# Vertically Autoscaling with Cassandra

**Karla Saur - Principal Research SDE**
*Microsoft - Gray Systems Lab (GSL)*

# About me

I am a Principal Research SDE in [Gray Systems Lab (GSL)](#) at Microsoft, an applied research lab within Azure Data.

Before Microsoft, I completed my PhD ~10 years ago, and worked as researcher on Telco/5G autoscaling
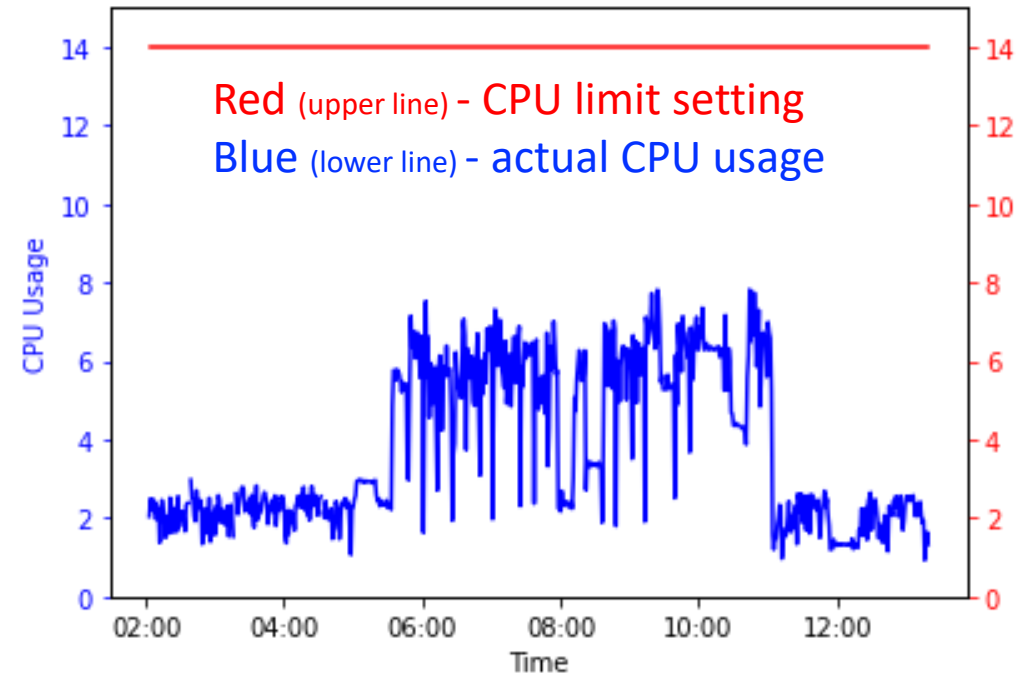
My main research focus is *optimizing cloud infrastructure for database and machine learning workloads* from a general perspective.

I hope that some of the techniques I mention are useful to you, and I am excited to learn more about Cassandra while at this conference!

# About me

I am a Principal Research SDE in [Gray Systems Lab (GSL)](#) at Microsoft, an applied research lab within Azure Data.

Before Microsoft, I completed my PhD ~10 years ago, and worked as researcher on Telco/5G autoscaling

My main research focus is *optimizing cloud infrastructure for database and machine learning workloads* from a general perspective.

I hope that some of the techniques I mention are useful to you, and I am excited to learn more about Cassandra while at this conference!

# Initial Project: Monolithic DB's on K8s

- Originally, our team was tasked with **optimizing deployments of monolithic databases** (ex: 1 primary, 2 secondaries, fixed) running on Kubernetes.

- We found that **many users were overprovisioned** in terms of CPU allocation, which is how we bill (#CPUs/hour)
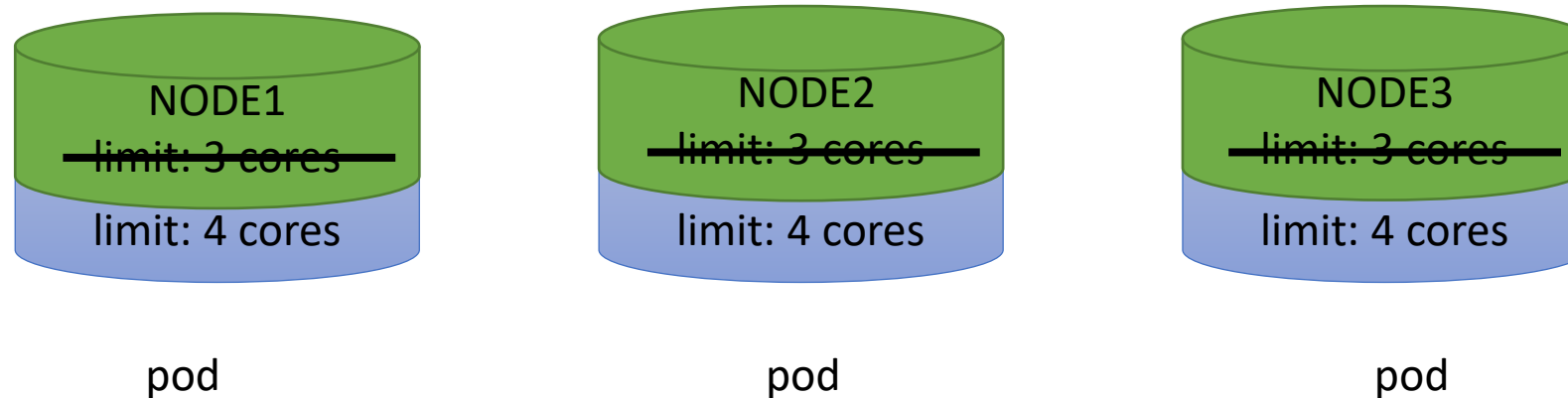
- This began our **vertical scaling journey**.



Red (upper line) - CPU limit setting
Blue (lower line) - actual CPU usage

# Vertically autoscaling with Cassandra

- Cassandra is famous for its linear scaling, seamlessly adds more nodes.
  - However, bringing up a new node (horizontal scaling) involves data movement and can take a significant amount of time

- Vertical scaling does not involve data movement.
  - It can provide an additional mechanism to right-size resources.

- There are many ways to run Cassandra (on VMs, on containers, on Kubernetes, K8ssandra, etc), and techniques are generally applicable.
  - We can often scale the cores in-place without restart

# Scenario: Vertically Scaling Generic DBs

- Example: a database runs as a set of Kubernetes Pods with $n$ cores each.
  - Users are billed based on a max CPU limit they specify



- Kubernetes excels at HORIZONTAL pod autoscaling, but our database use case is a fixed number of replicas. But we can scale VERTICALLY!

# Generic Vertical Autoscaling (end-to-end)

Application
limit: 3 cores

Application
limit: 3 cores

Application
limit: 3 cores

Could be Kubernetes,
some other container
orchestrator, etc.

scales

Could be Kubernetes, or
other program capable of
making changes to the
application as-needed.

controller

triggers action

Scaler
Healthy?
Resources?
Do decision.

metrics
server

Could be rest API,
recommendation
service, etc.

reads

Could be any (csv, open
telemetry, other standards).

select
metrics
data

reads

Pluggable
recommender
algorithms

writes

decisions

# Outline

- Mechanism
  - Vertically scaling in-place
  - Cassandra perf impacts?
- Policy
  - CaaSPER: Proactive/Reactive algorithm for balancing price-perf trade-off
- Code/demo: VASIM - Vertical Autoscaling SIMulator
  - Try your own autoscaling algorithm!
  - Autotuning: parameter tune your own algorithm
  - How to get started with Cassandra

# Mechanism: scaling in-place

Application
limit: 3 cores

Application
limit: 3 cores

Application
limit: 3 cores

scales

controller

Could be Kubernetes,
some other container
orchestrator, etc.

triggers action

Scaler

Healthy?
Resources?
Do decision.

metrics
server

Could be rest API,
recommendation
service, etc.

reads

Could be any (csv, open
telemetry, other standards).

select
metrics
data

reads

Pluggable
recommender
algorithms

writes

decisions

# Restarts hinder scaling nimbly in stateful workloads
*(even with containers!)*

Rolling restart HA process (~10-15 min) makes scaling perf much worse than necessary due to delay.





*Rolling upgrade process*

With no restarts, we could minimize throttling and optimize scaling further and more safely by reacting faster



*No restart!*

# In-place/no-restart scaling of CPUs

- Docker: `docker update some-cassandra --cpus` **or** `--cpu-quota`

- K8s: In-place updates officially an alpha feature in the 1.27 [release](release) (~April 2023) under the feature gate `InPlacePodVerticalScaling`
  - Changed simply by patching the running pod spec
  - Default behavior is "in-place" unless `resizePolicy` **is set to** `RestartContainer`

```
resizePolicy:
- resourceName: memory
  restartPolicy: RestartContainer
- resourceName: cpu
  restartPolicy: NotRequired
```

# What is the perf impact if *the DB thinks it has n cores, but we actually give it m cores?*

- 3 replica SQL (**on Linux**) deployed on *32 core* K8s nodes

- Start scaling op every 200 seconds
  - Because the machine has 32 cores, SQL Server thinks it has 32 cores
  - Scale from 2 → 4 → 6 → 8 → 10 → 12 → 16 → 24
  - Plot the average throughput during each segment

- 2 tests
  - Scaling only (i.e. SQL activates 32 SqlOS schedulers)
  - Coordination via affinity changes prior to scale

**Scaling only. Overall avg txn/s: 3960**

32 schedulers fighting for n < 32 cores
worth of time, results in lock waits,
priority inversions, etc.

Only 24 cores were
available for user workload

Pod core limit
(full node
schedulers)

2/32    4/32    6/32    8/32    10/32    12/32    16/32    24/32

Time Group (10 secs)

**Coordination via affinity tweaks. Overall avg txn/s: 4102**

Affinity tweaks especially useful when #
cores significantly lower than advertised
(ex: 2-6 cores out of 32)

Perf change
when
schedulers
match cores
limit

+31%    +30%    +30%    +1%    +6%    -3%    0%    0%

2/32    4/32    6/32    8/32    10/32    12/32    16/32    24/32

Time Group (10 secs)

IMPORTANT: Many
SQL instances are
often idle and could
be scaled down to 2-
4 cores. The longer
we can keep them
there, the longer we
have cores to use in
other places. But we
don't want to lose
perf!

# Perf impact of restart-free core scaling with Cassandra?



Average CPU Usage

- For this customer provisioned at 32 cores, we could scale down by ~20 cores each night
  - But what is impact to the locks/buffers/threads/etc if the DB thinks it has 32 cores but only has 12?



container

nproc=32
cpu-shares=12

# Testing perf impact of restart-free core scaling

## Quick experiment with Cassandra:

- Ran tried matched/mismatched on a customer's workload: *pleasantly boring.*



container
nproc=4
shares=4
4 core VM (E16-4ds_v5)

container
nproc=16
shares=4
16 core VM (E16ds_v5)

*calls/second diff:*
*in the noise/nearly identical*

container
nproc=8
shares=8
8 core VM (E16-8ds_v5)

container
nproc=16
shares=8
16 core VM (E16ds_v5)

*calls/second diff:*
*in the noise/nearly identical*

container
nproc=8
shares=8
8 cores available

container
nproc=96
shares=8
96 cores available

*calls/second diff:*
*in the noise/nearly identical*

- Java's `public int availableProcessors()` - apps must poll this explicitly
  - Cassandra doesn't, so we expected to see some mismatch due to threadpools, etc.
  - However: to use the new cores, Cassandra needs to be aware of the max cores at startup.
- Some JVM-weirdness related to List of Processors...*to be continued!*

# What about memory??

- Resizing memory without a restart is challenging, regardless of platform or environment (Python/C/JVM/etc).

- Memory resizing likely requires application changes.

- For now, restart/rolling-update when memory needs to be resized.

# Outline

- Mechanism
  - Vertically scaling in-place
  - Cassandra perf impacts?
- Policy
  - CaaSPER: Proactive/Reactive algorithm for balancing price-perf trade-off

# Vertical autoscaling: CaaSPER Algorithm



**Vertically Autoscaling Monolithic Applications with CaaSPER:**

Scalable Container-as-a-Service Performance Enhanced Resizing Algorithm for the Cloud

Anna Pavlenko, Joyce Cahoon, Yiwen Zhu, Brian Kroth, Michael Nelson,
Andrew Carter, David Liao, Travis Wright, Jesús Camacho-Rodríguez, Karla Saur
{firstname}.{lastname}@microsoft.com
Microsoft

*Vertically Autoscaling Monolithic Applications with CaaSPER (Pavlenko et al. SIGMOD 2024)*

NODE3
limit: 3 cores

triggers action

Scaler
healthy?
Resources?
Do decision.

controller

reads

metrics

decisions.txt

reads → Pluggable recommender algorithms → writes

metrics data

# Why not K8s built-in Vertical Pod Autoscaler (VPA)?

- For billing, we scale only at whole-cores, with `limits=requests`
  - In Kubernetes, `requests` and `limits` define guaranteed and burstable CPU resource allocation for applications. Setting these equal 'breaks' VPA.

- Must consider customer preferences when scaling
  - Most existing VPA tools are for optimal scheduling, not cost-perf preferences



Kubernetes VPA Default Algo
Scales initially, then never again

CaaSPER
Right-sizes quickly and adapts

Red - CPU limit setting
Blue - actual CPU usage

# Vertical auto-scaling approach:

## *2-part approach*



**Reactive:**
Handle the initial Pod size fit and adjust to spikes

**Proactive:**
Combine CaaSPER with existing algorithms to handle cyclical/ predictable loads over time

metrics data

Pluggable recommender algorithms

reads

writes

reads

# Vertical auto-scaling: Part 1
*Reactive/initial right-sizing with no data*

NODE1
limit: 3 cores

NODE2
limit: 3 cores

NODE3
limit: 3 cores

**Reactive:**
Handle the initial Pod size
fit and adjust to spikes

Proactive:
Combine CaaSPER with existing
algorithms to handle cyclical/
predictable loads over time

reads

decisions.txt

CPU Usage
10
5
0

reads

metrics
data

Pluggable
recommender
algorithms

writes

# Reactive CaaSPER

Doppler (prior work) provides initial SKU (#cores/#mem) selection *offline* for SQL Server based on personalized price-perf curve

- We adapted this price-perf curve for our container scenario by monitoring the change in slope over time, instead of focusing on a static price-perf curve for migration



if > threshold, scale up



if < threshold scale down

*Doppler: Automated SKU Recommendation in Migrating SQL Workloads to the Cloud. PVLDB 15, 12 (2022).*

# Reactive CaaSPER

## Steepness of price-perf curve determines how MUCH to scale

$$\text{Scaling Factor}(s) = \log(bs + c)$$

- $s$: Slope of the PP curve at the existing number of cores

- $b$: Skew estimate of the distribution of existing slopes

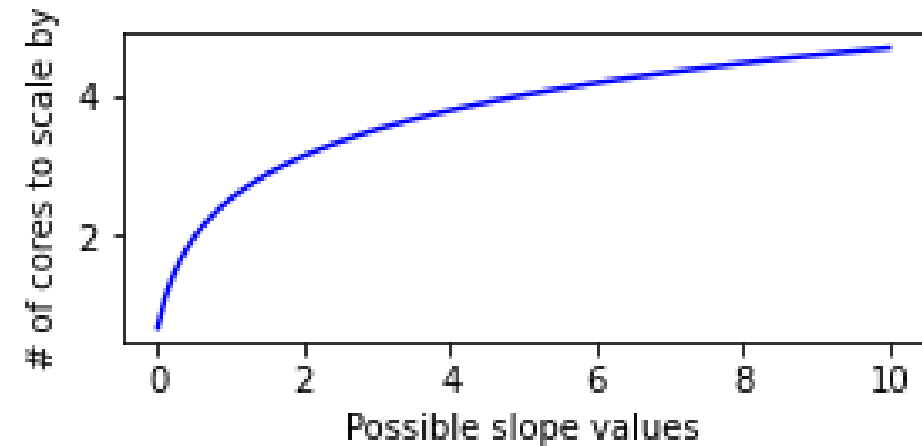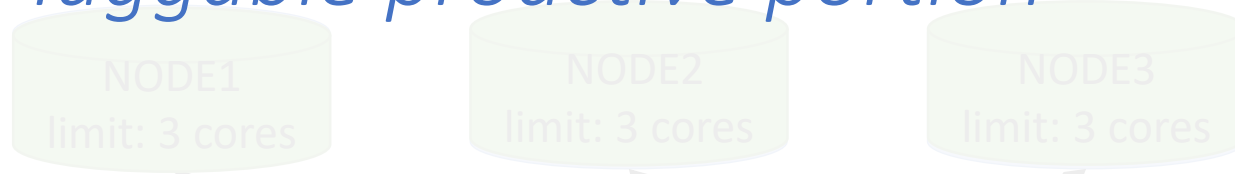- $c$: Minimum number of cores needed to operate



**Figure 6: Example shape of scaling-factor function** $SF(s)$ **of** *PvP-curve* **slope** $s$**. Scale-ups happen more aggressively for large** $s$ **(more throttling), than small** $s$ **(less throttling).**

# Vertical auto-scaling approach:
## *Pluggable proactive portion*

NODE1
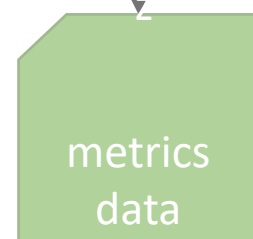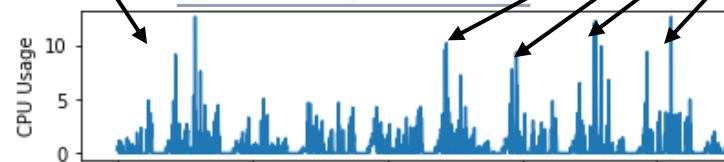limit: 3 cores

NODE2
limit: 3 cores

NODE3
limit: 3 cores

Reactive:
Handle the initial Pod size
fit and adjust to spikes

Proactive:
Combine CaaSPER with existing
algorithms to handle cyclical/
predictable loads over time

reads

CPU Usage

10

5

0

reads

decisions.txt

writes

metrics
data

Pluggable
recommender
algorithms
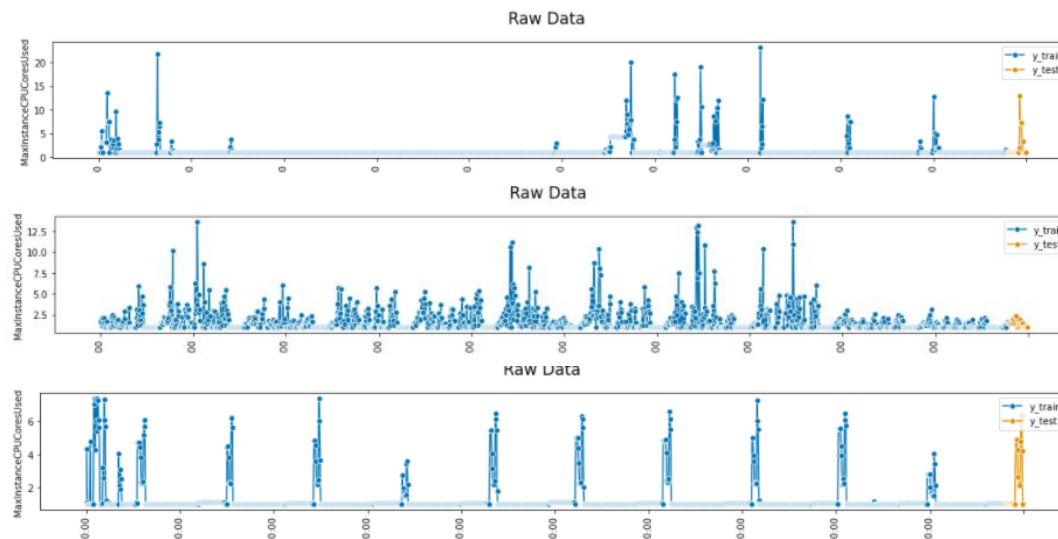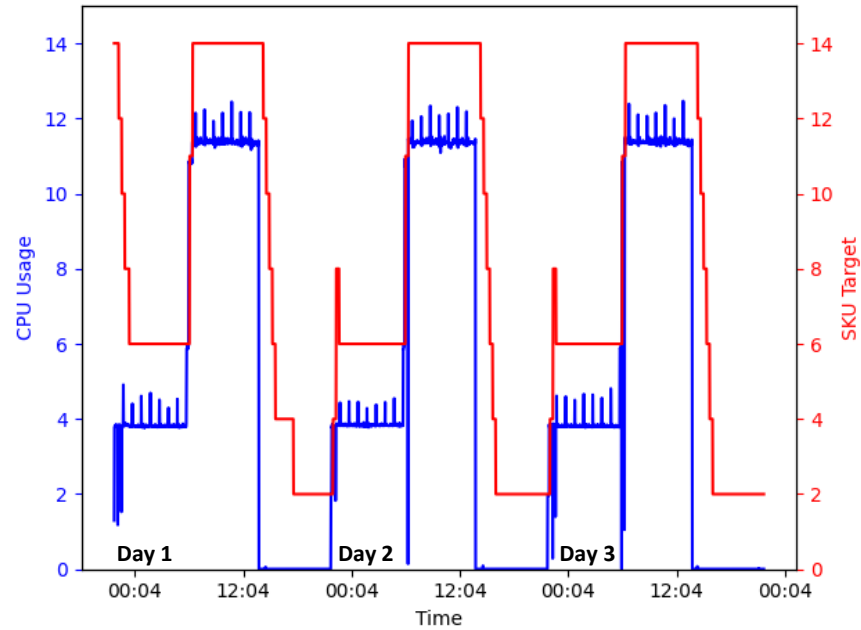
# Proactive:
*Real cyclical workload + Time series*

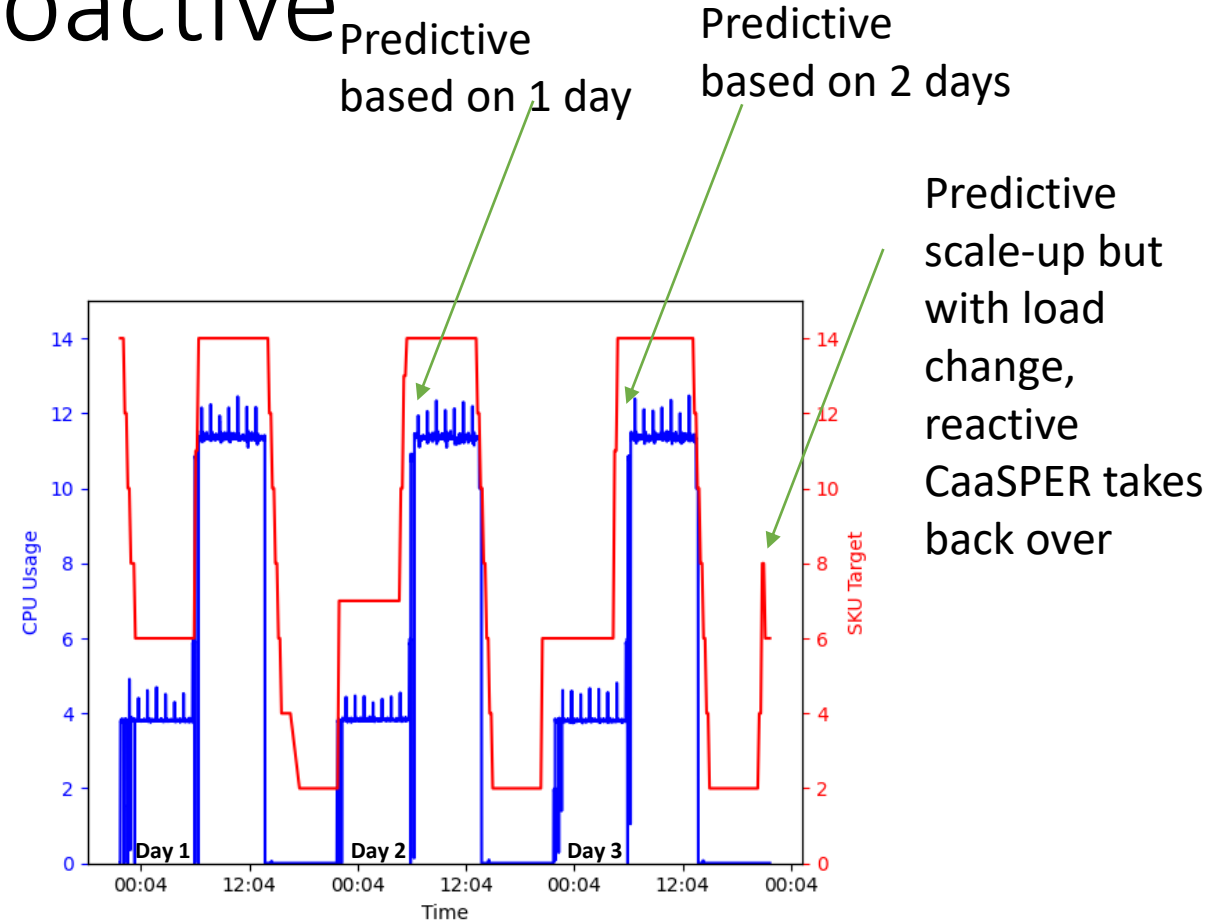- ## Started by looking at simulated + real workload CPU traces



- Experimented with many different algorithms and data preprocessing for prediction and measured throttling/fit/etc
- Naïve worked well for most of our scenarios, but can easily swap out
- Trade-off: complexity/robustness+debuggability

# Combining Reactive + Proactive



Predictive based on 1 day

Predictive based on 2 days

Predictive scale-up but with load change, reactive CaaSPER takes back over
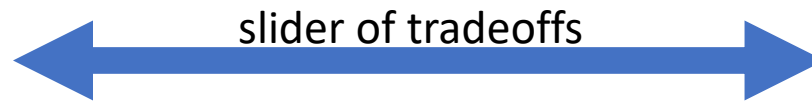
Reactive only

Reactive + Proactive
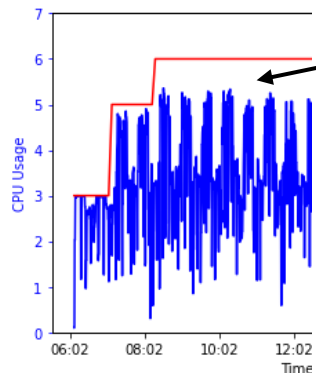
Red - CPU limit setting

Blue - actual CPU usage

# CaaSPER parameters

Users can specify preferences on a slider, or we can autotune in our simulator (next):
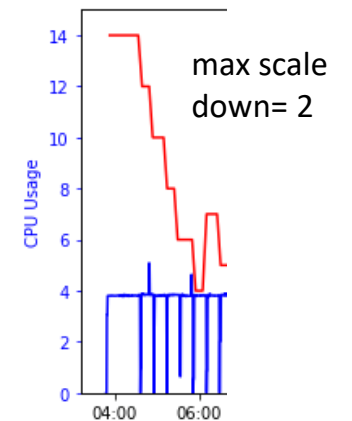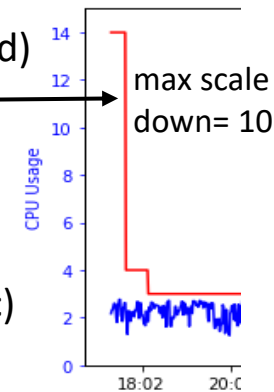
**More performant/More expensive**  slider of tradeoffs  **Less performant/More cost-effective**

## Impact of our parameters:



"slack" parameter is configurable. Here: 15%

- Slack (buffer between resources used and resources allocated)
- Max amount to auto-scale up/down
- How frequently to scale
- How much historical usage data to look at
- Balance between reactive vs proactive algorithm
- How early to be proactive (scale up 5 min early or 1 hour, etc)
- How frequently to scale
- Guardrails (ex: giant burst of traffic, how to behave)



max scale down= 10

max scale down= 2

# At this point: panicking.

- We had a paper deadline. We built an awesome algorithm, but tuning the 20+ parameters was challenging

- We needed to demonstrate our autoscaling algorithm for about 30 7-day long experiments to run, but we only had 3 functioning K8s clusters, and 10 days.

- Enter: VASIM

---

**Algorithm 1** CaaSPER autoscaling decision algorithm.

**Require:** $x_c$: CoreCount$_{cur}$
**Require:** $\{X_t\}$: Vector of workload CPU usage indexed by time (observed and/or predicted)
**Require:** $R$: System inputs (e.g., resource limit such as max CPU, price per core, granularity per core)
**Require:** $s_h$: High slope threshold
**Require:** $s_l$: Low slope threshold
**Require:** $m_h$: High slack threshold as percentage of capacity
**Require:** $m_l$: Low slack threshold as percentage of capacity
**Require:** SF$_h$: Maximum single step scale-up amount
**Require:** SF$_l$: Maximum single step scale-down amount
**Require:** $c_{min}$: Minimum resource requirements (scale-down lower bound)

1: **function** AUTOSCALE($x_c$, $\{X_t\}$)
2:     normalized cpu ← PREPROCESS CPU($\{X_t\}$)
3:     PvP curve ← SKU RECOMMENDATION TOOL(normalized CPU, $R$)
4:     PvP slopes ← CALCULATE SLOPES(PvP curve)
5:     skew ← CALCULATE SKEW(PvP slopes)
6:     $s$ ← GET CURRENT SLOPE(PvP slopes, $x_c$)
7:     SF ← CALCULATE SCALING FACTOR, SF($s$, skew)
8:     **if** $s \geq s_h$ or Quantile($\{X_t\}$) $\geq (1 - m_h) * x_c$ **then**
9:         SF ← min (SF, SF$_h$)
10:    **else if** $s \leq s_l$ or Quantile($\{X_t\}$) $\leq m_l * x_c$ **then**
11:        SF ← max ($-$SF, $-$SF$_l$)
12:    **else if** $s == 0$ and $x_c$ at top of PvP curve **then**
13:        SF ← UPDATE SCALING FACTOR(PvP curve, $x_c$)
14:    SF ← APPLY GUARDRAILS(SF, SF$_h$, SF$_l$, $c_{min}$, $R$)
15:    **return** SF

# VASIM: Vertical Autoscaling Simulator

VASIM replicates common components found in autoscaler architectures and replays CPU traces (real and estimated) with tunable parameters



*VASIM: Vertical Autoscaling Simulator Toolkit.*
*In IEEE International Conference on Data Engineering (ICDE 2024*

# VASIM: Vertical Autoscaling Simulator

You need 3 things: CPU Data, Autoscaling Algo, Parameters

```
TIMESTAMP,CPU_USAGE_ACTUAL
2023.04.02-00:09:00:000,7.2
2023.04.02-00:10:00:000,7.04
2023.04.02-00:11:00:000,6.88
2023.04.02-00:12:00:000,6.72
2023.04.02-00:13:00:000,6.48
2023.04.02-00:14:00:000,6.50
2023.04.02-00:15:00:000,6.52
2023.04.02-00:16:00:000,6.54
2023.04.02-00:17:00:000,6.56
```

```python
class SimpleAdditiveRecommender(Recommender):
    def __init__(self, cluster_state_provider, save_metadata=True):
        # Copy the code at the top of this function as-is.

        # Put your parameters here hard-coded, or pass them in to your
        # `metadata.json` file in the `algo_specific_config` section.
        self.my_param = self.algo_params.get("myparam", 2)

    def run(self, recorded_data):
        """
        This method runs the recommender algorithm and returns the new number of
            cores to scale to (new limit).

        Inputs:
            recorded_data (pd.DataFrame): The recorded metrics data for the current time window to
simulate
        Returns:
            latest_time (datetime): The latest time of the performance data.
            new_limit (float): The new number of cores to scale to.
        """

        # Your logic goes here! Look at the data in the `recorded_data` dataframe,
        #   do a calculation, and return the number of cores to scale to.

        return new_limit
```
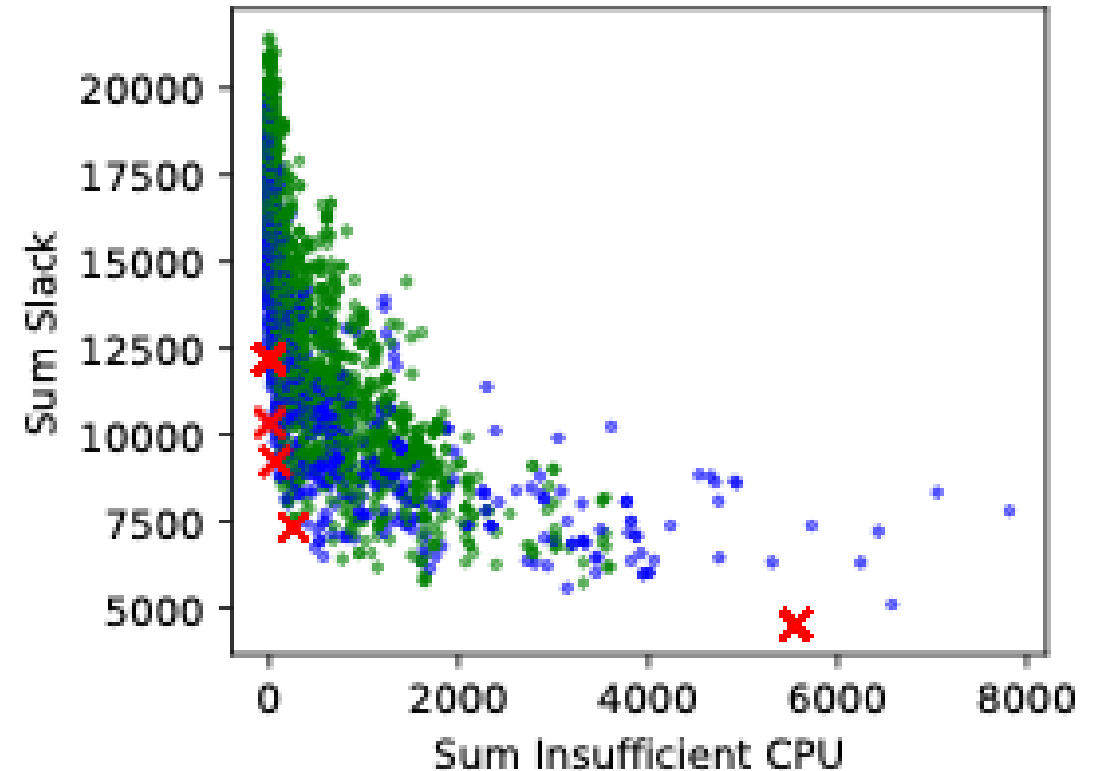
```json
"algo_specific_config": {
    "addend": 2
},
"general_config": {
    "window": 20,
    "lag": 10,
    "max_cpu_limit": 25,
    "min_cpu_limit": 2.0
},
```

# Simulating & Tuning parameters

When selecting parameters, we must find the ideal balance between:

- slack (resources wasted)
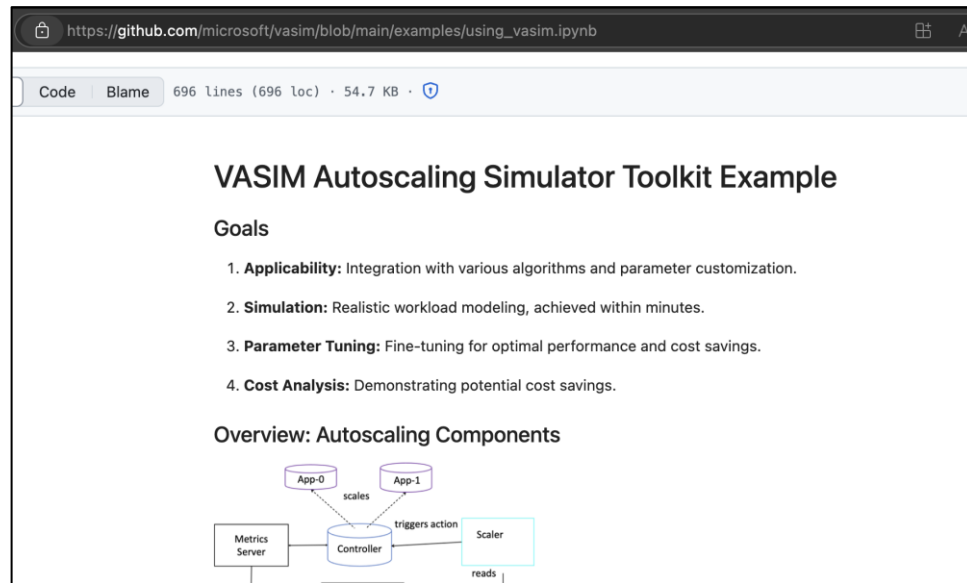- insufficient CPU (throttling)

# Outline

- Mechanism
  - Vertically scaling in-place
  - Cassandra perf impacts?
- Policy
  - CaaSPER: Proactive/Reactive algorithm for balancing price-perf trade-off
- Code/demo: VASIM - Vertical Autoscaling SIMulator
  - Try your own autoscaling algorithm!
  - Autotuning: parameter tune your own algorithm
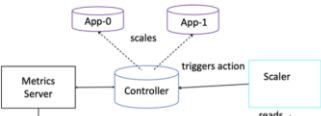  - How to get started with Cassandra

# (Go to GitHub…)

https://github.com/microsoft/vasim

# VASIM web demo

- Our [notebook](https://github.com/microsoft/vasim/blob/main/examples/using_vasim.ipynb)
https://github.com/microsoft/vasim/blob/main/examples/using_vasim.ipynb

- Together with Cassandra
https://github.com/microsoft/vasim/tree/kasaur/e2e-livedemo/examples/cassandra

- And the web front-end
https://github.com/microsoft/vasim/tree/main/examples/streamlit

# References

- Code repo: https://github.com/microsoft/vasim
  - Simulator demo: examples -> streamlit
  - Cassandra demo: examples -> cassandra

- Papers:
  - **VASIM: Vertical Autoscaling Simulator Toolkit** Anna Pavlenko, Karla Saur, Yiwen Zhu, Brian Kroth, Joyce Cahoon, Jesús Camacho Rodríguez. ICDE, 2024. [pdf]

  - **Vertically Autoscaling Monolithic Applications with CaaSPER: Scalable Container-as-a-Service Performance Enhanced Resizing Algorithm for the Cloud** Anna Pavlenko, Joyce Cahoon, Yiwen Zhu, Brian Kroth, Michael Nelson, Andrew Carter, David Liao, Travis Wright, Jesús Camacho Rodríguez, Karla Saur. SIGMOD, 2024. [pdf]