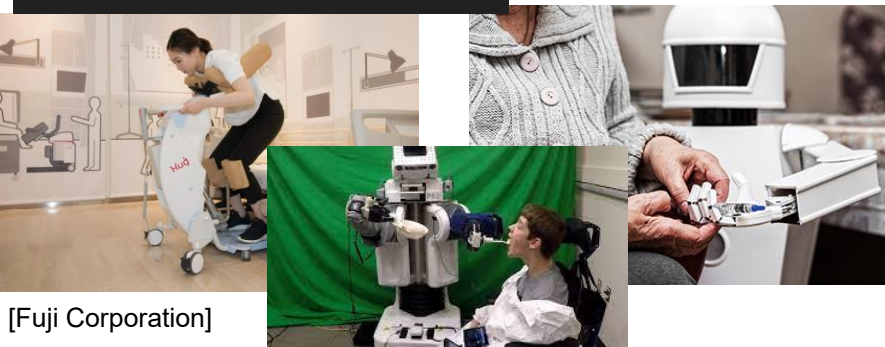


Minimal Mobile Systems via Cloud-based Adaptive Task Processing

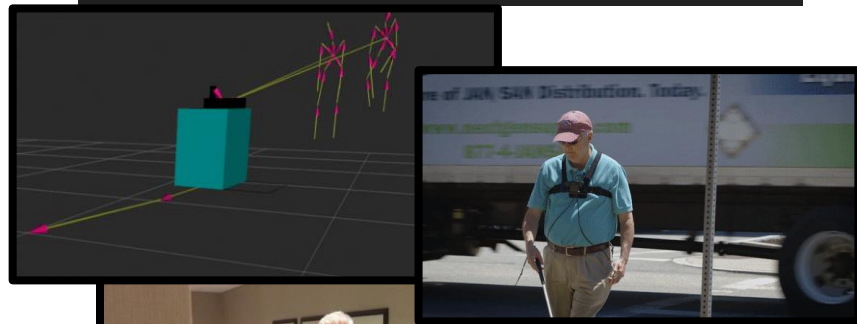
Assistive Care



[Fuji Corporation]

[Park et al., 2020]

Personalized Mobility



Delivery



[Zipline, Dispatch]

Rehabilitation



[ReWalk]



Renato Mancuso

Sanja Arora

Hee Jae Kim

Lei Lai

Bassel Mabsout

Eshed Ohn-Bar

Boston University

Minimal Mobile Systems via Cloud-based Adaptive Task Processing



Minimal Mobile Systems via Cloud-based Adaptive Task Processing

Video Model Training

OCCUPANCY NETWORK RECIPE

Pick 1.44B frames
Train for 100,000 GPU-hours at 90°C

14K GPUs

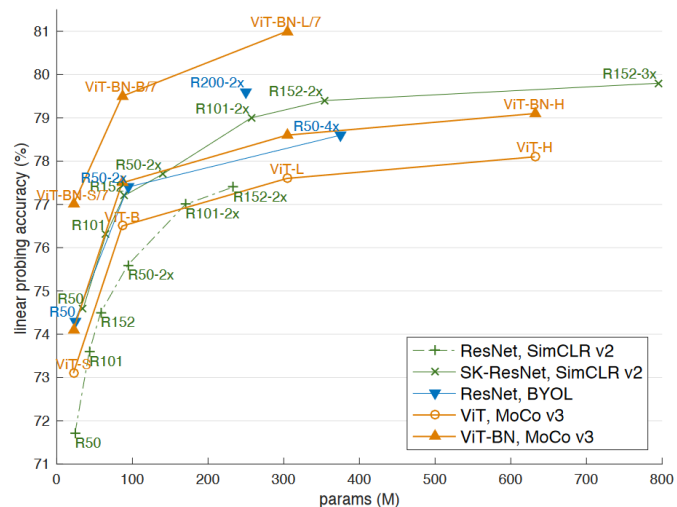
4K for auto labeling
10K for training

30PB DISTRIBUTED VIDEO CACHE

160B frames
500K videos rotating through cache/day
400K video instantiations per second: `Tclip(clip_id)`



14K GP
4K for a
10K for
30PB T



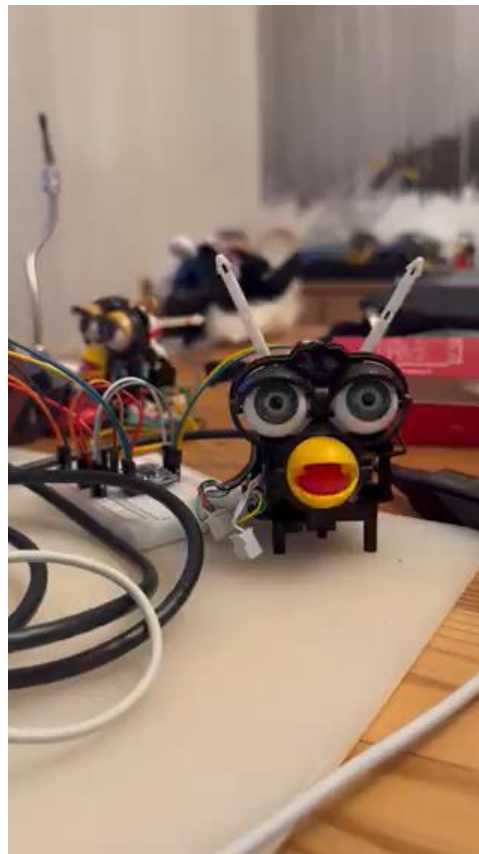
[Tesla AI Day, 2022]

[Chen et al, 2022]

Cloud-based Inference?



[BeMyEyes, OpenAI]



[@jessicacard]

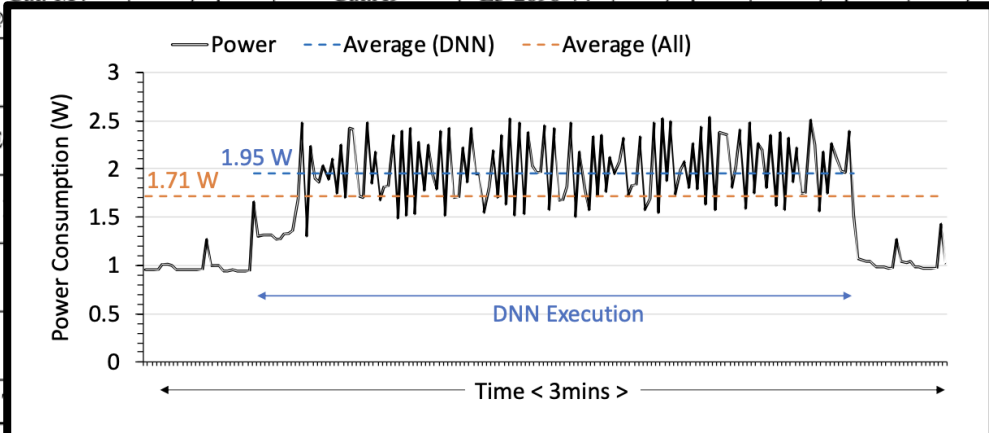
Edge Device DNN Inference is Costly

Category	IoT/Edge Devices	GPU-Based Edge Devices		Custom-ASIC Edge Accelerators	FPGA Based	CPU	HPC Platforms GPU			
Platform	Raspberry Pi 3B [34]*	Jetson TX2 [69]	Jetson Nano [36]	EdgeTPU [35]	Movidius NCS [37]†	PYNQ-Z1 [64]	Xeon	RTX 2080	GTX Titan X	Titan Xp
CPU	4-core Ctx.A53 @1.2 GHz*	4-core Ctx.A57 2-core Denver2 @2 GHz	4-core Ctx.A57 @1.43 GHz	4-core Ctx.A53 & Ctx.-M4 @1.5 GHz	N/Ap	4-core Ctx.A9 @650 MHz	2x 22-core E5-2696 v4 @2.20GHz	N/Ap*	N/Ap	N/Ap
GPU	No GPGPU	256-core Pascal μ A	128-core Maxwell μ A	N/Ap	N/Ap	N/Ap	N/Ap	2944-core Turing μ A	3072-core Maxwell μ A	3840-core Pascal μ A
Accelerator	N/Ap	N/Ap	N/Ap	EdgeTPU	Myriad 2 VPU	ZYNQ XC7Z020	N/Ap	N/Ap	N/Ap	N/Ap
Memory‡	1 GB LPDDR2	8 GB LPDDR4	4 GB LPDDR4	N/Av*	N/Av	630 KB BRAM 512 MB DDR3	264 GB DDR4	8 GB GDDR6	12 GB GDDR5	12 GB GDDR5X
Idle Power‡	1.33	1.90	1.25	3.24	0.36	2.65	\approx 70	\approx 39	\approx 15	\approx 55
Average Power‡	2.73	9.65	4.58	4.14	1.52	5.24	300 TDP	\approx	\approx 100	\approx
Platform	All	All	All	TFLite	NCSDK	TVM/FINN	All	All	All	All

† Effective memory size used for acceleration/execution of DNNs, e.g., GPU/CPU/Accelerator memory size. * Ctx.: Arm Cortex. N/Ap: Not applicable. N/Av: Not available.
‡ : Measured idle and average power while executing DNNs, in Watts. *: Raspberry Pi 4B [70], with 4-core Ctx.A72 and maximum of 4 GB LPDDR4, was released after this paper acceptance. With better memory technology and out-of-order execution, Raspberry Pi 4B is expected to perform better. † Intel Neural Compute Stick 2 [61] with a new VPU chip and support for several frameworks was announced during paper submission, but the product was not released.

Edge Device DNN Inference is Costly

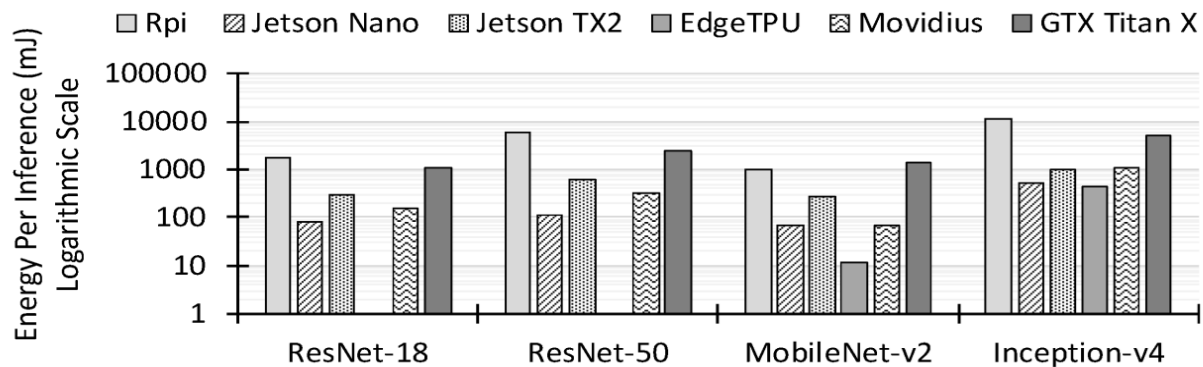
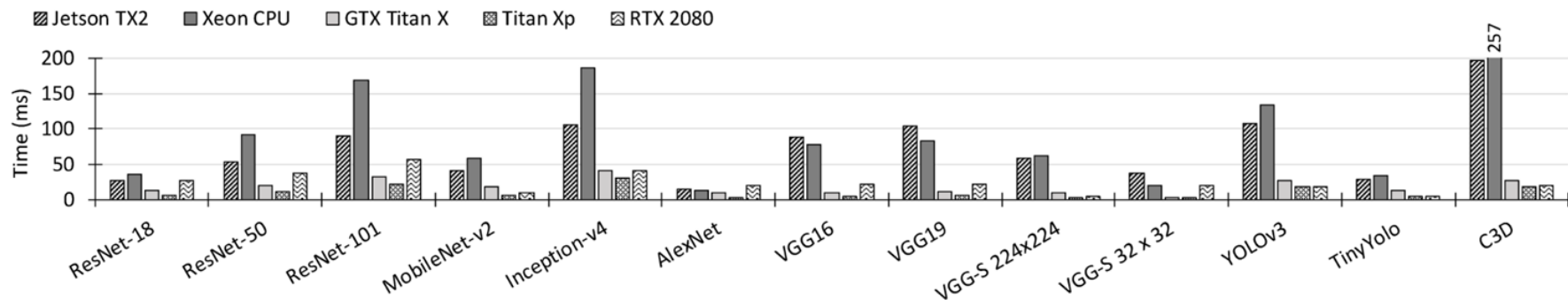
Category	IoT/Edge Devices	GPU-Based Edge Devices	Custom-ASIC Edge Accelerators	FPGA Based	CPU	HPC Platforms GPU				
Platform	Raspberry Pi 3B [34]*	Jetson TX2 [69]	Jetson Nano [36]	EdgeTPU [35]	Movidius NCS [37]†	PYNQ-Z1 [64]	Xeon	RTX 2080	GTX Titan X	Titan Xp
CPU	4-core Ctx.A53 @1.2 GHz*	4-core Ctx.A57 2-core Denver2 @2 GHz	4-core Ctx.A57 @1.43 GHz	4-core Ctx.A53 & Ctx.-M4 @	N/Ap	4-core Ctx.A9	2x 22-core E5-2696 v4	N/Ap*	N/Ap	N/Ap
GPU	No GPGPU	256-core Pascal μ A	128-core Maxwell μ A							
Accelerator	N/Ap	N/Ap	N/Ap	E						
Memory†	1 GB LPDDR2	8 GB LPDDR4	4 GB LPDDR4							
Idle Power‡	1.33	1.90	1.25							
Average Power‡	2.73	9.65	4.58							
Platform	All	All	All							



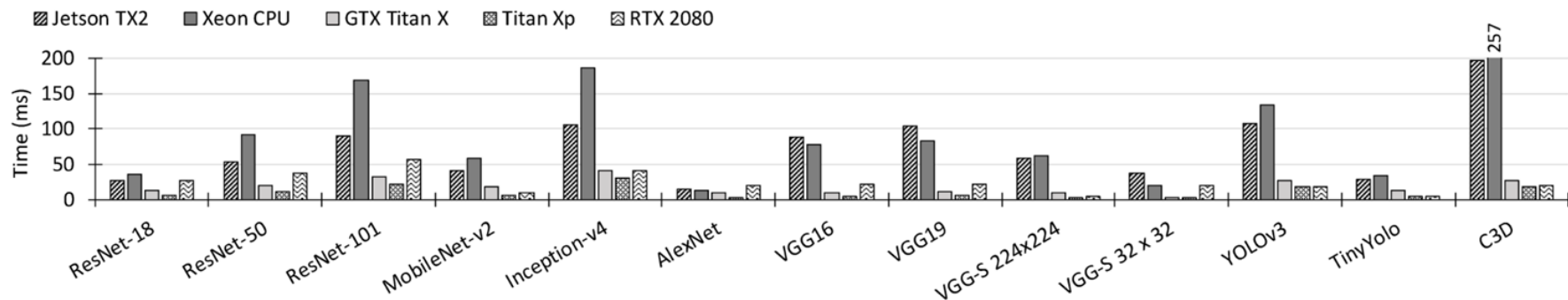
† Effective memory size used for acceleration/execution of DNNs.
 ‡: Measured idle and average power while executing DNNs, for acceptance. With better memory technology and out-of-order execution and support for several frameworks was announced during p

iRobot: 25-50+% Power increase in DNN Inference

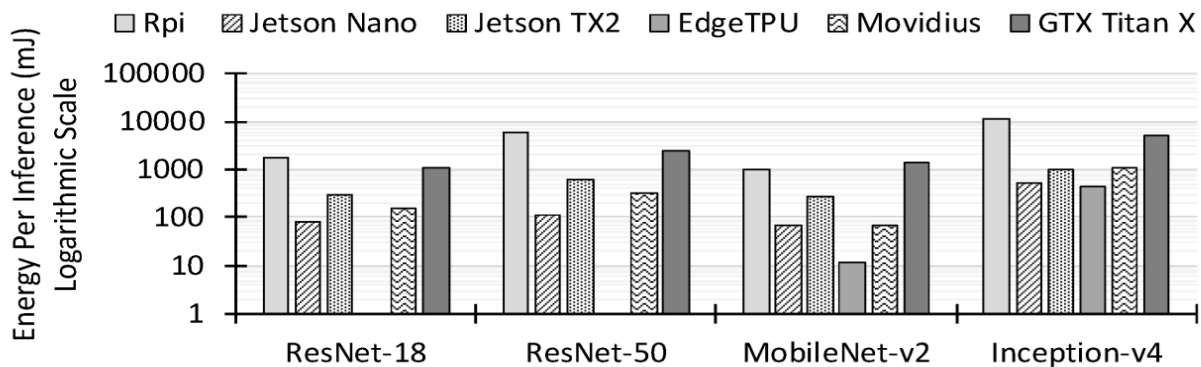
Edge Device DNN Inference is Costly



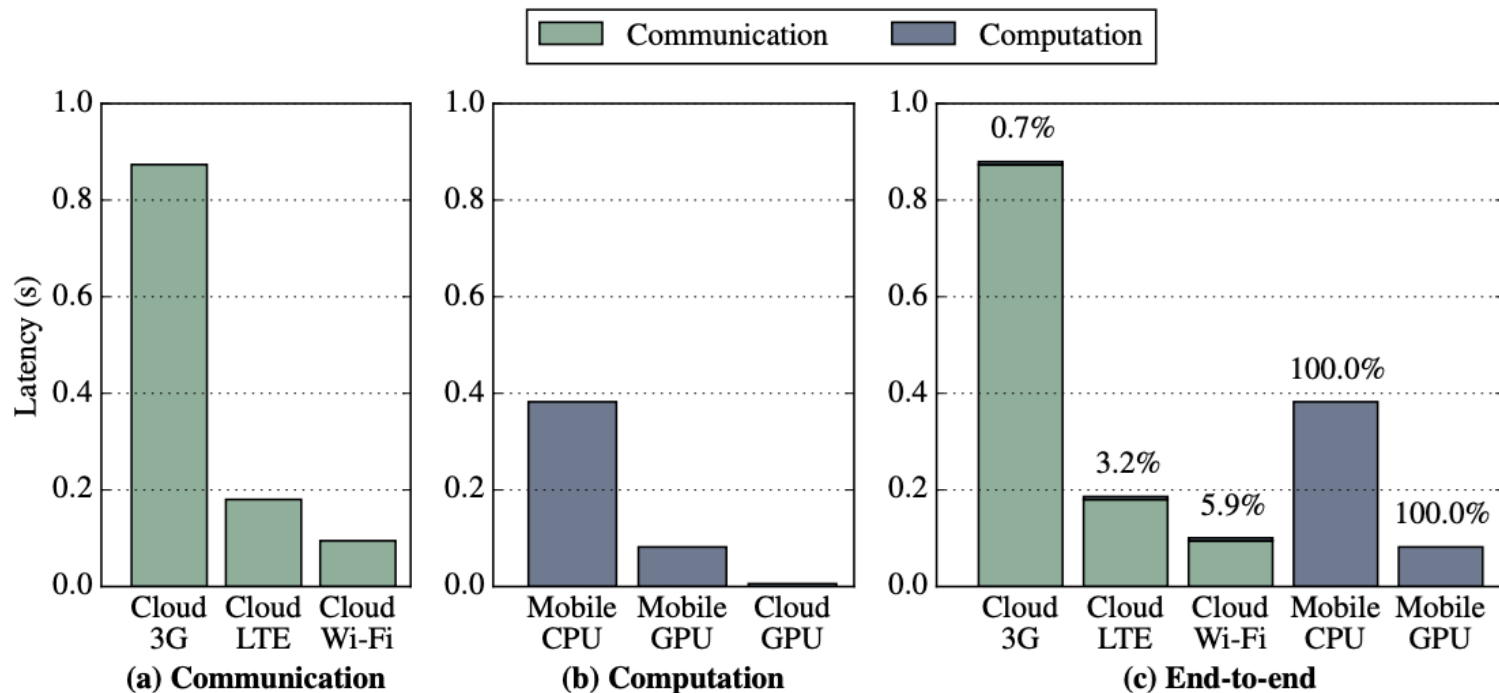
Edge Device DNN Inference is Costly



	mIOU	Params	MAdds
MNet V2*	75.32	2.11M	2.75B
	77.33	2.11M	152.6B
ResNet-101	80.49	58.16M	81.0B
	82.70	58.16M	4870.6B



Communication Latency is Decreasing Rapidly

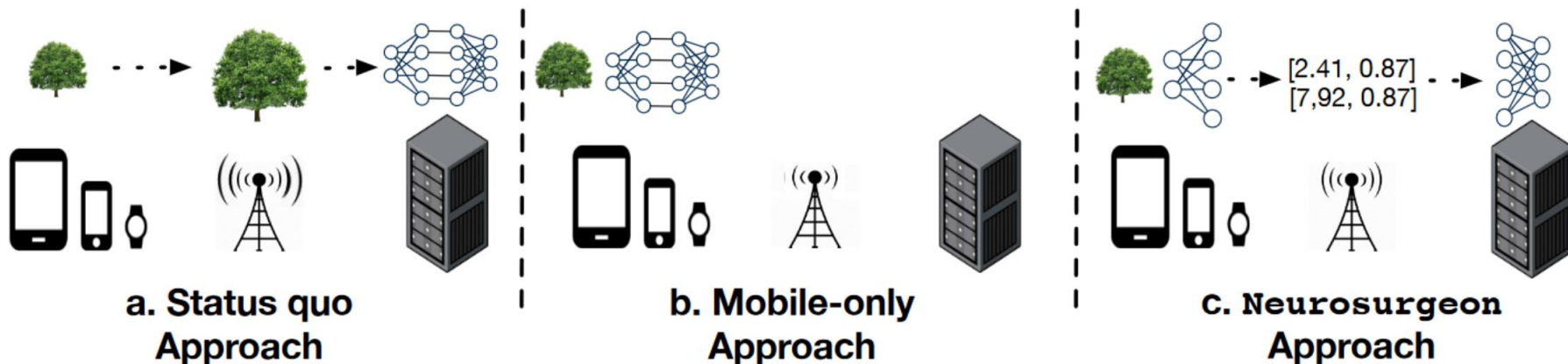


Transmitting
Images

Execution

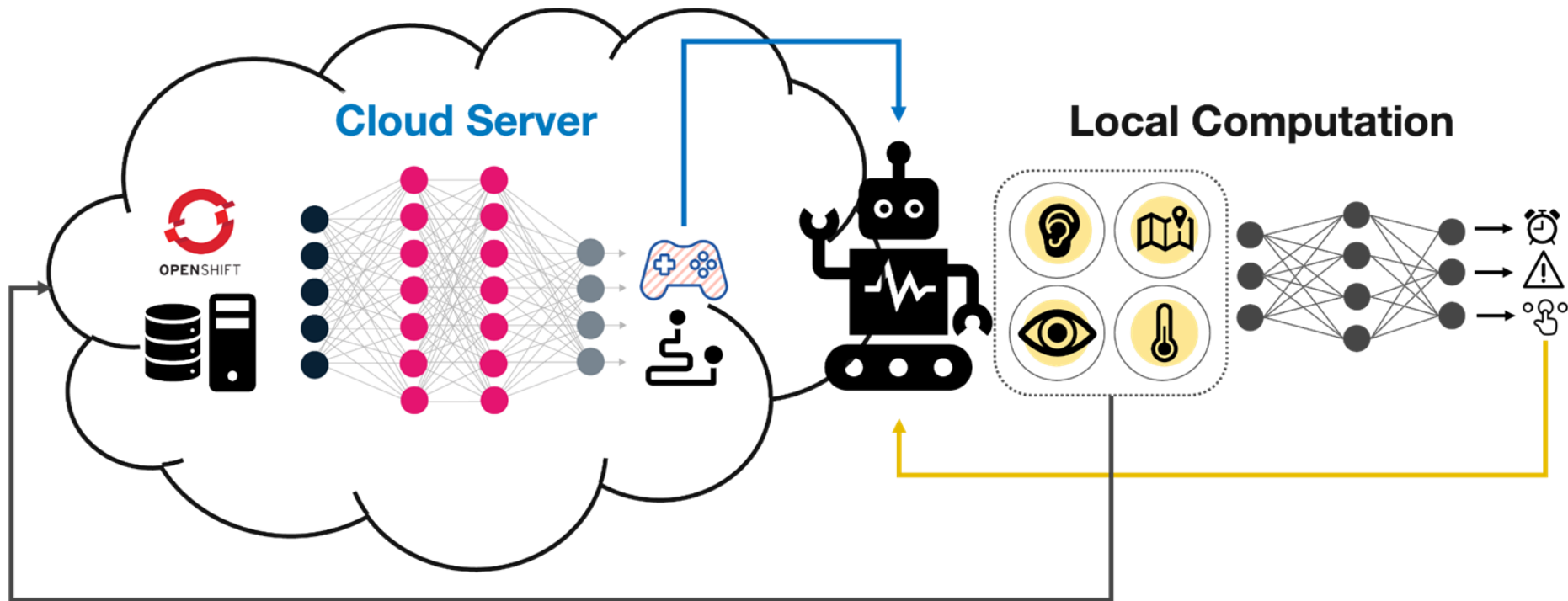
Total

Prior Work - Example



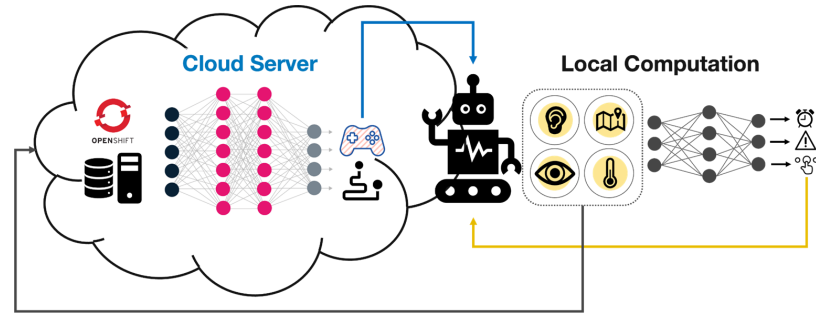
- May still have unacceptable latency
- Not ***task nor context-dependent***
- Cloud can run bigger and better models

How to Combine the Best of Both Worlds?

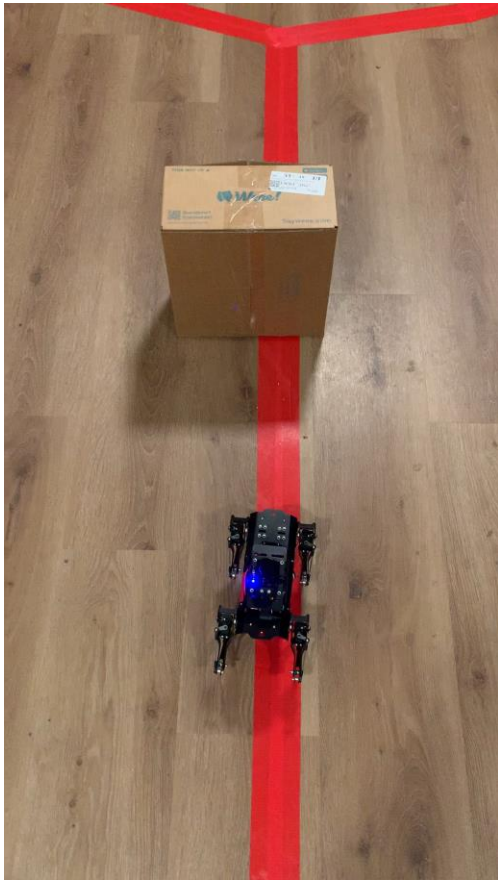


Task-Aware Markov Decision Process Formulation

- **States:** Data coming from the robot (e.g., *sensor data*, battery state, safety), potentially cloud server state (e.g., queues, prioritization, running tasks).
- **Actions:** The agent must select *the task* to perform and whether *locally or on the cloud*.
- **Task Reward:** Collision avoidance, path adherence, social disruption.

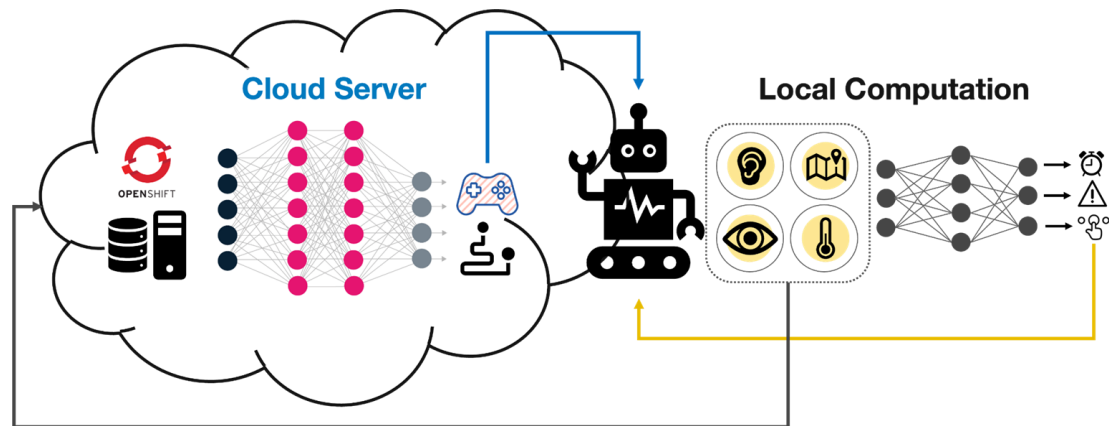


Preliminary Results



Open Challenges

- Standardized benchmarking?
- Latency model?
- Cloud processes?
- Local adaptation?



Standardized Benchmark with Streaming Perception



Accessibility, Vision, and Autonomy Challenge @ CVPR 2023



Win **\$500** to detect pedestrians and mobility aids:
[accessibility-cv.github.io](https://github.com/accessibility-cv)

Toyota suspends all self-driving vehicles at Paralympic Games after collision with athlete

Scooter 561 - Aug 27th 2021 10:58 am PT



Pitt suspends delivery robots after wheelchair user reports safety hazard

JAMIE MARTINES | Thursday, Oct. 24, 2019 7:27 p.m. SUPPORT LOCAL JOURNALISM

