# Henge: Intent-Driven
# Multi-Tenant Stream Processing

# Faria Kalim, Le Xu, Sharanya Bathey, Richa Meherwal, Indranil Gupta

## Distributed Protocols Research Group

## Department of Computer Science
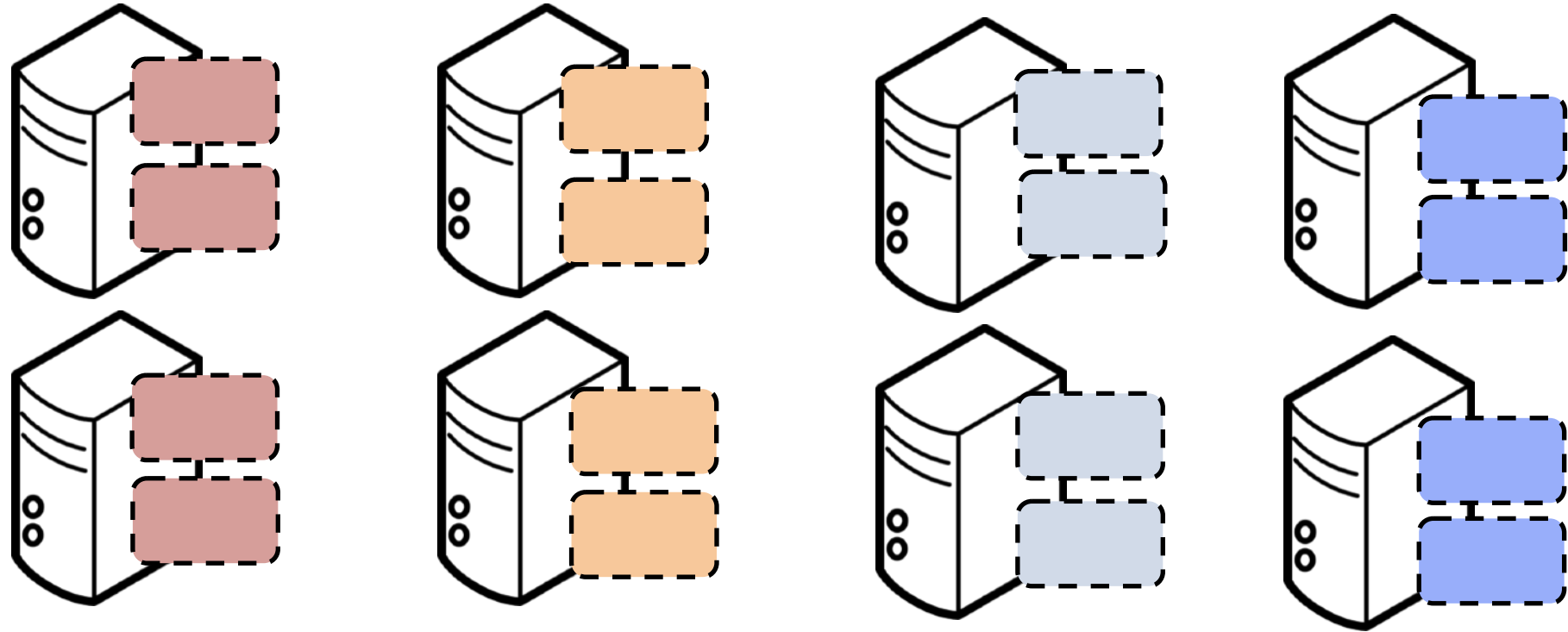## University of Illinois at Urbana Champaign

**Henge** allows stream processing jobs to satisfy **user-specified** performance requirements

while **reducing** costs

by performing online **resource reconfigurations** in a multi-tenant environment.
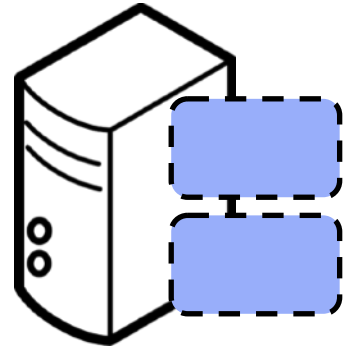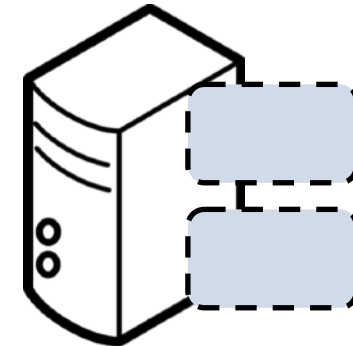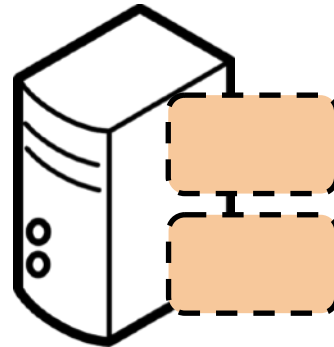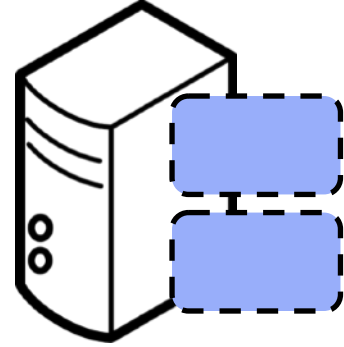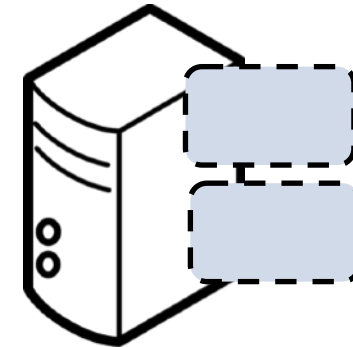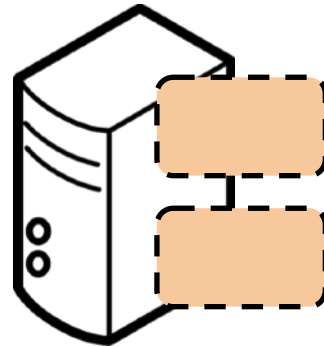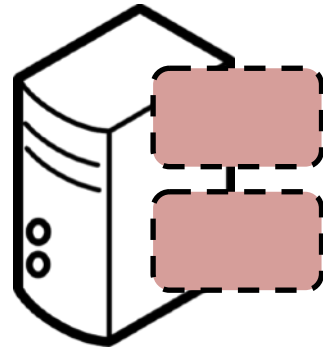
# A Typical Deployment



Per-job clusters →
overprovisioning

# A Typical Deployment

Low level metrics e.g., queue sizes, CPU load as performance indicators

# A Typical Deployment
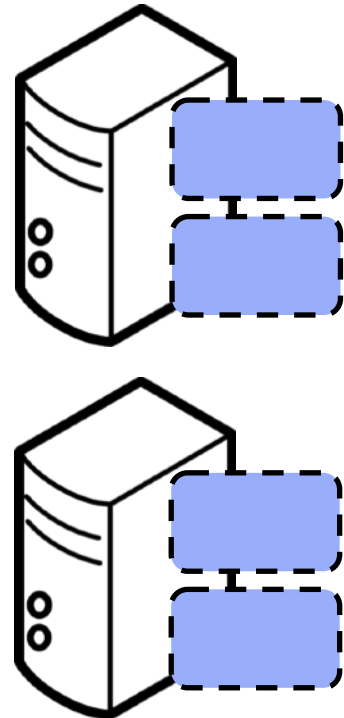
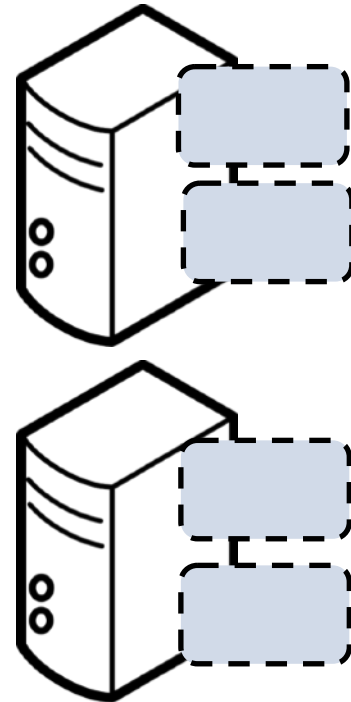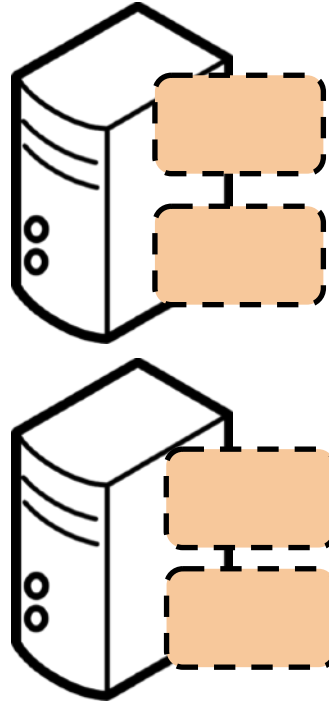Low level metrics e.g., queue sizes, CPU load as performance indicators
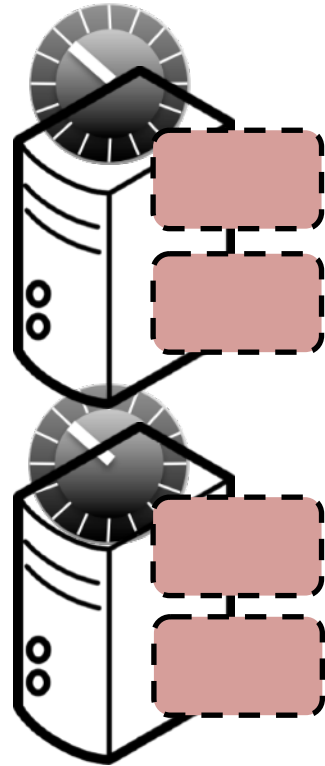


Job 1
Job 2
Job 3
Job 4

# A Typical Deployment

Low level metrics e.g., queue sizes, CPU load as performance indicators



Job 1

Job 2

Job 3

Job 4

Manual scaling

# Intent-Driven Multi-Tenancy

# Intent-Driven Multi-Tenancy

Efficient resource usage across multiple users
➔ Multi-tenancy

# Intent-Driven Multi-Tenancy

Efficient resource usage across multiple users
➜ Multi-tenancy

Application-aware adaptation to user requirements
➜ Intent-driven Multi-tenancy

# Intent-Driven Multi-Tenancy

Efficient resource usage across multiple users
➔ Multi-tenancy

Application-aware adaptation to user requirements
➔ Intent-driven Multi-tenancy

CPU Load, ❌ Queue Sizes ...

| Job | Description | Service Level Objective (SLO) |
|-----|-------------|-------------------------------|
| 1 | Finding ride price | Latency < 5 s |
| 2 | Analyzing earnings over | Throughput > 10K/hr. |

# Problem

How can we
**achieve user-facing service level objectives**
for stream processing jobs
on multi-tenant clusters?

# Problem

How can we
**achieve user-facing service level objectives**
for stream processing jobs
on multi-tenant clusters?

Latency,
Throughput

# Absolute Throughput SLOs are not Useful

# Absolute Throughput SLOs are not Useful



8

# Absolute Throughput SLOs are not Useful

# Absolute Throughput SLOs are not Useful



8

# Absolute Throughput SLOs are not Useful



8

# Absolute Throughput SLOs are not Useful

**Job Operations**

# Absolute Throughput SLOs are not Useful

**Job Operations**



**Juice**: fraction* of the input data processed by the job per unit time.

# Jobs benefit even below SLO threshold

Job Utility Functions

# Jobs benefit even below SLO threshold

Job Utility Functions



**Utility function for a single job**

# Jobs benefit even below SLO threshold

Job Utility Functions



Expected Utility

Current Utility

Latency (Seconds)

10

8

2

0

0    20    40    60    80    100

**Henge's goal → Maximize the total utility of the cluster**

**Utility function for a single job**

# Background: Stream Processing Topologies (Jobs)



**Logical DAG for a Word Count Job**

**Diamond Topology**

**Star Topology**

# Background: Stream Processing Jobs



["So it goes..."]

["So"]
["it"]
["goes"]
...

["So"]

["goes"]

["it"]

**Spout**

**Spout**

**Splitter**

**Splitter**

**Count**

**Count**

**Count**

**Count**

**Executors (Threads)**

# Background: Stream Processing Jobs

["So"]
["it"]
["So it goes..."] ["goes"]
...

["So"]

Spout

["goes"]

Splitter

Count

Count

Spout

Splitter

Count

Count

Parallelism → 2

["it"]

Executors (Threads)

13

# Background: A Physical Deployment

# Background: A Physical Deployment



Spout  Splitter
Count  Count

Workers

# Henge's Cluster-Wide State Machine



Not Converged

Converged

Total System Utility < Total Expected Utility

# Henge's Cluster-Wide State Machine



Reconfiguration

Reversion or Reconfiguration

Not Converged

Converged

Reduction

Total System Utility < Total Expected Utility

# Reconfiguration

De-congest operator by increasing parallelism level of executors

1) Reconfiguration

2) Reconfiguration



Not Converged

Converged

# Reconfiguration

De-congest operator by increasing parallelism level of executors
3) **Black-list** topologies that show less than Δ% improvement

1) Reconfiguration

2) Reconfiguration



Not Converged

Converged

# Bottlenecks



Workers

# Reconfiguration

Workers

18

Bottlenecks

Reconfigs.

High
Load

Spout

Splitter

Splitter

Splitter

Count

Count

Workers

19

Bottlenecks

Reconfigs.

High
Load

**SLO-Satisfying Job**

Spout

Splitter

Splitter

Splitter

Count

Count

Workers

Bottlenecks

Reconfigs.

High
Load

Reduction

# Reduction

20

# Reduction

Reconfigurations → drop in utility

**Reduction**

Not
Converged

# Reduction

Reconfigurations → drop in utility

If high CPU load on majority of machines, **reduce** parallelism for operators that

    a) are in topologies that satisfy their SLO

    b) are not congested

**Reduction**

Not Converged

# Reversion

Reconfigurations → drop in utility *and reduction is not possible*
**Revert** to a past configuration that provided best utility



**Reversion**

Not Converged

Converged

# Evaluation

Real-world workloads:

 Yahoo!

 Twitter

 Web log traces

Experimental Setup:

 10-40 node Emulab cluster

# Reducing cost and achieving high utilities

# Reducing cost and achieving high utilities



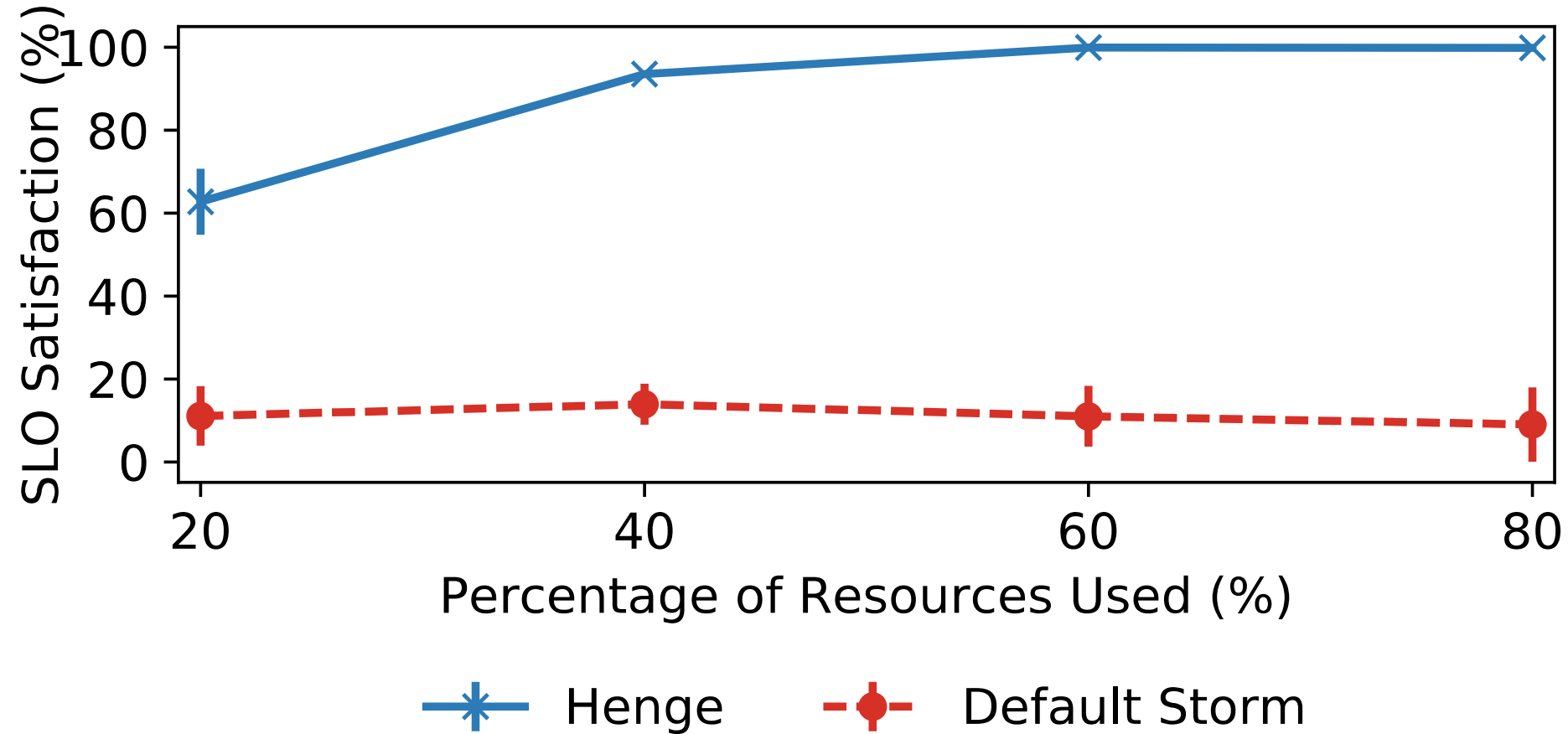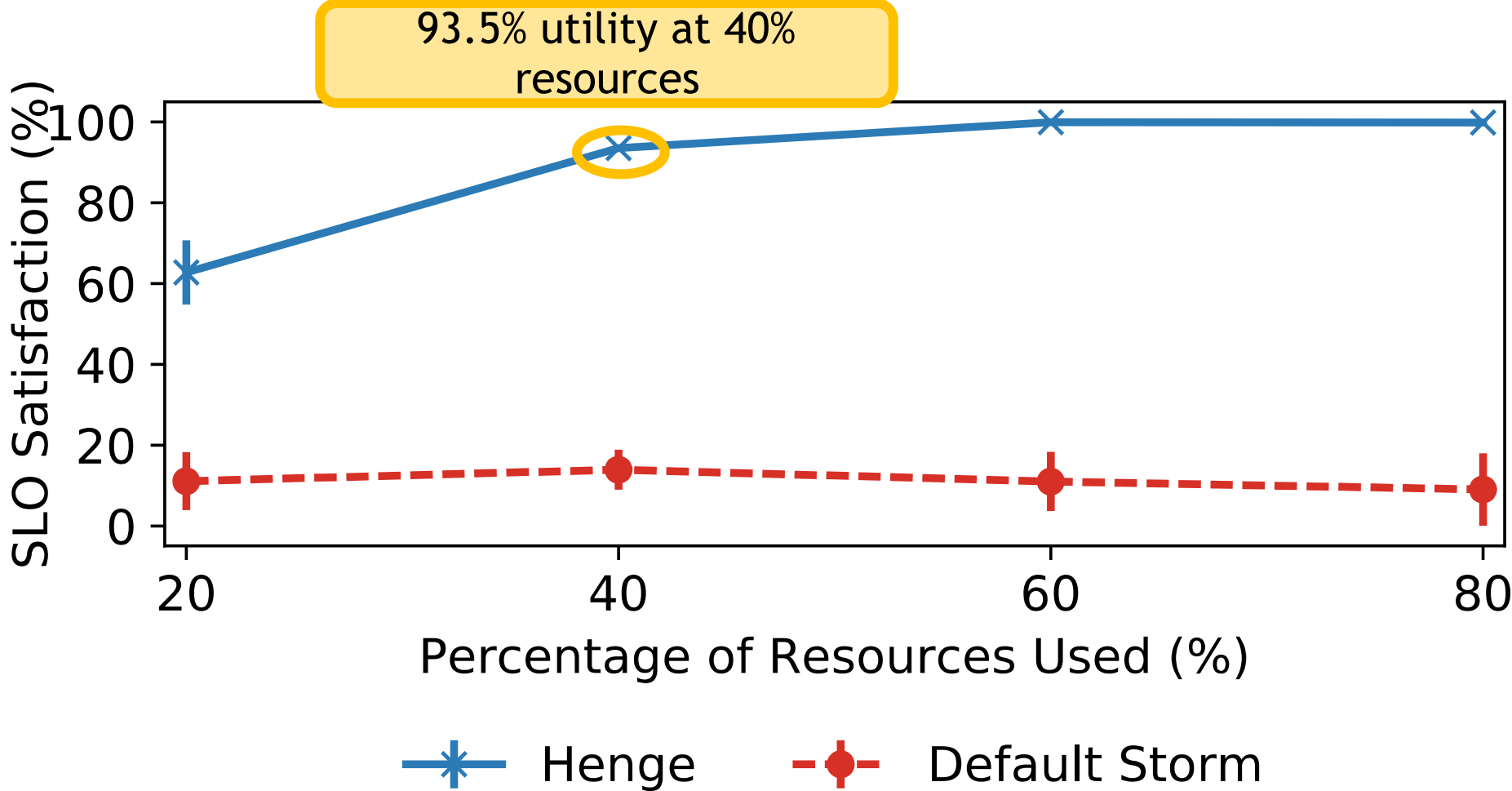93.5% utility at 40% resources

# Reducing cost and achieving high utilities



93.5% utility at 40% resources

100% utility at 60% resources

SLO Satisfaction (%)

Percentage of Resources Used (%)

Henge    Default Storm

# Adapting to a Diurnal Pattern

Legend: SDSC Input, SDSC Output, EPA Input, EPA Output

Top panel: No. of Tuples (x1000) vs time (Day 1, Day 2)

Bottom panel: Utility vs time (Day 1, Day 2), with Max. Utility and Reconfigurations annotations

25

Fewer reconfigurations are required once a job has adjusted to max load

Max. Utility

Reconfigurations

# Can Henge do better than manual configuration?

# Can Henge do better than manual configuration?



Henge does better in the 15th to 45th percentile, and is comparable later.

# Scaling Cluster Size

# Scaling Cluster Size



Limited resources entail more reconfigurations to reach max. utility

# More Results

Henge can:

    handle dynamic workloads

        abrupt e.g., spikes & natural fluctuations

        gradual e.g., diurnal patterns

    satisfy hybrid SLOs

    scale with number of jobs & cluster size

    gracefully handle failures

# Summary

- Henge allows users to specify **performance intents** for their jobs

- Henge's goal is to **maximize cluster-wide utility**

- The scheduler performs fine-grained **reconfigurations** to allow stream processing jobs to meet user-specified intents