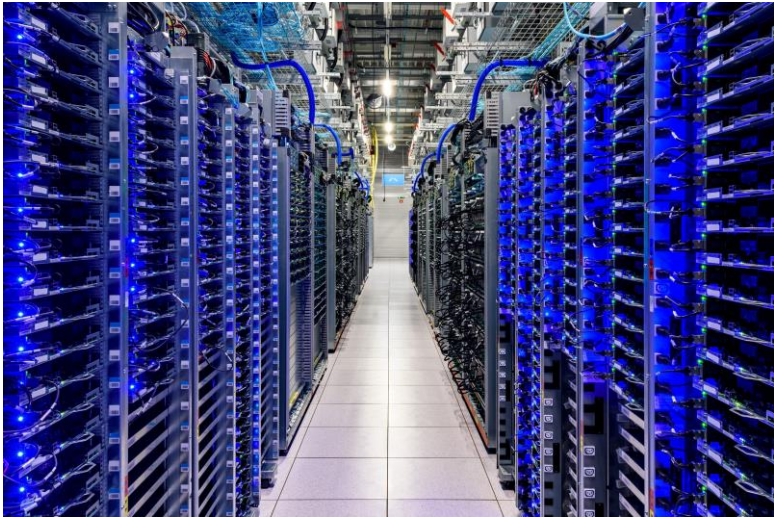


# Occam: A Programming System for Reliable Network Management

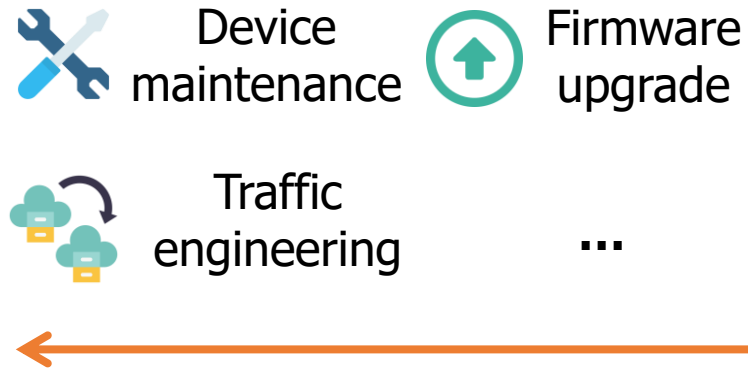
Jiarong Xing, Kuo-Feng Hsu, **Yiting Xia**, Yan Cai,  
Yanping Li, Ying Zhang, Ang Chen



# Problem: Managing large-scale networking services/devices



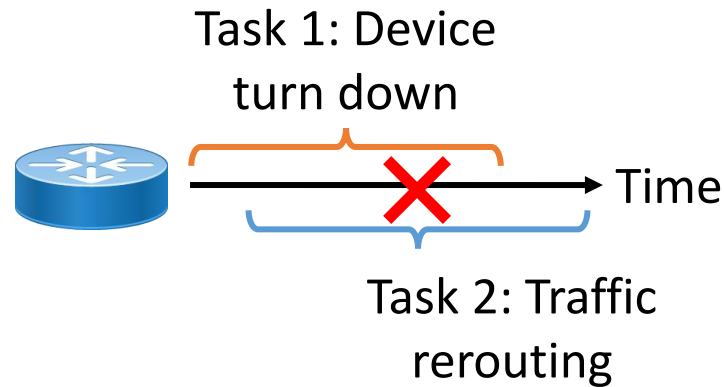
Large-scale data centers



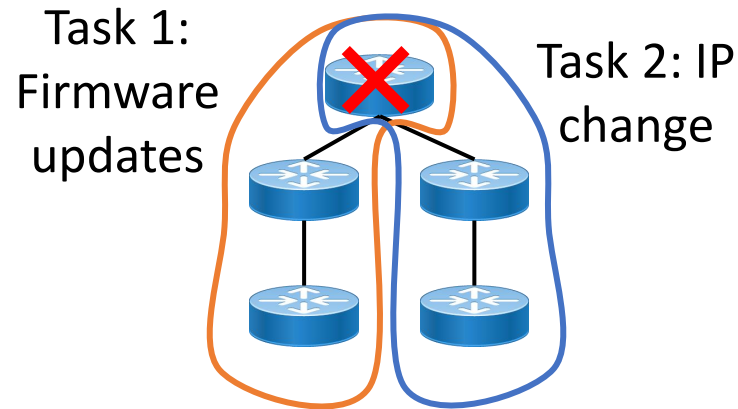
Network management team

- Modern data centers are constructed with large-scale networks
- Networks are constantly changing for evolving demands
  - e.g., some part of the network is being upgraded every day\*
- Network management becomes critical and challenging

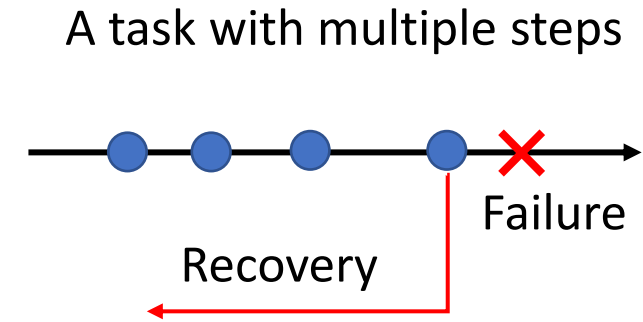
# Key challenge: Reliability!



EX1: Network state conflicts



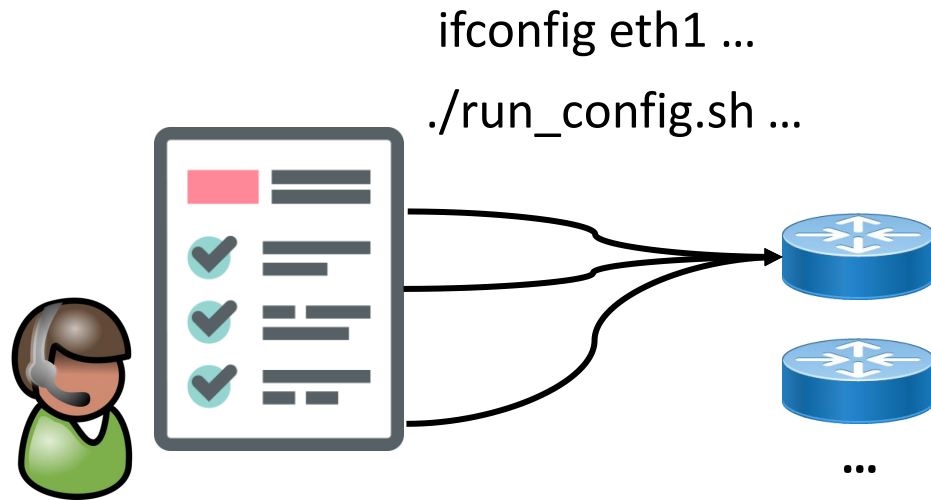
EX2: Device operation conflicts



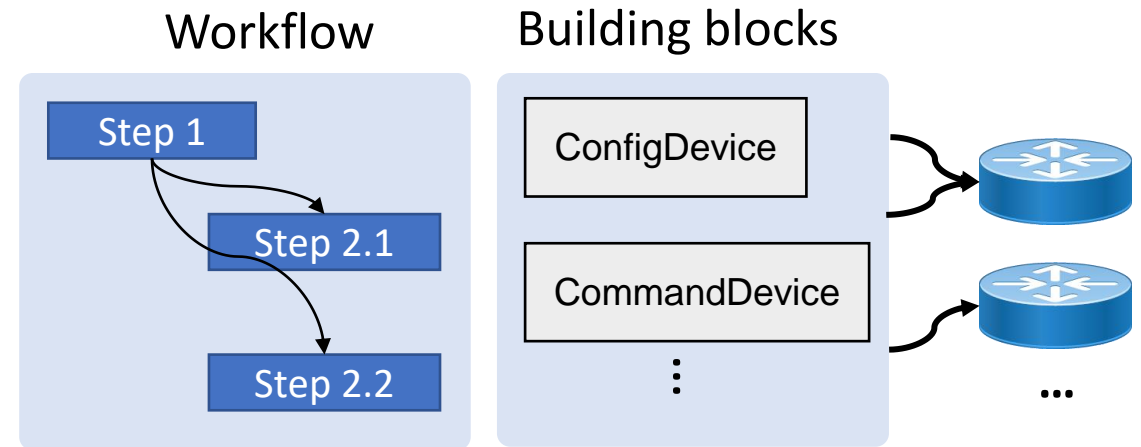
EX3: Failure handling

- Reliability is a key challenge in network management
  - e.g., net. state conflicts, device operation conflicts, failure handling
- Peripheral concerns, but make management very challenging

# Existing solutions and limitations



Method of procedure (MOP)



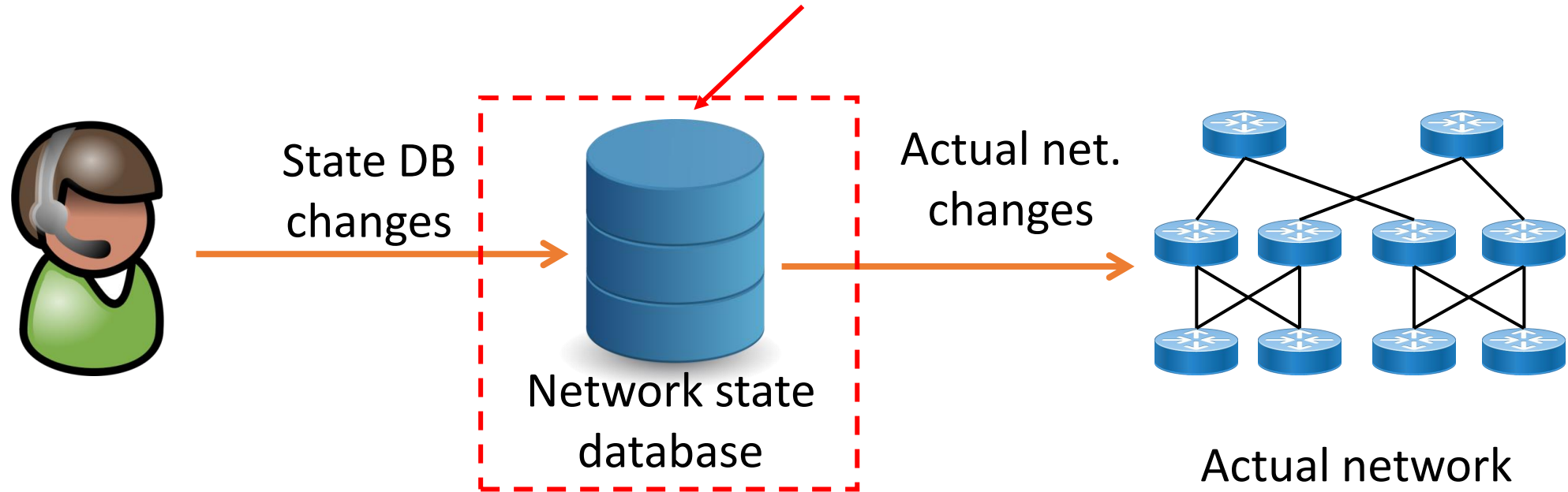
Workflow systems

- Method of procedure (MOP): A guide to manually run the commands
  - Manual operations, inefficient
- Workflow systems with reusable building blocks for automation
  - Not designed to address network-specific issues, e.g., reliability support

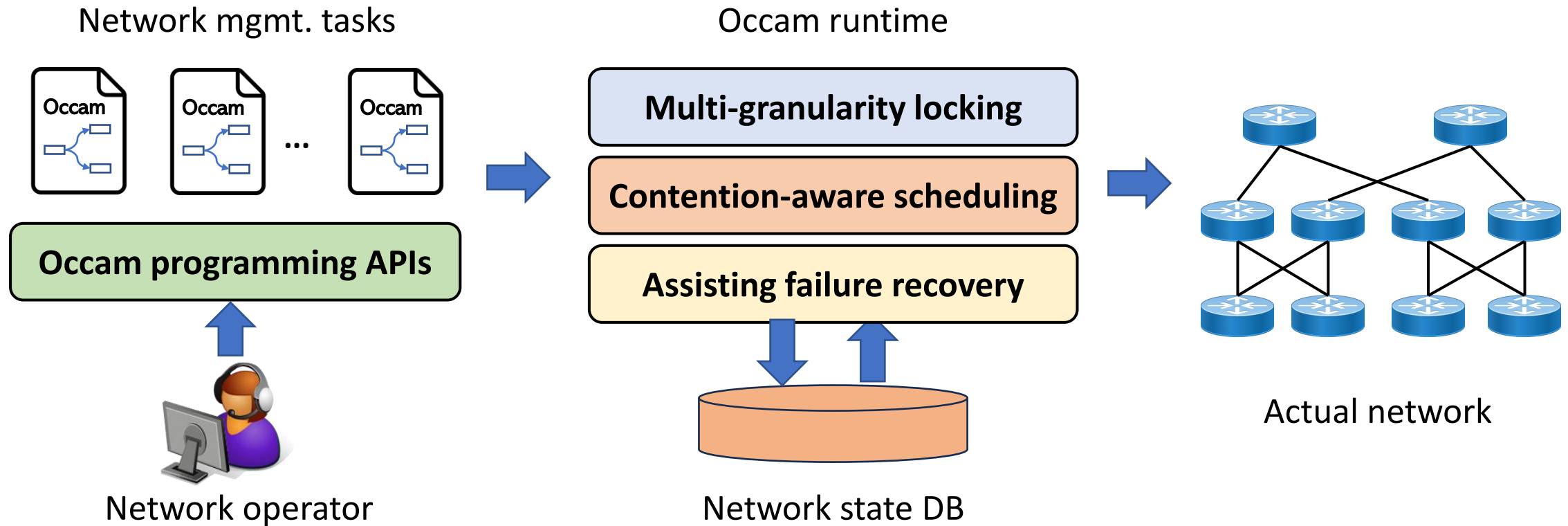
# Key question and idea

How to achieve reliable and efficient network management in modern data centers?

Handling reliability concerns using customized DB techniques



# Occam: Enabling reliable and efficient network management



- **Programming APIs:** Hide operators from peripheral concerns
- **Multi-granularity locking:** Ensures consistent device states
- **Contention-aware scheduling:** Maximizes task execution efficiency
- **Assisting failure recovery:** Facilitates management resilience

# Programming management tasks using Occam APIs

API	Description
→ <b>net(regex)</b>	Create and scope a network object
→ <b>net.get([attr_key])</b>	Get device attribute with key
→ <b>net.set([attr_key, val])</b>	Set device attribute with key
→ <b>net.apply(func)</b>	Execute func on physical device
<b>net.close()</b>	Commit state changes

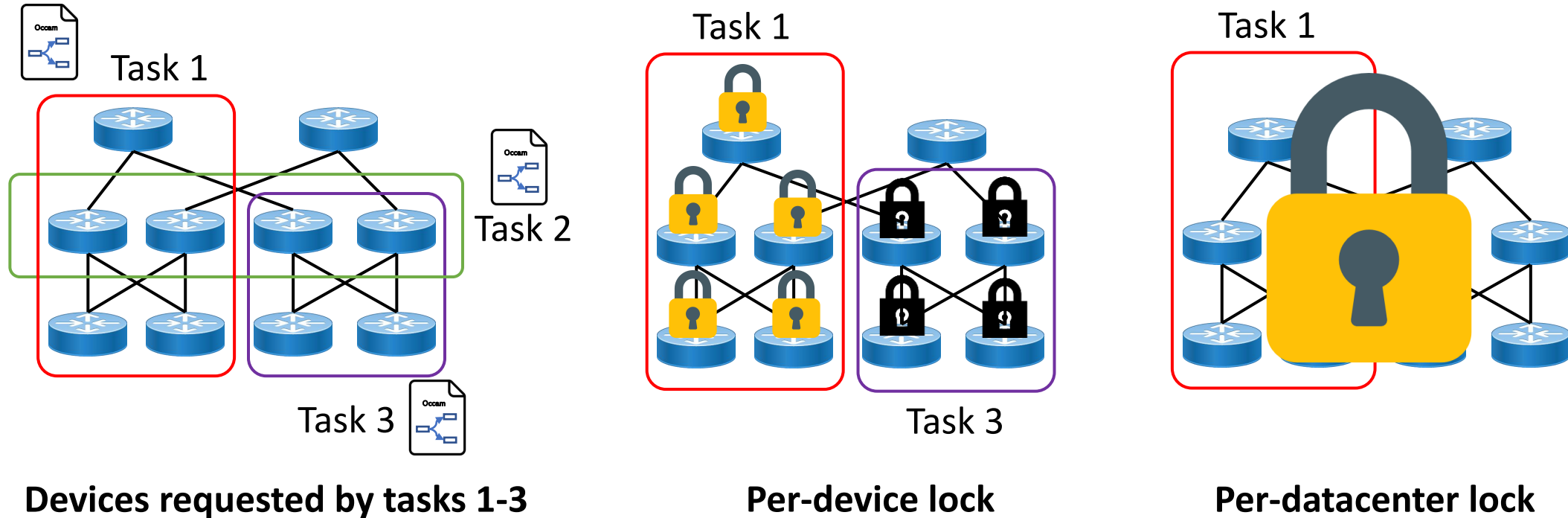
```
→ # device maintenance  
→ dc1pod3=net("dc1.pod3.*")  
→ dc1pod3.set("DEVICE_STATUS",  
             "UNDER_MAINTENANCE")  
→ dc1pod3.apply(f_push)  
dc1pod3.close()
```

Occam programming APIs

Ex: Device maintenance for a pod

- Regular expression (regex) to construct network object
- Capture the database read/write access with get()/set()
- Capture the device function execution with apply()
- Provides task semantics for Occam runtime

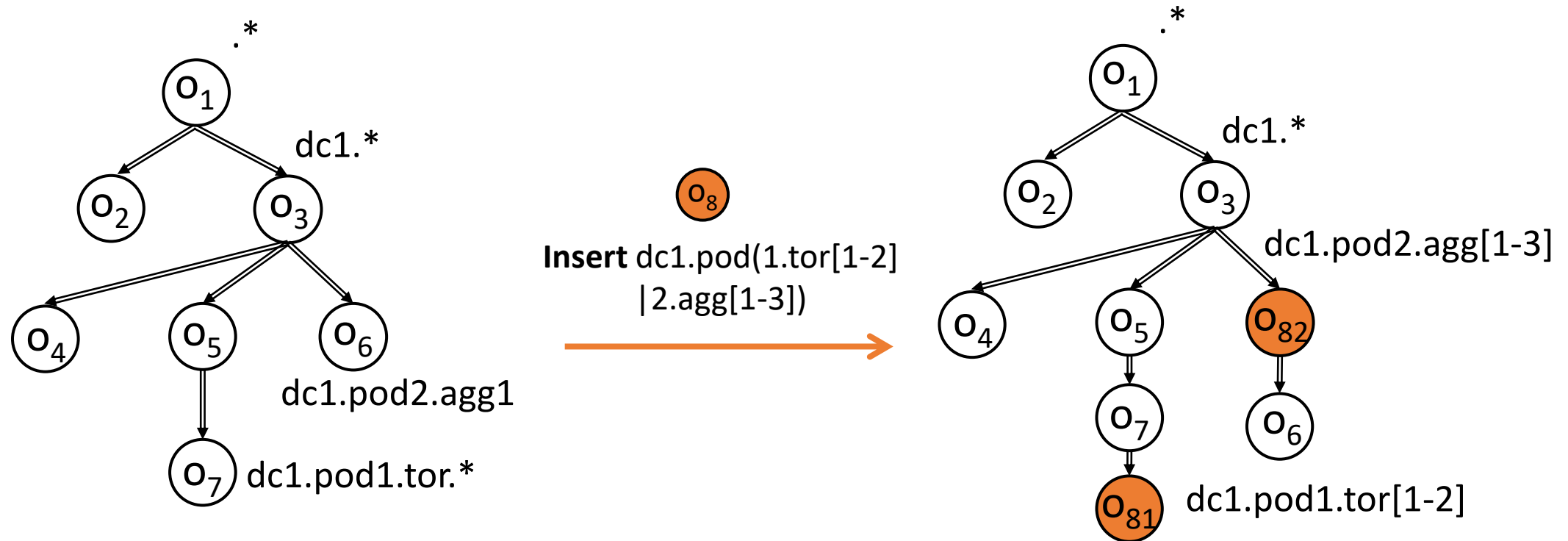
# How to ensure consistent device state changes?



- Different tasks can request to change the same set of devices
  - Could cause inconsistent state changes
- Per-device lock: Too fine-grained, high lock management overhead
- Per-datacenter lock: Too coarse-grained, limited parallelism

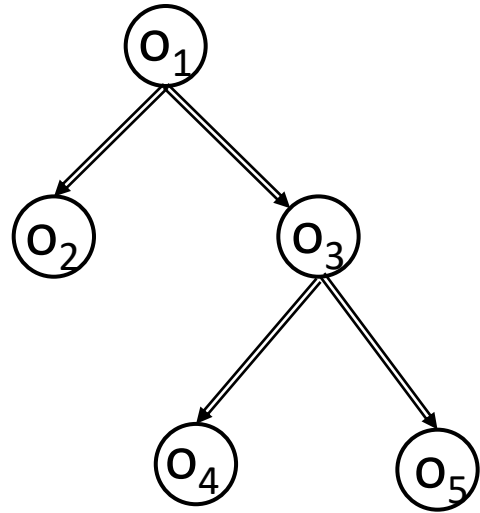


# Regex-based multi-granularity locking



- The locking unit is the **network object (scoped by regex)**
  - Varied sizes based on the management task
- A regex tree is used to encode containment relationships

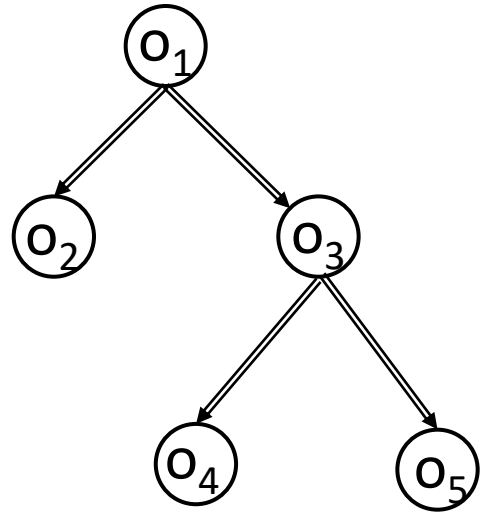
# Regex-based multi-granularity locking (cont'd)



○ Object

○<sub>1</sub> → ○<sub>2</sub> o<sub>1</sub> contains o<sub>2</sub>

# Regex-based multi-granularity locking (cont'd)

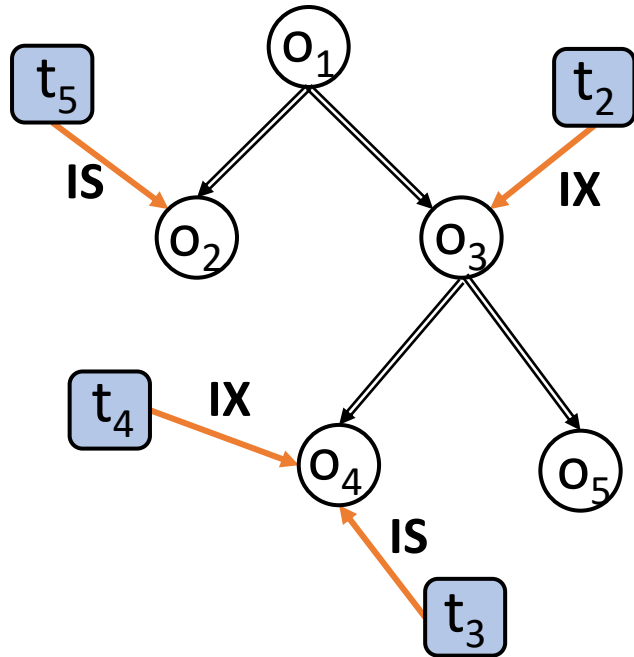


○ Object

○<sub>1</sub> → ○<sub>2</sub> o<sub>1</sub> contains o<sub>2</sub>

- The regex tree is enhanced with lock dependency
  - X/S: Exclusive/share locks
  - IX/IS: Intentional exclusive/share locks

# Regex-based multi-granularity locking (cont'd)



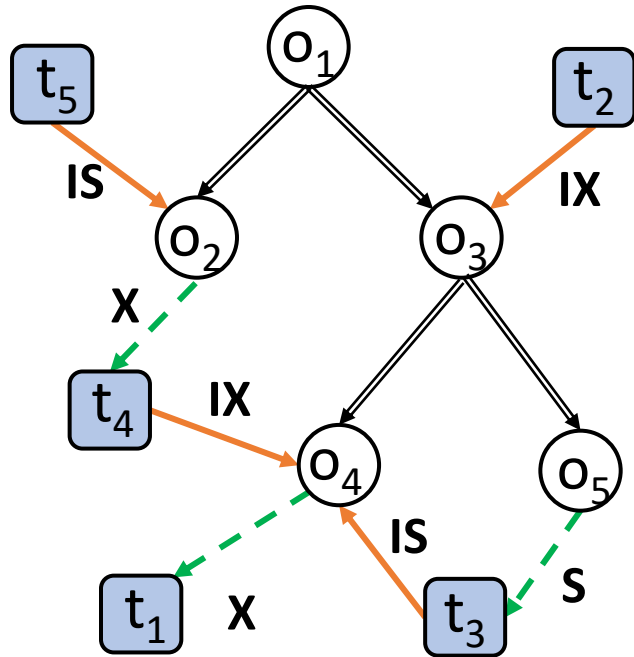
○ Object    □ Task

□  $t_1$  → ○  $o_1$      $t_1$  is waiting for the lock on  $o_1$

○  $o_1$  → ○  $o_2$      $o_1$  contains  $o_2$

- The regex tree is enhanced with lock dependency
  - X/S: Exclusive/share locks
  - IX/IS: Intentional exclusive/share locks

# Regex-based multi-granularity locking (cont'd)



○ Object    □ Task

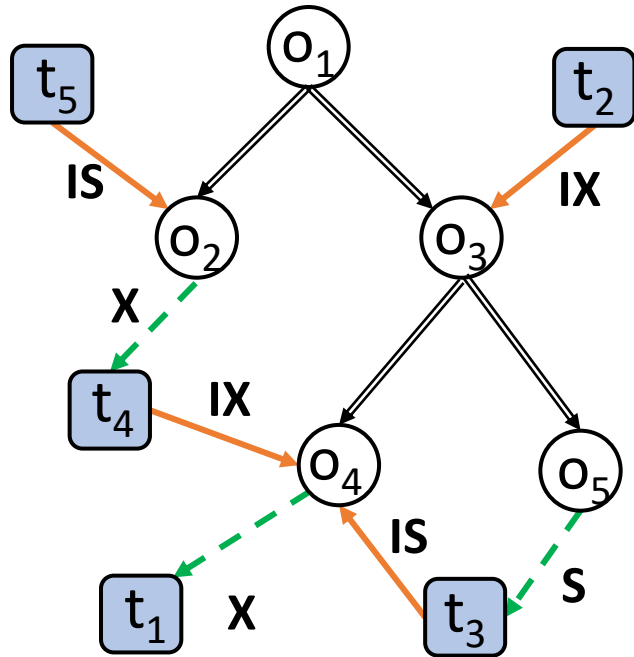
□  $t_1$  → ○  $o_1$   $t_1$  is waiting for the lock on  $o_1$

○  $o_1$  - - -> □  $t_1$  The lock on  $o_1$  has been granted to  $t_1$

○  $o_1$  ⇒ ○  $o_2$   $o_1$  contains  $o_2$

- The regex tree is enhanced with lock dependency
  - X/S: Exclusive/share locks
  - IX/IS: Intentional exclusive/share locks

# Regex-based multi-granularity locking (cont'd)



○ Object    □ Task

□ t<sub>1</sub> → ○ o<sub>1</sub>    t<sub>1</sub> is waiting for the lock on o<sub>1</sub>

○ o<sub>1</sub> - -> □ t<sub>1</sub>    The lock on o<sub>1</sub> has been granted to t<sub>1</sub>

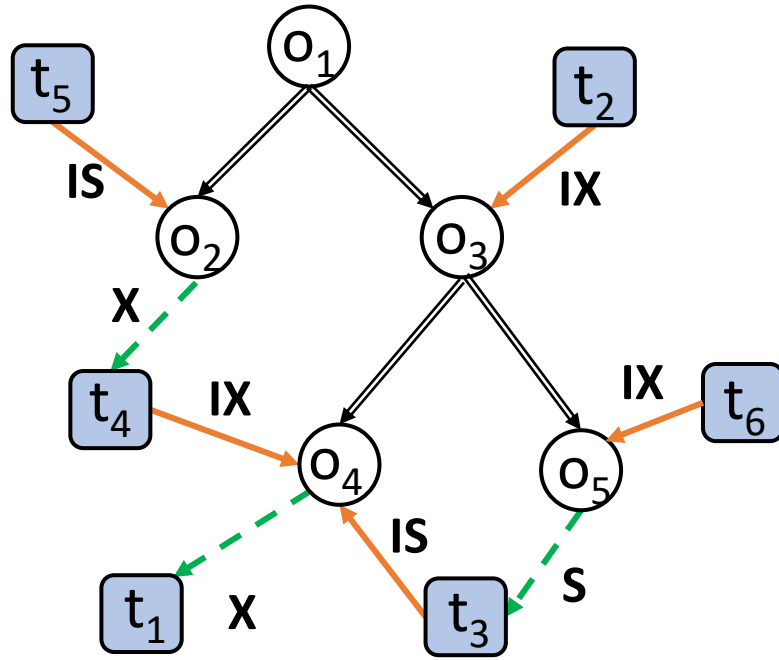
○ o<sub>1</sub> ⇒ ○ o<sub>2</sub>    o<sub>1</sub> contains o<sub>2</sub>

- The regex tree is enhanced with lock dependency
  - X/S: Exclusive/share locks
  - IX/IS: Intentional exclusive/share locks
- Tasks are scheduled based on this dependency graph

# Contention-aware scheduling

- Occam supports Largest Dependency Set First (LDSF) scheduling
  - Certain tasks may be blocking more tasks than others
  - Prioritizing them could simultaneously enable more tasks to proceed

# Contention-aware scheduling

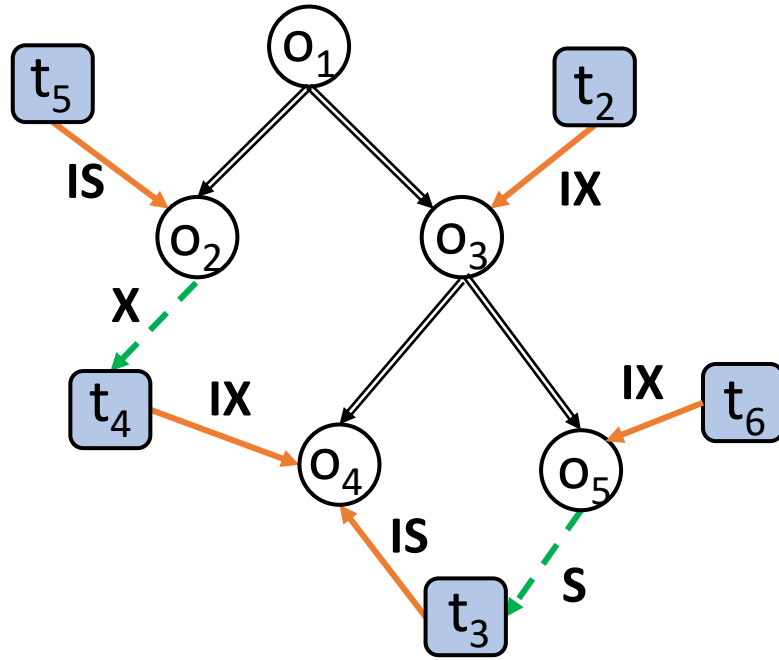


When **t1** releases **o4**, should **t3** or **t4** be scheduled?

- Occam supports Largest Dependency Set First (LDSF) scheduling
  - Certain tasks may be blocking more tasks than others
  - Prioritizing them could simultaneously enable more tasks to proceed



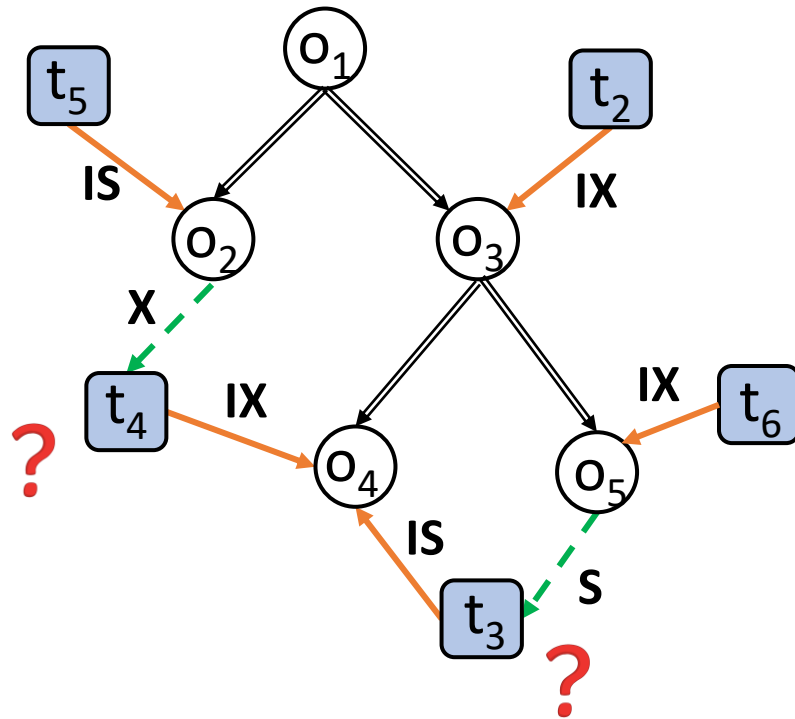
# Contention-aware scheduling



When **t1** releases **o4**, should **t3** or **t4** be scheduled?

- Occam supports Largest Dependency Set First (LDSF) scheduling
  - Certain tasks may be blocking more tasks than others
  - Prioritizing them could simultaneously enable more tasks to proceed

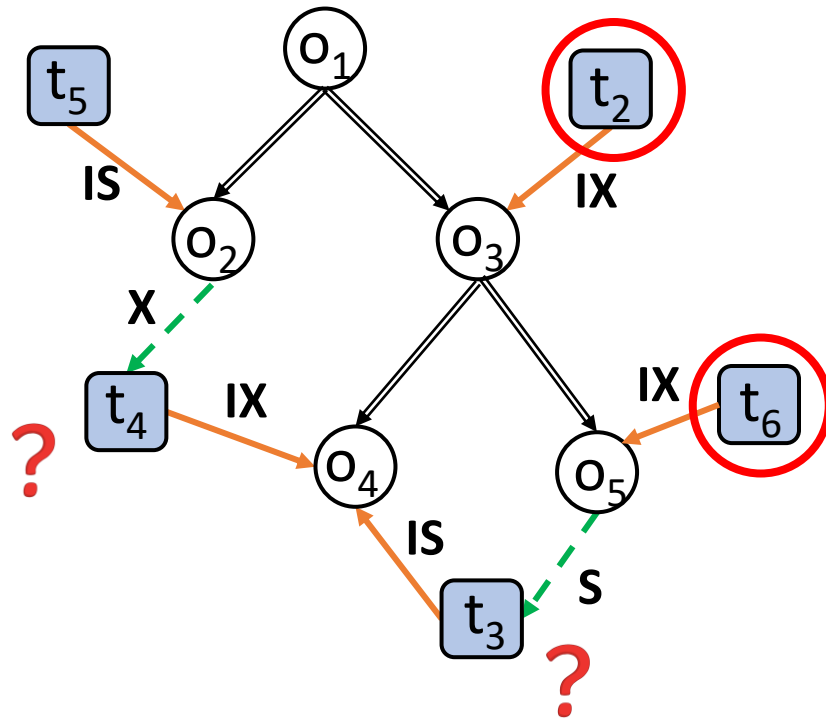
# Contention-aware scheduling



When **t1** releases **o4**, should **t3** or **t4** be scheduled?

- Occam supports Largest Dependency Set First (LDSF) scheduling
  - Certain tasks may be blocking more tasks than others
  - Prioritizing them could simultaneously enable more tasks to proceed

# Contention-aware scheduling

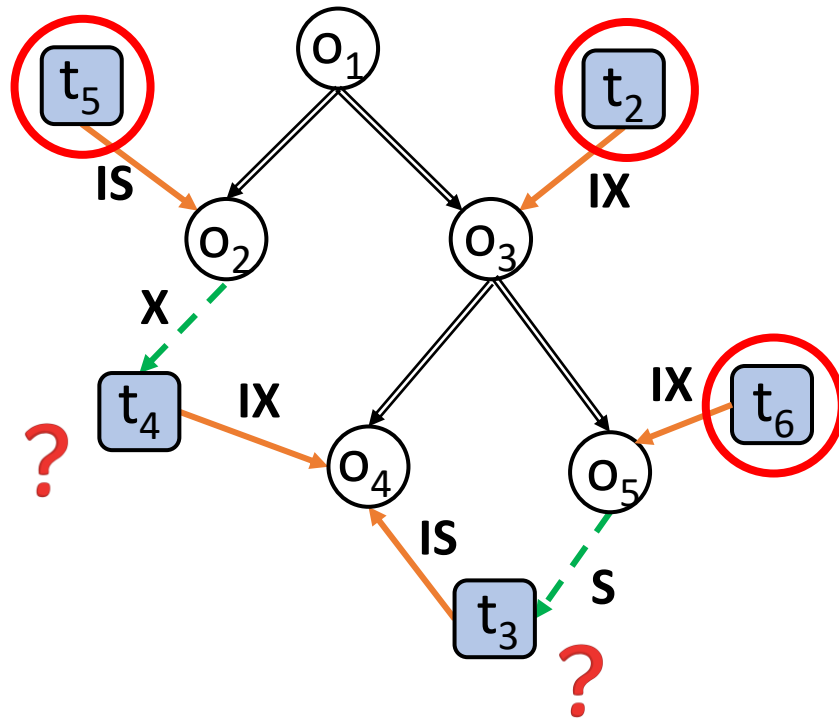


When **t1** releases **o4**, should **t3** or **t4** be scheduled?

**t3** blocks **t2** and **t6** because of **S** lock on **o5**

- Occam supports Largest Dependency Set First (LDSF) scheduling
  - Certain tasks may be blocking more tasks than others
  - Prioritizing them could simultaneously enable more tasks to proceed

# Contention-aware scheduling



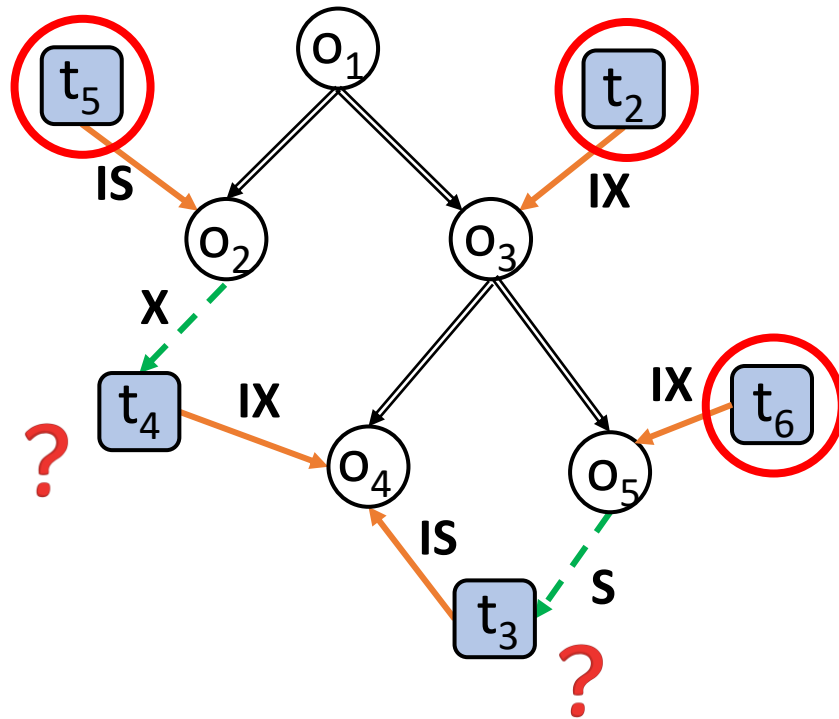
When **t1** releases **o4**, should **t3** or **t4** be scheduled?

**t3** blocks **t2** and **t6** because of **S** lock on **o5**

**t4** blocks **t5** because of **X** lock on **o2**

- Occam supports Largest Dependency Set First (LDSF) scheduling
  - Certain tasks may be blocking more tasks than others
  - Prioritizing them could simultaneously enable more tasks to proceed

# Contention-aware scheduling



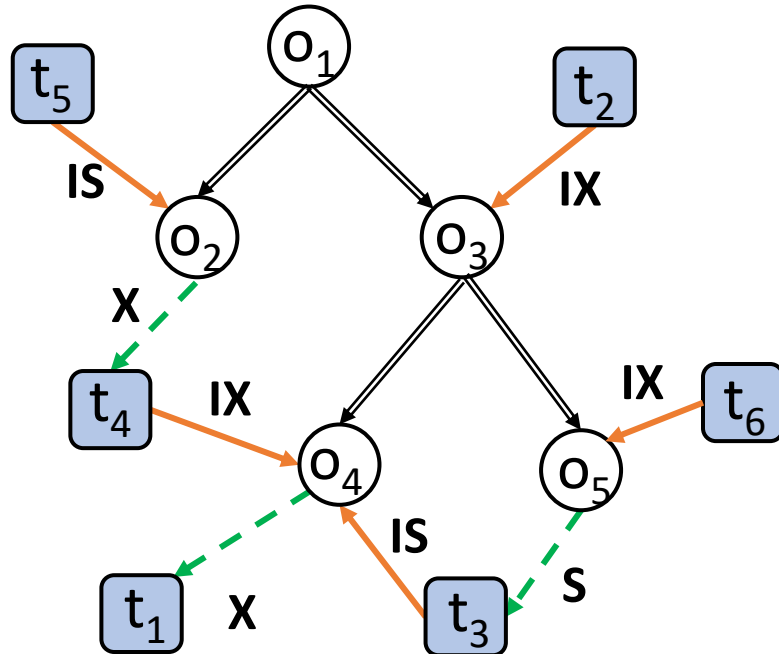
When  $t_1$  releases  $o_4$ , should  $t_3$  or  $t_4$  be scheduled?

$t_3$  blocks  $t_2$  and  $t_6$  because of S lock on  $o_5$

$t_4$  blocks  $t_5$  because of X lock on  $o_2$

- Occam supports Largest Dependency Set First (LDSF) scheduling
  - Certain tasks may be blocking more tasks than others
  - Prioritizing them could simultaneously enable more tasks to proceed

# Contention-aware scheduling



When **t1** releases **o4**, should **t3** or **t4** be scheduled?

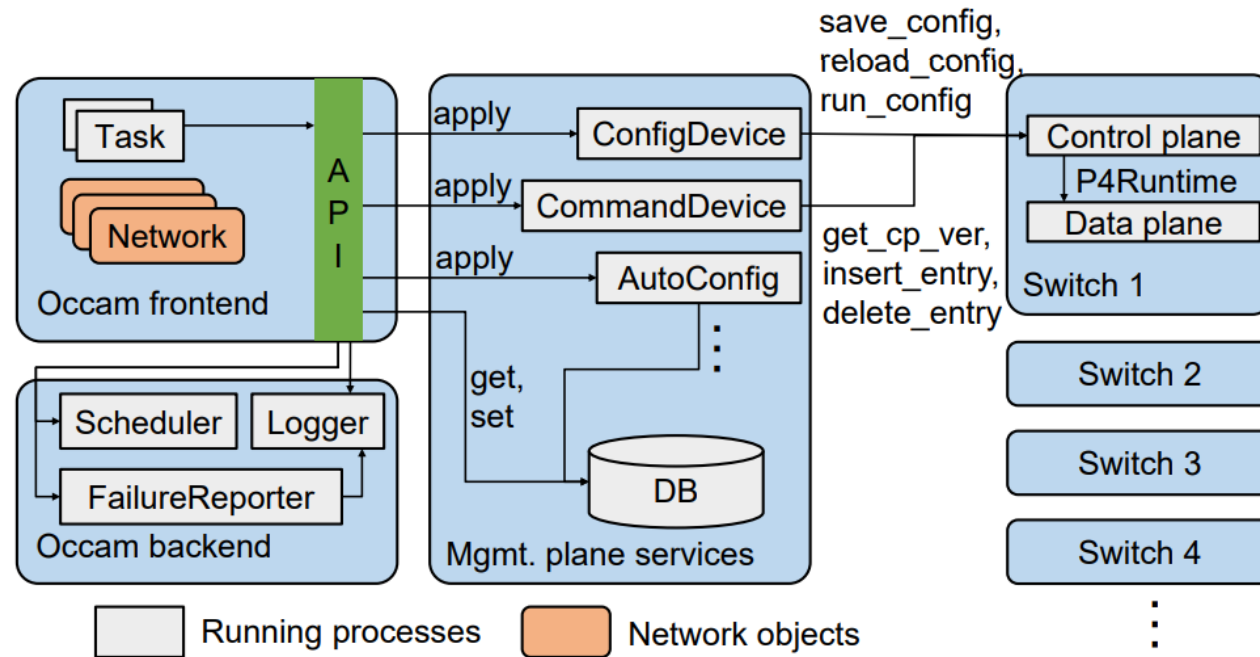
**t3** blocks **t2** and **t6** because of **S** lock on **o5**

**t4** blocks **t5** because of **X** lock on **o2**



- Occam supports Largest Dependency Set First (LDSF) scheduling
  - Certain tasks may be blocking more tasks than others
  - Prioritizing them could simultaneously enable more tasks to proceed
- Occam also assists failed tasks recovery (see §6 for details)

# Implementation and evaluation setup

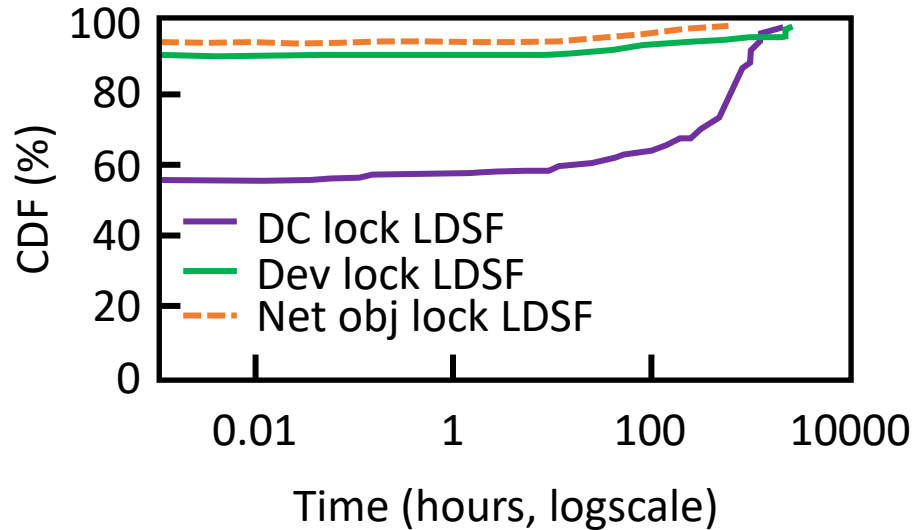


- Mininet-based network management emulation
  - Bmv2 P4 software switches
  - P4Runtime API for state changes
- Large-scale simulation
  - 5-month trace from Meta

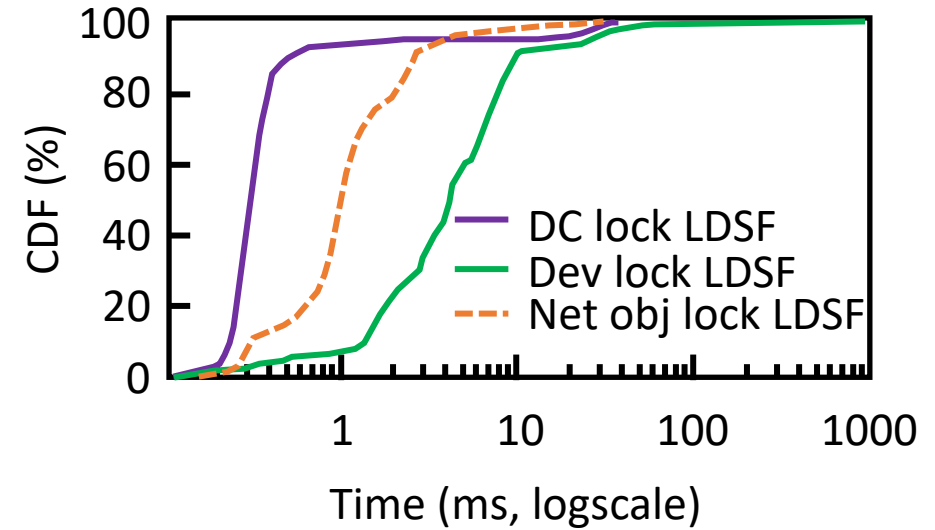


# Locking and scheduling performance

Higher is better



**Task waiting time**



**Scheduling computation time**

- Occam's multi-granularity locking and LDSF scheduling
  - Have minimal waiting time compared to per DC/device locking
  - With reasonable schedule computation time



# Summary

- Motivation: Reliable network management for large-scale data centers
- Occam provides
  - **Unified API:** Enable network-specific design
  - **Multi-granularity locking:** Avoid conflicts with multi-granularity locks
  - **Contention-aware scheduling:** Minimize waiting time via LSDF scheduling
  - **Assisting failure recovery:** Generate plans when failure happens
- Implementation and evaluation:
  - Open-sourced Occam emulator and simulator
  - Minimized waiting time with reasonable scheduling overhead

