

Turning an incident report into a design issue with TLA+

Finn Hackett (UBC) & Markus A. Kuppe (MSFT)

By Some Not-SREs



A. Finn Hackett

PhD Student at U. of British Columbia

Programming Languages Researcher

→Intern at Microsoft Research S'22



Markus A. Kuppe

Principal Research Engineer (MSFT)

TLA+ eng, Postmortem Reviewer

→Finn's internship mentor S'22



Joshua Rowe

Principal Engineer (MSFT)

Azure CosmosDB

→Domain expert

Thanks to postmortem owner **Ben Pannell** (MSFT)!

Highlights



We learned about a high-profile 28-day incident at Microsoft Azure.



Mitigation was a feature revert (full repair a costly multi-year redesign).



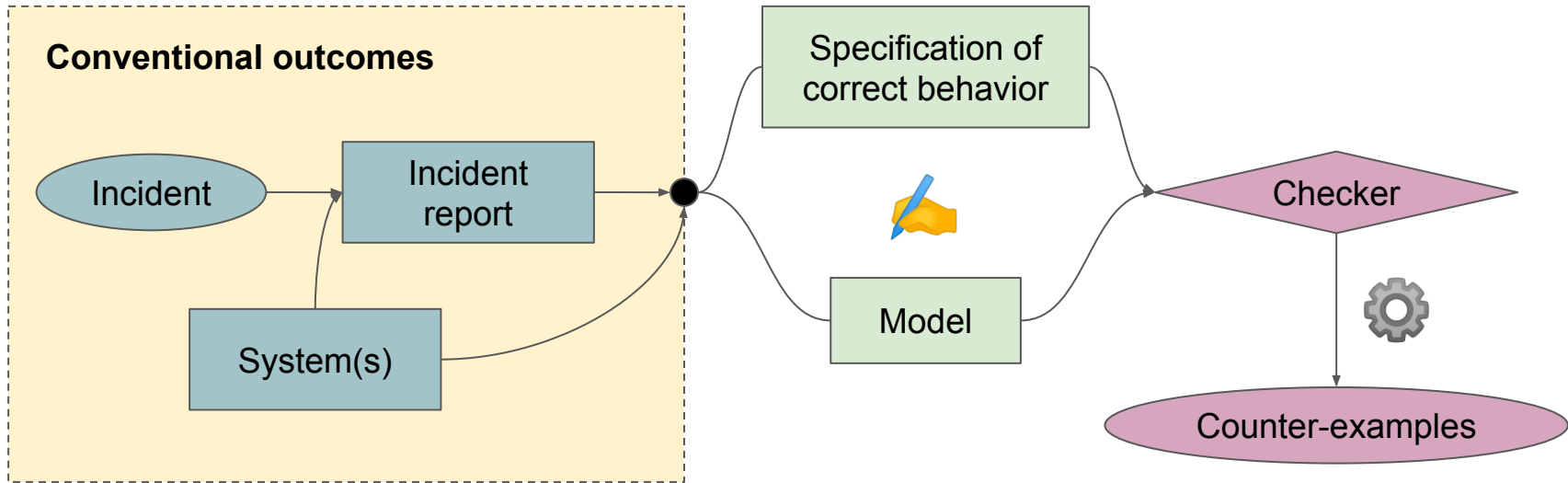
Reproducing the issue at small scale is impractical...



Idea: use modeling to unambiguously document the problem.

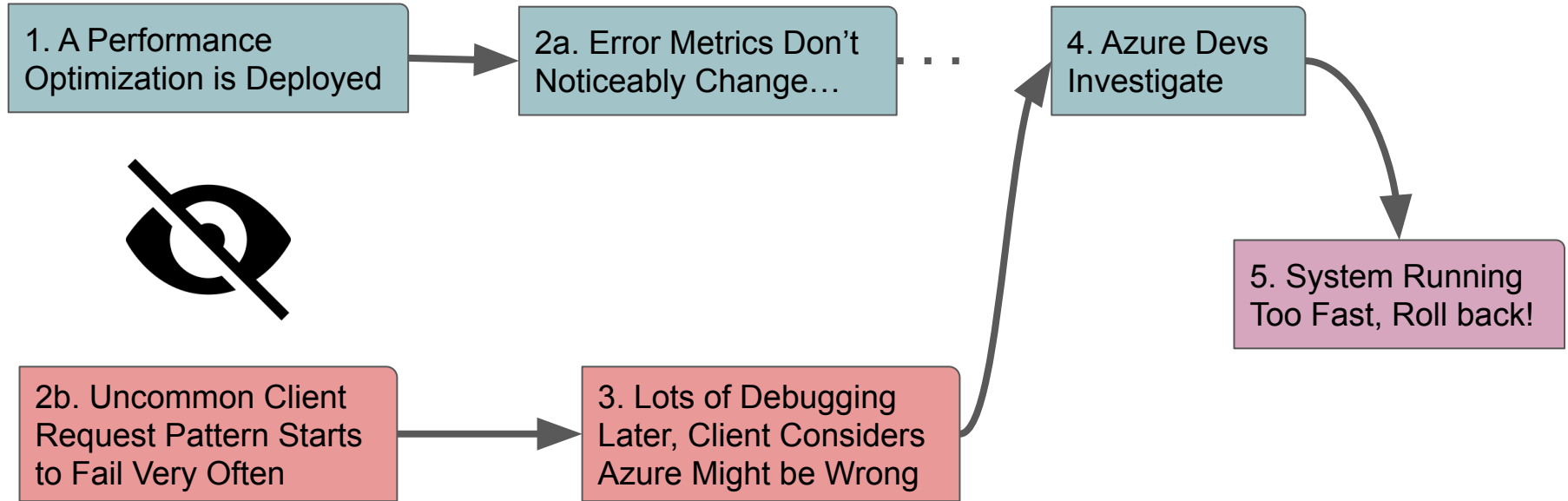
We use TLA+ here, but this applies to similar tools as well.

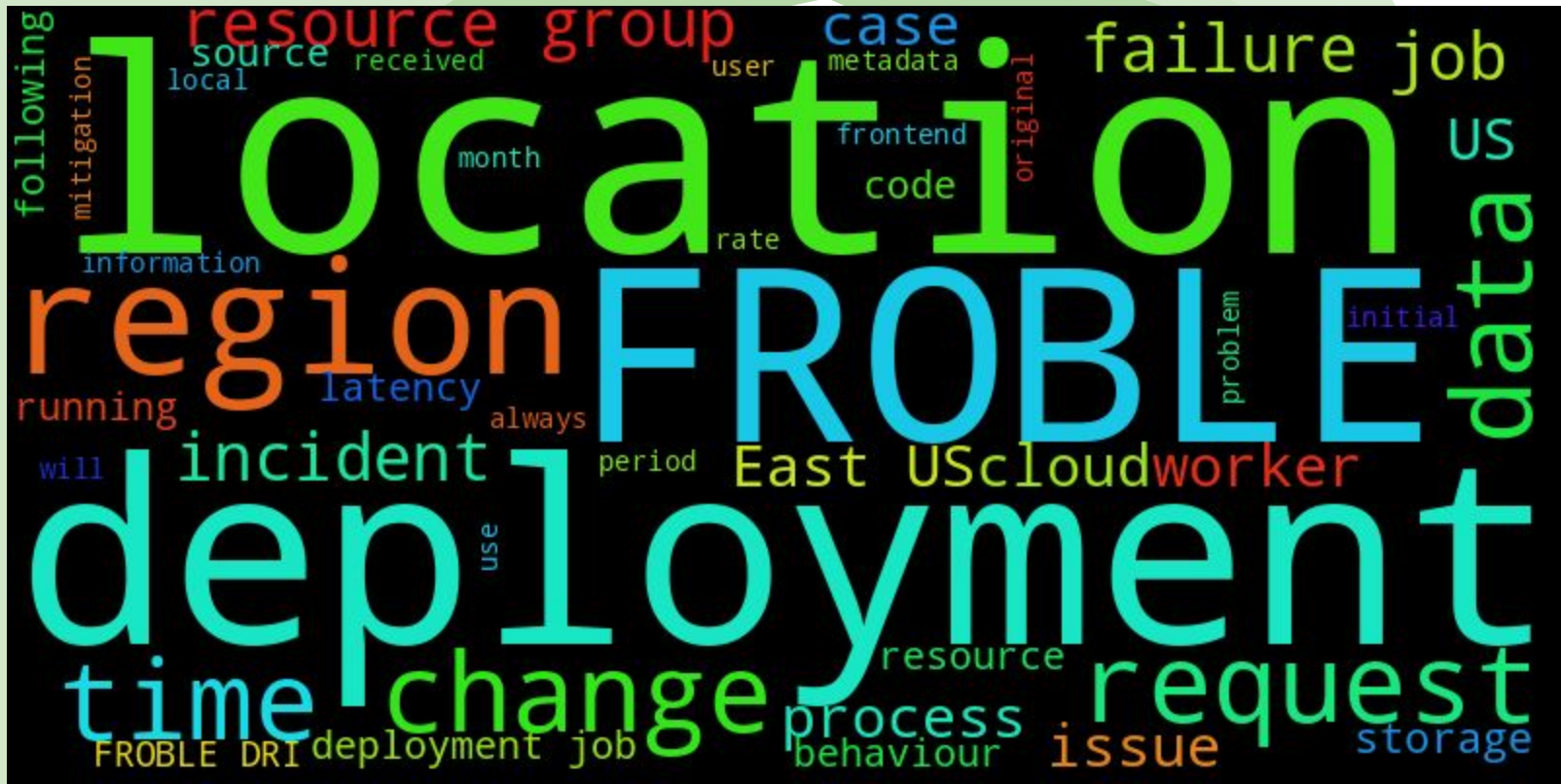
A Workflow Beyond Incident Reports



💡 ... which summarize the issue at a high level.

Story Time: an Incident





s/\${real product name}/FROBLE

What's the Big Deal?

Client is not seeing errors anymore, but...

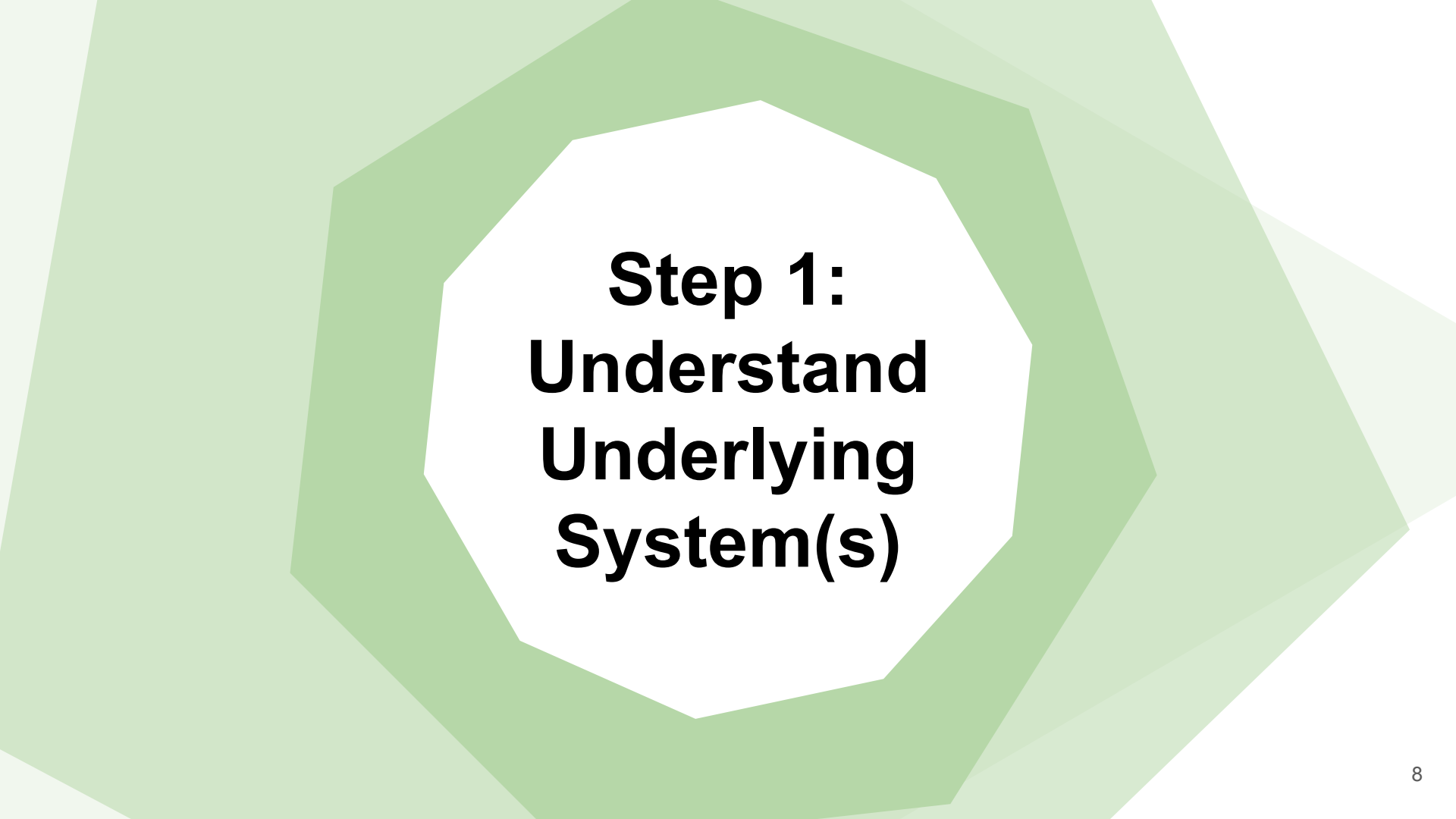
😞 That took a long time to figure out.

(28 days to fix)

😞 We know why we rolled back, but design-level insight is missing.

🤔 What about the design made a component being *too fast* a problem?

We want to clearly model and understand what happened at the conceptual level.



**Step 1:
Understand
Underlying
System(s)**

What is Azure CosmosDB?

⚙️ **Need to understand what the system is supposed to do to model it.**
What we present next is based on what we learned by drafting our model.

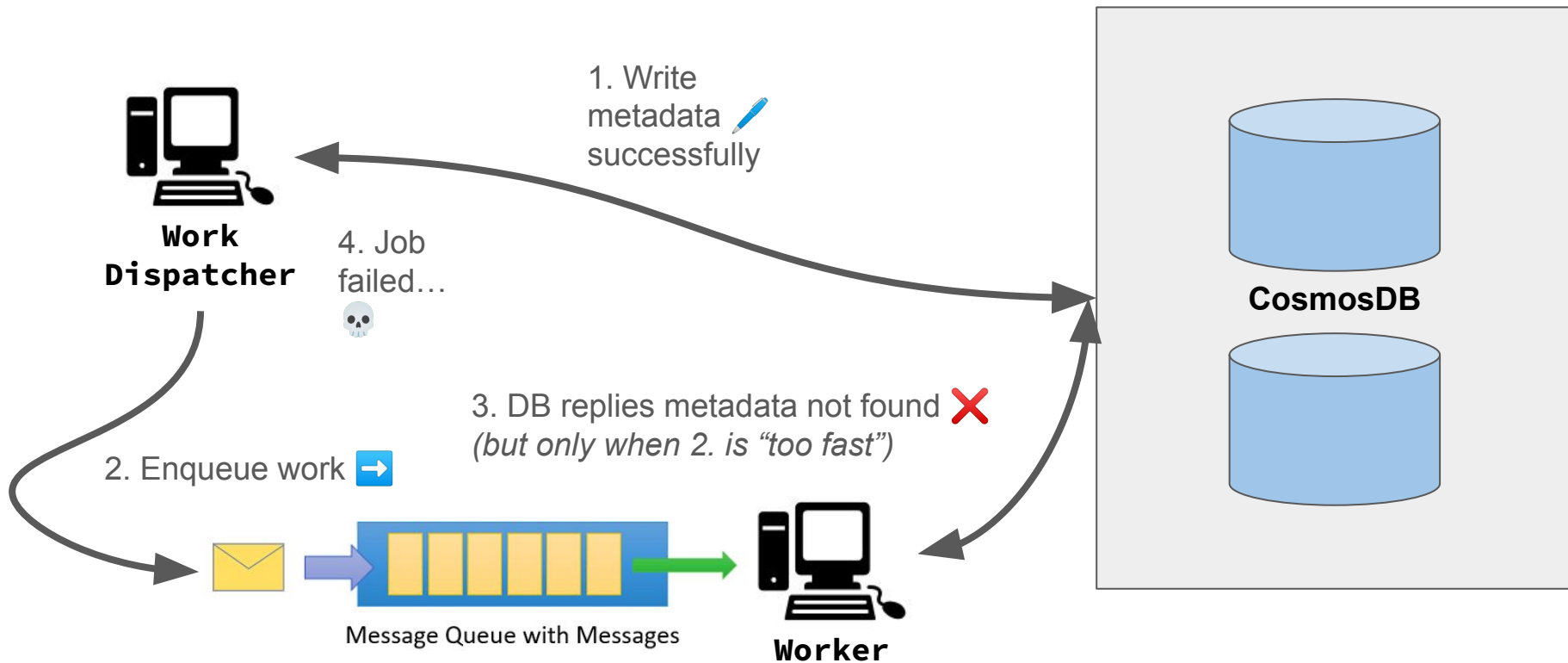
🌍 A planet-scale key-value store - a big distributed system.

📁 Stores your data as key-value mappings, useful for app/service data.

🤔 Similar to? Amazon DynamoDB, Google Firestore...

↩️ *Also using TLA+*

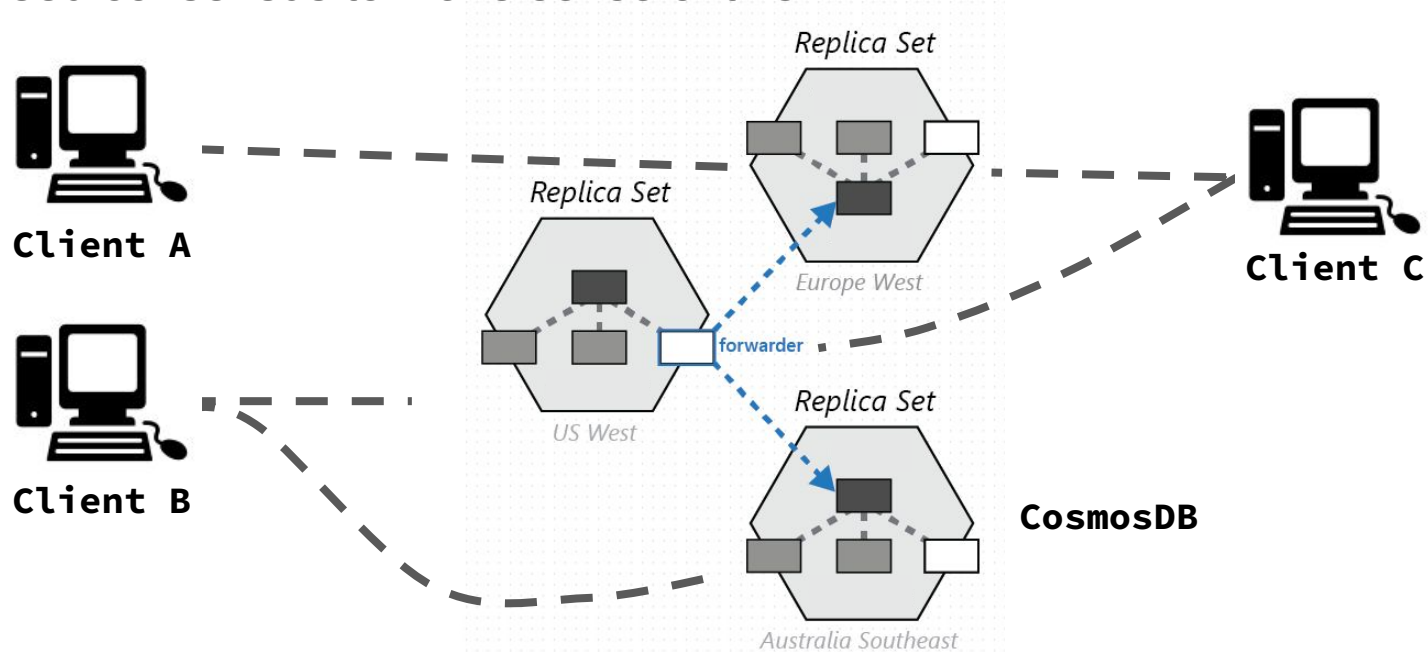
Key Detail: Communication with Azure Cosmos DB



CosmosDB: Consensus Machine

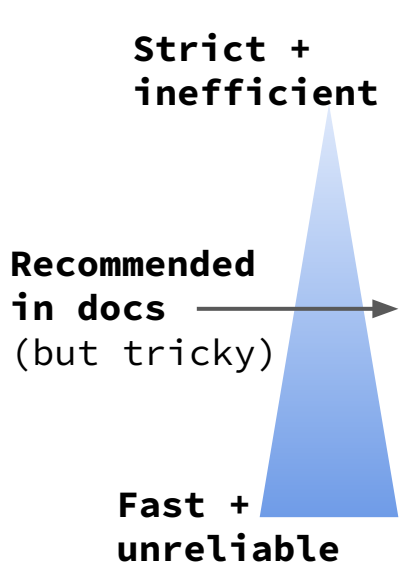
📧 Clients exchange get/set messages with different servers.

👉 Need consensus to make sense of this.

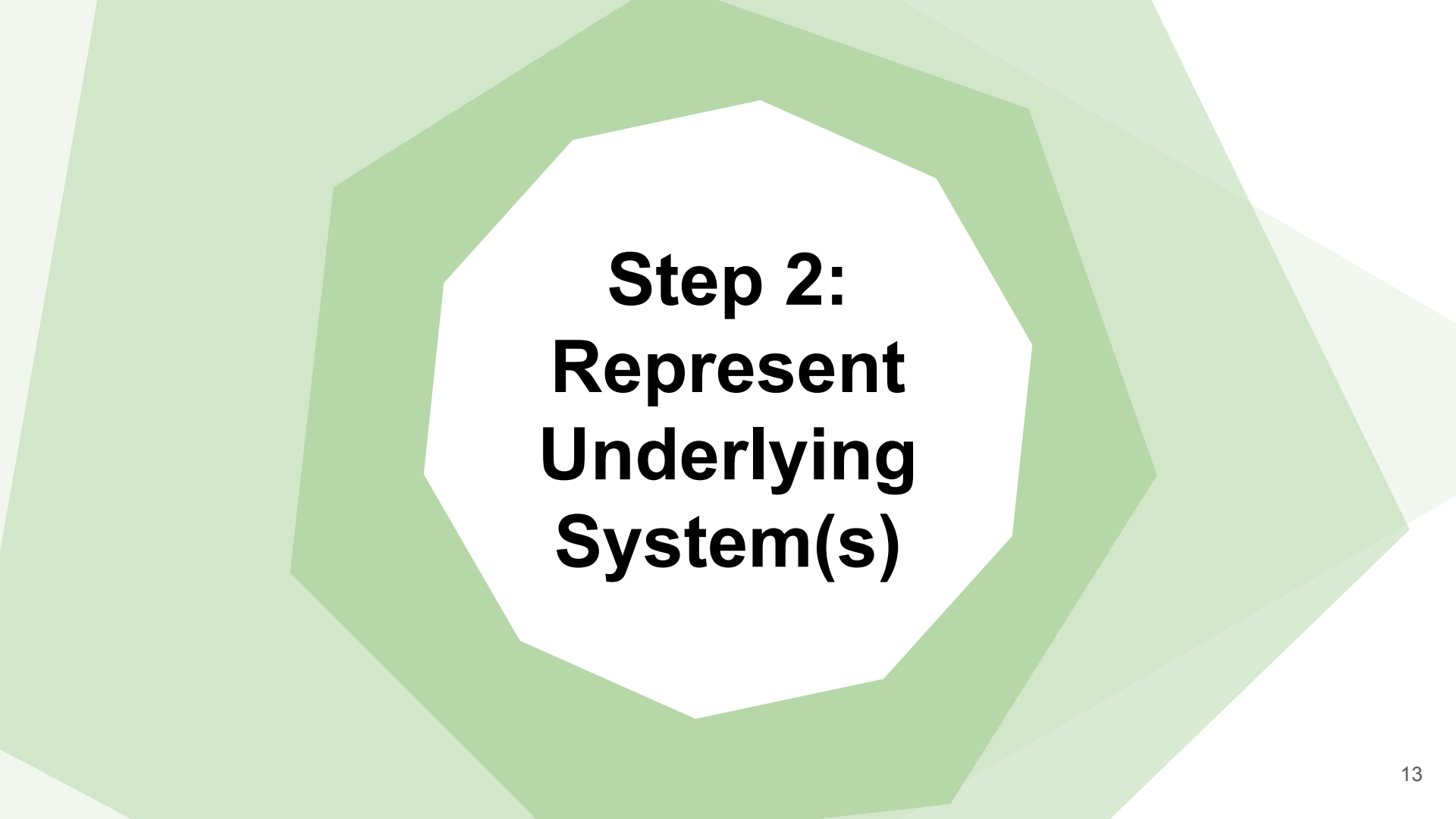


CosmosDB: Consistency Levels

🔧 5 consistency levels control how strictly CosmosDB maintains consensus.



Consistency Level	Effect on Client View
Strong Consistency	Global order; clients always see latest versions
Bounded Staleness	Old values visible for limited (<i>bounded</i>) time
Session Consistency	Synchronize only between clients with shared token
Consistent Prefix	Generally similar to Eventual Consistency
Eventual Consistency	Old values should eventually stop showing up



**Step 2:
Represent
Underlying
System(s)**

CosmosDB: What to Model for Our Purpose?

"The purpose of abstraction is not to be vague, but to create a new semantic level in which one can be absolutely precise."

– Edsger W. Dijkstra

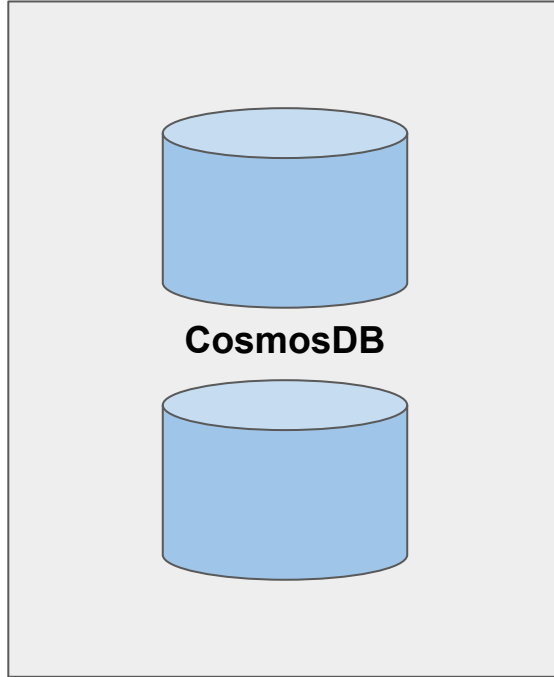
No

- ✗ Detailed client API syntax
- ✗ Server management
- ✗ Anything client can't see

Yes

- ✓ Plain key-value r/w
- ✓ Client view of servers
- ✓ Anything client can see
(especially unusual things)

... But Have a Domain Expert Model It



CosmosDB: the TLA+ Model

5✓ Simulate concurrent key-value reads and writes at all 5 consistency lvls.

🕒 Reusable, complete model: took 3 months to build w/ dev input.

🌟 But it took 1 day to use in this postmortem.

GitHub: <https://github.com/tlaplus/azure-cosmos-tla/tree/master/simple-model>

Understanding Inconsistency in Azure Cosmos DB with TLA+

Titus Barik University of Illinois Chicago titus@uic.edu	Johns Borek Microsoft johnsb@redhat.com	Matias Alvarez Kreppe Microsoft matias@redhat.com
---	--	--

Abstract—Formal implementation correctness of a distributed system is a complex problem to understand, verify, and maintain. This paper presents a formal model of the Cosmos DB replication and consistency logic. The model is implemented in TLA+ and verified using the TLC model checker. The model captures the state of the system, the state of the replication, and the state of the consistency. The model is used to verify the correctness of the Cosmos DB replication and consistency logic. The model is also used to generate test cases for the Cosmos DB replication and consistency logic. The model is implemented in TLA+ and verified using the TLC model checker. The model captures the state of the system, the state of the replication, and the state of the consistency. The model is used to verify the correctness of the Cosmos DB replication and consistency logic. The model is also used to generate test cases for the Cosmos DB replication and consistency logic.

I. INTRODUCTION

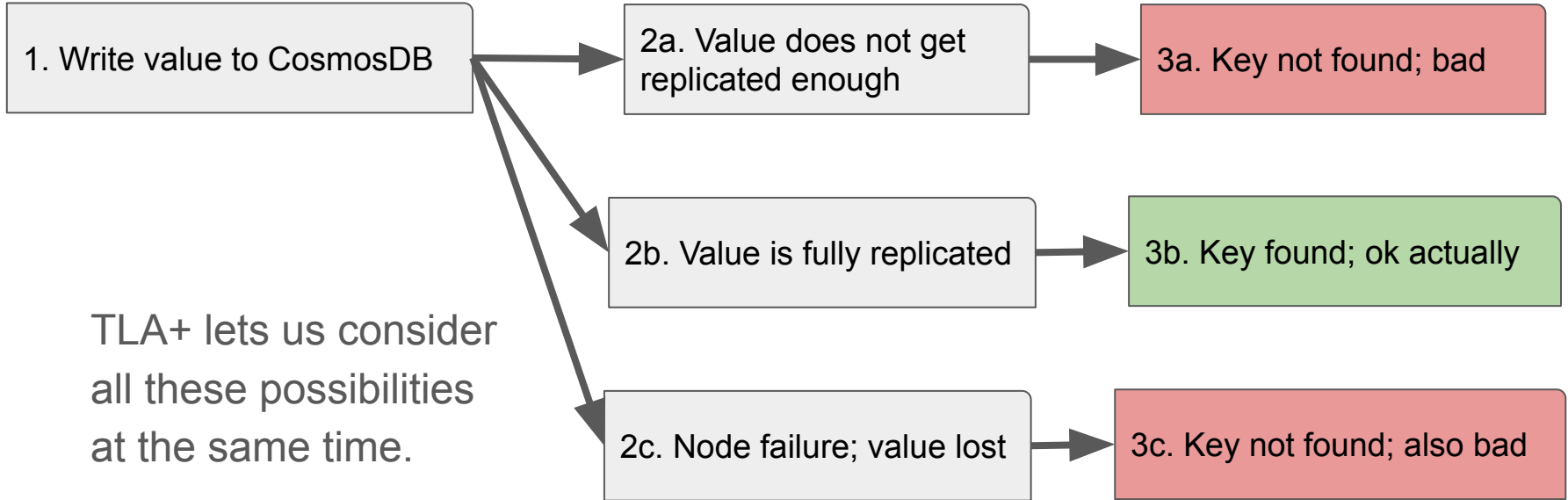
Formal verification is a powerful technique for verifying the correctness of distributed systems. In this paper, we present a formal model of the Cosmos DB replication and consistency logic. The model is implemented in TLA+ and verified using the TLC model checker. The model captures the state of the system, the state of the replication, and the state of the consistency. The model is used to verify the correctness of the Cosmos DB replication and consistency logic. The model is also used to generate test cases for the Cosmos DB replication and consistency logic.

See our conference paper for all the details:
<https://doi.org/10.48550/arXiv.2210.13661> (preprint)

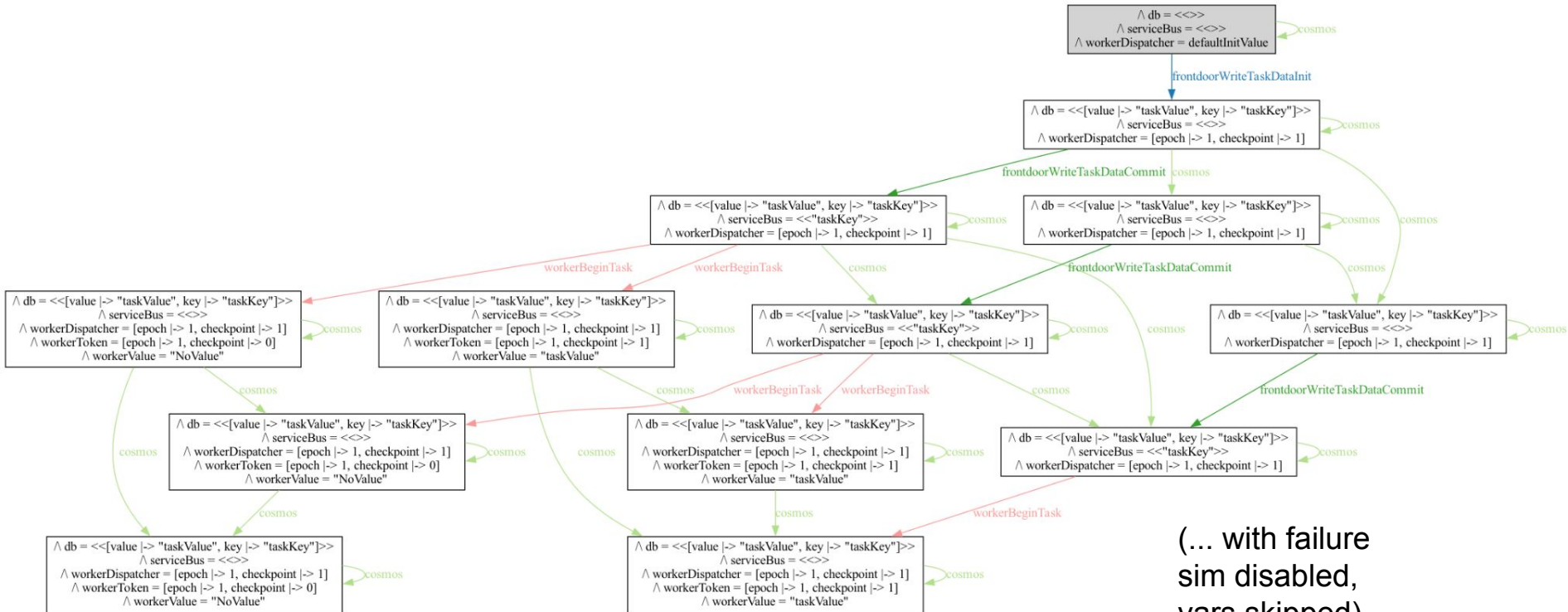


**Step 3: Model
the Incident
Demo Time!**

Thinking in States: Branching Possibilities

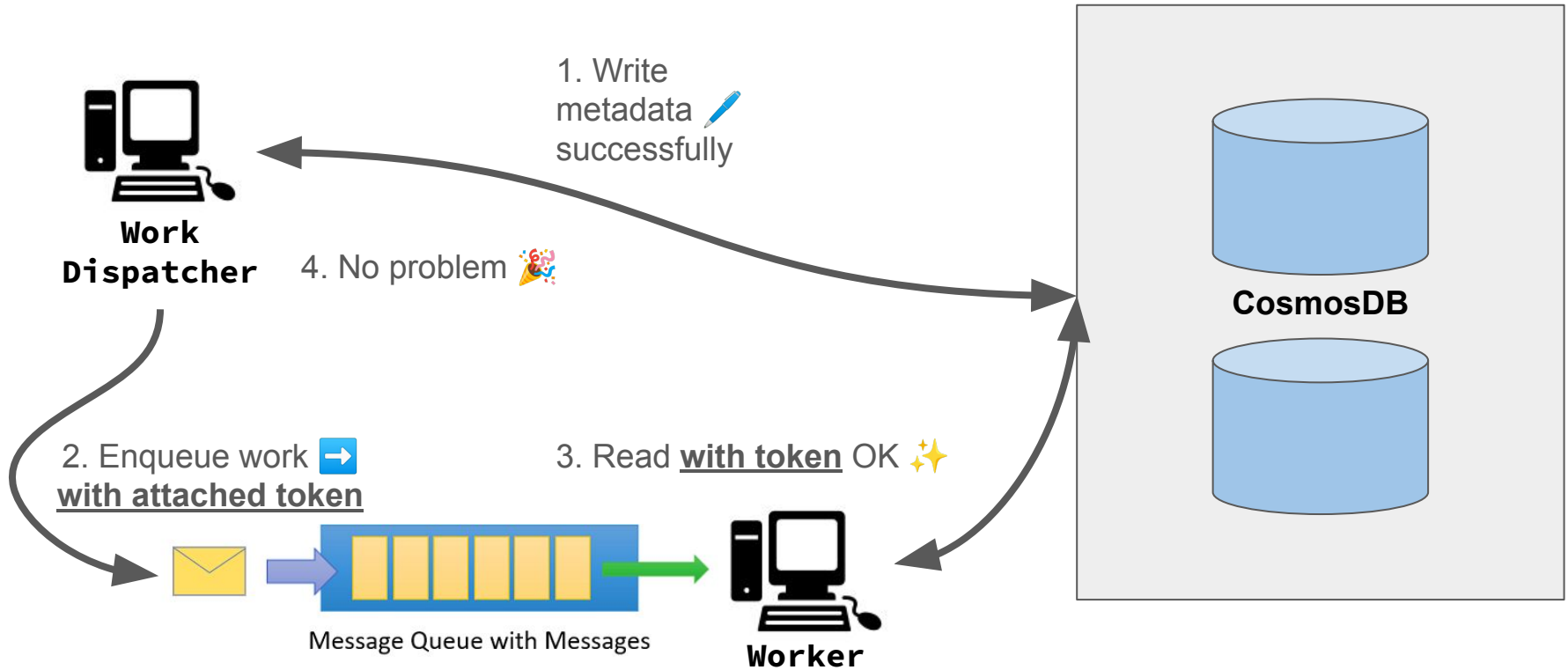


Thinking in States: the Full




(... with failure sim disabled, vars skipped)

The Corrected Design




What We Learned

 Problem is using session consistency without sharing tokens.

 Original mitigation changed the probability of the issue, but wasn't a fix. The issue was always there, just very rare.

 Clear when presented this way, but issue was hidden behind APIs.

 Engineers were not surprised by our results - we confirmed what was suspected but could not be demonstrated.

 Modeling helped us think and investigate.

Summary

Reproduced incident spanning multiple foundational systems. 🎉

Our modeling workflow can go beyond current incident reports.

You can do this too. Learn at <http://tlapl.us>.

TLA+ model (linked from official docs):

<https://learn.microsoft.com/en-us/azure/cosmos-db/consistency-levels>.

Paper: <https://doi.org/10.48550/arXiv.2210.13661>.

Any Questions?

