# On Optimizing Operational Efficiency in Storage Systems via Deep Reinforcement Learning

Sunil Srinivasa[1], Girish Kathalagiri[1], Julu Subramanyam Varanasi[2], Luis Carlos Quintela[1], Mohamad Charafeddine[1], ChiHoon Lee[1]

[1] Samsung SDS America

[2] Samsung Semiconductors Inc

# Motivation

- Data Centers as large sinkholes of energy.

- US data centers consumption is about 70 billion Kwh, 1.8% total US energy consumption (Lawrence Berkeley Lab study 2016), ~ $8B
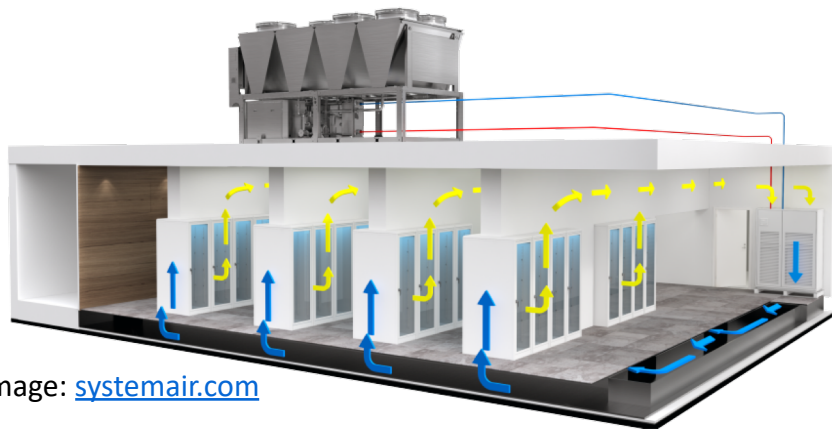


Image: systemair.com



Image: www.aashe.org

**Building-level**: Data Center energy usage optimization, optimizing water pumps, cooling towers, ..

**Server-level**: energy usage optimization per server, optimizing fans speed inside servers
**focus of this talk**

**SAMSUNG SDS** **SAMSUNG**

# Problem setup: operational efficiency of a storage server



Speed setting for the fans

24 temperature reading for each SSD drive

1 2 3 ... 24

24 SSDs solid state drives for storage

SAMSUNG SDS SAMSUNG

4

# Problem setup: operational efficiency of a storage server



Speed setting for the fans

24 temperature reading for each SSD drive

- High temperatures lower the performance of reading/writing
- High temperatures shorten the lifespan of the SSD

- Temperature varies based on the workload: read or write, how many KB, and how often.

- Current status quo control is rule-based.

We would like to use Deep Reinforcement Learning to represent the state of operation and learn the optimum policy for controlling the fan speeds, as a proxy of energy usage for cooling.

SAMSUNG SDS  SAMSUNG

# Contributions

- **Model-free approach**, no need to understand to SSD server behavior dynamics of workload and temperature

- We **trained on the real environment** (the server), not through a simulator

- We **designed a Reward function for this problem** to guide the algorithm towards the desired operation behavior (servers temperature & fans speeds)

**SAMSUNG SDS** SAMSUNG

# Reinforcement Learning

$$s_t \mid s_{t+1} \mid s_{t+2} \mid \ldots s_{t+H-1}$$

Action $a_t \mid a_{t+1} \mid a_{t+2} \mid \ldots a_{t+H-1}$

**RL Agent**

**Environment**

Perceived **State** of the environment $s_t$

*Policy $\pi(a_t|s_t)$*
to optimize expected
*return* $R = E\left(\sum_{k=0}^{H-1} \gamma^k r_{t(s_t, at)}\right)$*, i.e.,*
*expected sum of discounted rewards*

**Rewards** $r_t$

$$r_t \mid r_{t+1} \mid r_{t+2} \mid \ldots r_{t+H-1}$$

**SAMSUNG SDS** SAMSUNG

# Reinforcement Learning: Action Value, State Value, Advantage functions

- **Action Value Function $Q^\pi(s_t, a_t)$**: is the scalar **expected long-term return** achieved for following **Policy $\pi$** when in **state $s_t$** by performing **action $a_t$**. → How good is this state-action pair.

- **State Value Function $V^\pi(s_t)$:** is the scalar **expected long-term return** achieved for following **Policy $\pi$** when in **state $s_t$** averaged over all possible actions {$a$} from that state. → How good is this state $s_t$.

- **Advantage Function $A^\pi(s_t, a_t) = Q^\pi(s_t, a_t) - V^\pi(s_t)$** is the advantage of taking action $a_t$ when in state $s_t$ compared to the average of all possible actions from that state.

The Action and Value functions are unknown, and we need to **learn them by interacting with the environment**. One approach is through training Deep Neural Network(s) which takes the state as an input and outputs an approximation of Q and V.
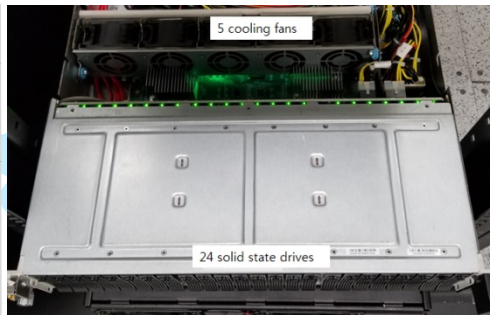
**SAMSUNG SDS** **SAMSUNG**

# Problem setup using Deep Reinforcement Learning

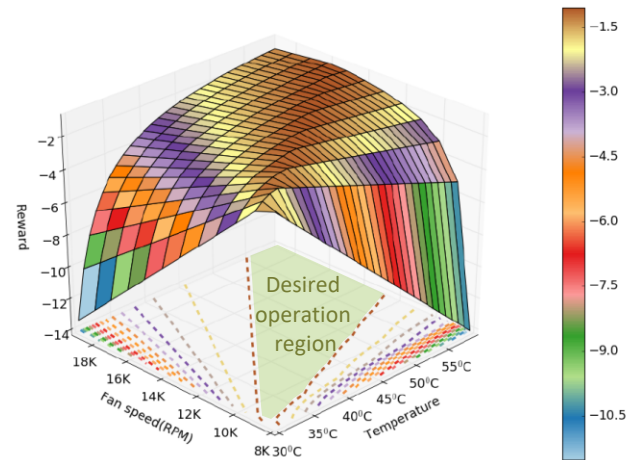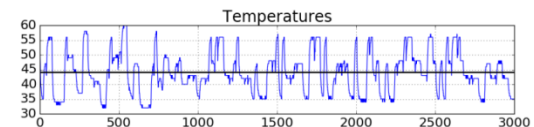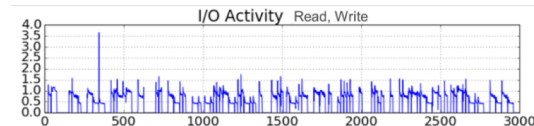**State**: function of workload (Reads or Writes) and temperature

**Actions:** Fan Speeds
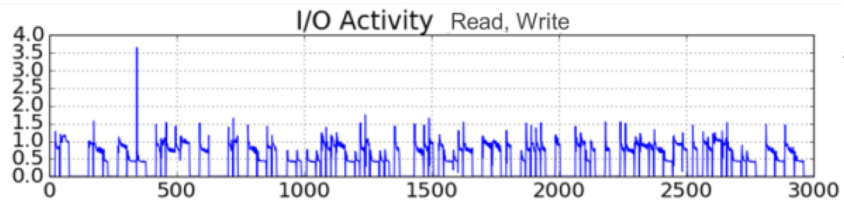


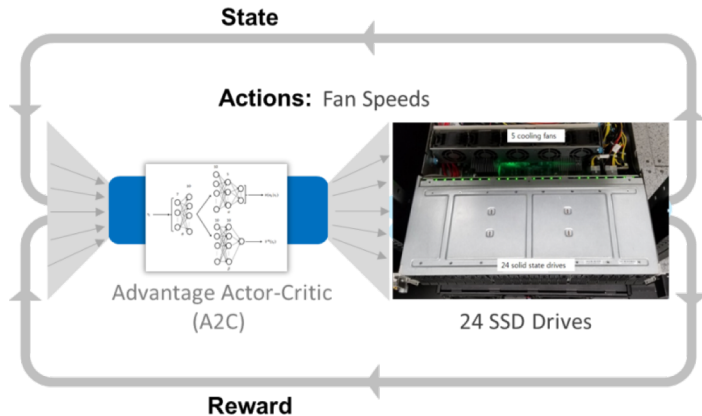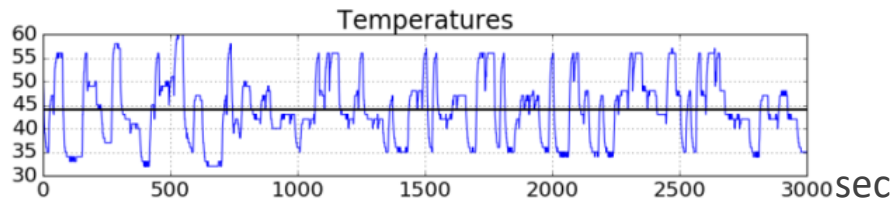Advantage Actor-Critic
(A2C)

24 SSD Drives

**Reward:** A function of Temperatures & Fans Speeds

# Problem Formulation: State
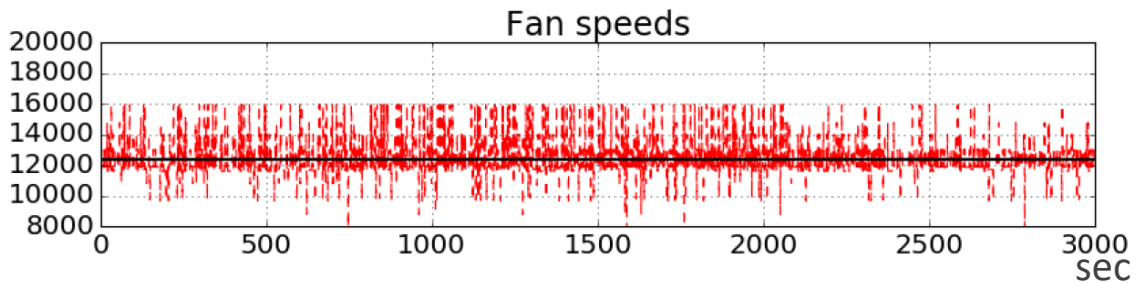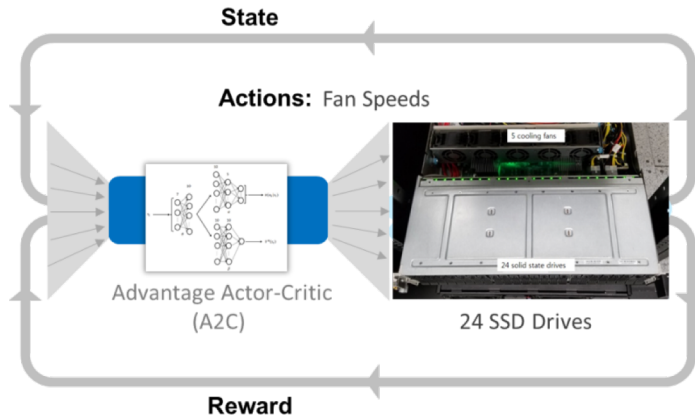


Time slot = 25 sec
Horizon = 10 slots

State =

1. $\Gamma$(**# of I/O requests** to the server **per second**)
2. $\Gamma$(**# of Kbytes read**, averaged over all SSDs per sec)
3. $\Gamma$(**# of Kbytes written**, averaged over all SSDs per sec)
4. $\Gamma$(**# of Kbytes read** in **previous time slot**, averaged over all SSDs)
5. $\Gamma$(**# of Kbytes written** in **previous time slot**, avg. over all SSDs)
6. $\Gamma$(**Mean Temperate**: averaged over all SSDs)
7. $\Gamma$(**Mean Fan Speed**: averaged over all 5 cooling fans)

Read via `iostat`

Read via `ipmi-sensors`

Where normalization operation per field is $\Gamma(x) = \frac{x - minX}{maxX - minX}$

# Problem Formulation: Actions



**Actions:** Fan Speeds

Advantage Actor-Critic (A2C)

24 SSD Drives

**State**

**Reward**

Fan speeds

We tried 2 approaches

## Raw Actions

Size of Actions space = **7**

A = {6, 8, 10, .., 18}
Kilo rpm (revolutions per minute)

## Incremental Actions

Size of Actions space = **3**

A = {Decrease by 1000,
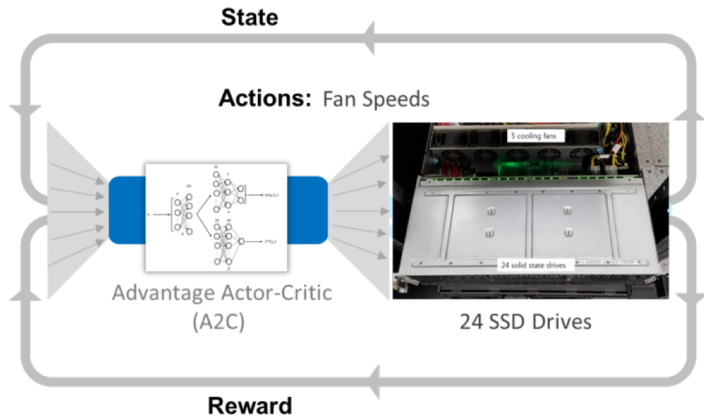
No change,

Increase by 1000} rpm

Allows for smoother transitions
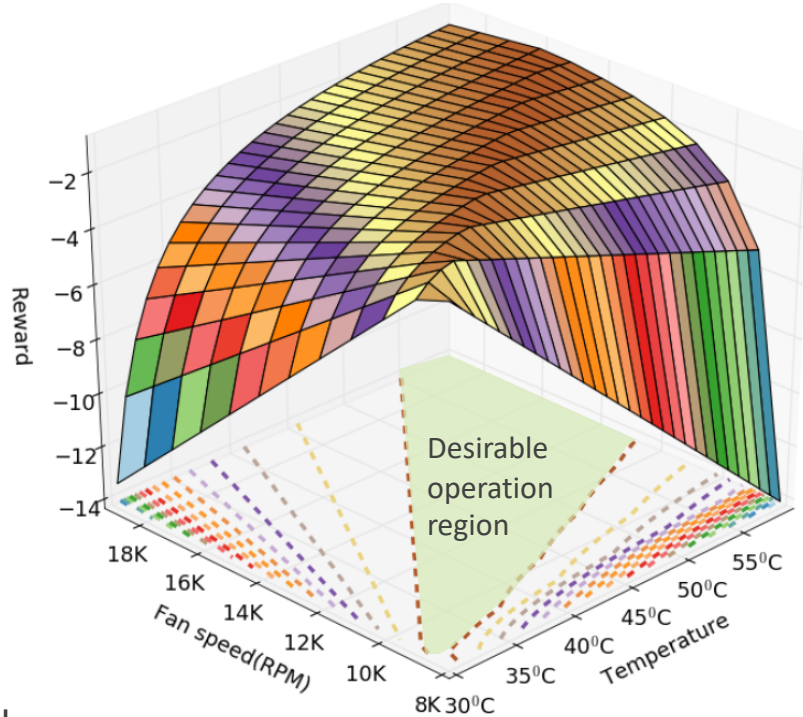
Action is applied to all 5 fans via the `ipmitool` command

# Problem Formulation: **Rewards**



**State**

**Actions:** Fan Speeds

Advantage Actor-Critic
(A2C)

24 SSD Drives

**Reward**



Desirable operation region

**Reward shaping** is an important design step which requires domain knowledge and guides the RL agent towards its optimal behavior

**For this RL application**, a good reward function should have:
1) High reward for low fan speeds and low temperature
2) Low reward for high fan speeds and low temperature – as this is wasting energy
3) Low reward for low fan speed and high temperature – as this will damage the SSDs

$$R = -max\left(\frac{\Gamma(T)}{\Gamma(F)}, \frac{\Gamma(F)}{\Gamma(T)}\right)$$

We chose this design above.
T is the mean temp across all SSDs.
F is the mean speed across fans.

**SAMSUNG SDS** **SAMSUNG**

# Actor Critic (A2C) Dueling Network Architecture

A NN with 2 branches: one for the Policy evaluation (what action to do: actor), and one for policy estimation (how good is it being in such state: critic)



*Policy Neural Network*

**Actor**: Approximates optimal policy, of choosing Best action for state $s$

$\pi(a_t|s_t)$

**State**:
Vector of length 7

$s_t$

*Value Function Neural Network*

$V^{\pi}(s_t)$

**Critic**: approximates the state-value function to estimate how attractive is it being at state $s$

Other option is to estimate the Actor in one NN, and the Critic in another NN. Our approach has the advantage of having: (1) fewer parameters, (2) mitigates overfitting one function over the other, (3) faster convergence, and (4) more accurate approximations for the actor & the critic.

# A2C algorithm pseudo-code

---

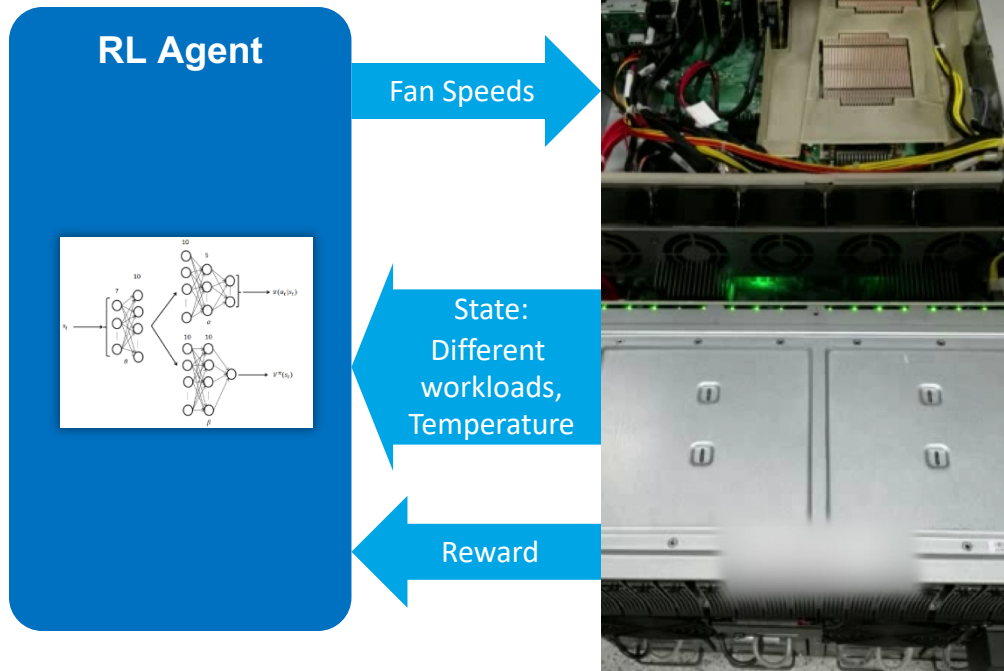**Algorithm 1** Advantage Actor-Critic (A2C) - pseudocode

---

// Notate shared parameters by $\theta$ and actor- and critic-specific parameters by $\alpha$ and $\beta$, respectively.

// Assume same learning rates $\eta$ for $\theta$, $\alpha$ as well as $\beta$. In general, they may all be different.

Initialize $\theta$, $\alpha$ and $\beta$ via uniformly distributed random variables.

**repeat**

    Reset gradients $d\theta = 0$, $d\alpha = 0$ and $d\beta = 0$.

    Sample $N$ trajectories $\tau_1, \ldots, \tau_N$ under the (current) policy $\pi(\cdot; (\theta, \alpha))$.

    $i = 1$

    **repeat**

        $t_{\text{start}} = t$

        Obtain state $s_t$

        **repeat**

            Perform action $a_t$ sampled from policy $\pi(a_t|s_t; (\theta, \alpha))$.

            Receive reward $r_t(s_t, a_t)$ and new state $s_{t+1}$

            $t \leftarrow t + 1$

        **until** $t - t_{\text{start}} = H$

        $i \leftarrow i + 1$

        Initialize $R$: $R = V(s_t; (\theta, \beta))$

        **for** $i \in \{t - 1, \ldots, t_{\text{start}}\}$ **do**

            $R \leftarrow r_i(s_i, a_i) + \gamma R$

            Sum gradients w.r.t $\theta$ and $\alpha$: // gradient ascent on the actor parameters

            $d\theta \leftarrow d\theta + \nabla_\theta \log \pi(a_i|s_i; (\theta, \alpha)) (R - V(s_i; (\theta, \beta)))$

            $d\alpha \leftarrow d\alpha + \nabla_\alpha \log \pi(a_i|s_i; (\theta, \alpha)) (R - V(s_i; (\theta, \beta)))$

            Subtract gradients w.r.t $\beta$ and $\theta$: //gradient descent on the critic parameters

            $d\theta \leftarrow d\theta - \nabla_\theta (R - V(s_i; (\theta, \beta)))^2$

            $d\beta \leftarrow d\beta - \nabla_\beta (R - V(s_i; (\theta, \beta)))^2$

        **end for**

    **until** $i = N$

    // Optimize parameters

    Update $\theta$, $\alpha$ and $\beta$: $\theta \leftarrow \theta + \eta d\theta$, $\alpha \leftarrow \alpha + \eta d\alpha$, $\beta \leftarrow \beta + \eta d\beta$.
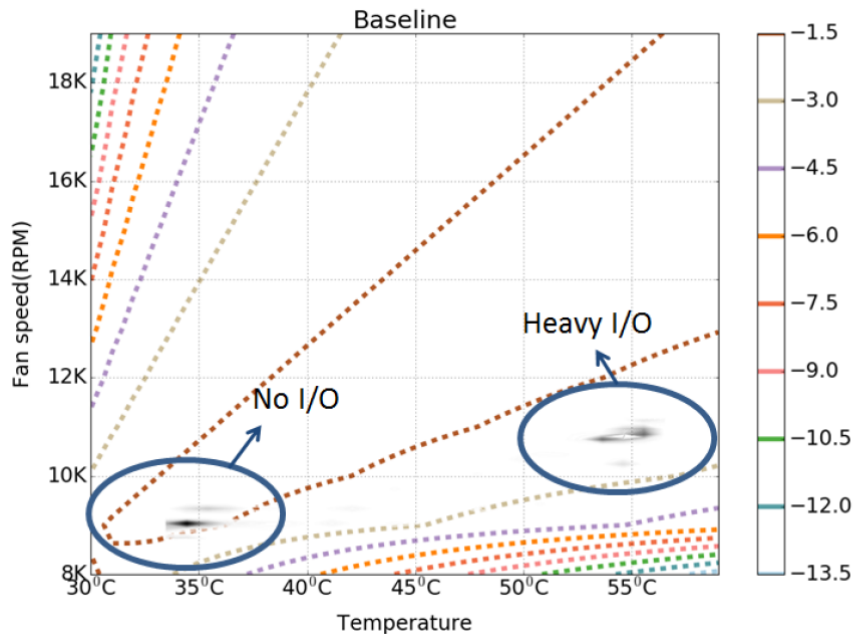
**until** convergence.

---

SAMSUNG SDS SAMSUNG

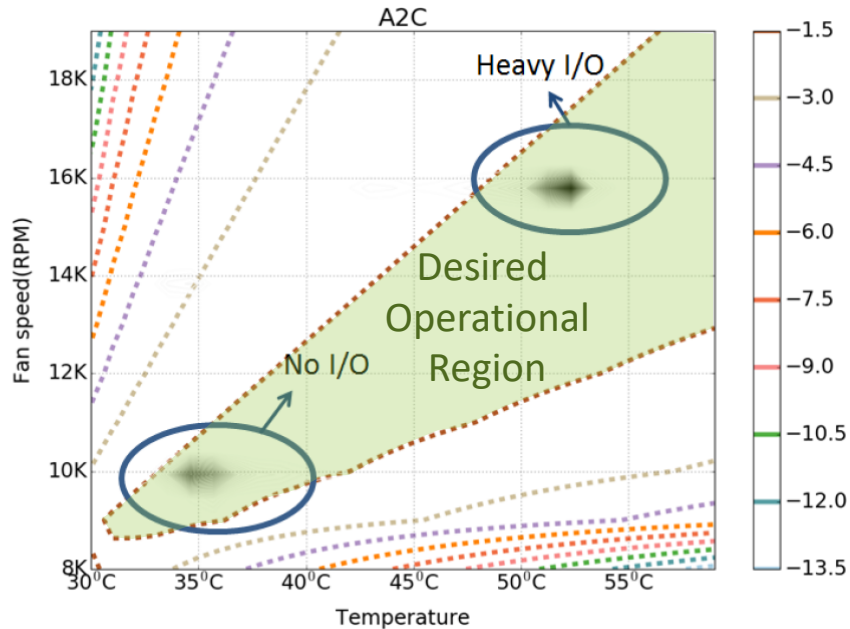# Learning… over few days



- Learning directly on the real environment (no simulator)

- Model-free: does not require any knowledge of the SSD server behavior dynamics

- Exposed to different stochastic workloads

# Performance for Idle Vs Heavy Periodic I/O workloads on the operational contours
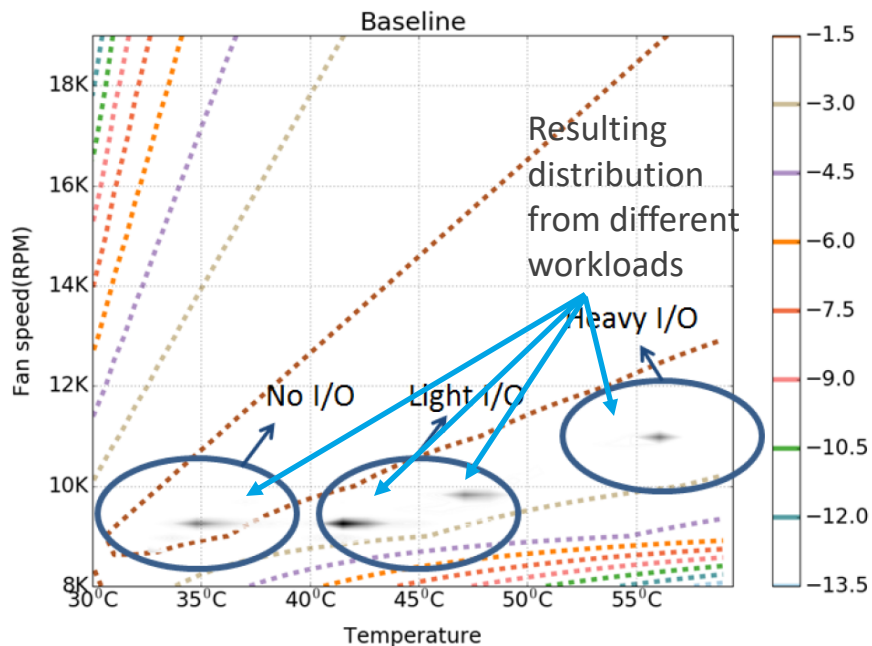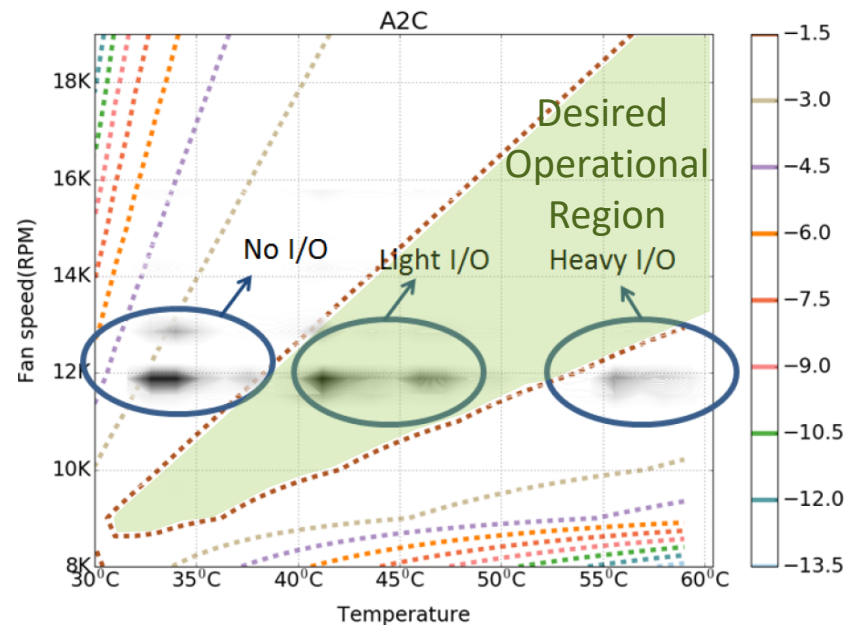


Status Quo controller

Using Deep Reinforcement Learning with Raw Actions

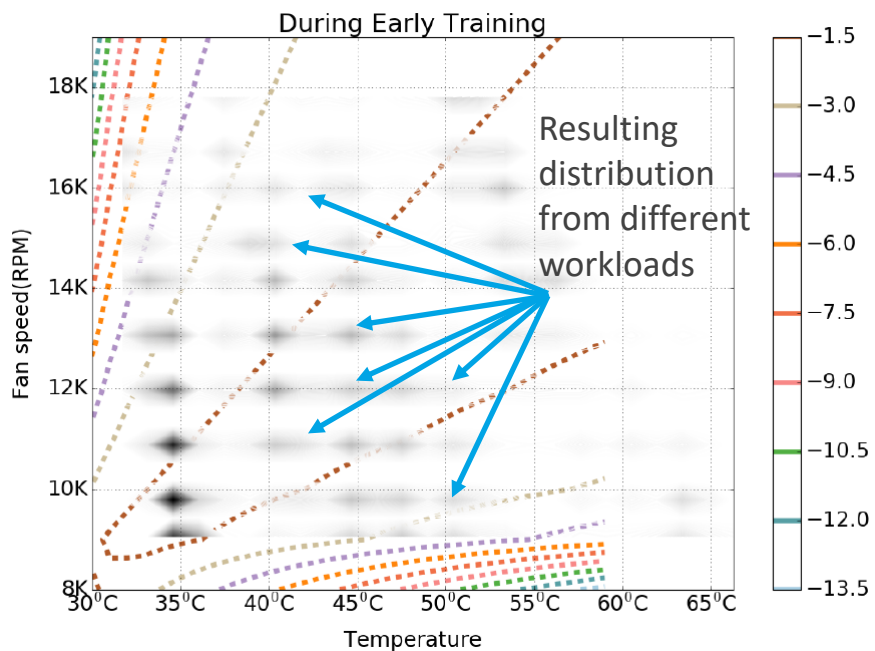# Performance for different stochastic workloads – Attempt 1



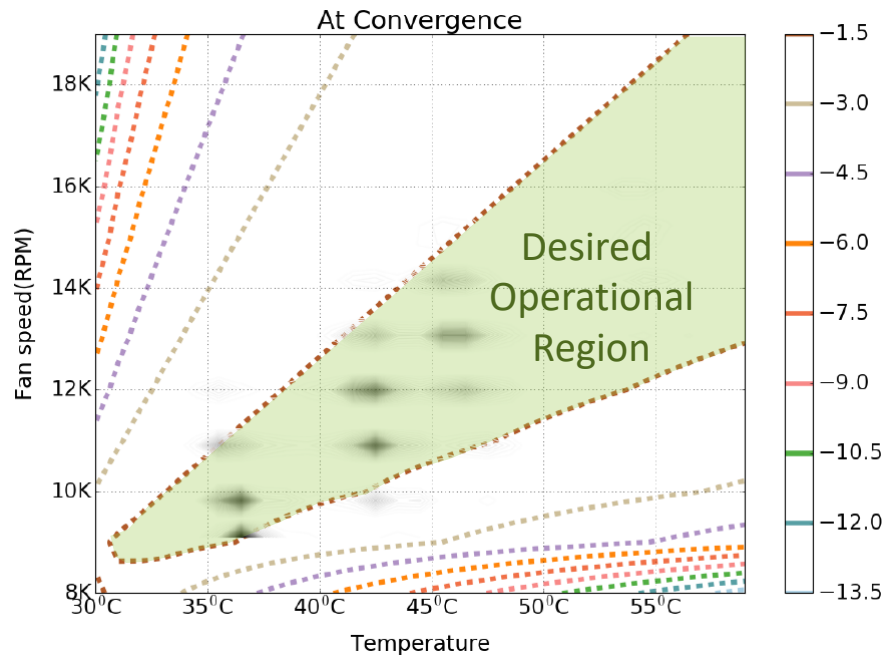At the beginning of training, algorithm is exploring and learning

Using Deep Reinforcement Learning with Raw Actions – Suboptimal performance, likely due to insufficient exploration

# Performance for different stochastic workloads – Attempt 2



**During Early Training**

Resulting distribution from different workloads

**At Convergence**

Desired Operational Region

At the beginning of training, algorithm is exploring and learning

Once finished learning right policy, operational behavior is within desired region. Used DRL with incremental actions.

insight to !nspiration

SAMSUNG    SAMSUNG SDS