# ESCHER

Expressive SCHeduling with Ephemeral Resources

**Romil Bhardwaj**, Alexey Tumanov, Richard Liaw, Stephanie Wang, Robert Nishihara, Philipp Moritz, Ion Stoica
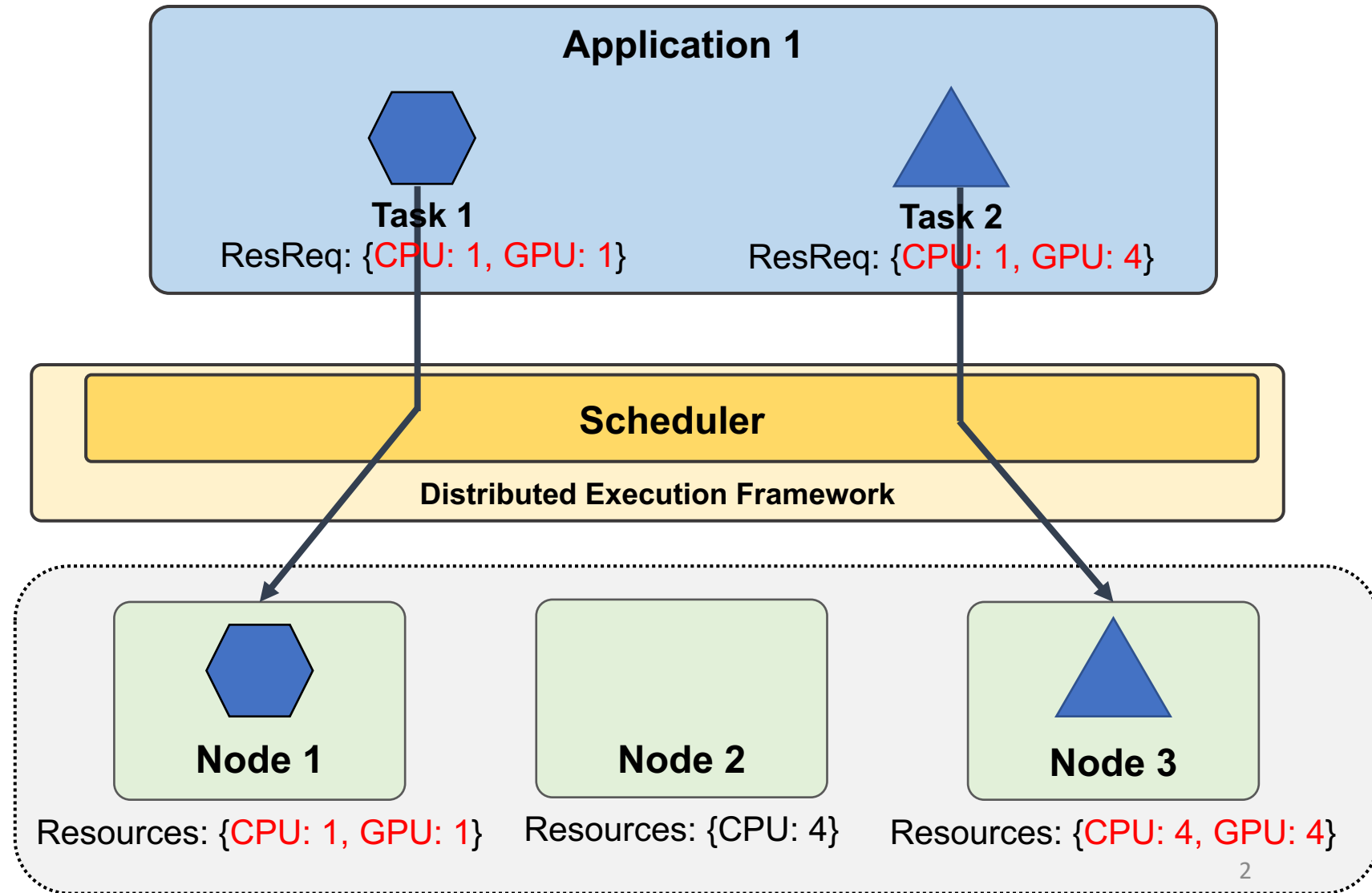
*Relativity* by M. C. Escher

# Typical Distributed Application

Application is composed of tasks with resource requirements

**Scheduler** matches task resource **requirements** to node resource **availabilities**

Cluster is composed of nodes with resource configurations

**Application 1**

**Task 1**
ResReq: {CPU: 1, GPU: 1}

**Task 2**
ResReq: {CPU: 1, GPU: 4}

**Scheduler**

**Distributed Execution Framework**

**Node 1**

**Node 2**

**Node 3**

Resources: {CPU: 1, GPU: 1}

Resources: {CPU: 4}

Resources: {CPU: 4, GPU: 4}

2

# Example - Distributed Training

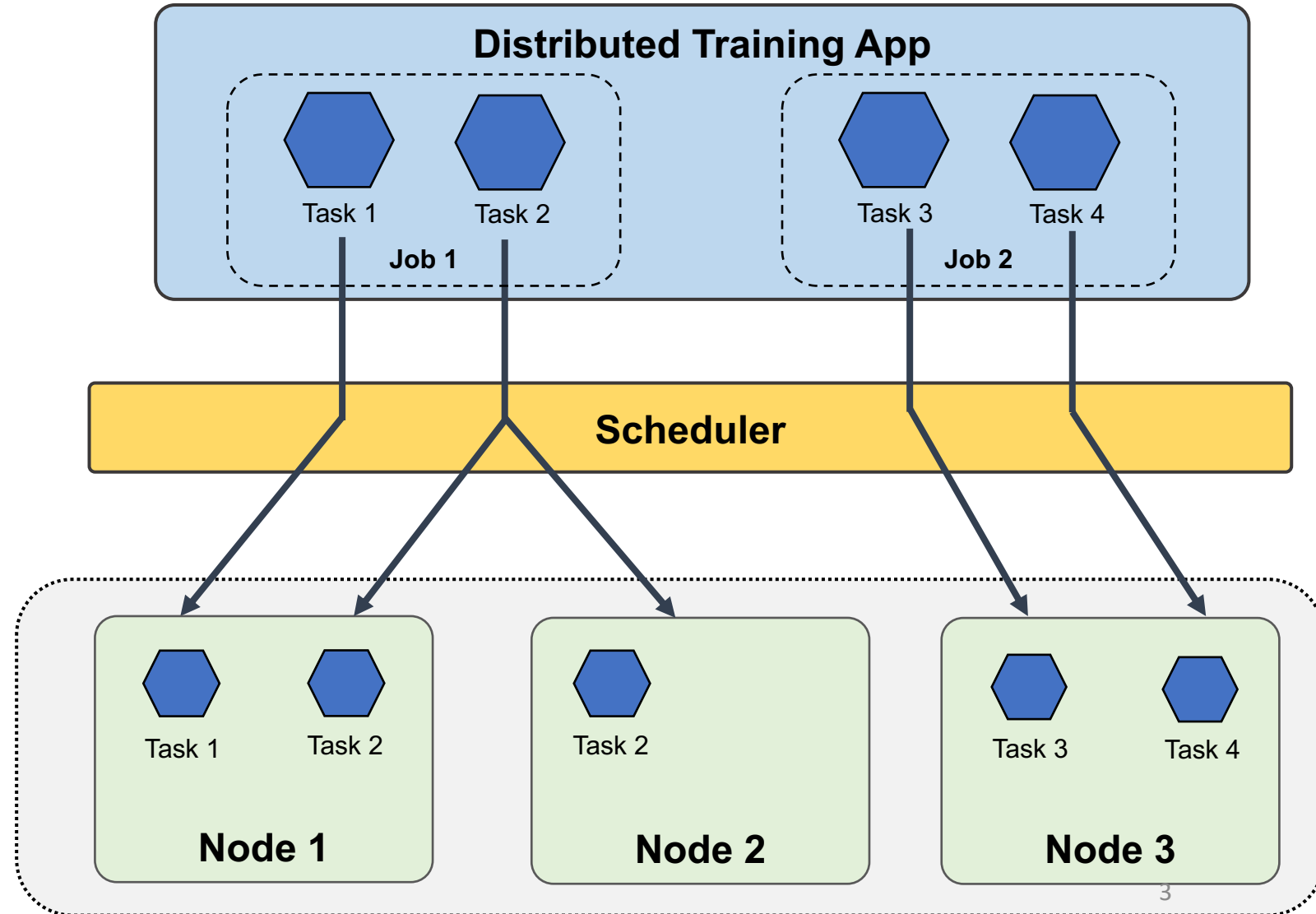**Scheduling Requirements**

**1. Gang Scheduling**
- Scheduling co-dependent tasks requires all-or-none semantics

**2. Co-location**
- Tasks of a job share parameter updates and must be placed on the same node for performance

**3. Anti-affinity**
- Avoid interference and resource contention by spreading jobs evenly spread across nodes
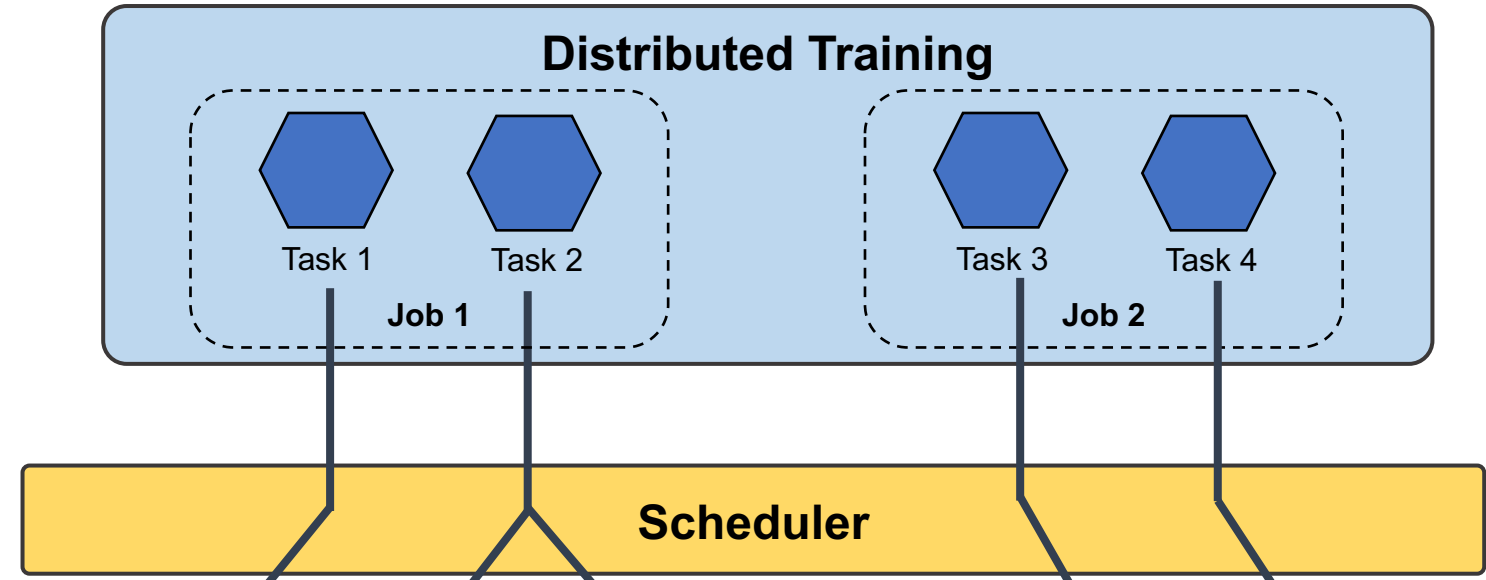
# Example - Distributed Training



**Scheduling Requirements**

**1. Gang Scheduling**
- Scheduling co-dependent tasks requires all-or-none semantics

**2. Co-location**
- Tasks of a job share parameter updates and must be placed on evenly spread across nodes

**Distributed Training**

Task 1  Task 2  **Job 1**

Task 3  Task 4  **Job 2**

**Scheduler**

Supporting custom scheduling constraints requires *evolvable* schedulers

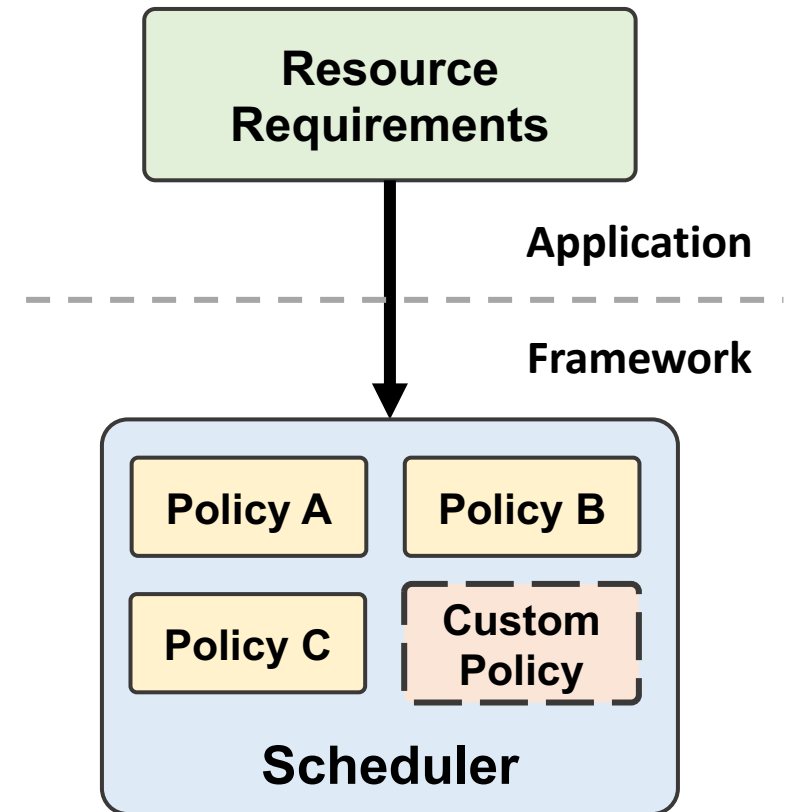Task 1  Task 2  **Node 1**

Task 2  **Node 2**

Task 3  Task 4  **Node 3**

# Evolvability in Monolithic Schedulers
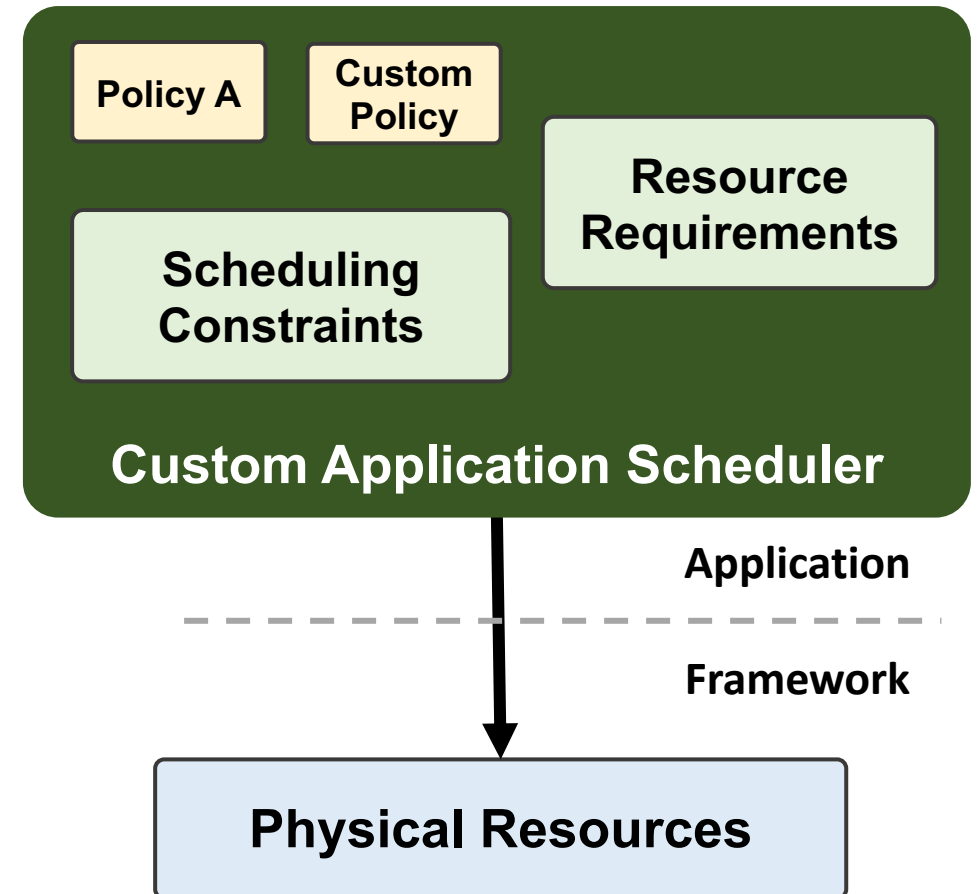
Kubernetes, YARN

- Applications state resource requirements
- Scheduler provides a fixed set of supported policies
  - E.g., Affinity, anti-affinity
- Challenging to evolve
  - Implementing custom policies requires modifying the core scheduler
  - Can take months to add support
  - Difficult to maintain - must commit to maintaining branch

# Evolvability in Two-level Schedulers

Mesos, Omega

- Physical resources are exposed to applications

- Applications implement end-to-end scheduling

- Highly flexible, but application must implement a scheduler:
  - Resource state tracking
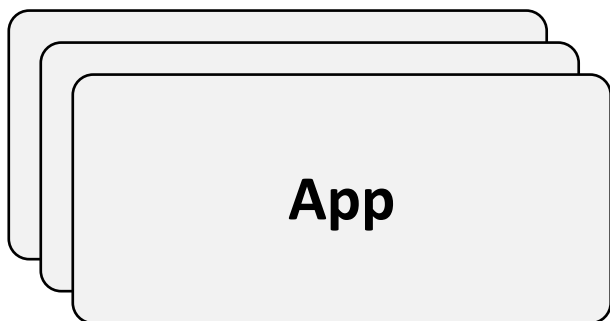  - Task queueing
  - Fault tolerance



Policy A | Custom Policy

Resource Requirements

Scheduling Constraints

**Custom Application Scheduler**

Application

Framework

**Physical Resources**

# Summary of solutions today

| Monolithic Schedulers | Two-level Schedulers | ESCHER |
|:---:|:---:|:---:|
| Simple, but hard to evolve | Highly evolvable, but complex | Simple and evolvable |



**Application Layer**

Monolithic: App

Two-level: App — Scheduling Policy + Scheduling Mechanism

ESCHER: App — Scheduling Policy

**Cluster Framework**

Monolithic: Scheduling Policy + Scheduling Mechanism

ESCHER: Scheduling Mechanism

# ESCHER Insights

**With the following two scheduling abstractions, frameworks can allow applications to express a wide range of scheduling policies:**
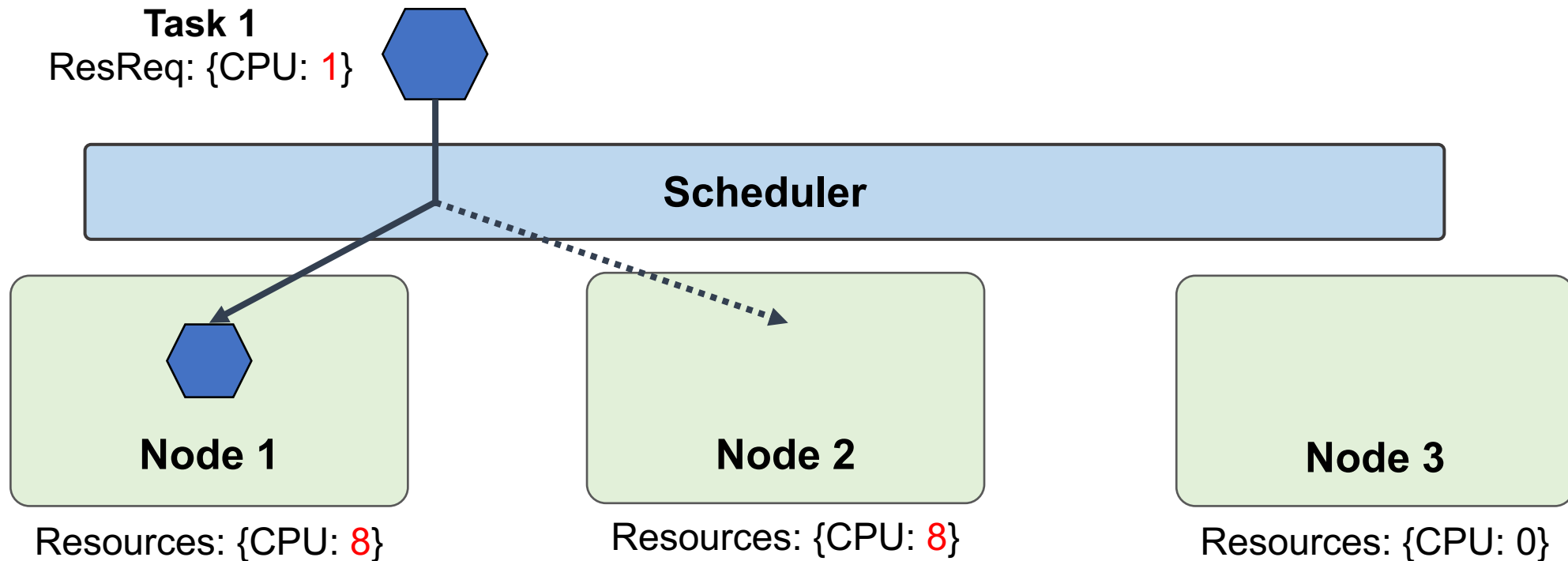
**1. A simple <span style="color:red">resource matching scheduler</span>**

**2. An API for applications to <span style="color:red">create resources at runtime</span>**

# Abstraction 1 - Resource Matching Scheduler



Scheduler matches tasks resource **requirements** to node resource **availabilities**

**Task 1**
ResReq: {CPU: 1}

**Scheduler**

**Node 1**

Resources: {CPU: 8}

**Node 2**

Resources: {CPU: 8}

**Node 3**

Resources: {CPU: 0}

# Abstraction 2 – Create Resources on-the-fly

**Applications should be able to <span style="color:red">create resources and get cluster state</span> at runtime through an API**
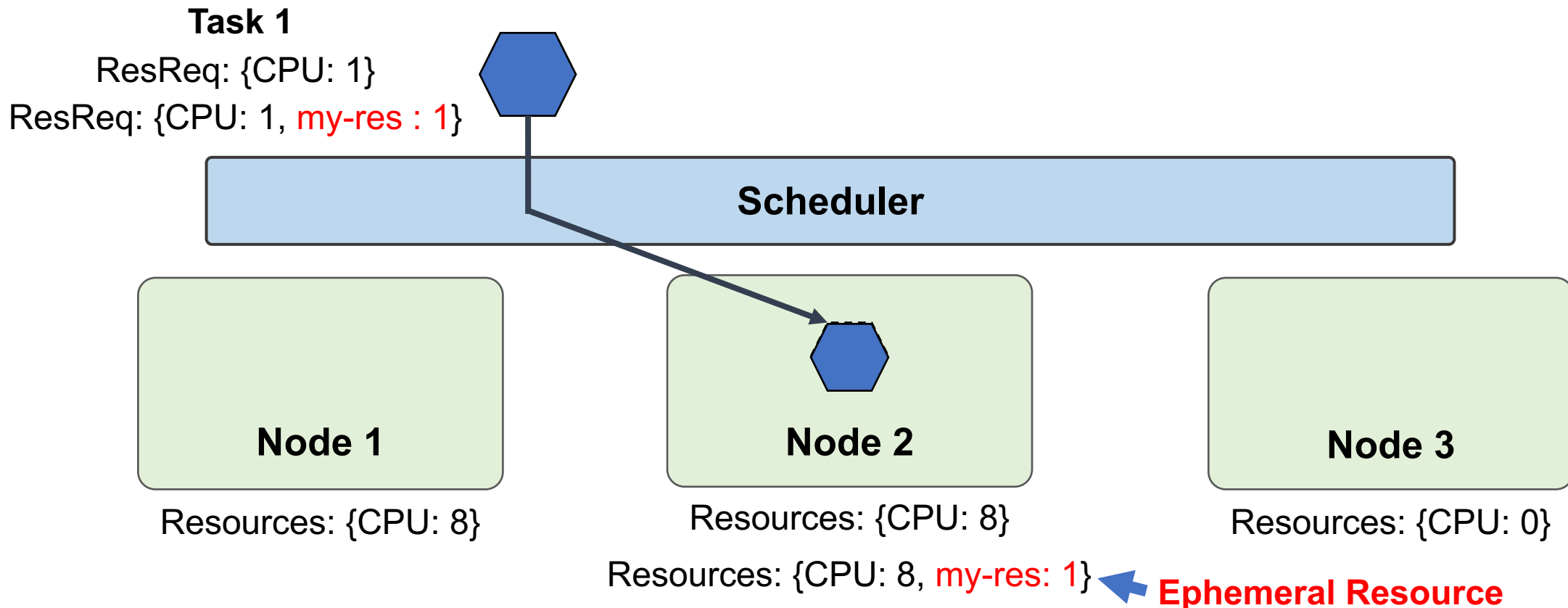
```python
def set_resource(resource_name, capacity, node_spec=None)
```

```python
def get_cluster_state()       # Returns a map of {node: resources}
```

- Can specify resource availability constraints for resource creation

- If not node_spec not specified, resource created locally
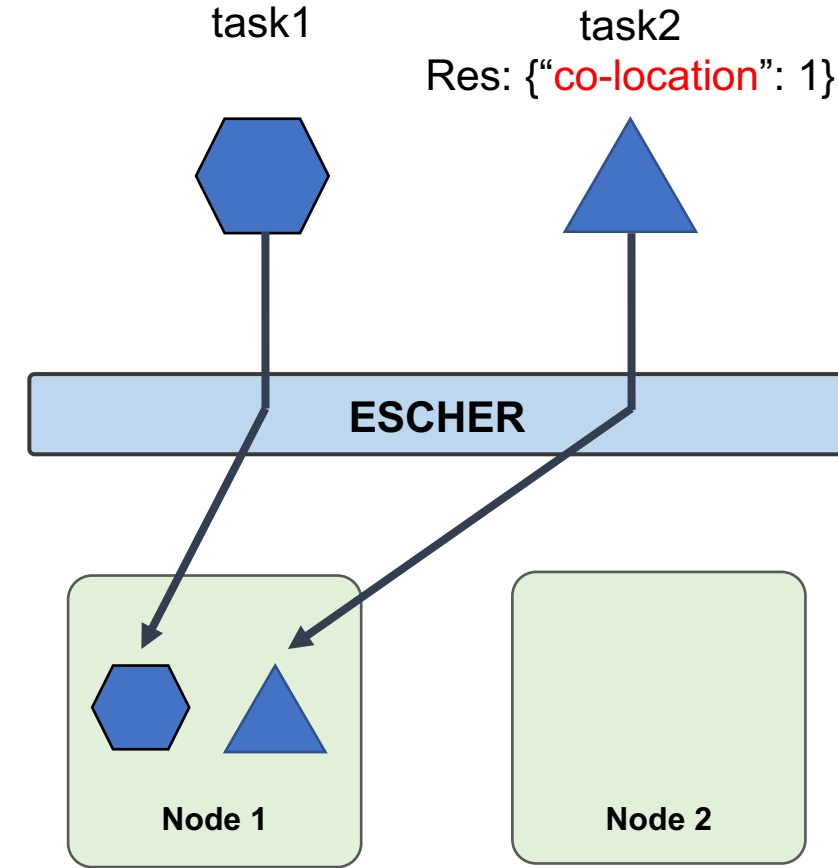
# Scheduling with Ephemeral Resources

A simple resource matching scheduler can be induced to make targeted placement decisions with short-lived *ephemeral* resources

**Task 1**

ResReq: {CPU: 1}

ResReq: {CPU: 1, my-res : 1}

**Scheduler**

**Node 1**

**Node 2**

**Node 3**

Resources: {CPU: 8}

Resources: {CPU: 8}

Resources: {CPU: 0}

Resources: {CPU: 8, my-res: 1}  **Ephemeral Resource**

# Example - Task co-location

Run tasks on the same node

```python
def task1():
    set_resource(label="co_location", capacity=1)
    ...

def task2():
    ...

def main():
    launch(task1, res = {})
    launch(task2, res = {"co_location": 1})
```
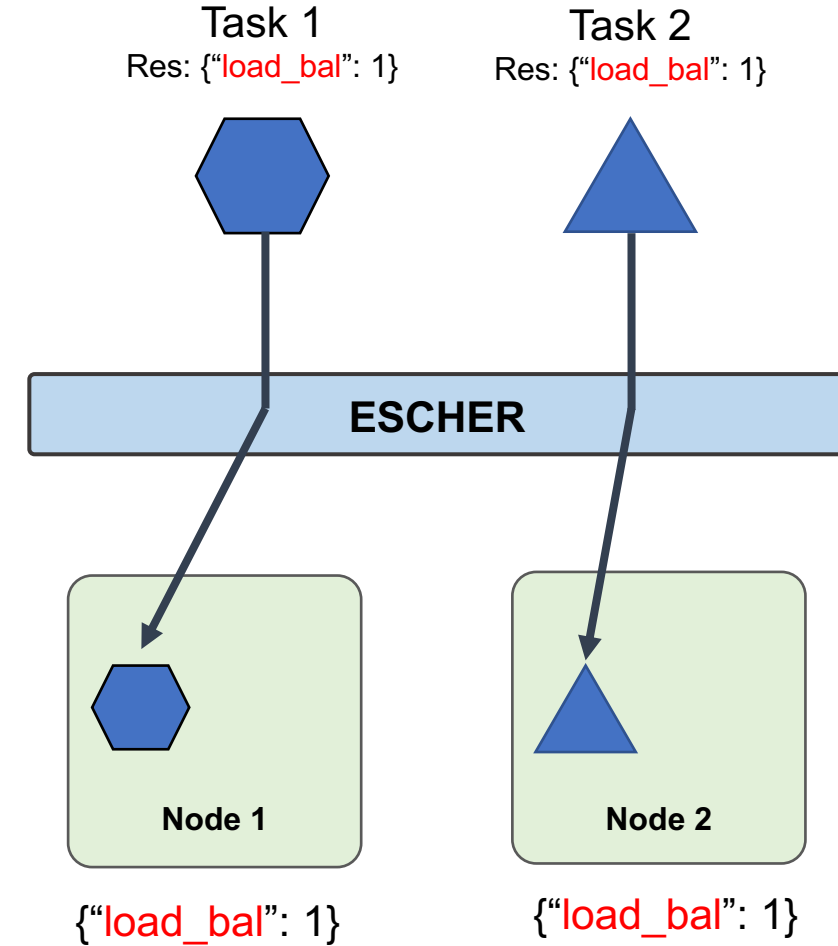
task1

task2
Res: {"co-location": 1}

**ESCHER**

Node 1

Node 2

{"co-location": 1}

**ESCHER allows declarative specification of scheduling policies by dynamically creating ephemeral resources**

# Example - Load Balancing

Spread tasks across machines

Task 1
Res: {"load_bal": 1}

Task 2
Res: {"load_bal": 1}

ESCHER

```python
def static_load_balancing(num_tasks, num_nodes):
    resource_capacity = ceiling(num_tasks/num_nodes)
    set_resource(label="load_bal", node_spec={},
                 capacity=resource_capacity)
    for task in tasks:
        task.resources = {"load_bal": 1}
        task.launch()
```
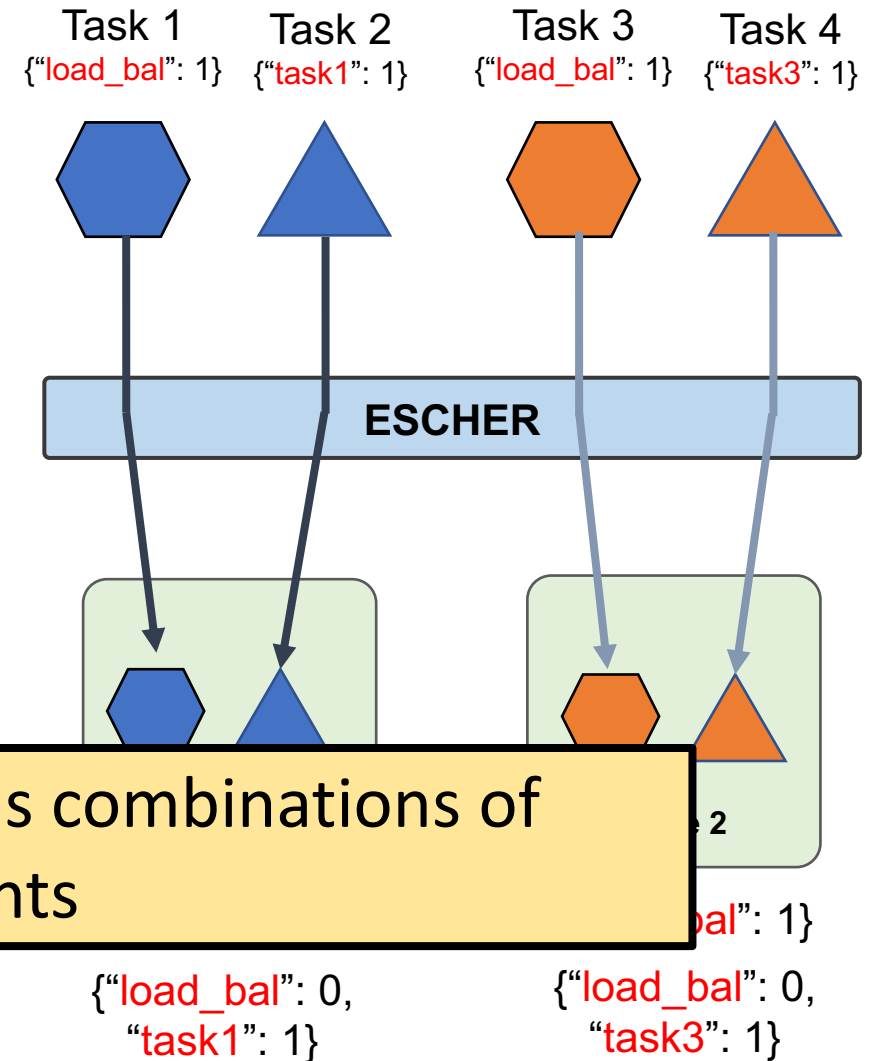
Node 1

Node 2

{"load_bal": 1}

{"load_bal": 1}

0}

Many more ESCHER policies (gang scheduling, bin-packing, anti-affinity, soft constraints, priorities, hierarchical fair sharing) in the paper!
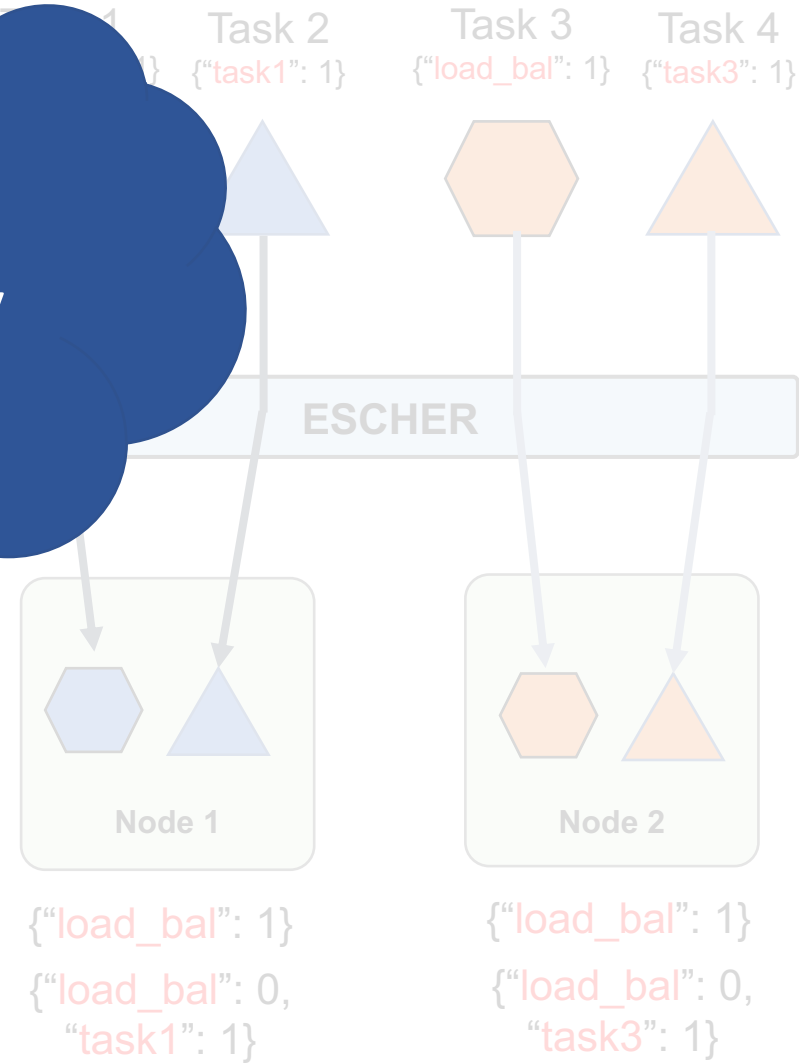
# Policy Composition: Load Balancing & Co-location

Co-locate two tasks and spread out pairs of tasks

```python
def task1(id):
    set_resource(label=id, node=None, capacity=1)
    ...


def task2():
    ...


def main():
    # Create load-balancing resources
    set_resource(label="load_bal", capacity=1, node_spec={})

    # Launch tasks
```

Task 1
{"load_bal": 1}

Task 2
{"task1": 1}

Task 3
{"load_bal": 1}

Task 4
{"task3": 1}

**ESCHER**

Compositions of policy can be represented as combinations of ephemeral resource constraints

{"load_bal": 0,
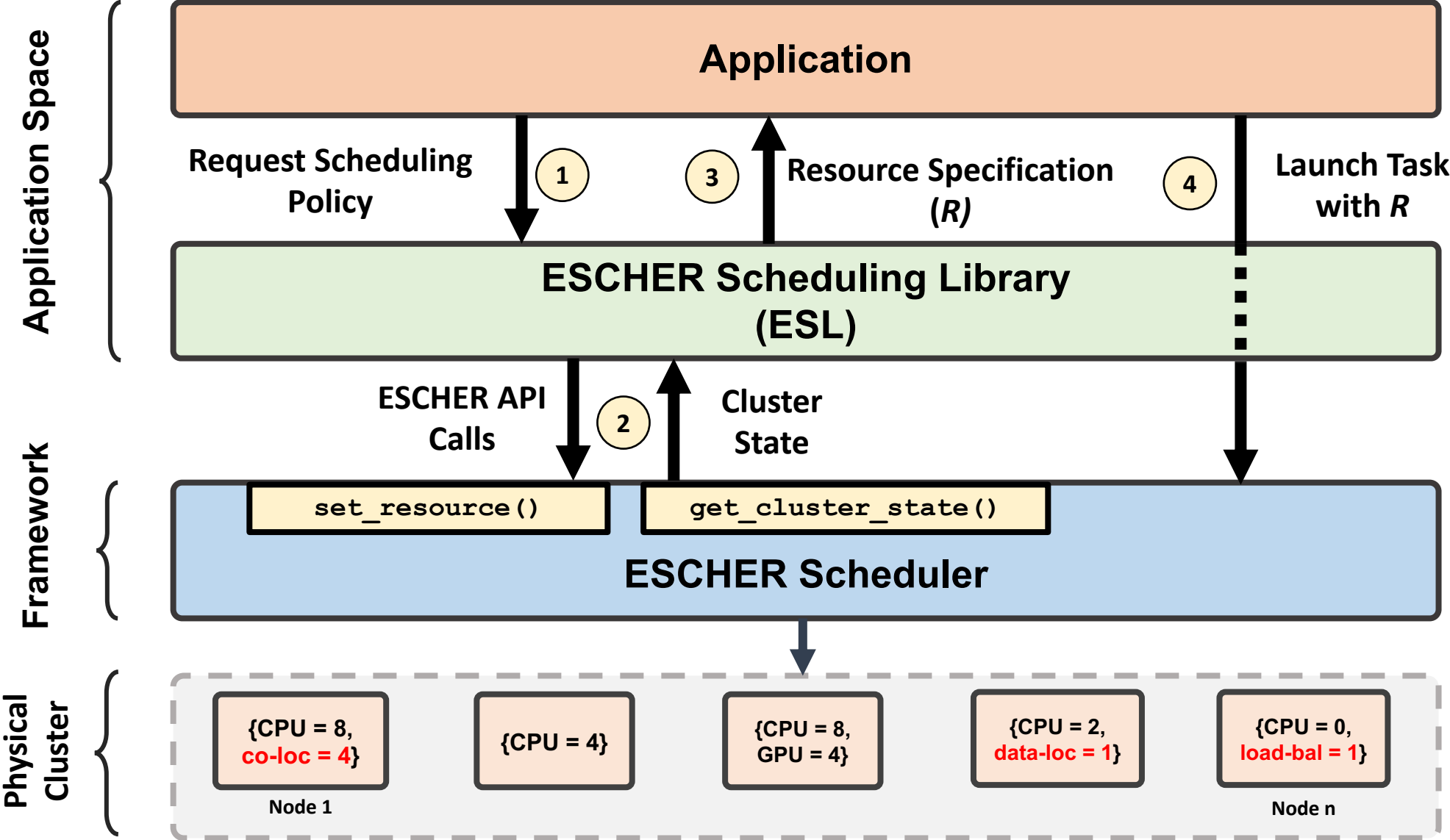"task1": 1}

{"load_bal": 0,
"task3": 1}

# ESCHER Scheduling Libraries (ESLs)

- An app-level library of scheduling policies which encapsulate all state management for ephemeral resources

- Encourage code-reuse and simplify application code

```python
def colocated_task():
    ...


def main():
    esl = CoLocationESL()
    coloc_res = esl.get_colocation_group("mygroup", res_req={gpu: 8})
    launch(colocated_tasks, res += coloc_res)
```
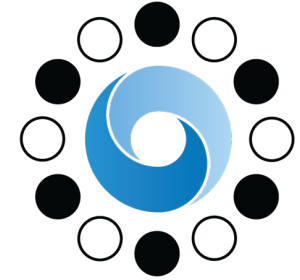
# ESCHER Workflow

# Implementation



Modified the Ray Scheduler to support online resource updates

No changes required in Kubernetes core – we reuse the extended resources API
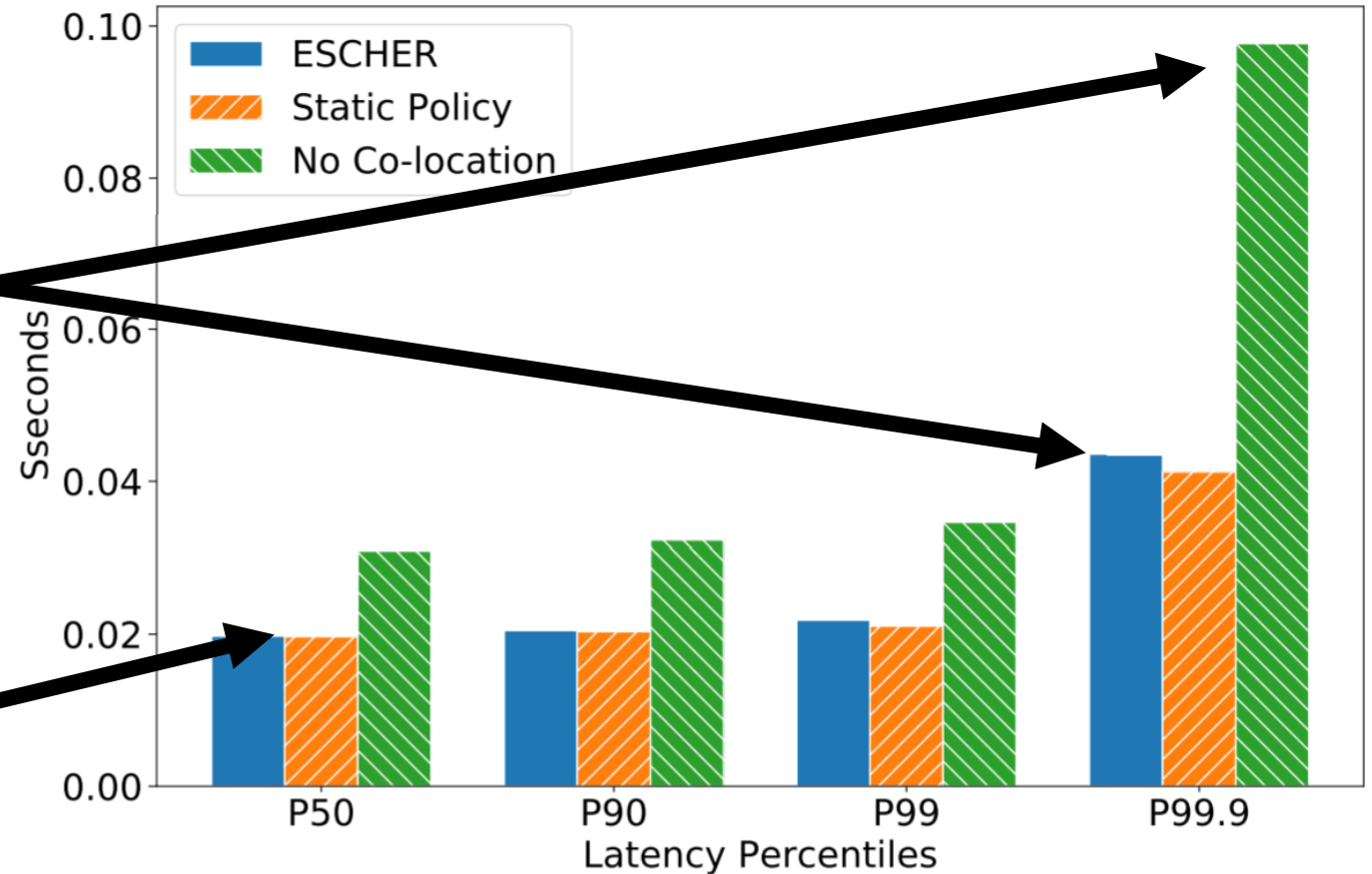
18

# Evaluation - AlphaZero

- AlphaZero trains an RL agent to play Go

- Training has two key processes:

  - **Board Generation**: CPU intensive generation of possible game states

  - **Board Evaluation**: A GPU agent predicts the "goodness" of the generated states and chooses an action

- These processes require both **co-location** and **load-balancing**

# AlphaZero on Ray

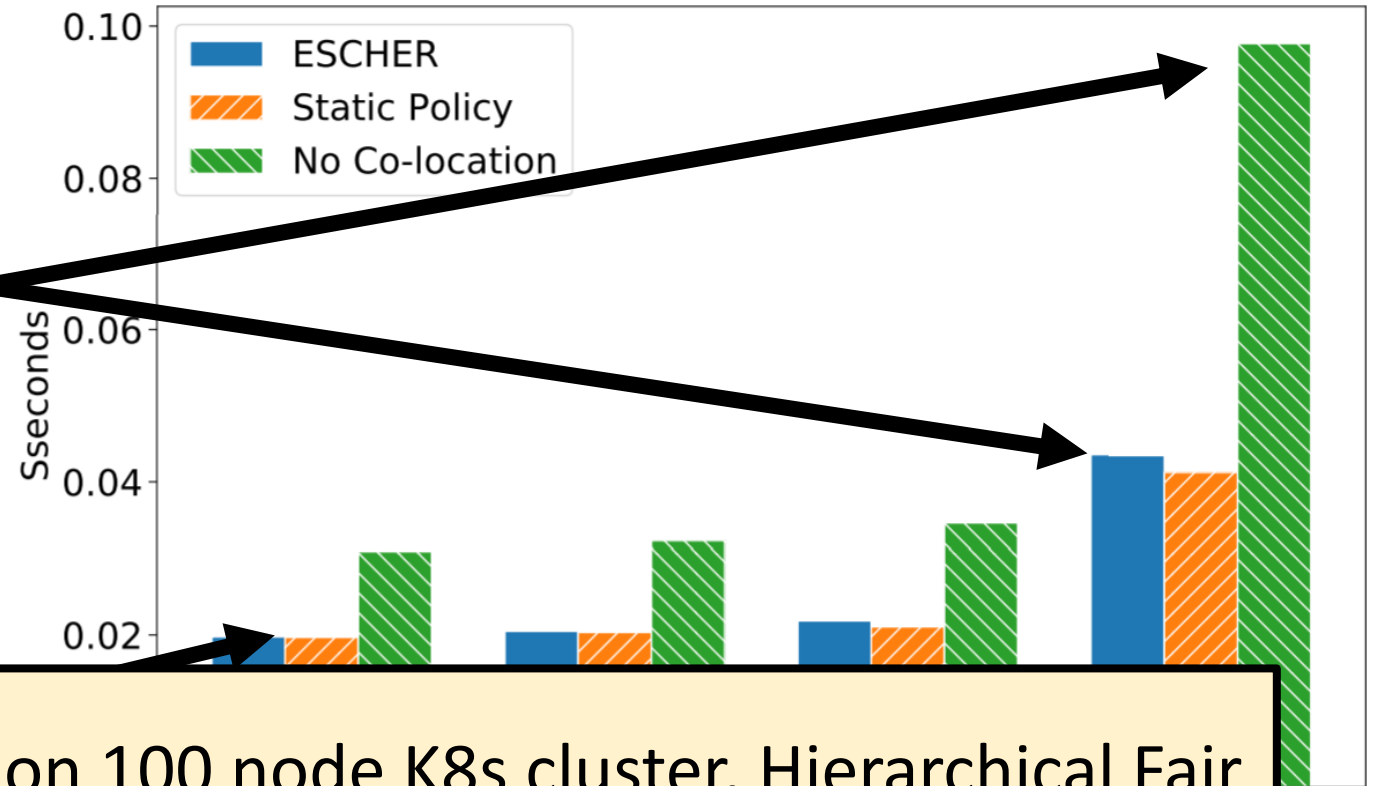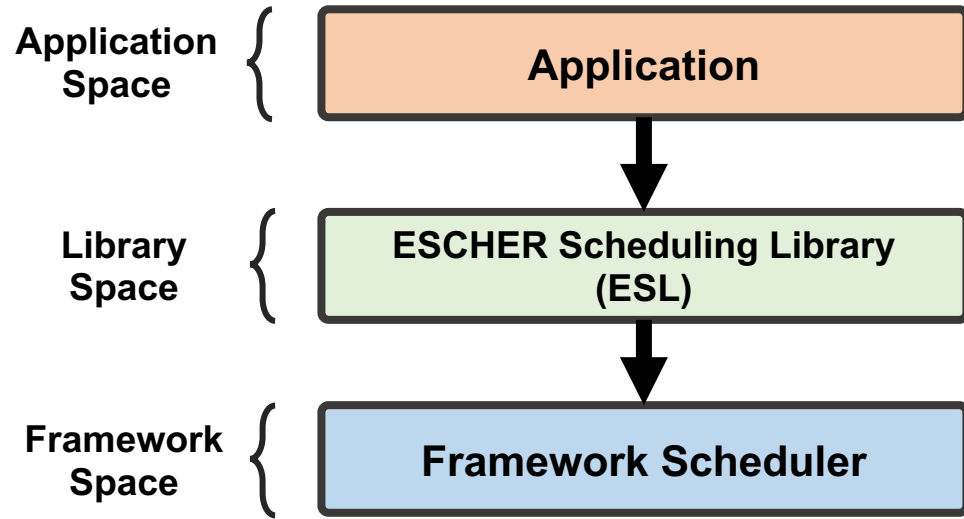ESCHER is **1.5-2x faster** in exploring board states than a locality-unaware scheduler

ESCHER performs **comparably with a static hard-coded policy** with just **5 lines of code changes**



**Board Exploration Latencies - 128 GPUs**

Legend:
- ESCHER
- Static Policy
- No Co-location

Y-axis: Sseconds (0.00, 0.02, 0.04, 0.06, 0.08, 0.10)
X-axis: Latency Percentiles (P50, P90, P99, P99.9)

# AlphaZero on Ray

**Board Exploration Latencies - 128 GPUs**

ESCHER is **1.5-2x faster** in exploring board states than a locality-unaware scheduler

Legend:
- ESCHER
- Static Policy
- No Co-location

Y-axis: Sseconds (0.02, 0.04, 0.06, 0.08, 0.10)
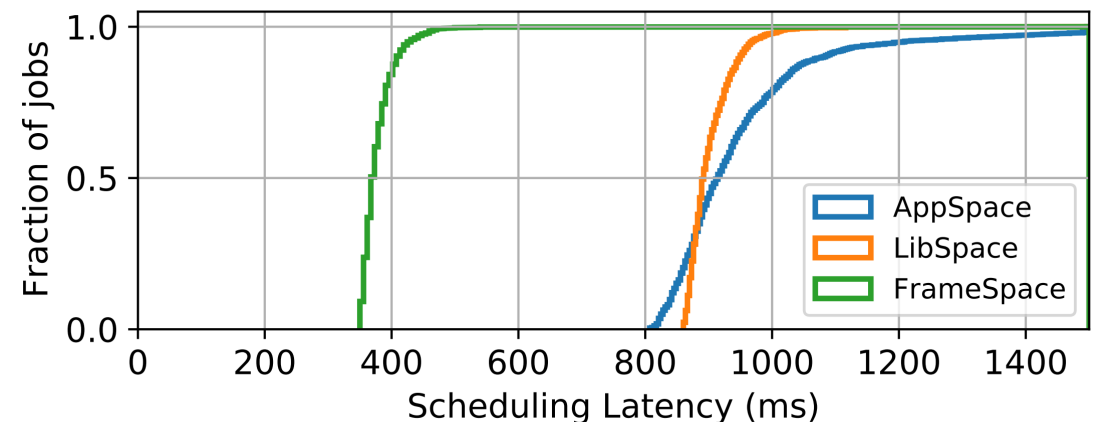
ESCHER performs **comparably**

More results (MapReduce on 100 node K8s cluster, Hierarchical Fair Sharing, Distributed Training, Microbenchmarks) in the paper!

# ESCHER Overheads vs Evolvability



Application Space { Application

ESCHER Scheduling Library (ESL)

Framework Space { Framework Scheduler

Library Space {

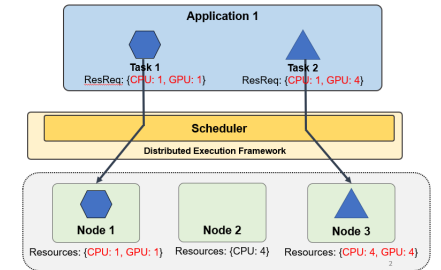| Gang Scheduling Implementation | Lines of Code | Median Scheduling Latency |
|---|---|---|
| **AppSpace** (ESCHER) | | |
| **LibSpace** (ESLs in ESCHER) | | |
| **FrameSpace** (Monolithic Scheduler) | | |

Using ESCHER adds latency for some policies such as gang scheduling, but significantly reduces the implementation burden.
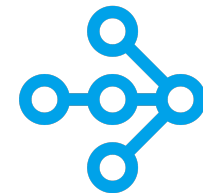
# ESCHER Summary

- Applications need **fine-grained scheduling** control without the complexity of implementing scheduling mechanisms.

- ESCHER presents an evolvable scheduler architecture with two key abstractions – a **resource matching scheduler** and **set_resource API**

- Ephemeral resources **are easily implemented in Ray and Kubernetes** and provide scheduling flexibility for a range of workloads with minimal overhead.



```
def set_resource()
```