**SJTU SAIL** (Software Architecture and Infrastructure Lab)

**SHANGHAI JIAO TONG UNIVERSITY**

# Insight of Medea and Neptune

Scheduling of Long Running Applications  &
Scheduling Suspendible Tasks for Unified Stream/Batch Applications

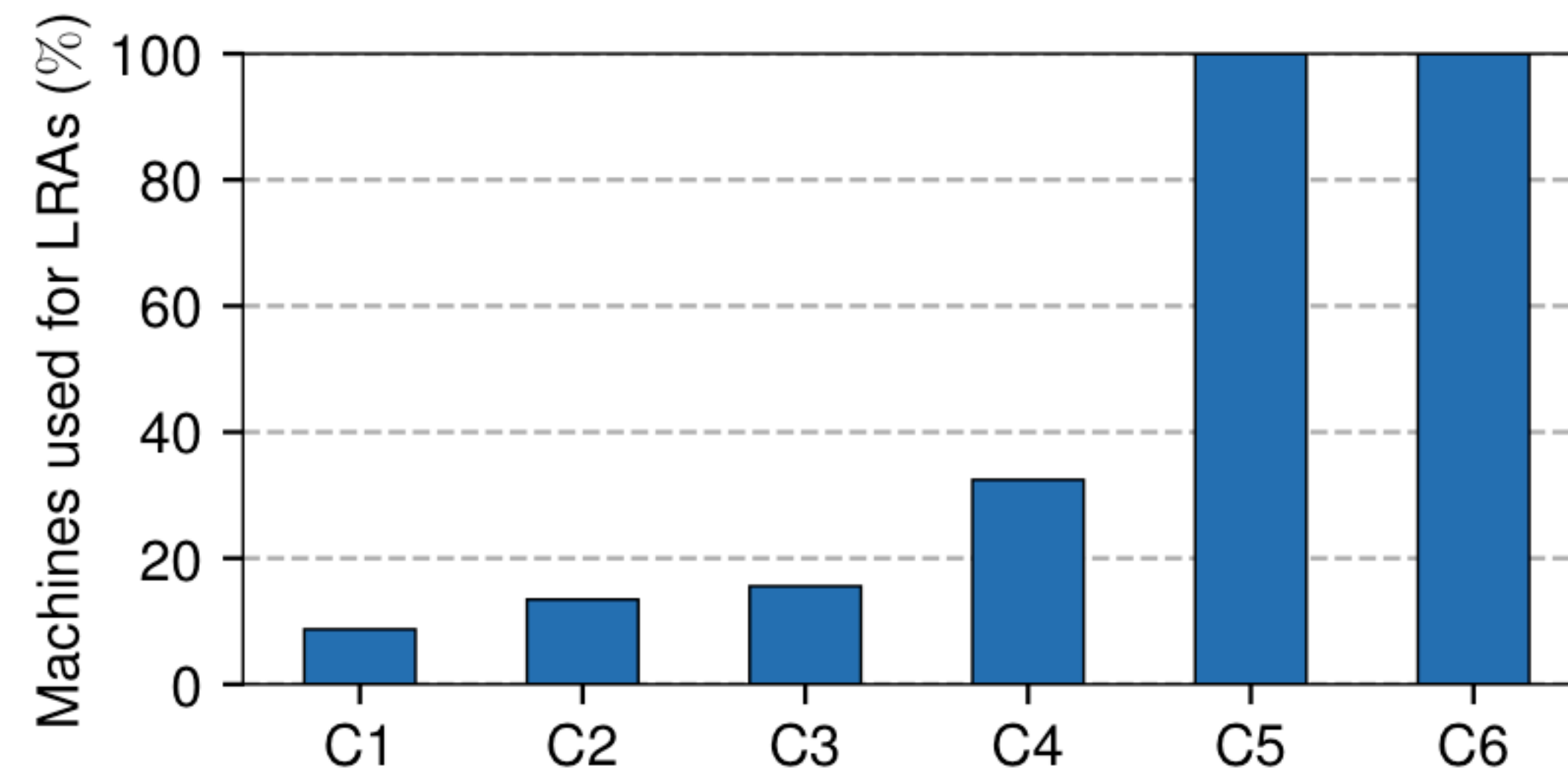Lin Hsu Yalun    2020/09/21

# Medea

## Scheduling of Long Running Applications in Shared Production Clusters

# Long-Running Applications (LRAs)

- **Interactive data-intensive** applications

  - Spark, Hive LLAP

- **Streaming** systems

  - Flink, Storm, SEEP

- **Latency-sensitive** applications

  - HBase, Memcached

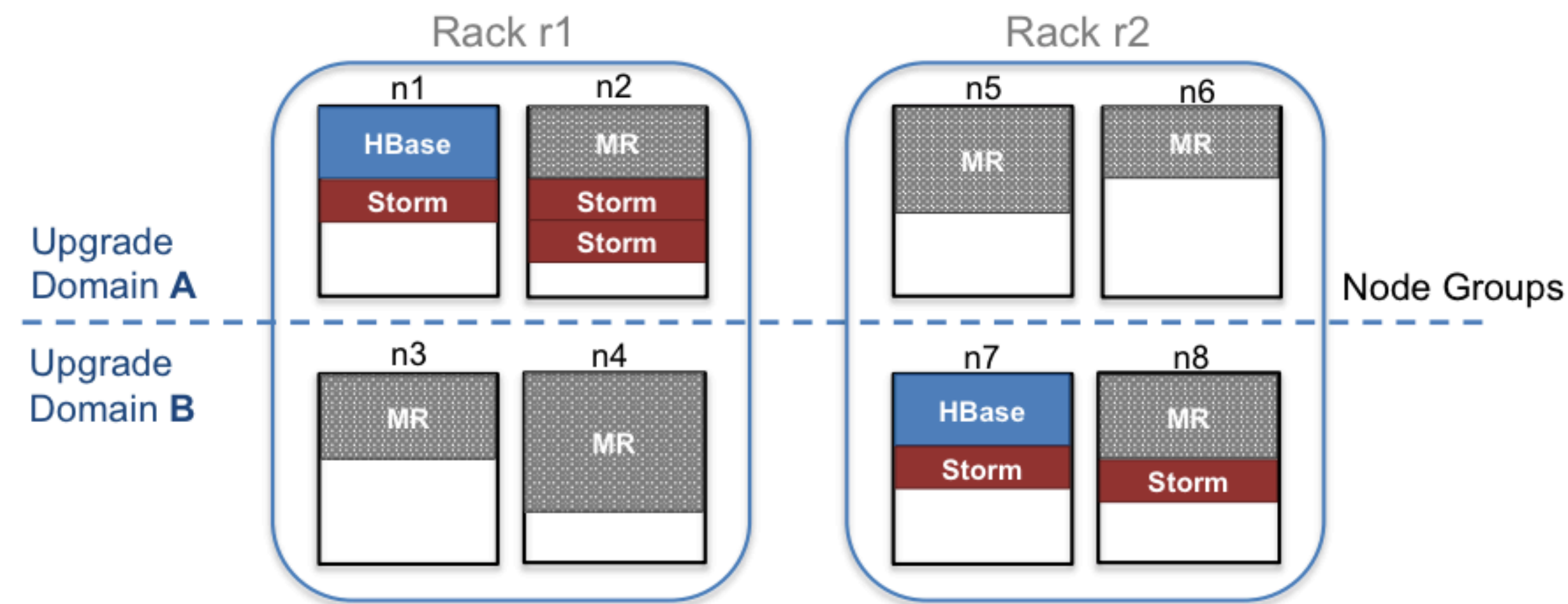- **ML** frameworks

  - TensorFlow, Spark ML-lib

**LRAs** = applications
with long-running containers
(running from hours to months)
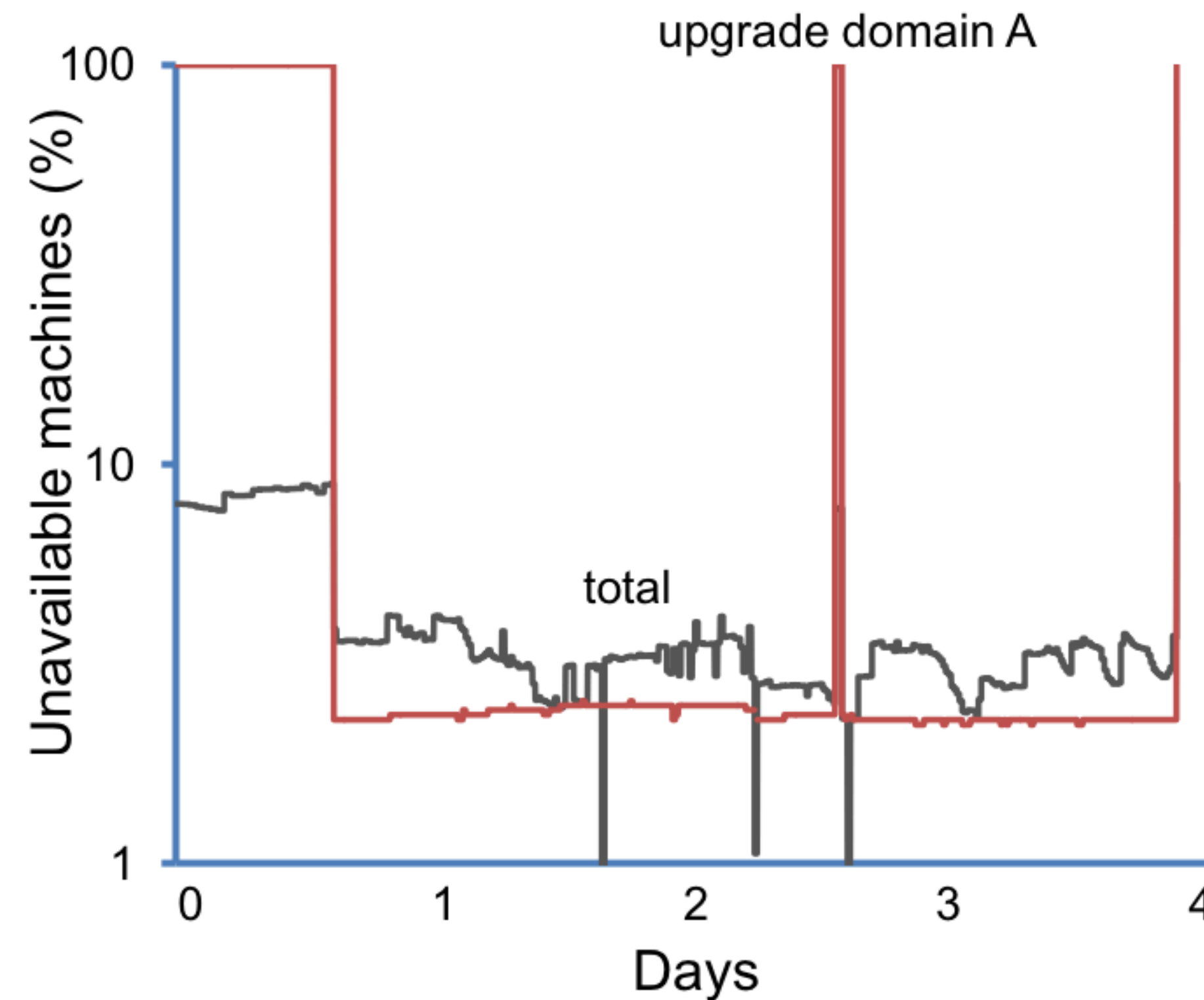
# LRAs in Microsoft's analytics clusters



**LRA placement is important**

# LRA scheduling problem



- **Performance**: "Place Storm containers in the same rack as HBase"

- **Cluster objectives**: "Minimize resource fragmentation"

- **Resilience**: "Place HBase containers across upgrade domains"

# Machine unavailability in a Microsoft cluster



With random placement, an LRA might lose all containers at once

# Challenges

- How to **relate containers** to node groups?

- How to **express** different types of constraints related to LRA containers?

- How to achieve **high quality placement** without affecting task-based jobs?

# Medea

- How to **relate containers** to node groups?

- Support container tags and logical node groups

- How to **express** different types of constraints related to LRA containers?

- Introduce expressive cardinality constraints

- How to achieve **high quality placement** without affecting task-based jobs?
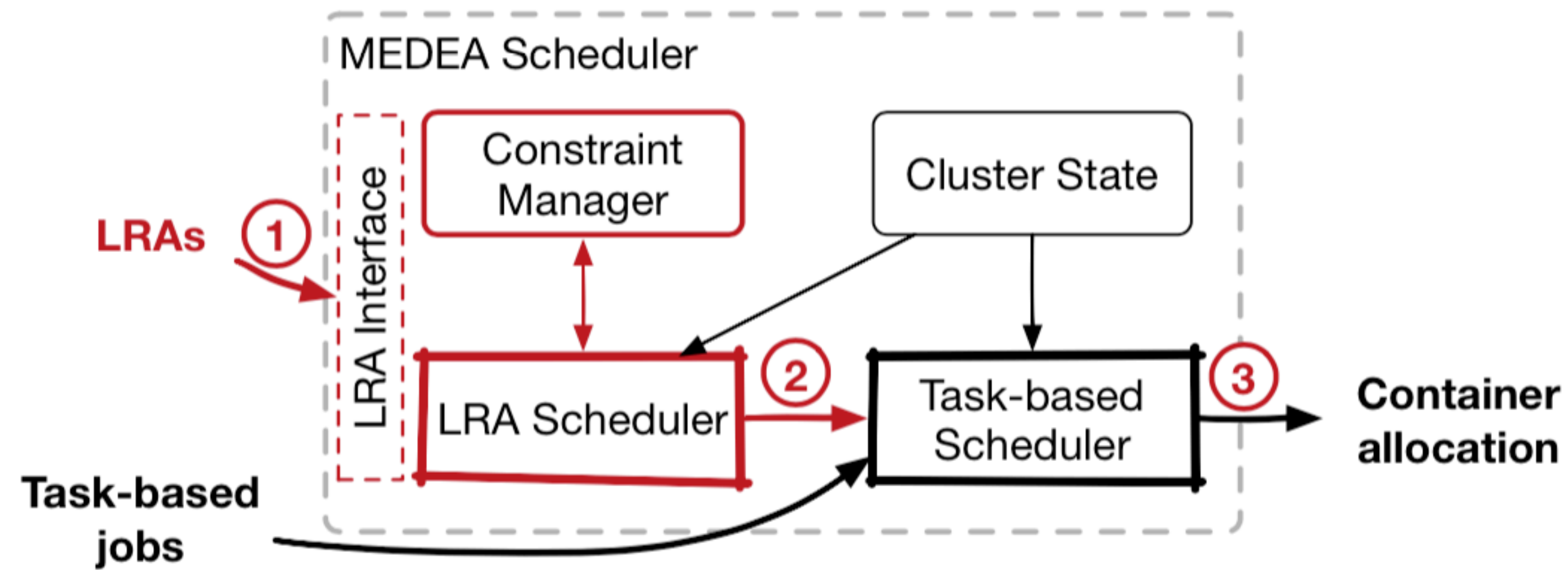
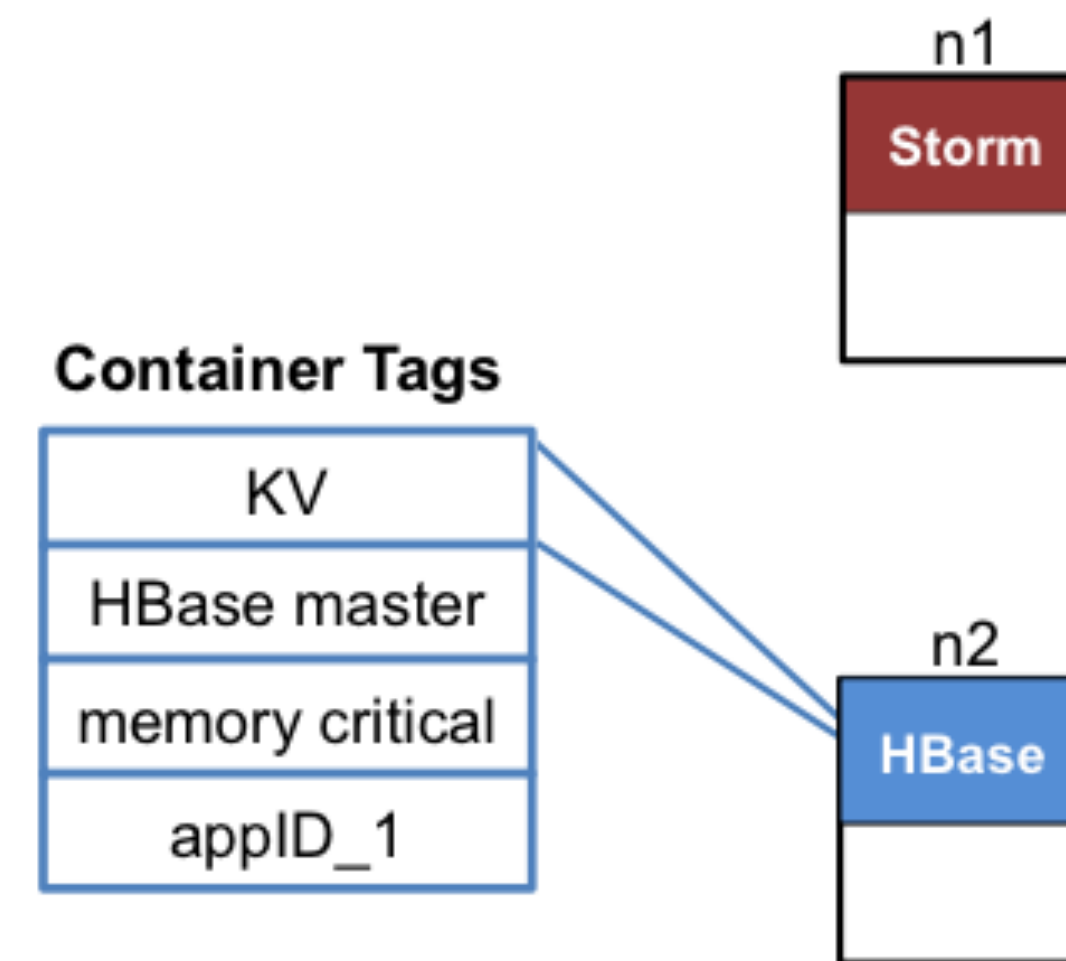- Follow a two-scheduler design

# Medea Design



Figure 4: MEDEA scheduler design

# Container tagging

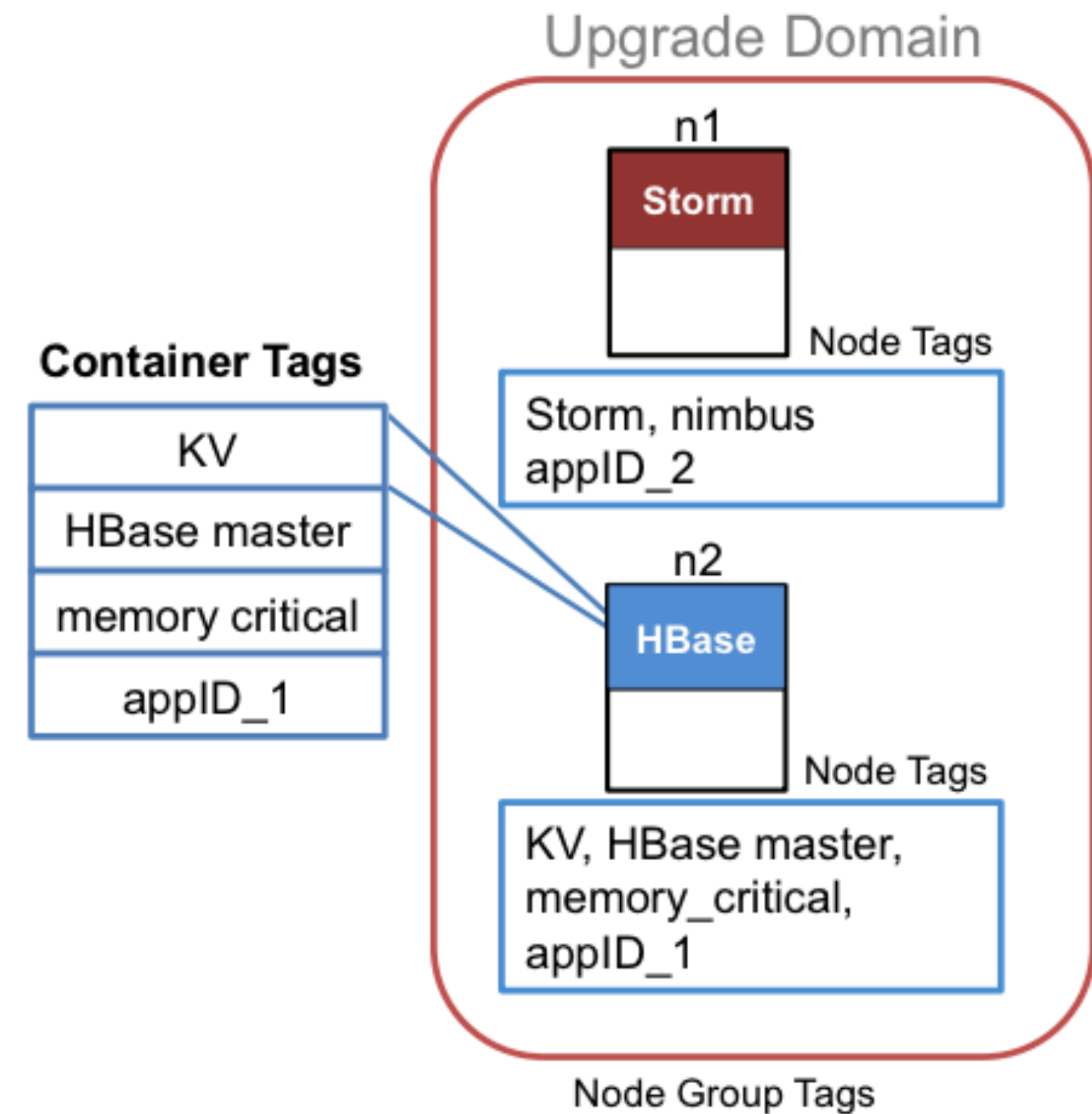- **Idea**: use container tags to refer to group of containers

  - Describe
    - application type
    - application role
    - resource specification
    - global application ID

**Container Tags**

| |
|---|
| KV |
| HBase master |
| memory critical |
| appID_1 |

n1
Storm

n2
HBase

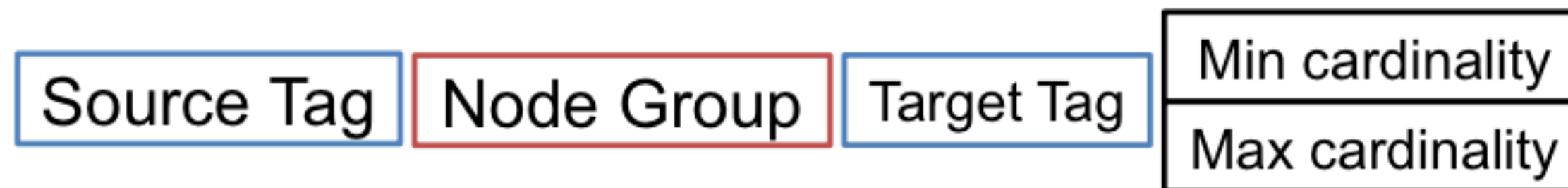- Can refer to any current or future LRA container

# Hierarchical grouping of nodes

- **Idea**: logical node groups to refer to dynamic node sets

  - E.g. node, rack. Upgrade domain

  - Associate nodes with all the container tags that live there
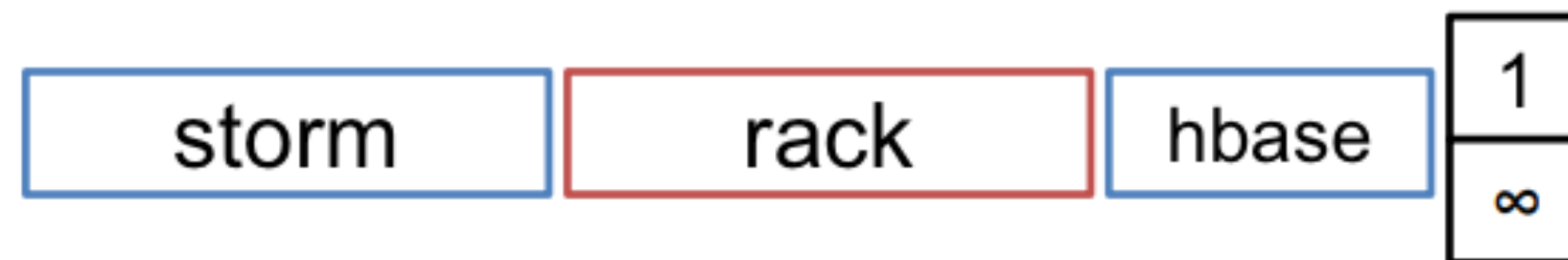
  - Hide infrastructure "spread across upgrade domains"

# Defining constraints

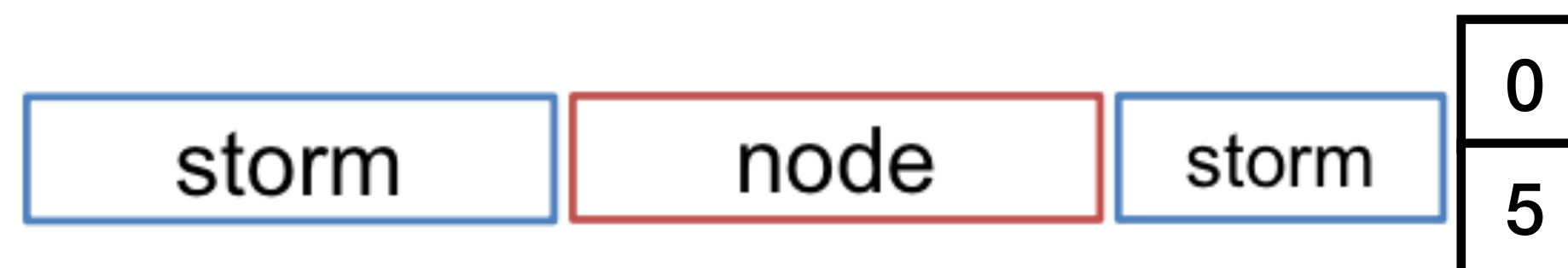- Generic constraints to capture a variety of cases

| Source Tag | Node Group | Target Tag | Min cardinality |
|---|---|---|---|
| | | | Max cardinality |

  - Min cardinality ≤ occurrences (Target Tag) ≤ Max cardinality

- Affinity "Place Storm containers in the same rack as HBase"

| storm | rack | hbase | 1 |
|---|---|---|---|
| | | | ∞ |

- Cardinality "Place up to 5 Storm containers in the same node"

| storm | node | storm | 0 |
|---|---|---|---|
| | | | 5 |

# Two-scheduler design

# ILP-based scheduling algorithm

$$\text{maximize} \quad \frac{w_1}{k} \sum_{i=1}^{k} S_i + \frac{w_2}{m} \sum_{l=1}^{m} v_c^l + \frac{w_3}{N} \sum_{n=1}^{N} z_n \tag{1}$$

subject to:

$$\forall i, j : \sum_{n=1}^{N} X_{ijn} \leq 1 \tag{2}$$

$$\forall n : \sum_{i=1}^{k} \sum_{j=1}^{T_i} r_{ij} \cdot X_{ijn} \leq R_n^f \tag{3}$$

$$\forall i : \sum_{n=1}^{N} \sum_{j=1}^{T_i} X_{ijn} - T_i S_i = 0 \tag{4}$$

$$\forall n : \sum_{i=1}^{k} \sum_{j=1}^{T_i} r_{ij} \cdot X_{ijn} - B_n(1 - z_n) \leq R_n^f - r_{min} \tag{5}$$

For each constraint $C_l = \{\texttt{s\_tag}, \{\texttt{c\_tag}, c_{min}^l, c_{max}^l\}, \texttt{G}\}$,
$\forall$ container $t_{i_s j_s} \in \texttt{s\_tag}$, $\forall$ node set $\mathcal{S} \in \texttt{G}$:

$$\sum_{n \in \mathcal{S}} \left( \sum_{\substack{i,j: \texttt{tag} \in t_{ij} \\ t_{ij} \neq t_{i_s j_s}}} X_{ijn} + D_n(1 - X_{i_s j_s n}) \right) - c_{min}^l + c_{min}^{l,v} \geq 0 \tag{6}$$
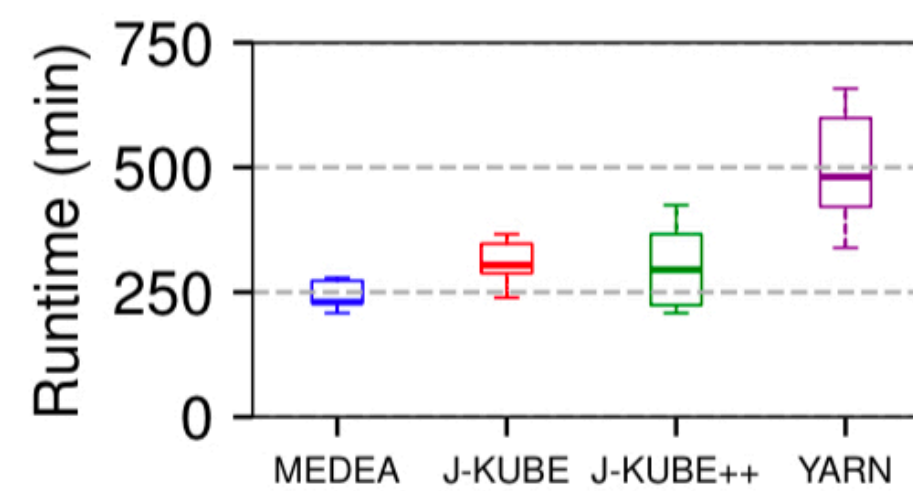
$$\sum_{n \in \mathcal{S}} \left( \sum_{\substack{i,j: \texttt{tag} \in t_{ij} \\ t_{ij} \neq t_{i_s j_s}}} X_{ijn} - D_n(1 - X_{i_s j_s n}) \right) - c_{max}^l - c_{max}^{l,v} \geq 0 \tag{7}$$

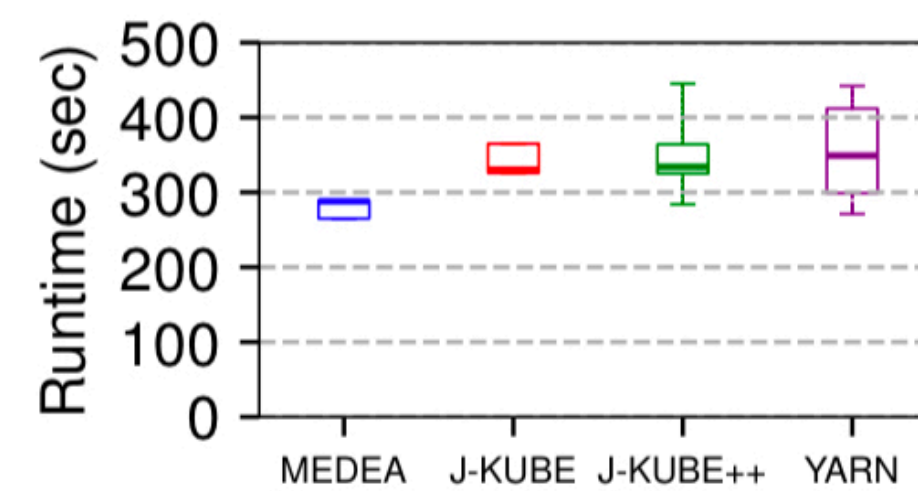$$v_c^l = \frac{c_{min}^{l,v}}{c_{min}^l} + \frac{c_{max}^{l,v}}{c_{max}^l} \tag{8}$$

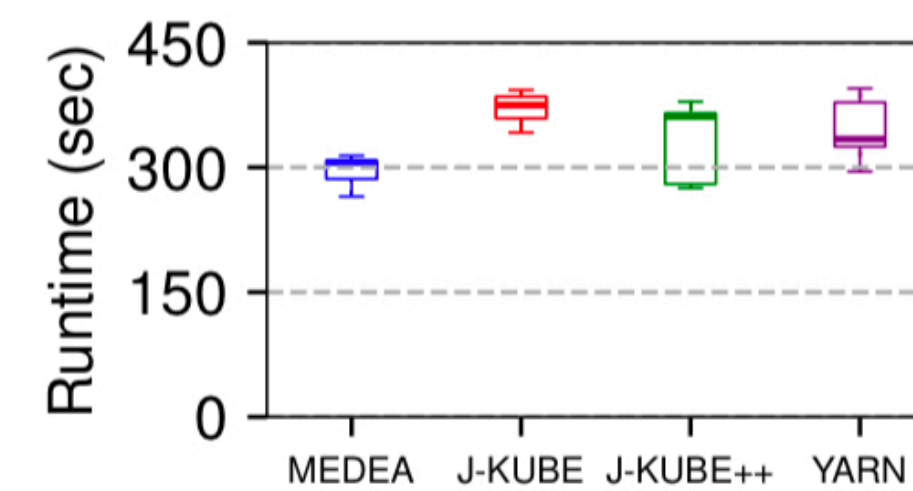| Symbol | Description |
|---|---|
| $k$ | Number of LRAs to be placed |
| $N$ | Number of cluster nodes |
| $T_i$ | Number of containers of LRA $i$ |
| $R_n^f, R_n^u$ | Free, used resources of node $n$[6] |
| $m$ | Total number of constraints |
| $w_i$ | Weights of components in objective function |
| $B_n, D_n$ | Sufficiently large integers, used in inequalities |
| $S_i$ | 1 if all containers of LRA $i$ are placed; 0 otherwise |
| $X_{ijn}$ | 1 if container $j$ of LRA $i$ placed at node $n$; 0 otherwise |
| $r_{ij}$ | Resource demand of container $j$ of LRA $i$ |
| $r_{min}$ | Minimum resource demand |
| $c_{min}^{l,v}, c_{max}^{l,v}$ | Violation of cardinalities $c_{min}, c_{max}$ for constraint $C_l$ |
| $v_c^l$ | Violation for constraint $C_l$ |
| $z_n$ | 1 if free resources $\geq r_{min}$ after placement; 0 otherwise |

# Evaluation

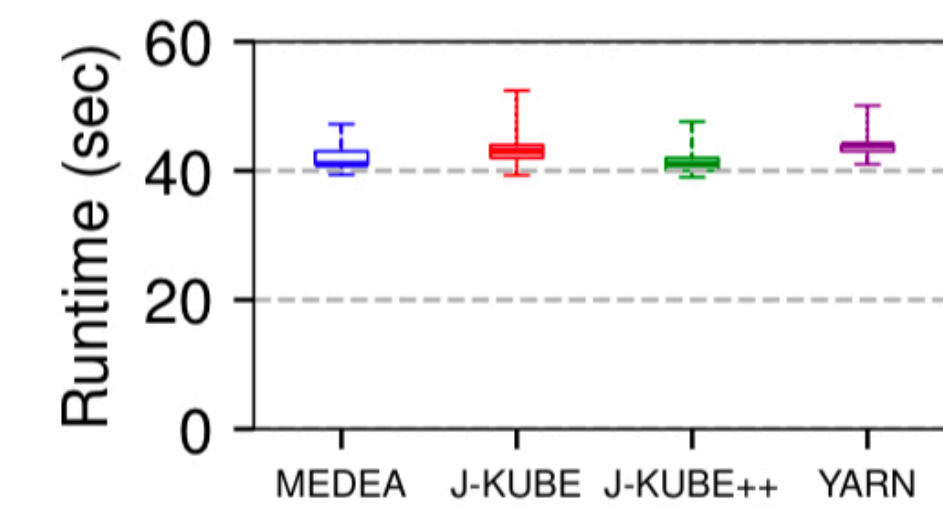**A 400 node pre-production cluster grouped into 10 racks, supplemented by simulation.**



Figure 7: **Application performance** (lower is better)

# Evaluation

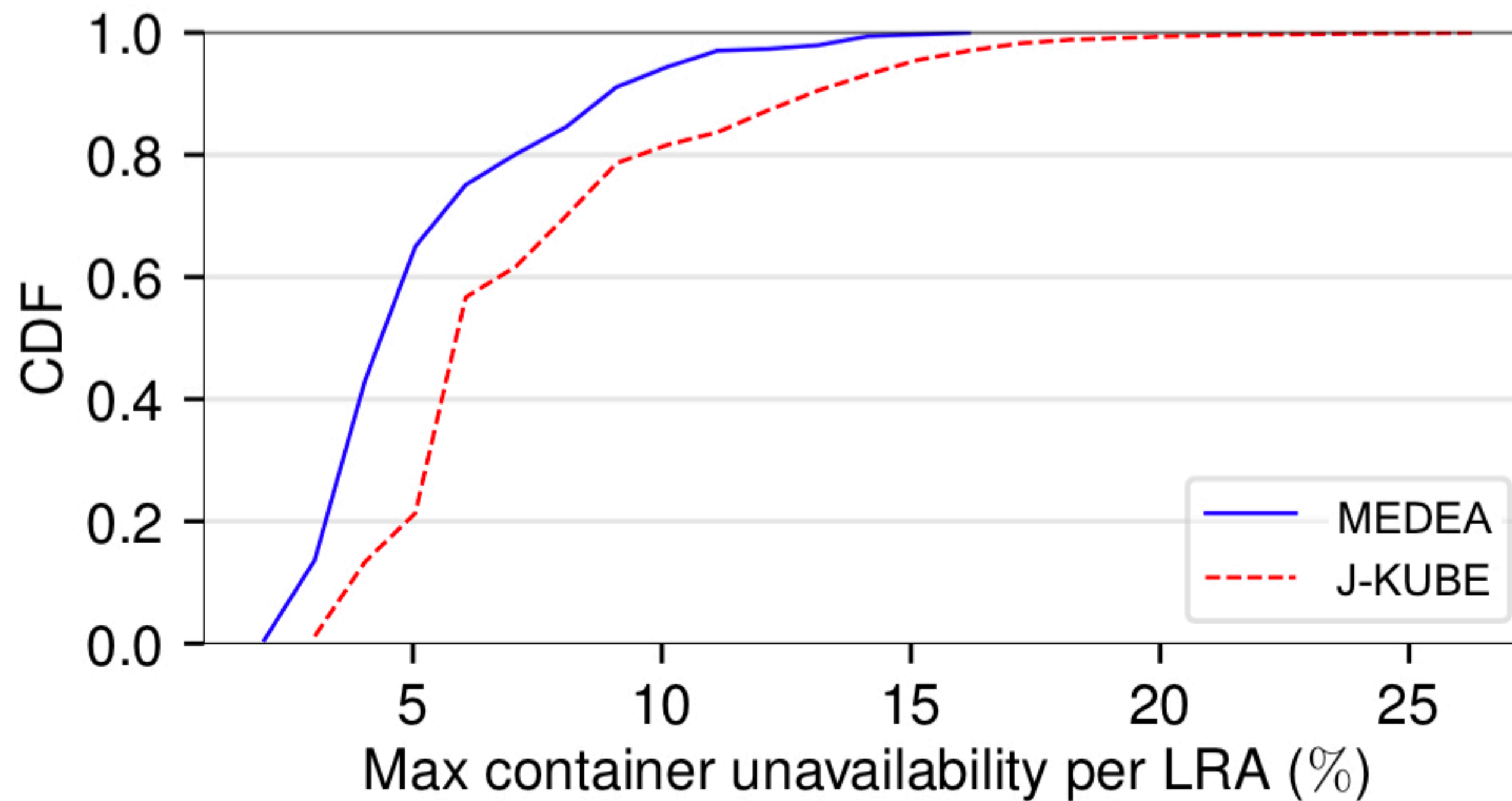**A 400 node pre-production cluster grouped into 10 racks, supplemented by simulation.**



**Figure 8: Application resilience over 15 days**

# Evaluation

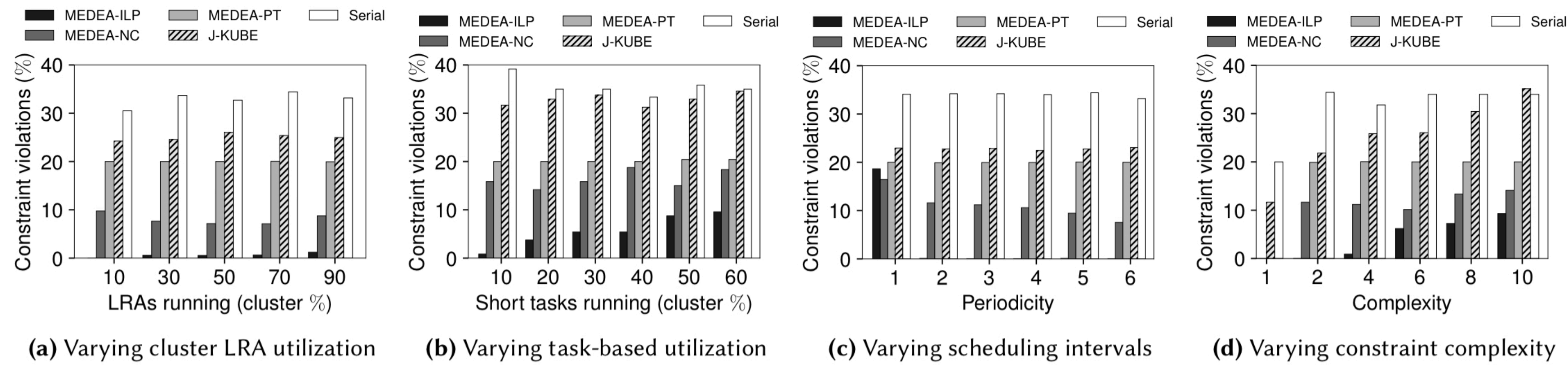**A 400 node pre-production cluster grouped into 10 racks, supplemented by simulation.**
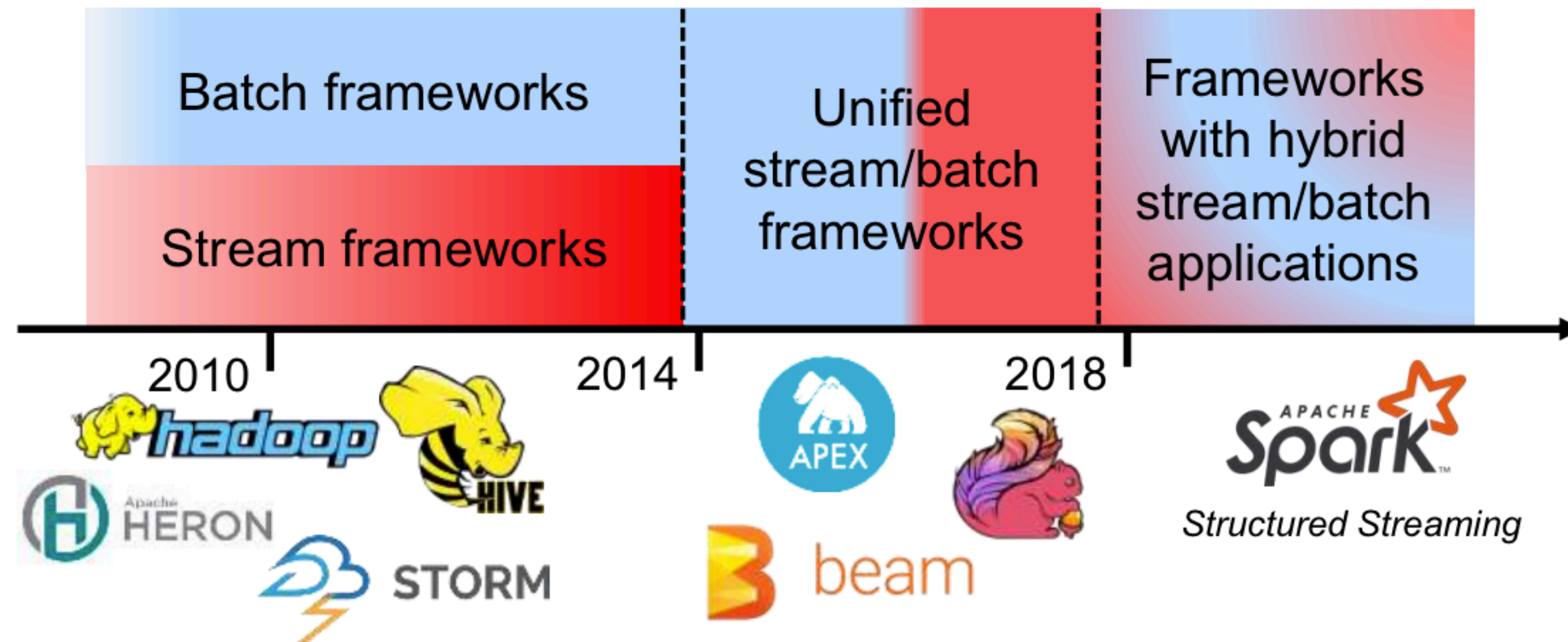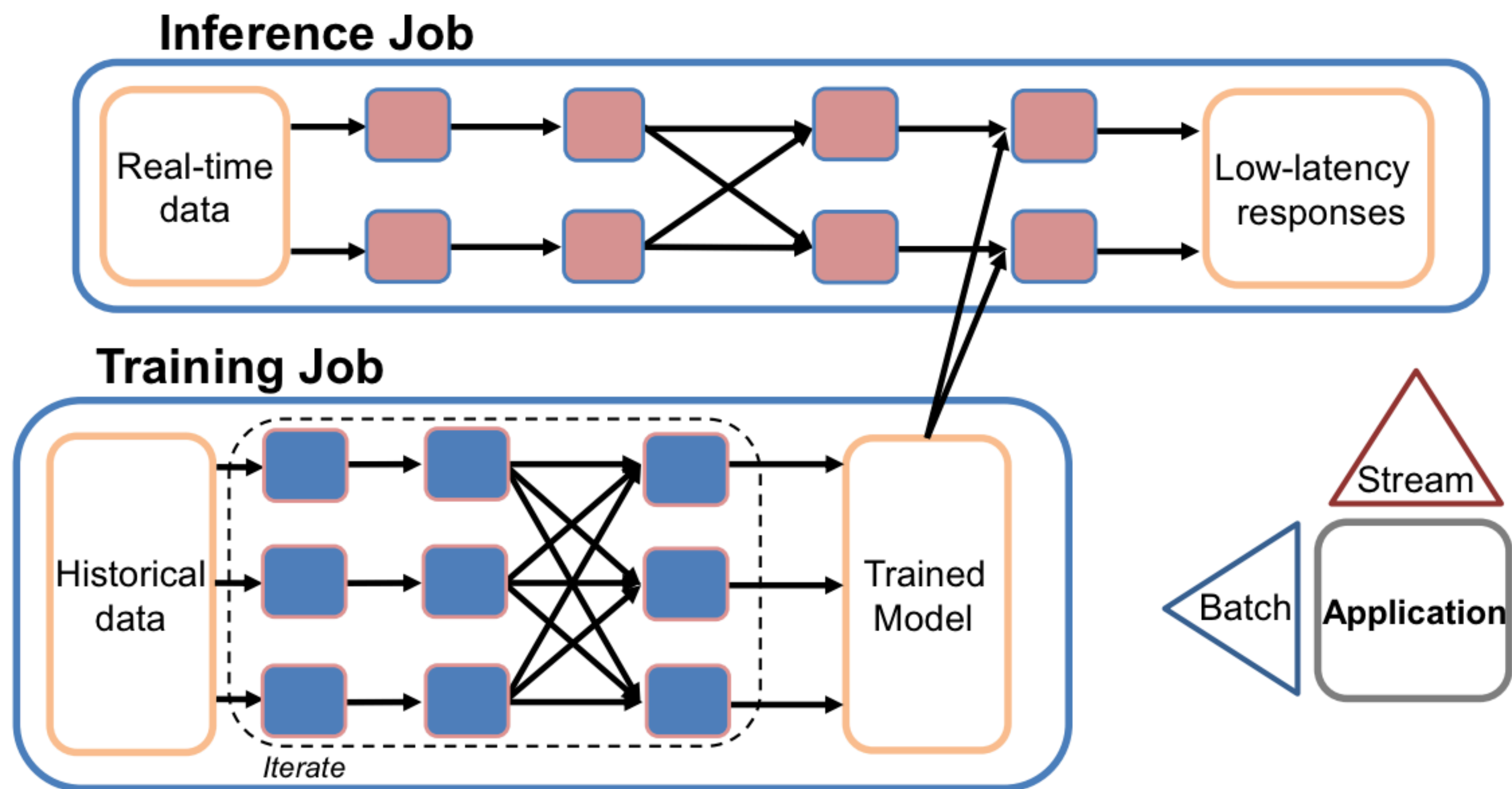


**Figure 9: Constraint violations**

# Neptune

Scheduling Suspendable Tasks
for Unified Stream/Batch Applications

# Evolution of analytics frameworks
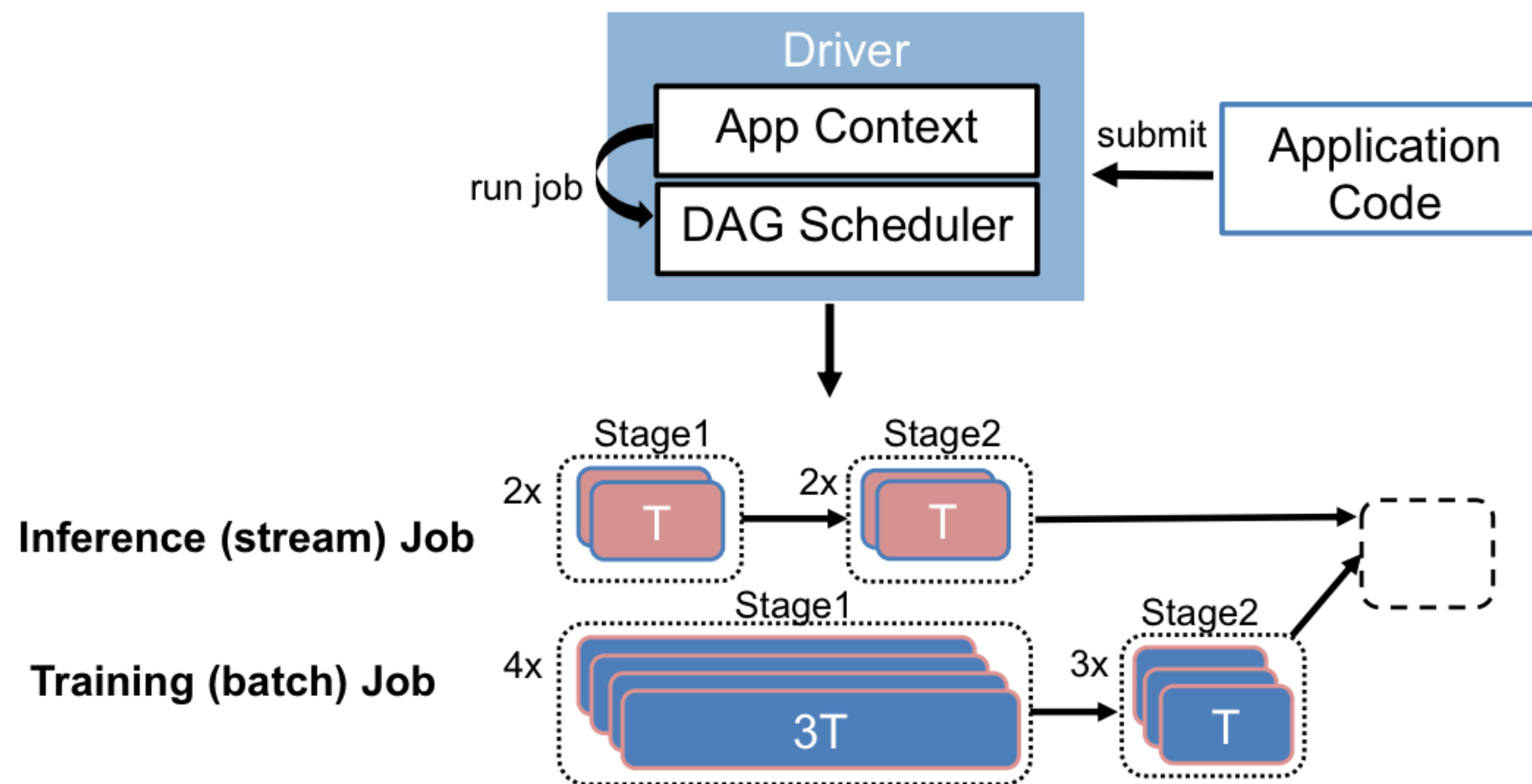


**UNIFICATION**

# Unified application example

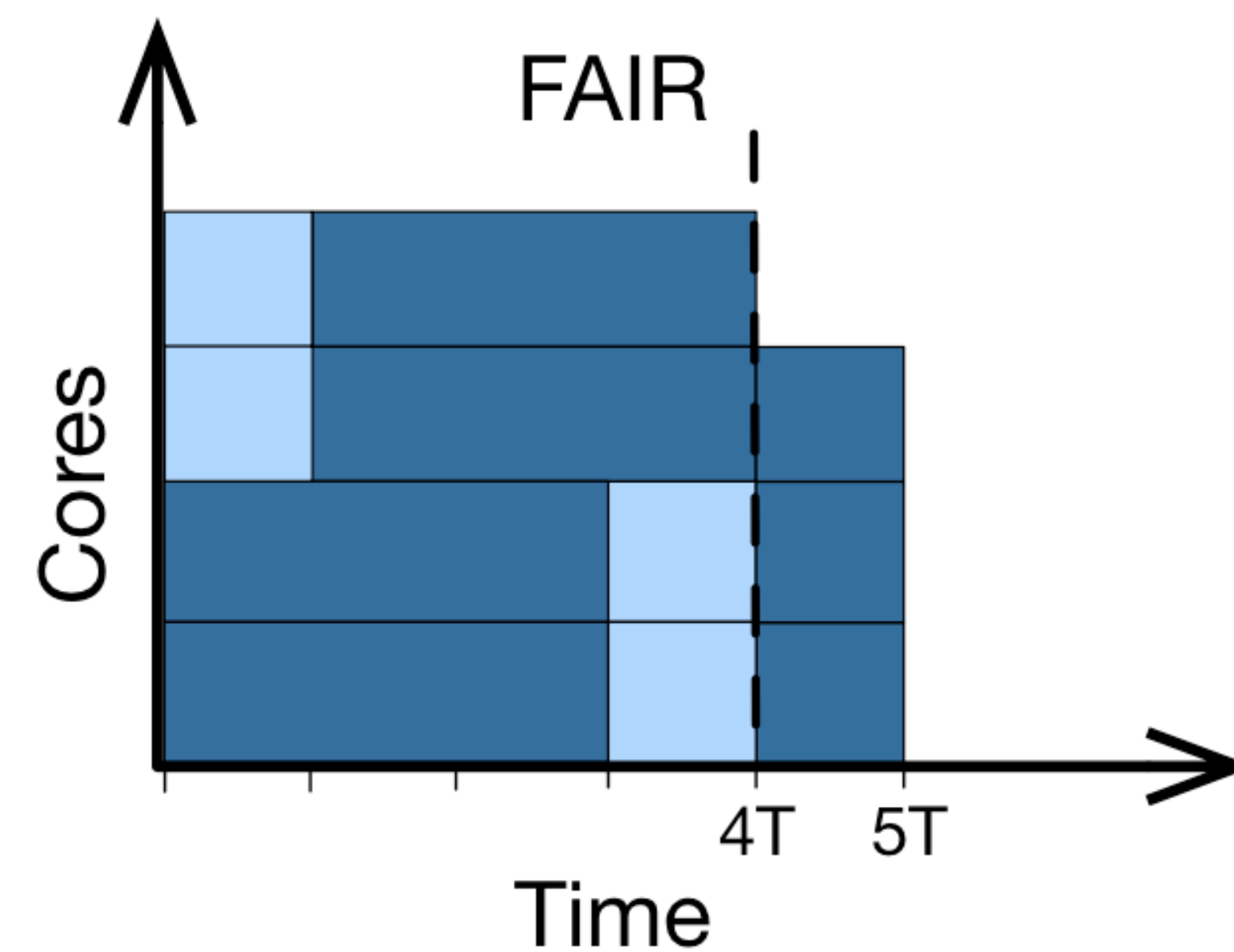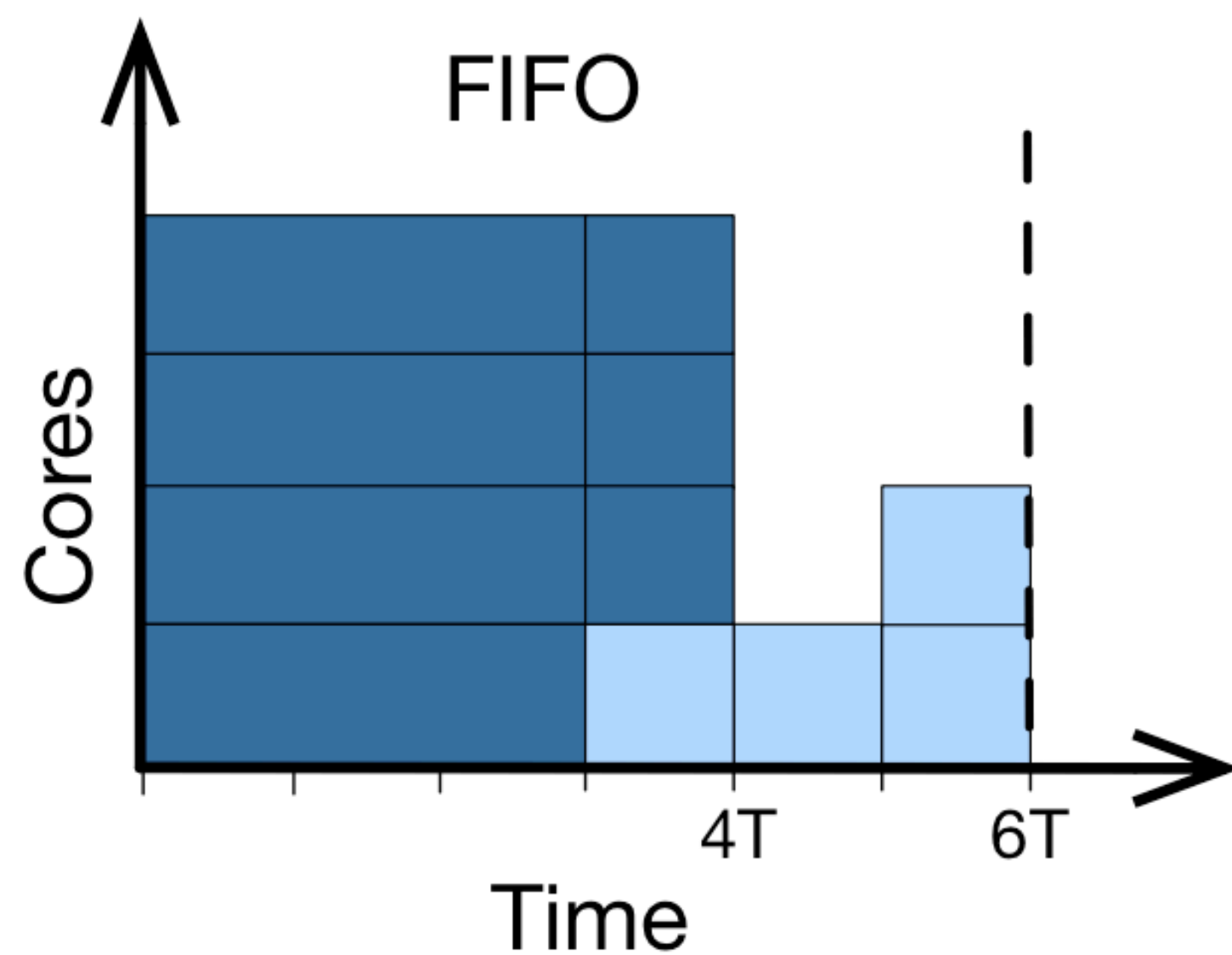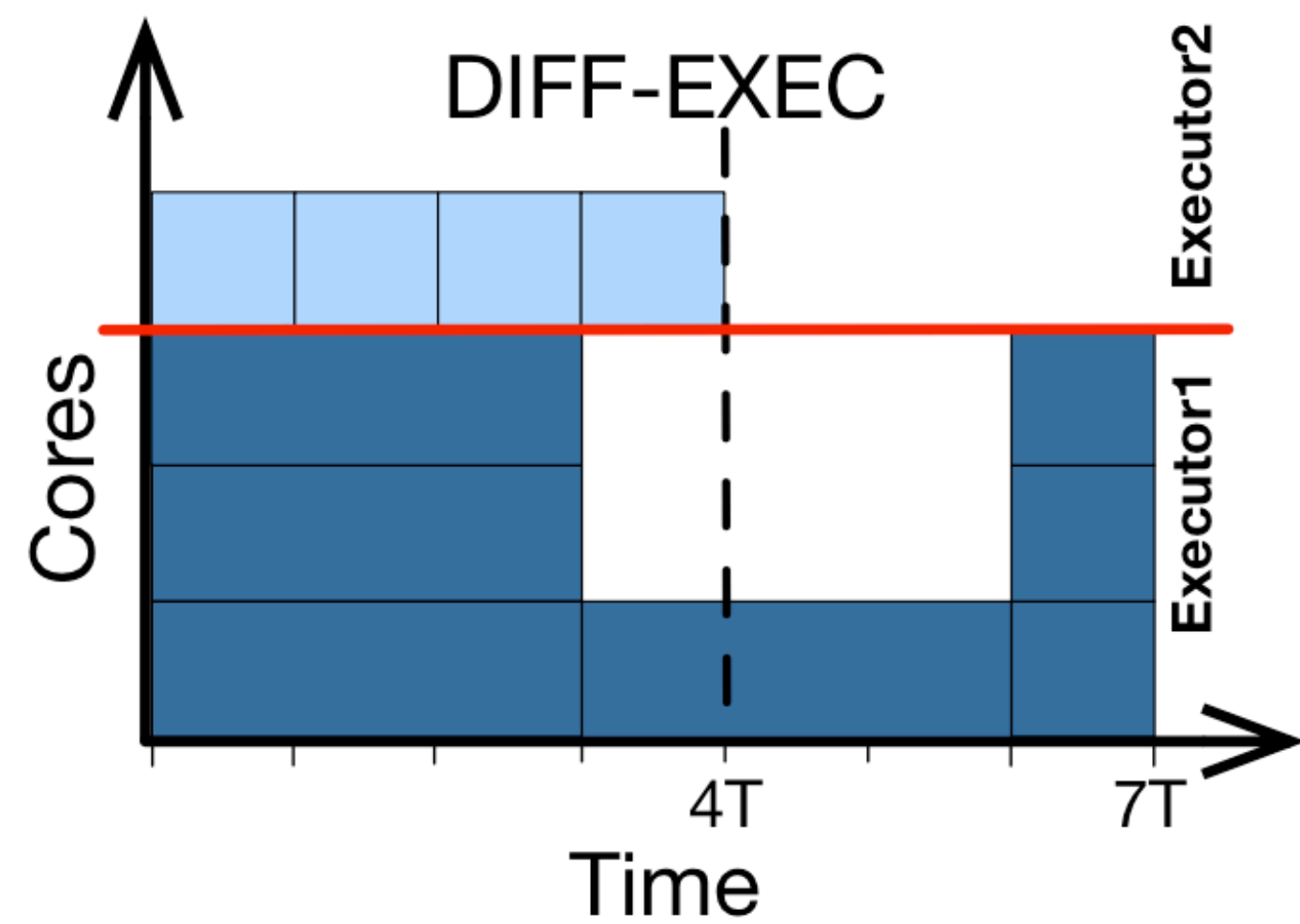# Stream/Batch application requirements
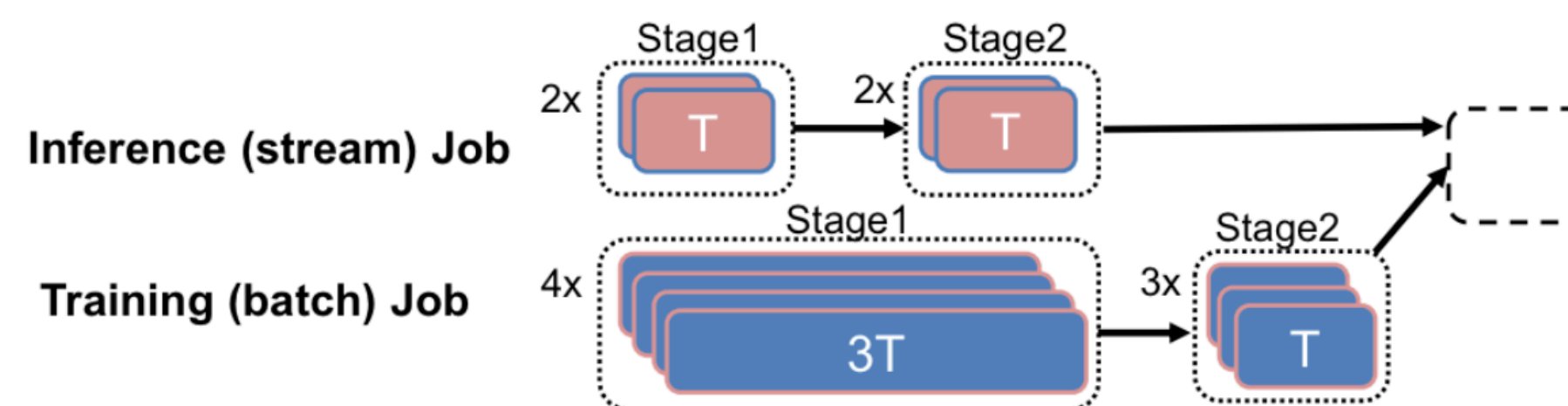
- **Requirements**
  - Latency: Execute inference job with minimum delay
  - Throughput: Batch jobs should not be compromised
  - Efficiency: Achieve high cluster resource utilization

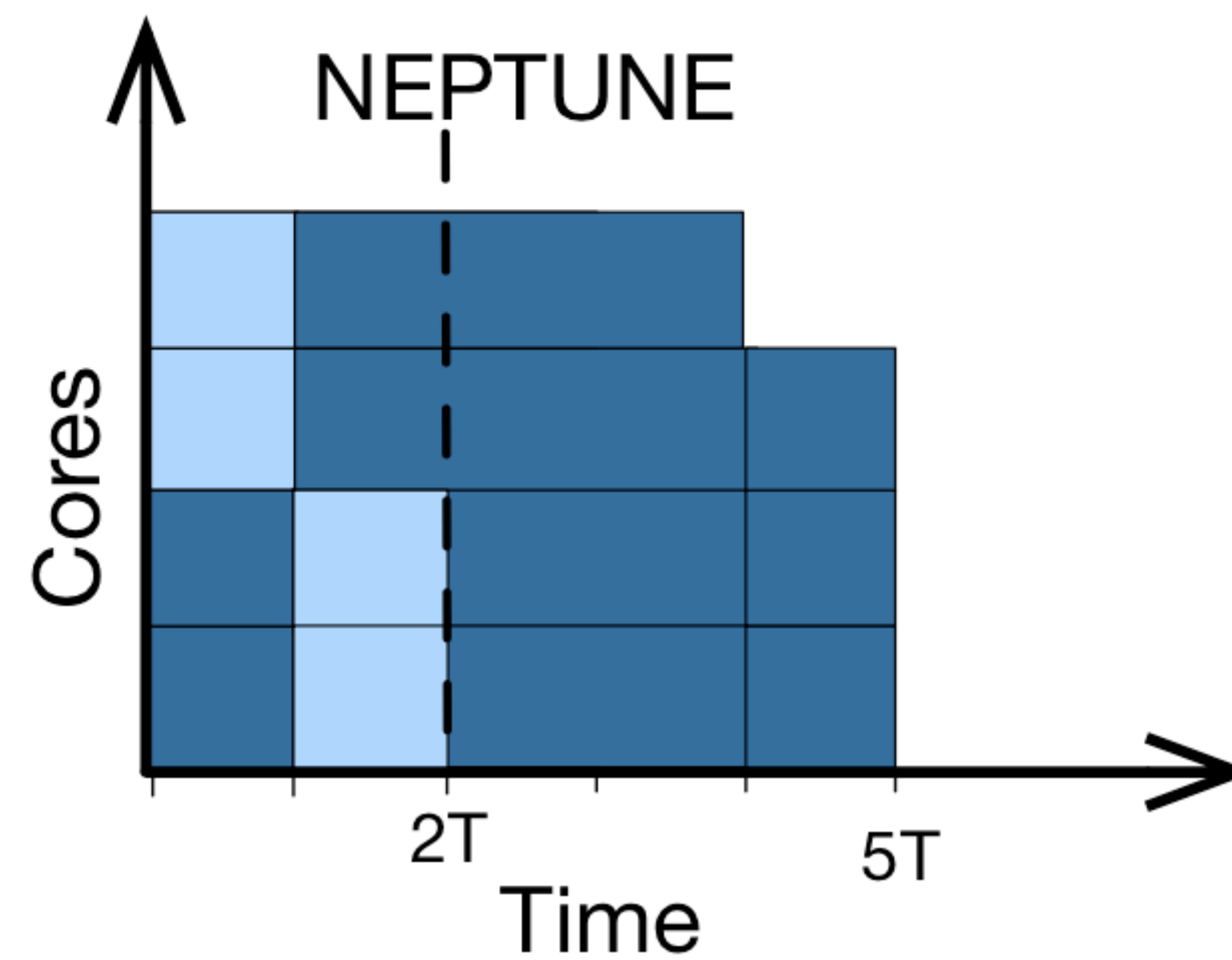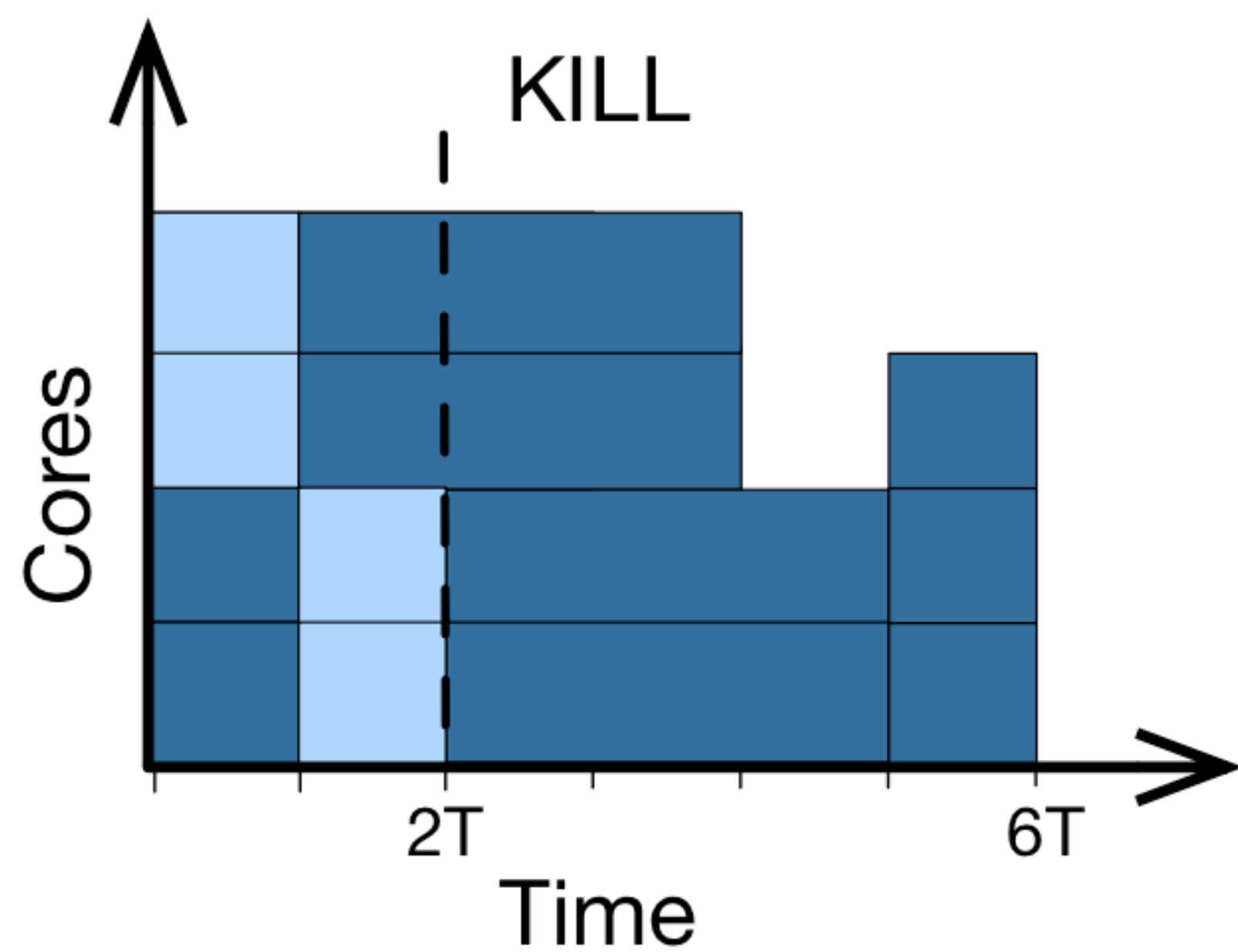**Challenge**: schedule stream/batch jobs to satisfy their diverse requirements
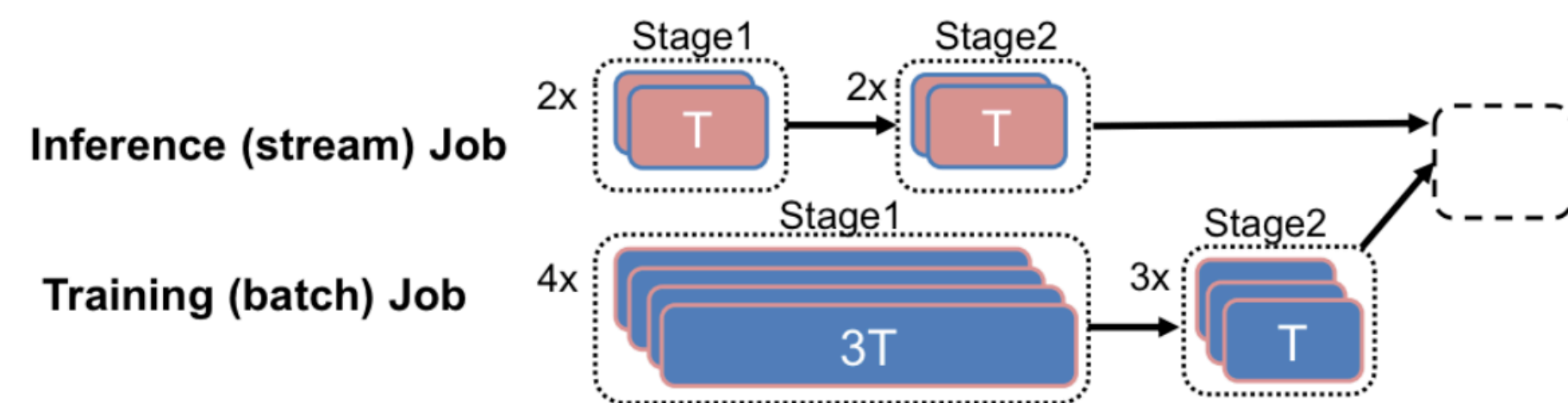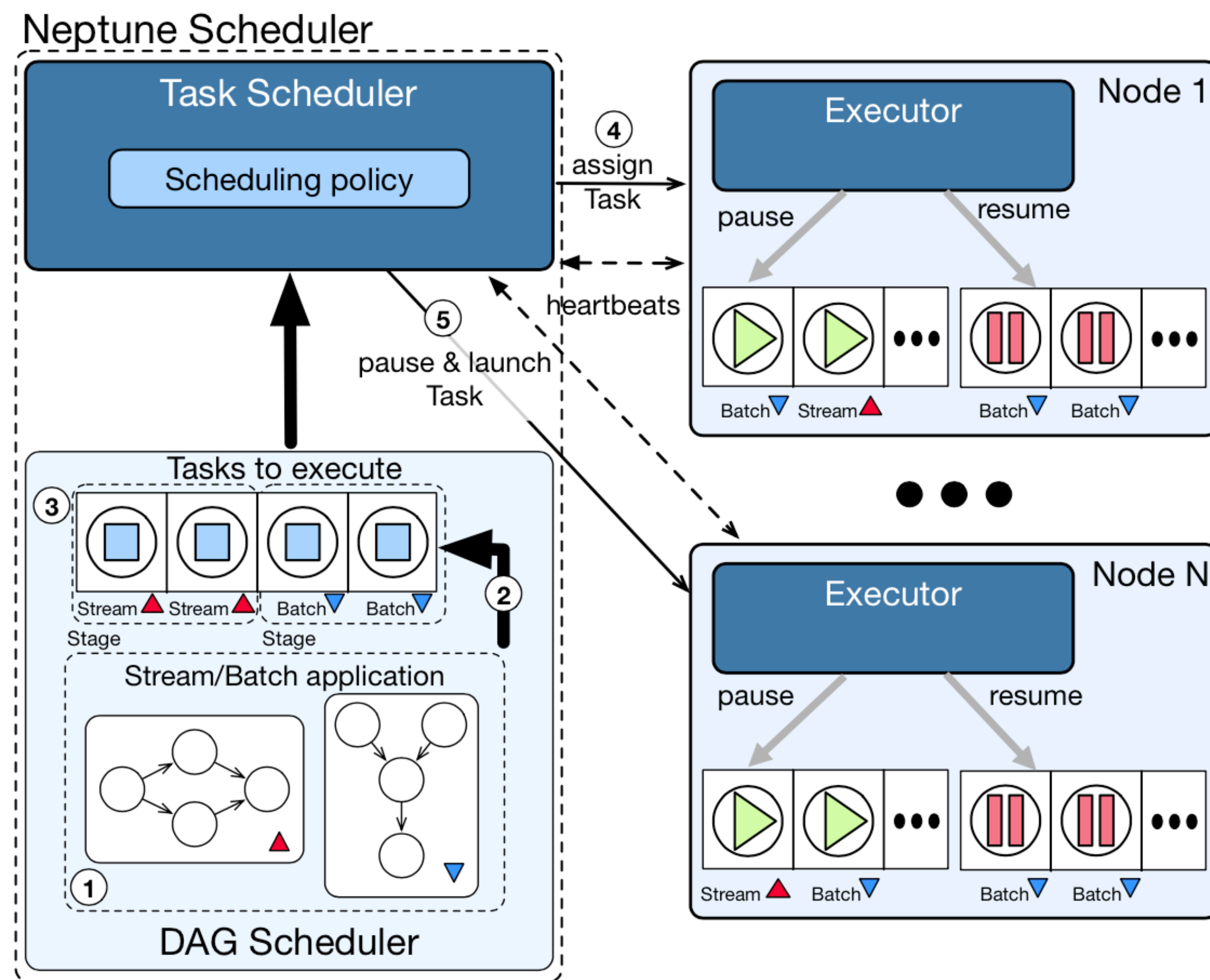
# Stream/Batch application scheduling

# Stream/Batch application scheduling

# Stream/Batch application scheduling

# Neptune Design
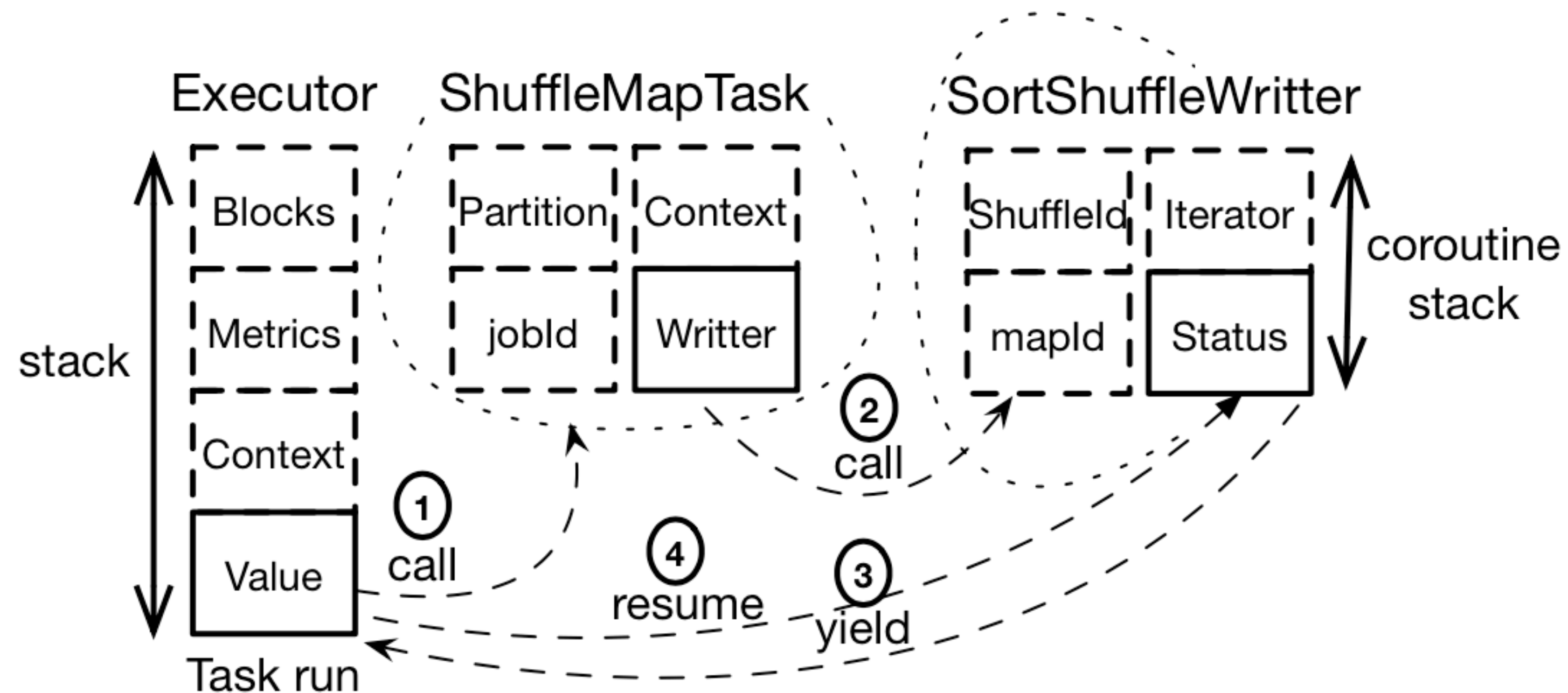
# Suspendable Tasks

**Neptune uses stackful coroutines to implement suspendable tasks, which have a yield point after the processing of each record.**

**Listing 2: Collect coroutine task**

```
1  val  collect : (TaskContext,  Iterator[T]) -> (Int, Array[T]) =
2      coroutine {(context: TaskContext, itr : Iterator[T]) =>
3        val  result  = new mutable.ArrayBuffer[T]
4        while ( itr .hasNext)   /* iterate  records */
5        | result .append(itr .next)   /* append record to dataset */
6        | if (context .isPaused())   /* check task context */
7        | | yieldval (0)  /* yield value to caller */
8      result .toArray /* return result dataset */
```

# Suspendable Tasks

**Neptune uses stackful coroutines to implement suspendable tasks, which have a yield point after the processing of each record.**

# Scheduling policies

- **Idea**: policies trigger task suspension and resumption
  - Guarantee that stream tasks <u>bypass</u> batch tasks
  - Satisfy <u>higher-level objectives</u> i.e. balance cluster load
  - Avoid <u>starvation</u> by suspending up to a number of times


- **Load-balancing** (LB): takes into account executors' memory conditions and equalize the number of tasks per node
- **Locality- and memory aware** (LMA): respect task locality preferences in addition to load-balancing
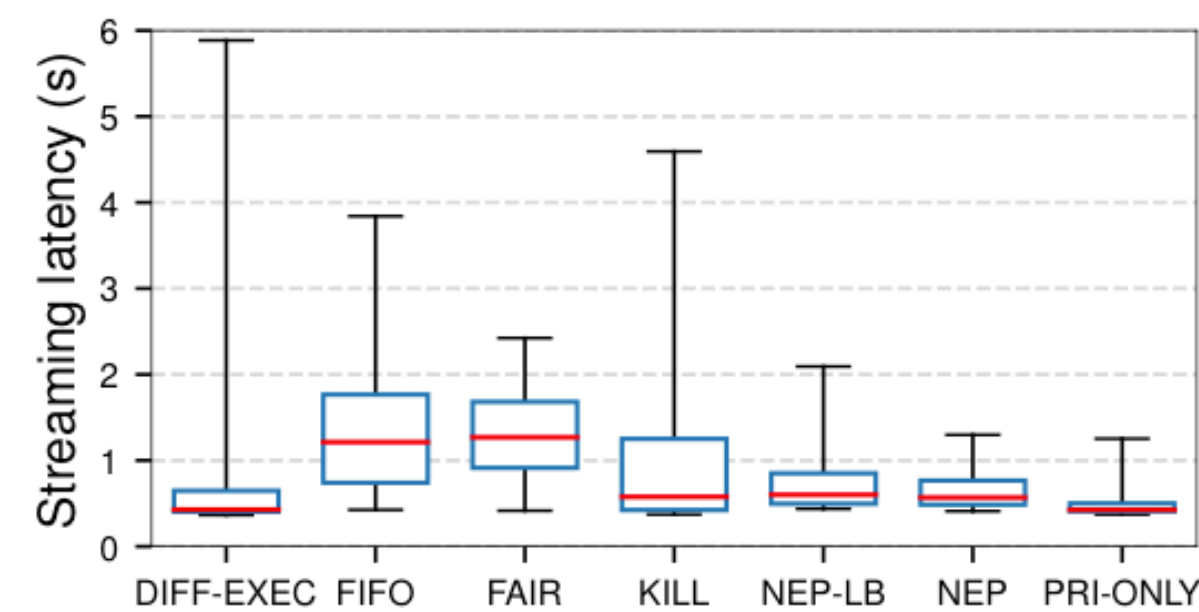
# LMA Scheduling policies

**Algorithm 1:** LMA scheduling policy

   **Input:** List Executors,       // In descending preference order

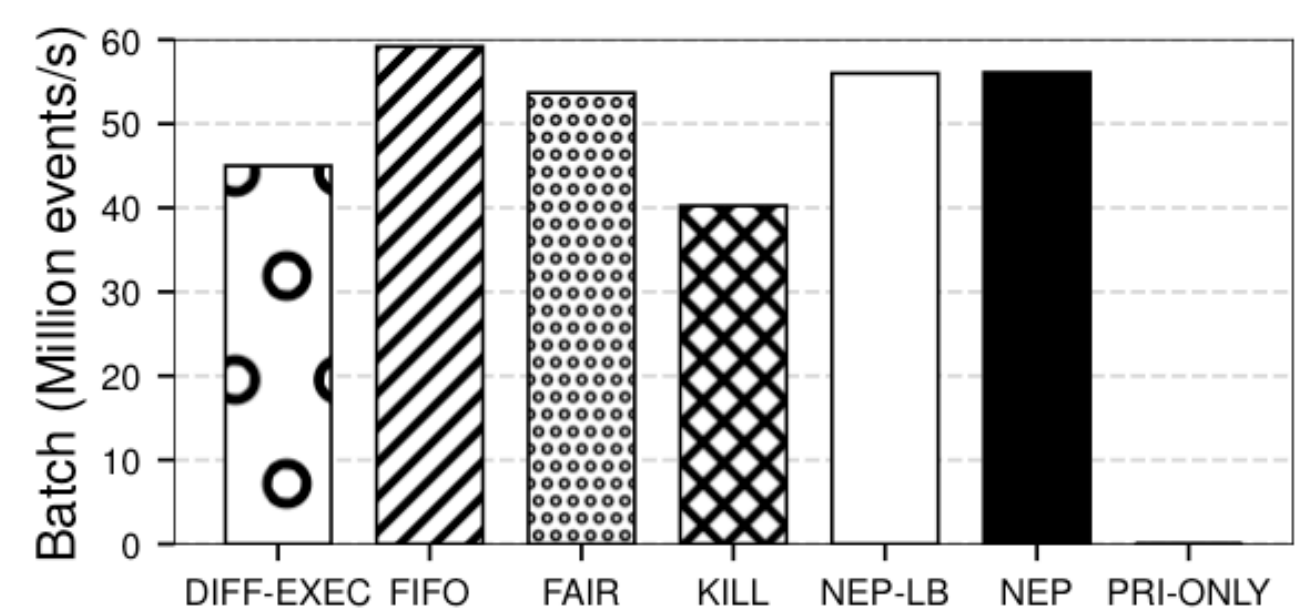1   List ExecutorsMetricsWindow         // Executor metrics

2   **Upon event** *HEARTBEAT hb from EXECUTOR e* **do**

3        ExecutorsMetricsWindow[e].push(hb.metrics)

4        Executors.updateOrdering

5   **return**

6   **Upon event** *SUBMIT Task t* **do**

        // Cache local executor

7        Executor $t_{exec}$ ← hostToExecutor.get(t.cacheLocation)

8        **if** $t_{exec}$ *is None* **or** $t_{exec}$.*freeMemory < threshold* **then**

            // Executor with the least pressure

9            $t_{exec}$ ← Executors.head

10      **if** $t_{exec}$ *has availableCores* **then**

            // Launch task *t* on free cores

11           $t_{exec}$.Launch(*t*)

12      **else**

            // Suspend task $t_p$ and launch *t*

13           Task $t_p$ ← $t_{exec}$.tasks.filter(LowPriority)

14           $t_{exec}$.PauseAndLaunch(*t*, $t_p$)
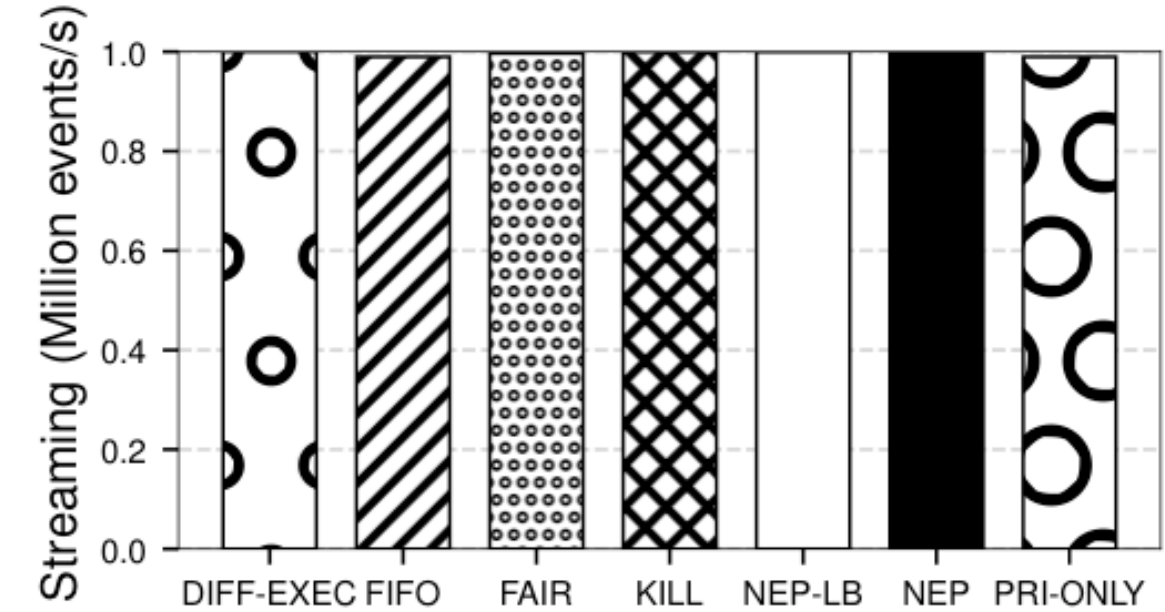
15  **return**

# Evaluation


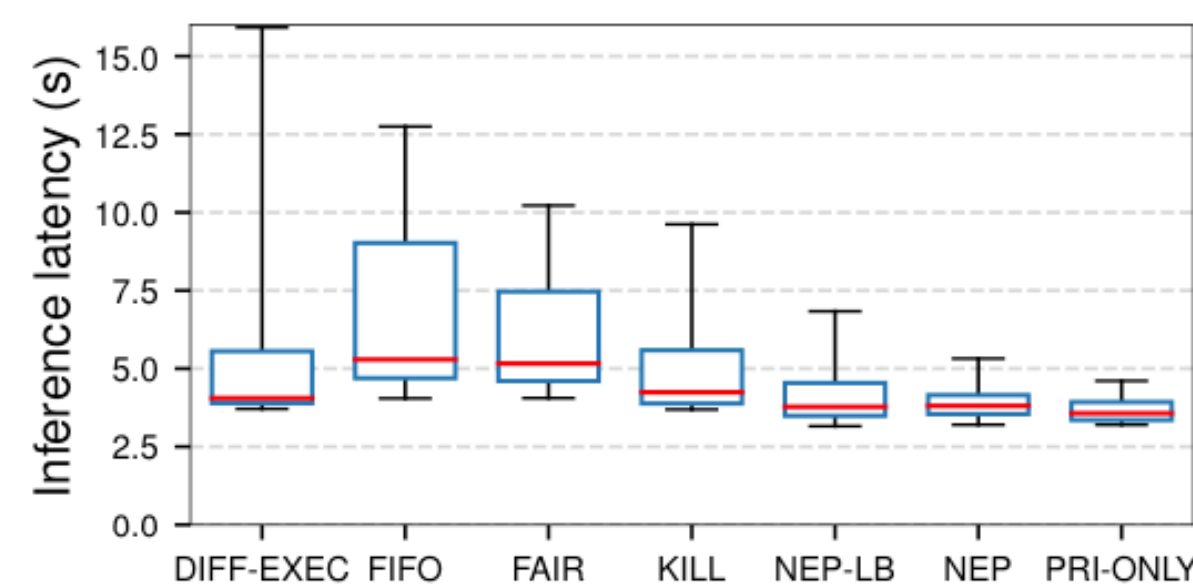
**(a)** Streaming query latency
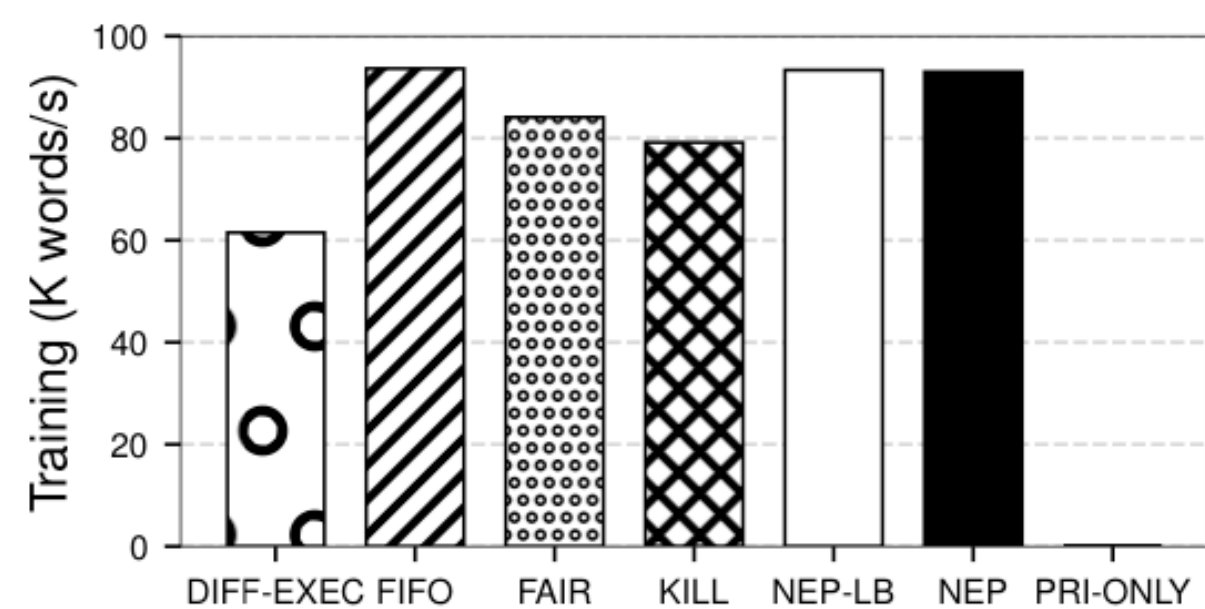
**(b)** Batch query throughput
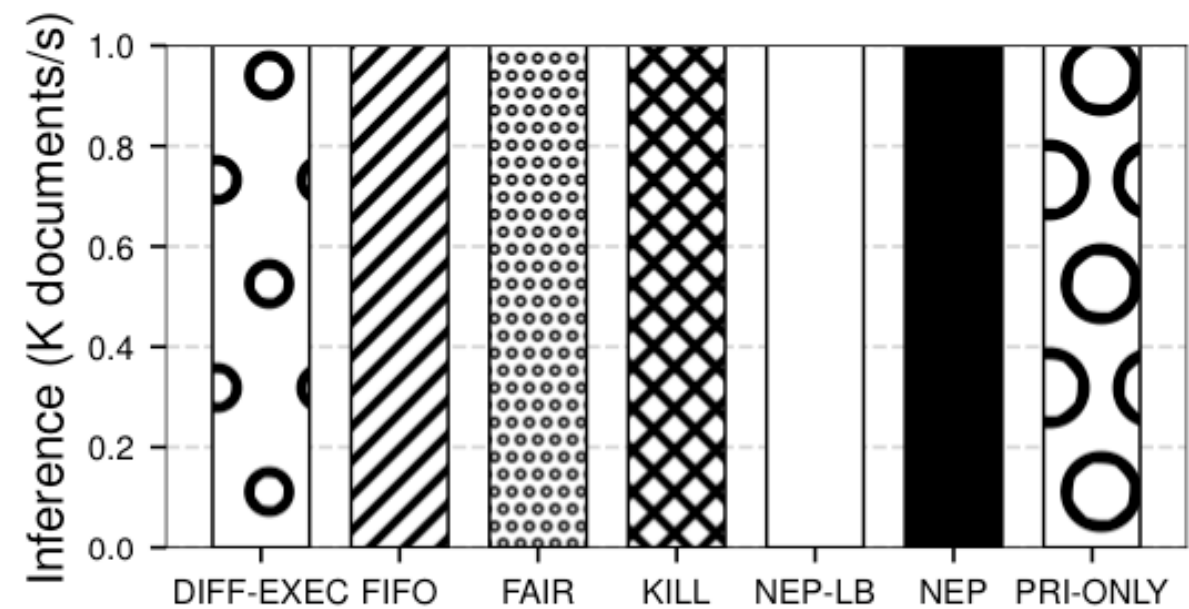
**(c)** Streaming query throughput

**Figure 9: Yahoo Streaming benchmark (streaming + batch)**
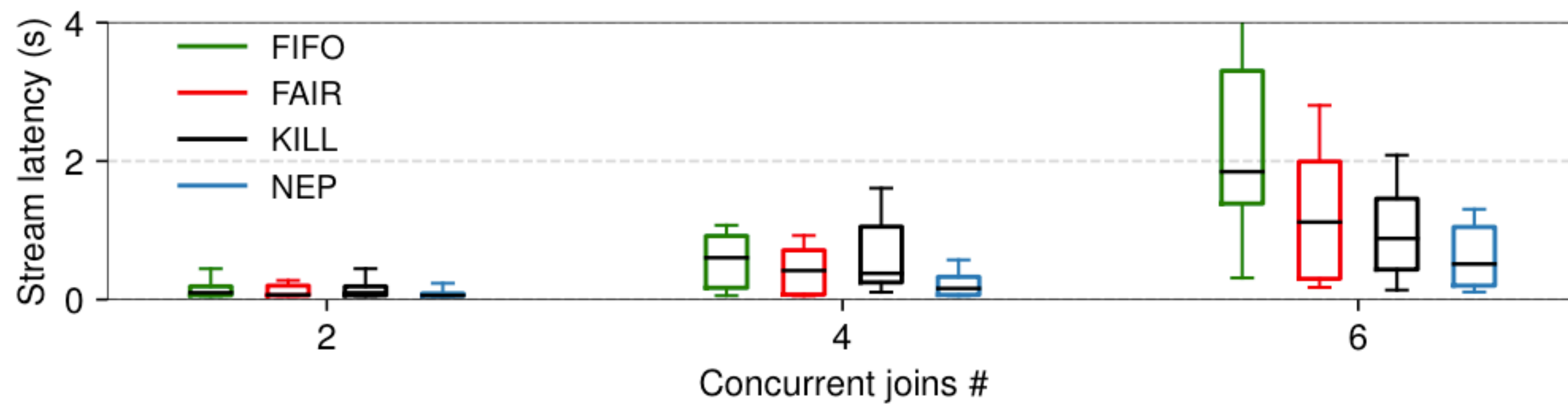
**(a)** Inference latency
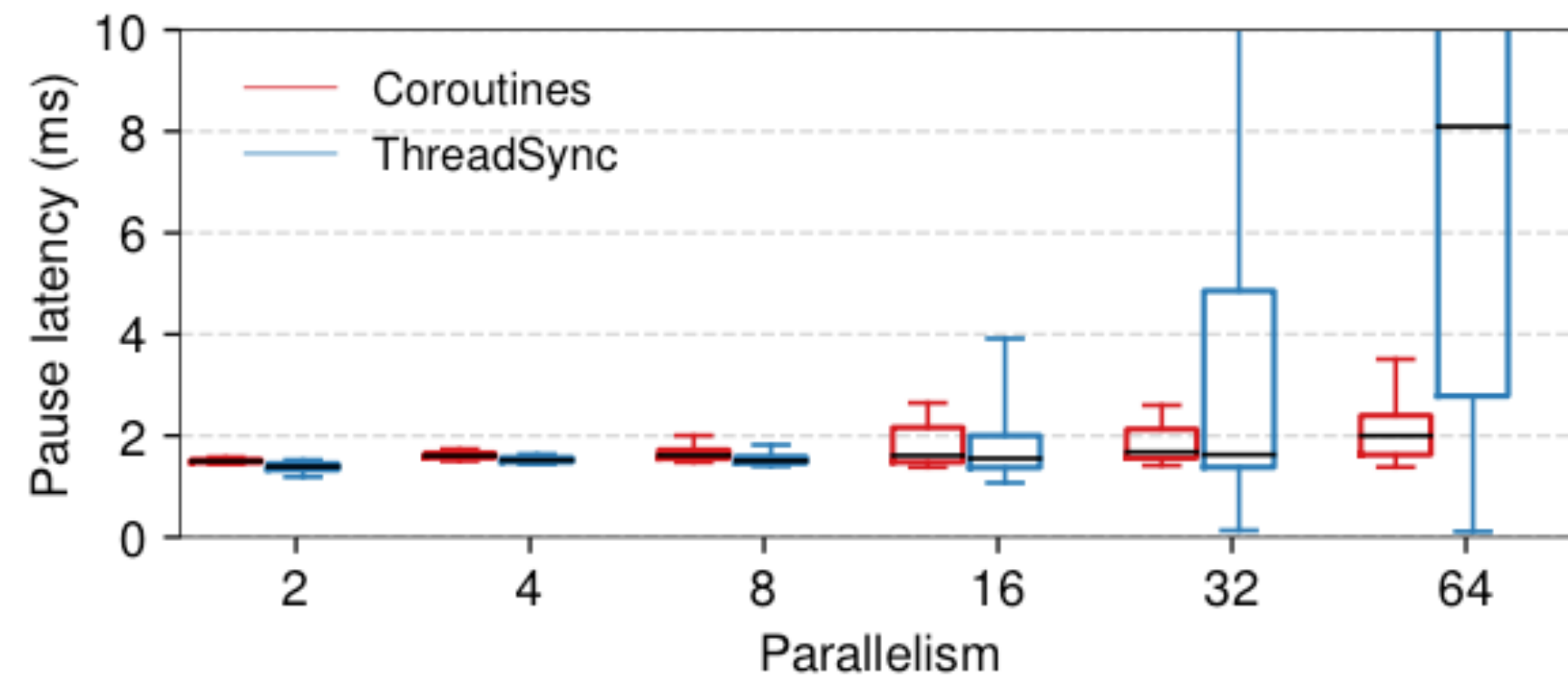
**(b)** Training query throughput

**(c)** Inference query throughput

# Evaluation

# Evaluation

# Summary

- **Medea**:
  - Support container tags and logical node groups
  - Expressive cardinality constraints
  - Two scheduler design

- **Neptune**
  - Suspendable Tasks
  - LMA Scheduling policies

# Thanks