
Dagster Labs

LLM Routing

LLM Routing

The best model on every input

As the number of models increases exponentially, selecting the right model for each request becomes increasingly challenging. LLM routing solves this problem by intelligently directing each query for you in real time, outperforming any individual model while reducing costs and latency by using smaller models when doing so doesn't degrade quality.



Why LLM routers are the future of prompting

Intelligent Model Usage

Always use the most effective model

Not every model is created equally. LLM routing ensures that your prompt is automatically directed to the model to produce the best result. This replaces suboptimal human heuristics with data-driven model selection.

Optimize for quality, cost, and speed

LLM performance varies significantly based on the type of request and costs can quickly escalate based on your model choice. Using an LLM router allows you to weigh your preferences between quality, cost, and speed.

AI Engineering & Data Engineering

AI engineering *is* data engineering

Data engineering tools lead to resilient AI pipelines

Working with LLMs and model inference fits nicely into the traditional tools of data pipeline building, allowing AI engineers to take advantage of a plethora of tools.

LLMs can empower data engineers

LLMs give data engineers the ability to make sense of unstructured organic data

It's easier than ever to make sense of messy and organic datasets. By introducing a step in your data processing pipelines that call out to AI services, engineers can add structure to data that was traditionally difficult to process.

How are data engineers using LLMs today?

Extracting insights

Finding trends and anomalies in data

Parsing unstructured data

Giving structure to non-deterministic organic data and classifying parsed data

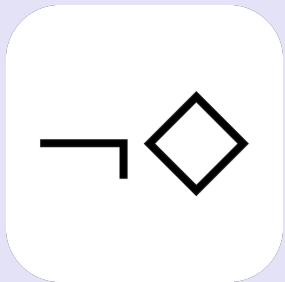
Documentation and code generation

Metadata for data cataloging, like generated asset descriptions.

Troubleshooting

Debugging of failures during operation, or in development

Not Diamond



The backstory of Not Diamond

- Not Diamond was founded on the bet that the future of AI won't have a single, giant model. Instead, we'll have a decentralized ecosystem of highly specialized models. This will lead to higher-performing, more computationally efficient, and safer AI.
- As we pursued this bet, we increasingly saw teams at both Global 500s and leading startups struggling to figure out when to use which model.
- The more developers we talked to, the more conviction we built that multi-model infrastructure would be critical to moving off traditional AI workflows that waterfall development through one model at a time and towards data-driven, multi-path development.
- Our team is made of serial founders, AI researchers, and ML engineers with decades of experience in the field, and we're backed by leading AI engineers and scientists like Jeff Dean, Julien Chaumond, Ion Stoica, and Tom Preston-Werner.

Features



Model Router

Leverage intelligent, real-time routing that dynamically recommends the best model for each query to maximize quality, cost-efficiency, and latency, delivering optimal performance at scale.



Model Gateway

Instead of orchestrating every request to each LLM provider separately, you can use Not Diamond's gateway to seamlessly leverage all of the most popular models with a unified API interface.

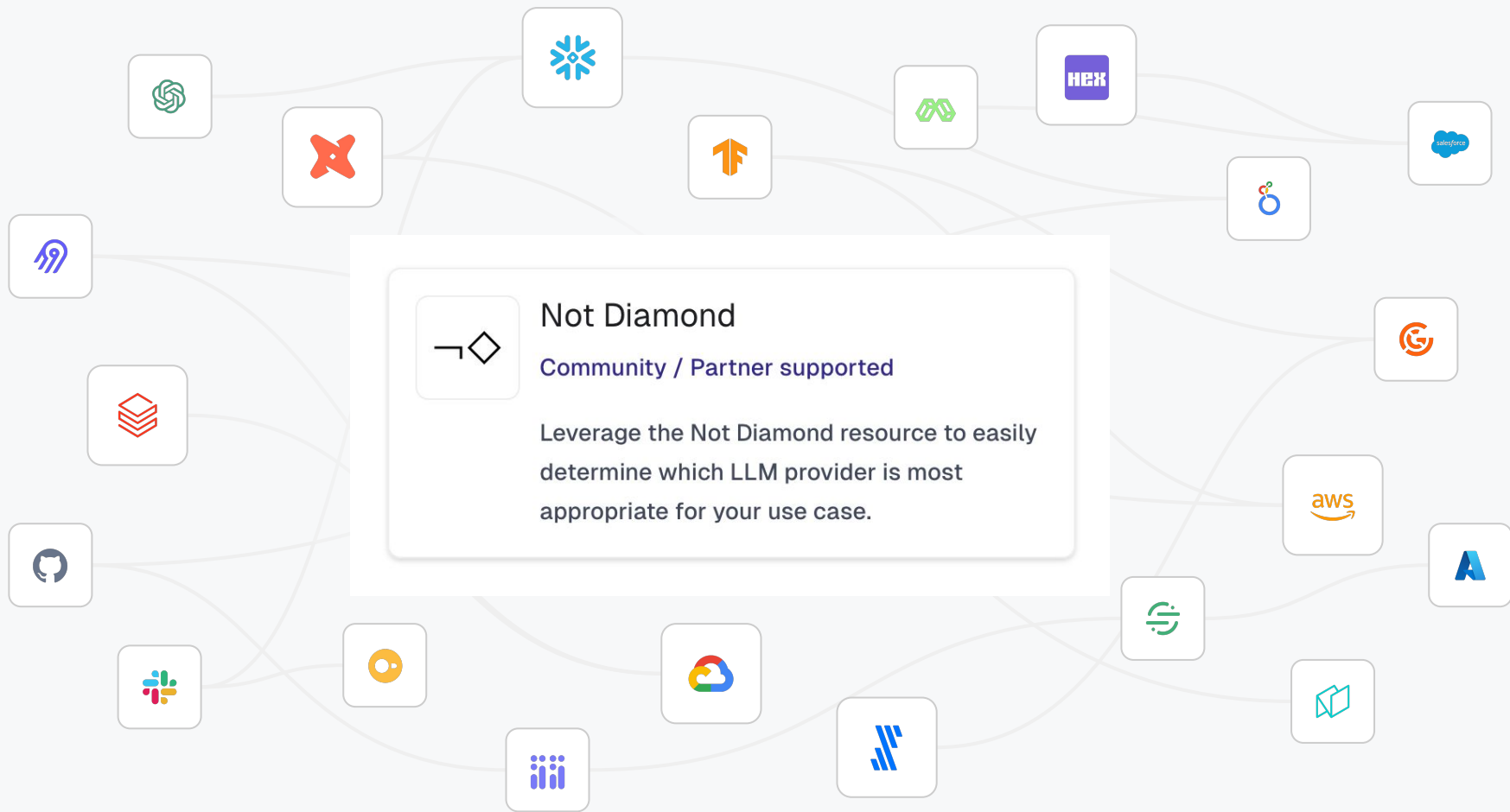


Beta

Prompt Adaptation

Model selection and prompting are very tightly coupled. Not Diamond can help you optimize your prompt to each model to ensure you're using the best model with the best prompt on every input.

The `dagster-notdiamond` Integration



pip install dagster-notdiamond



The Not Diamond Resource

The `dagster-notdiamond` integration provides a resource that wraps the Not Diamond client.

This allows you to easily leverage the power of Not Diamond model routing from your assets, determining the best model for a prompt, and using that in tandem with integrations like `dagster-openai`, `dagster-anthropic`, or other providers.

However, Not Diamond can also be used directly from the OpenAI client through the *model gateway*.



```
1 import dagster_notdiamond as nd
2
3
4 notdiamond_resource = nd.NotDiamondResource(
5     api_key=dg.EnvVar("NOTDIAMOND_API_KEY")
6 )
7
8 defs = dg.Definitions(
9     assets=[...],
10    resources={
11        "notdiamond": notdiamond_resource
12    }
13 )
14
15
16
17
18
19
20
21
22
23
24
```

Integration usage: **Model routing**

Model routing

Both the `NotDiamondResource` and `OpenAIResource` are registered for use in our assets.



```
1 import dagster as dg
2 import dagster_notdiamond as nd
3 import dagster_openai as oai
4
5
6 defs = dg.Definitions(
7     assets=[book_review_data, book_reviews_summary],
8     resources={
9         "notdiamond": nd.NotDiamondResource(
10             api_key=dg.EnvVar("NOTDIAMOND_API_KEY")
11         ),
12         "openai": oai.OpenAIResource(
13             api_key=dg.EnvVar("OPENAI_API_KEY")
14         ),
15     },
16 )
17
18
19
20
21
22
23
24
```

Model routing

One approach to using Not Diamond is in tandem with integrations with other providers, for example [dagster-openai](#).

Let's imagine we have a dataset of book reviews that we would like to summarize.



```
1 @dg.asset(kinds={"python"})
2 def book_review_data(context: dg.AssetExecutionContext) -> dict:
3     data = {
4         "title": "Cat's Cradle",
5         "author": "Kurt Vonnegut",
6         "genre": "Science Fiction",
7         "publicationYear": 1963,
8         "reviews": [
9             {
10                 "reviewer": "John Doe",
11                 "rating": 4.5,
12                 "content": "A thought-provoking satire...",
13             },
14             {
15                 "reviewer": "Jane Smith",
16                 "rating": 5,
17                 "content": "An imaginative and darkly...",
18             },
19             {
20                 "reviewer": "Alice Johnson",
21                 "rating": 3.5,
22                 "content": "Intriguing premise but...",
23             },
24         ],
25     }
26     context.add_output_metadata(metadata={"num_reviews": len(data.get('reviews', []))})
27     return data
```

Model routing

One approach to using Not Diamond is in tandem with integrations with other providers, for example [dagster-openai](#).

Let's imagine we have a dataset of book reviews that we would like to summarize.



```
1 @dg.asset(  
2     kinds={"openai", "notdiamond"}, automation_condition=dg.AutomationCondition.eager()  
3 )  
4 def book_reviews_summary(  
5     context: dg.AssetExecutionContext,  
6     notdiamond: nd.NotDiamondResource,  
7     openai: oai.OpenAIResource,  
8     book_review_data: dict,  
9 ) -> dg.MaterializeResult:  
10  
11     prompt = f"""  
12     Given the book reviews for {book_review_data["title"]}, provide a detailed summary:  
13  
14     {'|'.join([r['content'] for r in book_review_data["reviews"]])}  
15     """  
16  
17     ...  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38
```

Model routing

We pass our `prompt` to a method called `model_select` provided by the `notdiamond` resource.

This returns the best model to use when optimizing for `cost` as we specified `tradeoff="cost"`.



```
1 @dg.asset(
2     kinds={"openai", "notdiamond"}, automation_condition=dg.AutomationCondition.eager
3 )
4 def book_reviews_summary(
5     context: dg.AssetExecutionContext,
6     notdiamond: nd.NotDiamondResource,
7     openai: oai.OpenAIResource,
8     book_review_data: dict,
9 ) -> dg.MaterializeResult:
10
11     ...
12
13     with notdiamond.get_client(context) as client:
14         start = time.time()
15         session_id, best_llm = client.model_select(
16             model=["openai/gpt-4o", "openai/gpt-4o-mini"],
17             tradeoff="cost",
18             messages=[
19                 {"role": "system", "content": "You are an expert in literature"},
20                 {"role": "user", "content": prompt},
21             ],
22         )
23         duration = time.time() - start
24
25     ...
26
27
28
29
30
31
32
33
```

Model routing

Then, we take the `best_llm.model` returned by Not Diamond, and supply that to our call to OpenAI.

This allows us full control over the result of the model.

And in our materialized result, we include metadata around the LLM chosen.



```
1 @dg.asset(
2     kinds={"openai", "notdiamond"}, automation_condition=dg.AutomationCondition.eager
3 )
4 def book_reviews_summary(
5     context: dg.AssetExecutionContext,
6     notdiamond: nd.NotDiamondResource,
7     openai: oai.OpenAIResource,
8     book_review_data: dict,
9 ) -> dg.MaterializeResult:
10
11     ...
12
13     with openai.get_client(context) as client:
14         chat_completion = client.chat.completions.create(
15             model=best_llm.model,
16             messages=[{"role": "user", "content": prompt}],
17         )
18
19     summary = chat_completion.choices[0].message.content or ""
20
21     return dg.MaterializeResult(
22         metadata={
23             "nd_session_id": session_id,
24             "nd_best_llm_model": best_llm.model,
25             "nd_best_llm_provider": best_llm.provider,
26             "nd_routing_latency": duration,
27             "summary": dg.MetadataValue.md(summary),
28         }
29     )
30
31
32
33
34
```

Model Routing

Example #1

Materialize selected

book_review_data

No description

Materialized Feb 5, 4:44 PM

Python

book_reviews_summary

No description

Materialized Feb 5, 4:44 PM

Open AI notdiamond

book_reviews_summary

[View in Asset Catalog](#)

Latest materialization

Run	Run 7ef4bc95	View logs
Timestamp	Feb 5, 4:44 PM	
nd_session_id	c56f8fbc-e232-44f8-8b0c-c8a91d40ef42	
nd_best_llm_model	gpt-4o-mini	
nd_best_llm_provider	openai	
nd_routing_latency	0.9765541553497314	
summary	[Show Markdown]	
openai.calls	1	
openai.total_tokens	357	
openai.prompt_tokens	122	
openai.completion_tokens	235	

Example #2

Integration usage: **Model gateway**

Model gateway

Only the `OpenAIResource` is required, but we configure the `api_key` and `base_url` to point to Not Diamond.



```
1 import dagster as dg
2 import dagster_openai as oai
3
4
5 defs = dg.Definitions(
6     assets=[book_review_data, book_reviews_summary],
7     resources={
8         "openai": oai.OpenAIResource(
9             api_key=dg.EnvVar("NOTDIAMOND_API_KEY"),
10             base_url="https://proxy.notdiamond.ai/v1/proxy",
11         ),
12     },
13 )
14
15
16
17
18
19
20
21
22
23
24
25
26
27
```

Model gateway

The model selection is handled for you!

You can provide `extra_body` parameters for additional options provided to Not Diamond.



```
1 @dg.asset(
2     kinds={"openai", "notdiamond"}, automation_condition=dg.AutomationCond
3 )
4 def book_reviews_summary(
5     context: dg.AssetExecutionContext,
6     openai: oai.OpenAIResource,
7     book_review_data: dict,
8 ) -> dg.MaterializeResult:
9
10     ...
11
12     with openai.get_client(context) as client:
13         chat_completion = client.chat.completions.create(
14             model="notdiamond",
15             extra_body={
16                 "models": ["gpt-4o", "claude-3-5-sonnet-20240620"],
17                 "tradeoff": "cost",
18                 "preference_id": "YOUR_PREFERENCE_ID"
19             },
20             messages=[{"role": "user", "content": prompt}],
21         )
22
23     ...
24
25
26
27
28
29
30
```

Discussion on LLM usage and best practices

“

Model routing has led to a 14% increase in our LLM output quality and a 72% reduction in our inference costs.

Tilen Babnik

CTO of Samwell

Summary

Summary



Not Diamond

Not Diamond provides a method for choosing the best LLM for your given use case.

The [dagster-notdiamond](#) integration makes it easy to leverage this service.



Model Routing

With the integration, you can provide your prompt, and Not Diamond will determine which model is most suitable, which you can then pass to your LLM provider of choice.



Model Proxying

Not Diamond also offers a proxying service, which automatically handles the routing, and provides requests for you behind the scenes.

References & Resources

01

Not Diamond quickstart

docs.notdiamond.ai/docs/quickstart

02

Integration reference

docs.dagster.io/integrations/libraries/notdiamond

03

Integration source code

[github.com/dagster-io/community-integrations/.../dagster-notdiamond](https://github.com/dagster-io/community-integrations/tree/main/dagster-notdiamond)

Q & A





Subtitle goes here

Thank you!



The advantage of data-driven LLM usage

Intelligent Model Usage

Optimize routing with your data

Like everything in AI, data is the driving force of every outcome. Create a custom router by providing your dataset to tune the various paths each request can take. Let the router do the hard work with some subtle guidance.

Real-time tuning for model routing

Part of being human is having preferences and current models don't do a great job of personalizing the outcome. Let your users' preferences influence the model selection by reacting to real-time feedback.