# Gray Failure: The Achilles' Heel of Cloud-Scale Systems

Ryan Huang, Chuanxiong Guo, Lidong Zhou, Jacob R. Lorch,
Yingnong Dang, Murali Chintalapati, Randolph Yao
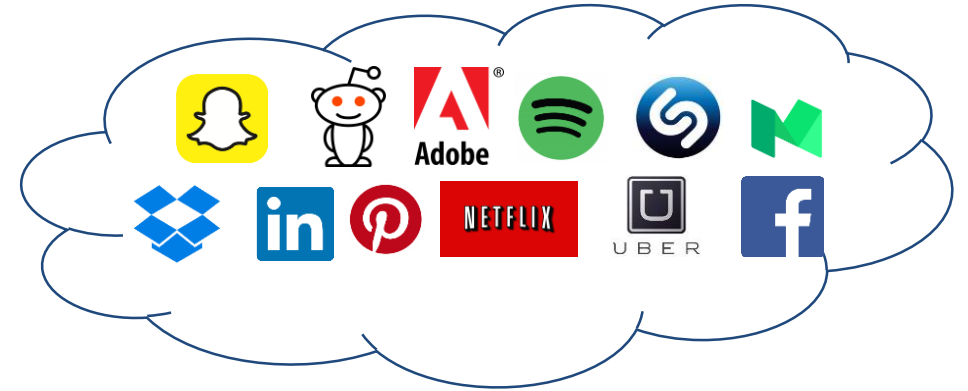
Microsoft

JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

# Outline

❑ **Background & the gray failure problem**

❑ **Real-world gray failure cases in Azure**

❑ **A model and a definition for gray failure**
   ❖ differential observability

❑ **Potential future directions**

# Rapid Growth of Cloud System Infra

» **Software user shift**
  - direct: e.g., office 365, Google Drive
  - indirect: e.g., Netflix on AWS

» **Workload diversity**
  - website, workflow, big data, machine learning

» **Internal composition**
  - more data centers, larger cluster, special h/w
  - containerization, micro-services

# Demanding Requirement on Availability

» **Users intolerant of service downtime**



» **Failure more costly**
- SLA violation, reputation hit, customer loss, engineering resource waste
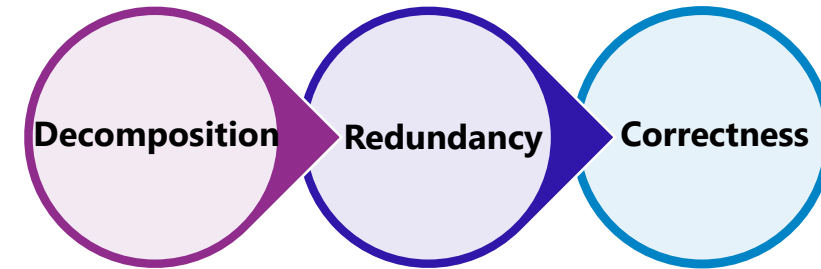
» **New availability bar**
- 3 nines to 5 or 6 nines

# Key: Embrace Fault-tolerance!
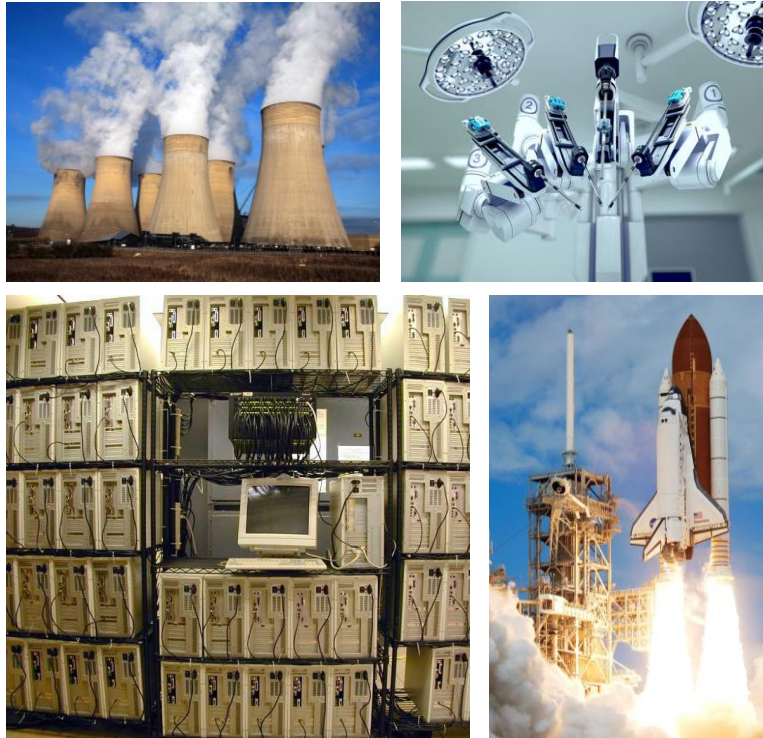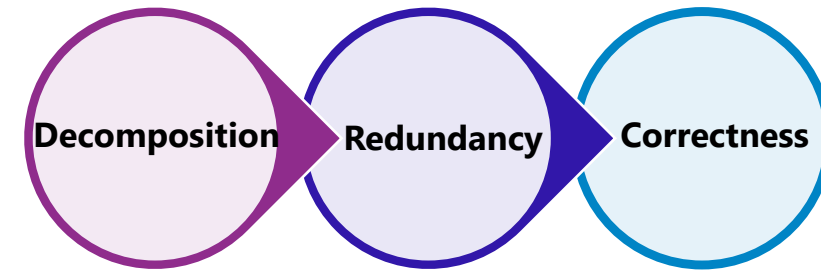
**Rich history since 1950s**

**Steps:**

Decomposition → Redundancy → Correctness

# Key: Embrace Fault-tolerance!

**Rich history since 1950s**

**Steps:**



Decomposition → Redundancy → Correctness

**Building block:**

| process pair | triple modular redundancy | Primary/backup |

| RAID | state machine replication | checkpoint | chain replication |

| transaction | 2PC | Gossip | N-version | erasure coding |

| virtual synchrony | Paxos | Zab | PBFT | Zyzzyva | ... |

# Status Quo

# Status Quo

» **By and large, the efforts paid off**

- many faults successfully detected, tolerated, and repaired every day
- few global outages
- 99.9% is achievable

» **But moving forward...**

- simplistic assumptions start to break
- reasoning about availability becomes hard
- frequent bizarre phenomenon in production
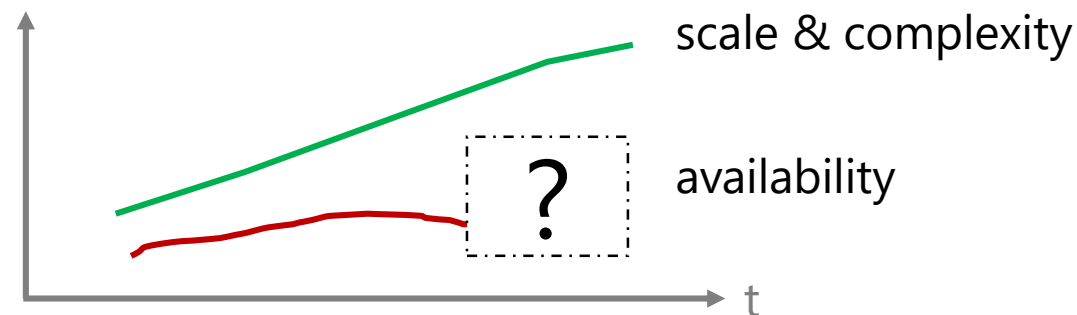- 99.999% and beyond is challenging

# Status Quo

» **By and large, the efforts paid off**

- many faults successfully detected, tolerated, and repaired every day
- few global outages
- 99.9% is achievable

» **But moving forward...**

- simplistic assumptions start to break
- reasoning about availability becomes hard
- frequent bizarre phenomenon in production
- 99.999% and beyond is challenging



scale & complexity

availability

?

t

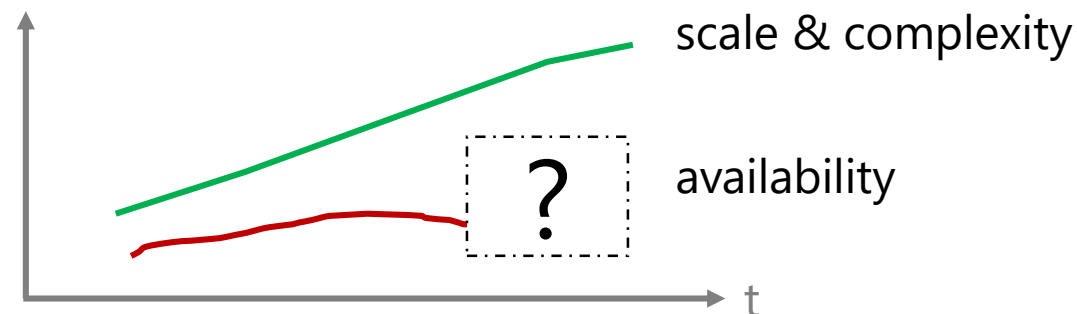# Status Quo

» **By and large, the efforts paid off**
- many faults successfully detected, tolerated, and repaired every day
- few global outages
- 99.9% is achievable

» **But moving forward...**
- simplistic assumptions start to break
- reasoning about availability becomes hard
- frequent bizarre phenomenon in production
- 99.999% and beyond is challenging

} A common theme: the overlooked gray failure problem

scale & complexity

? availability

t

# The Elephant in the Cloud - Gray Failure

A component either works correctly or stops

**Fail-stop**

process pair    TMR

RAID    primary backup

2PC    chain replication

Paxos    erasure coding

virtual synchrony    Zab

...

# The Elephant in the Cloud - Gray Failure

A component either works correctly or stops

**Fail-stop**

process pair   TMR

RAID   primary backup

2PC   chain replication

Paxos   erasure coding

virtual synchrony   Zab

...

A component may behave arbitrarily

**Byzantine**

PBFT   Zyzzyva

Q/U   BAR Gossip

Aliph   UpRight

...

# The Elephant in the Cloud - Gray Failure

A component either works correctly or stops

A component appears to be still working but is in fact experiencing severe issue

A component may behave arbitrarily

**Fail-stop**

**Gray failure**

**Byzantine**

process pair | TMR

RAID | primary backup

2PC | chain replication

Paxos | erasure coding

virtual synchrony | Zab

...

?

PBFT | Zyzzyva

Q/U | BAR Gossip

Aliph | UpRight

...

# The Elephant in the Cloud - Gray Failure

A component either works correctly or stops

A component appears to be still working but is in fact experiencing severe issue

A component may behave arbitrarily

**Fail-stop**                    **Gray failure**                    **Byzantine**

# The Elephant in the Cloud - Gray Failure

A component either works correctly or stops

A component appears to be still working but is in fact experiencing severe issue

A component may behave arbitrarily

**Fail-stop**

**Gray failure**

**Byzantine**

**symptom**

- subtle and ambiguous: e.g., switch random packet loss, non-fatal exceptions, memory thrashing, flaky disk I/O, overload..
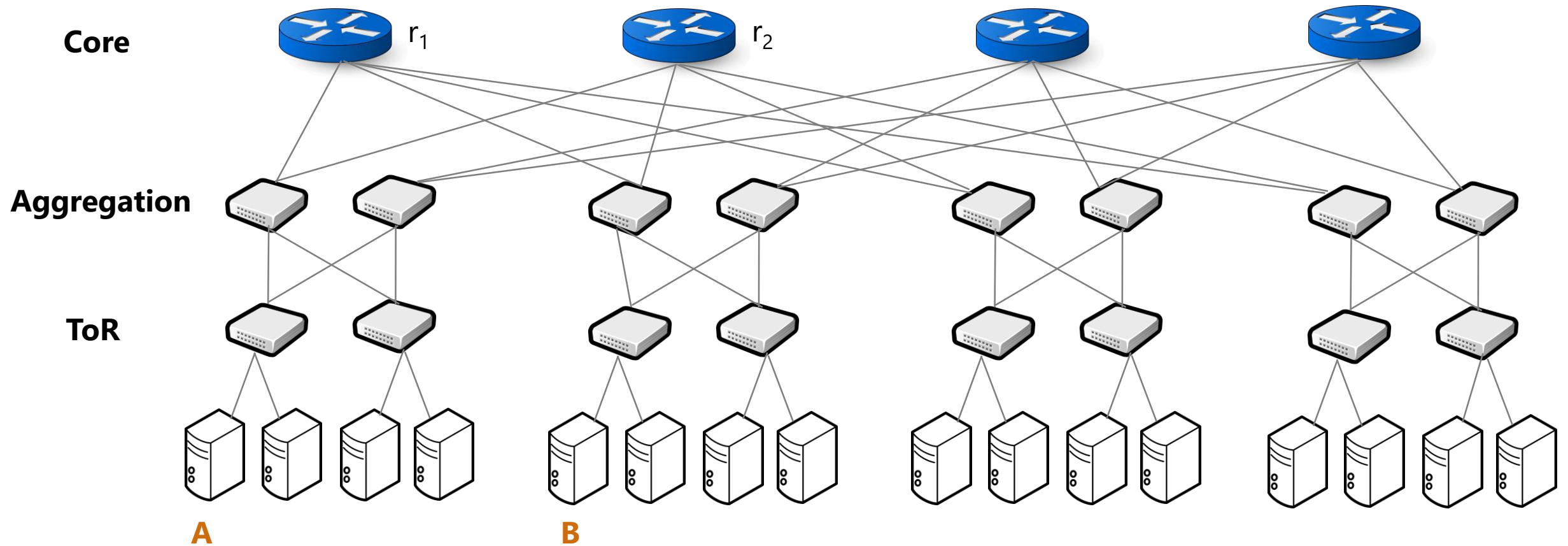
# The Elephant in the Cloud - Gray Failure

A component either works correctly or stops

A component appears to be still working but is in fact experiencing severe issue

A component may behave arbitrarily

**Fail-stop**             **Gray failure**             **Byzantine**

## symptom

- subtle and ambiguous: e.g., switch random packet loss, non-fatal exceptions, memory thrashing, flaky disk I/O, overload..

## occurrence

- across s/w and h/w stack in the infra due to various defects
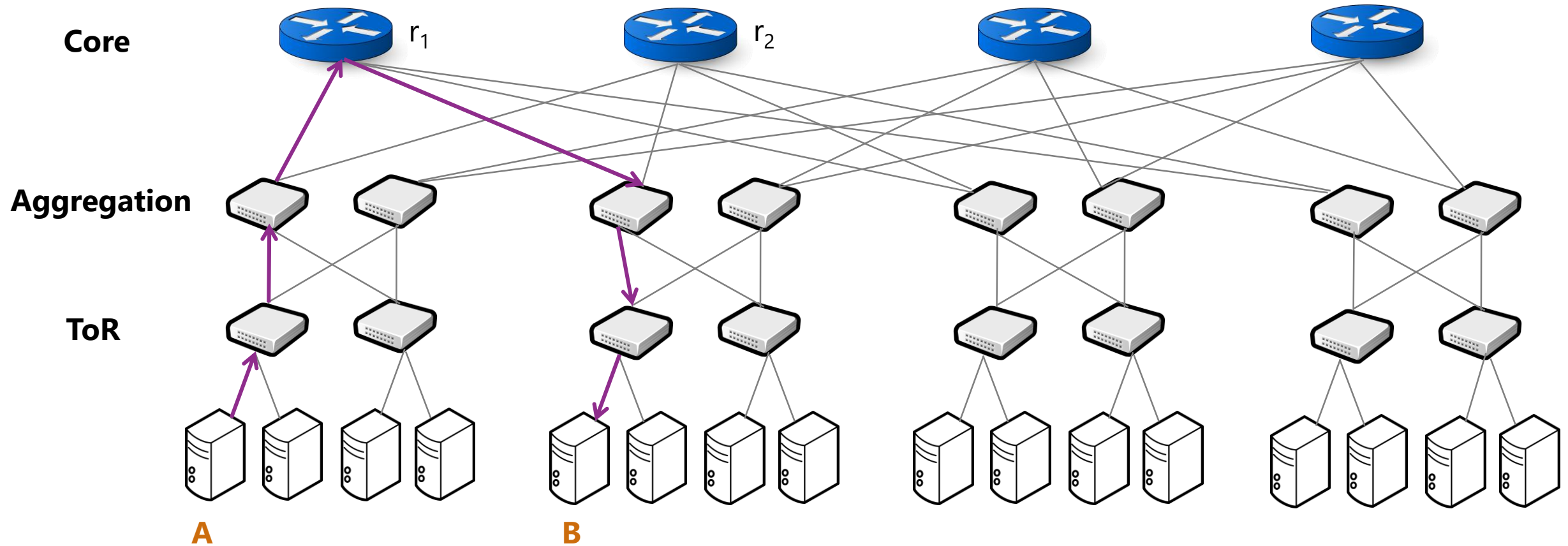- behind most service incidents we've seen in Azure

# The Elephant in the Cloud - Gray Failure

A component either works correctly or stops

A component appears to be still working but is in fact experiencing severe issue

A component may behave arbitrarily

**Fail-stop** | Gray failure | **Byzantine**

**symptom**
- subtle and ambiguous: e.g., switch random packet loss, non-fatal exceptions, memory thrashing, flaky disk I/O, overload..

**occurrence**
- across s/w and h/w stack in the infra due to various defects
- behind most service incidents we've seen in Azure

**danger**
- fault-tolerance ineffective or counterproductive
- faults take engineers & designers huge efforts to nail down
- teams play the blame game with each other

# Real-world Gray Failure Cases in Azure

# Case I: Redundancy in Datacenter Network (1)



Core

Aggregation

ToR

$r_1$     $r_2$

A     B

# Case I: Redundancy in Datacenter Network (1)



**Core**

**Aggregation**

**ToR**

A    B

# Case I: Redundancy in Datacenter Network (1)



crash

Core    $r_1$    $r_2$

Aggregation

ToR

A    B

# Case I: Redundancy in Datacenter Network (1)



Core

Aggregation

ToR

crash

$r_1$

$r_2$

A

B

46

# Case I: Redundancy in Datacenter Network (1)

increasing # of core switches helps with availability



crash

Core    $r_1$    $r_2$

Aggregation

ToR

A    B

# Case I: Redundancy in Datacenter Network (2)



**Core**

random packet drop

$r_1$     $r_2$

**Aggregation**

**ToR**

**A**     **B**

**Workload: single round trip**

# Case I: Redundancy in Datacenter Network (2)

- packets will not be re-routed ➡ application glitches or increased latency
- increasing # of core switches may not affect chance of being affected $p$



**Core**

random packet drop

$r_1$     $r_2$

**Aggregation**

**ToR**

A     B

**Workload: single round trip**

# Case I: Redundancy in Datacenter Network (3)

- high chance to involve every core switches for each front-end request
- gray failure at *any* core switch will cause delay
- more core switches ➡ worse tail latencies

$$1 - (1 - p)^n$$



**Core**

random packet drop

$r_1$     $r_2$     $r_3$     $r_4$

**Aggregation**

**ToR**

A     B     C     D     E

**Workload: send multiple requests wait for all to finish (e.g., search)**

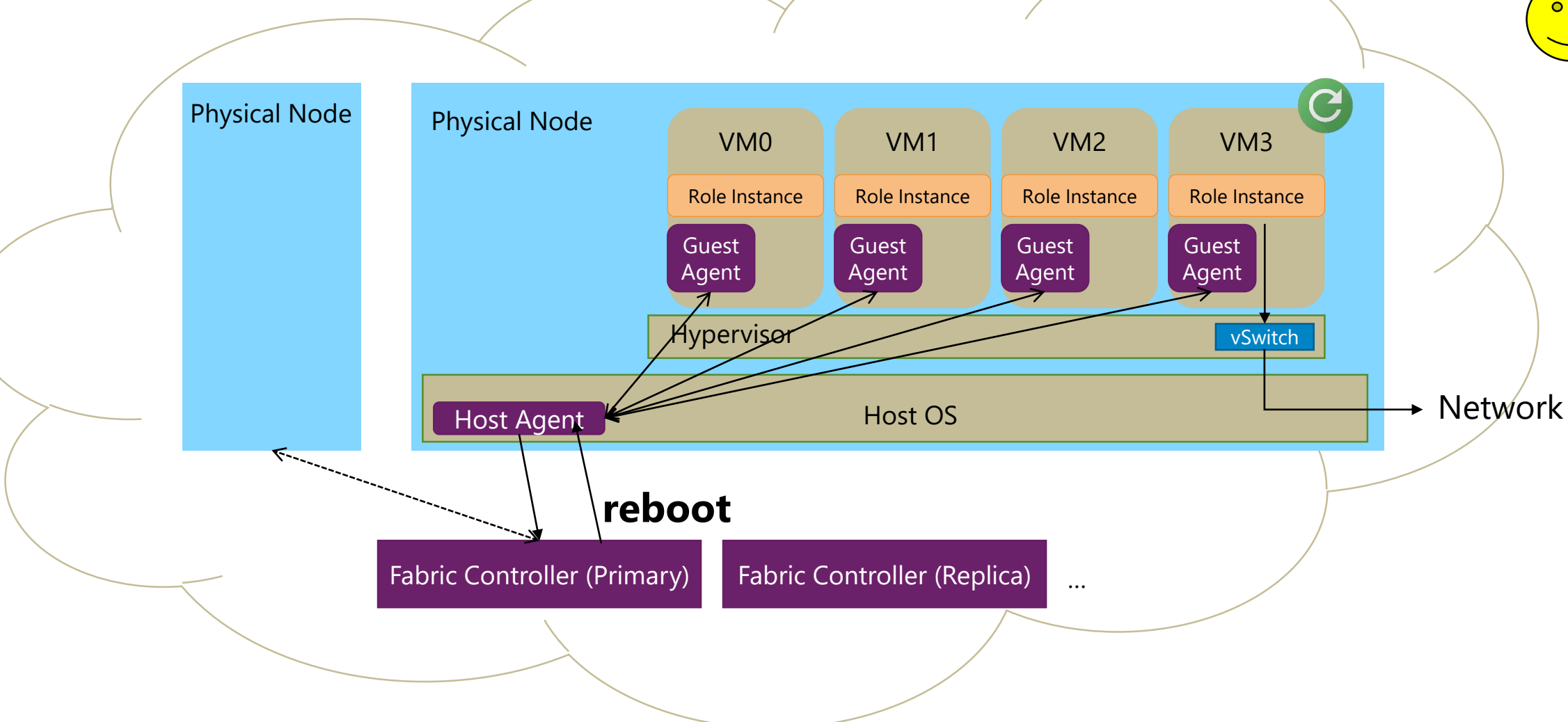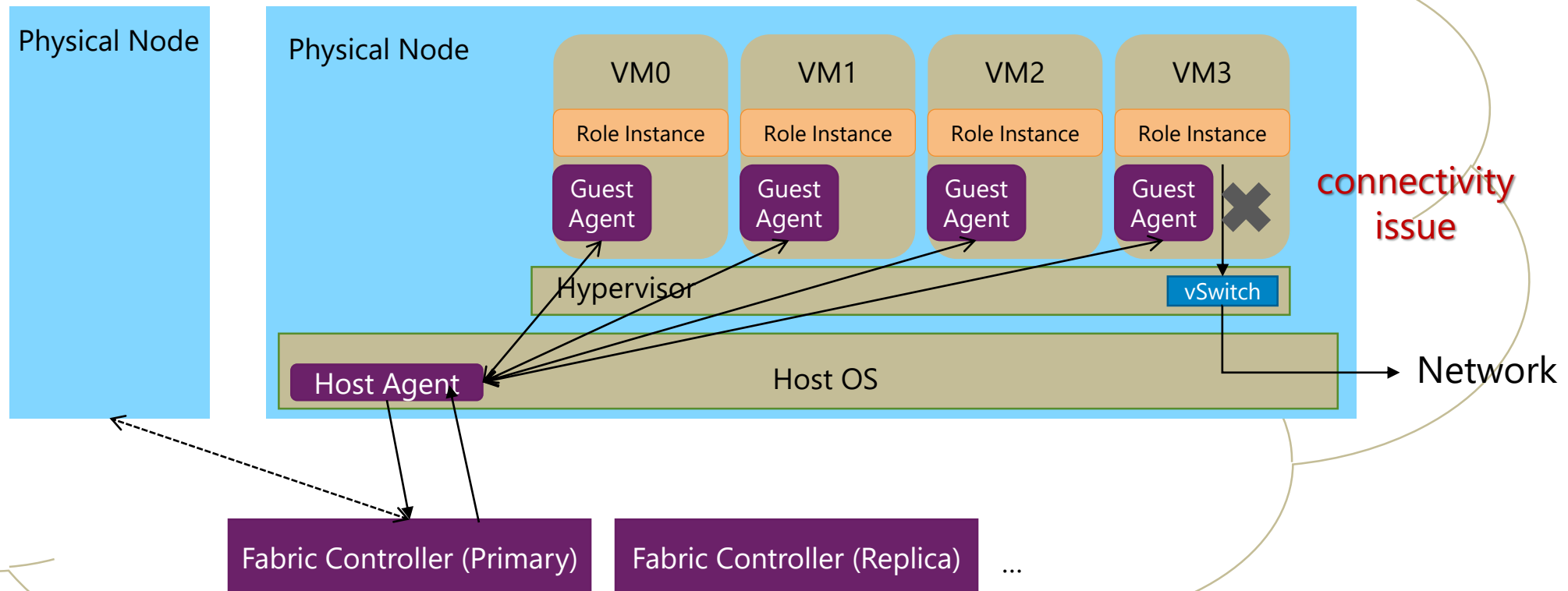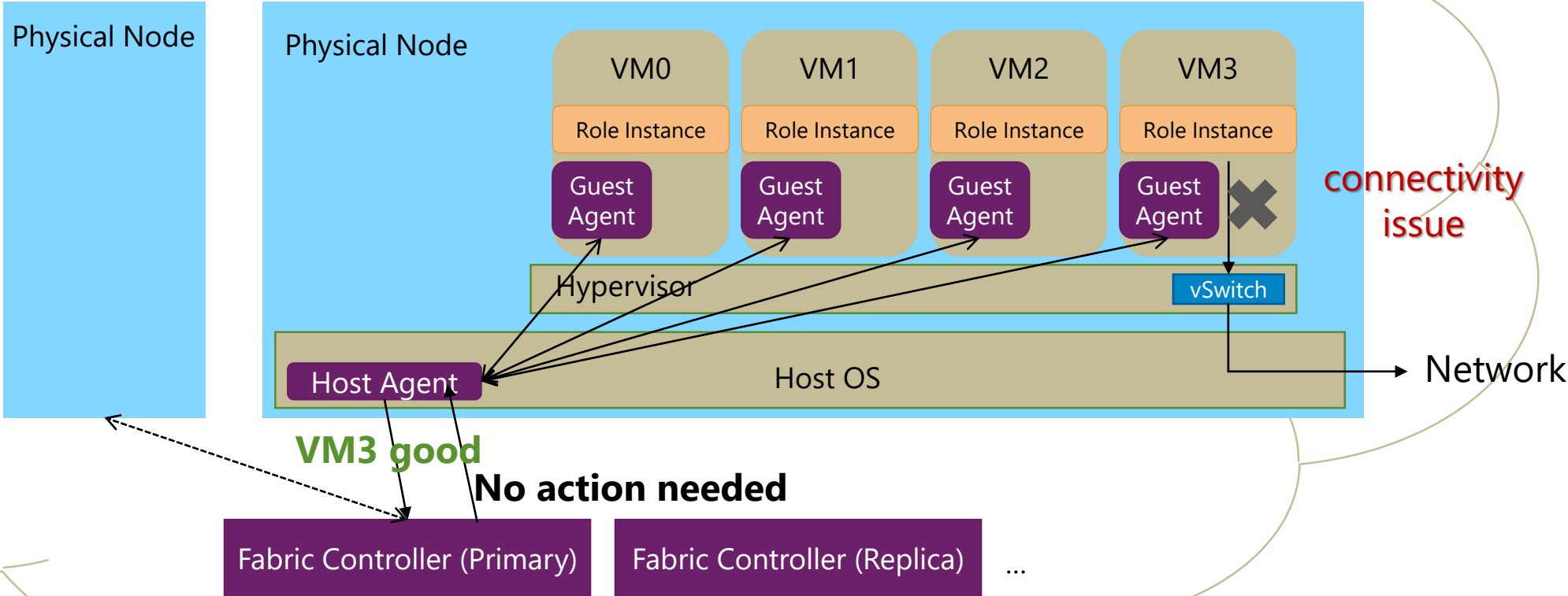# Case II: Failure Detector in Compute Service

# Case II: Failure Detector in Compute Service



Hierarchical agents to catch failure in different layers

# Case II: Failure Detector in Compute Service



Hierarchical agents to catch failure in different layers

# Case II: Failure Detector in Compute Service



Hierarchical agents to catch failure in different layers

# Case II: Failure Detector in Compute Service
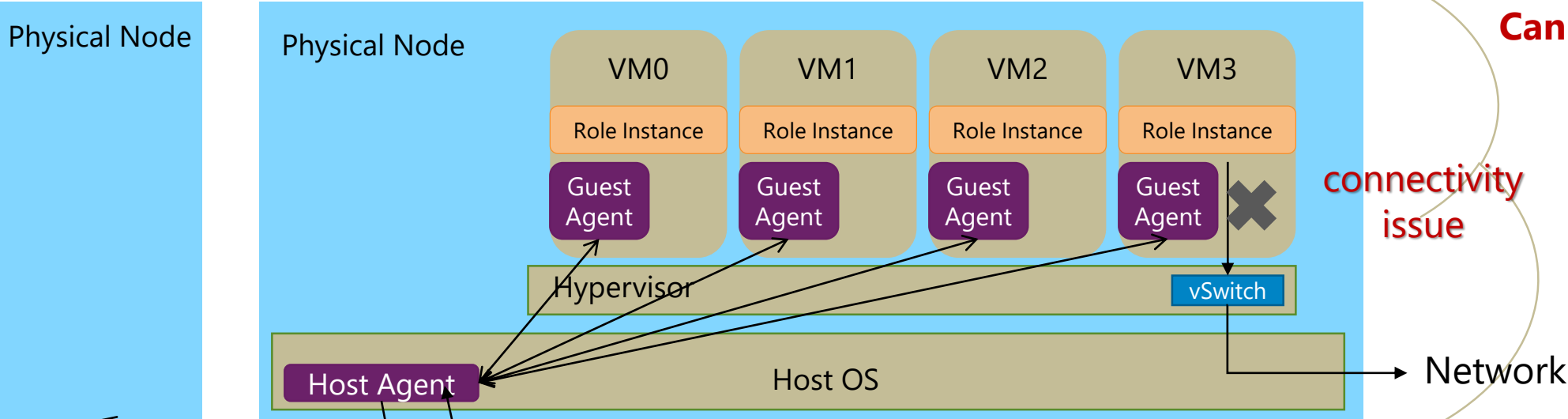


Hierarchical agents to catch failure in different layers

Physical Node

Physical Node

VM0 | VM1 | VM2 | VM3

Role Instance | Role Instance | Role Instance | Role Instance

Guest Agent | Guest Agent | Guest Agent | Guest Agent

connectivity issue

Hypervisor

vSwitch

Host OS

Host Agent

Network

Fabric Controller (Primary)

Fabric Controller (Replica) ...

# Case II: Failure Detector in Compute Service

# Case II: Failure Detector in Compute Service

# Case III: Recovery in Storage Service

# Case III: Recovery in Storage Service



**Front End**　　**Front End**　　**Front End**

**Stream Manager**

**low free blocks**

EN₁　EN₂　EN₃　EN₄　EN₅　...

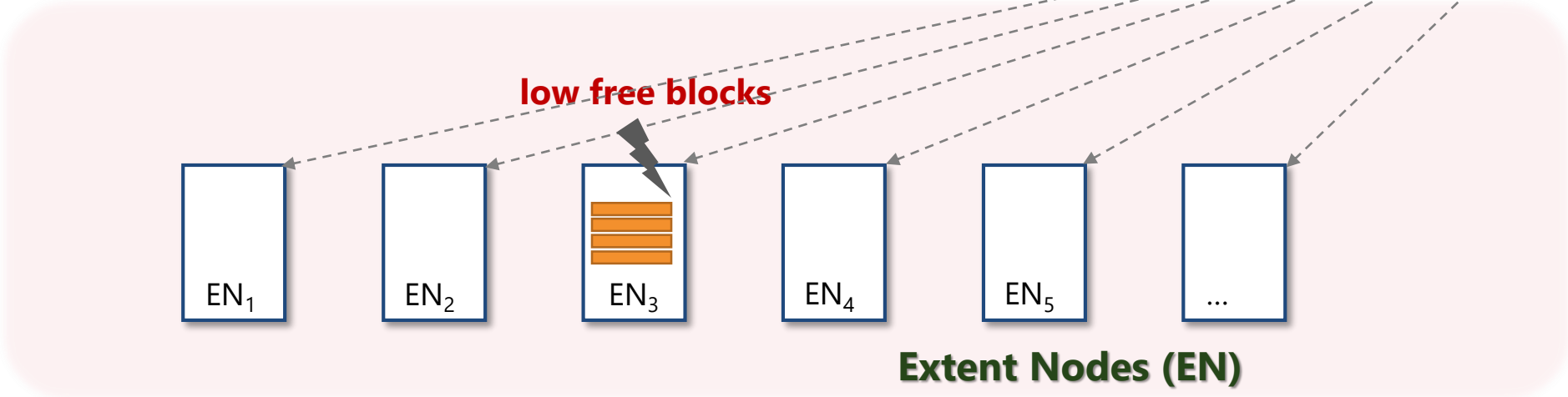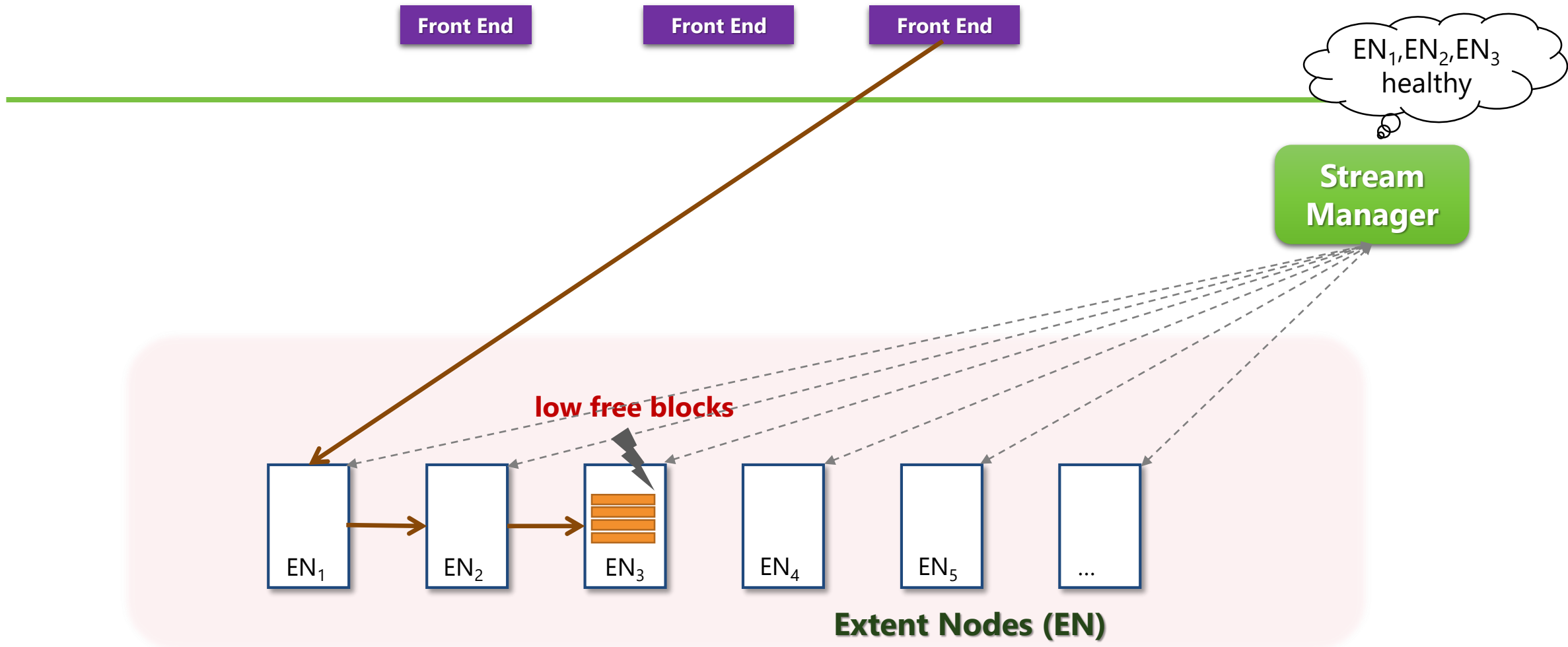**Extent Nodes (EN)**

# Case III: Recovery in Storage Service

# Case III: Recovery in Storage Service



Front End    Front End    Front End

EN$_1$,EN$_2$,EN$_3$ healthy

Stream Manager

low free blocks

EN$_1$    EN$_2$    EN$_3$    EN$_4$    EN$_5$    ...

**Extent Nodes (EN)**

# Case III: Recovery in Storage Service

# Case III: Recovery in Storage Service



Front End   Front End   Front End

EN$_3$ is down

**Stream Manager**

**crash**

EN$_1$   EN$_2$   EN$_3$   EN$_4$   EN$_5$   ...

**Extent Nodes (EN)**

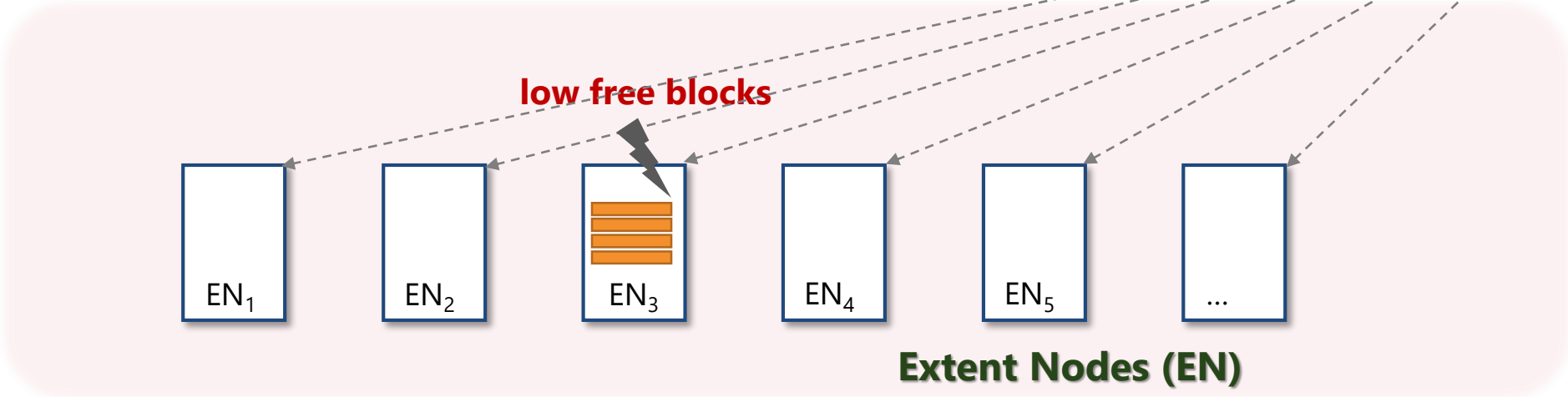# Case III: Recovery in Storage Service

# Case III: Recovery in Storage Service



Front End   Front End   Front End

Stream Manager

low free blocks

$EN_1$   $EN_2$   $EN_3$   $EN_4$   $EN_5$   ...

**Extent Nodes (EN)**

# Case III: Recovery in Storage Service

# Case III: Recovery in Storage Service

# Case III: Recovery in Storage Service

# Case III: Recovery in Storage Service

# Case III: Recovery in Storage Service

# Case III: Recovery in Storage Service

# Case III: Recovery in Storage Service

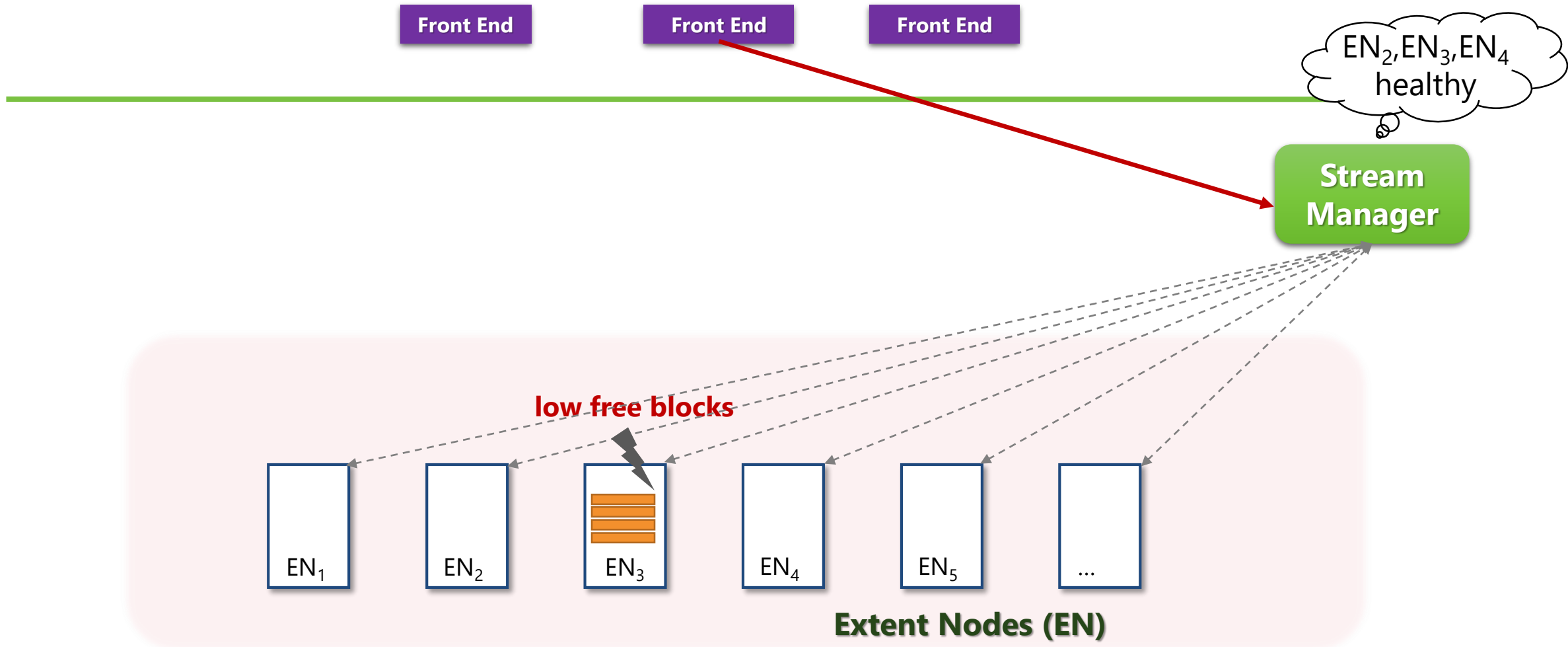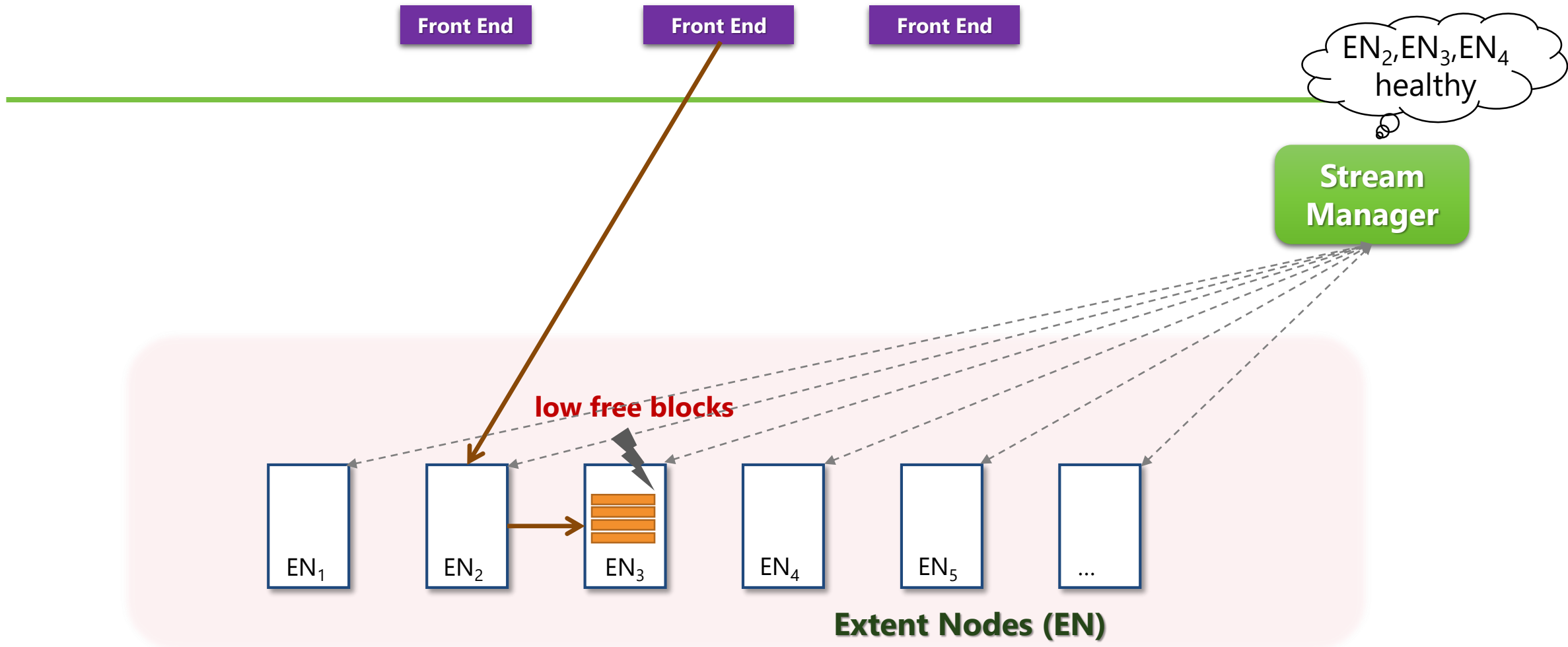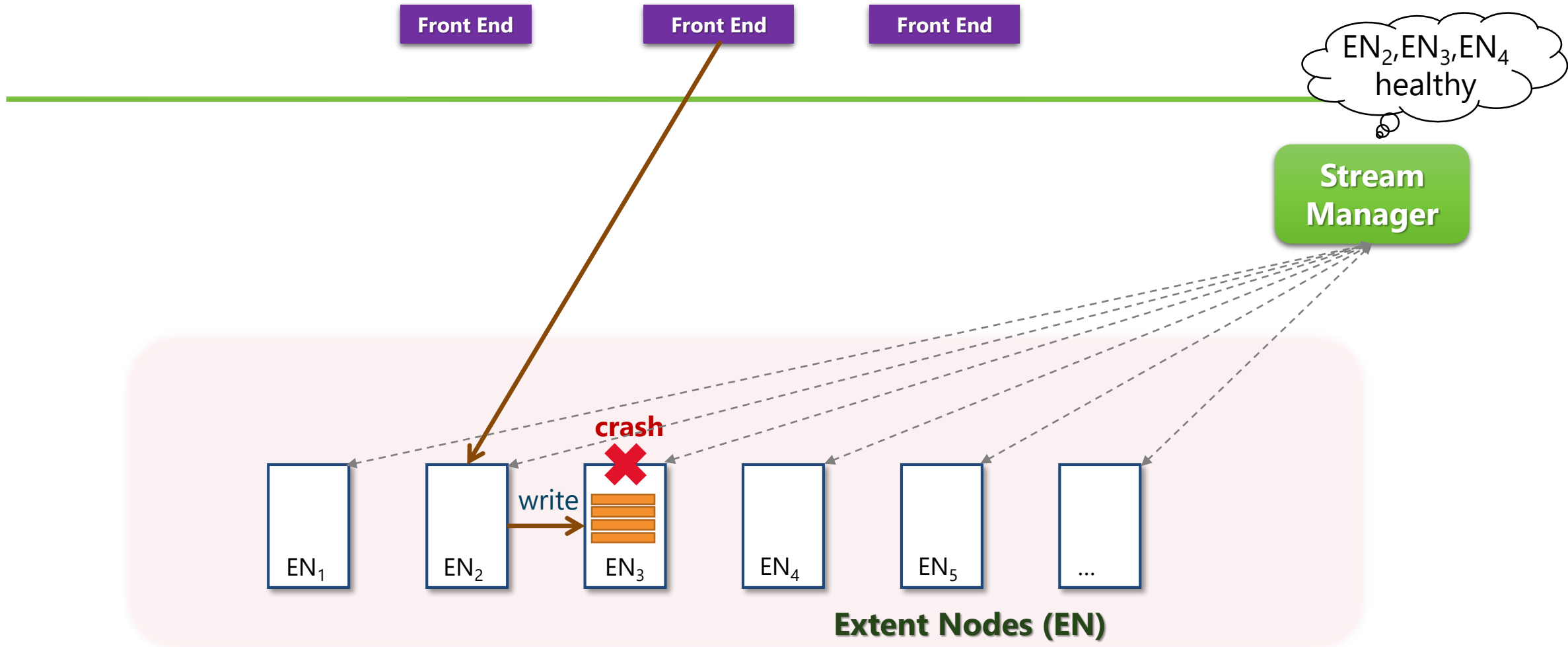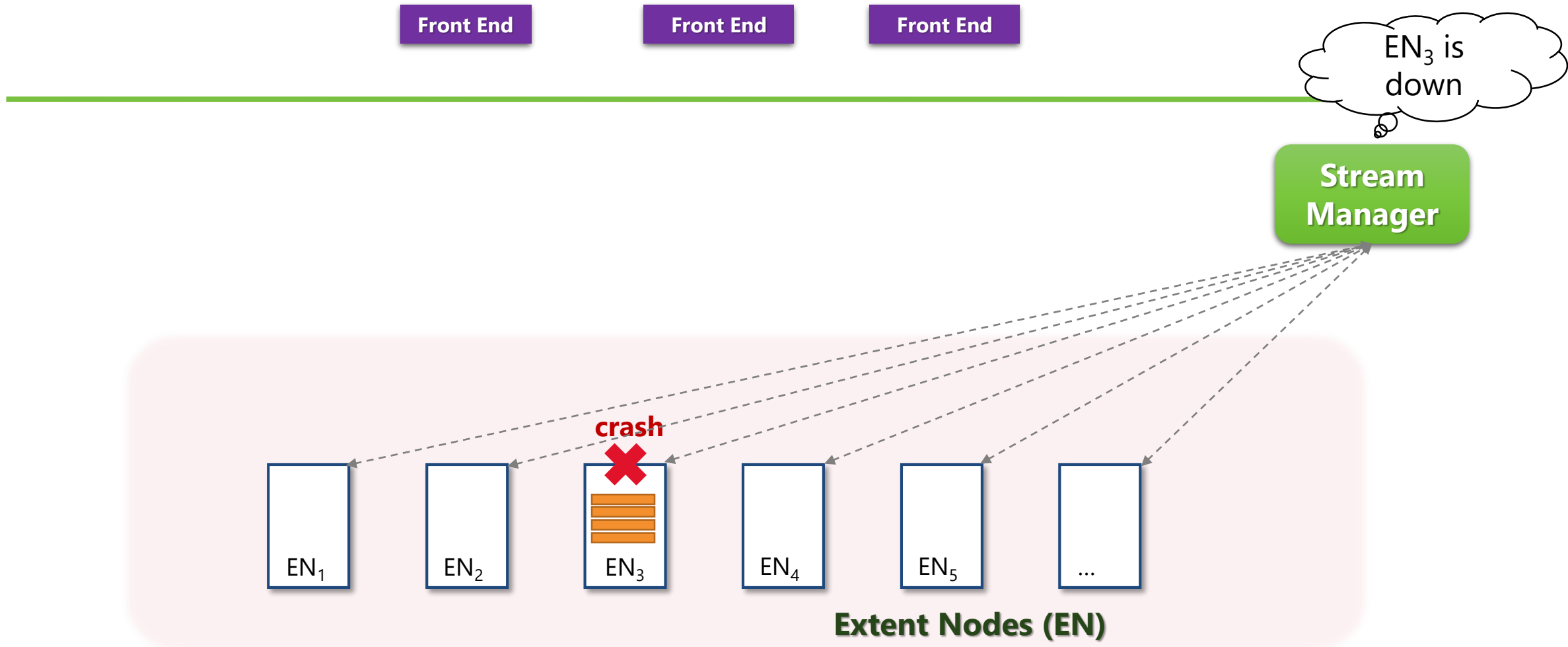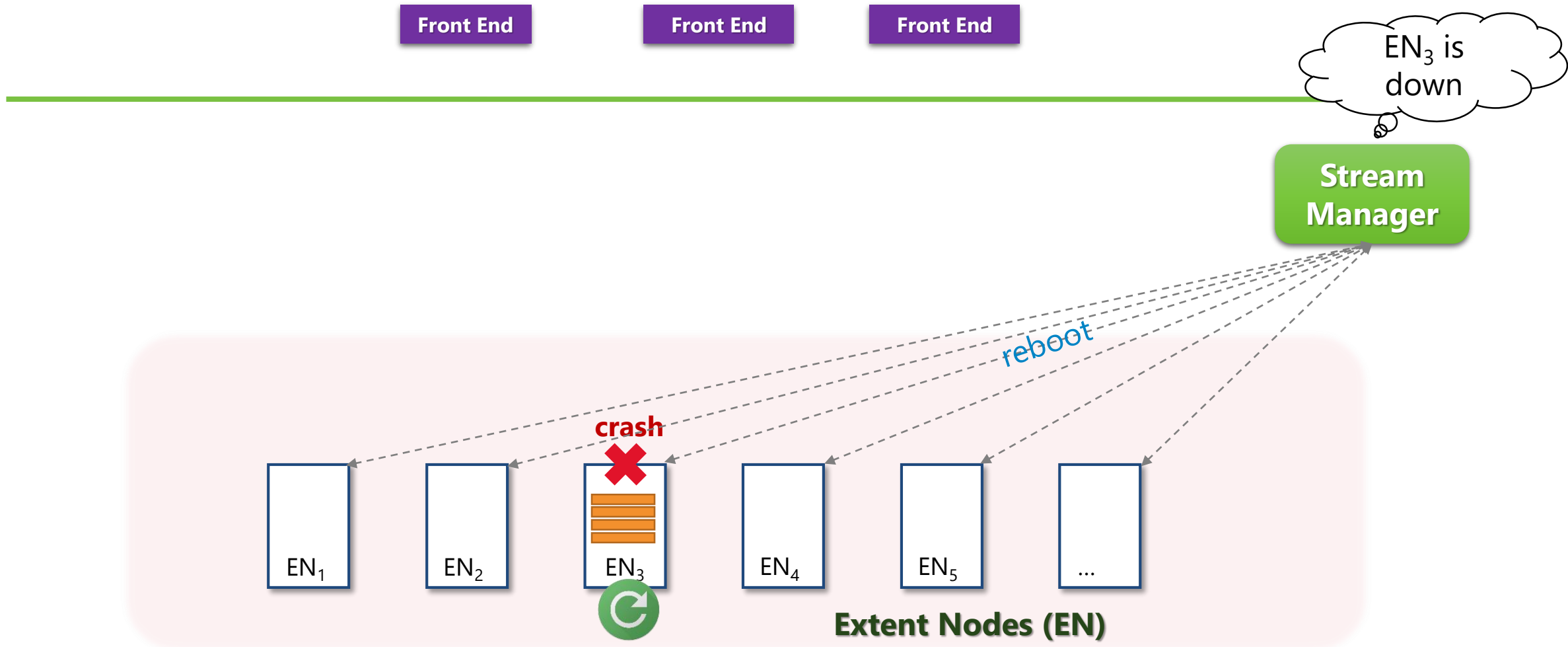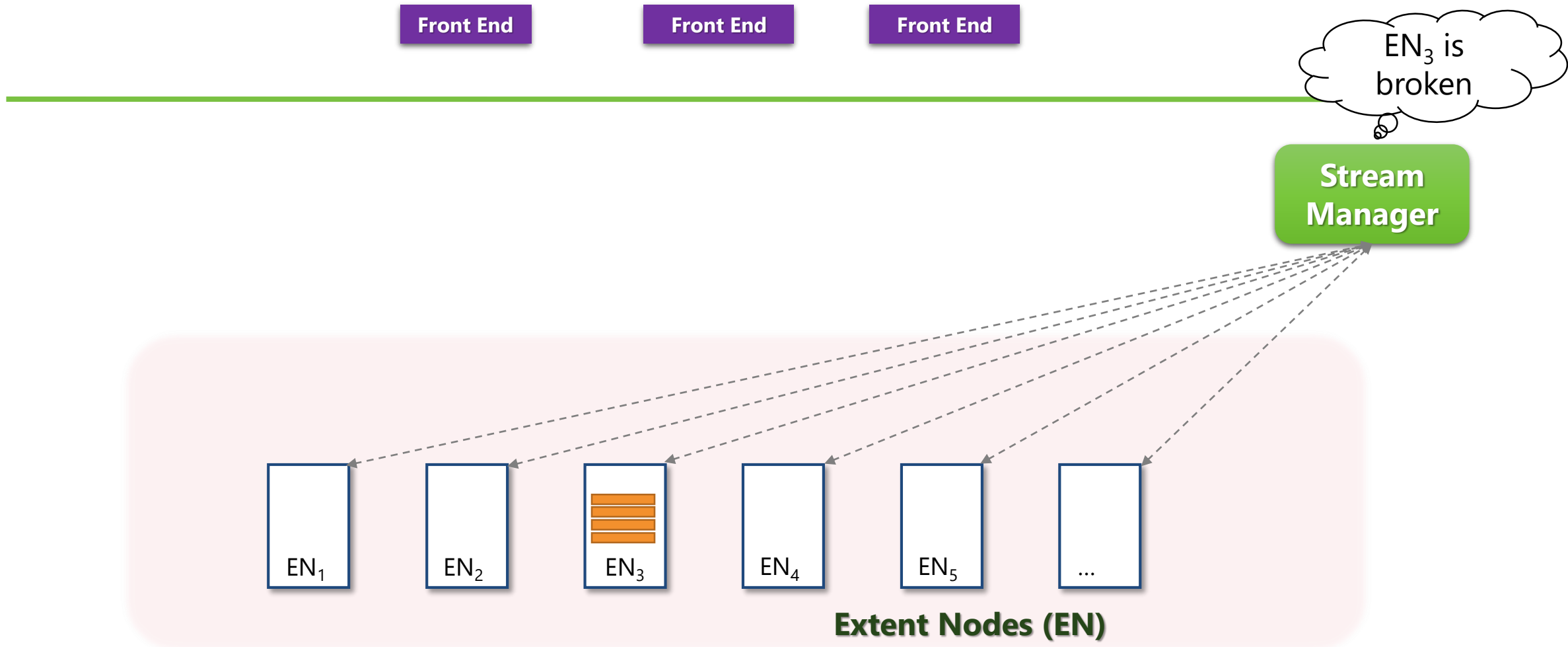# Case III: Recovery in Storage Service

# Case III: Recovery in Storage Service

# Understanding Gray Failure

# The Many Faces of Gray Failure

So, what is a gray failure?

# The Many Faces of Gray Failure

So, what is a gray failure?

" A performance issue. "

# The Many Faces of Gray Failure

So, what is a gray failure?

" A performance issue. "

" A problem that some thinks is a failure but some thinks is not, e.g., a 2% packet loss. The ambiguity itself defines gray failure. If everyone agrees it is a problem, it is not a gray failure. "

# The Many Faces of Gray Failure

So, what is a gray failure?

" A performance issue. "

" A problem that some thinks is a failure but some thinks is not, e.g., a 2% packet loss.  The ambiguity itself defines gray failure. If everyone agrees it is a problem, it is not a gray failure. "

" A Heisenbug, sometimes it occurs and sometimes it does not. "

# The Many Faces of Gray Failure

**So, what is a gray failure?**

" A performance issue. "

" A problem that some thinks is a failure but some thinks is not, e.g., a 2% packet loss.  The ambiguity itself defines gray failure. If everyone agrees it is a problem, it is not a gray failure. "

" A Heisenbug, sometimes it occurs and sometimes it does not. "

" The system is failing slowly, e.g., memory leak. "

# The Many Faces of Gray Failure

So, what is a gray failure?

66 A performance issue. 99

66 A problem that some thinks is a failure but some thinks is not, e.g., a 2% packet loss. The ambiguity itself defines gray failure. If everyone agrees it is a problem, it is not a gray failure. 99

66 A Heisenbug, sometimes it occurs and sometimes it does not. 99

66 The system is failing slowly, e.g., memory leak. 99

66 There is an increasing number of transient errors in the system, which results in reduced system capacity even if the system still manages to continue working. 99

# The Many Faces of Gray Failure

So, what is a gray failure?

" A performance issue. "

" A pro...                    ...ks is not,
e.g., a...                     ...ay failure.
If eve...                      ...re. "

" A He...gr...                 ...oes not. "

" The system is failing slowly, e.g., memory leak. "

" There is an increasing number of transient errors in the system,
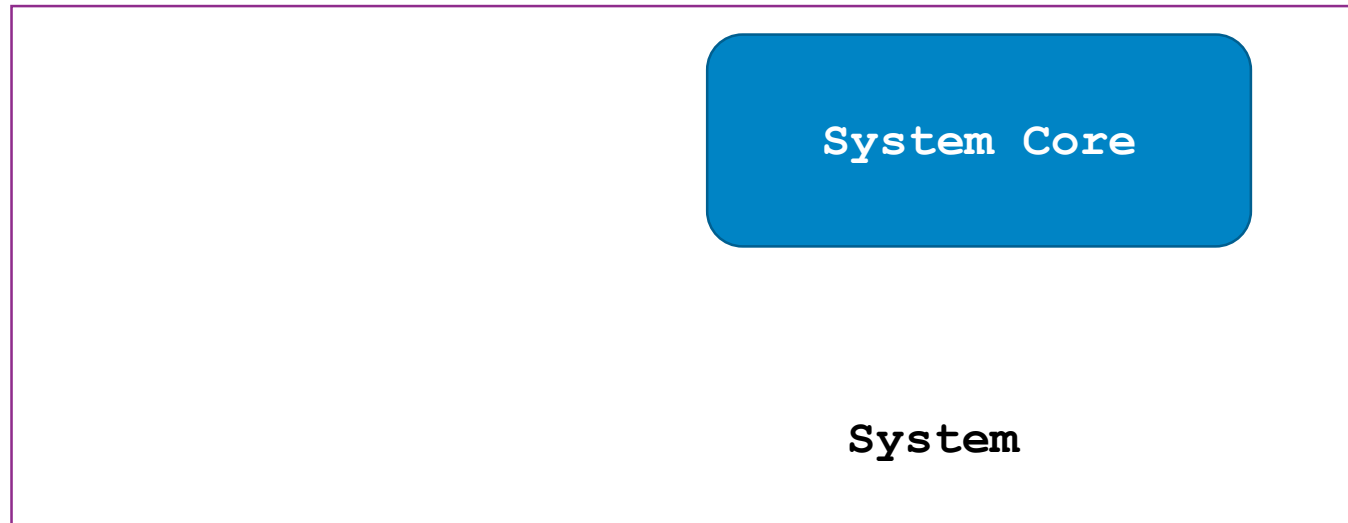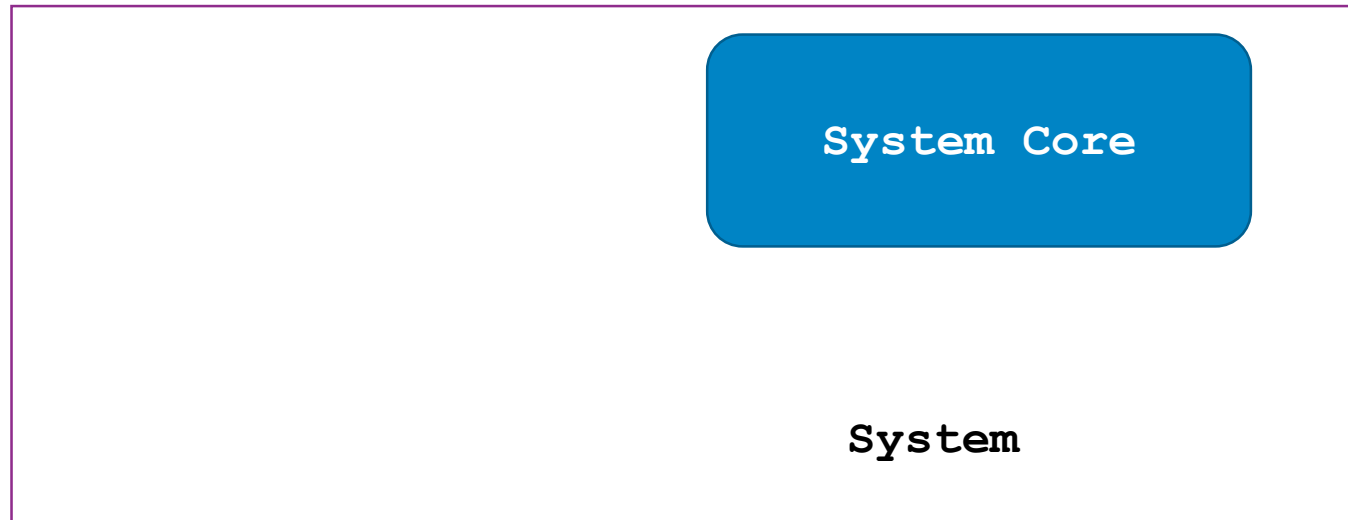which results in reduced system capacity even if the system still
manages to continue working. "

What is a formal way to define and
study gray failure, one that potentially
sheds light on how to address it?

# An Abstract Model

```
                    ┌─────────────────────┐
                    │                     │
                    │     System Core     │
                    │                     │
                    └─────────────────────┘



                           System
```

*Note:* these are logical entities

# An Abstract Model

**System Core**

**System**

- distributed storage system
- IaaS platform
- data center network
- search engine
- ...

*Note:* these are logical entities
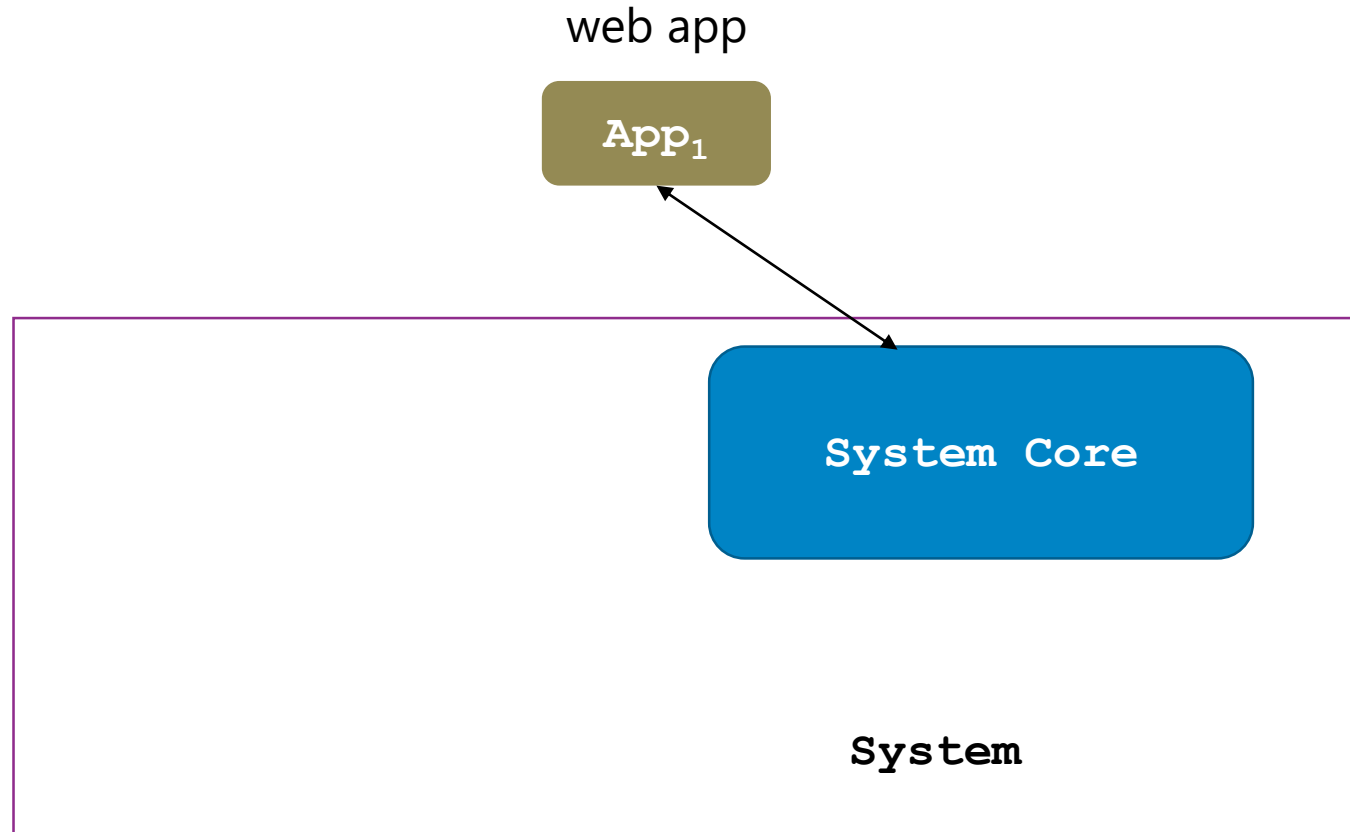
# An Abstract Model

web app

App$_1$

System Core

System

- distributed storage system
- IaaS platform
- data center network
- search engine
- ...

*Note:* these are logical entities

# An Abstract Model
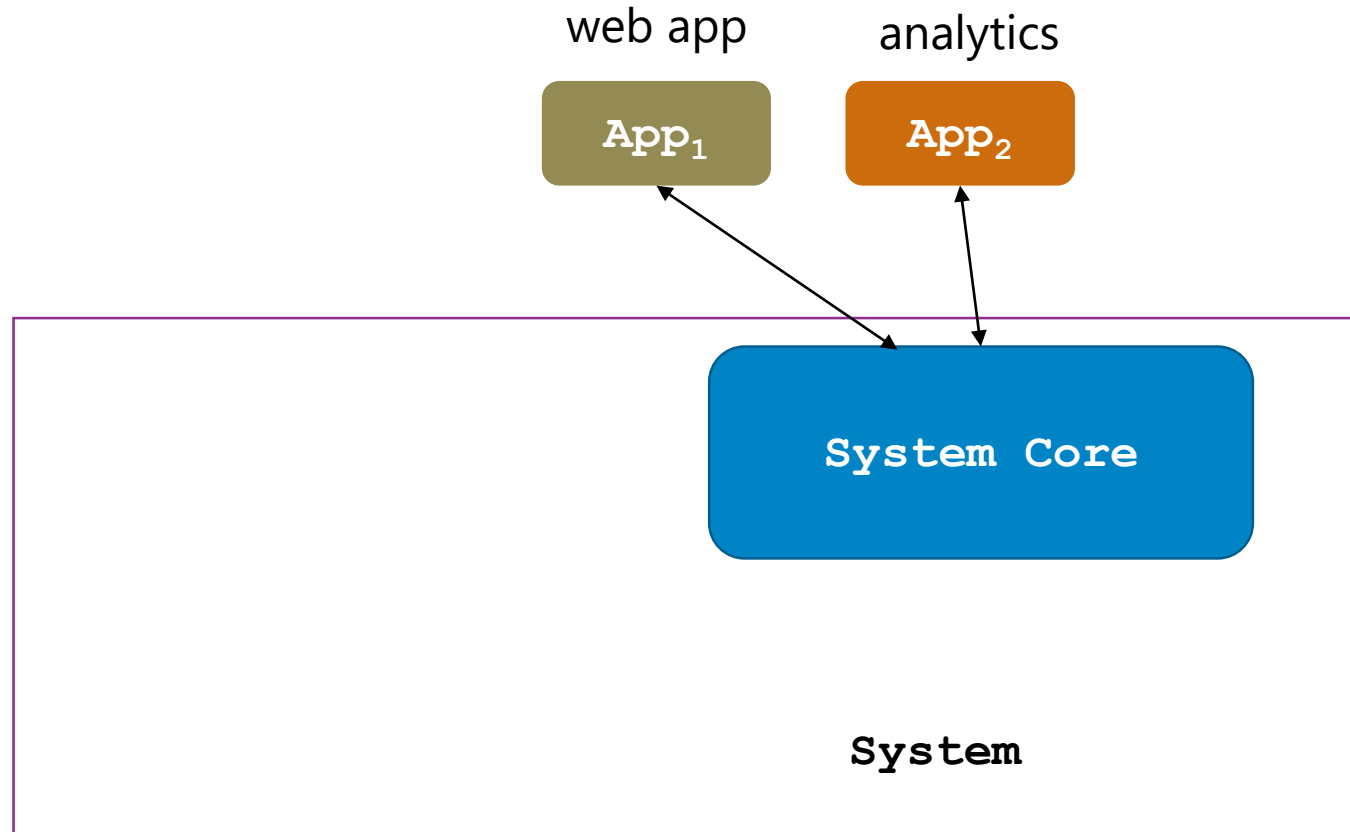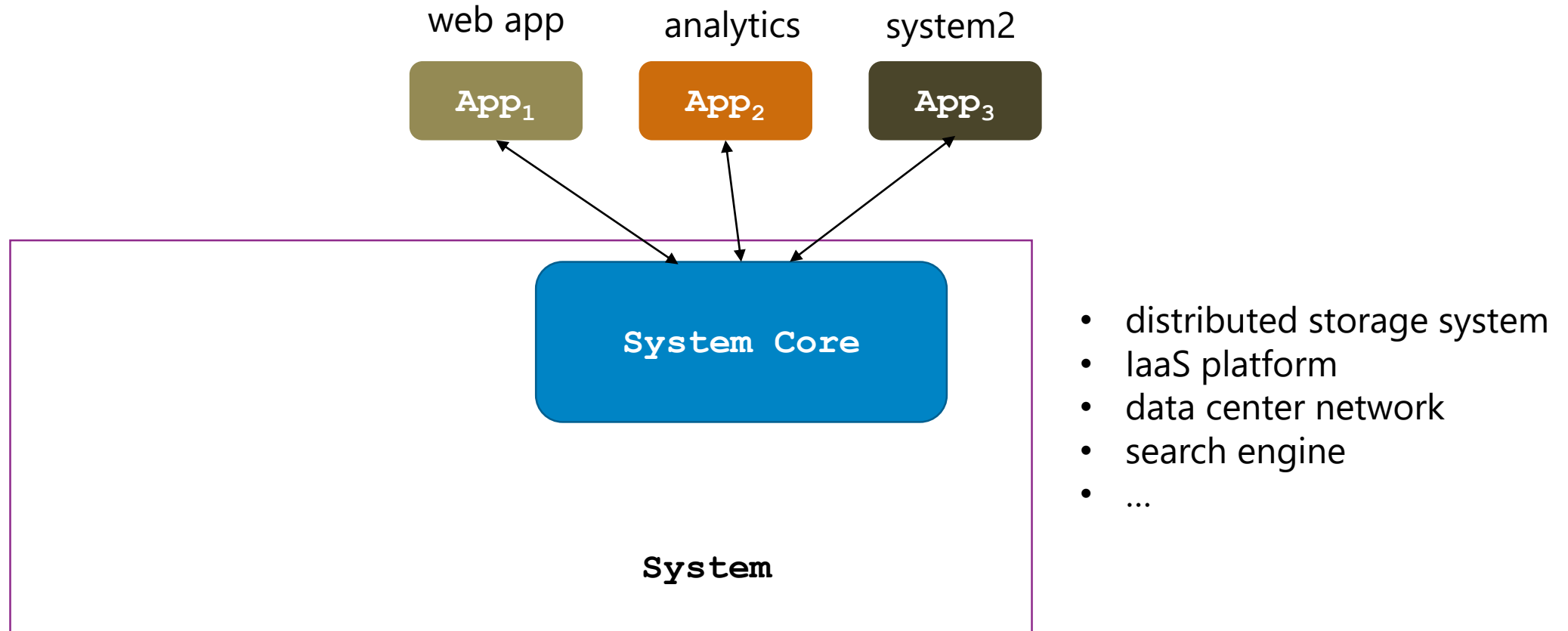
web app   analytics

App₁   App₂

System Core

- distributed storage system
- IaaS platform
- data center network
- search engine
- ...

System

*Note:* these are logical entities

# An Abstract Model

web app     analytics     system2

**App$_1$**     **App$_2$**     **App$_3$**

**System Core**

- distributed storage system
- IaaS platform
- data center network
- search engine
- ...

**System**

*Note:* these are logical entities

# An Abstract Model

web app     analytics     system2     user/operator

**App$_1$**     **App$_2$**     **App$_3$**     **App$_n$**     ...

**System Core**

- distributed storage system
- IaaS platform
- data center network
- search engine
- ...

**System**

*Note:* these are logical entities
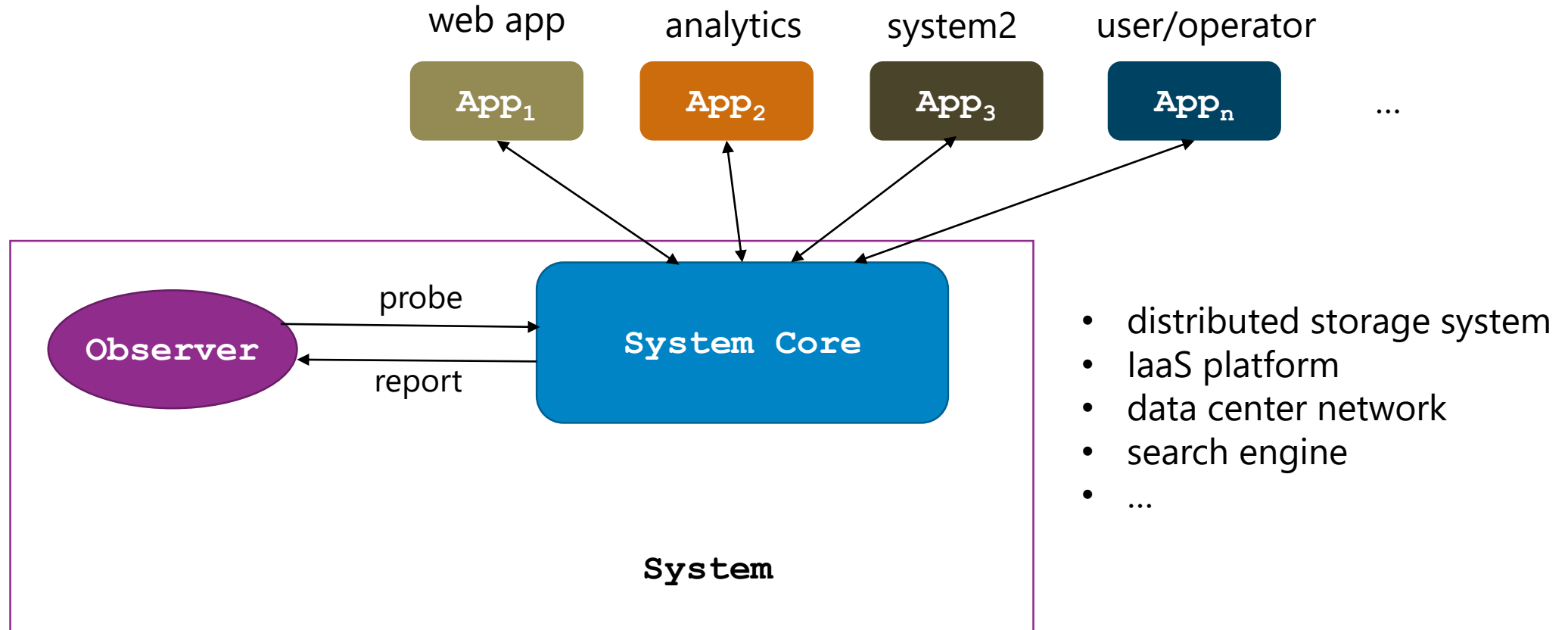
# An Abstract Model



web app     analytics     system2     user/operator

**App$_1$**     **App$_2$**     **App$_3$**     **App$_n$**     ...

**Observer**

probe

**System Core**

report

- distributed storage system
- IaaS platform
- data center network
- search engine
- ...

**System**

*Note:* these are logical entities

# An Abstract Model

web app     analytics     system2     user/operator

$App_1$     $App_2$     $App_3$     $App_n$     ...

**Observer**

probe →

← report

**System Core**

**Reactor**

**System**

- distributed storage system
- IaaS platform
- data center network
- search engine
- ...

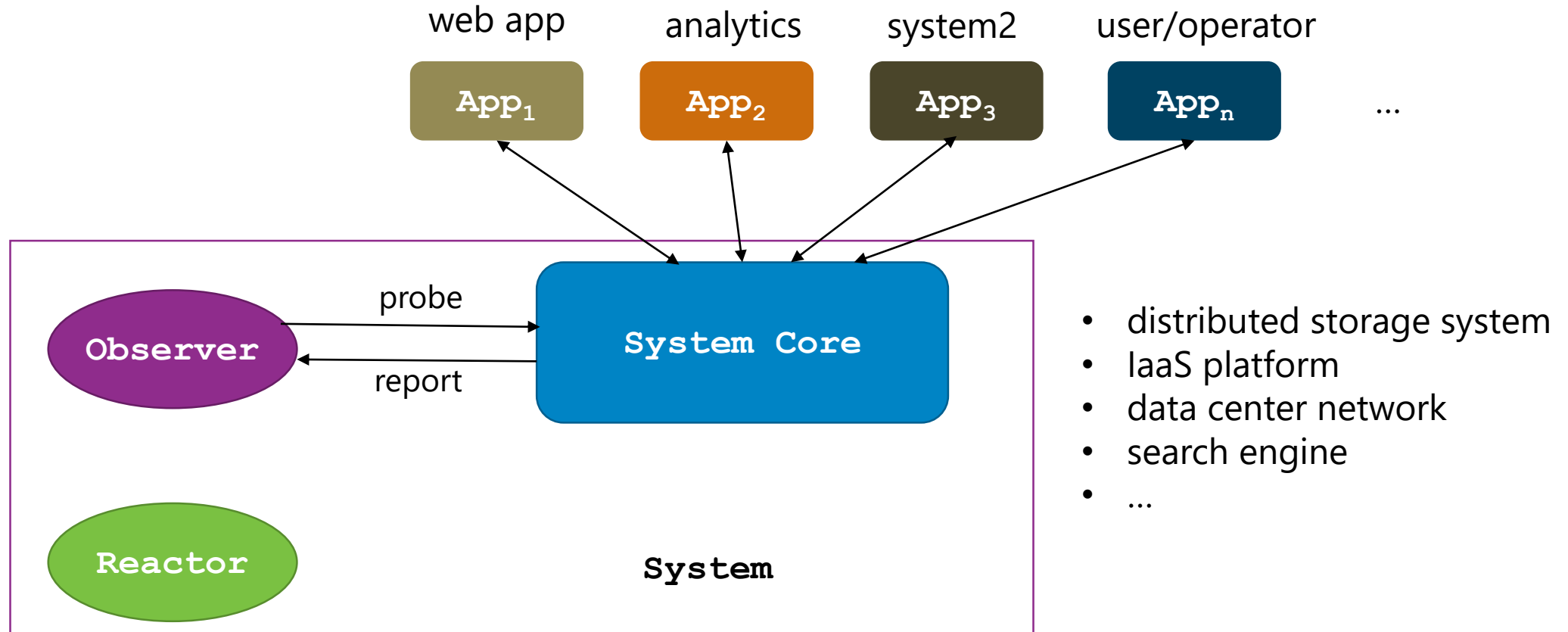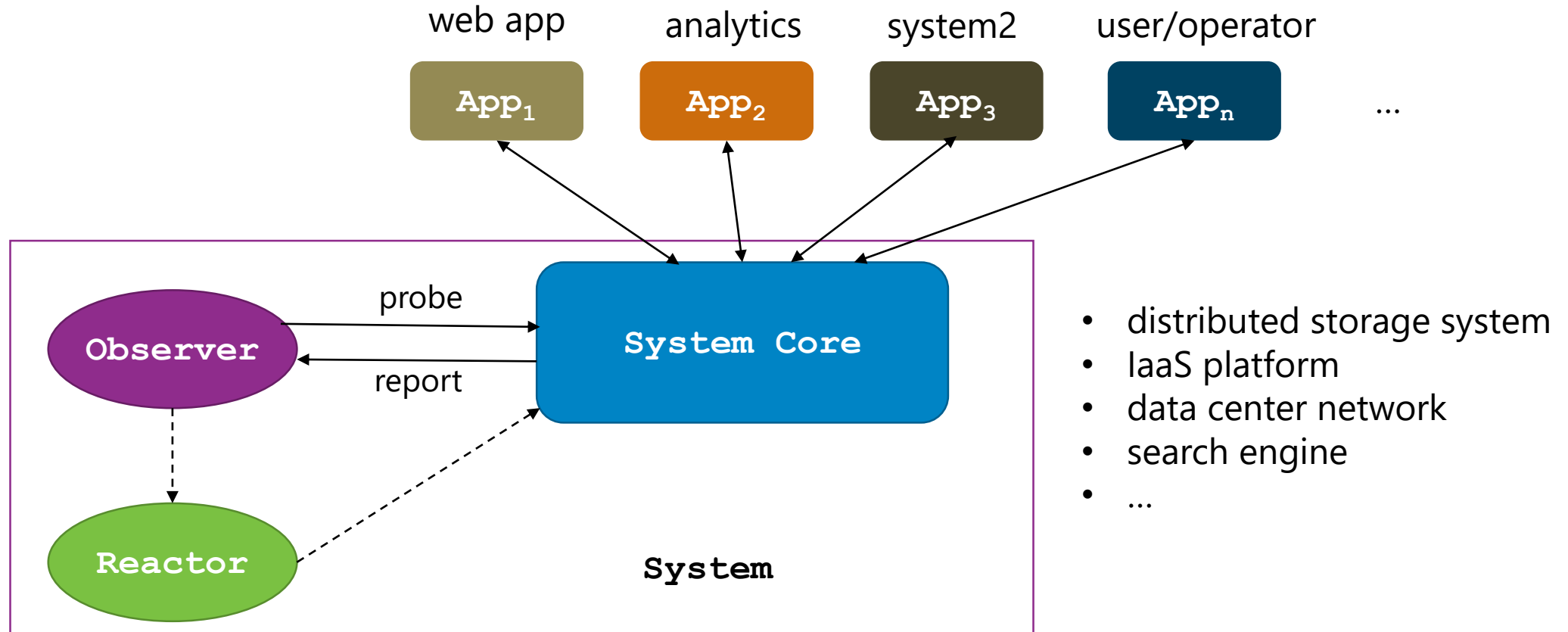*Note:* these are logical entities

# An Abstract Model
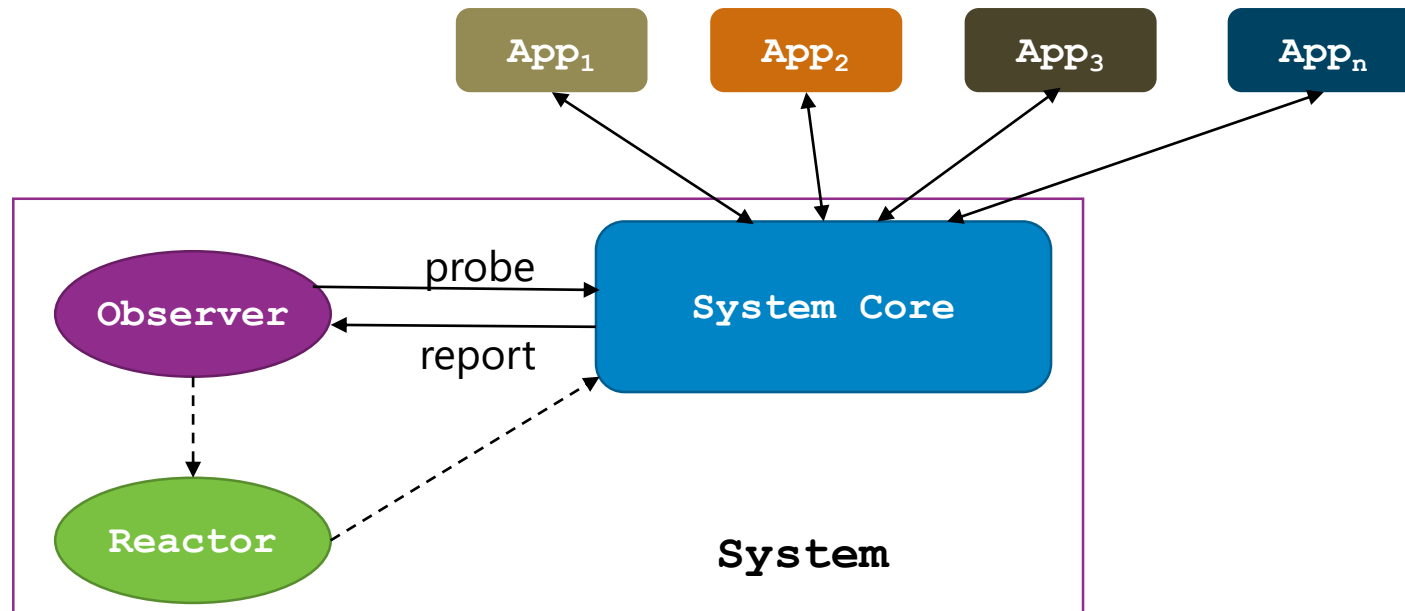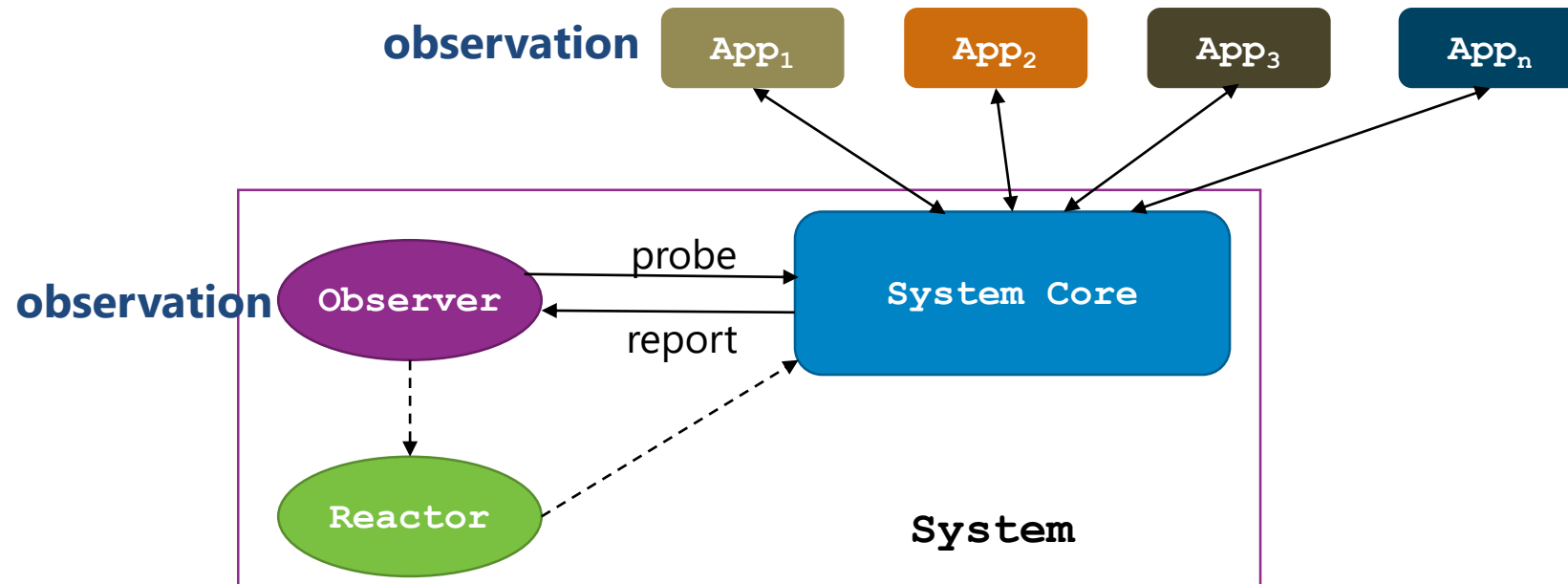
# An Abstract Model



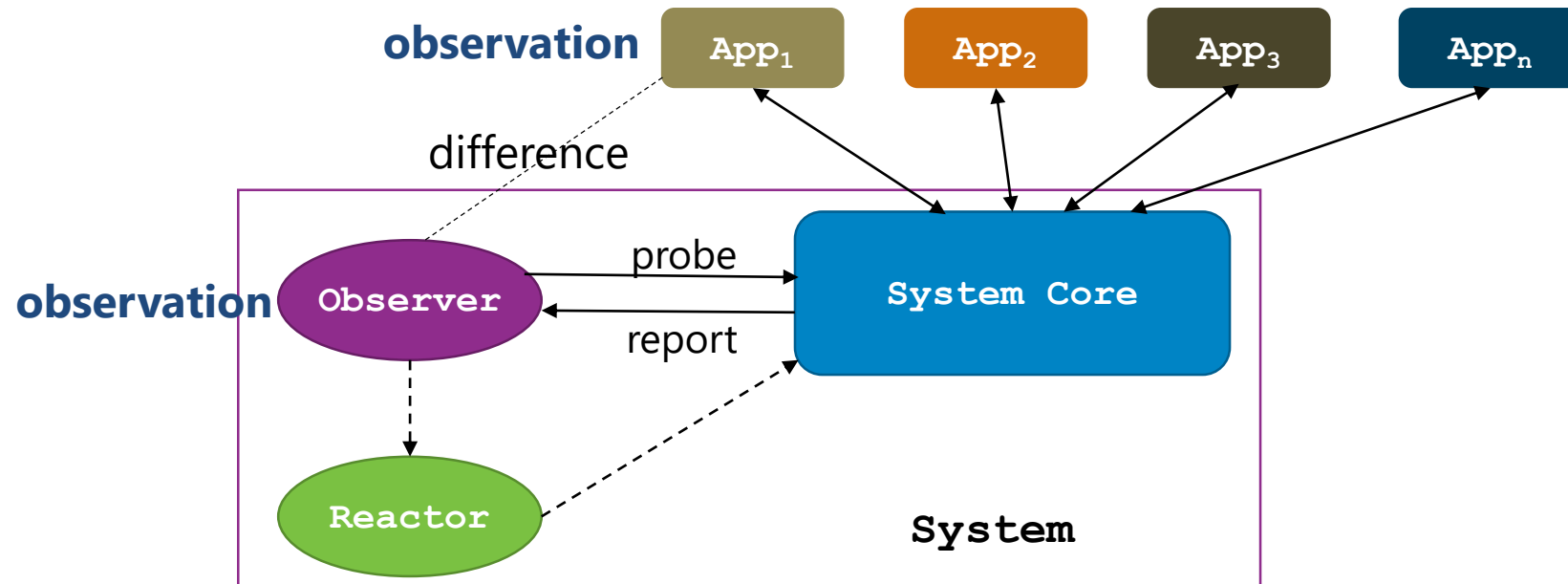*Note:* these are logical entities

# Gray Failure Trait: *Differential Observability*

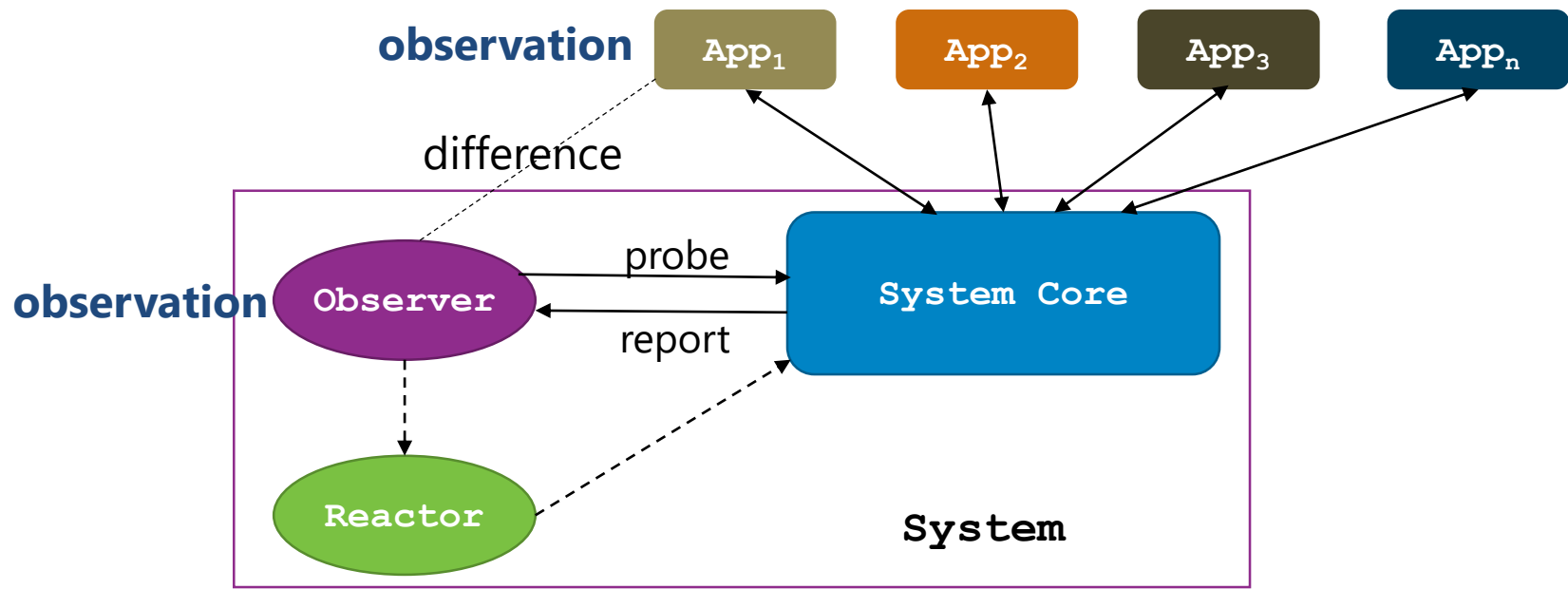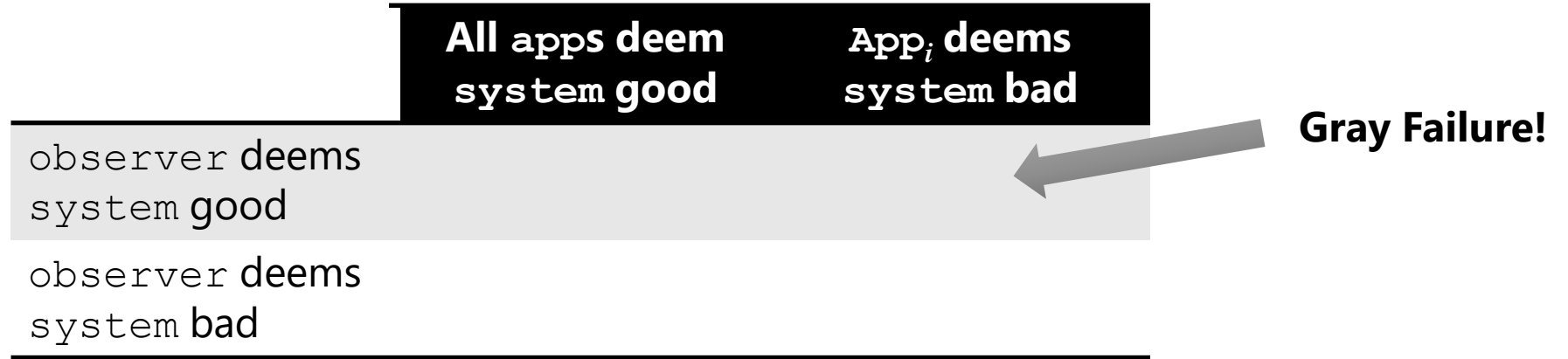# Gray Failure Trait: *Differential Observability*

# Gray Failure Trait: *Differential Observability*

different entities come into different conclusions about whether a system is working or not

# Gray Failure Trait: *Differential Observability*

| | All apps deem system **good** | App$_i$ deems system **bad** |
|---|---|---|
| observer deems system **good** | | ← **Gray Failure!** |
| observer deems system **bad** | | |

**observation** App$_1$    App$_2$    App$_3$    App$_n$

difference

**observation**    Observer → probe → System Core

report

Reactor

System

# Gray Failure Trait: *Differential Observability*



**Healthy or w/ minor latent fault**

| | All apps deem system **good** | App$_i$ deems system **bad** |
|---|:---:|:---:|
| observer deems system good | ❶ | ❷ |
| observer deems system bad | | |

**Gray Failure!**

observation

App$_1$  App$_2$  App$_3$  App$_n$

difference

observation

Observer — probe / report — System Core

Reactor

System

# Gray Failure Trait: *Differential Observability*

**Healthy or w/ minor latent fault**

**Gray Failure!**

| | All apps deem system good | App$_i$ deems system bad |
|---|---|---|
| observer deems system good | ❶ | ❷ |
| observer deems system bad | ❸ | |

**Fault tolerance at play, or a false positive**

**observation** App$_1$  App$_2$  App$_3$  App$_n$

difference

**observation**

Observer — probe → System Core

Observer ← report —

Reactor

System

# Gray Failure Trait: *Differential Observability*

**Healthy or w/ minor latent fault**

**Gray Failure!**

**Fault tolerance at play, or a false positive**

**Crash, fail-stop**

|  | All apps deem system **good** | App$_i$ deems system **bad** |
|---|---|---|
| observer deems system **good** | ❶ | ❷ |
| observer deems system **bad** | ❸ | ❹ |

**observation**  App$_1$   App$_2$   App$_3$   App$_n$

difference

**observation**  Observer  —probe→  System Core  ←report—

Reactor

System

# Applying the Model

## Case I

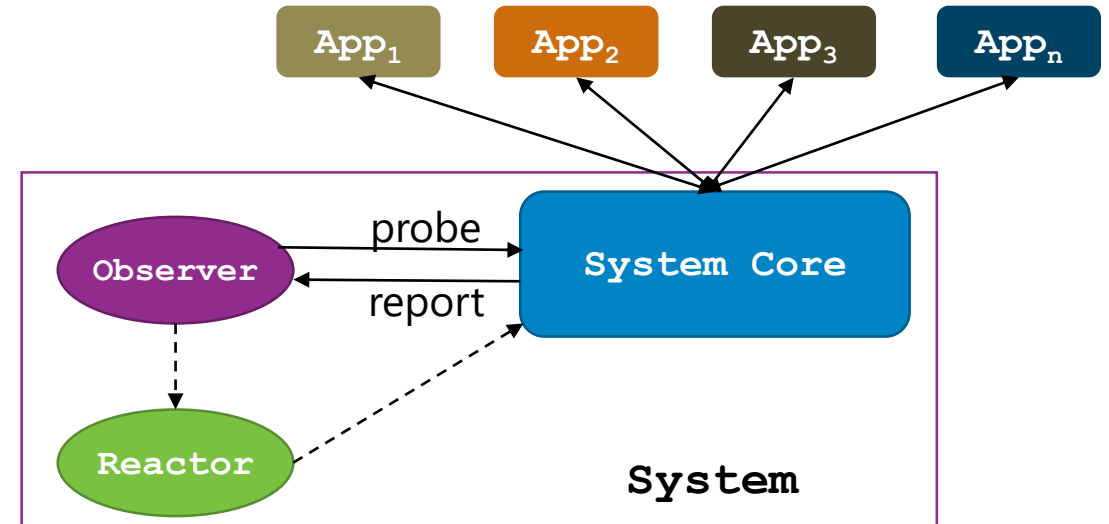| term | example |
| --- | --- |
| system | data center network |
| observer | switch peers |
| reactor | routing protocols |
| $app_1$ | simple web server |
| $app_2$ | search engine |
| gray failure | app2 observed glitches but neighbors of the bad switch (and app1) didn't |

# Applying the Model

## Case III

| term | example |
|------|---------|
| system | Azure storage service |
| observer | storage master |
| reactor | storage master |
| $app_{1..n}$ | Azure VMs |
| gray failure | some VMs hit remote I/O exceptions while storage master deems ENs healthy |

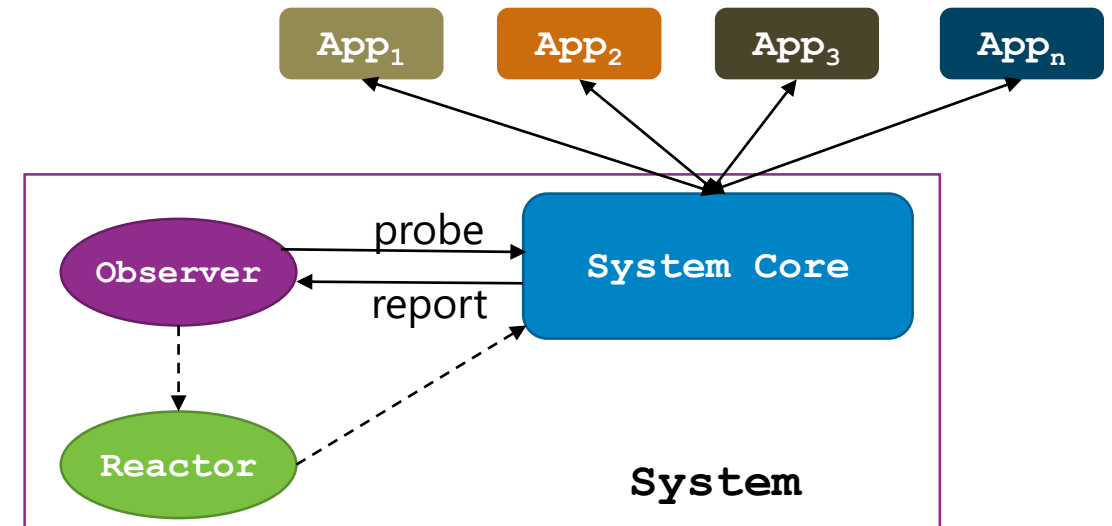# Applying the Model

## Case III

| term | example |
|------|---------|
| system | Azure storage service |
| observer | storage master |
| reactor | storage master |
| $app_{1..n}$ | Azure VMs |
| gray failure | some VMs hit remote I/O exceptions while storage master deems ENs healthy |

# Applying the Model

## Case III

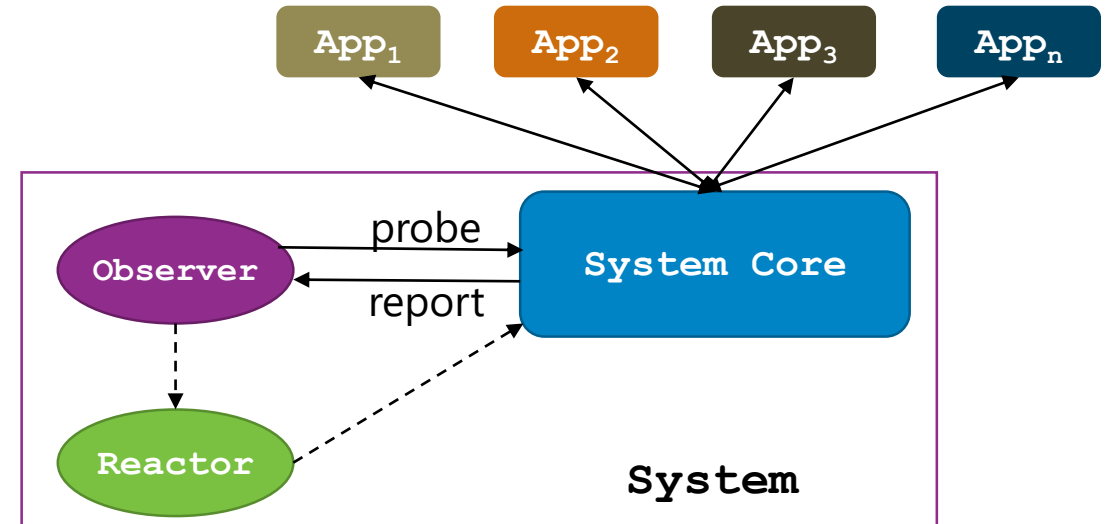| term | example |
|---|---|
| system | Azure storage service |
| observer | storage master |
| reactor | storage master |
| $app_{1..n}$ | Azure VMs |
| gray failure | some VMs hit remote I/O exceptions while storage master deems ENs healthy |



- $EN_3$ degraded
- Observation difference

# Applying the Model

## Case III

| term | example |
|------|---------|
| system | Azure storage service |
| observer | storage master |
| reactor | storage master |
| $app_{1..n}$ | Azure VMs |
| gray failure | some VMs hit remote I/O exceptions while storage master deems ENs healthy |

App₁ App₂ App₃ Appₙ

Observer — probe → System Core
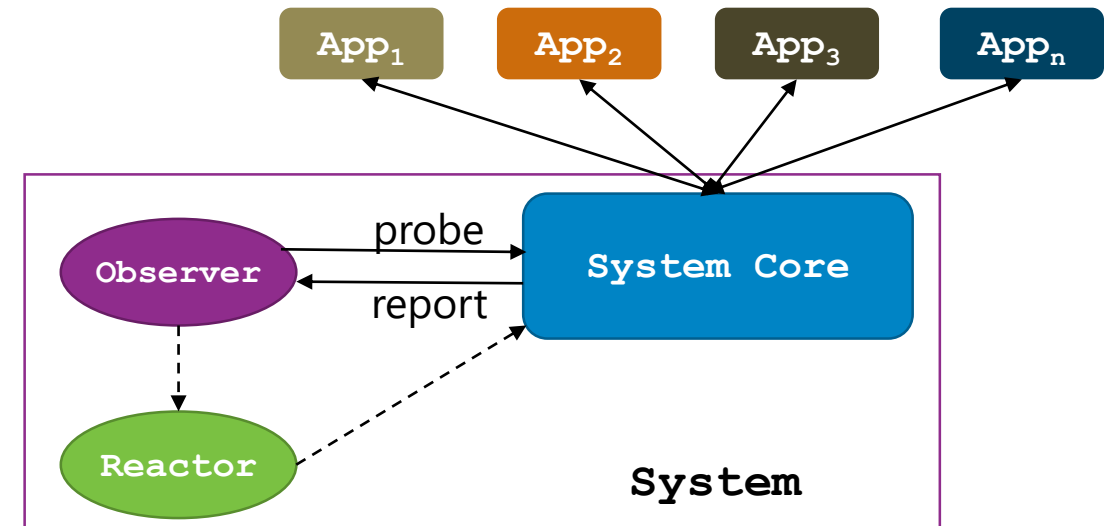
Observer ← report — System Core
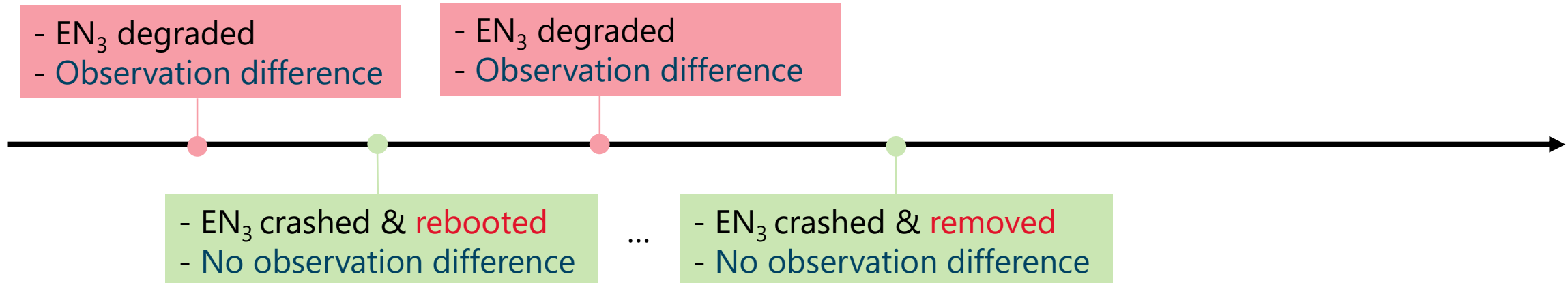
Reactor

System

- EN₃ degraded
- Observation difference

- EN₃ crashed & rebooted
- No observation difference

# Applying the Model

## Case III

| term | example |
|------|---------|
| system | Azure storage service |
| observer | storage master |
| reactor | storage master |
| $app_{1..n}$ | Azure VMs |
| gray failure | some VMs hit remote I/O exceptions while storage master deems ENs healthy |

App₁  App₂  App₃  Appₙ

Observer — probe → System Core
Observer ← report — System Core

Reactor

**System**

- $EN_3$ degraded
- Observation difference

- $EN_3$ degraded
- Observation difference

- $EN_3$ crashed & rebooted
- No observation difference

# Applying the Model

## Case III

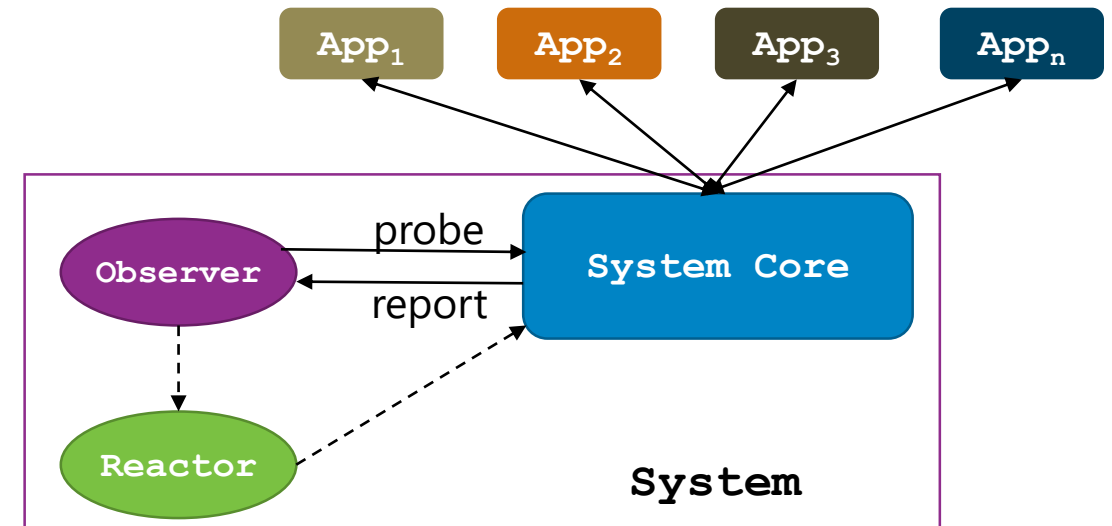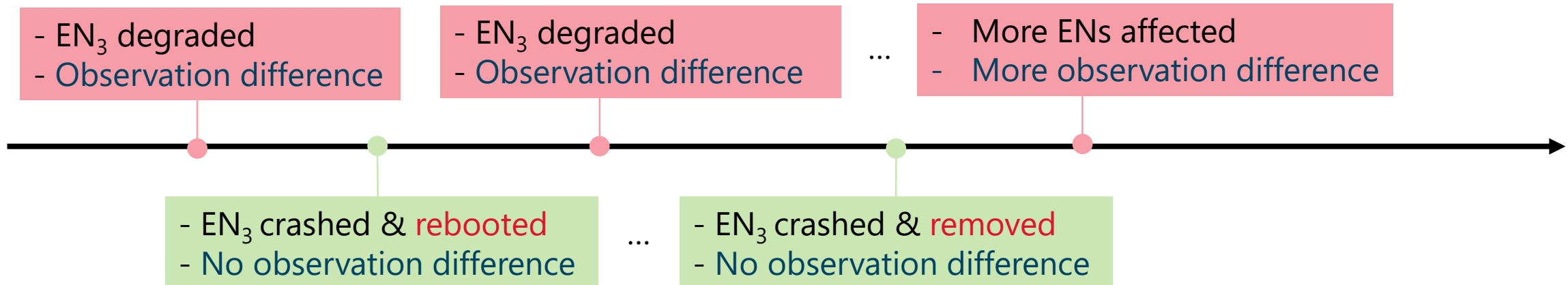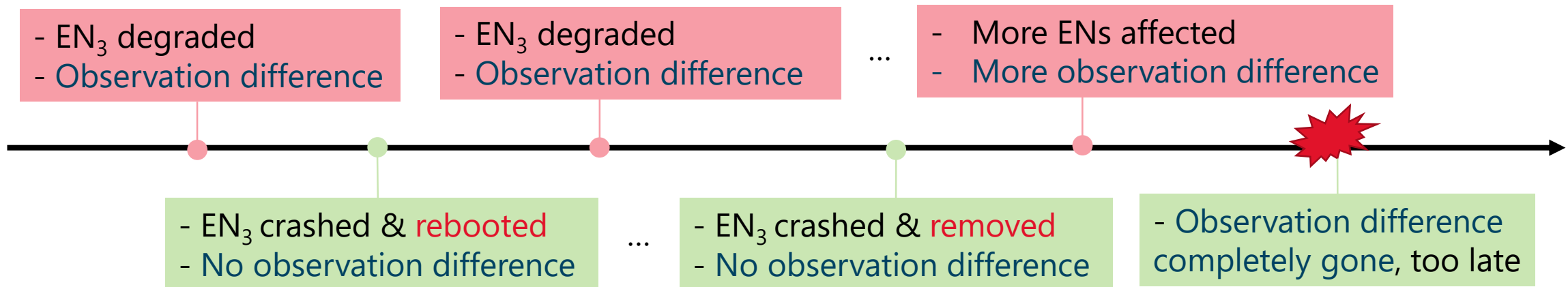| term | example |
|------|---------|
| system | Azure storage service |
| observer | storage master |
| reactor | storage master |
| $app_{1..n}$ | Azure VMs |
| gray failure | some VMs hit remote I/O exceptions while storage master deems ENs healthy |



- EN$_3$ degraded
- Observation difference

- EN$_3$ degraded
- Observation difference

- EN$_3$ crashed & rebooted
- No observation difference

...

- EN$_3$ crashed & removed
- No observation difference

# Applying the Model

## Case III

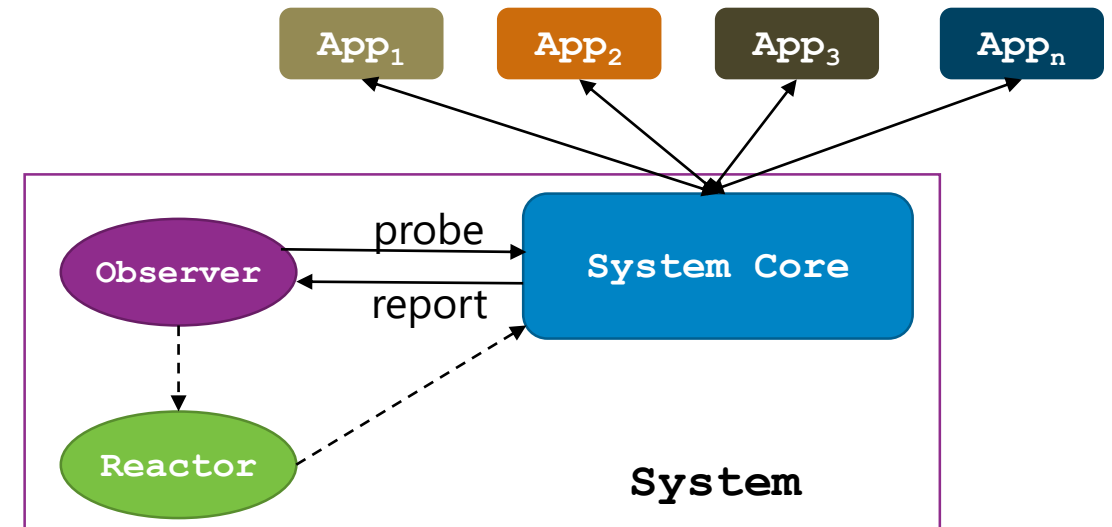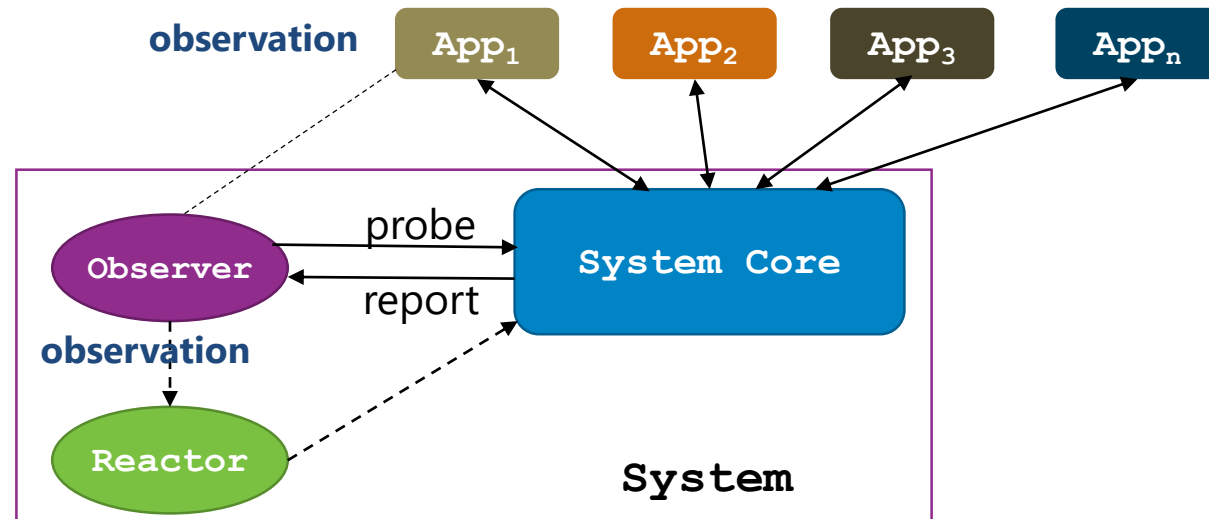| term | example |
|---|---|
| system | Azure storage service |
| observer | storage master |
| reactor | storage master |
| $app_{1..n}$ | Azure VMs |
| gray failure | some VMs hit remote I/O exceptions while storage master deems ENs healthy |



- $EN_3$ degraded
- Observation difference

- $EN_3$ degraded
- Observation difference

...

- More ENs affected
- More observation difference

- $EN_3$ crashed & rebooted
- No observation difference

...

- $EN_3$ crashed & removed
- No observation difference

# Applying the Model

## Case III

| term | example |
|------|---------|
| system | Azure storage service |
| observer | storage master |
| reactor | storage master |
| $app_{1..n}$ | Azure VMs |
| gray failure | some VMs hit remote I/O exceptions while storage master deems ENs healthy |



App₁ App₂ App₃ Appₙ

Observer — probe / report — System Core

Reactor

**System**

- EN$_3$ degraded
- Observation difference

- EN$_3$ degraded
- Observation difference

...

- More ENs affected
- More observation difference

- EN$_3$ crashed & rebooted
- No observation difference

...

- EN$_3$ crashed & removed
- No observation difference

- Observation difference completely gone, too late

108

# Direction 1: Close Observation Gap

Traditional failure detector ➡ multi-dimensional health monitor

# Direction 1: Close Observation Gap
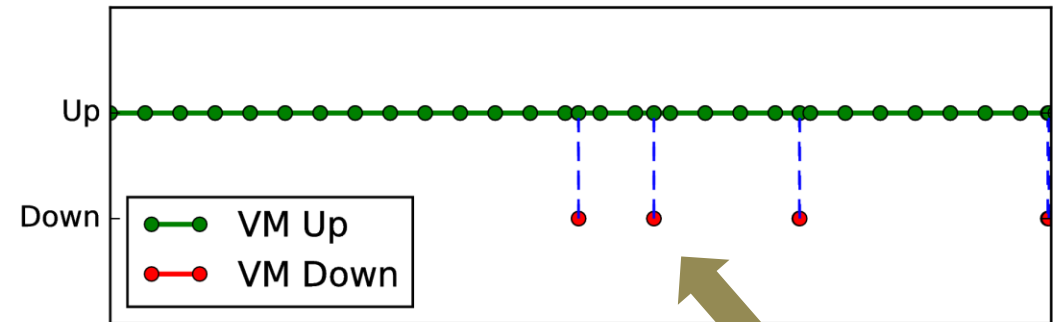
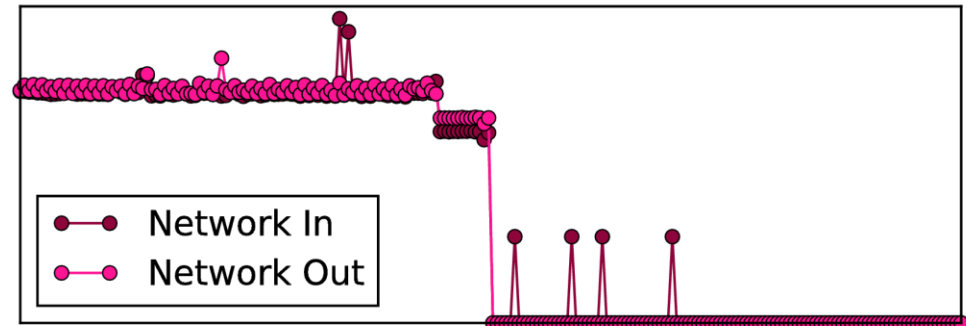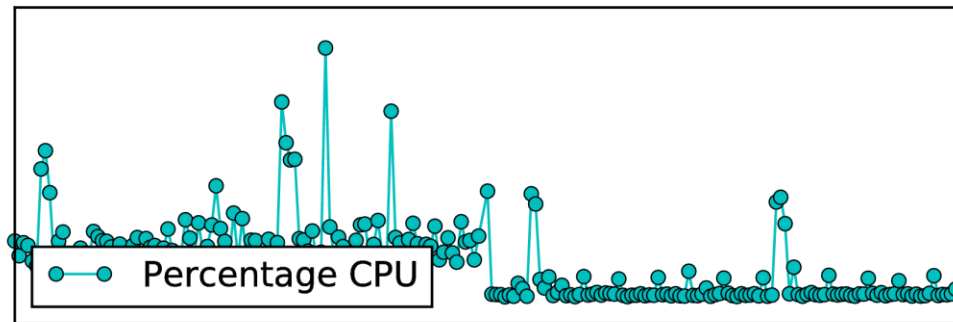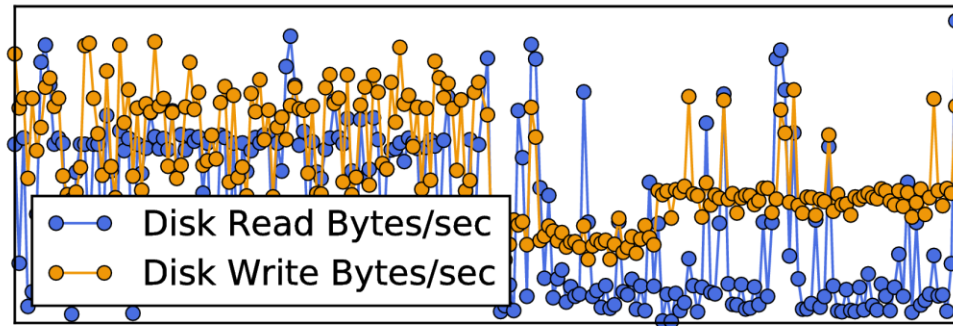Traditional failure detector ➡ multi-dimensional health monitor

Doctors can't use heartbeat as the only signal of a patient's health
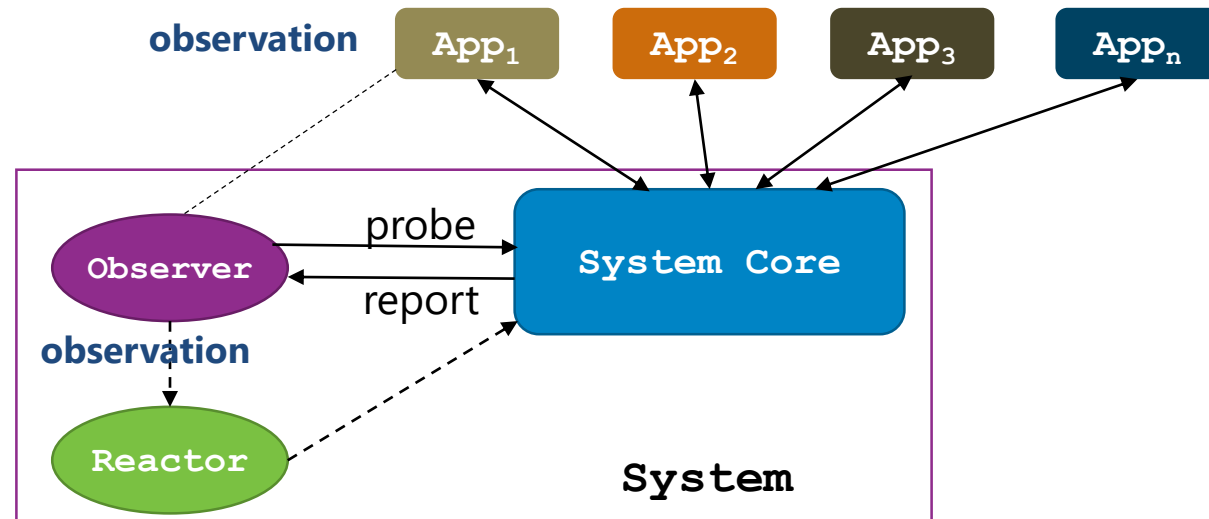
# Direction 1: Close Observation Gap

Traditional failure detector ➡ multi-dimensional health monitor

Doctors can't use heartbeat as the only signal of a patient's health



In-VM performance counters before and during the gray failure incident (case II)

Heartbeat-based hierarchical failure detector
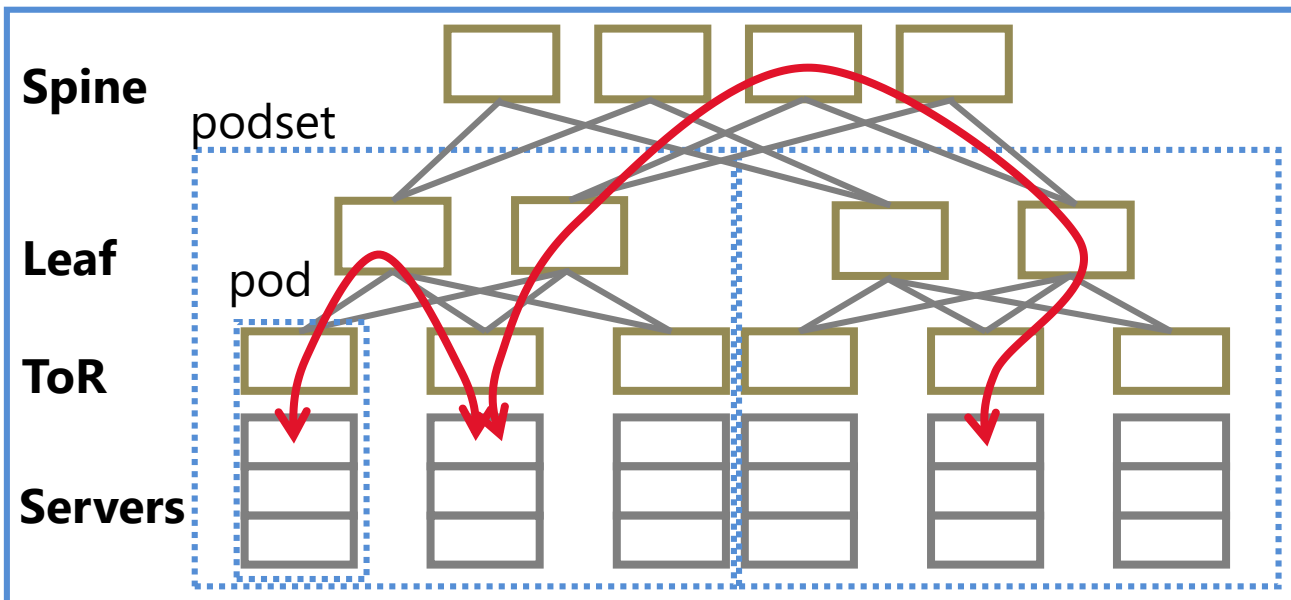
# Direction 2: Approximate Application View

- Infeasible to completely eliminate differential observability due to multi-tenancy and modularity constraints
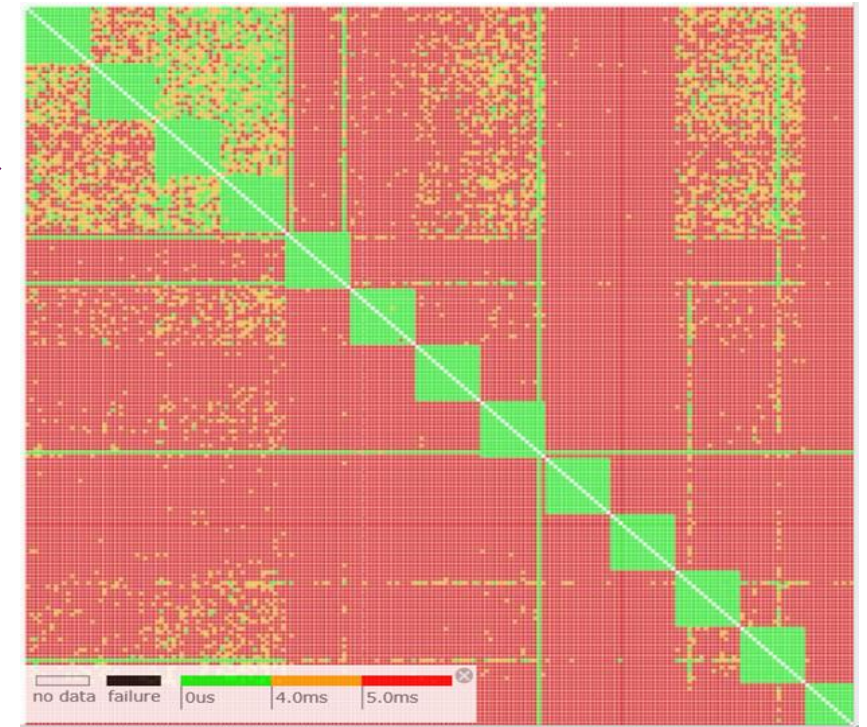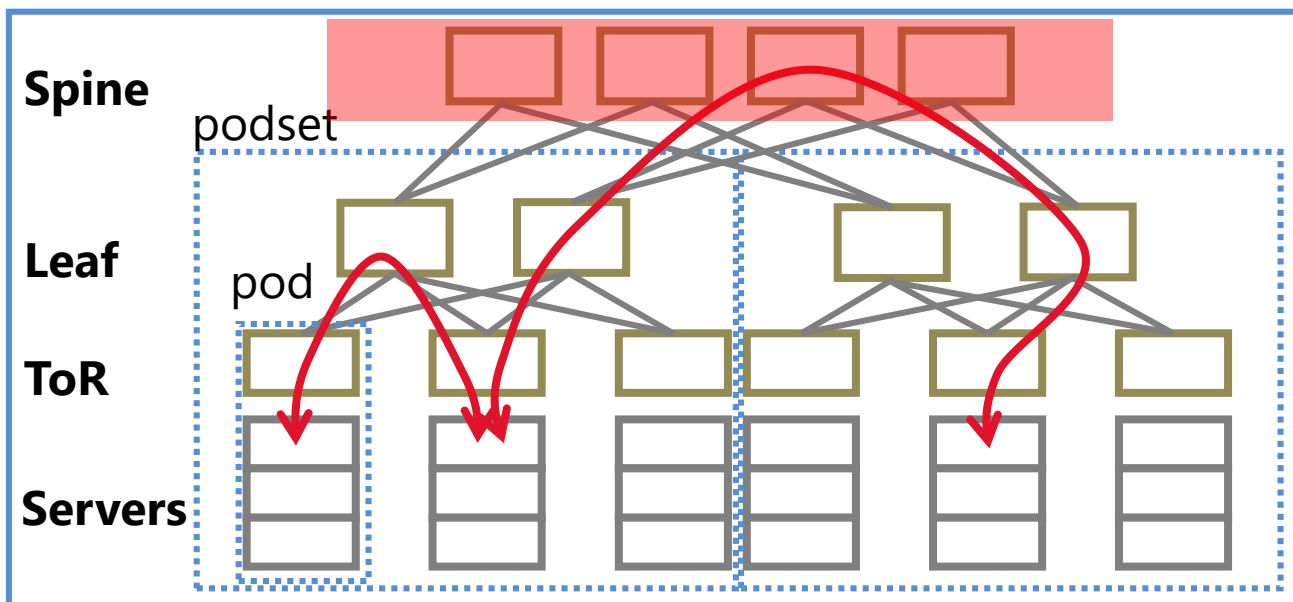- System sends probes to emulate common application usage

# Direction 2: Approximate Application View

PingMesh [SIGCOMM '15]

# Direction 2: Approximate Application View

PingMesh [SIGCOMM '15]

Spine

podset

Leaf

pod

ToR

Servers

no data | failure | 0us | 4.0ms | 5.0ms
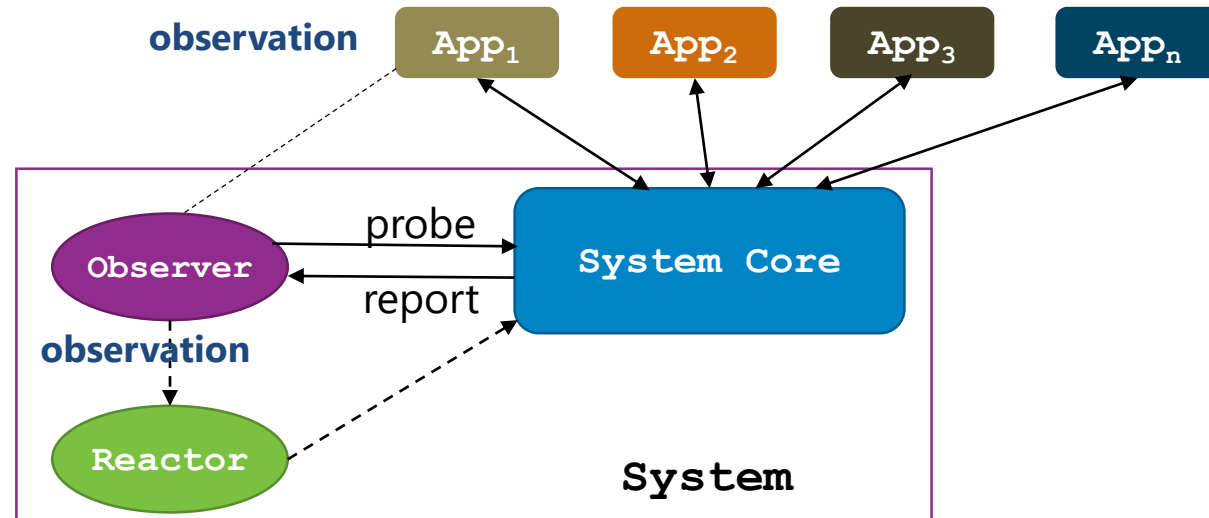
Failure in Spine

114

# Direction 3: Leverage Power of Scale

Break the observation silos to complement each other

# Direction 3: Leverage Power of Scale

Break the observation silos to complement each other

## Observable vs. Detectable

- although end-to-end probe may expose differential observability, it may not detect *who* is responsible for the difference
- example solution: infer signals from many network devices to localize fault

## Blame game

- *A* thinks *B* is problematic, *B* thinks *A* is problematic
- example solution: correlate VM failure events with topology info to judge

# Direction 3: Leverage Power of Scale

Break the observation silos to complement each other
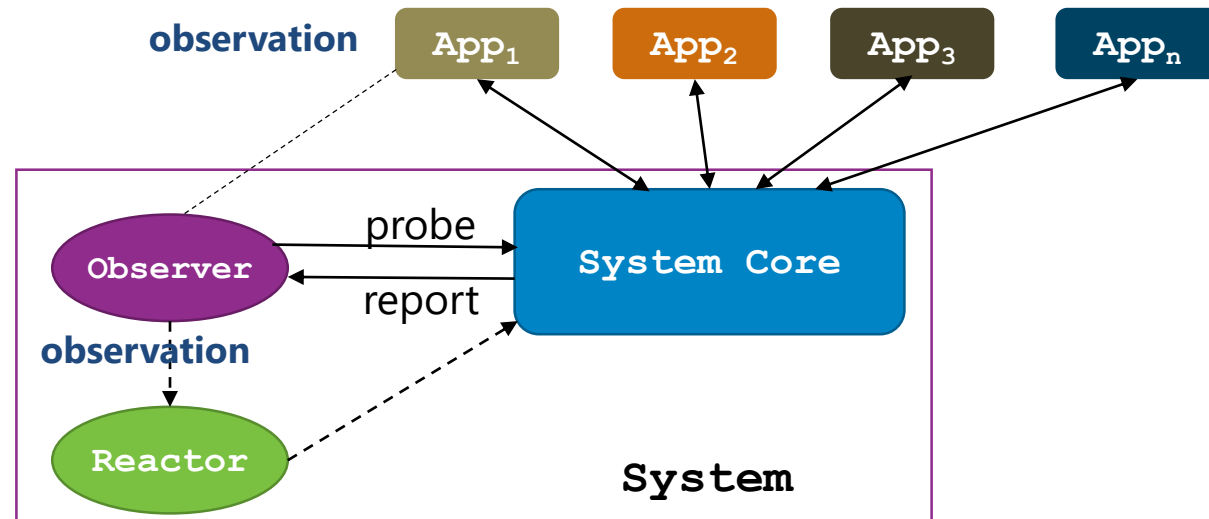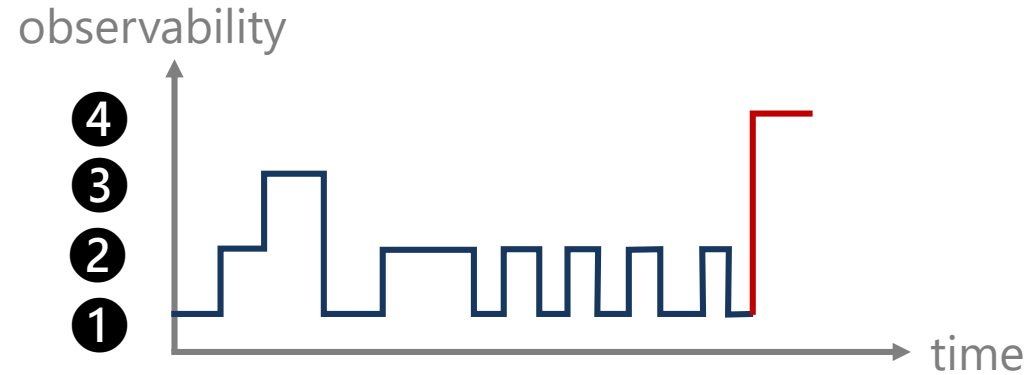
## Observable vs. Detectable

- although end-to-end probe may expose differential observability, it may not detect *who* is responsible for the difference
- example solution: infer signals from many network devices to localize fault

## Blame game

- *A* thinks *B* is problematic, *B* thinks *A* is problematic
- example solution: correlate VM failure events with topology info to judge
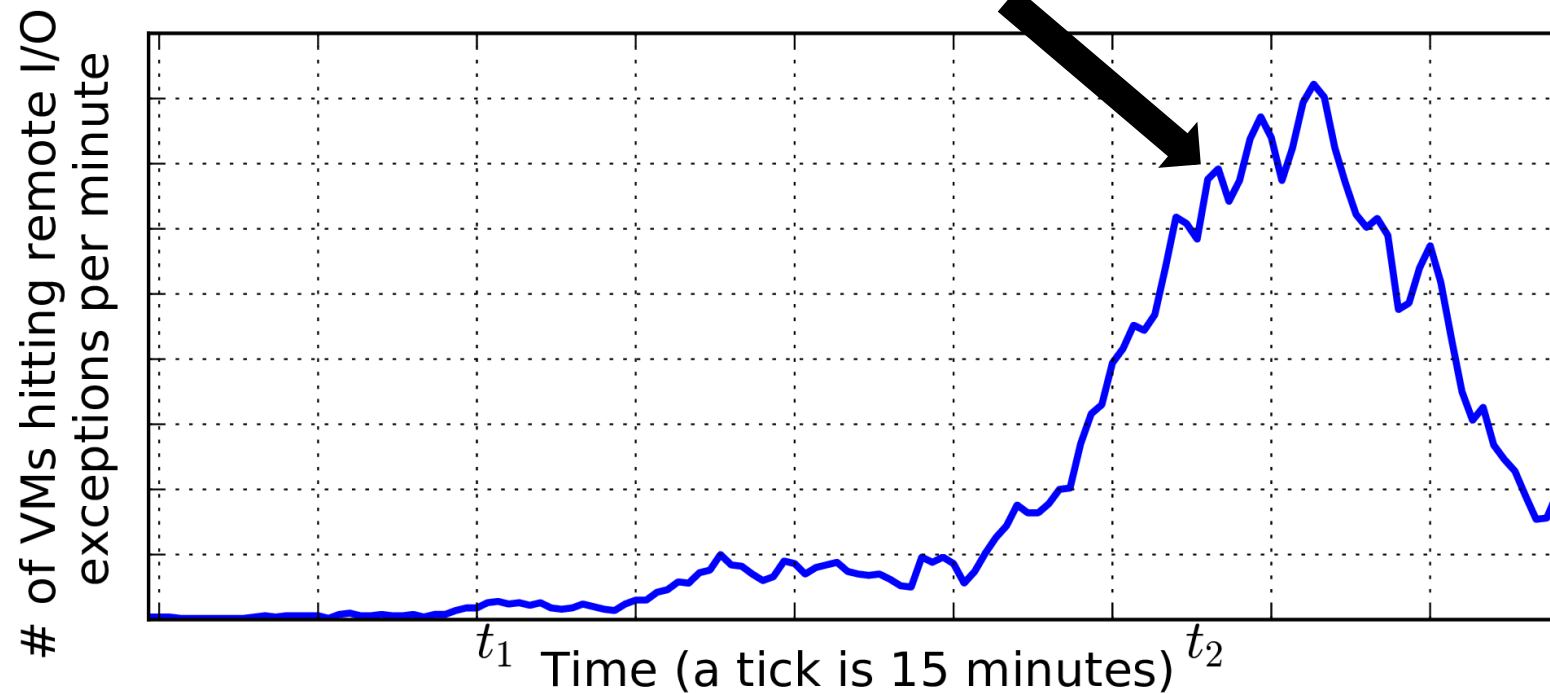
# Direction 4: Harness Temporal Patterns

Time series and trend analysis of key health signals

# Direction 4: Harness Temporal Patterns

Time series and trend analysis of key health signals



**Storage side observed availability issue**

# of VMs hitting remote I/O exceptions per minute

$t_1$  Time (a tick is 15 minutes)  $t_2$

# Conclusion

**Cloud system are adept at handling crash and fail-stop failures**

» decades of efforts and research advances have paid off

**Gray failure is a major challenge moving forward**

» behind many service incidents
» an acute pain for system designers and engineers
» fault tolerance 1.0 �disp 2.0

**A first attempt to define and explore this problem domain**

» differential observability is a fundamental trait
» addressing this trait is key to tackle gray failure
» potential future directions with open challenges

# Discussion

## Why does differential observability occur?

- simplistic model and assumption about component behavior
- modularity principle, unexpected dependency, improper error handling
- focus on narrow point availability but dismissed *app* availability

## Should (can) academia work on this problem?

- practitioners have been troubled by these issues for quite a while, relying on intuition, workaround, resources and process
- hungry for principled approach to understand and tackle the problem
- many issues exist in open-source distributed system stack as well