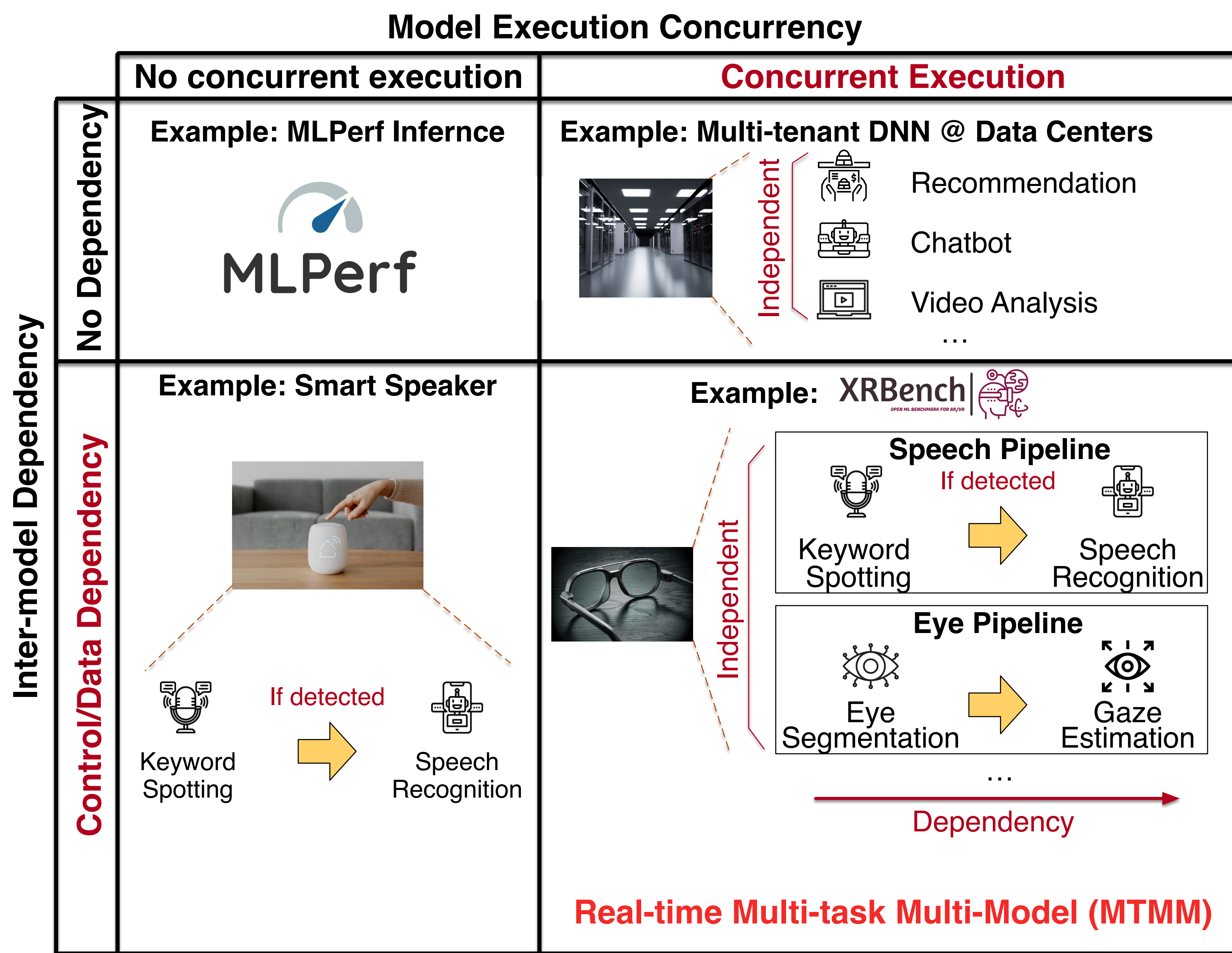




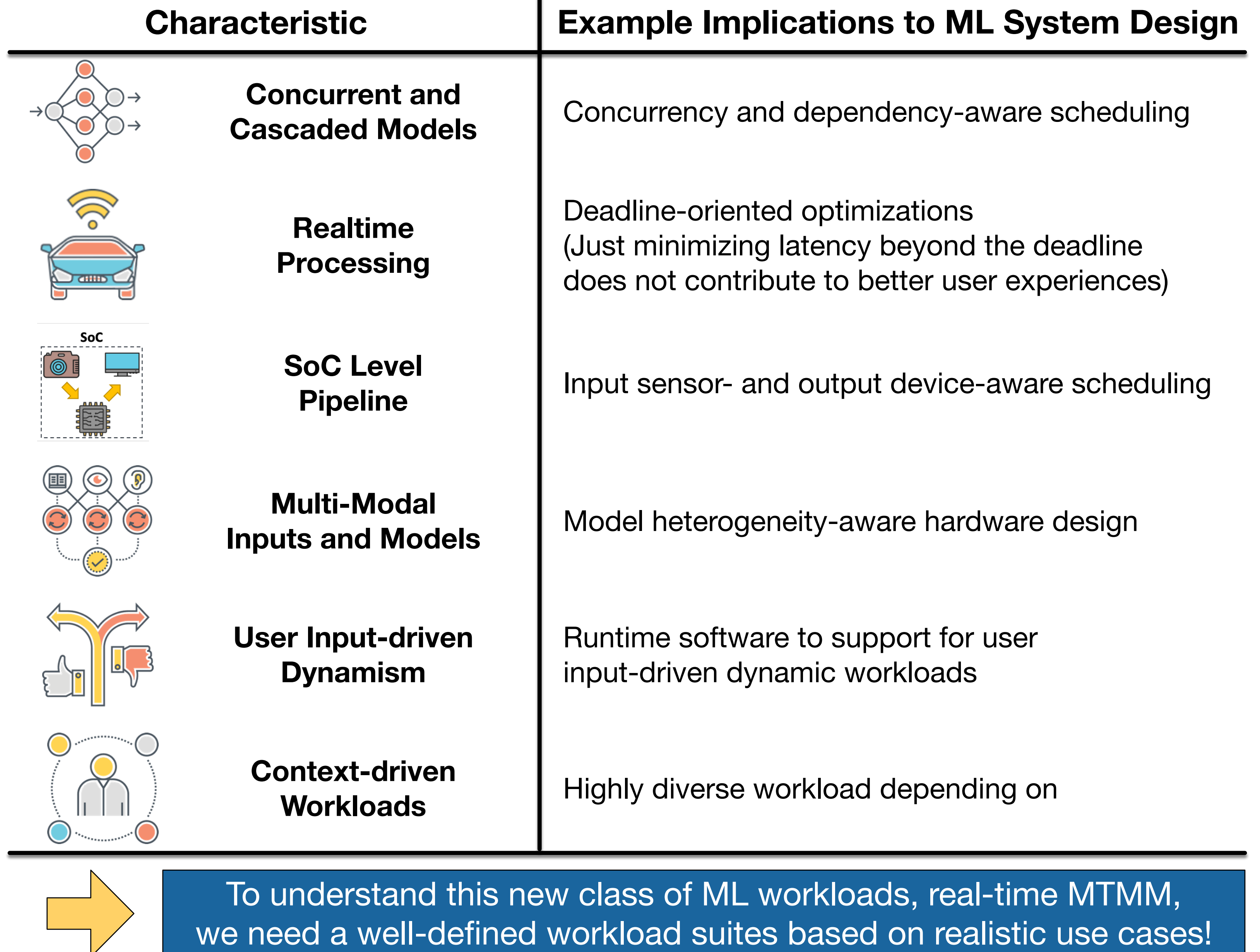
Hyoukjun kwon,<sup>1,2</sup> Krishnakumar Nair,<sup>2</sup> Jamin Seo,<sup>3</sup> Jason Yik,<sup>4</sup> Debabrata Mohapatra,<sup>2</sup> Dongyuan Zhan,<sup>2</sup> Jinook Song,<sup>2</sup> Peter Capak,<sup>2</sup> Peizhao Zhang,<sup>2</sup> Peter Vajda,<sup>2</sup> Colby Banbury,<sup>4</sup> Mark Mazumder,<sup>4</sup> Liangzhen Lai,<sup>2</sup> Ashish Sirasao,<sup>2</sup> Tushar Krishna,<sup>3</sup> Harshit Khaitan,<sup>2</sup> Vikas Chandra,<sup>2</sup> Vijay Janapa Reddi<sup>4</sup>



## ML Workload Taxonomy



## Characteristics of Real-time MTMM ML Workloads



## XRBench v0.1

### Unit Models

Category	Task	Model	Dataset	Model Perf. Requirement
Interaction	Hand Tracking (HT)	Hand Shape/Pose (Ge et al., 2019)	Stereo Hand Pose (Zhang et al., 2017)	AUC PCK, GT 0.948
	Eye Segmentation (ES)	RITNet (Chaudhary et al., 2019)	OpenEDS 2019 (Garbin et al., 2019)	mIoU, GT 90.54
	Gaze Estimation (GE)	Eyecod (You et al., 2022)	OpenEDS 2020 (Palmero et al., 2021)	Angular Error, LT 3.39
	Keyword Detection (KD)	Key-Res-15 (Tang & Lin, 2018)	Google Speech Cmd (Google, 2017)	Accuracy, GT 85.60
	Speech Recognition (SR)	Emformer (Shi et al., 2021)	LibriSpeech (Panayotov et al., 2015)	WER (others), LT 8.79
Context Understanding	Semantic Segmentation (SS)	HRViT (Gu et al., 2022)	Cityscape (Cordts et al., 2016)	mIoU, GT 77.54
	Object Detection (OD)	D2Go (Meta, 2022b)	COCO (Lin et al., 2014)	boxAP, GT 21.84
	Action Segmentation (AS)	TCN (Lea et al., 2017)	GTEA (Fathi et al., 2011)	Accuracy, GT 60.8
	Keyword Detection (KD)	Key-Res-15 (Tang & Lin, 2018)	Google Speech Cmd (Google, 2017)	Accuracy, GT 85.60
	Speech Recognition (SR)	Emformer (Shi et al., 2021)	LibriSpeech (Panayotov et al., 2015)	WER (others), LT 8.79
World Locking	Depth Estimation (DE)	MidAS (Ranftl et al., 2020)	KITTI (Geiger et al., 2012)	$\delta > 1.25$ , LT 22.9
	Depth Refinement (DR)	Sparse-to-Dense (Ma & Karaman, 2018)	KITTI (Geiger et al., 2012)	$\delta_1$ , GT 85.5(100 samples)
	Plane Detection (PD)	PlaneRCNN (Liu et al., 2019)	KITTI (Geiger et al., 2012)	$AP^{0.5}$ , GT 0.37

### Selection Criteria

- Recommendation from industry ML researchers and engineers
- Model performance (e.g., accuracy) and efficiency

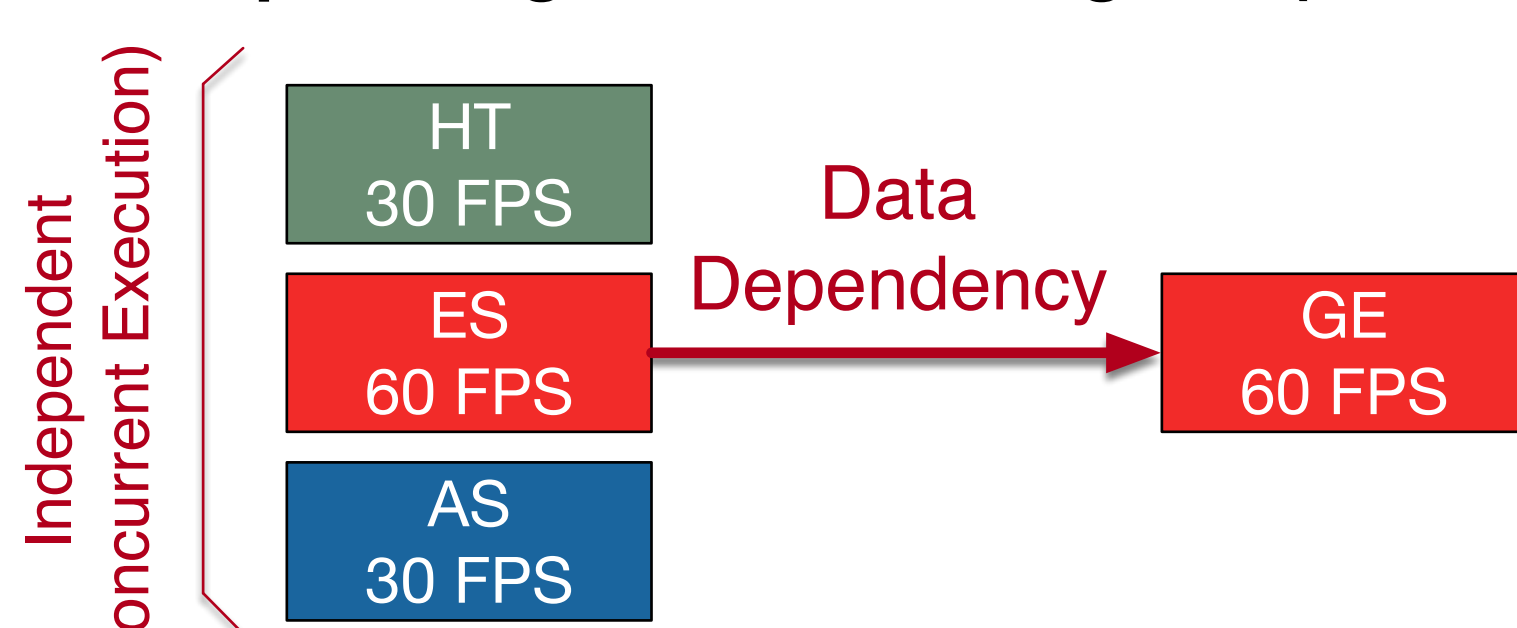


Envisioning On-device AI

### Usage Scenarios

Usage Scenario	HT	Eye Pipeline ES → GE (dep: D)	Speech Pipeline KD → SR (dep: C)	SS	OS	AS	DE	DR	PD	Example Usage Scenario Description
Social Interaction A	30	60	60						30	AR messaging with AR object rendering
Social Interaction B		60	60							In-person interaction with AR glasses
Outdoor Activity A				3	3				10 30	Hiking with smart photo capture
Outdoor Activity B				3	3				30	Rest during hike
AR Assistant				3	3				10 10	Urban walk with informative AR objects
AR Gaming	45								30 30	Gaming with AR object
VR Gaming	45	60	60							Highly-interactive Immersive VR gaming

### Example Usage Scenario Diagram (Social Interaction B)

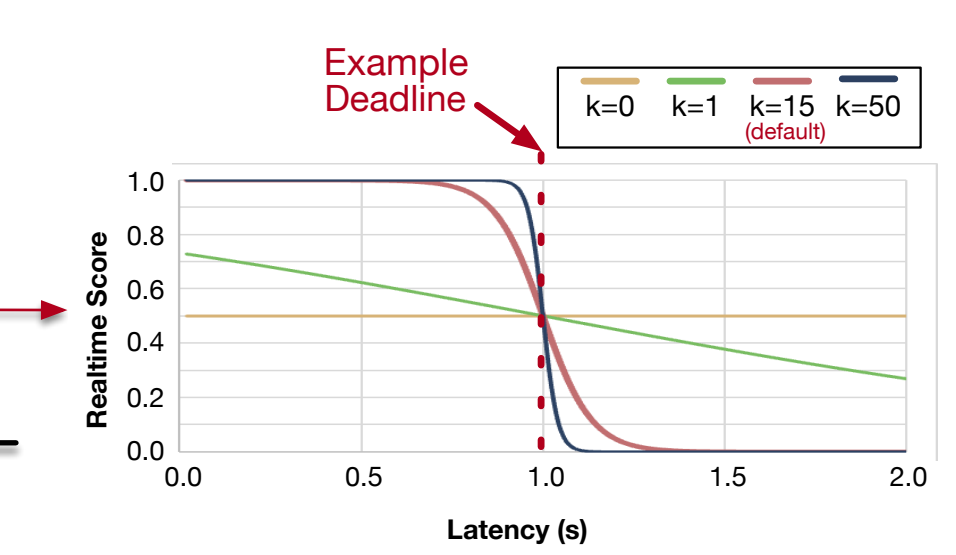


- **ML pipeline:** Cascade multiple models to implement complicated functionalities (e.g., ES → GE in this example; eye pipeline)
- **High FPS eye pipeline:** Enables low latency human-device interaction
- **Dependency in eye pipeline:** ES results are used as inputs of GE

\* HT: hand tracking, ES: eye segmentation  
AS: action segmentation

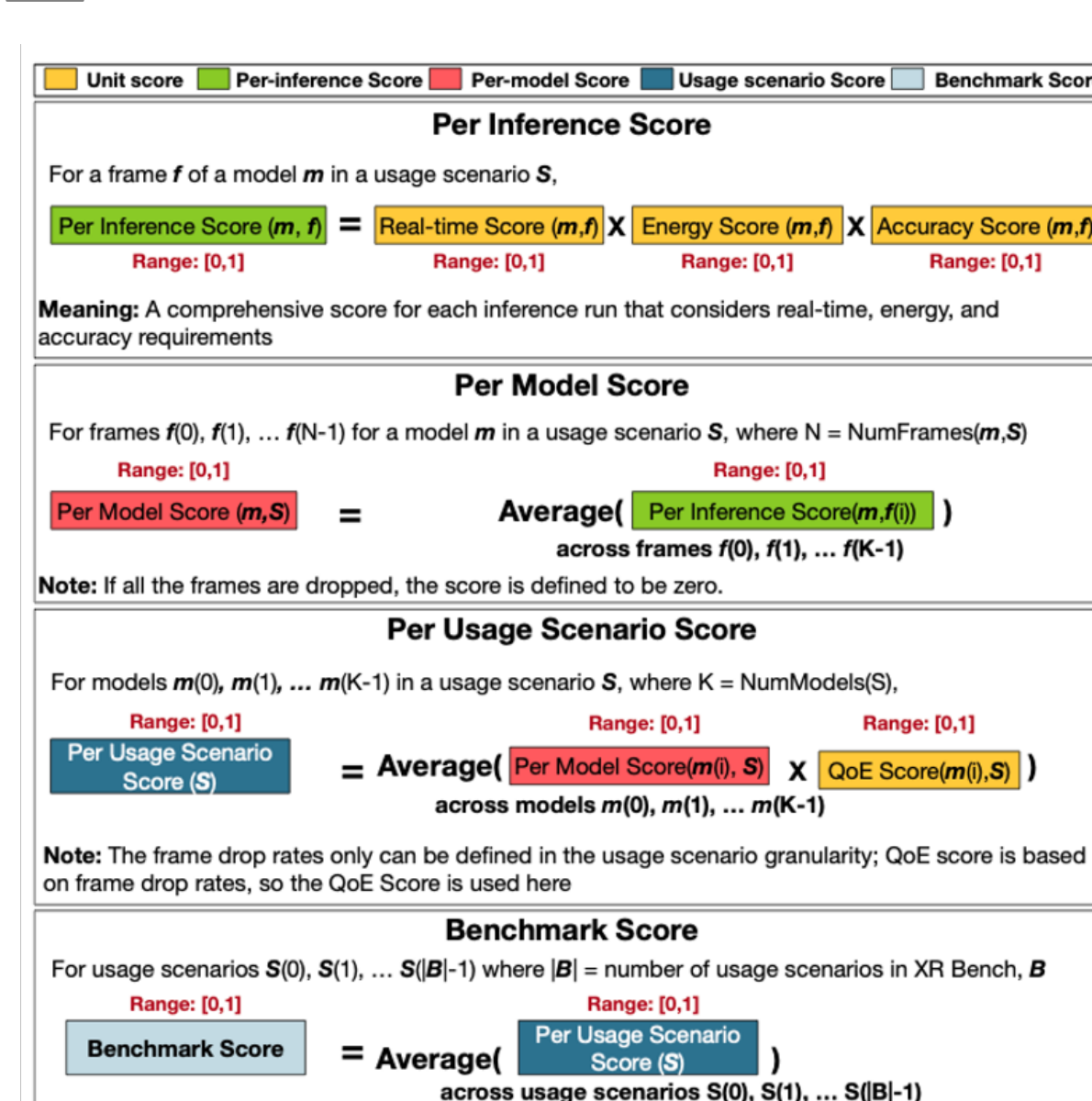
## Score Metric

### Four Unit Scores

Unit Score	What does it measure?
Real-time	Degree of deadline violations (Not absolute latency!) 
Energy	Energy consumption
Accuracy	Relative model performance compared to reported numbers in original papers → mIoU, accuracy, mAP, etc.
Quality of Experience (QoE)	Frame drop rate

All formulated to be higher-is-better metric within [0,1] range

### Benchmark Score



### Hierarchical Formulation

- Define scores from fine-grained to coarse-grained execution units (each inference → all the inference runs for a model → usage scenario → benchmark)

### Composable Formulation

- Range of unit scores: [0,1]
- Product of unit scores ⇒ benchmark score

Single-metric: Provides comprehensive insights and facilitates industry score submissions

(Break-down scores are still available from the benchmark)

Full detailed formulation is available in the paper!

## Case Studies and Conclusion

### Evaluation: Various ML Accelerator Systems

Acc. ID	Acc. Style	Dataflow
A	WS	WS
B	FDA	OS
C		RS
D		WS + WS (1:1 partitioning)
E		OS + OS (1:1 partitioning)
F		RS + RS (1:1 partitioning)
G	SFDA <sup>1</sup>	WS + WS + WS + WS (1:1:1:1 partitioning)
H		OS + OS + OS + OS (1:1:1:1 partitioning)
I		RS + RS + RS + RS (1:1:1:1 partitioning)
J		WS + OS (1:1 partitioning)
K		WS + OS (3:1 partitioning)
L	HDA	WS + OS (1:3 partitioning)
M		WS + OS + WS + OS (1:1:1:1 partitioning)

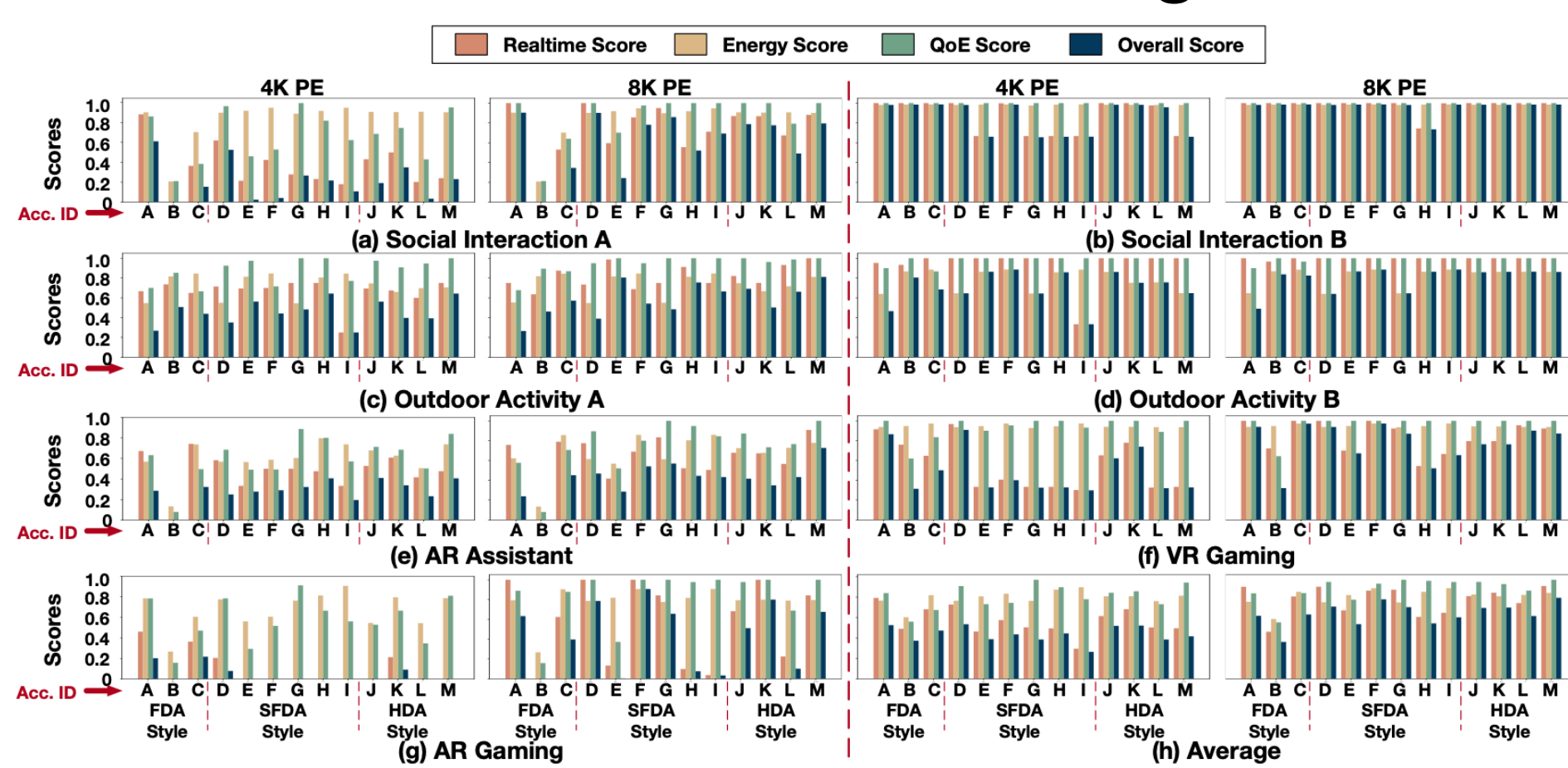
### Accelerator Styles

- FDA: Fixed dataflow accelerator
- SFDA: Scaled-out FDA
- HDA: Heterogeneous dataflow accelerator

### Accelerator Dataflow

- WS: Weight-stationary
- OS: Output-stationary

### Main Evaluation Results and Insights

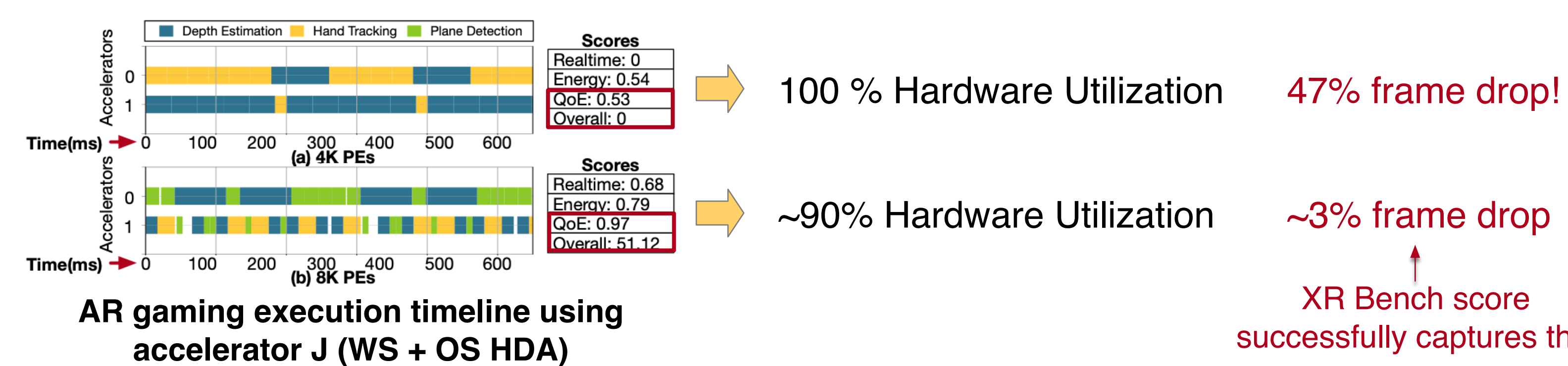


ML Systems for XR needs to be co-designed with usage scenarios

Optimal accelerator style depends on the chip scale

Multi-accelerator systems are friendly to XR systems

### More insights: Example implication to ML system design for XR



AR gaming execution timeline using accelerator J (WS + OS HDA)

Hardware utilization is an incorrect metric for XR ML system design!

### Conclusion

Real-time MTMM workloads have unique characteristics and new implications to ML system design

We developed XRBench to invite everyone for this new problem domain: ML system for real-time MTMM workloads

XRBench is an open project

We look forward to working with the community!

