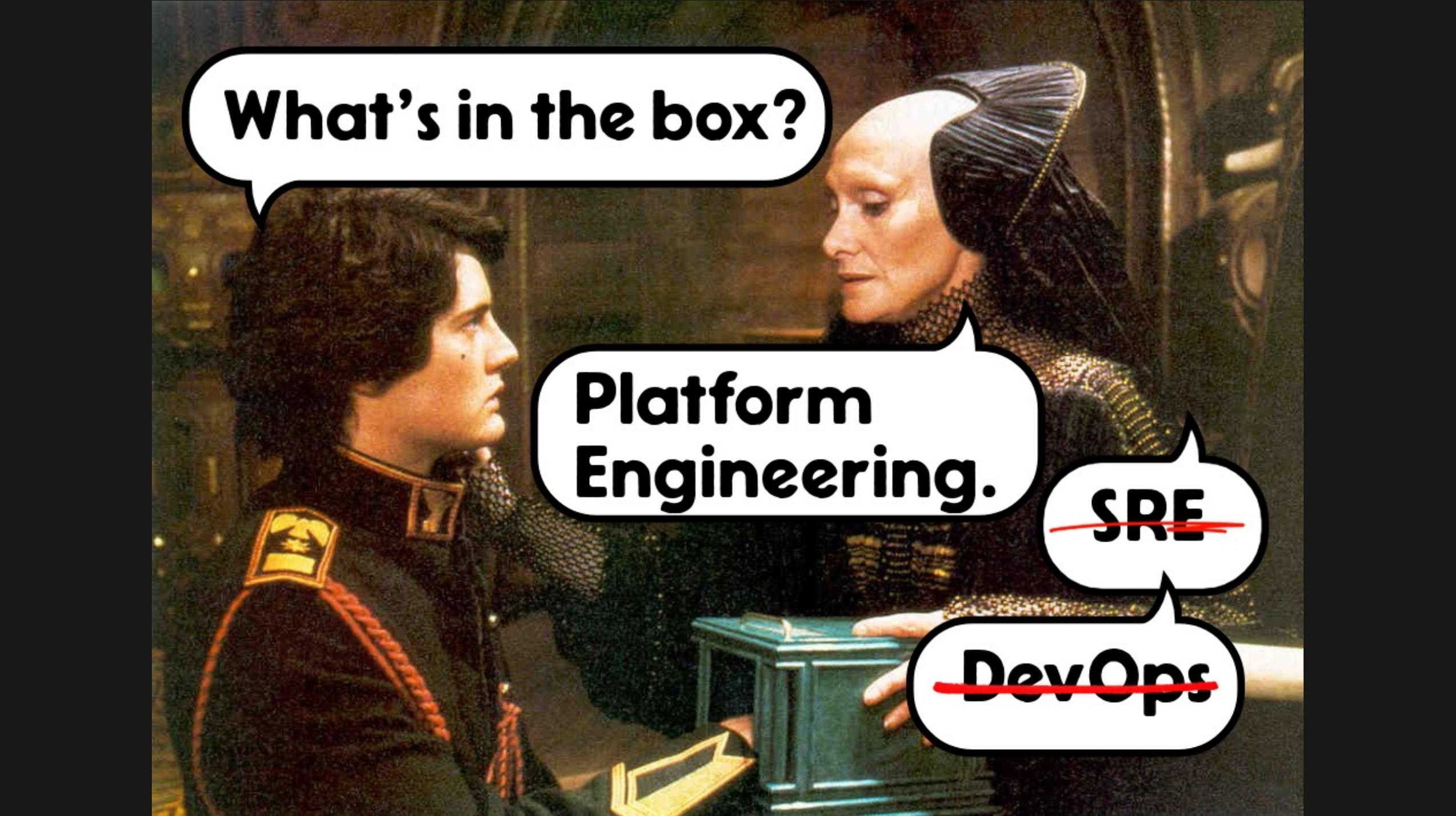


Platform Engineering for Private Cloud

Coté – Civo Navigate London, September 30th, 2025



What's in the box?

**Platform
Engineering.**

~~**SRE**~~

~~**DevOps**~~

50%

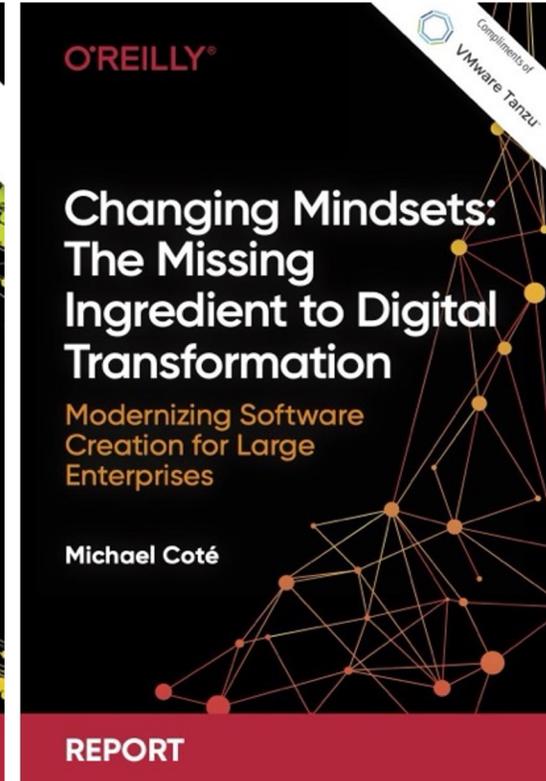
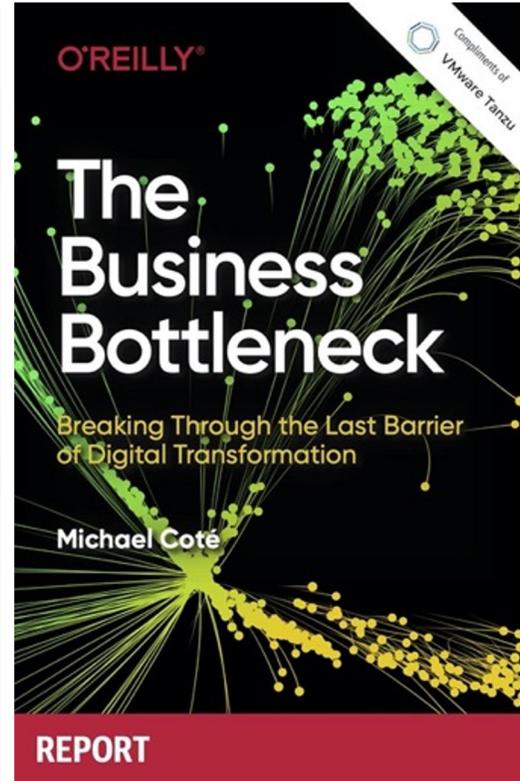
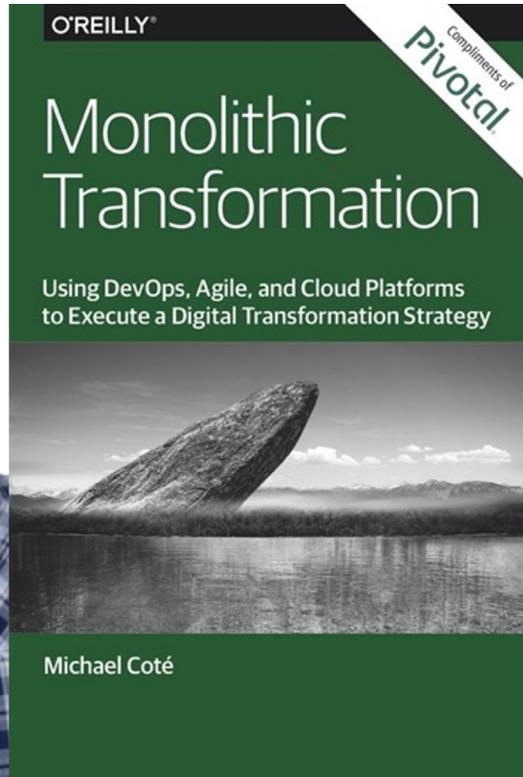
of enterprise apps
run on private cloud

AI:

- Needs apps (maybe)
- Brings new platform customers

Coté

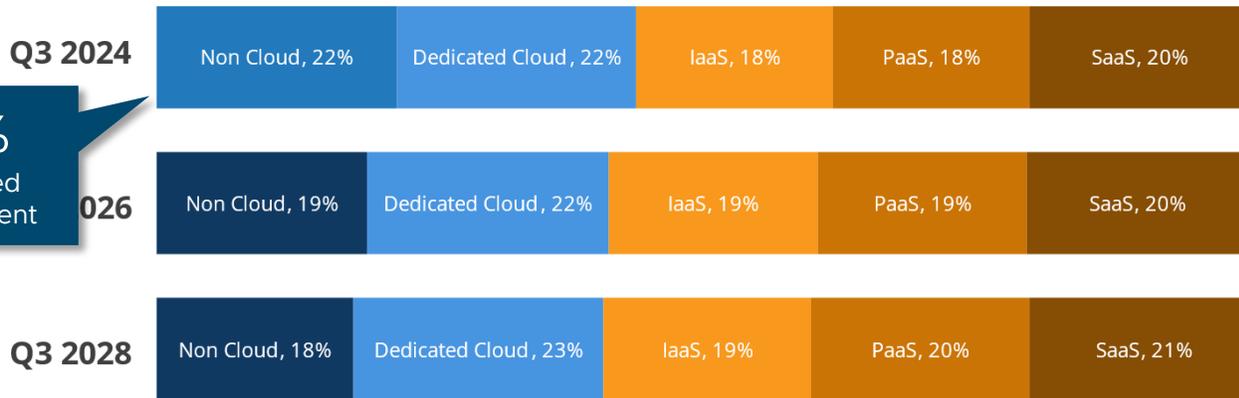
<https://newsletter.cote.io/> | cote@broadcom.com



Where are the apps?

Cloud buyers increased their level of cross-platform interoperability in 2024. They also relied more upon providers for the management of their IT environments

Where are all your organizations applications deployed?



44%
dedicated environment

Level of Interoperability Between Clouds

Public Cloud	41%	Hosted Dedicated Cloud
On Prem Dedicated Cloud	39%	Public Cloud
Hosted Dedicated Cloud	38%	On Prem Dedicated Cloud
Public Cloud	38%	Public Cloud

HOW THE EDGE IS MANAGED

IAAS EDGE

- 49% of IaaS deployments are considered as remote of edge deployments

DEDICATED EDGE

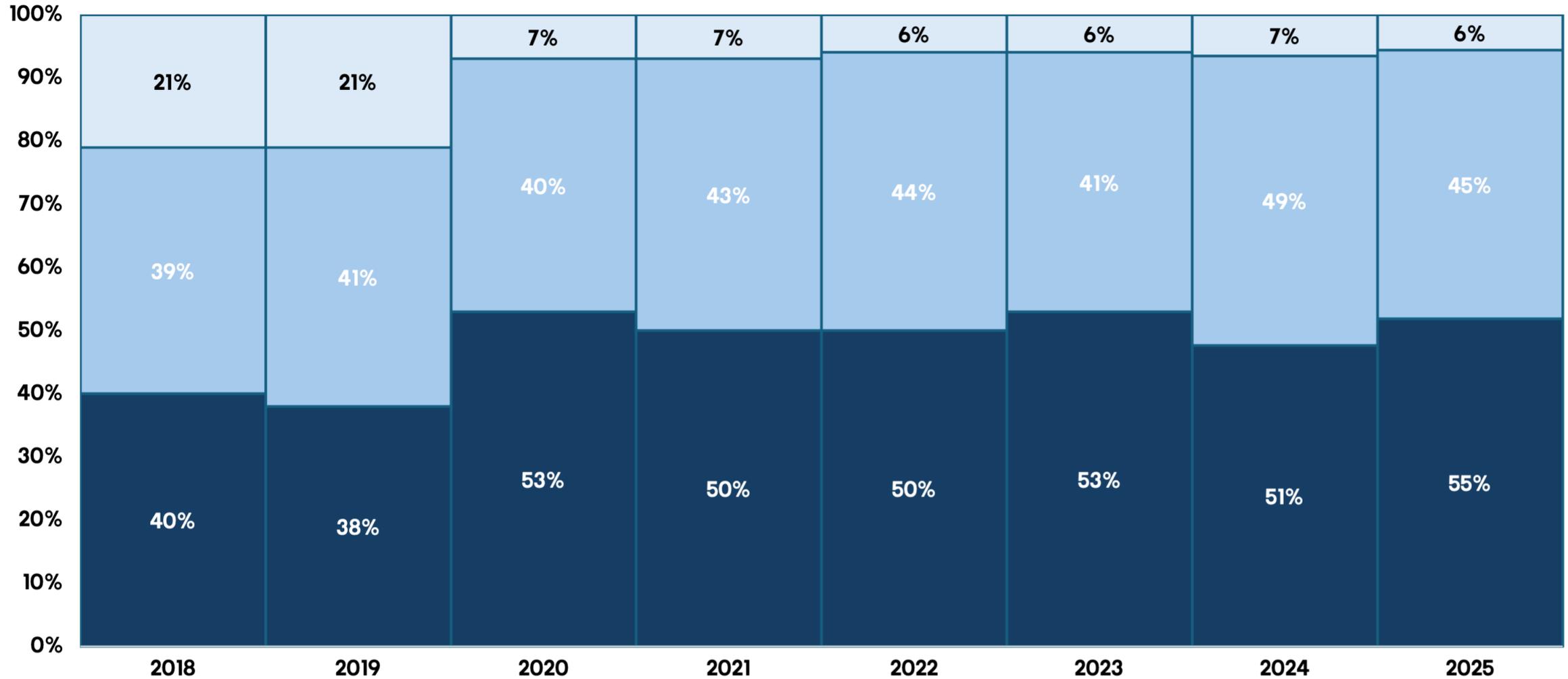
- 48% of dedicated cloud environments are classified as remote or edge deployments.



N = 1,724, QC6 Source: QDV1 3Q24 Cloud Pulse Survey, October 2024, IDC

Workload placement over the years (Flexera)

■ Public Cloud ■ Private Cloud ■ Other



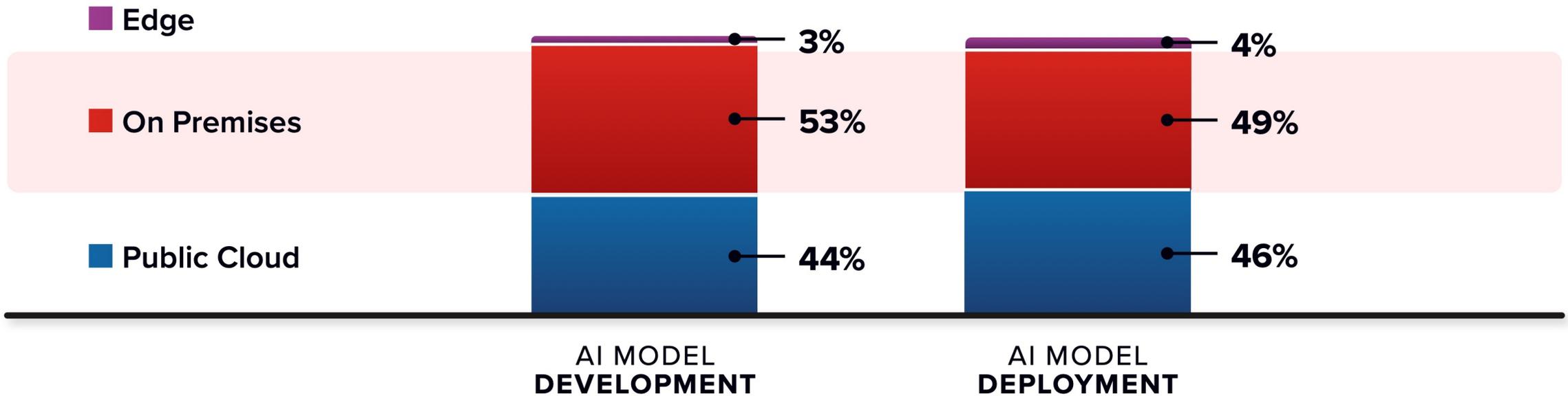
"Other" is "Non-cloud" in 2018 and 2019, "Additional workloads in public cloud in 12 months" in remaining years.

FIGURE 5

Deployment Location for the Development and Deployment of AI Models

Where does your organization primarily develop and deploy AI models?

(Percentage of respondents)

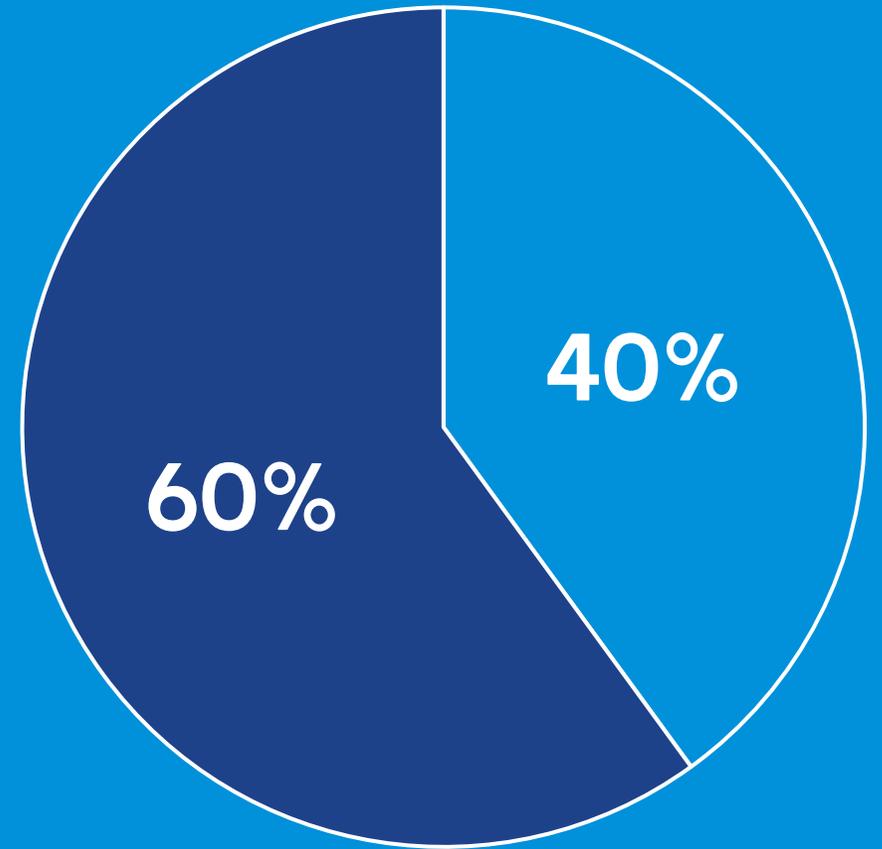
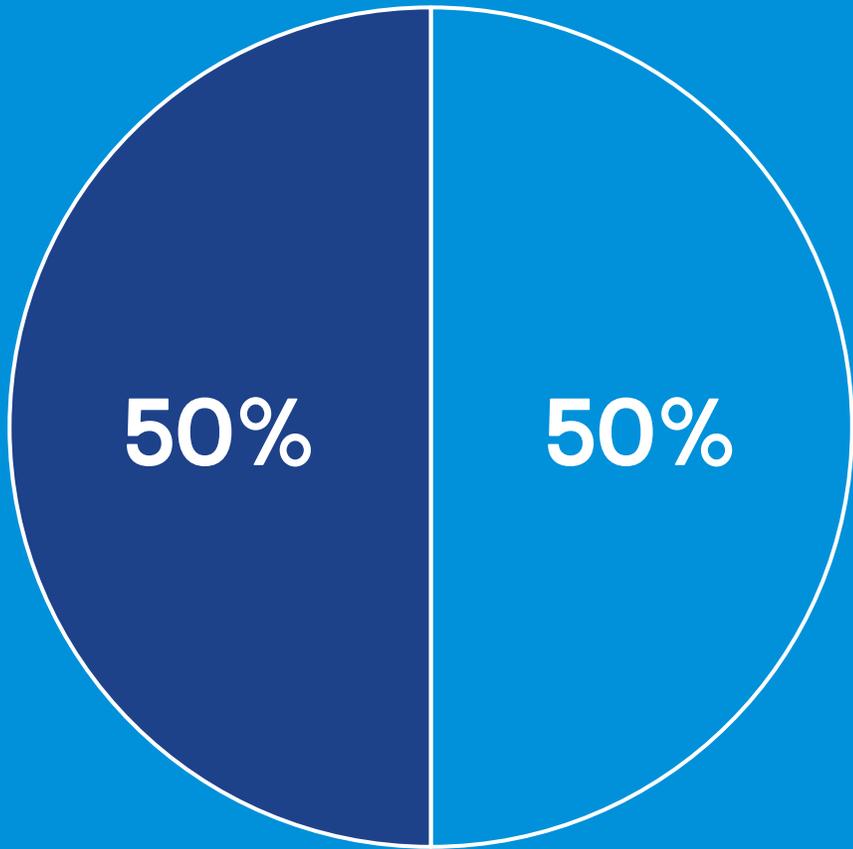


Note: Public cloud includes public cloud provider, managed services provider, and special-purpose cloud. On premises includes private cloud, traditional infrastructure, and colocation provider. n = 411; Source: IDC's *AI Infrastructure Survey*, July 2024.

For an accessible version of the data in this figure, see [Figure 5 Supplemental Data](#) in Appendix 1.



Where the workloads live, rough estimates



What is a platform?

2007?

2019?

Today?



~~PaaS~~
Platform



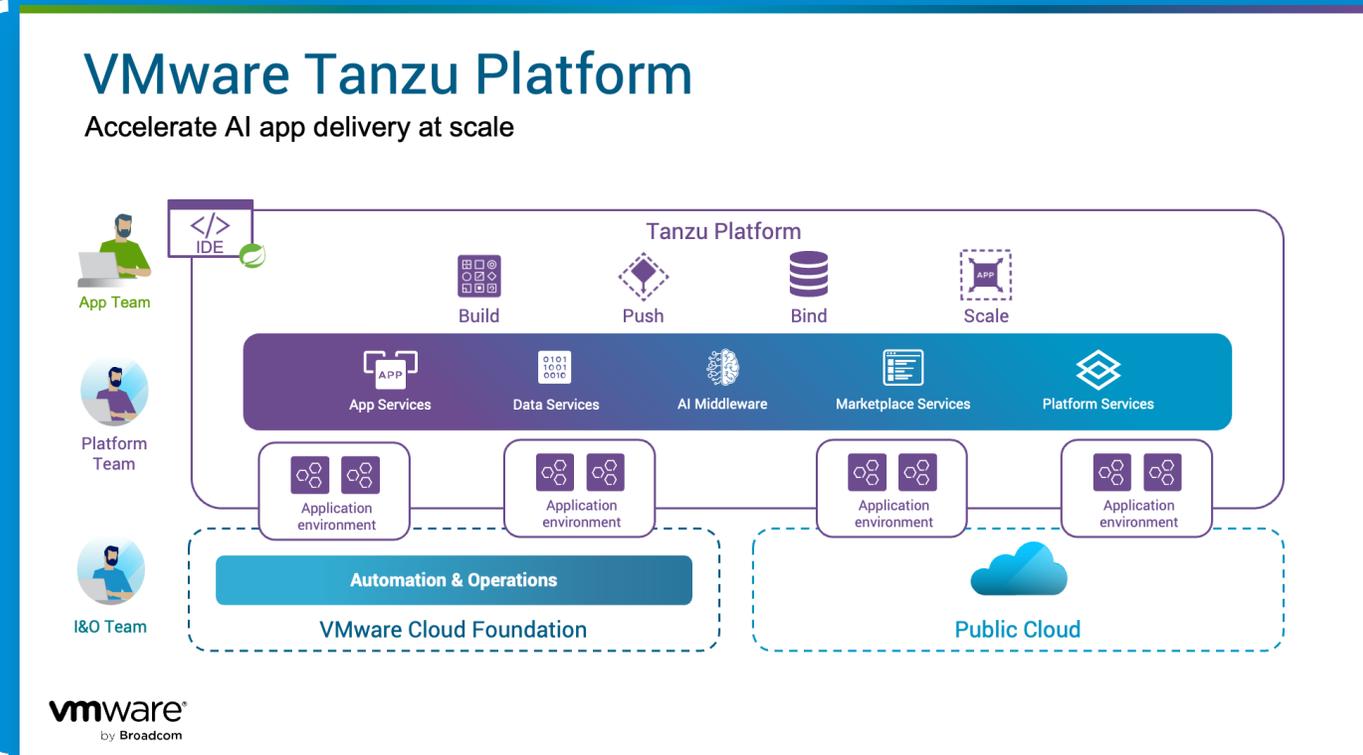
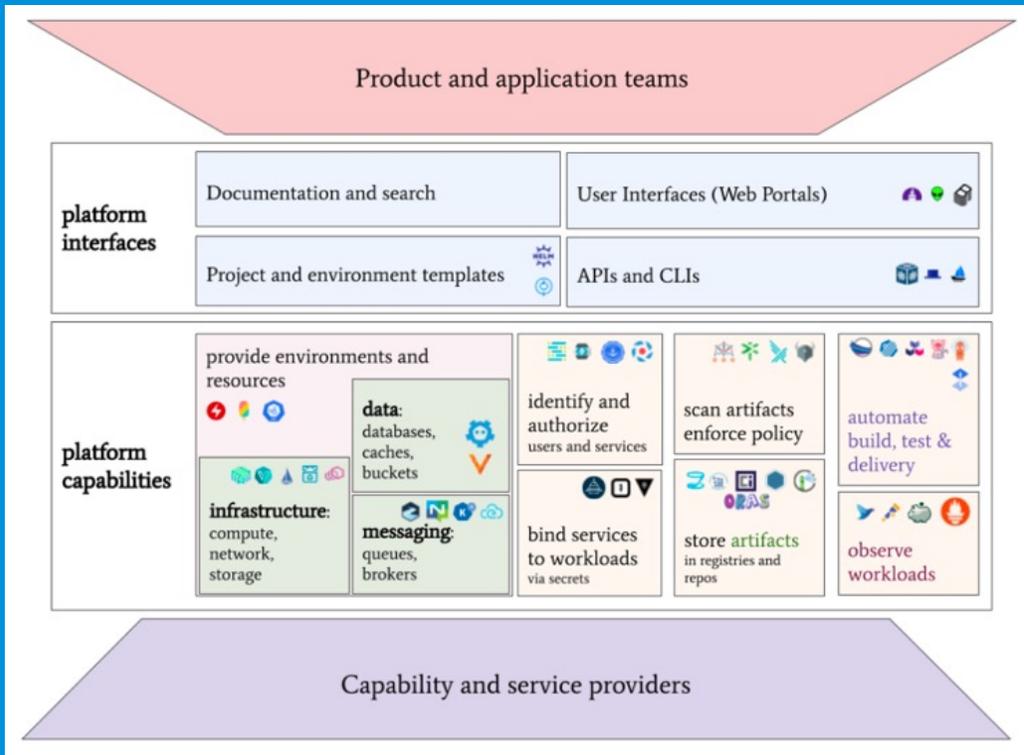
A digital platform is a foundation of self-service APIs, tools, services, knowledge and support which are arranged as a compelling internal product.

[SO THAT] Autonomous delivery teams can make use of the platform to deliver product features at a higher pace, with reduced co-ordination.

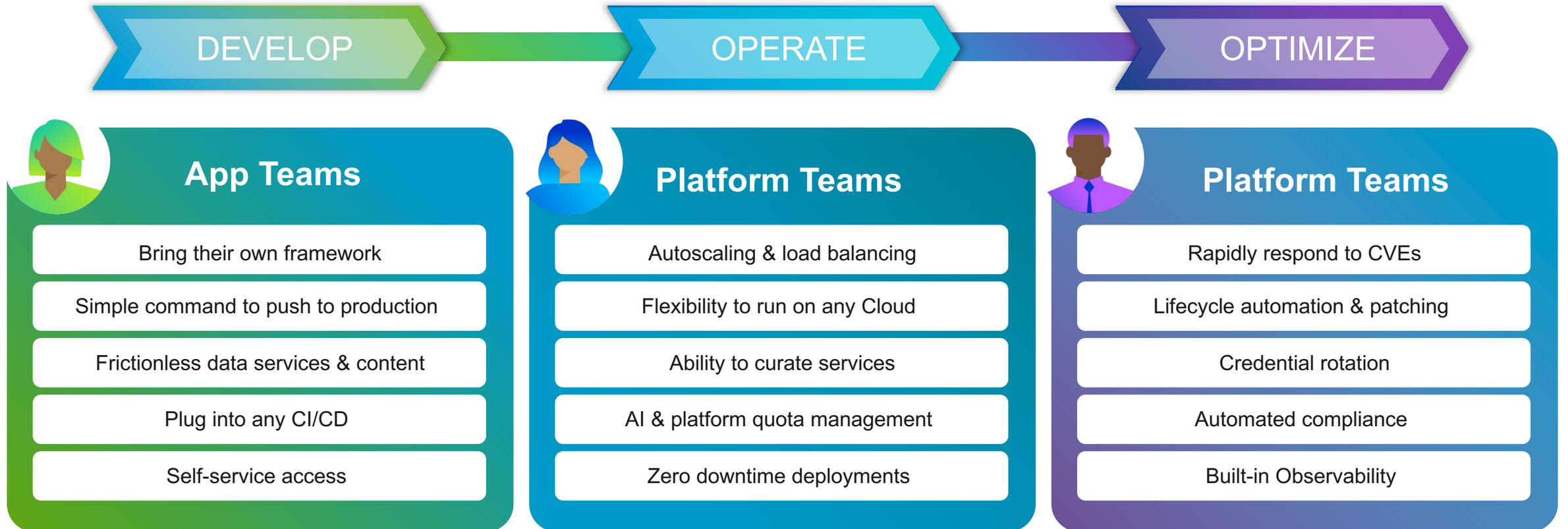
[Evan Bottcher](#), March, 2018

What is a *platform*?

Centralized, standardized stack for building, running, and managing in-house apps.



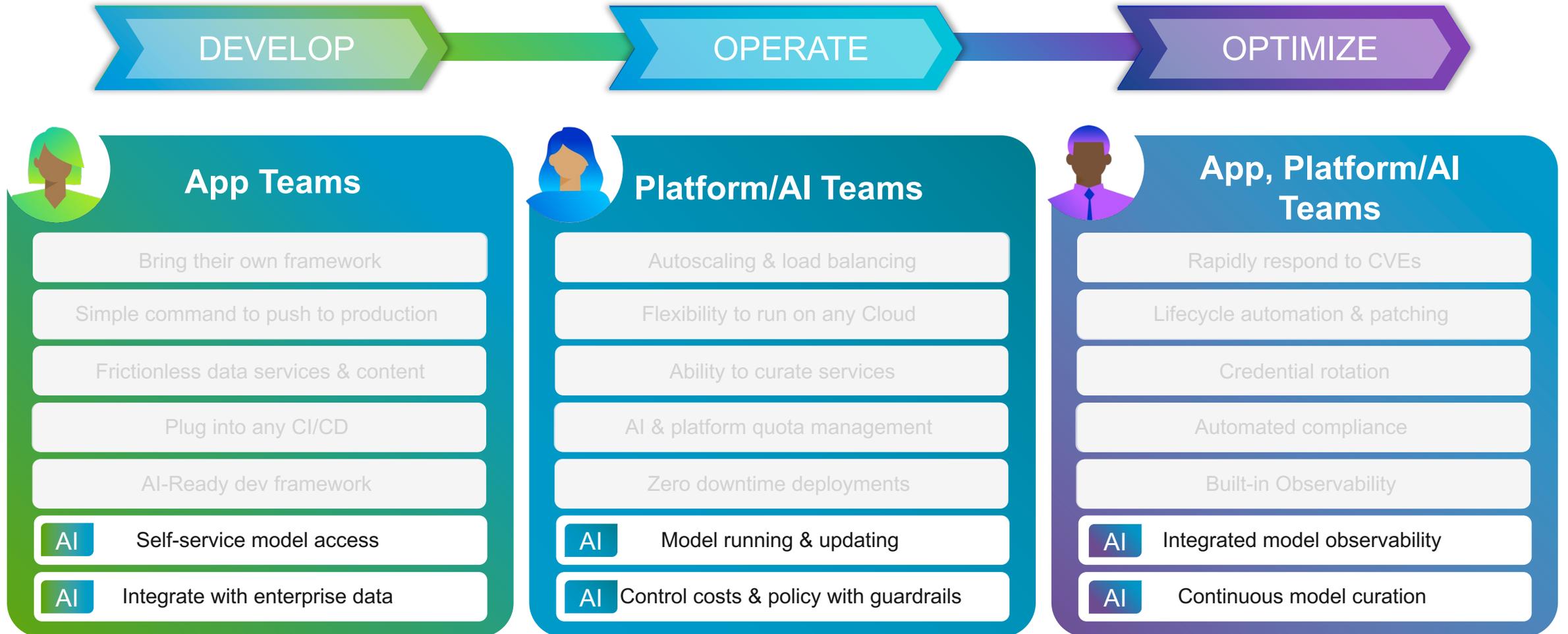
Developer productivity & all the ops -ilities



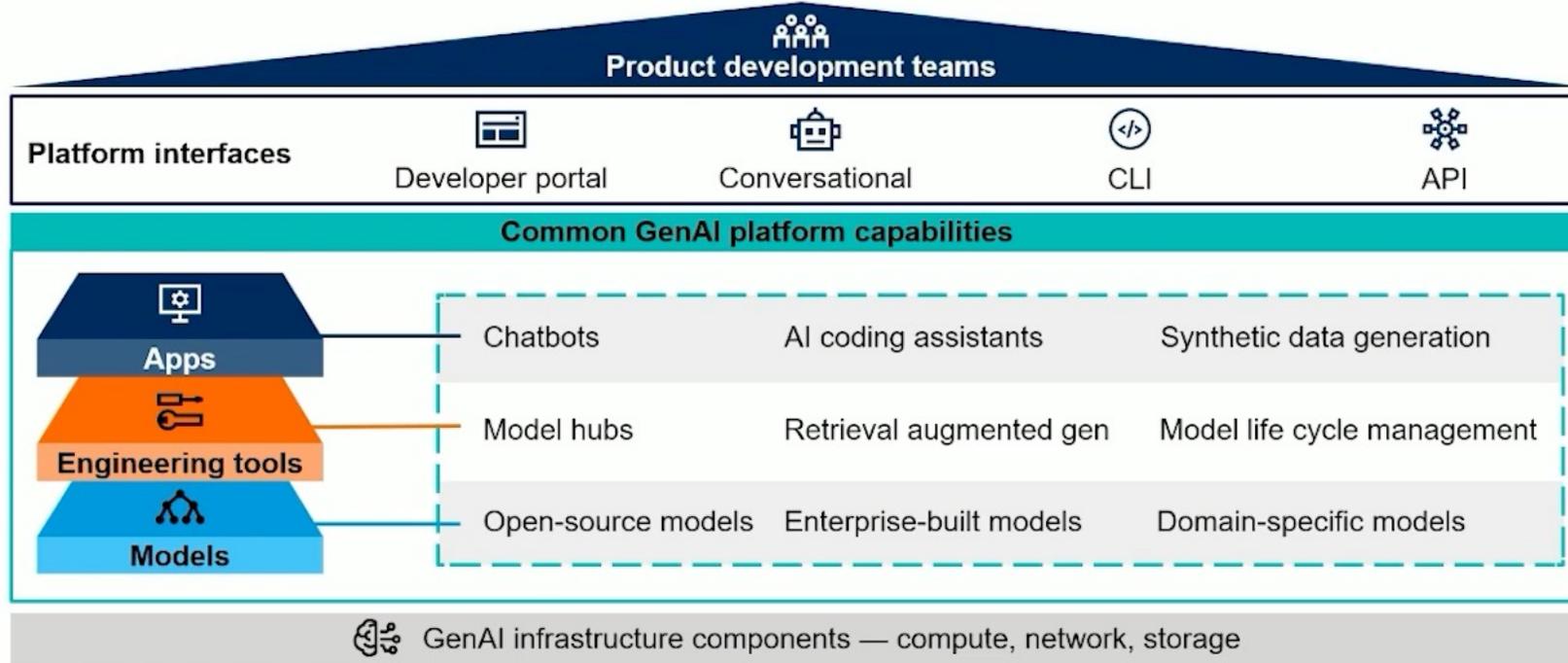


Adding AI...?

A platform treats AI like any other service, adding AI middleware & focusing on new models & frameworks



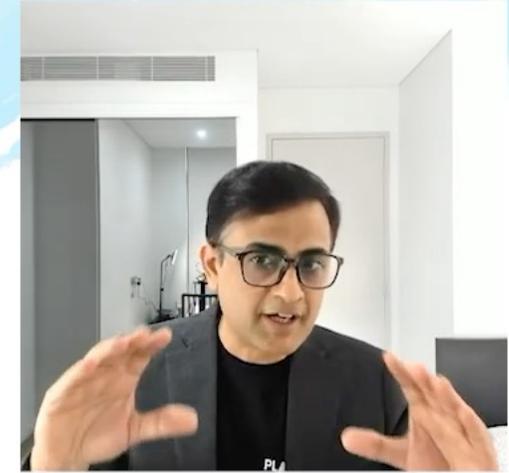
Platform Teams Provide Internal Platform Services to Support Common GenAI Needs



*TRISM — AI trust, risk and security management
Source: Gartner

20 © 2025 Gartner, Inc. and/or its affiliates. All rights reserved. Gartner is a registered trademark of Gartner, Inc. and its affiliates.

Gartner



Manjunath Bhat

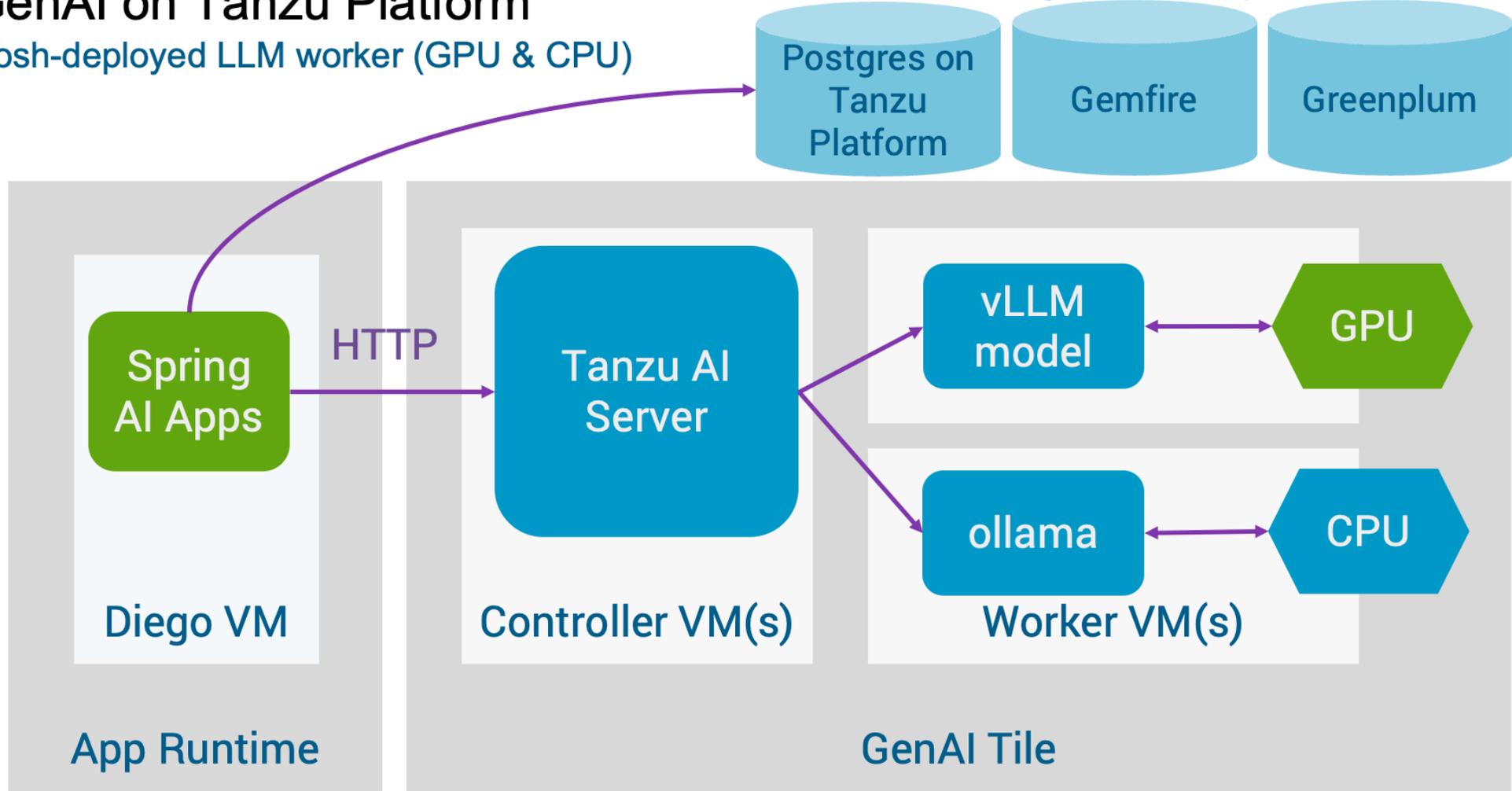
Research VP, Gartner

IMPACT

GenAI on Tanzu Platform

Bosh-deployed LLM worker (GPU & CPU)

“Embeddings/Data Pipelines”



EXPLORE

Broadcom Proprietary and Confidential. Copyright © 2025 Broadcom. All Rights Reserved. The term “Broadcom” refers to Broadcom Inc. and/or its subsidiaries.

8

How do you run a platform?

(in private cloud)

“We are building this platform not for us, we are building it for Mercedes-Benz developers.”

Thomas Müller, Mercedes-Benz



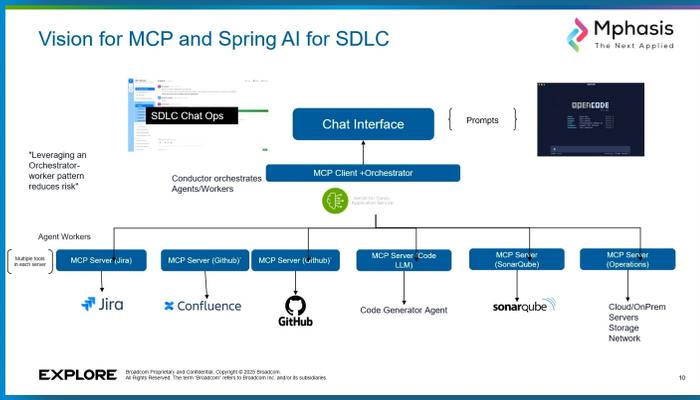
Find the Developer Toil, Confusion, Blockers

Find the Developer Toil, Confusion, Blockers

- What are we making?
- We have a strong vision for our product, and we're doing important work together every day to fulfill that vision.
- I have the context I need to confidently make changes while I'm working.
- I am proud of the work I have delivered so far for our product.
- I am learning things that I look forward to applying to future products.
- My workstation seems to disappear out from under me while I'm working.
- It's easy to get my workstation into the state I need to develop our product.
- What aspect of our workstation setup is painful?
- It's easy to run our software on my workstation while I'm developing it.
- I can boot our software up into the state I need with minimal effort.
- What aspect of running our software locally is painful? What could we do to make it less painful?
- It's easy to run our test suites and to author new ones.
- Tests are a stable, reliable, seamless part of my workflow.
- Test failures give me the feedback I need on the code I am writing.
- What aspect of production support is painful?
- We collaborate well with the teams whose software we integrate with.
- When necessary, it is within my power to request timely changes from other teams.
- I have the resources I need to test and code confidently against other teams' integration points.
- What aspect of integrating with other teams is painful?
- I'm rarely impacted by breaking changes from other tracks of work.
- We almost always catch broken tests and code before they're merged in.
- What aspect of committing changes is painful?
- Our release process (CI/CD) from source control to our story acceptance environment is fully automated.
- If the release process (CI/CD) fails, I'm confident something is truly wrong, and I know I'll be able to track down the problem.
- What aspect of our release process (CI/CD) is painful?
- Our team releases new versions of our software as often as the business needs us to.
- We are meeting our service-level agreements with a minimum of unplanned work.
- When something is wrong in production, we reproduce and solve the problem in a lower environment.

Platform engineers play a critical role in supporting generative AI initiatives within companies. Here are some key responsibilities and tasks they might undertake:

- 1. Infrastructure Development:**
 - Design and build robust cloud or on-premises infrastructure to host AI models and related services.
 - Implement scalable solutions to handle AI and machine learning workloads efficiently.
 - Ensure high availability and resilience of AI systems.
- 2. Model Deployment:**
 - Create and manage pipelines for deploying AI models from development to production.
 - Automate tasks related to model updates, scaling, and rollback if necessary.
- 3. Monitoring and Logging:**
 - Set up monitoring systems to track performance and health of AI models in production.
 - Ensure proper logging mechanisms to record model predictions, usage patterns, and error messages for later analysis.
- 4. Security and Compliance:**
 - Implement security measures to protect data and models, including encryption and access controls.
 - Ensure compliance with data protection regulations and industry standards.



IMC Insurance - Telematics Dashboard

Tracking 15 Drivers | Connected

ID	Status	Location
400001	0.0 mph - BREAK_TIME	Juniper Rd
400002	0.0 mph - BREAK_TIME	Dogwood Dr & East Way
400003	0.0 mph - BREAK_TIME	Marietta Way
400004	22.2 mph - DRIVING	Azalia Way
400006	0.0 mph - BREAK_TIME	Downtown Dr
400007	0.0 mph - BREAK_TIME	Piedmont Dr
400008	0.0 mph - BREAK_TIME	West Peachtree Way & Azalia Dr
400009	0.0 mph - BREAK_TIME	West Peachtree Pkwy
400010	33.9 mph - DRIVING	Juniper Rd
400011	0.0 mph - BREAK_TIME	Peachtree St
400013	0.0 mph - BREAK_TIME	Piedmont Rd

Vehicle details for 400015:
 Vehicle ID: 300015
 Speed: 0.0 mph
 State: POST_CRASH_IDLE
 Street: Buckhead Way
 Route: East St + Six Flags Over Georgia
 G-Force: 0.97g
 Location: 34.025096, -83.453720
 Updated: 2:25:25 PM

Chat (for normals)...

- Your own ChatGPT.
- Customer service.
- Better search.
- Chat-as-UI

Programming...

- New & old code.
- SDLC juicing.
- Trad'l "data science."
- Making pptx?

AI in apps...

- Sales assistants.
- Sloppy integration.
- Science-ing.
- ???

Note: audio, pictures, video are omitted.

Sources: Taznu customers; AI at Goldman, FT, September 14th, 2025; "Leverage Generative AI to Streamline the Software Development Lifecycle," Banu Parasuraman, Andrew Berenato, Explore 2025, August, 2025.

Run and customize your PaaS, don't build it

350 apps supported by **7** platform engineers

300 apps supported by **8** platform engineers

1,200 developers supported by **6** platform engineers

2,500 developers supported by **5** platform engineers

45 app teams supported by **1** ops team

Platform marketing

Organizational Learning



Focus on ways of working.....



Developer centric platform

Internal Resources

- Social Intranet
- GitHub Pages

Developer Trainings

- 3 Day Mercedes Benz specific developer training for Cloud Foundry and K8s
- Udemy Business
- Developer workshops and deepdives

News and Updates

- Regular newsletter
- Internal tech talks
- Quarterly Update Call

Feedback

- Internal Rating Portal with 5-Star Rating and Net Promoter Score
- Follow-Up meetings after go-lives

Communication

- Support via Mattermost, Teams and email
- Requests via GitHub issues
- Changes and incidents via ServiceNow

Consulting

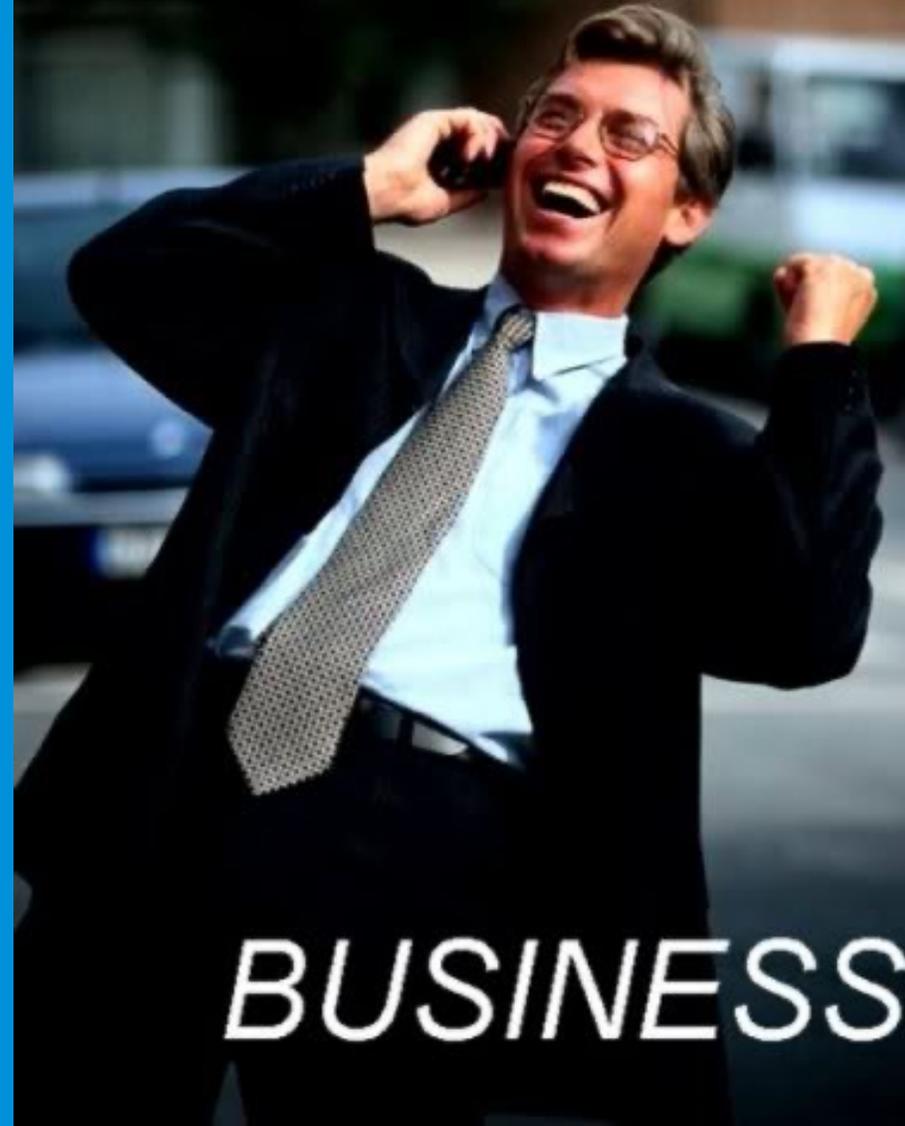
- Weekly consulting and QnA Session
- App cloud readiness scan
- App transformation consulting



“What have
you done for
me lately?”

Tales of ROI, or,
Metrics == Money

HA HA!



BUSINESS

Metrics for the LINE OF BUSINESS

Speed

Velocity is a vector comprised of speed and direction.

We bring a raw speed of advantage to the LOBs and also enable them to rapidly and reliably respond to changes in direction in the service of the business based on user feedback loops.

MEASUREMENTS

- ❑ Time to value (cycle time)
- ❑ Frequency of customer feedback
- ❑ Time between bug identification and fix
- ❑ Time from feedback to deployment of change
- ❑ Customer satisfaction (NPS)
- ❑ Business satisfaction

Stability

Reality is a complex landscape of changing priorities, emergent bugs, evolving architectures, and staffing changes.

We help the LOB achieve resiliency and low volatility as they deliver customer value in the face of this complex reality.

MEASUREMENTS

- ❑ Volatility (std dev in velocity / mean velocity)
- ❑ # of defects generated per developer - year
- ❑ % of software launches / upgrades delayed due to defects
- ❑ Employee satisfaction (ENPS)

Scalability

LOBs need to scale across two dimensions:

People - LOBs strive to attract developers and ramp productivity linearly with personnel.

Apps - LOBs need to rapidly scale their applications and their complexity to handle demand.

MEASUREMENTS

- ❑ # of products in development
- ❑ # of products measuring business success
- ❑ Investment ratios: spend developing software vs operating and systems
- ❑ Disruption caused by doubling workload
- ❑ Ability to attract and retain talent (# of internal referrals)

Security

To move rapidly the team needs to feel secure in making code changes aggressively. Automated test coverage provides this safety net.

To rapidly search for customer value LOBs must adopt a learning culture that fosters psychological safety necessary to fail and learn from failure.

MEASUREMENTS

- ❑ % teams using CI
- ❑ % teams doing TDD
- ❑ Time from commit to deployment

Savings

Teams must reduce risk and waste through small batch delivery and fast consumer feedback.

This drives significant savings as use of the product grows and is key to maintaining their trust and enabling them to go fast, forever.

MEASUREMENTS

- ❑ Fraction of developer time spend writing code and delivering value
- ❑ Product:dev ratio
- ❑ Business satisfaction
- ❑ # of go/no-go decisions based on business success

Metrics for the IT

Speed

IT can efficiently upgrade, patch, and manage the platform.

They rapidly onboard new application teams and provide the necessary services to quickly unblock teams and enable them to deliver consumer value.

MEASUREMENTS

- ❑ # prod/dev deploys per month
- ❑ # platform upgrades per month
- ❑ Platform upgrade speed
- ❑ # of new apps onboarded/month
- ❑ Team distribution of skills

Stability

Our customers entrust us with their production workloads and their developer productivity.

We must provide adequate SLOs to meet their needs and earn their trust by ensuring compatibility and uptime across platform upgrades.

MEASUREMENTS

- ❑ Minutes of prod outage per year
- ❑ Minutes of dev outage per year
- ❑ Mean time to recovery
- ❑ Mean time between failures
- ❑ # of upgrade-related failures

Scalability

IT needs to provide an “at-scale” service on-demand at the whim of the business.

They need to explore all options with minimal friction as they grapple with the mix of workloads on-premise and in the cloud.

MEASUREMENTS

- ❑ Queries per second
- ❑ # of AIs per foundation
- ❑ # of SIs per foundation
- ❑ # of foundations
- ❑ # of teams using the platform
- ❑ Does increasing workload on existing

Security

Security is a paramount concern for our customers. We earn their trust by providing a platform that is secure by default.

We solve for security and reduce security-related friction and toil in order to enable our customers to go fast, forever.

MEASUREMENTS

- ❑ Time between identifying and patching a CVE
- ❑ Cost in person-hours or dollars of leaked credential
- ❑ Fraction of operator time spent on security configuration
- ❑ # of disruptions/suspensions due to security concerns

Savings

IT must meet the needs of thousands of developers within tight budgetary constraints.

We provide a platform that simultaneously reduces complexity and sprawl and improves the ops:dev ratio.

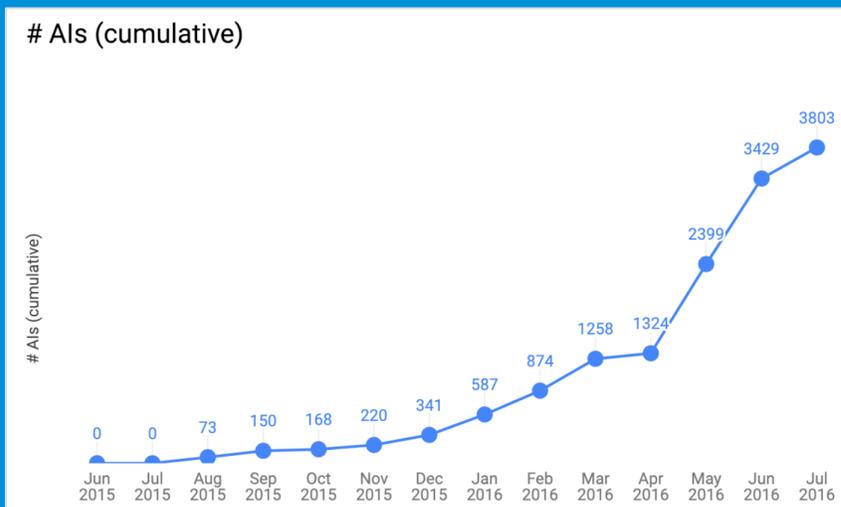
MEASUREMENTS

- ❑ Operator:developer ratio
- ❑ # of apps per operator
- ❑ # of foundations per operator
- ❑ Degree of automation for provisioning, build, test, change approval governance, deployment, perf

Scaling Phase – Pairing & Seeding to build trust & training



1. Create platform marketing program.
2. Find two to five more apps.
3. Pair & seed from first dev & platform team to new teams.
4. "Shift Left" – build golden paths for governance, security, etc.
5. Add more infrastructure staff with pairing & seeding.
6. Do this for three months.
7. Repeat, growing number of apps as pairing & seeding allows.



Running a platform: Platform Engineering Maturity Model

Aspect		Provisional	Operational	Scalable	Optimizing
<u>Investment</u>	<i>How are staff and funds allocated to platform capabilities?</i>	Voluntary or temporary	Dedicated team	As product	Enabled ecosystem
<u>Adoption</u>	<i>Why and how do users discover and use internal platforms and platform capabilities?</i>	Erratic	Extrinsic push	Intrinsic pull	Participatory
<u>Interfaces</u>	<i>How do users interact with and consume platform capabilities?</i>	Custom processes	Standard tooling	Self-service solutions	Integrated services
<u>Operations</u>	<i>How are platforms and their capabilities planned, prioritized, developed and maintained?</i>	By request	Centrally tracked	Centrally enabled	Managed services
<u>Measurement</u>	<i>What is the process for gathering and incorporating feedback and learning?</i>	Ad hoc	Consistent collection	Insights	Quantitative and qualitative

Getting ready for a long journey

Questions, topics, & analysis

1. **Expect 3 to 5 years: “Small organizations might swim through level three, while large enterprise organizations with lots of heritage and legacy, might be here for years.”**
2. **Do you have and/or need microservices and 12 factor apps? How many new apps do you have?**
3. **How will you modernize your existing apps? Would lift-and-shift or leave alone be better?**
4. **How will you standardize & centralize infrastructure, [dev|ops|security|compliance] practices?**
5. **How can you make continuous learning stick?**
6. **How will you change the organization and team definition, structures, and incentives?**
7. **How will you get The Business to work weekly with the developers?**
8. **Costs do not decrease until you decommission old infrastructure and legacy apps.**

Get an AI App Up and Running in Under 2 Weeks

What is it:

Developed by innovators of Cloud Foundry, the AI Starter Kit provides a seamless path to your first AI application on Tanzu Platform.

Who is it for:

Platform engineers, application developers and data teams can utilize AI Starter Kit to develop patterns for scaling AI apps.



Why it's needed:

- Take AI out of the shadows and support internal innovation with minimal friction
- De-risk AI exploration and prepare for scaled AI app rollout
- Ensure your models are ready for primetime with performance, accuracy and ROI

TryTanzu.ai

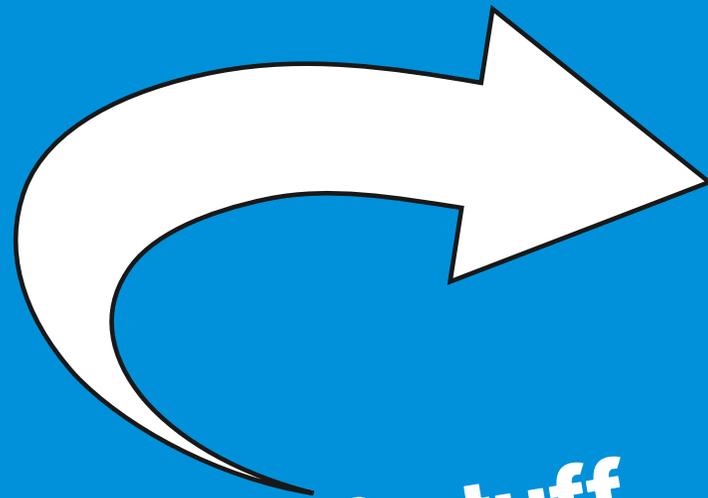
Expected results:

- Frictionless AI Prototyping
- Forecast AI needs more accurately
- Identify optimal delivery strategies
- Optimize production performance
- Reduce GPU dependency

Thanks!

 <https://newsletter.cote.io/>

 cote@broadcom.com



Slides & stuff

