# 11/5/25: Evaluating Launch Darkly AI Configs

## AI Model Config

Config for a particular model deployment

<mark>endpoint</mark>

<mark>deployment name</mark>

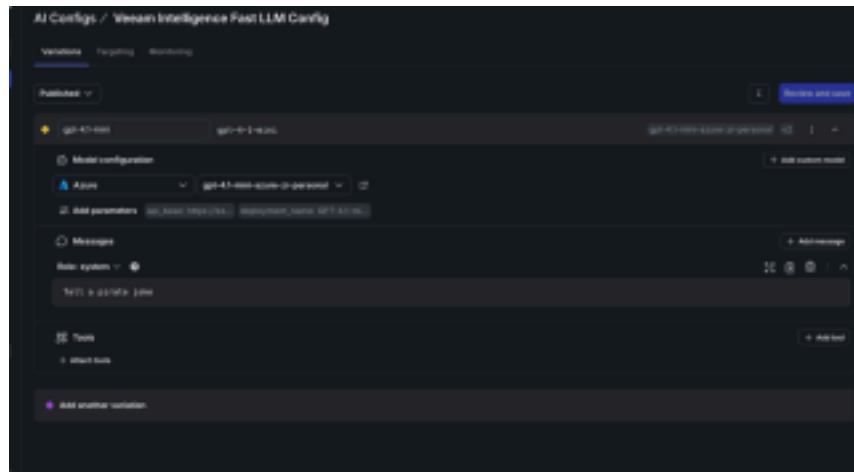> <mark>temperature</mark> (maybe – we have the option to override the temperature depending on the use case)



## AI Config

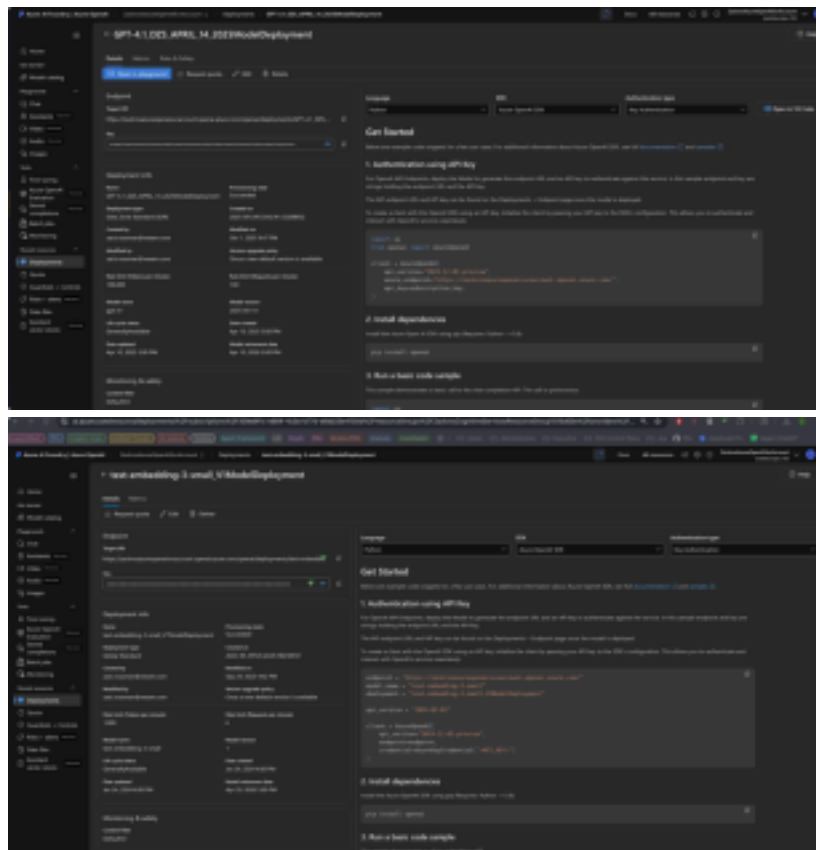Config for a particular LLM-based query

Which AI model config to use

Prompts to be fed to the AI model config

Creating AI Model Configs

First, let's create AI model configs that correspond to the model deployments in our Azure Open AI Service account.

Example model deployment details from Azure Open AI Service





In this it seems AWS Bedrock has a standard endpoint per region and the region needs to be pulled in from the example code env vars (it's not part of the Launch Darkly Provider class. Azure Open AI Service does things different than AWS here Bedrock – each Azure Open AI Service account has its own endpoint.

We have model deployments in each VDC environment (personal, dev, stage, prod).See Azure OpenAI Configuration | Quot

a Allocations for Each Environment

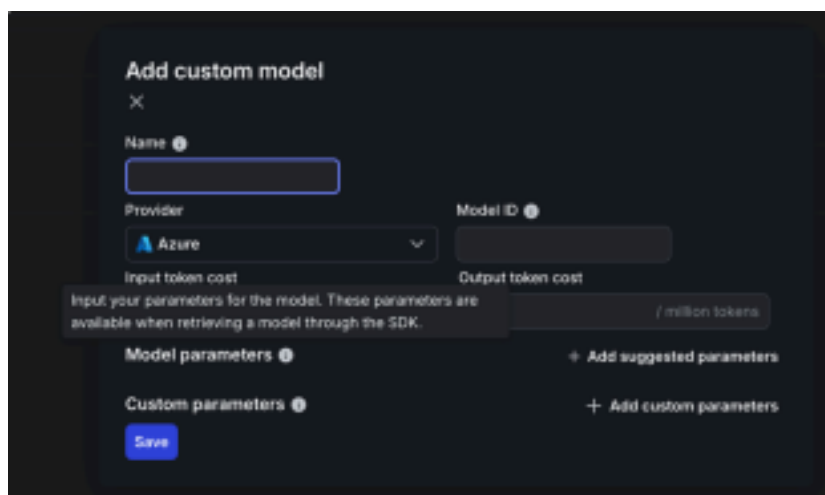for more details. We'll need to create AI model configs for the model deployments in each

env because Azure Open AI Service creates a unique endpoint for each Azure Open AI Service account. We'll add the api base endpoint, the deployment name, and maybe even the api version as a Custom Parameter.

The endpoint and deployment_name params would probably be a better fit in Model Params rather than Custom Params because these are core params for all of our Azure LLM deployments.
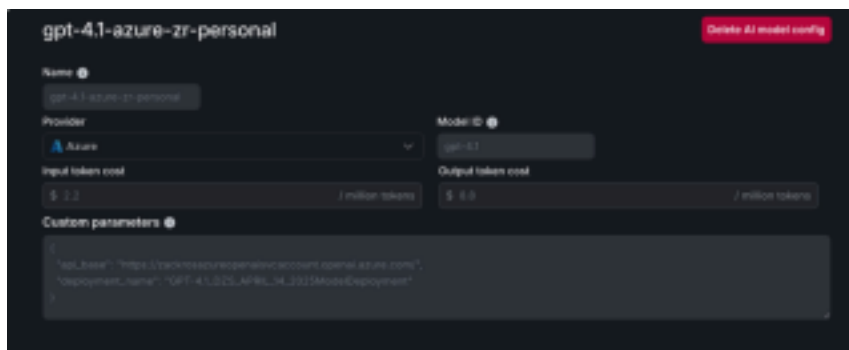
After talking to Launch Darkly solutions engineers, they said this was because AWS Bedrock has a consistent endpoint for all models deployed in the same region (see this ).

[example code](#)

Not ideal but at the same time, it's hard to generalize across all model providers and the Custom Params option gives us the flexibility to add the info we need for each AI Model Config.
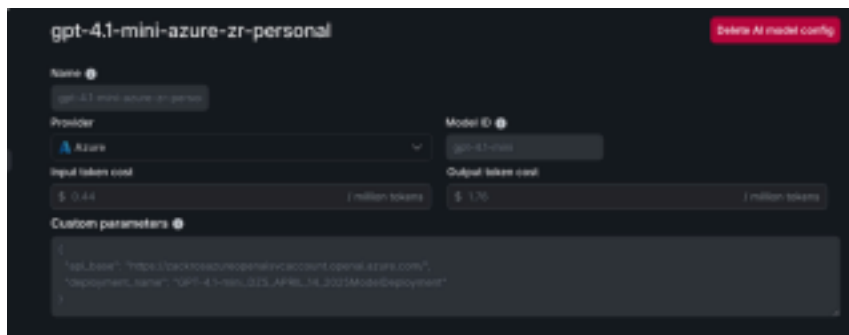






Some AI model providers support "model params" and "custom params" (e.g. Azure) while others only support "custom params" (e.g. Open AI)?

(minor feedback) why can't we edit the Model Config after creation? Have to delete and recreate if we want to change the name, token input/output cost, or the custom params



Creating AI Configs

I created two "completion based" AI config (rather than "AI agent" AI config) because this is the best fit for the Veeam Intelligence workflow.

Details on "Completion-Based" AI configs vs. "Agent-Based" AI configs

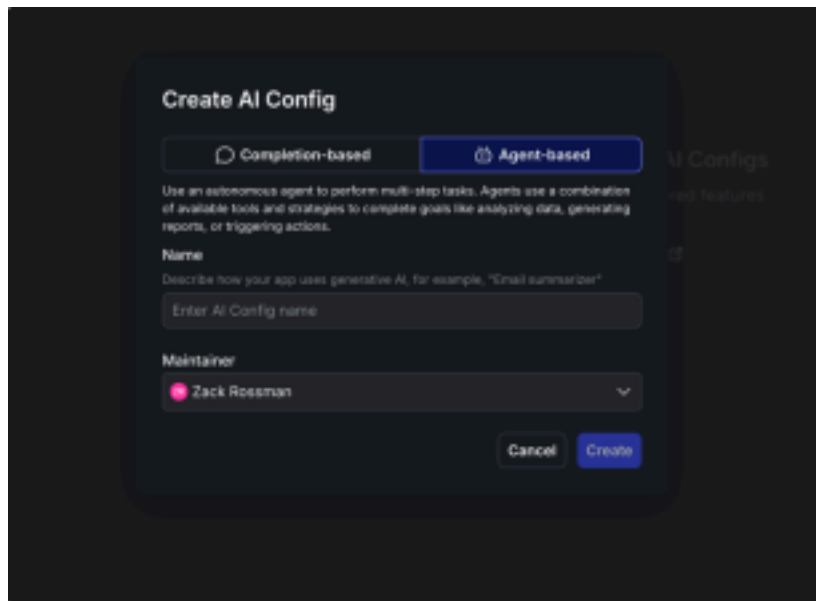Quickstart for AI Configs | LaunchDarkly | Documentation

AI Configs can be created in two modes:
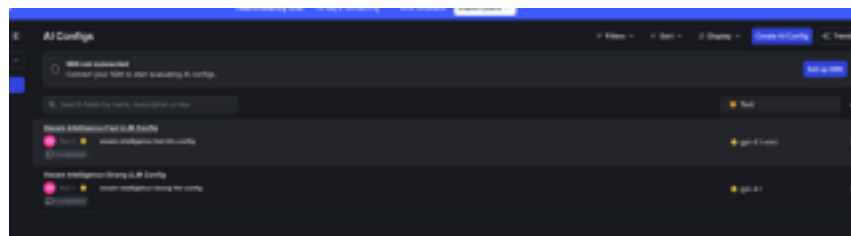
Prompt-based: Configure prompts with messages and roles.

Agent-based: Configure multi-step workflows with instructions and tools. To learn more, read Agents in AI Configs.

Agents in AI Configs | LaunchDarkly | Documentation

Use agent-based AI Configs to define, monitor, and
adapt structured model workflows in real time. Agents
use the same configuration model as standard AI
Configs, including environment targeting, variation
management, and monitoring. The difference is that
agents define how a model performs reasoning and takes
actions across multiple steps instead of responding to a
single prompt.



One AI config for the "strong" model for handling more complex queries, and another for the "fast" model for handling simpler queries.

Get SDK Key

This was tricky. At first I thought the application code needed to use an API key to fetch the AI configs. But actually we need to use an SDK key to fetch the AI configs.
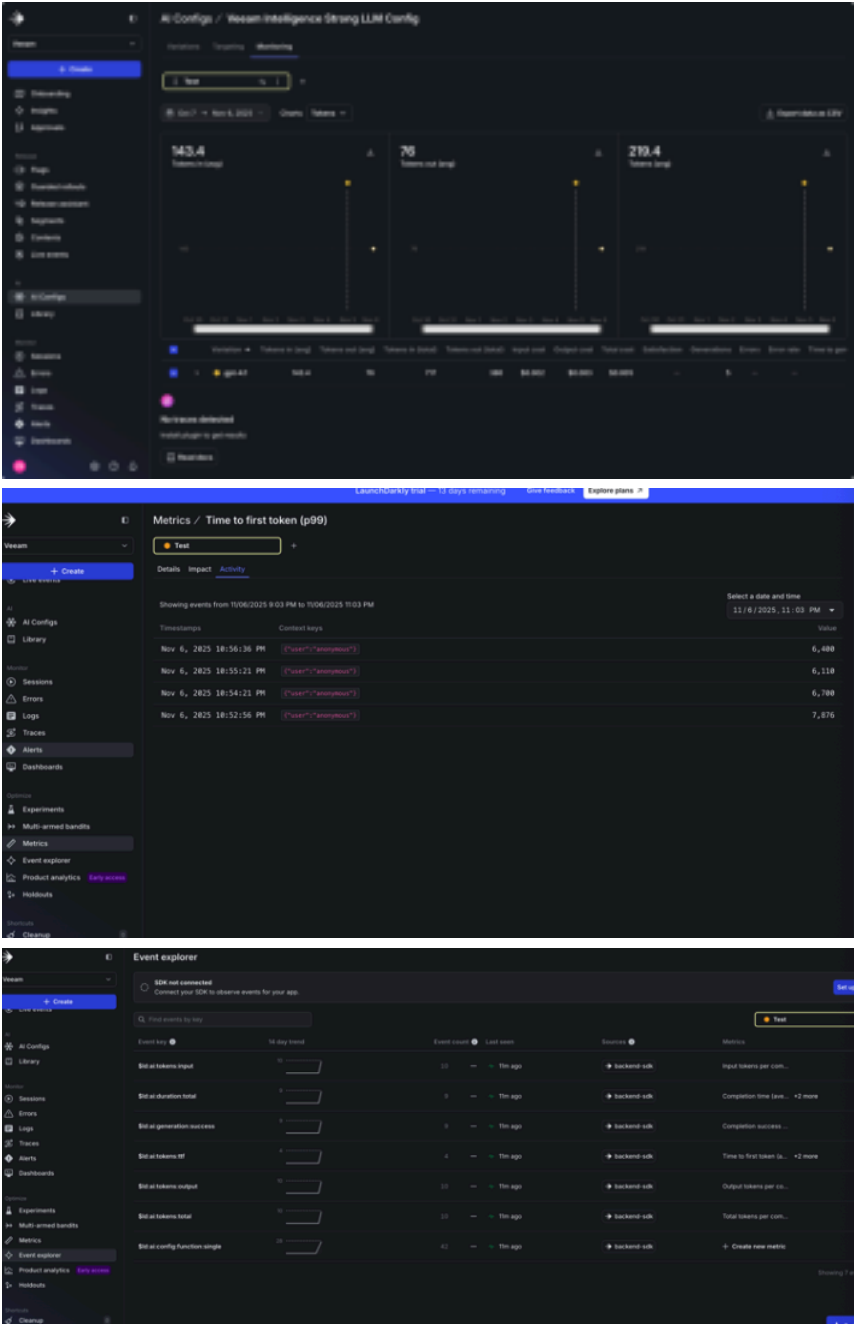


Tracking Metrics

I added a wrapper to Veeam Intelligence to track various metrics
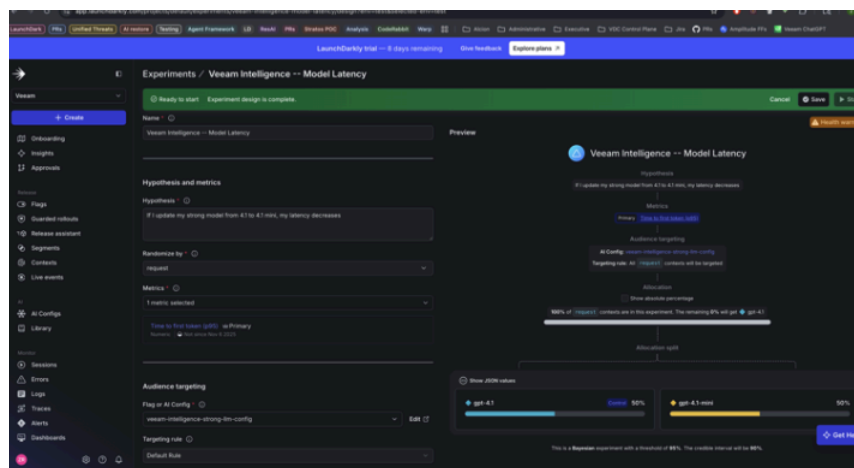
Generation count - Number of successful completions

Input tokens - Tokens sent to model

Output tokens - Tokens generated by model

Duration - Time to generate

Time to first token - Latency metric
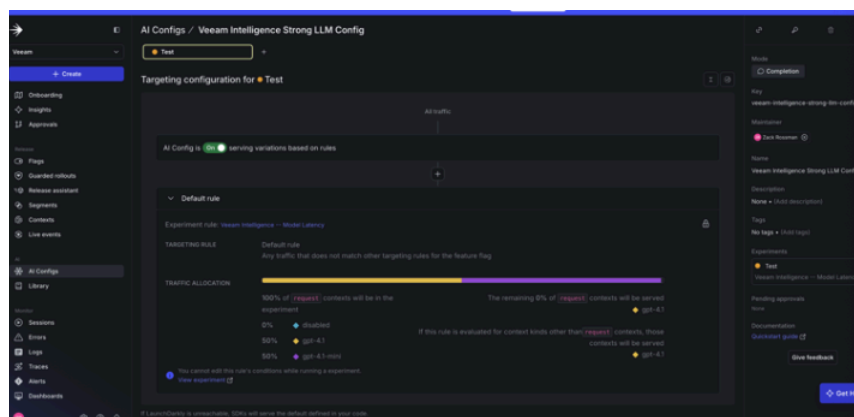
Errors - Failed generations







Experiments

I created an experiment to understand – can we quantify the difference in "time to first token" metric if we use a lighter weight
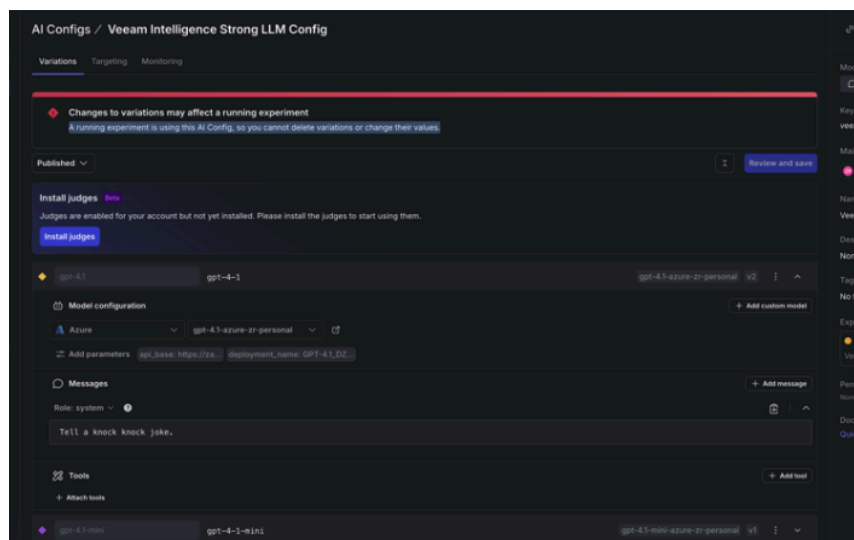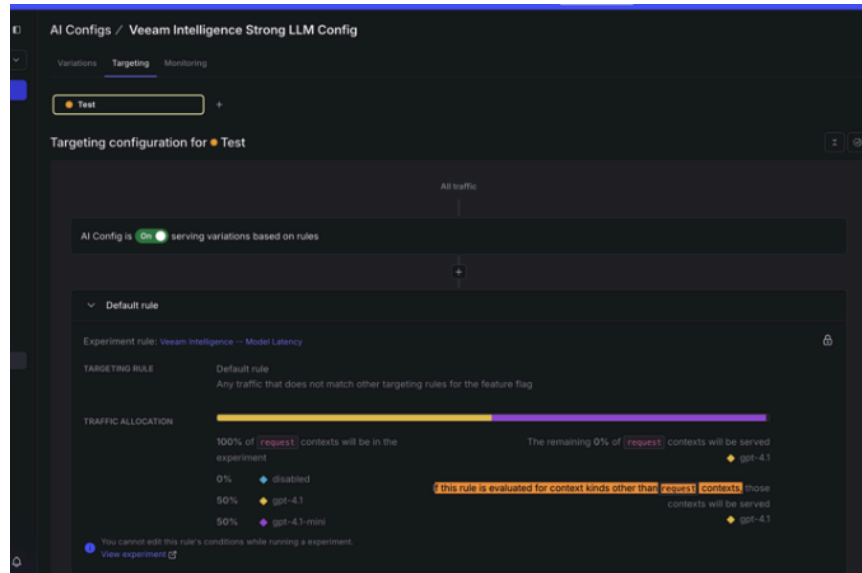
model like GPT 4.1-mini rather than GPT 4.1?



Once the experiment was started, the targeting rules for my "strong" AI config were automatically updated according to the 50/50 allocation split for my experiment.



Note – when an AI config is being used as part of an experiment, you can't delete variations or change their value. This is smart, don't want to subvert the results



At first, I noticed that the experiment wasn't tracking the results for requests and the AI config variations weren't being distributed 50/50. Rather, all 100% of requests were getting the gpt-4.1 model config. The issue was nuanced – my code wasn't using the correct "context kind"
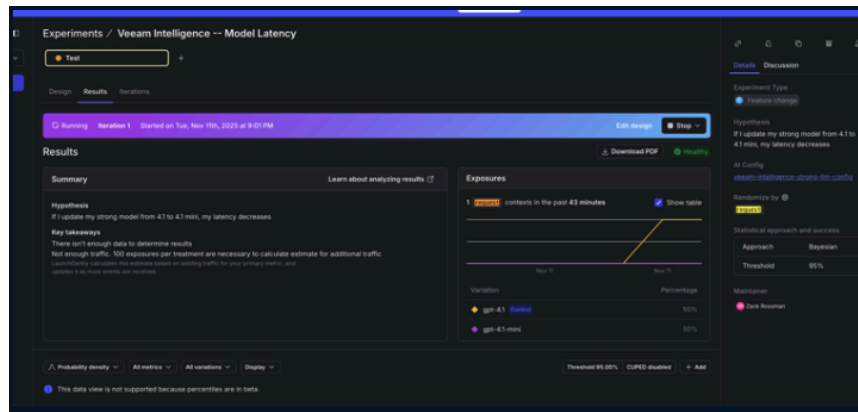
I updated my code to use the `request` context kind



Hmmm… but even with the context `kind` set to `request` , Launch Darkly was still serving the default `gpt-4.1`☐

Debugging the issue with Warp AI agent

I see the issue! The context key is always "anonymous" (or the same user_id), which means Launch Darkly sees it as the same context across all requests. For per-request randomization in experiments, you need a unique key per request.

The experiment's traffic allocation uses the context key to determine which variation to serve. If the key is always the same, it will always hash to the same variation (gpt-4.1 in your case).

Here's the fix - generate a unique key per request:
Perfect! Now each request will have a unique UUID as the context key, which allows Launch Darkly's traffic allocation algorithm to properly randomize between the experiment variations.

What changed:

1. Unique key per request: Now using uuid.uuid4() to generate a unique key for each request
2. user_id as attribute: Moved user_id from the key to an attribute, so you can still target by user if needed in the future
3. Proper randomization: With unique keys, Launch Darkly can properly hash and distribute traffic according to your 50/50 experiment split

Why this fixes it:

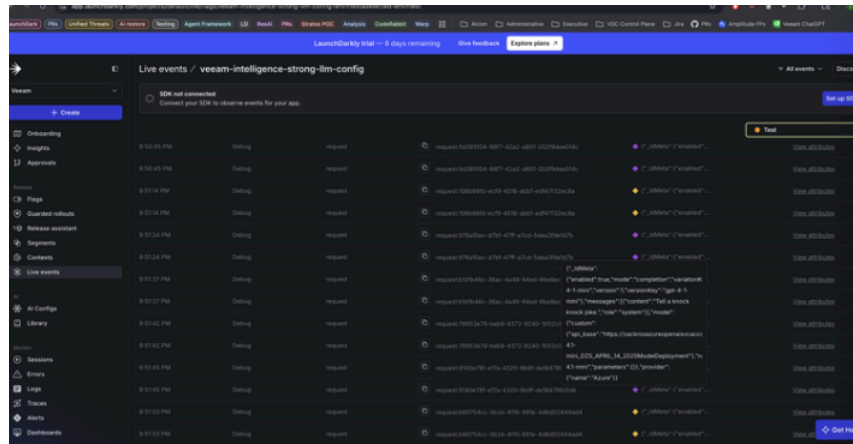Launch Darkly's traffic allocation uses the context key

to hash and determine which variation to serve

When you were using the same key ("anonymous") for all requests, it always hashed to the same bucket → always gpt-4.1

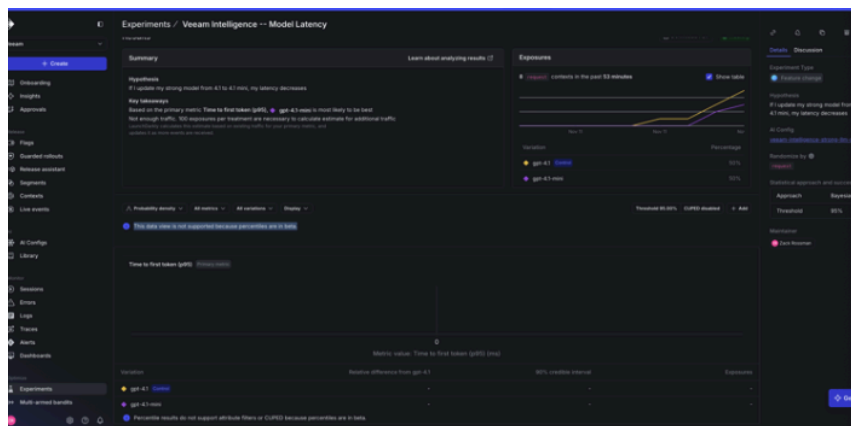With unique keys per request, each request gets independently randomized → true 50/50 split

Now your experiment should properly distribute traffic between gpt-4.1 and gpt-4.1-mini!

The issue was I needed to add a unique requestID for each AI config evaluation. This fixes it and, indeed, now we can see an even split between requests☐



Unfortunately, it seems the result view and the CUPED adjustment isn't supported for percentile analysis like p95☐

Creating feature change experiments | LaunchDarkly | Documentation

So I recreated the experiment where we're measuring average time to first token (rather than p95 time to first token)



I tried to create a new "iteration" of the existing experiment, but I had to start a new experiment altogether because the "primary single metric must match a metric in the metrics list". So the issue was that I created the initial experiment with a percentage-based metric as the primary metric.

When you stop an experiment, Launch Darkly asks to "ship" a single variation. Is this a "feature" or a "bug"? I like it because it forces us to be decisive.

This is one example of how the Launch Darkly experimentation platform is strongly opinionated.



I recreated the experiment with time to first token (average) as the primary metric

They have a nice "health check" feature in on the right hand side of the experiment UI. This is a really thoughtful feature – helps us build confidence in our experiment setup and operations☐



After 115 requests, 59 requests were handled by the gpt-4.1 variation and 54 requests were handled by the gpt-4.1- mini variation. We can see there's a 62% chance that gpt-4.1-mini is the superior variation if we want to pick the variation that reduces the average time to first token☐



Evaluations

Online evals use AI judges, which are a type of AI Config that evaluates another AI Config's output in real time. Judges apply a consistent evaluation prompt and scoring framework to measure each response on key metrics such as accuracy, relevance, and toxicity.

Scores appear automatically on the Monitoring tab for each variation, alongside latency, cost, and satisfaction metrics. This provides a continuous signal of model performance with real users and data, enabling teams to detect regressions, improve reliability, and apply guardrails within the AI Config workflow.

Online evals differ from offline or pre-deployment testing. Instead of running evaluations manually in a sandbox or against datasets, they measure AI Config quality continuously in production.

Online evals work alongside observability. Observability helps you view model responses and routing details, while online evals provide quality scores you can use to trigger alerts, manage guarded rollouts, or run experiments.
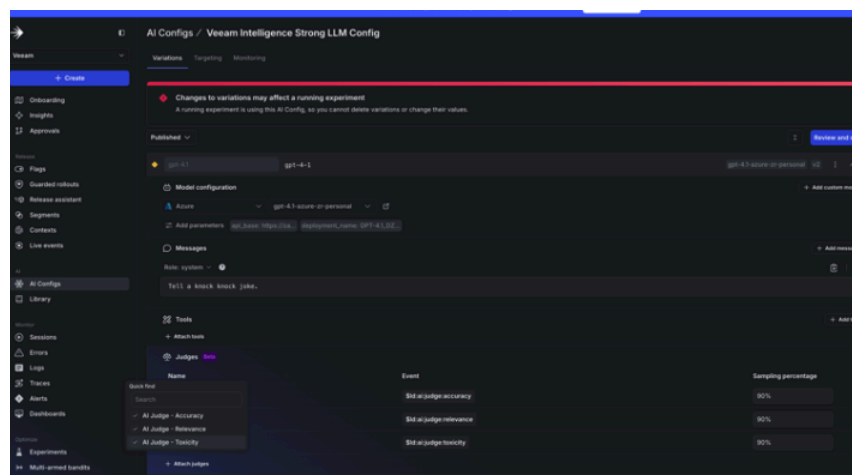
Use online evals to:

Continuously monitor AI quality in production

Detect regressions immediately after a rollout

Automate rollbacks or alerts based on evaluation metrics

Compare prompt or model variations using live performance data

I attached all three judges to my AI configs



The Launch Darkly judge evaluation capability is only supported in the_. This is a bummer because we primarily [JS SDK](#) use Python and Go for our backend AI services. But this is a temporary limitation, estimated that the Python SDK will have support within 3 weeks (by Dec 8, 2025 ).

Observability

[LLM observability | LaunchDarkly | Documentation](#)
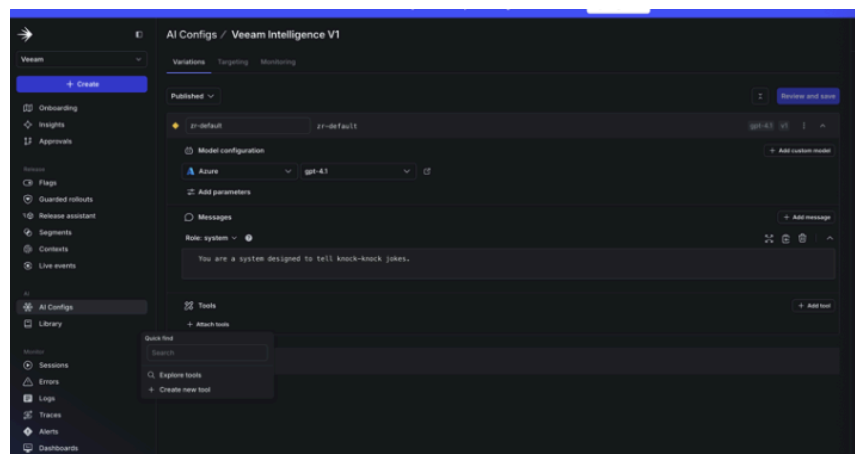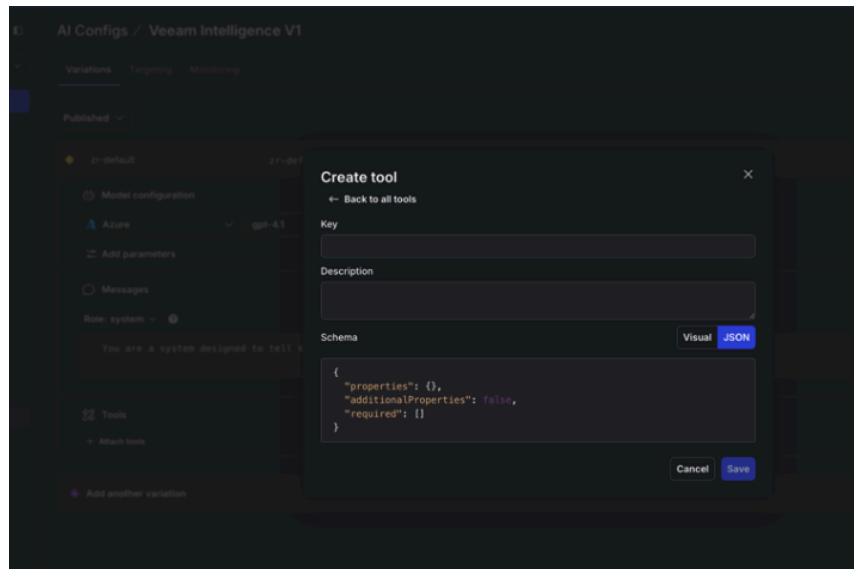
☐ TODO

Tracing

☐ TODO

Creating Tools

[Tools in AI Configs | LaunchDarkly | Documentation](#)

IMO this feel like they're trying to do too much. Trying to apply "dynamic configuration" to every aspect of AI workflows, including agents and tools. Not sure if they're well-positioned in the stack to dictate agentic workflows and tool orchestration, probably better suited for dedicated AI agent framework code/platform. But I like that Launch Darkly has a pulse on AI/agentic trends.

TODO

Test dynamic prompts with placeholder values in the AI config prompt

POC Code

Here's the PR with the updates needed to integrate AI configs into Veeam Intelligence for VDC ☐

https://github.com/veeam-ai/veeam-intelligence/pull/742 RESTRICTED CONTENT

Analysis

TODO Zack