



SHEEO: Continuous Energy Efficiency Optimization in Autonomous Embedded Systems

Presenter: Xinkai Wang (Shanghai Jiao Tong University)

Advisor: Chao Li, Minyi Guo

ICCD 2024, Milan, Italy



Content



1. Background and Motivation

Sensing	Perception	Actuation
Software Gather &Sync	2D Perception Move 3D Perception ments Object Tracking Localization Positions	Movement ↑ Control ↑ Planning
Cams LiDA	CUDA کے - TensorRT ای Optimizer R Middleware Efficiency Resource Optimizer B Optimizer	DenCV Variance
Radar IMU	Hardware 🖉 CPU 🗔 GPU 📋 AS	nc 🕺 🖈 🗍

2. Design of SHEEO



3. Evaluation Results



4. Conclusion and Vision





Content



1. Background and Motivation

Sensing	Perception	Actuation
Software	2D Perception Move 3D Perception ments Object Tracking Localization Positions	Movement ▲ Control ◆ Planning
Cams LiDA	CUDA کَظُنُ TensorRT اللَّلَى Of Middleware اللَّلَي Efficiency اللَّقَانِ Resource Optimizer Optimizer	DenCV Runtime Variance
Radar IML	, Hardware 🖉 CPU 🗔 GPU 📋 AS	ic 大 🗍

2. Design of SHEEO



3. Evaluation Results



4. Conclusion and Vision





Autonomous Embedded System (AES)

AES is promising to minimize human intervention in various tasks.





Autonomous Embedded System (AES)

AES is promising to minimize human intervention in various tasks.





Background: Present AES Pipeline

















Background: Present AES Pipeline





Background: AES Efficiency Optimization (EO)

AES faces varied environments and complex hardware to manage.

T	emporal Variance: External Dynamics.							
-	AES 📇 / Urban Road / Traffic Jam / Overtake / Highway							
	Dynamics Case 1 Case 2 Case 3 Case 4							
	Vehicle Speed	Low	Low	Fluctuated	High			
	Surrounding	Low	Fluctuated	Low	Low			
	Task Complexity	Low	High	High	Low			
- 				•				
	Current Method	\	//F ↓	V/F	= 1			
	Expected Method	V/F↓	Adapting V/	F with time	V/F ↓			
	Requir	re Cor	ntinuou	is EO				





Require Intelligent EO

AES expects EO in 100ms granularity within a large optimization space.

[1] Kim, Young Geun, and Carole-Jean Wu. "Autoscale: Energy efficiency optimization for stochastic edge inference using reinforcement learning." MICRO 2020

Background: Review of AES Optimization

Researchers explore resource and energy optimization separately.



- [1] Yu, Bo, et al. "Building the computing system for autonomous micromobility vehicles: Design constraints and architectural optimizations." MICRO 2020
- [2] Krishnan, Srivatsan, et al. "Automatic domain-specific soc design for autonomous unmanned aerial vehicles." MICRO 2022
- [3] Qiu, Haoran, et al. "FIRM: An intelligent fine-grained resource management framework for SLO-Oriented microservices." OSDI 2020
- [4] Mishra, Nikita, et al. "Caloree: Learning control for predictable latency and low energy." ASPLOS 2018



Motivation: EO is Costly for AES



Desirable efficiency optimization interfere with AES pipeline.





Motivation: EO is Costly for AES



Desirable efficiency optimization interfere with AES pipeline.

Complex ML models introduces great training / inference overheads.

Continuous EO disturbs execution on resource-constrained AES.

10% slowdown due to efficiency optimization is unacceptable for safety-critical AES.

ing	Collecting Model		CALOREE, ASPLOS'18
ol Lear			Transfer Learning
Contre		Learned Control Model	Overhead-500ms / 2ms
			Overneau-Joomis / Zins

Frame 1=100ms	Frame 2=100ms		Frame 1=100ms	EO=10ms	Frame 2=100ms
		\Box		EO	

Typical EO using normal computing power causes ~10% slowdown of each iteration.



AES presents underutilized shadow cycles for deploying EO.



The ignored resources in AES are promising for EO without slowdown.





1. Background and Motivation

Sensing	Perception	Actuation
Gather &Sync	2D Perception Move Perception Fusion 3D Perception ments Object Tracking Localization Positions Prediction	Movement Control Planning
Cams LiDAF	Middleware CUDA CLUDA CL	nc Runtime

2. Design of SHEEO



3. Evaluation Results



4. Conclusion and Vision





SHEEO Design Consideration









SHEEO exploits shadow cycles for continuous and intelligent EO.



> SHEEO has two cooperative components: real-time monitor and continuous learning manager.

> SHEEO works at OS, enhancing **power management facility** without intrusion into AES pipeline.



SHEEO Observation: Environment and Workload

SHEEO observes internal and external status via real-time monitor.



- Runtime Monitor observes the start and end of AES tasks for recording shadow cycles.
- It observes the environment changes and system status and records them into the 7-dimensional state table.

State		Descriptions	Discrete Values
Workload	SCONV	Number of CONV layers	L (<30), M (<50), H (≥50)
Features	S _{FC}	Number of FC layers	L (<10), H (≥10)
reatures	S _{RC}	Number of RC layers	L (<10), H (≥10)
	S _{Com}	Computing units' utilization	L (<25%), M (<75%), H (≥75%)
Runtime	S _{Mem}	Memory utilization	L (<25%), M (<75%), H (≥75%)
Dynamics	S _{Speed}	Current vehicle speed	L (<10mph), M (<40mph), H (≥40mph)
	Svar	Variation from last execution	L (<10%), M (<30%), H (≥30%)

(a) Multi-dimensional state-related features

- State Table for Internal Runtime Status
 - Model software and hardware features.
 - Adjustable for actual AES workloads.
- State Table for External Environment
 - Model the external stochastic variance.

Use discrete ranges to reduce complexity.



SHEEO Observation: Environment and Workload

SHEEO observes internal and external status via real-time monitor.



- Runtime Monitor observes the start and end of AES tasks for recording shadow cycles.
- It observes the environment changes and system status and records them into the 7-dimensional state table.

State		Descriptions	Discrete Values
Workload	SCONV	Number of CONV layers	L (<30), M (<50), H (≥50)
Features	S _{FC}	Number of FC layers	L (<10), H (≥10)
reatures	S _{RC}	Number of RC layers	L (<10), H (≥10)
	S _{Com}	Computing units' utilization	L (<25%), M (<75%), H (≥75%)
Runtime	SMem	Memory utilization	L (<25%), M (<75%), H (>75%)
Dynamics	S _{Speed}	Current vehicle speed	L (<10mph), M (<40mph), H (≥40mph)
	S _{Var}	Variation from last execution	L (<10%), M (<30%), H (≥30%)

(a) Multi-dimensional state-related features

- State Table for Internal Runtime Status
 - Model software and hardware features.
 - > Adjustable for actual AES workloads.
- State Table for External Environment
 - Model the external stochastic variance.
 - Use discrete ranges to reduce complexity.



SHEEO Optimization: Learning and Control

SHEEO optimizes efficiency via continuous learning manager.

Environment	1 Obse	erve	State	Action	Reward
	Real-time		<i>S</i> ¹	A ¹	R_e^1, R_l^1
Workloads			S	tore Tab	ole
			\bullet S^t	A ^t	R_e^t, R_l^t
	4 Feedb	pack	VSC Pe ② Update	eriods Q-Table	States, Rewards
Hardware	HSC Periods ③ Select	Action State	n A ₁		A _m
CPU (Internete)	Action	<i>S</i> ₁	$Q(S_1, A$	1)	$Q(S_1, A_m)$
GPU memory	Power Management			Q-Tabl	e
	Decisions	S _n	$Q(S_n, A$	1)	$Q(S_n, A_m)$

Input: Pre-trained Q-table Q(S, A), learning rate α , discount factor γ , exploration probability ϵ Output: Fine-tuned O-table and Action A for each iteration while \exists VSC periods and \exists Store Table do Calculate recorded reward R_{energy} and $R_{latency}$; For consecutive S and S' and chosen action A; Choose action A' with the largest Q(S', A'); $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)];$ $S \leftarrow S'$ 8 foreach HSC periods do Gather S and locate in Q(S, A); if rand() $< \epsilon$ then 10 Choose action A randomly; 11 else 12 13 Choose action A with the largest Q(S, A); end 14 Store Table \leftarrow (A, S') ; 15 16 end

Algorithm 1: Algorithm for continuous learning

- Utilizing preemptive and variable-sized VSC. Two-fold Reward
- Minimize performance slack
- Maximize energy efficiency
- Used for updating Q-table

```
\begin{split} R &= aR_{latency} + bR_{energy} \text{ s.t. } R_l \leq \text{Constraint} \\ R_{latency} &= \max_{t \in \text{Perception}} T^t_{end} - T_{start} \\ R_{energy} &= \sum_i E^i_{Unit} = \sum_i \int_{T_b}^{T_e} P_f \delta t + P_{idle} \times t_{idle} \end{split}
```

- Continuous Learning Manager adopts reinforcement
 learning (Q-learning) to optimize energy efficiency.
- It embeds continuous learning into VSC and action control into HSC to reduce power management facility costs.
- > SHEEO learns optimal action of given states.
 - Consider status of consecutive iterations.
 - Adjustable for Q-learning, DDPG, DRL, etc.



SHEEO Optimization: Learning and Control

SHEEO optimizes efficiency via continuous learning manager.

Environment	1 Obse	erve	State	Action	Reward
	Real-time	Monitor	<i>S</i> ¹	A ¹	R_e^1, R_l^1
Workloads			S	tore Tab	le
			S ^t	A ^t	R_e^t, R_l^t
	(4) Feedb	ack (2	VSC Pe	riods Q-Table	States, Rewards
Hardware	HSC Periods ③Select	Action State	A ₁		A _m
CPU (Internet)	Action	<i>S</i> ₁	$Q(S_1, A_1)$)	$Q(S_1, A_m)$
GPU Integrate memory	Power Management			Q-Tabl	e
	Decisions	S _n	$Q(S_n, A_1)$)	$Q(S_n, A_m)$

- Continuous Learning Manager adopts reinforcement
 learning (Q-learning) to optimize energy efficiency.
- It embeds continuous learning into VSC and action control into HSC to reduce power management facility costs.

Al	gorithm 1: Algorithm for continuous learning		
I	nput: Pre-trained Q-table $Q(S, A)$, learning rate α ,		
	discount factor γ , exploration probability ϵ		
C	Output: Fine-tuned Q-table and Action A for each		
	iteration		
1 W	Thile \exists VSC periods and \exists Store Table do		
2	Calculate recorded reward R_{energy} and $R_{latency}$;		
3	For consecutive S and S' and chosen action A ;		
4	Choose action A' with the largest $Q(S', A')$;	\triangleright	
5	$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)];$		
6	$S \leftarrow S'$		
7 e	nd		(
8 fc	preach HSC periods do		
9	Gather S and locate in $Q(S, A)$;		
10	if $rand() < \epsilon$ then		
11	Choose action A randomly;		
12	else		
13	Choose action A with the largest $Q(S, A)$;		
14	end		
15	Store Table \leftarrow (A, S'):		

Utilizing short-lived and periodic HSC.

7-dimensional Action

- Divided into discrete levels.
- Consider heterogeneous HW.

Actions	Descriptions
SX	Turn On/Off Component X
NCPU	Number of Active Core
FCPU	V/F Level of CPU Core
F _{GPU}	V/F Level of GPU
F _{DLA1}	V/F Level of DLA1
F _{DLA2}	V/F Level of DLA2
FMEM	V/F Level of Memory

> SHEEO controls the action in next iteration.

- > Use ϵ to control exploitation & exploration.
- > Agilely make decisions with largest Q(S, A).





1. Background and Motivation

Sensing	Actuation	
Gather &Sync	2D Perception Move Perception Fusion 3D Perception ments Object Tracking Localization Positions Prediction	Movement Control
Cams LiDAF	Middleware CUDA CUDA CUDA CUDA CUDA CUDA CUDA CUDA	Runtime Variance

2. Design of SHEEO



3. Evaluation Results



4. Conclusion and Vision





Evaluation Methodology and Settings

We evaluate SHEEO extensively on commercial edge platforms.

Evaluated Platform Specifications			
Device	Nvidia Jetson AGX Orin Module		
CPU	8-core ARM Cortex-A78 v8.2		
GPU	1792-core Ampere GPU with 56 tensor cores		
Memory	32GB 256-bit LPDDR5		
Accelerator	2x NVDLA v2, 1x PVA v2		
System	Linux 5.19.104-tegra with Jetpack 5.1.1		
Software	CUDA 11.4 and TensorRT 8.5.2		

Evaluated Workloads					
SSD	ssd-mobilenet-v1 for object detection.				
YOLO	yolov3-tiny-416 for image detection.				
SRCNN	super-resolution-bsd500 for image reconstruction.				
Evaluated Baselines					
Workload-Aware Control [1]		Use static model to profile workloads.			
Learning-Based Control [2]		Use ML model and historical statistics.			
Oracle		Offline optimal energy efficiency scheme			

> We evaluate SHEEO on Nvidia Jetson AGX Orin and consider shadow cycles on GPU and NVDLA.

> We set learning rate as 0.9, discount factor as 0.1, and exploration probability as 0.1.

We mimic the perception stage of AES pipeline with three networks with different combinations.

[1] Bateni, Soroush, and Cong Liu. "NeuOS: A Latency-Predictable Multi-Dimensional Optimization Framework for DNN-driven Autonomous Systems." USENIX ATC 20 [2] Mishra, Nikita, et al. "Caloree: Learning control for predictable latency and low energy." ASPLOS 2018



Harvesting Heterogeneous Resources

SHEEO harvests shadow cycles in AES under various scenarios.



- Q-learning is too lightweight to occupy all the shadow cycles, but more advanced facilities could utilize better.
- The harvesting ratio is higher under higher vehicle speed and higher runtime variance scenarios.



Reducing Energy Consumptions



SHEEO improves energy efficiency of AES under various scenarios.

SHEEO provides a more agile and

accurate power management solution.



SHEEO achieves better energy efficiency under higher runtime variance and medium vehicle speed.

SHEEO achieves better energy-delay

product (EDP) compared with baselines.



SHEEO improves energy efficiency with a little sacrifice of execution performance.





1. Background and Motivation

Sensing	Perception	Actuation
Software	2D Perception Move Perception Fusion 3D Perception ments Object Tracking Localization Positions Prediction	Movement Control
Cams LiDAR	Middleware CUDA CLUCA CL	penCV Variance

2. Design of SHEEO



3. Evaluation Results



4. Conclusion and Vision





Vision: Future of Autonomous Systems

Towards better accelerator usage for efficient system operations.



Moving Towards Efficient Autonomous Systems!

[1] Xinkai, Wang, et al. "Not All Resources are Visible: Exploiting Fragmented Shadow Resources in Shared-State Scheduler Architecture." SoCC 2023

- [2] Lingyu, Sun, et al. "A2: Towards Accelerator Level Parallelism for Autonomous Micromobility Systems." TACO 2024
- [3] Lingyu, Sun, et al. "Jigsaw: Taming BEV-centric Perception on Dual-SoC for Autonomous Driving." RTSS 2024



Conclusion



Initialization Localization Actuation Decisions We analyze shadow cycles in autonomous embedded systems Sensor Buffer Intermediate Buffer Senso 4Frame1 HSC Detection Recognition VSC HSC and the need to **continuously optimizing** energy efficiency. HSC HSC Depth Estimation Tracking Sensing Perception Stage Actuation

- SHEEO is an intelligent continuous energy efficiency optimizer that first embeds the costly facility into shadow cycles.
 - SHEEO harvests the underutilized hetero-computing resources and improves AES energy efficiency under various scenarios.

Environment	1 Obse	rve	State	Action	Reward
	Real-time	Monitor	S1	A^1	R_e^1, R_l^1
	🤌 👰 🔹		Store Table		
VVorkioads			St	A^t	R_e^t, R_l^t
	④Feedb	ack (2	VSC Pe	eriods Q-Table	States, Rewards
Hardware	3Select	Action State	A1		Am
	Action	<i>S</i> ₁	$Q(S_1, A)$	1)	$Q(S_1, A_m)$
GPU Integrate memory	Power Management		Q-Table		le
	Decisions	S _n	$Q(S_n, A$	1)	$Q(S_n, A_m)$



By exploiting "free" shadow cycles within more heterogeneous accelerators, AES can enable a wide range of intelligent management facilities and autonomous workloads.



Thank You! Q & A

SHEEO: Continuous Energy Efficiency Optimization in Autonomous Embedded Systems

Discussion: Xinkai Wang (Presenter), unbreakablewxk@sjtu.edu.cn



