# BECOMING AI ENGINEERS

Lawrence Jones

LDX London (June 2025)

# Hi!
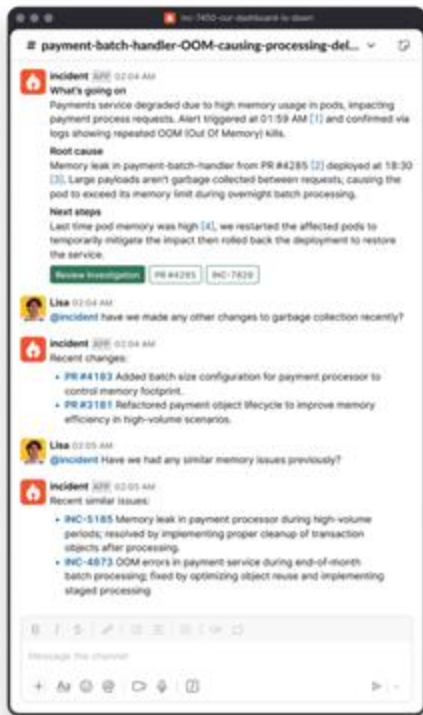
➔ Lawrence Jones (@lawrjones)

➔ Product Engineer at incident.io, leading AI efforts

➔ Joined three years ago as first hire

➔ Now 100 people, dual-HQ in London and SF
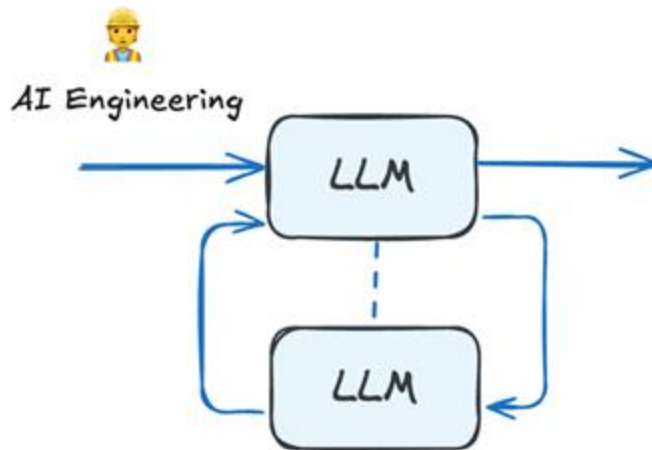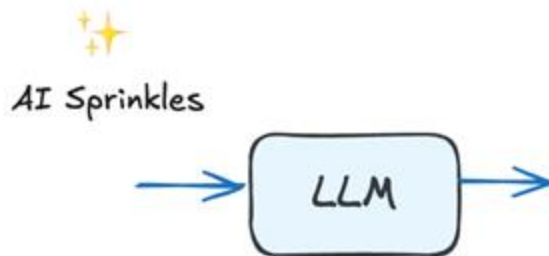
# The timeline to *fully* *automated* incident response.

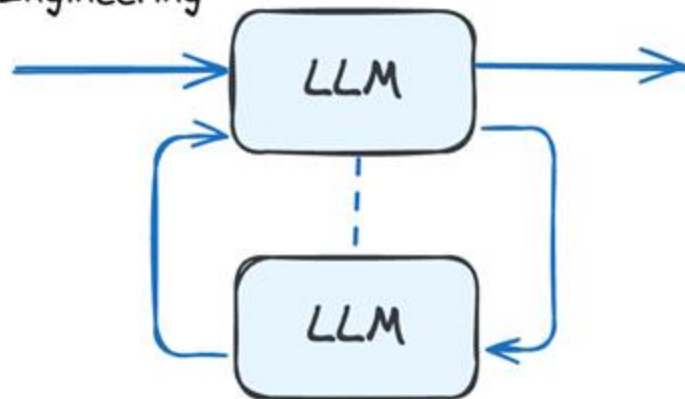# Automated incident investigation and resolution

# Spectrum of AI in product development

➔ No AI

➔ AI Sprinkles

➔ AI Engineering

➔ Foundational model engineering

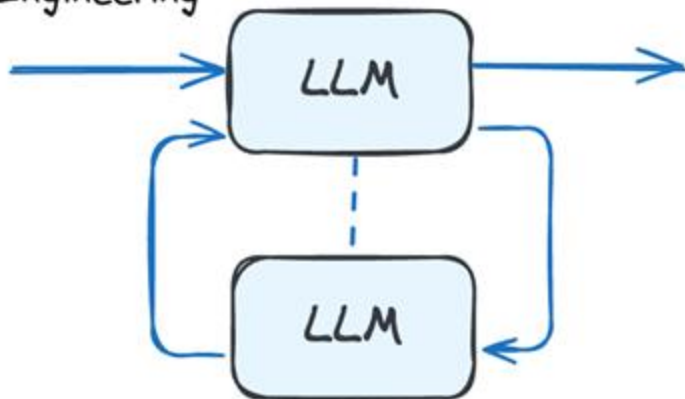AI Sprinkles

LLM

AI Engineering

LLM

LLM

🧑‍🏭 AI Engineering

**Maturity phases**

# Development
# Running in production
# Evolving the system

# Maturity phases

**Maturity phases**

# Development
## Running in production
## Evolving the system

**You need to...**

# ...call LLMs with your data, and receive structured output

# Prompt

➜ Abstraction around LLM calls

➜ Integrate type-system to structured output where possible

➜ Support various LLM providers (OpenAI, Anthropic, Google)

➜ Apply instrumentation and cost tracking

```go
type PromptYarr struct{}

type PromptYarrInput struct {
    Message string `json:"text" description:"Text to translate into pirate"`
}

type PromptYarrResult struct {
    Message string `json:"message" description:"Resulting pirate-speak"`
}

func (p PromptYarr) Model() string {
    return aimodel.GPT41_20250414
}

func (p PromptYarr) Result() PromptYarrResult {
    return PromptYarrResult{}
}

func (p PromptYarr) Prompt(input PromptYarrInput) []goopenai.ChatCompletionMessage {
    return []goopenai.ChatCompletionMessage{
        {
            Role: "system",
            Content: `
# Context

You are a pirate translator who helps translate normal text into pirate speak.

# Task

Translate the given text into pirate speak, making sure to:
* Use common pirate phrases and terminology
* Keep the core meaning of the message intact
* Be consistent with pirate dialect throughout
* Add nautical references where appropriate
`,
        },
        {
            Role: "user",
            Content: input.Message,
        },
    }
}
```

**You need to...**

# ...iterate on a prompt to test it works for some set of inputs

# Evals

➜ These are LLM 'unit tests'

➜ Each test case has prompt input data and expected result

➜ Eval suites define matchers on result and expected

➜ Non-determinism is the challenge

**Matchers** — prompt_yarr.go

```go
func (p PromptYarr) Evals() (
    *EvalSuite[PromptYarrInput, PromptYarrResult], error,
) {
    type (
        Input  = PromptYarrInput
        Result = PromptYarrResult
    )

    return EvalLoad(p,
        EvalCheck("PirateSpeech", EvalCheckPromptV2[Input](
            func(res Result) string {
                return res.Message
            },
Grade how well the translation captures pirate speech patterns on a scale of 1-4:

1: Not pirate-like at all - sounds like normal modern English
2: Uses some pirate words but inconsistently or incorrectly
3: Consistent pirate dialect with common phrases and good nautical flavor
4: Perfect pirate speech with period-appropriate terms and authentic feel

Key aspects to evaluate:
* Use of classic pirate phrases (ahoy, avast, ye, matey, etc.)
* Nautical terminology where appropriate
* Consistent dialect throughout
* Maintains original meaning while sounding authentically pirate

The test should pass if the grade is >= 3.
        `,
        )),
        EvalCheck("MeaningPreservation", EvalCheckPromptV2[Input](
            func(res Result) string {
                return res.Message
            },
Grade how well the pirate translation preserves the original meaning on a scale of 1-4:

1: Completely different meaning or missing key information
2: Some of the original meaning preserved but important details lost
3: Core meaning intact with minor details potentially altered for pirate flavor
4: Perfect preservation of original meaning while adding pirate flair

The test should pass if the grade is >= 3.
        `,
        )),
        EvalCheck("NotEmpty", EvalCheckMatch(func(ctx context.Context, g Gomega, tc
EvalCase[Input], expected, actual EvalCheckOutcome[Result]) {
            g.Expect(actual.Result.Message).NotTo(BeEmpty(), "Pirate translation should not be
empty")
            g.Expect(len(actual.Result.Message)).To(BeNumerically(">", 5), "Pirate translation
should be substantial")
        })),
    )
}
```

**Cases** — prompt_yarr.yaml

```yaml
...
cases:
  - name: Simple greeting
    input:
      text: "Hello, how are you today?"
    expected:
      message: "Ahoy there, matey! How be ye this fine day?"

  - name: Warning message
    input:
      text: "Be careful, there's danger ahead!"
    expected:
      message: "Avast! Thar be danger ahead, ye scurvy dogs!"

  - name: Question about location
    input:
      text: "Where is the treasure located?"
    expected:
      message: "Where be the treasure hidden, ye landlubber?"

...
```

**Results** — bash

```
$ go run cmd/copilot/main.go eval --prompt PromptYarr

# Aggregate results by case

+------------------------+------+------+--------+--------+--------+
| CASE                   | PASS | FAIL | ACTUAL | TARGET | RESULT |
+------------------------+------+------+--------+--------+--------+
| Warning message        | 2    | 0    | 100%   | 100%   | PASS   |
+------------------------+------+------+--------+--------+--------+
| Command to move        | 2    | 0    | 100%   | 100%   | PASS   |
+------------------------+------+------+--------+--------+--------+
| Expression of anger    | 2    | 0    | 100%   | 100%   | PASS   |
+------------------------+------+------+--------+--------+--------+
| Question about location| 2    | 0    | 100%   | 100%   | PASS   |
+------------------------+------+------+--------+--------+--------+
| Expression of excitement | 2  | 0    | 100%   | 100%   | PASS   |
+------------------------+------+------+--------+--------+--------+
| Polite request         | 2    | 0    | 100%   | 100%   | PASS   |
+------------------------+------+------+--------+--------+--------+
| Statement about weather | 2   | 0    | 100%   | 100%   | PASS   |
+------------------------+------+------+--------+--------+--------+
```

**incident-io-bot** commented 2 minutes ago                                ...

# Copilot Evals

*Total cost: $9.8178 |* *Branch output* *|* *Master output*

### 🎯 Regressions

Prompts that are failing more than on master:

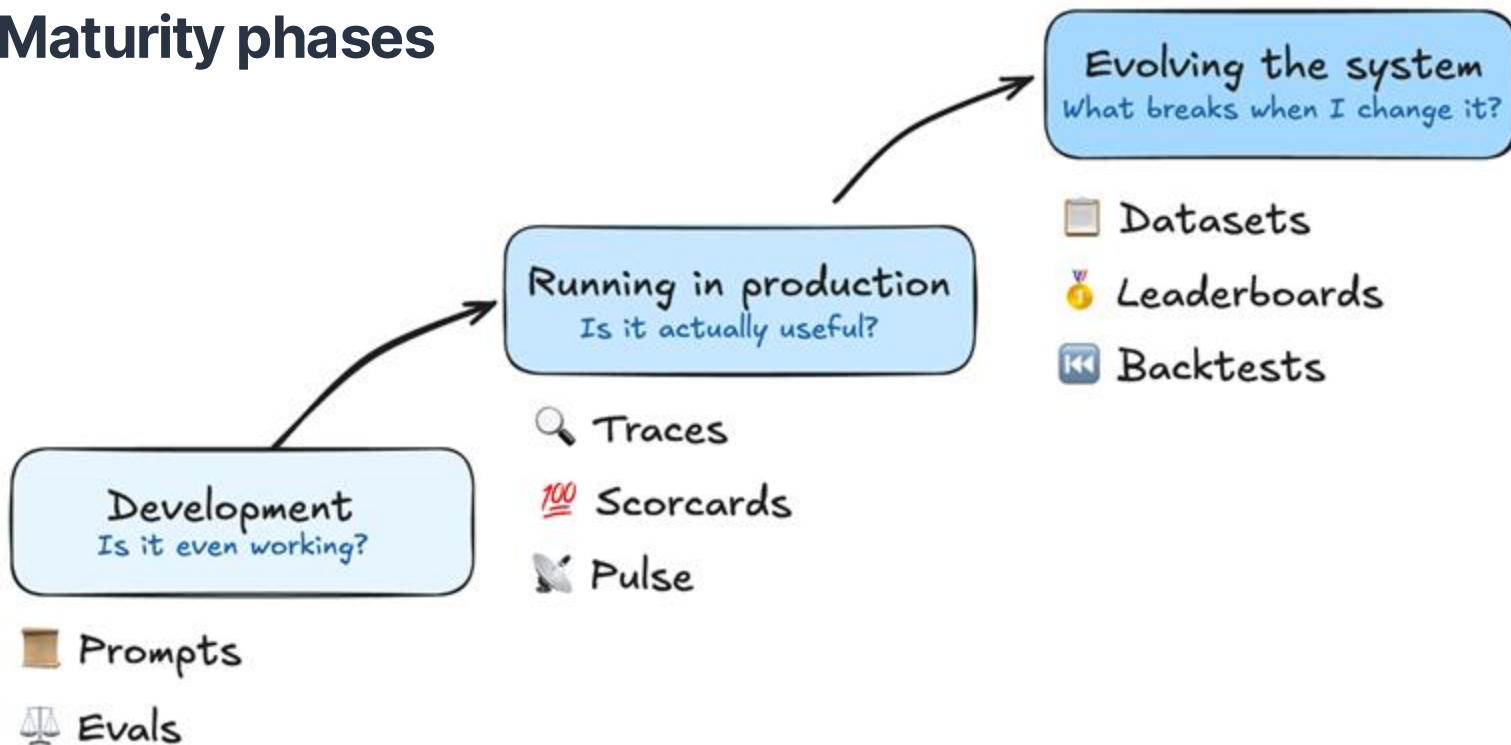| Prompt | Master | Branch | Diff |
|---|---|---|---|
| PromptActionsDraft | 100% (7/7) | 86% (6/7) | ▼ -1 case |
| PromptCopilotAreYouAskingMe | 100% (29/29) | 97% (28/29) | ▼ -1 case |
| PromptCopilotScoreThread | 100% (10/10) | 80% (8/10) | ▼ -2 cases |

### ✅ Broken on master

There are no prompts which are always failing.

▶ All prompts

☺

# Maturity phases



**Evolving the system**
What breaks when I change it?

📋 Datasets
🥇 Leaderboards
⏮️ Backtests

**Running in production**
Is it actually useful?

🔍 Traces
💯 Scorcards
📡 Pulse

**Development**
Is it even working?

📜 Prompts
⚖️ Evals

**Maturity phases**

Development
# Running in production
Evolving the system

# payment-batch-handler-OOM-causing-processing-del... ⌄

**incident** APP 12:28 AM

**What's going on**
The payments service is experiencing serialization errors when handling custom fields, causing failures in processing payments.

**Root cause**
Improper handling of numeric custom fields due to changes introduced in PR #35285 [1]. Specifically, the changes failed to validate or resolve placeholder values in numeric fields during serialization.

| Review Investigation | PR #35285 | Error logs |

**Milly** 12:29 AM
@incident have we made any other changes to the serialization logic in the payment service?

**incident** APP 12:29 AM
Recent changes:

- **PR #34289** Add validation for dynamic custom field types in payment routes
- **PR #31384** Fix: Broken serialization on legacy payment configs with deprecated fields

**Lisa** 12:29 AM
@incident Have we had any similar issues previously?

**incident** APP 12:29 AM

**Lisa** 7:52 PM
@incident create a follow-up to understand why web/web CPU was climbing - looks like a memory leak

✅ 1   😊+

**incident** APP 7:52 PM
Lisa Karlin Curtis (via @incident) added a follow-up:

**Investigate web tier high CPU and memory usage for potential memory leak**

🔗 ONC-6792

Assigned to @Lisa   Priority: **High**   Status: **Open**

Trace • Evals • $0.14 • 10.423s

**You need to...**

**...see exactly what we sent and received, in order to debug**

# Traces

➔ Not standard o11y traces, these contain rich data

➔ See exactly what was sent and received from the model

➔ Allow you to re-run and compare models

➔ 'Steal' a trace for an eval

'Steal' an eval

**You need to...**

**...score AI interactions, marking them as good/bad and why**

**Lisa** 7:52 PM
@incident create a follow-up to understand why web/web CPU was climbing - looks like a memory leak
✅ 1

**incident** APP 7:52 PM
Lisa Karlin Curtis (via @incident) added a follow-up:

**Investigate web tier high CPU and memory usage for potential memory leak**    🔗 ONC-6792

Assigned to @Lisa    Priority: **High**    Status: **Open**

Trace • Evals • $0.14 • 10.423s

Scored 2 days ago    📊 Show trace    🔄 Refresh    ▼ Hide events

**Copilot thread**    thread

✓ **Overall score** 100.0%  2 / 2    ✓ **Efficiency** 100.0%  5 / 5
The overall score of the thread    Direct path to user's goal without unnecessary steps

✓ **Alignment** 100.0%  5 / 5
Appropriate tool selection and operation execution

Intents ⓘ
create-followup

Error Type ⓘ
none

Supported Action ⓘ
true

Source ⓘ
user

User Frustration ⓘ  Colour by dimension

Specificity ⓘ  Colour by dimension

Actions Taken ⓘ
follow_ups_create

# This is our first objective definition of 'good' for AI interactions

# #copilot-interactions-pulse

💚 Copilot thread | incident-io | ✍️

Classification: action • Type: create_followup • Score: high

Message • Evals • $0.1397 • 11.934s

> 👤 Lisa
> @incident create a follow-up to understand why web/web CPU was climbing - looks like a memory leak
> Posted in # inc-13251-intermittent-google-cloud-failures | Jun 12th | View message

💩 Copilot thread | incident-io | ✍️

Classification: action • Type: update_summary • Score: low

Message • Evals • $0.0585 • 6.723s

> 👤 Liz
> @incident update summary to say the impact was that one user would have seen an error when trying to filter the page with the outdated ID
> Posted in 🗂 inc-13076-scim-group-not-found-causing-incident-loading-failure-for-hiya (archived) | May 29th
> | View message

🔥 incident APP 12:47 PM
🤭 Copilot thread | incident-io

Classification: action • Type: draft_name, accept_incident • Score: medium • Tags: incident_management, incident_decline, action_performed, manual_solution, expectation_mismatch, efficiency_issues

Message • Evals • $0.1750 • 10.458s

> 👤 Isaac
> @incident better name than that, come on. also accept it please
> Posted in 🗂 inc-12880-master-ci-builds-failing-due-to-buildkite-cache-and-node_modules (archived) | May
> 14th | View message

# Maturity phases

# Maturity phases

**Development**
**Running in production**
**Evolving the system**

**You need to...**

**...make a change and know how it's impacted the system for real customers**

# Backtests

➔ Datasets + Scorecards = Backtests

➔ Re-run interactions, score them, aggregate scores

➔ Make a change and see in ~15m how it changes dataset scores

➔ This is how you avoid terrible regressions

# Filter unrelated code changes before topk #36559

💬 Conversation 7 | ◦ Commits 1 | ☑ Checks 5 | ⊡ Files changed 1

**millyleadley** commented on May 12                                    ···

Previously we were passing all high confidence items into the topk analysis, even those which had a causality of `unrelated` .

Also, if there weren't any high confidence items but there were some medium ones, we were just dropping them on the floor.

Instead, let's:

1. Drop findings that had a causality of unrelated or uncertain
2. Still consider potentialyl causing code changes with a `medium` confidence if there were no high confidence ones

☺

◦— 🐥 Filter unrelated code changes before topk ···                  ✓ af9d03c

# Filter unrelated code changes before topk #36559

**Merged** millyleadley merged 1 commit into `master` from `milly/filter_unrelated_code_changes_before_topk` on May 12

| Conversation 7 | Commits 1 | Checks 5 | Files changed 1 |

**millyleadley** commented on May 12

Previously we were passing all high confidence items into the topk analysis, even those which had a causality of `unrelated`.

Also, if there weren't any high confidence items but there were some medium ones, we were just dropping them on the floor.

Instead, let's:

1. Drop findings that had a causality of unrelated or uncertain
2. Still consider potentially causing code changes with a `medium` confidence if there were no high confidence ones

🙂

Filter unrelated code changes before topk ✓ af9d83c



**milhouse** May 12th at 4:41 PM

Nice win from that change I made today which changes the filtering before topk in code changes:

- 25% faster, because we'll be skipping topk when its not needed
- 10% cheaper, because less calls
- 15% better recall, because we'll be passing some `medium` confidence code changes into topk when there aren't any `high`'s

Screenshot 2025-05-12 at 16.38.02.png ▼



🎉 6

**You need to…**

**…track existing baselines and motivate people to hit 'high scores'**

# Leaderboards

➔ Evaluation criteria, baseline, hill-climb

➔ Pin the best scorecard to a leaderboard

➔ This is how you reestablish the dopamine hit

➔ Make your team feel good when they succeed!

# Workbench

## Workbench

**Features**

- Copilot threads
- Copilot messages
- Investigations
- Suggestions
- Searches

**Evaluation**

- Leaderboards
- Backtests
- Datasets
- Search ground truths

**Processed resources**

- Onboarding
- Code changes

**Platform**

- Backfills
- Health reports
- Processors
- Requests

## 🏆 Leaderboards

[Create leaderboard]

### 💡 Investigation (code-changes)

**Causing code changes**

**Cost**
- 🔵 Total cost `$0.39`

**Code changes**
- ✅ Primary precision `85.0%` `17 / 20`
- 🟠 Precision `66.7%` `20 / 30`
- 🟠 Recall `64.5%` `20 / 31`
- 🟠 F1 `60.3%`

**Causing code changes (including negatives)**
- 🟠 Precision `57.7%` `15 / 26`
- 🟠 Recall `48.4%` `15 / 31`
- ✅ F1 `76.8%`

### 🟠 Incident facts

**Incident learnings**
- ✅ Recall `77.6%` `45 / 58`
- ✅ Precision `75.0%` `45 / 60`
- ✅ Evidence Quality `96.7%` `58 / 60`

### 💡 Investigation (messages)

**Investigation (messages)**
- ✅ Precision `80.6%` `25 / 31`

# That's it!

# Maturity phases



**Evolving the system**
What breaks when I change it?

📋 Datasets
🥇 Leaderboards
⏮️ Backtests

**Running in production**
Is it actually useful?

🔍 Traces
💯 Scorcards
📡 Pulse

**Development**
Is it even working?

📜 Prompts
⚖️ Evals
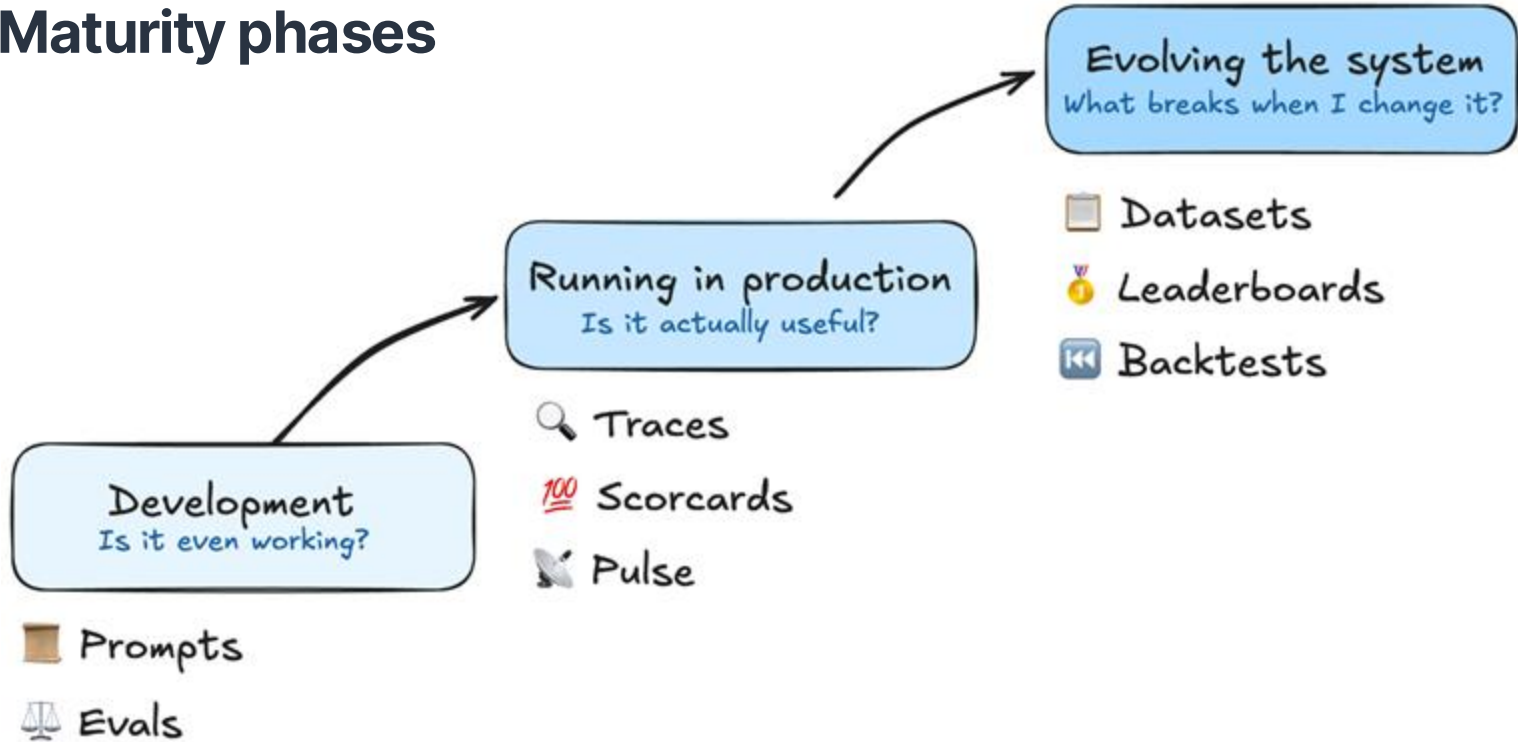
# We've discussed tools, but they solve for very real human problems.

➔ Non-deterministic systems hurt

➔ Working with AI is a rollercoaster, it can be exhausting

➔ Objective measures are the answer

➔ Help celebrate the small wins

# A note from the team

We're building an agent that will **investigate** the incident **with** you, show you what's wrong and why it's wrong. Eventually, we'll show you how to fix it or offer to fix it on your behalf. It will feel like we're right there with you, helping you resolve the incident just like your best & most experienced colleague would.

To get there, we're on the sharp edge of what's possible with AI. It's pushed us to become 'AI Engineers', build internal tools to tame the chaos of non-deterministic systems, and rethink our product process for AI systems which are much harder to evaluate. We've built resilience, learned to navigate the uncertainty and stay motivated through tough, ambitious R&D.

These are the people and stories behind AI at incident.io.

This could be you!

incident.io raises *$62M*
To build AI agents that resolve incidents with you