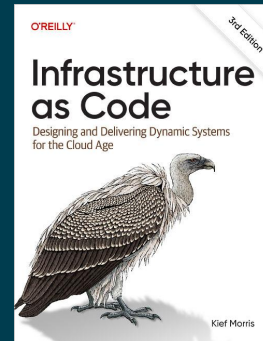


Pulling Continuous Delivery inside the agentic loop

How do we use GenAI to build software to last?



Kief Morris

Distinguished Engineer
Technology Advisory Services



<https://bit.ly/4dS28pT>



I see value in
hands-off, "all-in"
agentic workflows.

A blue industrial robotic arm is the central focus, positioned in a factory environment. The arm is mounted on a white base and is reaching towards the right. In the background, there are various pieces of industrial machinery, including a white machine with a black sign that reads 'T 07'. The scene is lit with bright, cool-toned lights, creating a high-tech atmosphere. Large, white, sans-serif text is overlaid on the image, reading 'But it means getting better at ensuring "good".'

But it means
getting better at
ensuring "good".



Which means
being clear about
how we define
"good".

/thoughtworks

Photo by [Homa Appliances](#) on [Unsplash](#)

I believe using CD pipelines
within agentic workflows can
help **build better software.**



All-in, "hands-off"
agentic workflows.



You don't touch code.

Even when you don't like the code that that the agents extrude.



You probably shouldn't look at
the code either.

You can look at the code if you
really want. (Although it feels
like a losing battle.)



So how do we avoid cognitive debt?

How do we avoid cognitive surrender?

A large industrial factory interior, possibly a steel mill, with a crane hook hanging from the ceiling. The scene is filled with complex machinery, metal structures, and a high ceiling with visible beams and lighting. The lighting is dramatic, with strong highlights and deep shadows, creating a sense of scale and industrial complexity.

Do we really need our code
to be good?

A large industrial factory interior, likely a steel mill, with a crane hook hanging from the ceiling. The scene is filled with complex machinery, metal structures, and a high ceiling with visible beams and lighting. The overall atmosphere is industrial and somewhat dimly lit, with some light streaming in from the background.

It's the LLM's problem now!

A scenic landscape featuring a vibrant rainbow arching over a calm lake. In the background, there are snow-capped mountains under a clear blue sky. The foreground is filled with a field of bright yellow wildflowers. The overall atmosphere is bright and cheerful.

Vibe-oriented approaches build software that is good in the ways that users know they care about.

"It does what I want it to!"

A scenic landscape featuring a vibrant rainbow arching over a calm lake. In the background, there are snow-capped mountains under a clear blue sky. The foreground is filled with a field of bright yellow wildflowers. The overall scene is bright and colorful, suggesting a beautiful day in nature.

But does GenAI build software
that is good in the ways that
users will care about later?

Not unless you make it do it.

A scenic landscape featuring a vibrant rainbow arching over a calm lake. In the background, snow-capped mountains rise against a clear blue sky. The foreground is filled with a field of bright yellow wildflowers. The overall scene is bright and cheerful, suggesting a positive outcome or a 'good' result.

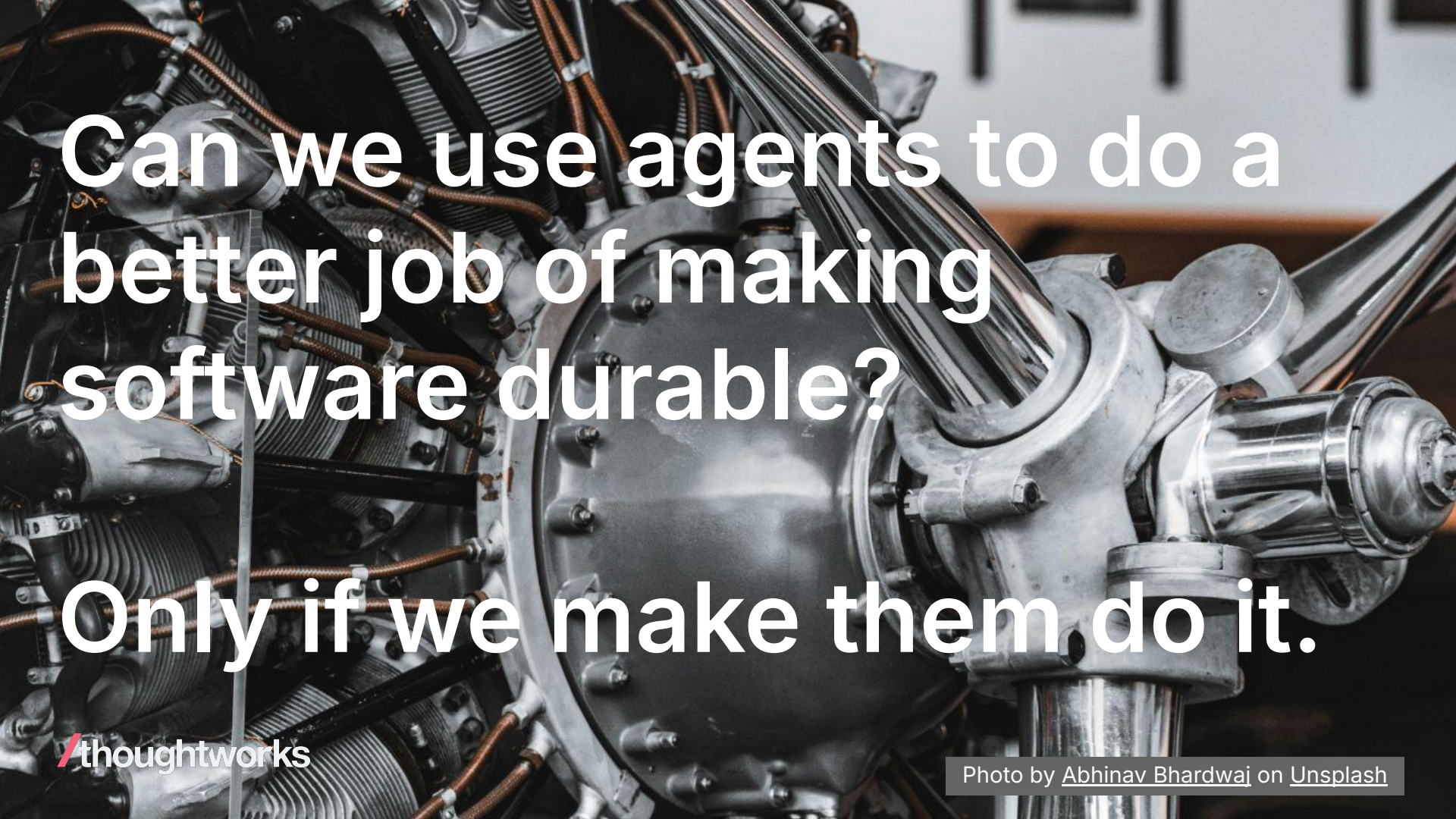
We've always had a hard
time getting users to care
about making software
"good" in ways that show
up later.

And maybe that's fine for some software.

Maybe for a lot of software made with AI.




But what if we're building
software that **needs to last?**



Can we use agents to do a better job of making software durable?

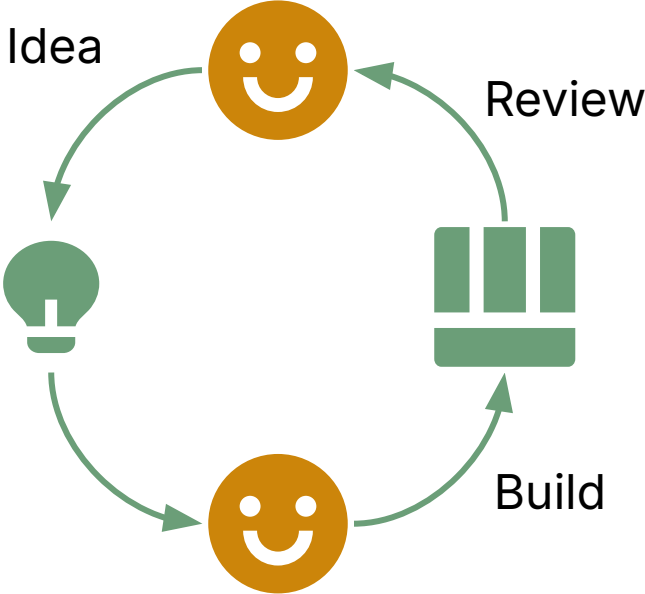
Only if we make them do it.



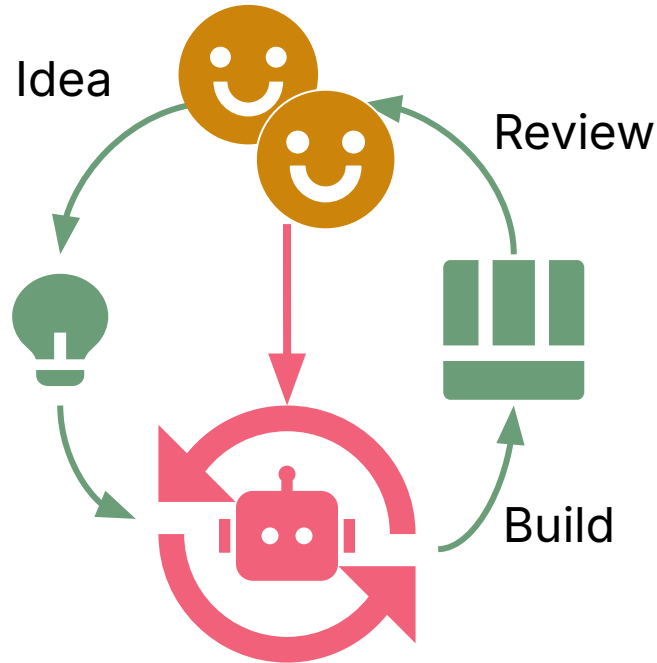
If do we get our agents to build software to last, that can have benefits even for software that we might not think needs it.

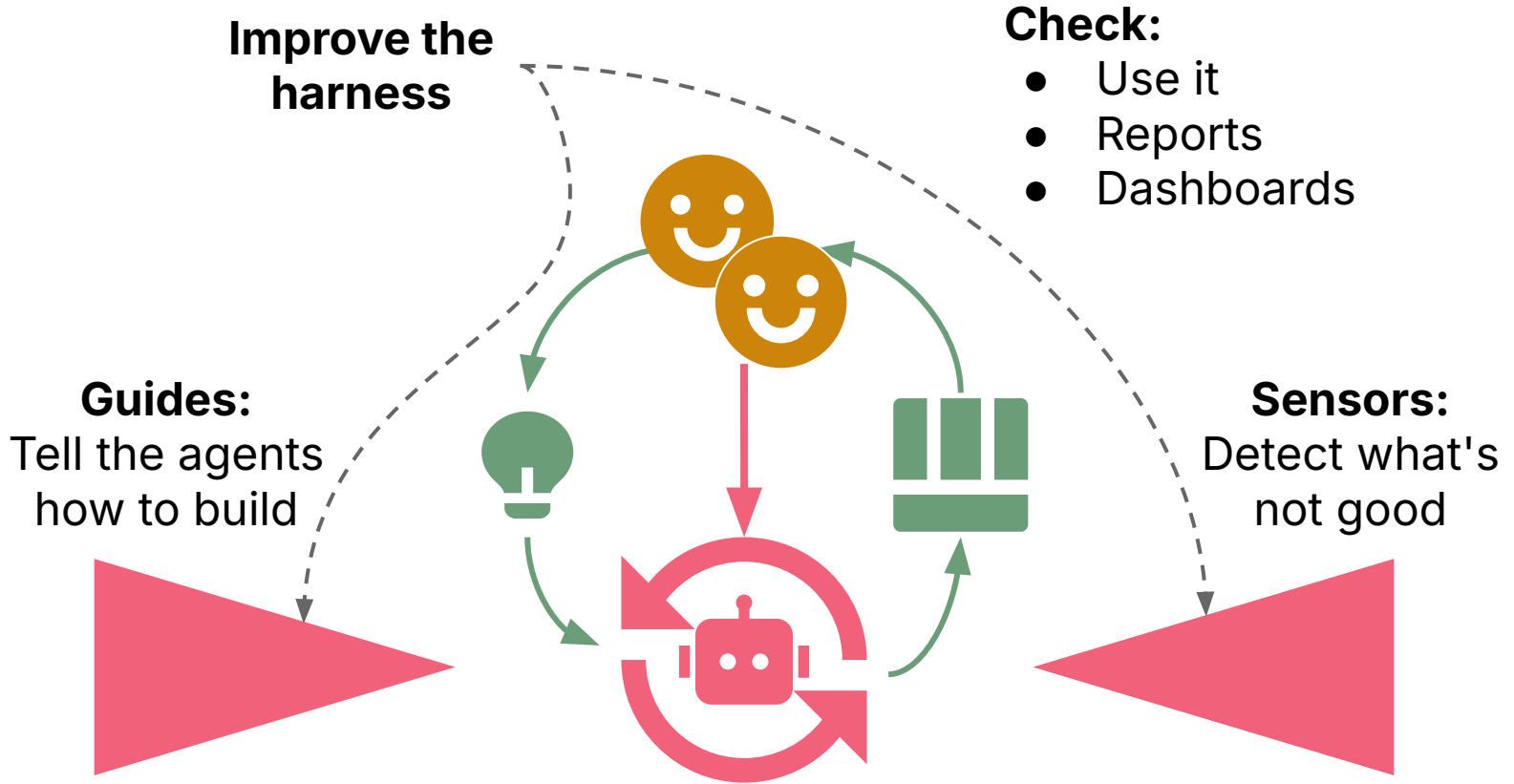
Agentic coding harness:
A system that drives the LLM
to make good code.

Legacy approach: Humans are the loop



Agentic loop. Work on the harness, not the code





Prerequisite:

We need to clearly and effectively articulate what "good" looks like.

Code quality?

Who cares!

It's the LLM's problem now!

Code quality impacts **DELIVERABILITY**

We need to:

Change it. Configure it. Improve it. Add to it.
Troubleshoot it. Fix it.

How do we measure deliverability?

Lagging measures

DORA and things like that.

Leading measures

Linting

Test coverage

Mutation testing, fuzz testing

Code complexity

Vulnerability and supply chain scanning

LLM code evaluations

Implementation affects **OPERABILITY**

We need to:

Know when it breaks.

Restore it quickly when it does.

We also need it to:

Be fast enough.

Have enough capacity.

Be reliable.

Not cost too much to run.

Handle data properly.

Be secure.

Comply with laws and standards.

How do we measure operability?

Lagging measures

Deliverability metrics (DORA)

User satisfaction

Business operations metrics
(transactions, cost)

Commercial metrics

Leading measures

Deployability

Performance and scalability testing

Failure scenario testing

Penetration testing

Compliance auditing

Continuous Delivery

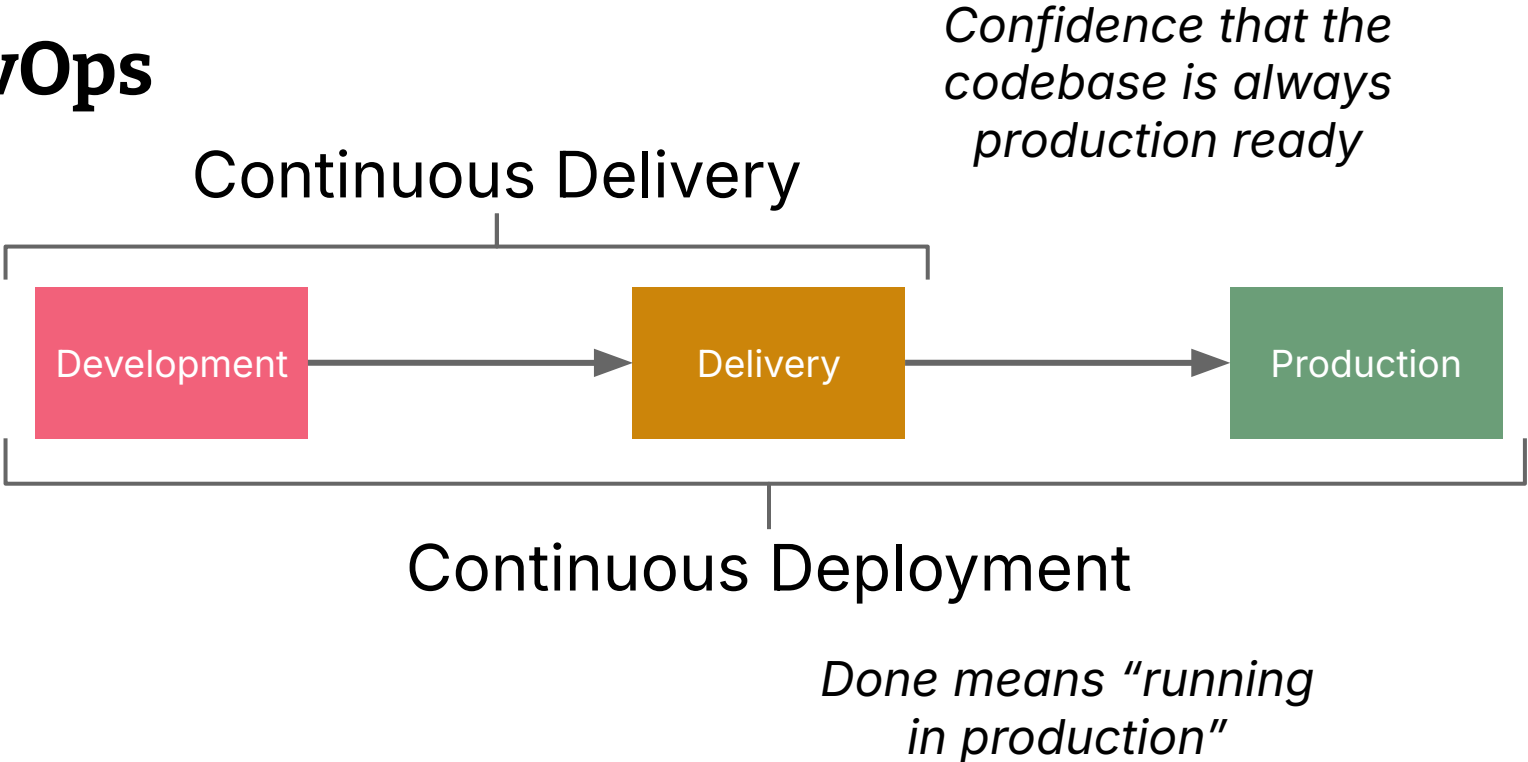
Legacy software delivery flow



*"Done" means
"code complete"*

*Make it "production
ready" after it's done.*

DevOps

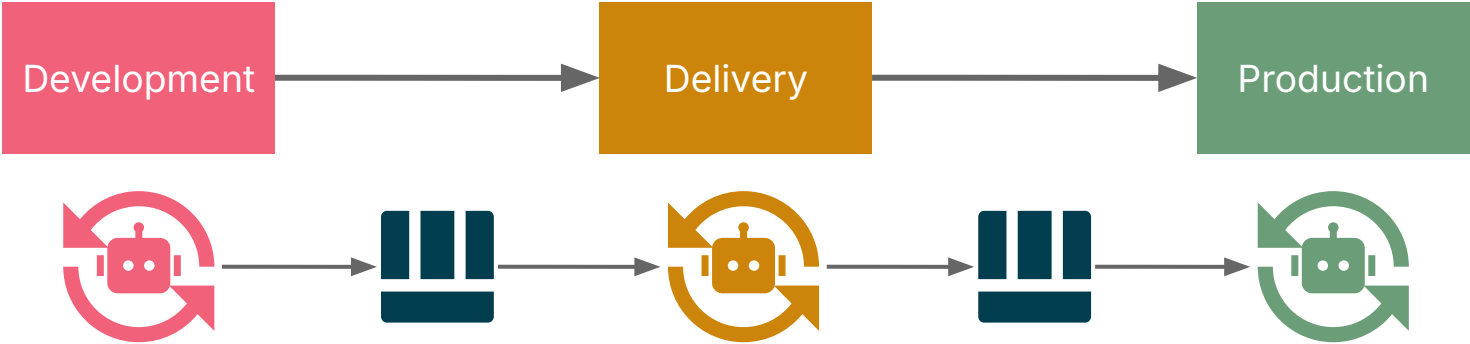


Maybe we don't need CD with agentic engineering?

Agents build the complete solution

Agents make the build production-ready

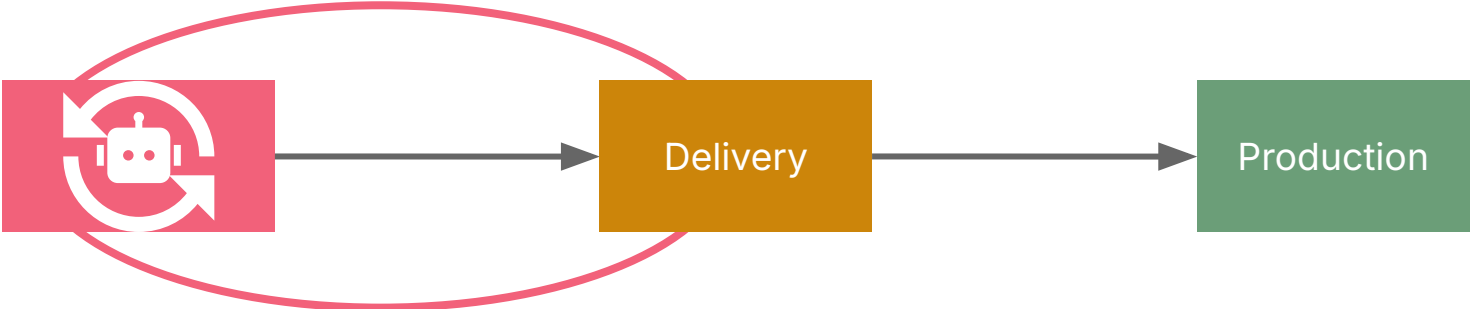
Agents fix issues in production



Operability: software + infrastructure + processes

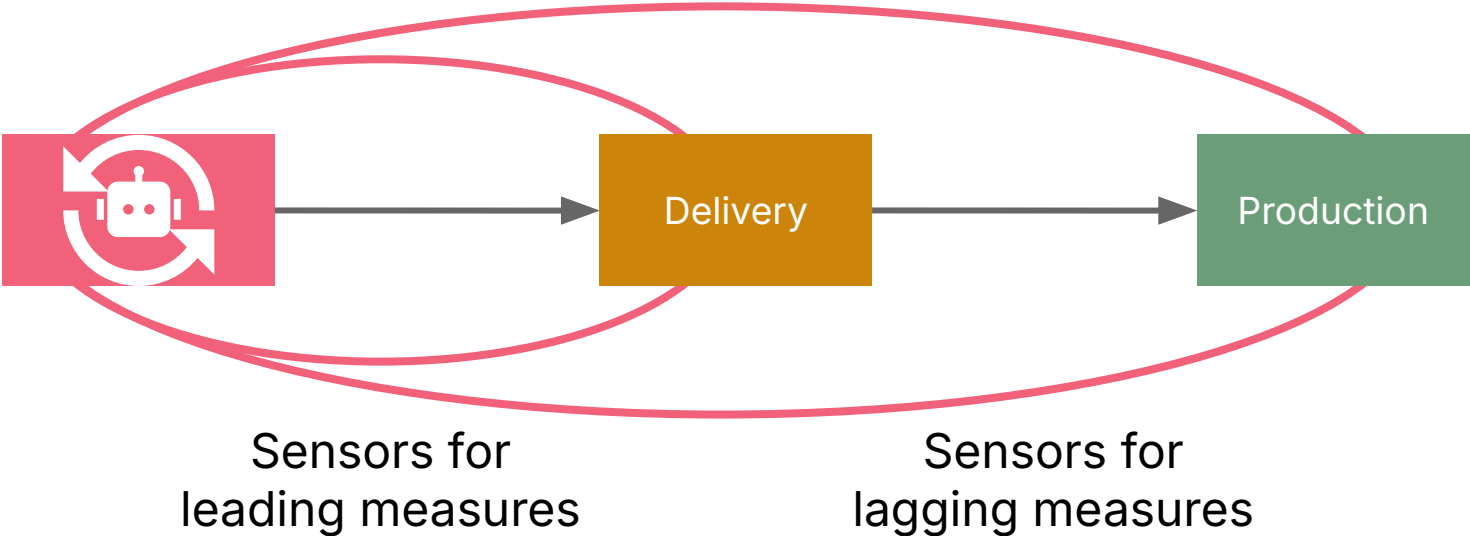
**We can use a Continuous
Delivery pipeline as an
operability harness**

Expanding the agentic loop to include deployment



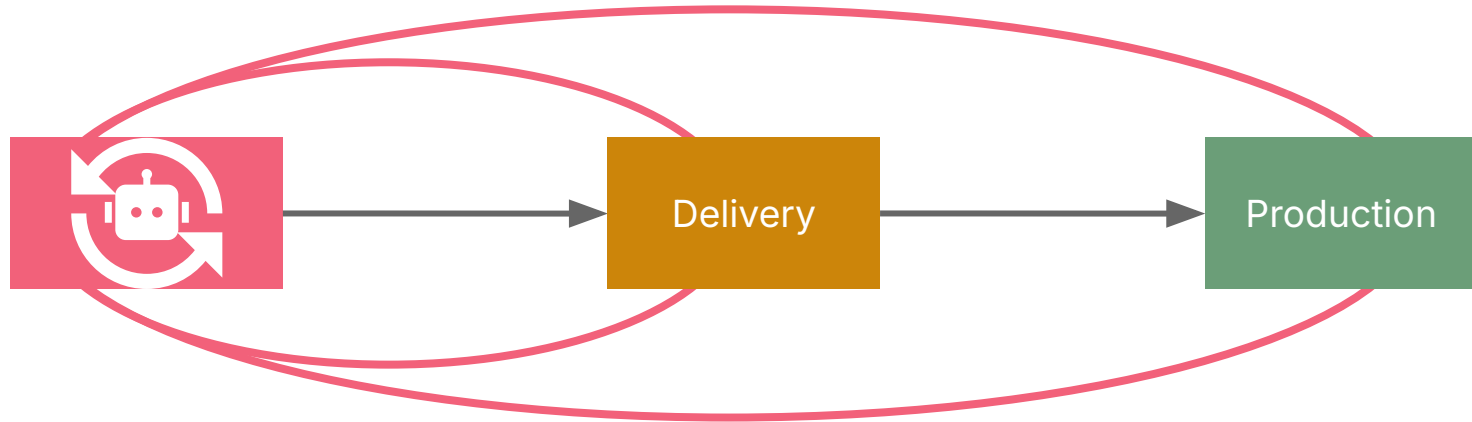
Sensors for
leading measures

Expanding the agentic loop to include production



Challenge: Agent makes fixes downstream

Solution: Prevent access to the deployed infrastructure



"Watch the job. If it fails,
replicate locally, fix the
code, push again

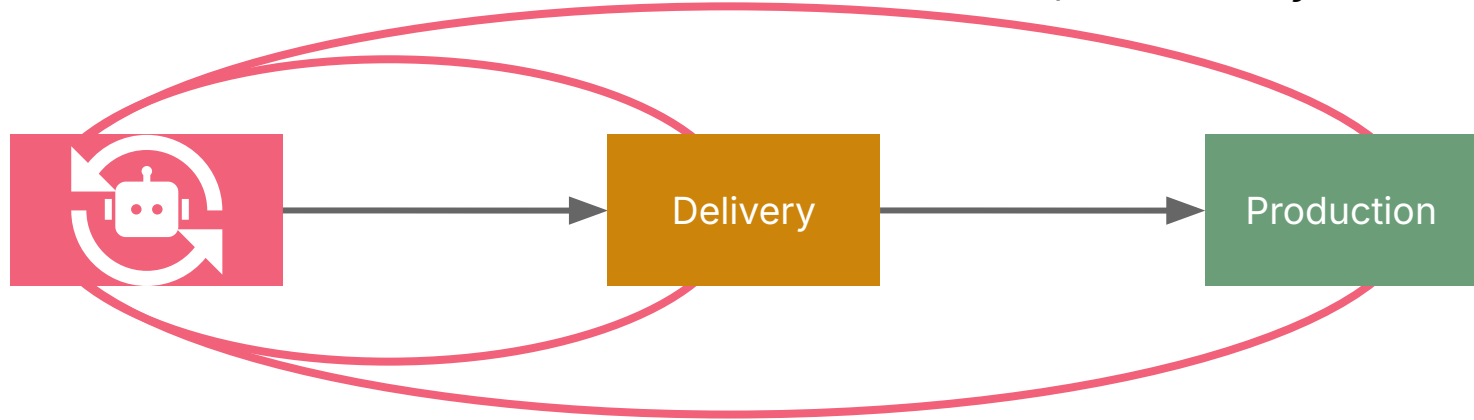
Challenge: Doesn't take responsibility

Agent: *"I fixed it!"*

Agent: *"It was red before my change."*

Me: *"The build is still red?"*

Me: *"Yeah, because you broke it!"*

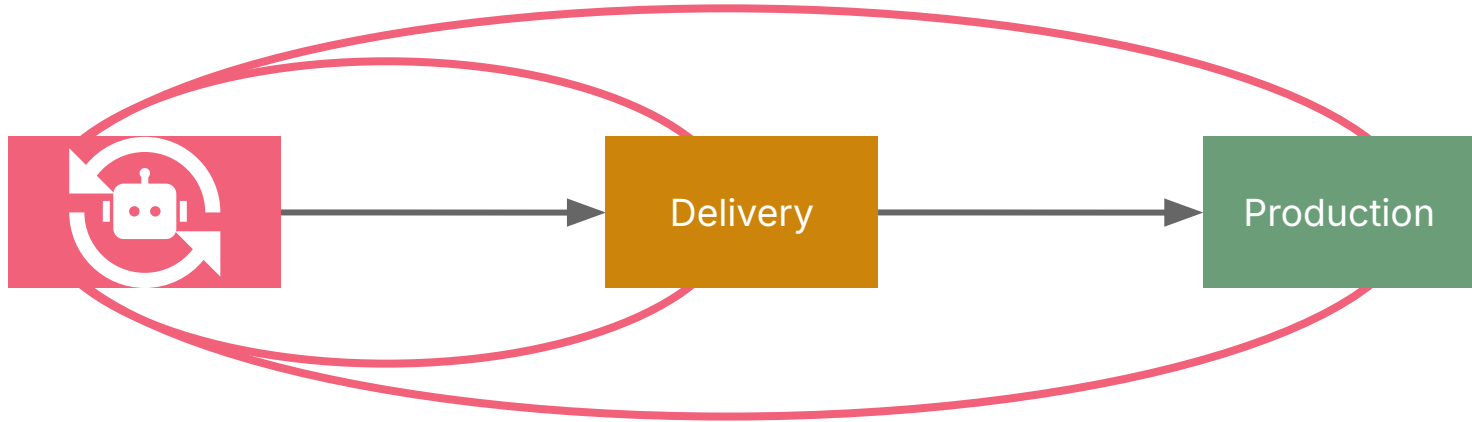


Solution:

Skills for build discipline (?)

Challenge: Worked on my machine!

Cause: Over-fitting

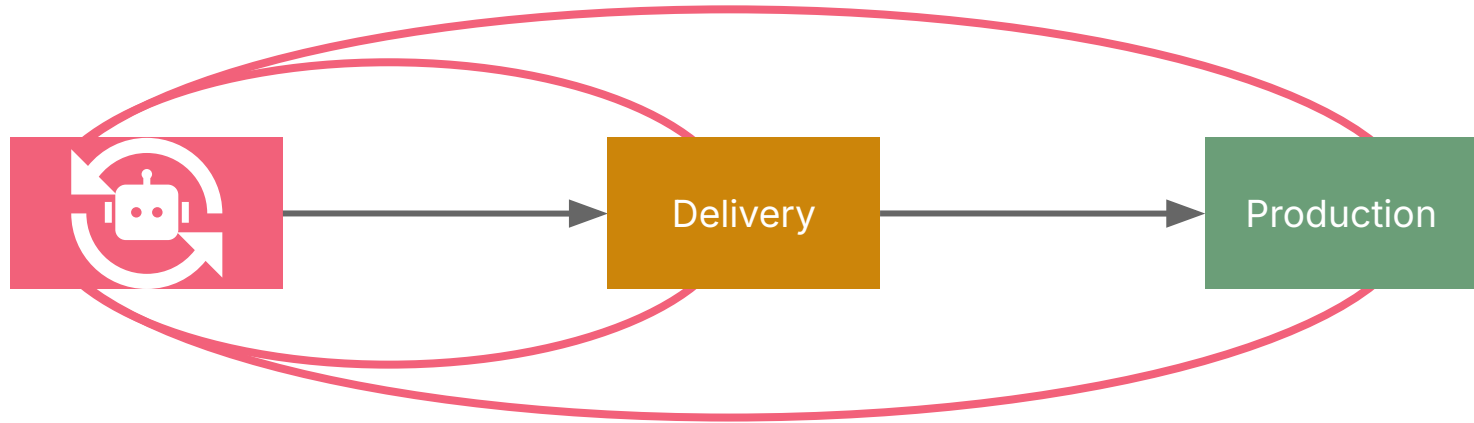


Solution:

Test environment
parameterization earlier

Challenge: Deploying is expensive

Plus: Agents get it wrong. A lot.



Solution:

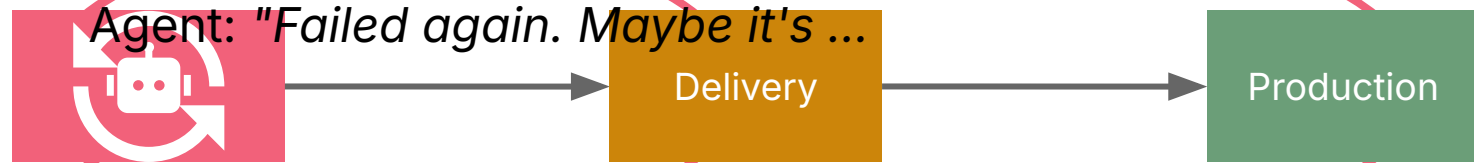
Better local testing

Challenge: Agents fail a lot

Agent: "Hmm. The deployment failed. Let me try something else."

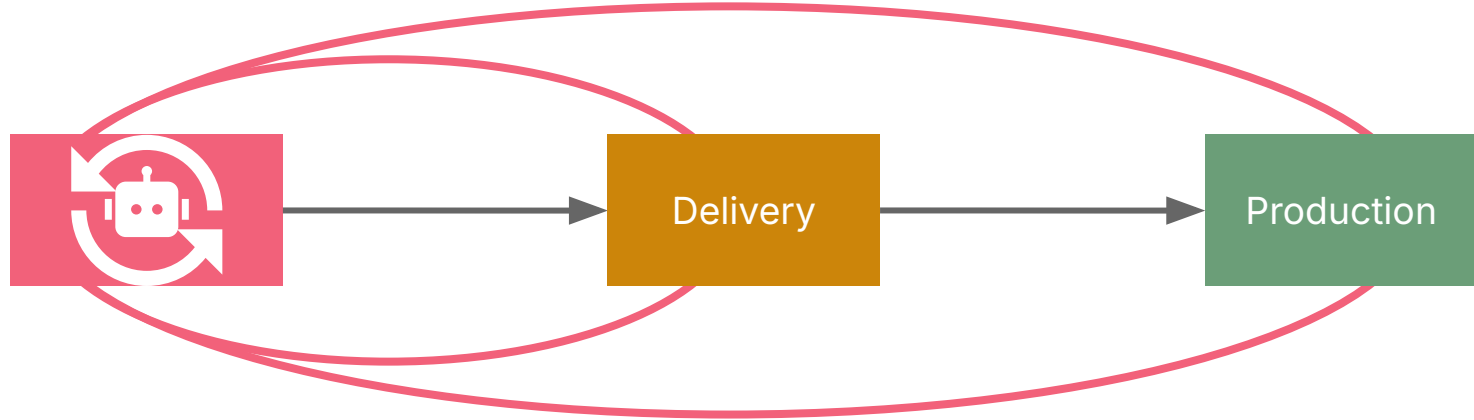
Agent: "That didn't work. Let me try something else."

Agent: "Still fails. I need to think about this. Let me try something else."



Challenge: Agents suck at deploying

Cause: LLMs don't actually *know* how to deploy.



Solution: Detailed skills. "Here's how to deploy. Step 1 ..."

Solution: Prefer scripts and tools over markdown.

Solution: Composability.

Solution: Better models. (Or something else?)

What I think

We can learn a lot with "all-in", hands-off agentic loops

LLMs don't care about building software to last

We need to figure out how to make them do it anyway

CD pipelines work pretty well for humans

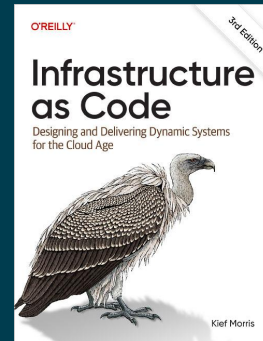
Put CD pipelines inside the agentic loop to learn

We will learn better ways

Pulling Continuous Delivery inside the agentic loop

How do we use GenAI to build software to last?

 /thoughtworks



Kief Morris

Distinguished Engineer
Technology Advisory Services



<https://bit.ly/4dS28pT>