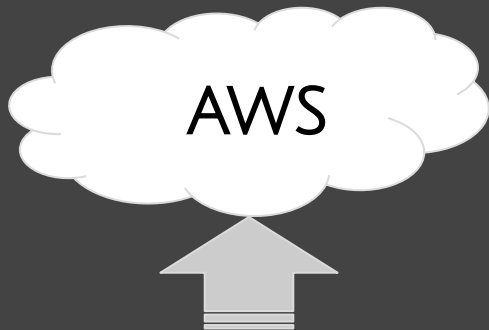
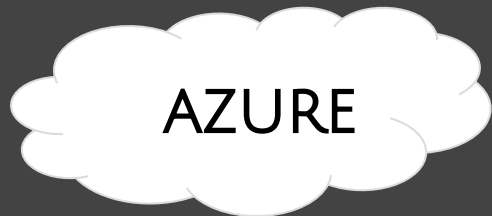
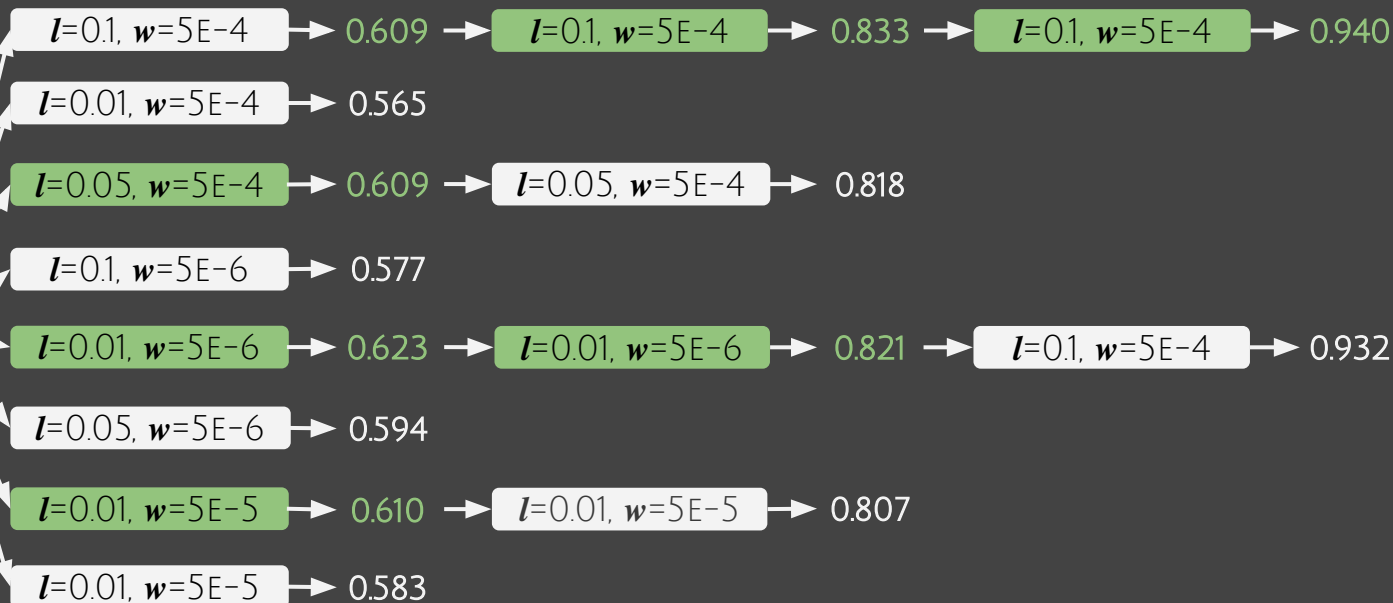


RUBBERBAND

UJVAL MISRA*, RICHARD LIAW*, LISA DUNLAP, ROMIL BHARDWAJ, KIRTHEVASAN
KANDASAMY, JOSEPH E. GONZALEZ, ION STOICA, ALEXEY TUMANOV

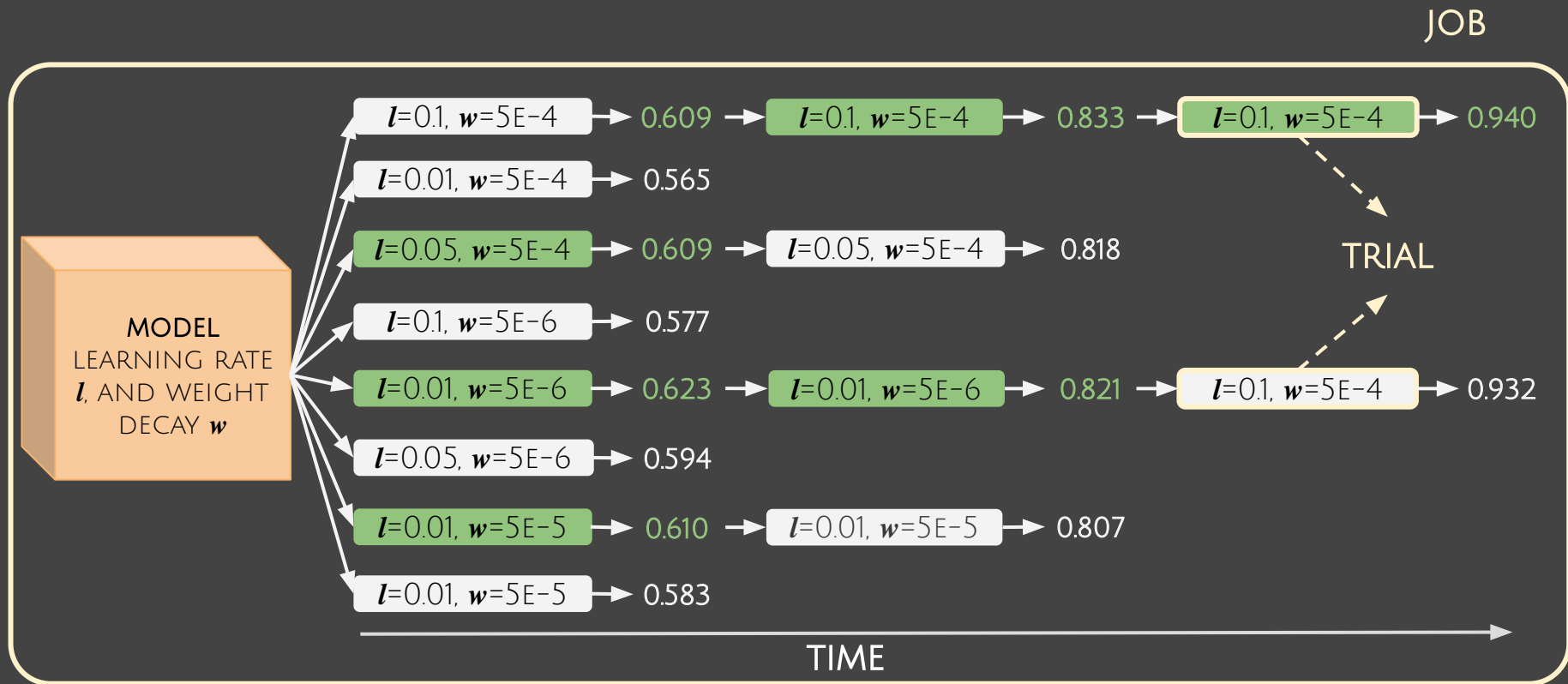


MODEL
LEARNING RATE
 l , AND WEIGHT
DECAY w

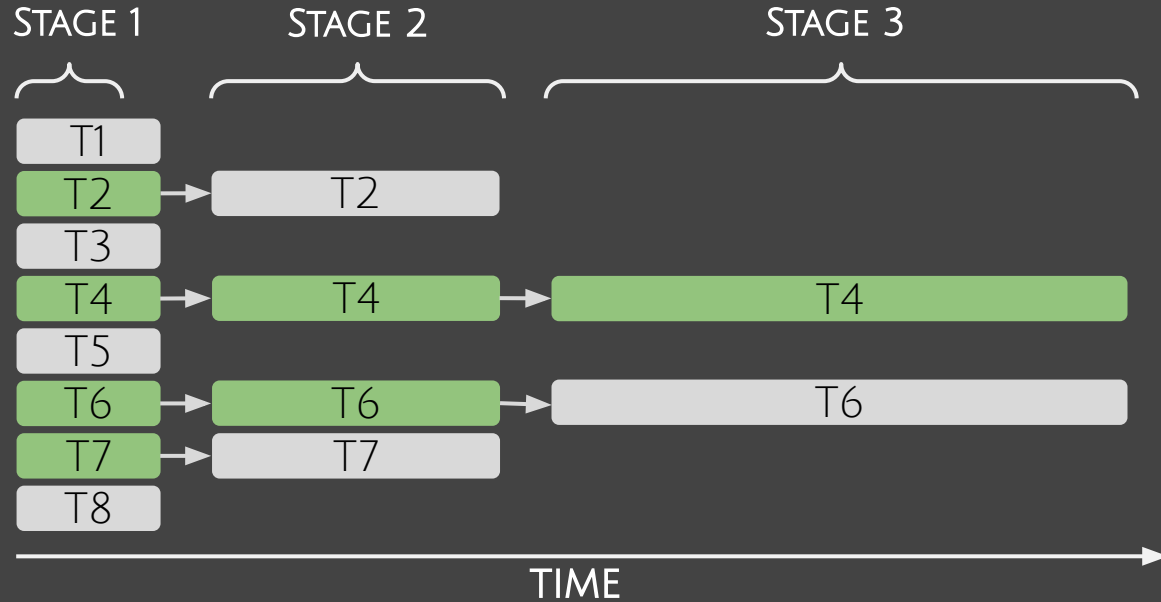


TIME

DISTRIBUTED, PARALLEL HYPERPARAMETER TUNING

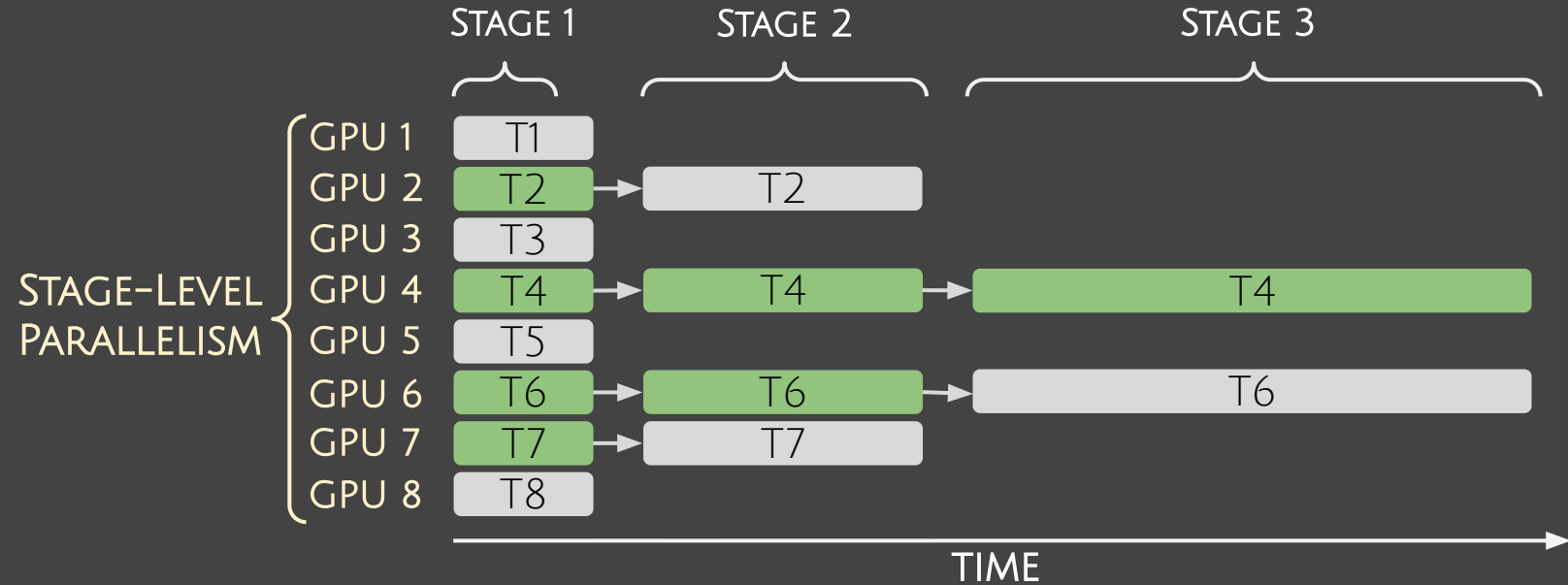


DISTRIBUTED, PARALLEL HYPERPARAMETER TUNING



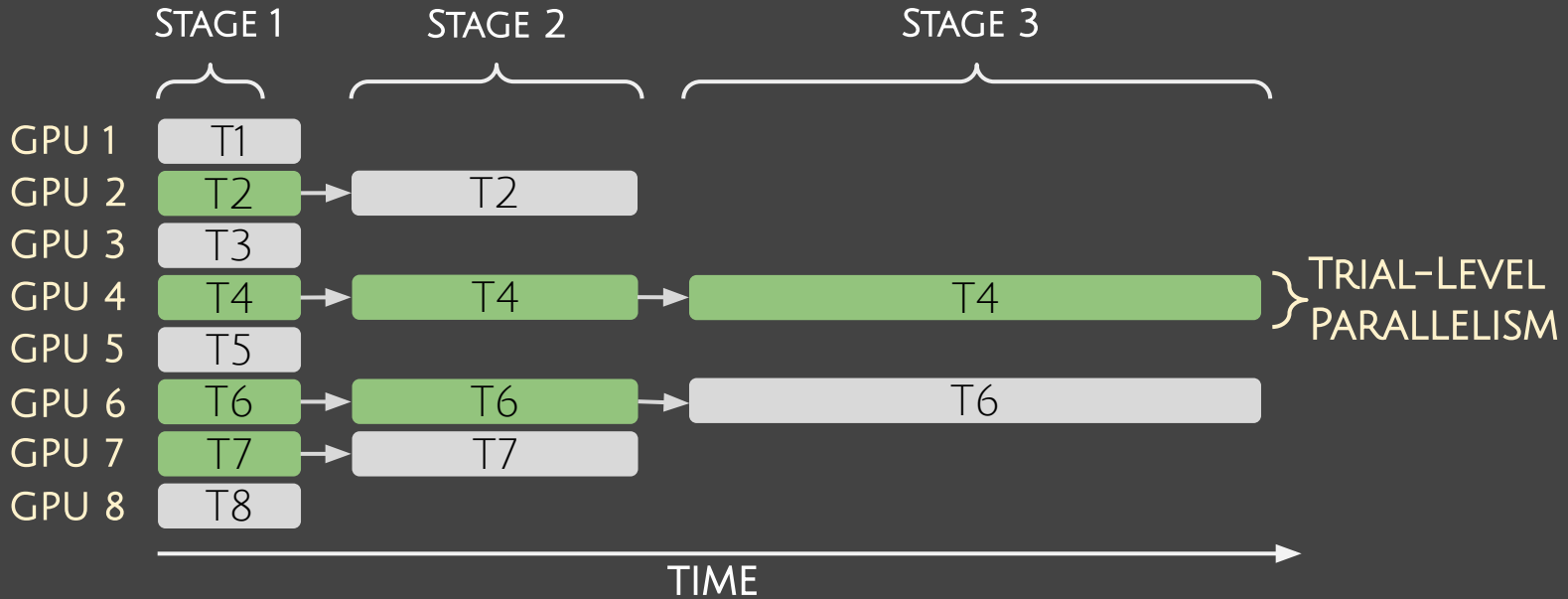
SUCCESSIVE HALVING ALGORITHM (SHA)

DISTRIBUTED, PARALLEL HYPERPARAMETER TUNING

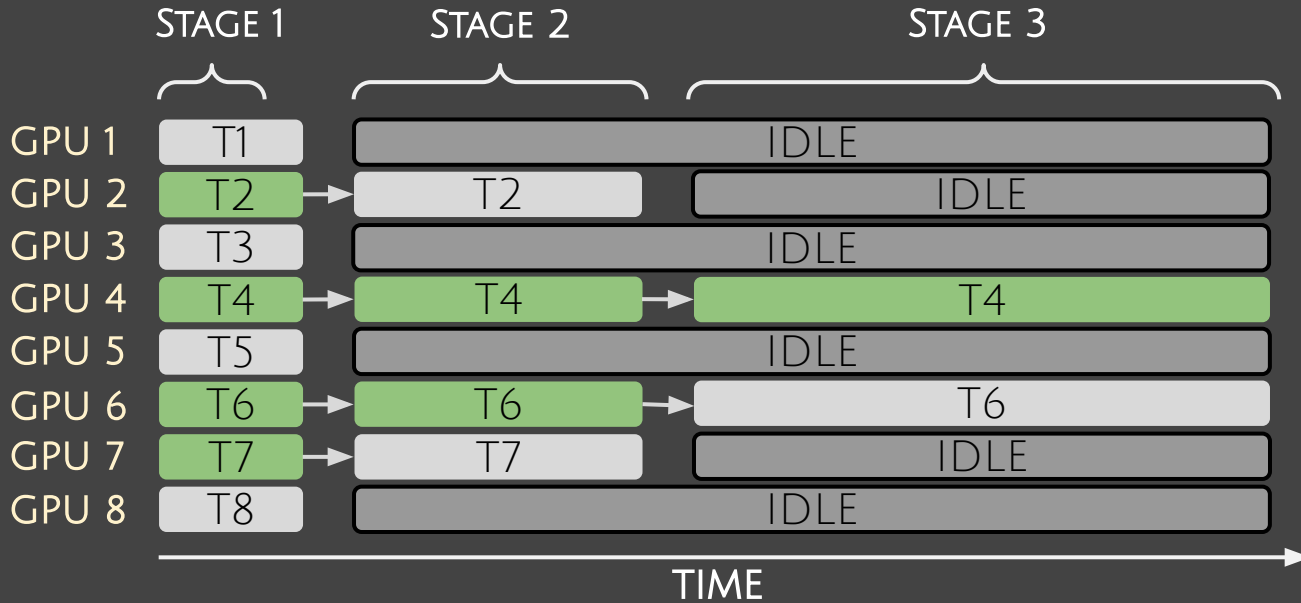


SUCCESSIVE HALVING ALGORITHM (SHA)

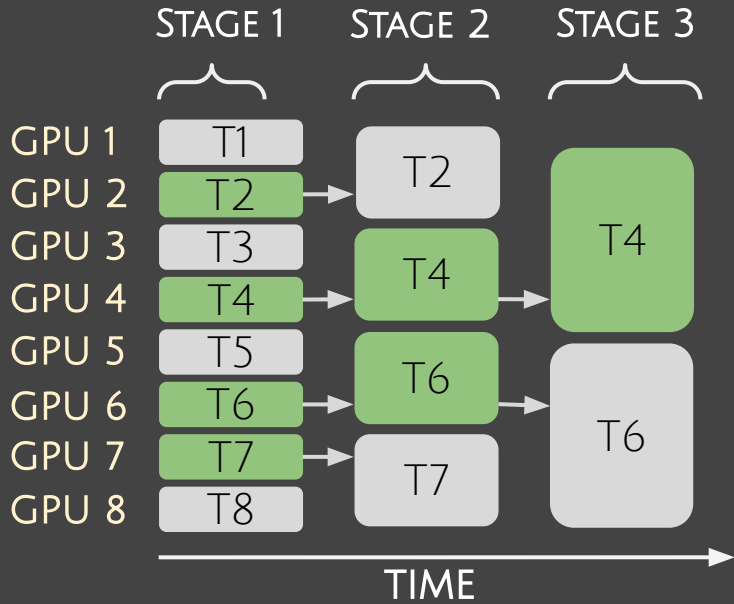
DISTRIBUTED, PARALLEL HYPERPARAMETER TUNING



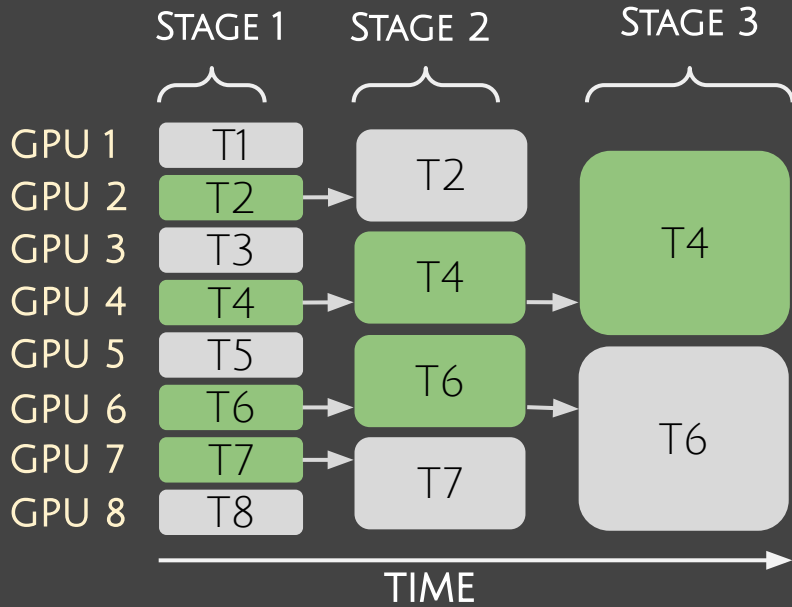
DISTRIBUTED, PARALLEL HYPERPARAMETER TUNING



DISTRIBUTED, PARALLEL HYPERPARAMETER TUNING

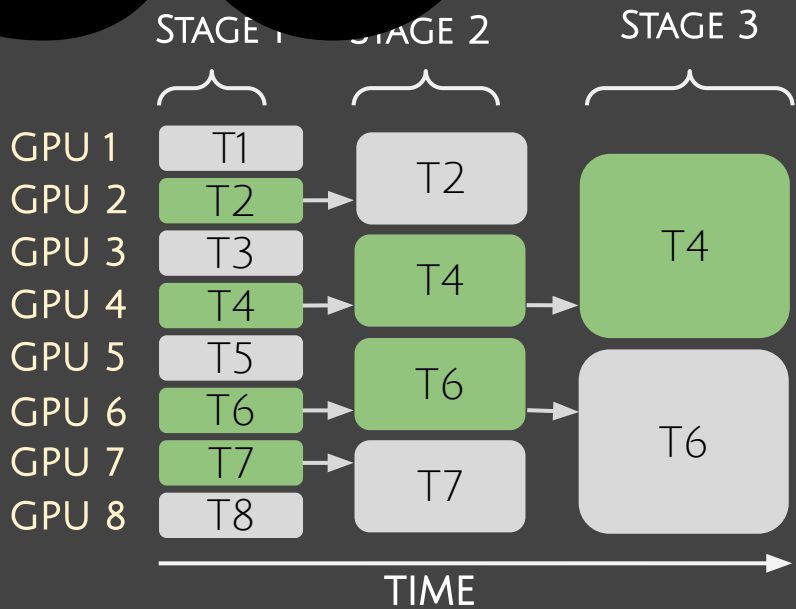


DISTRIBUTED, PARALLEL HYPERPARAMETER TUNING



SO

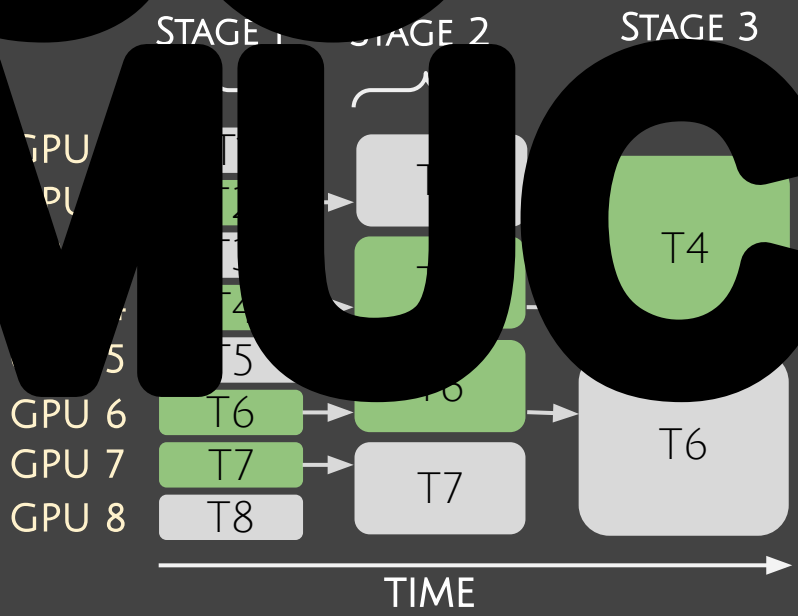
SHRINKED, PARALLEL HYPERPARAMETER TUNING



SO

RESTRICTED, PARALLEL HYPERPARAMETER TUNING

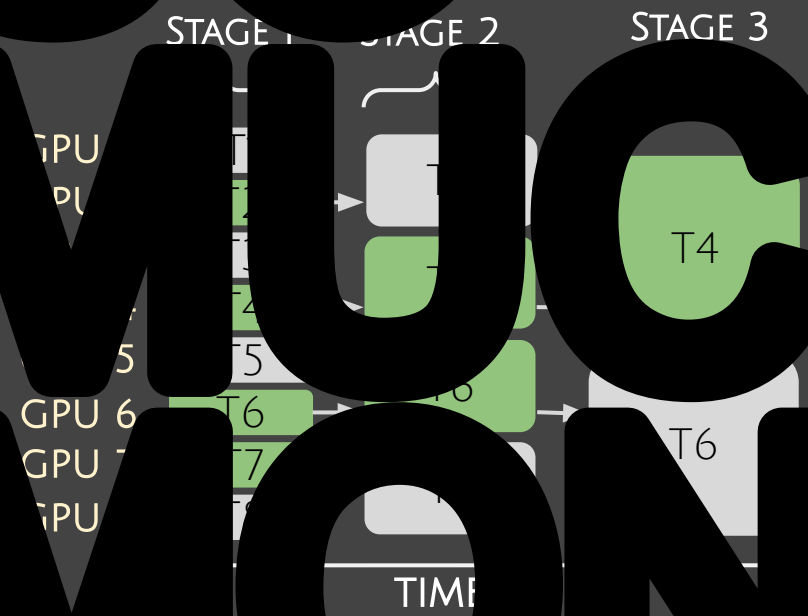
MUCH



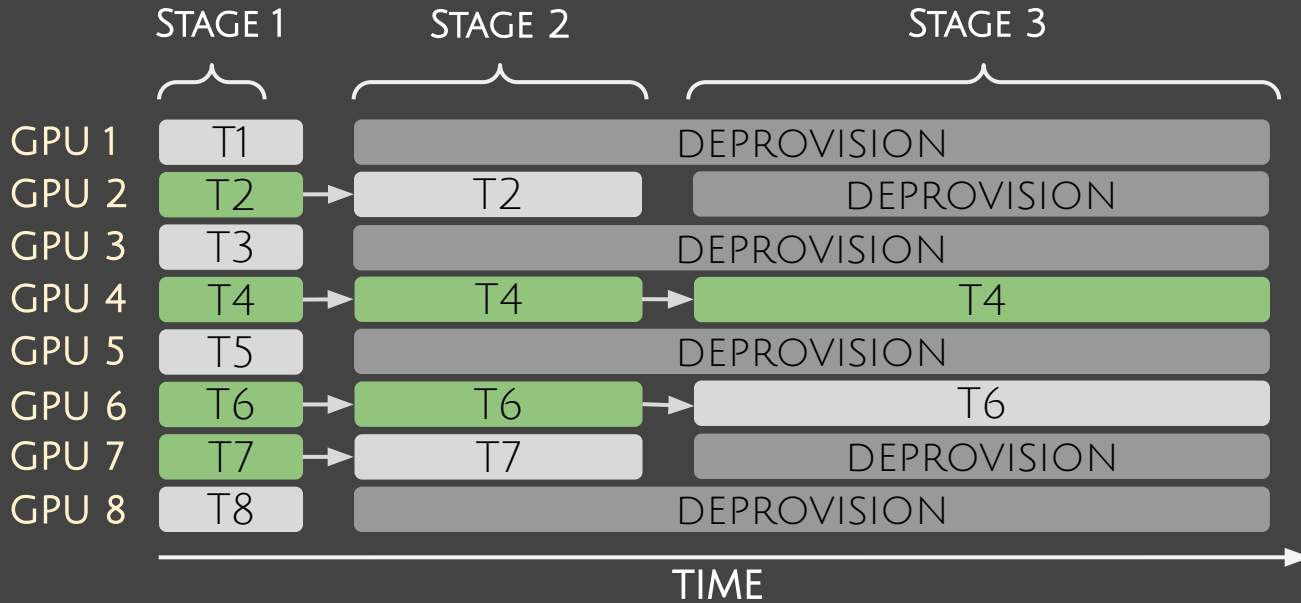
SO

DISTRIBUTED, PARALLEL HYPERPARAMETER TUNING

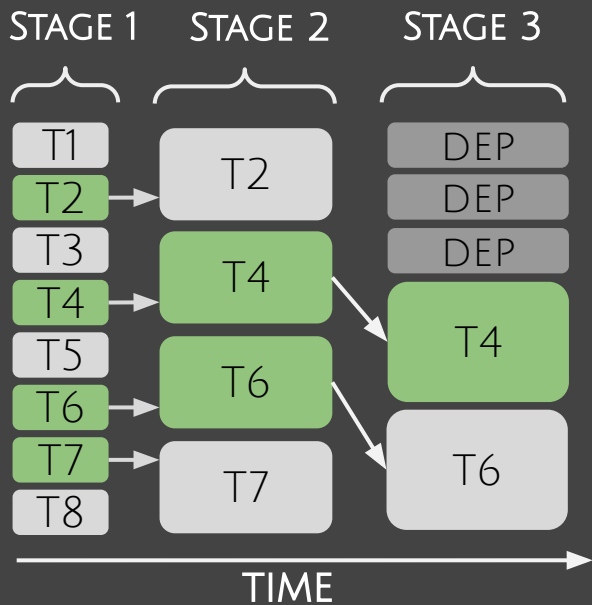
MUCH MONEY



DISTRIBUTED, PARALLEL HYPERPARAMETER TUNING



GIVEN A TIME CONSTRAINT, MINIMIZE THE COST OF EXECUTING A HYPERPARAMETER TUNING JOB.



STAGE	EPOCHS	TRIALS	GPUS/TRIAL
1	0-4	8	1
2	5-12	4	2
3	13-28	2	3

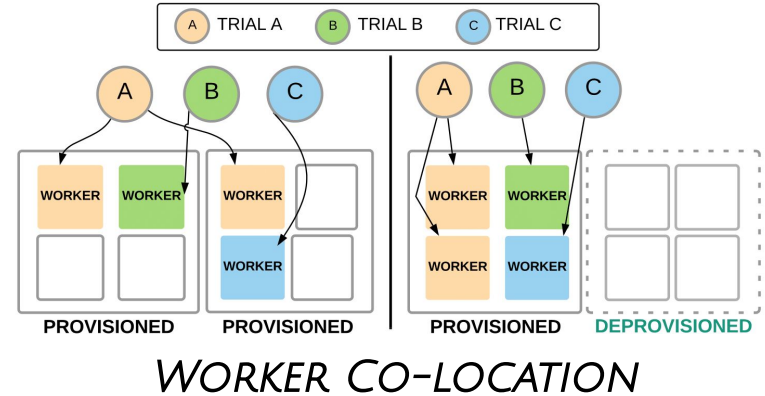
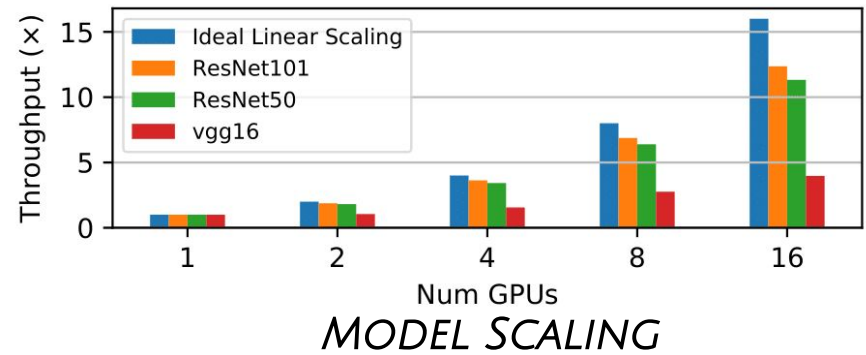
ALLOCATION PLAN

CHALLENGES

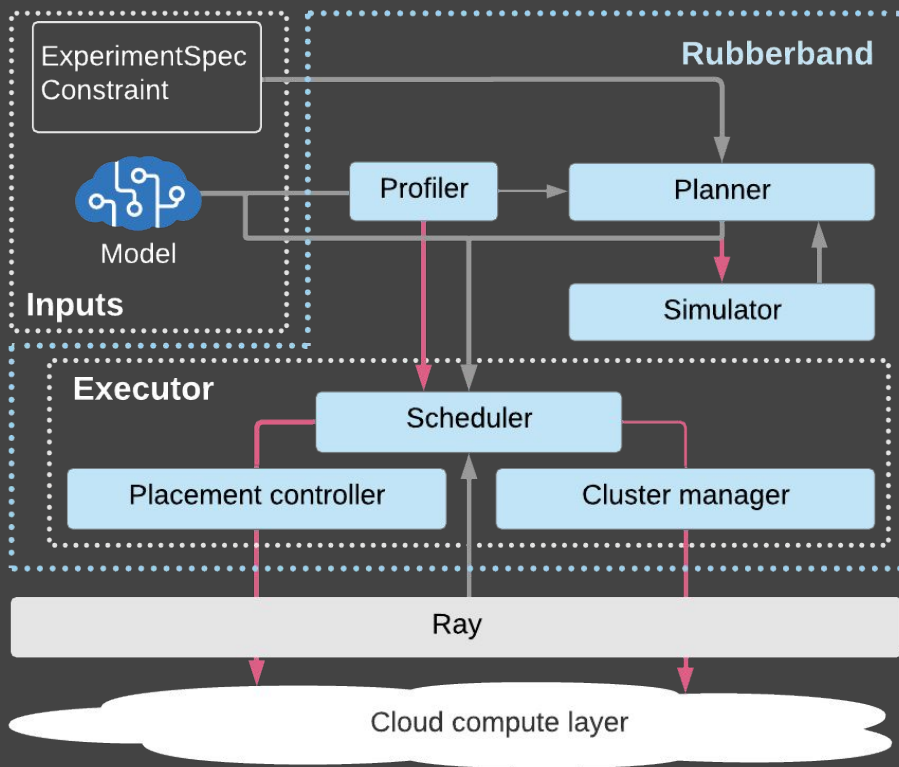
HOW CAN WE MODEL THE JOB COMPLETION TIME AND COST OF THE GIVEN ALLOCATION PLAN?

HOW CAN WE GENERATE A LOW COST ALLOCATION PLAN THAT COMPLETES ON TIME?

HOW CAN WE SCHEDULE SAID ALLOCATION TO OPTIMIZE WORKER CO-LOCATION + CLUSTER UTILIZATION?



RUBBERBAND



CHALLENGES

HOW CAN WE MODEL THE JOB COMPLETION TIME AND COST OF THE GIVEN ALLOCATION PLAN?

HOW CAN WE GENERATE A LOW COST ALLOCATION PLAN THAT COMPLETES ON TIME?

HOW CAN WE SCHEDULE SAID ALLOCATION TO OPTIMIZE WORKER CO-LOCATION + CLUSTER UTILIZATION?

RUBBERBAND

COST/PERFORMANCE MODEL VIA PROFILING DL MODEL TRAINING LATENCY AND PROVISIONING OVERHEADS

DAG-BASED EXECUTION MODEL WHICH FINDS FEASIBLE AND COST-EFFICIENT RESOURCE ALLOCATIONS

FULL-STACK SYSTEM FOR PLACEMENT, SCHEDULING, AND SCALING

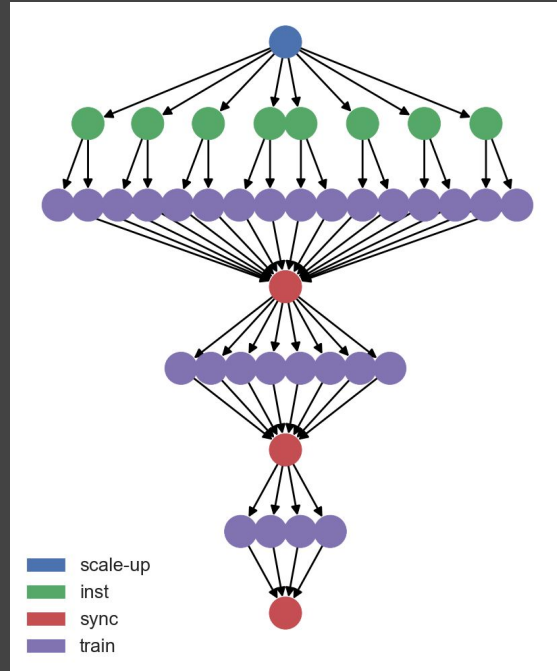
ISSUE 1: MODELING JOB COMPLETION TIME AND COST

PERFORMANCE MODELING

- TRAINING LATENCY
- PROVIDER QUEUING DELAY
- INSTANCE INITIALIZATION LATENCY

COST MODELING

- COMPUTE PRICE
- BILLING GRANULARITY
- DATA PRICE



ISSUE 2: FINDING A LOW-COST ALLOCATION PLAN

ISSUE 2: FINDING A LOW-COST ALLOCATION PLAN

STEP 1: GENERATE CANDIDATES

ISSUE 2: FINDING A LOW-COST ALLOCATION PLAN

STEP 1: GENERATE CANDIDATES

STEP 2: USE SIMULATOR TO PREDICT JOB COMPLETION TIME

ISSUE 2: FINDING A LOW-COST ALLOCATION PLAN

STEP 1: GENERATE CANDIDATES

STEP 2: USE SIMULATOR TO PREDICT JOB COMPLETION TIME

STEP 3: GREEDILY SELECT BEST CANDIDATE

ISSUE 2: FINDING A LOW-COST ALLOCATION PLAN

STEP 1: GENERATE CANDIDATES

STEP 2: USE SIMULATOR TO PREDICT JOB COMPLETION TIME

STEP 3: GREEDILY SELECT BEST CANDIDATE

MAXIMIZE *COST-MARGINAL BENEFIT*:

$$M = \frac{\text{COST OF CURRENT BEST PLAN} - \text{COST OF PROPOSED PLAN}}{\text{JCT OF PROPOSED PLAN} - \text{JCT OF CURRENT BEST PLAN}}$$

ISSUE 2: FINDING A LOW-COST ALLOCATION PLAN

STEP 1: GENERATE CANDIDATES

STEP 2: USE SIMULATOR TO PREDICT JOB COMPLETION TIME

STEP 3: GREEDILY SELECT BEST CANDIDATE

STEP 4: ITERATE WITH NEW BEST CANDIDATE

MAXIMIZE *COST-MARGINAL BENEFIT*:

$$M = \frac{\text{COST OF CURRENT BEST PLAN} - \text{COST OF PROPOSED PLAN}}{\text{JCT OF PROPOSED PLAN} - \text{JCT OF CURRENT BEST PLAN}}$$

ISSUE 3: EFFECTIVELY EXECUTE ALLOCATION PLAN

END OF STAGE



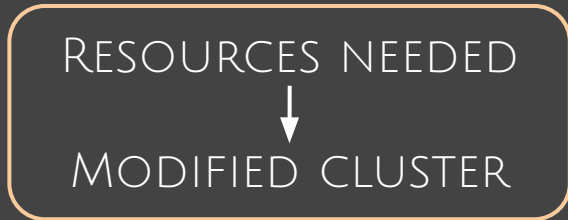
SCHEDULER



PLACEMENT CONTROLLER



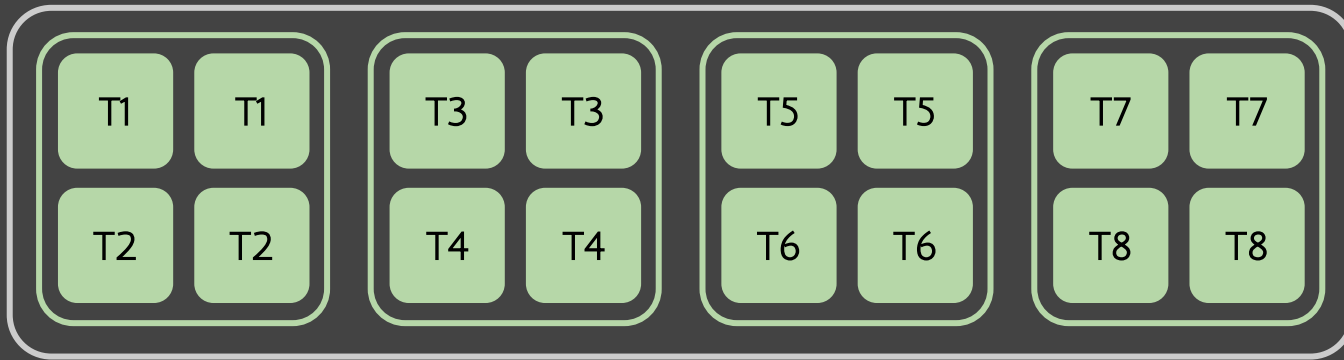
CLUSTER
MANAGER



ISSUE 3: EFFECTIVELY EXECUTE ALLOCATION PLAN

STAGE	EPOCHS	TRIALS	GPUS/TRIAL	CLUSTER SIZE
1	0-4	8	2	4
2	5-12	4	3	3

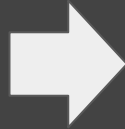
← END OF STAGE



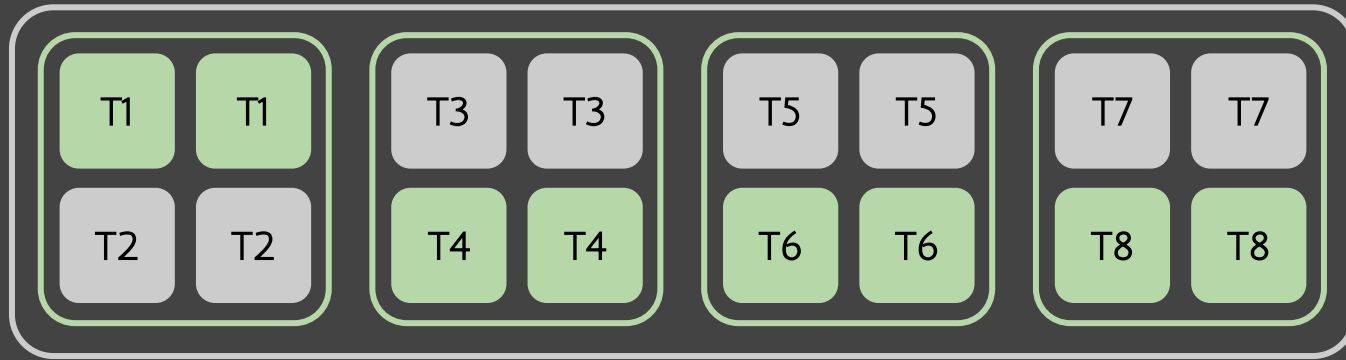
CLUSTER

ISSUE 3: EFFECTIVELY EXECUTE ALLOCATION PLAN

SCHEDULER



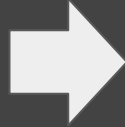
STOP - T2, T3, T5, T7
CONTINUE - T1, T4, T6, T8



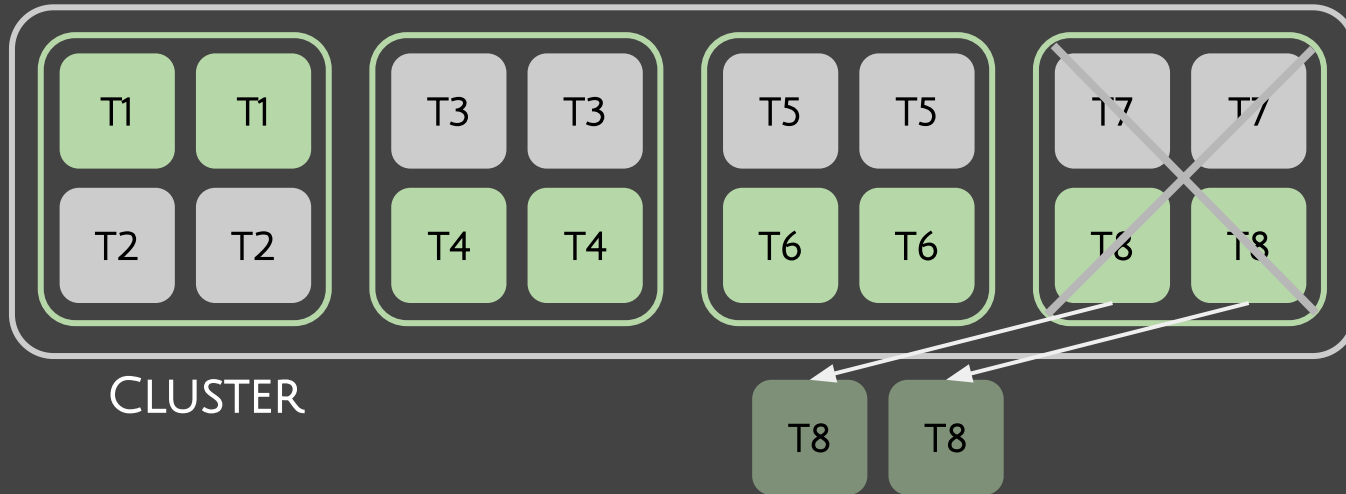
CLUSTER

ISSUE 3: EFFECTIVELY EXECUTE ALLOCATION PLAN

CLUSTER
MANAGER



DEPROVISION NODE 4

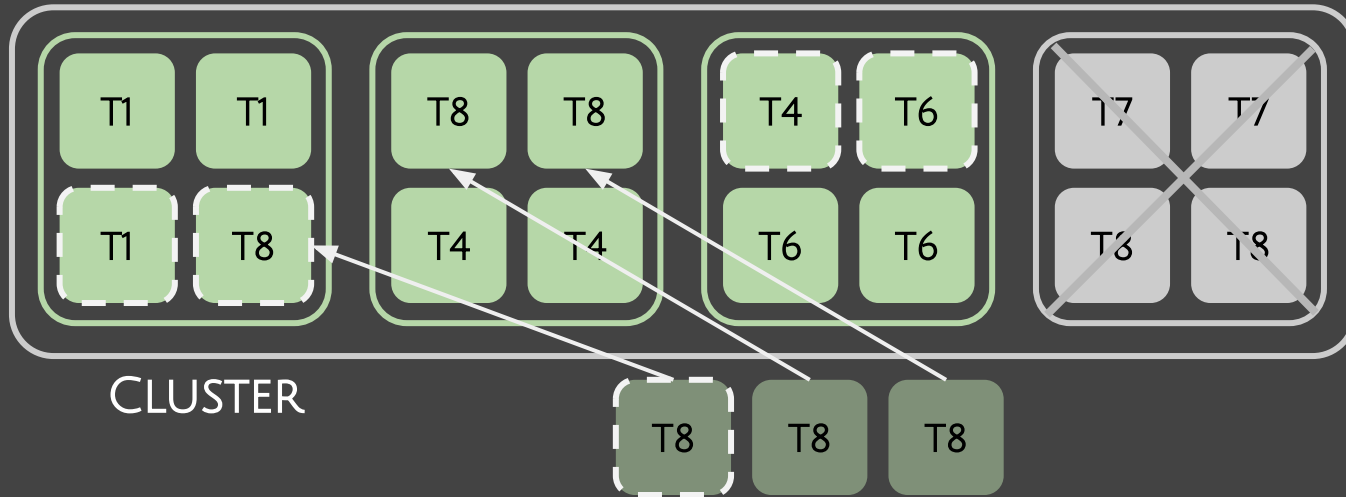


ISSUE 3: EFFECTIVELY EXECUTE ALLOCATION PLAN

PLACEMENT
CONTROLLER

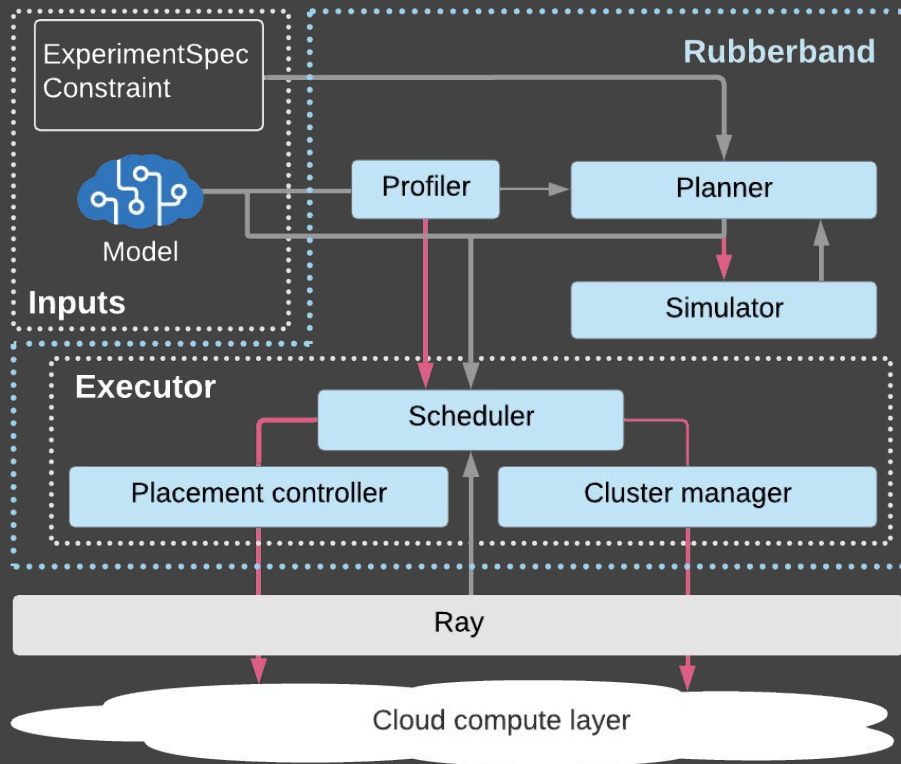


MOVE T8 TO NODE 1 AND 2



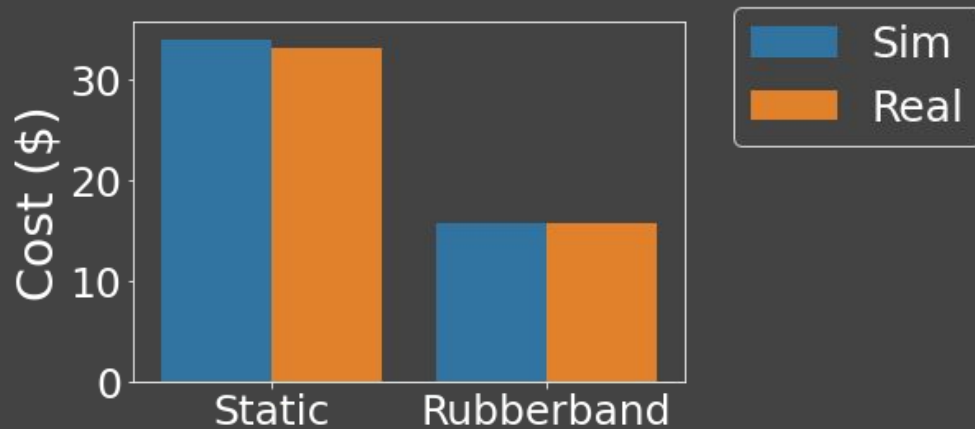
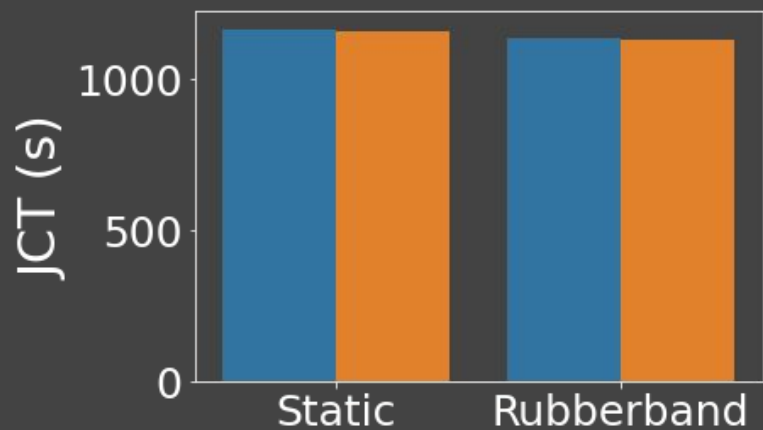
SYSTEM

1. **PROFILER** AND **SIMULATOR** MODEL JOB COMPLETION TIME + COST OF POTENTIAL ALLOCATIONS
2. **PLANNER** GENERATES A LOW COST ALLOCATION PLAN THAT COMPLETES ON TIME
3. **SCHEDULER**, **PLACEMENT CONTROLLER**, AND **CLUSTER MANAGER** EXECUTES THE ALLOCATION PLAN SUCH THAT WORKER CO-LOCATION AND CLUSTER UTILIZATION ARE MAXIMIZED

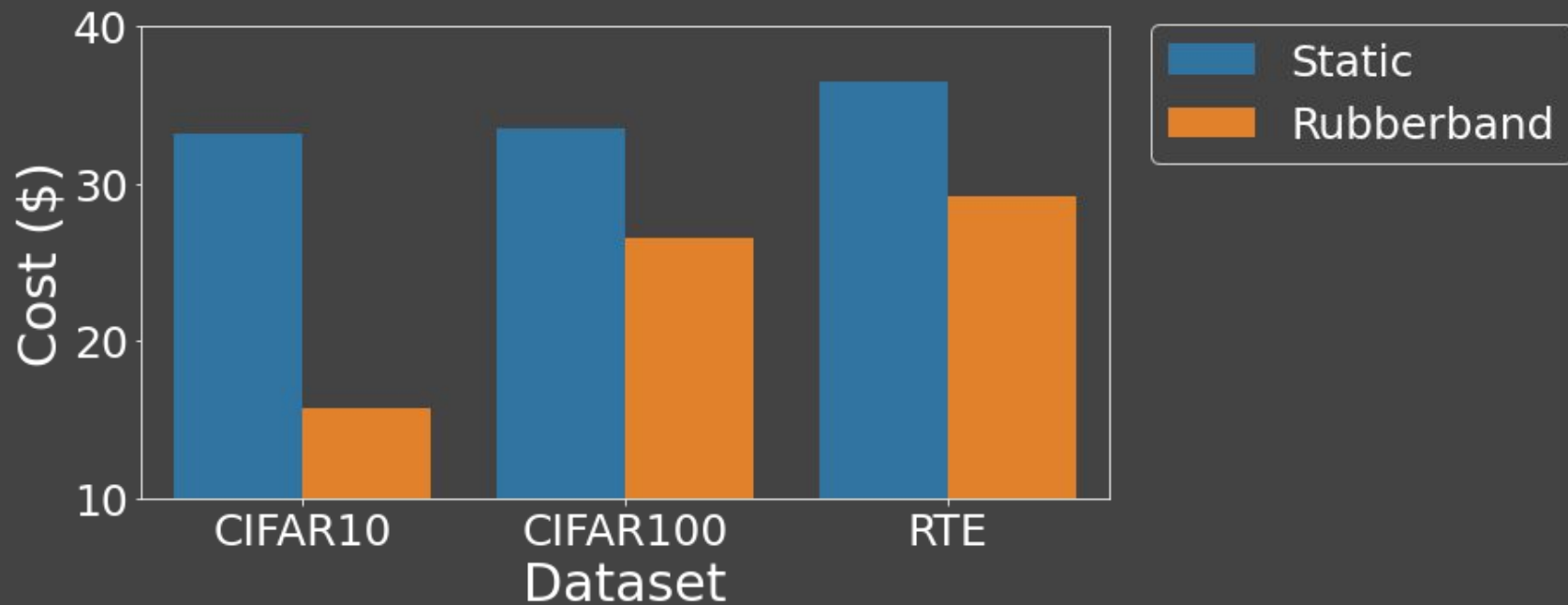


END-TO-END RESULTS

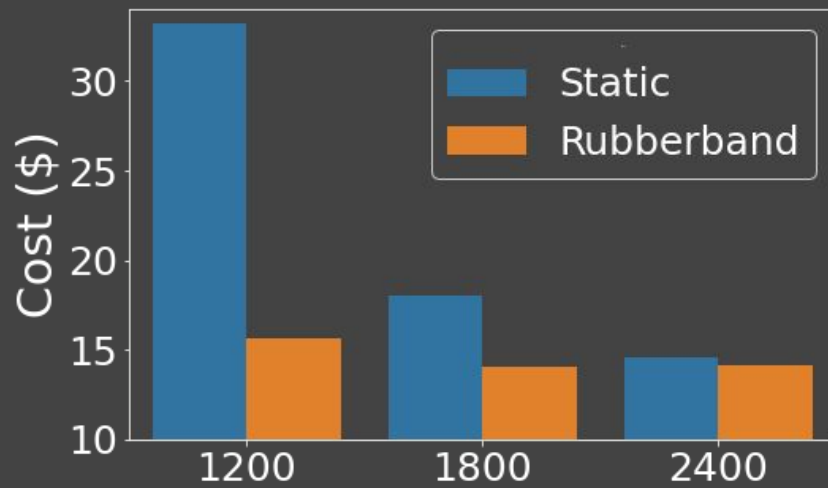
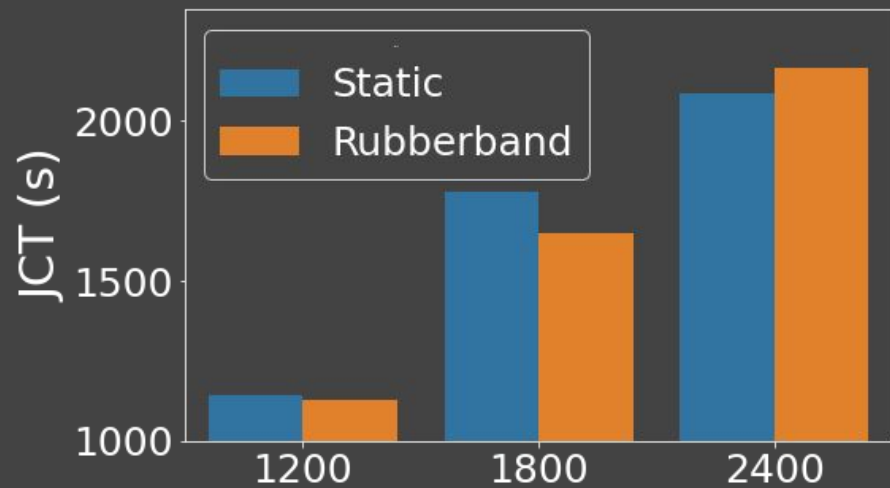
SIMULATION QUALITY



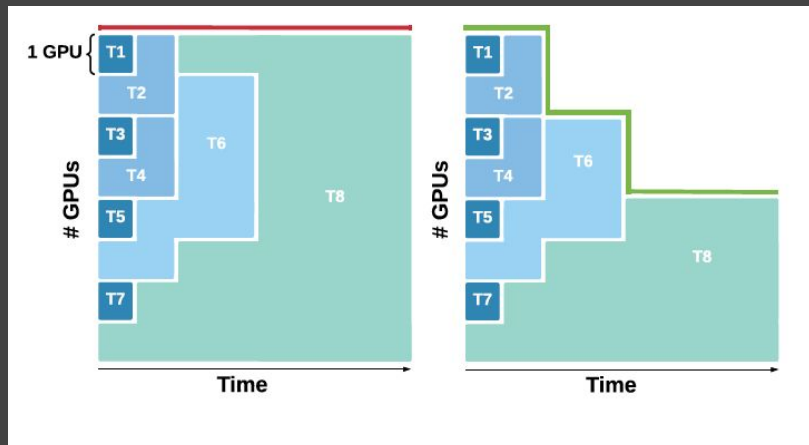
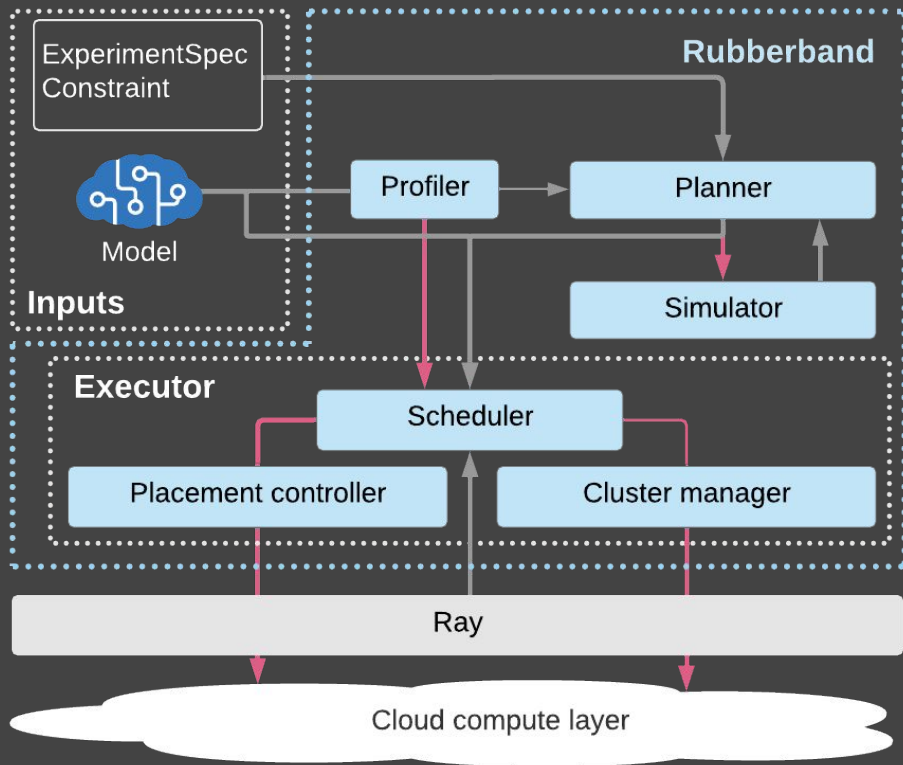
ACROSS DATASETS



ACROSS DEADLINES



RUBBERBAND



THANK YOU!