

# Low-Latency Proactive Continuous Vision

**Yiming Gan**

Department of Computer Science,  
University of Rochester

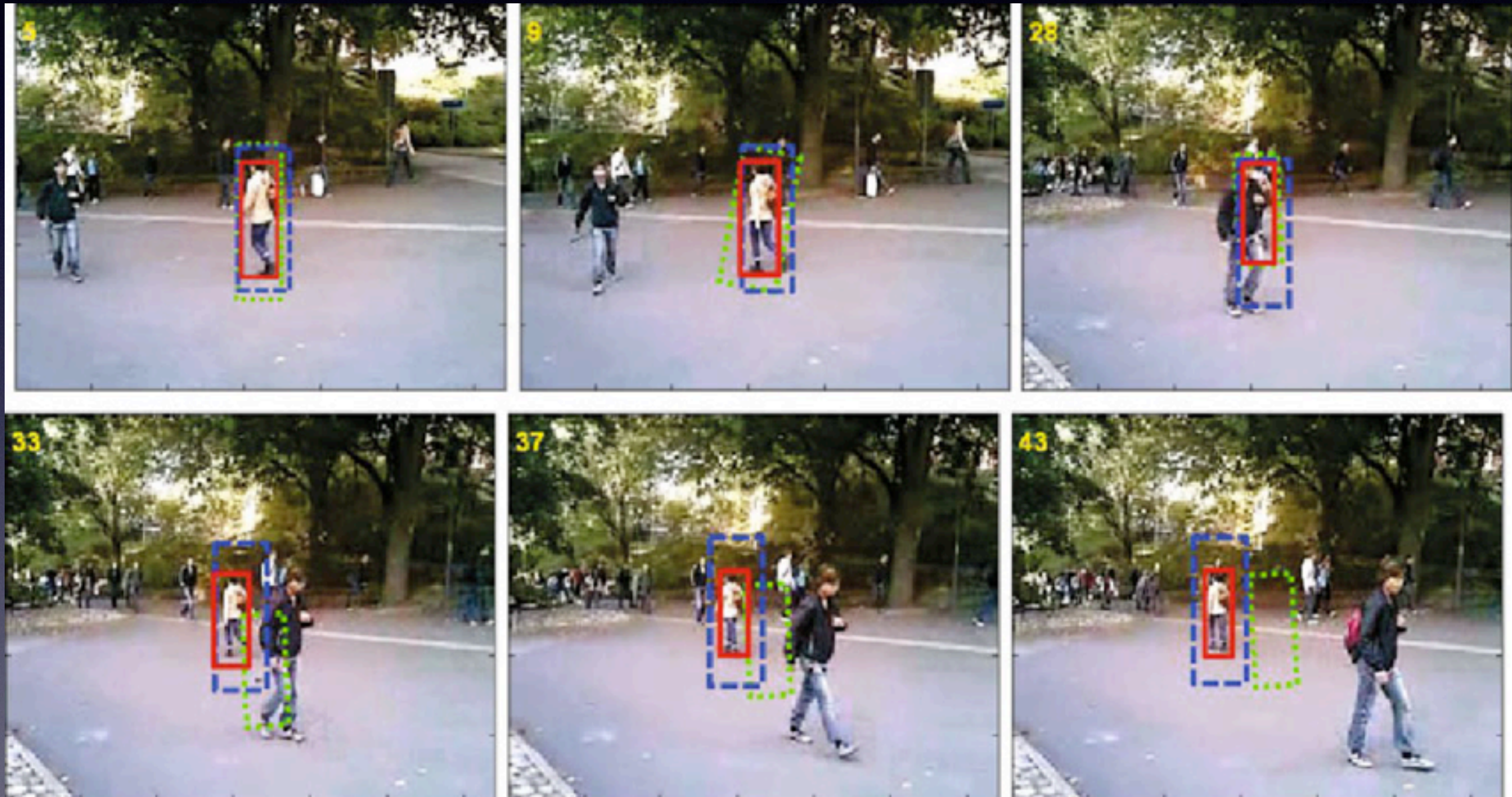
with

Yuxian Qiu,	Shanghai Jiao Tong University
Lele Chen,	University of Rochester
Jingwen Leng,	Shanghai Jiao Tong University
Yuhao Zhu	University of Rochester





# Continuous Vision: Long Frame Latency





# Bottleneck: Serialization

Sensor



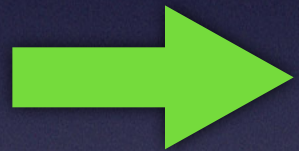


# Bottleneck: Serialization

**Sensor**

**Image Signal Processor**

Light



Raw Pixels



RGB



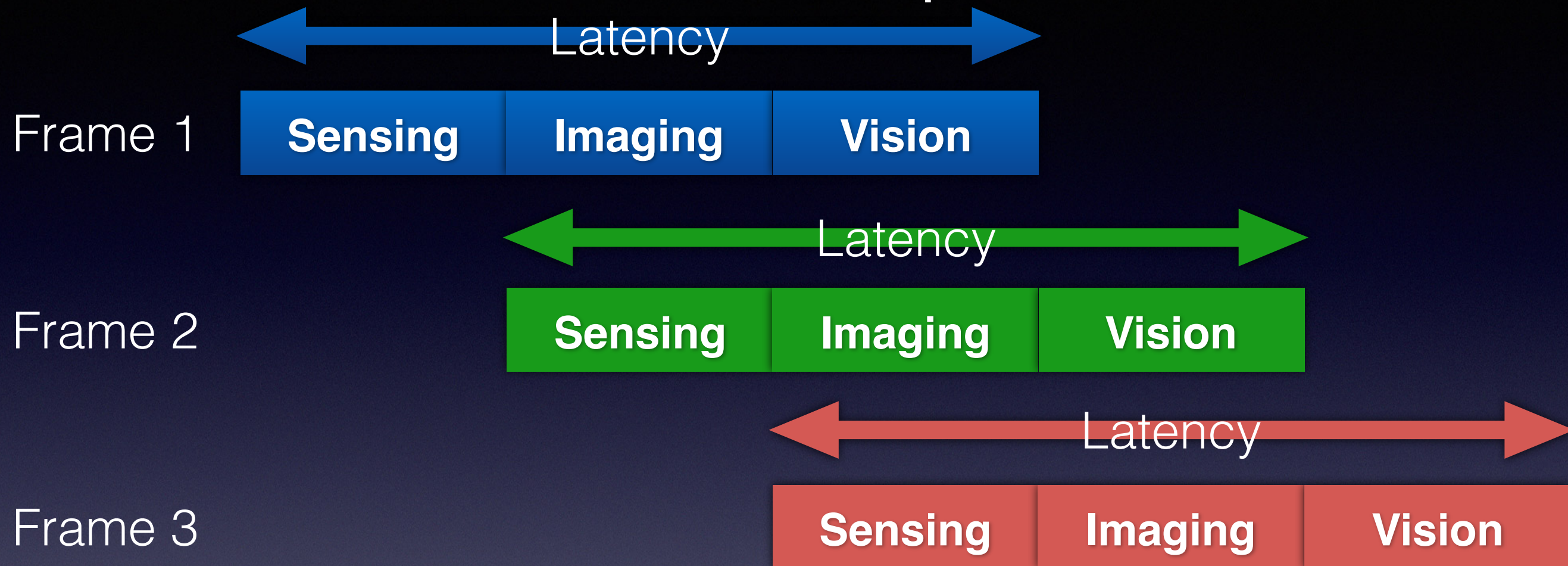
# Bottleneck: Serialization

**Sensor      Image Signal Processor      DNN Accelerator**





# Traditional Pipeline





# Proactive Pipeline



Frame 1

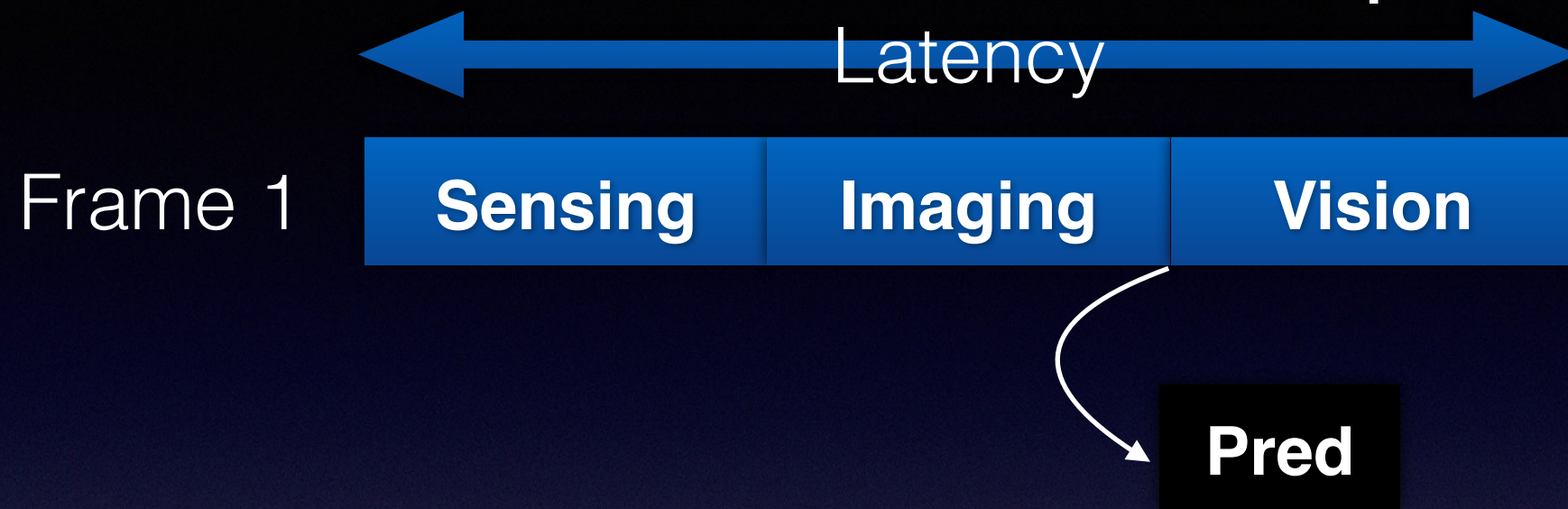
**Sensing**

**Imaging**

**Vision**

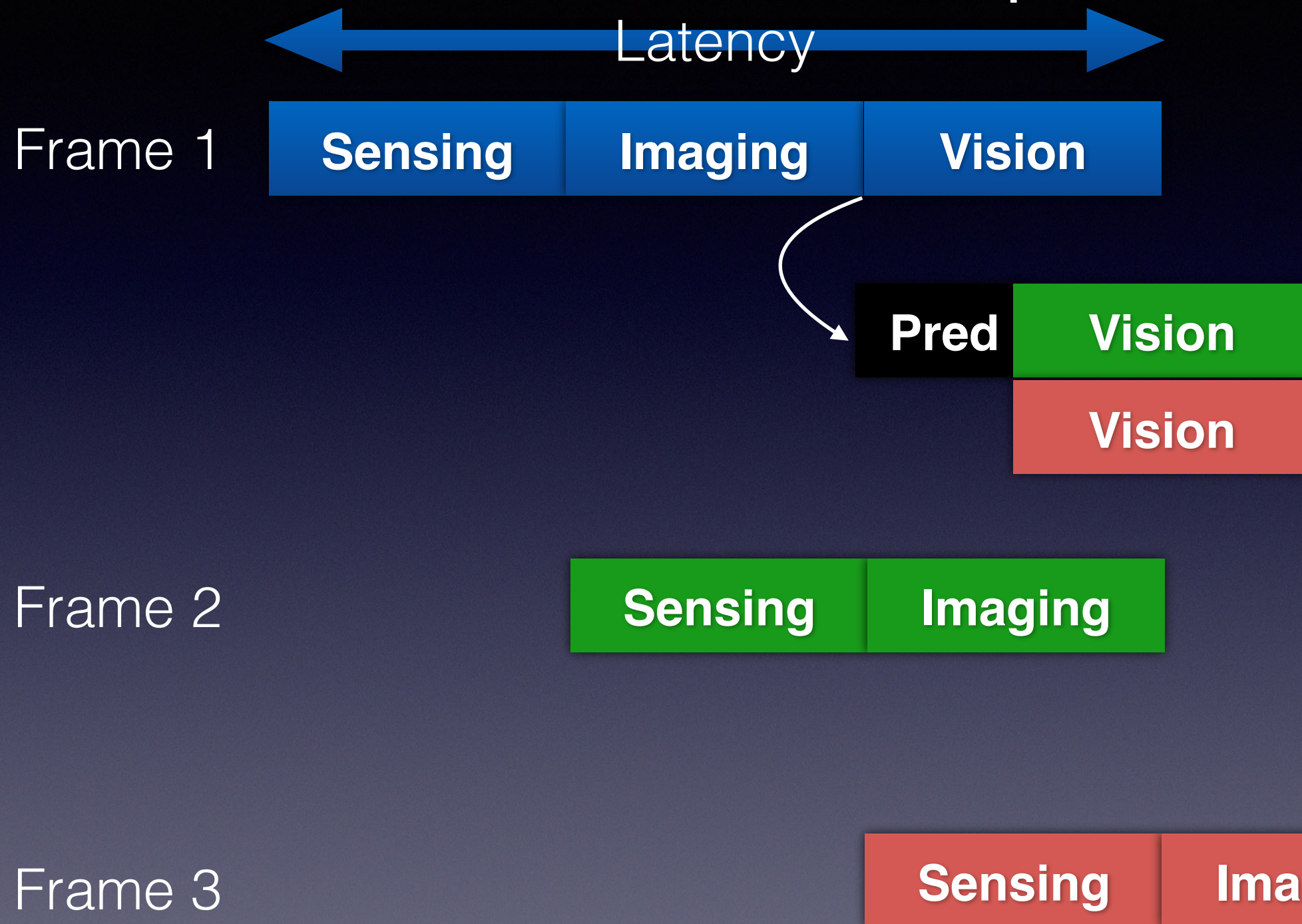


# Proactive Pipeline



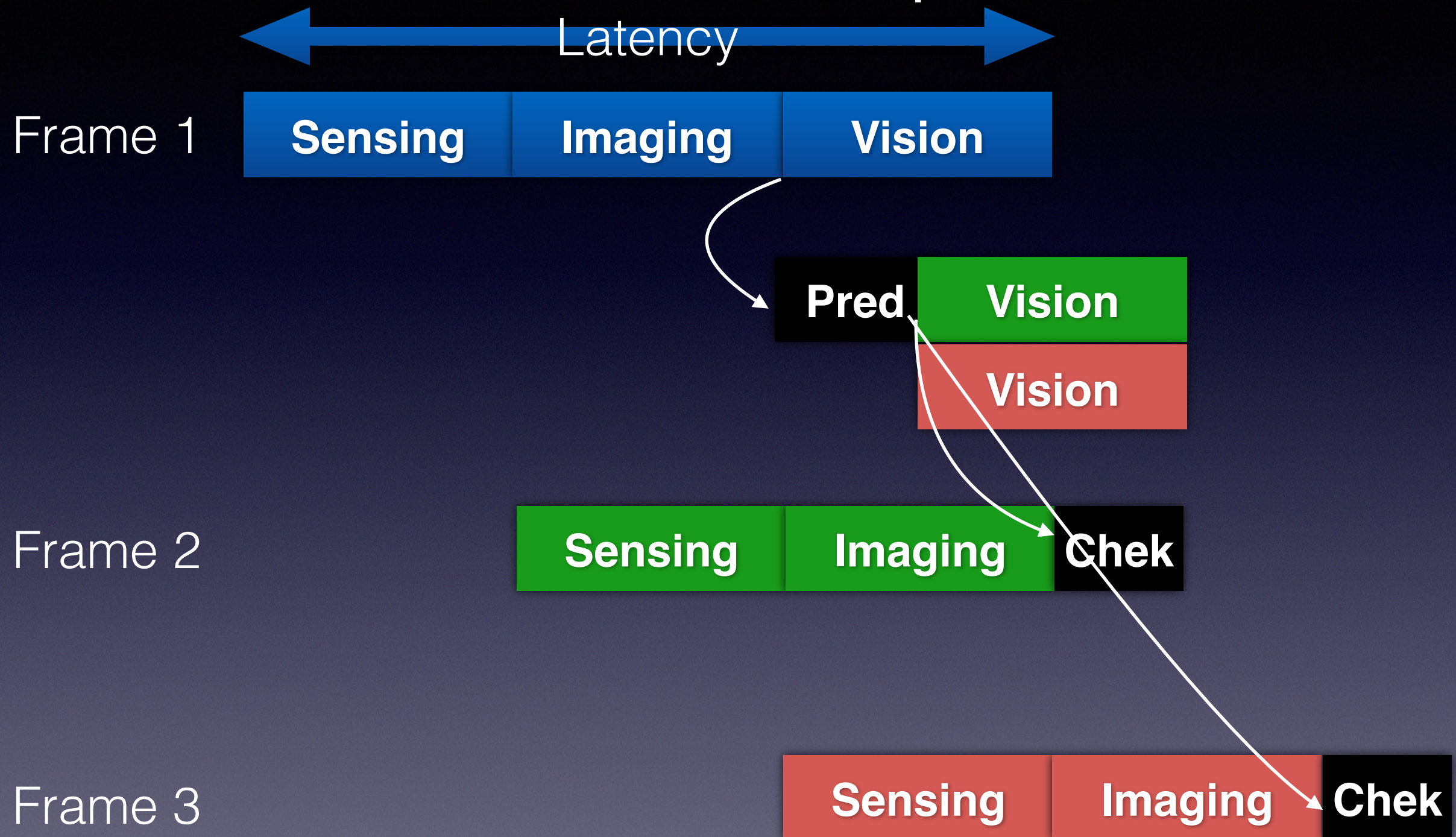


# Proactive Pipeline



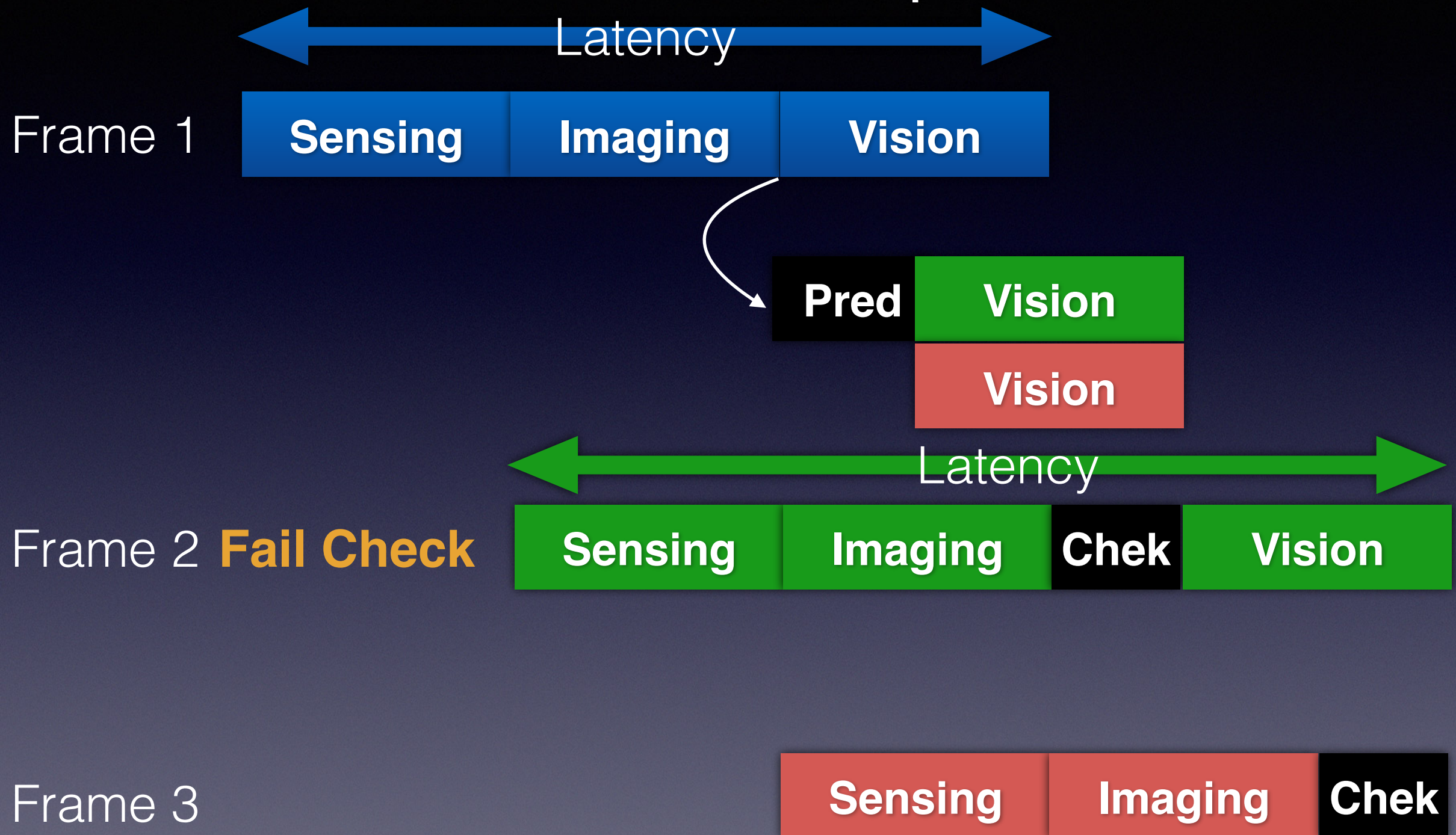


# Proactive Pipeline



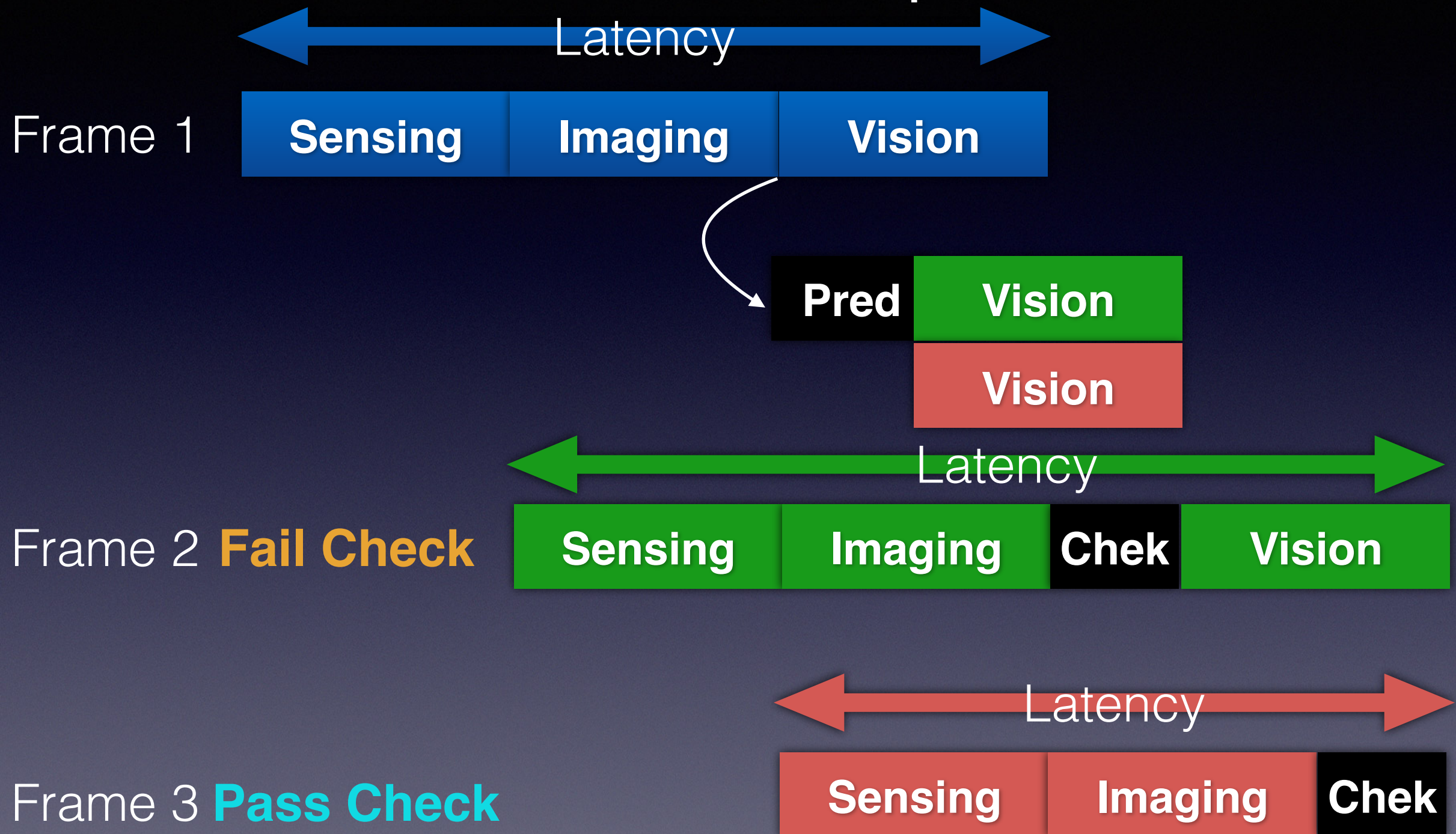


# Proactive Pipeline



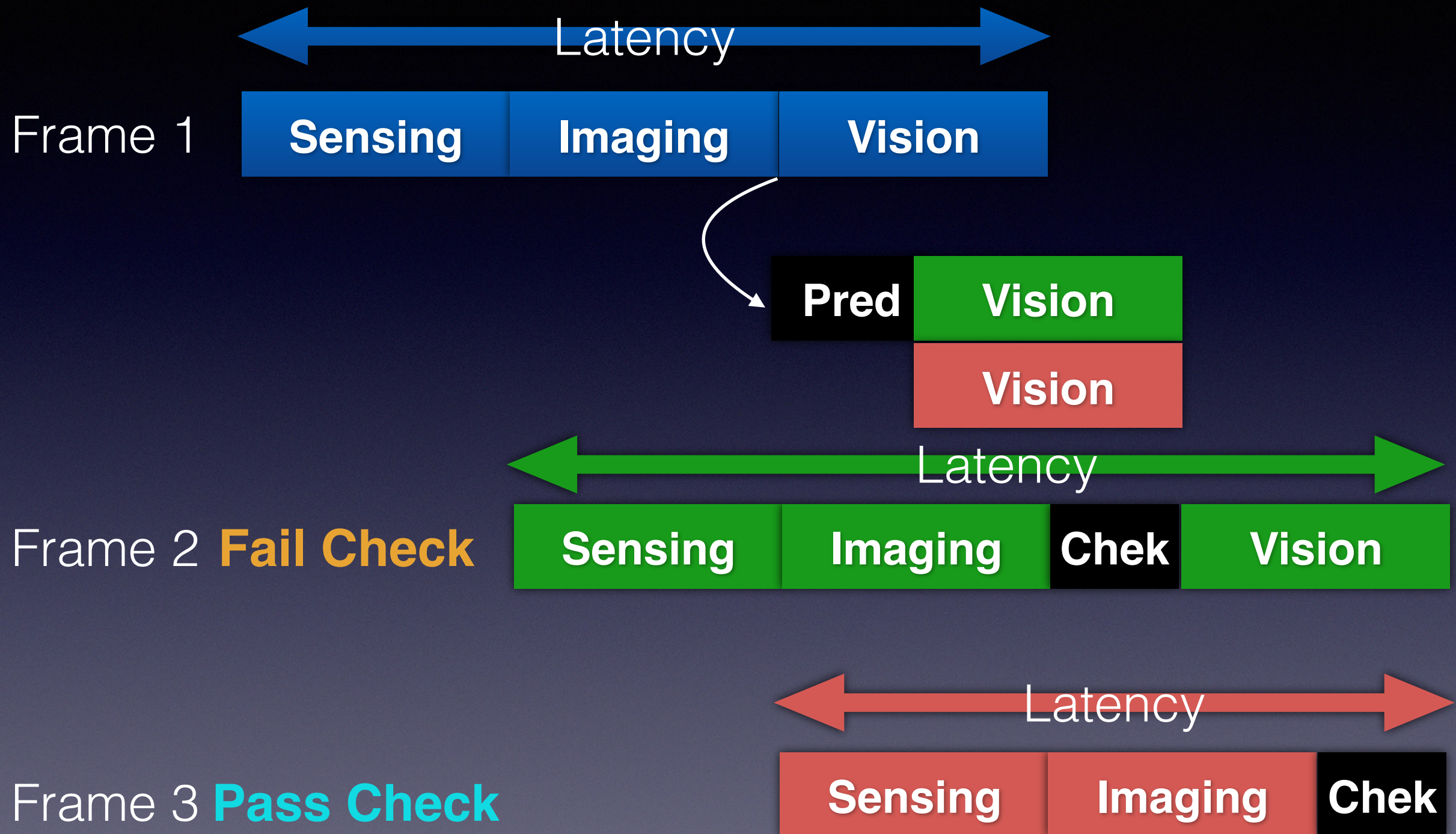


# Proactive Pipeline



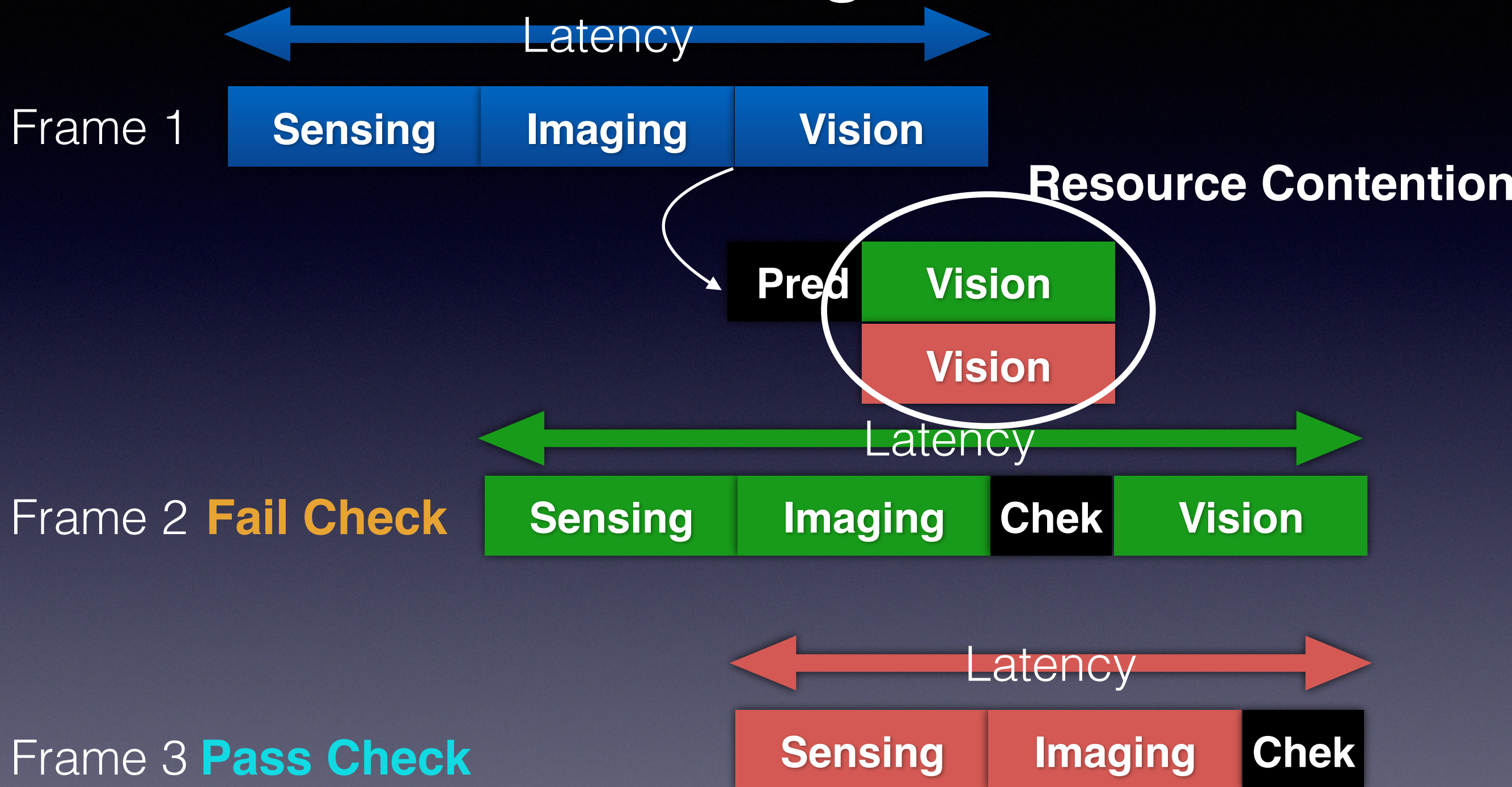


# Gains





# Challenges





# Solutions

## Snapdragon 675

### CPU

4<sup>th</sup> gen Kryo CPU  
Performance @ 2.0GHz  
Efficiency @ 1.7GHz

### Artificial Intelligence

3<sup>rd</sup> gen Qualcomm® AI Engine

### Modem

Snapdragon X12  
Cat 12 DL, up to 600 Mbps

### Audio

Qualcomm Aqstic™ audio  
Qualcomm® aptX™ audio

### Display

Up to FHD+ display



### GPU

6<sup>th</sup> gen Adreno GPU

### DSP

6<sup>th</sup> gen Hexagon DSP  
DSP Security

### Camera

2<sup>nd</sup> gen Qualcomm Spectra ISP  
Up to 25 Megapixels @ 30fps ZSL  
48 Megapixels snapshot

### Charging

Qualcomm® Quick  
Charge™ 4+ technology

Commercial devices expected in Q1 2019



# Solutions

## Snapdragon 675

### CPU

4<sup>th</sup> gen Kryo CPU  
Performance @ 2.0GHz  
Efficiency @ 1.7GHz

### Artificial Intelligence

3<sup>rd</sup> gen Qualcomm® AI Engine

### Modem

Snapdragon X12  
Cat 12 DL, up to 600 Mbps

### Audio

Qualcomm Aqstic™ audio  
Qualcomm® aptX™ audio

### Display

Up to FHD+ display



### GPU

6<sup>th</sup> gen Adreno GPU

### DSP

6<sup>th</sup> gen Hexagon DSP  
DSP Security

### Camera

2<sup>nd</sup> gen Qualcomm Spectra ISP  
Up to 25 Megapixels @ 30fps ZSL  
48 Megapixels snapshot

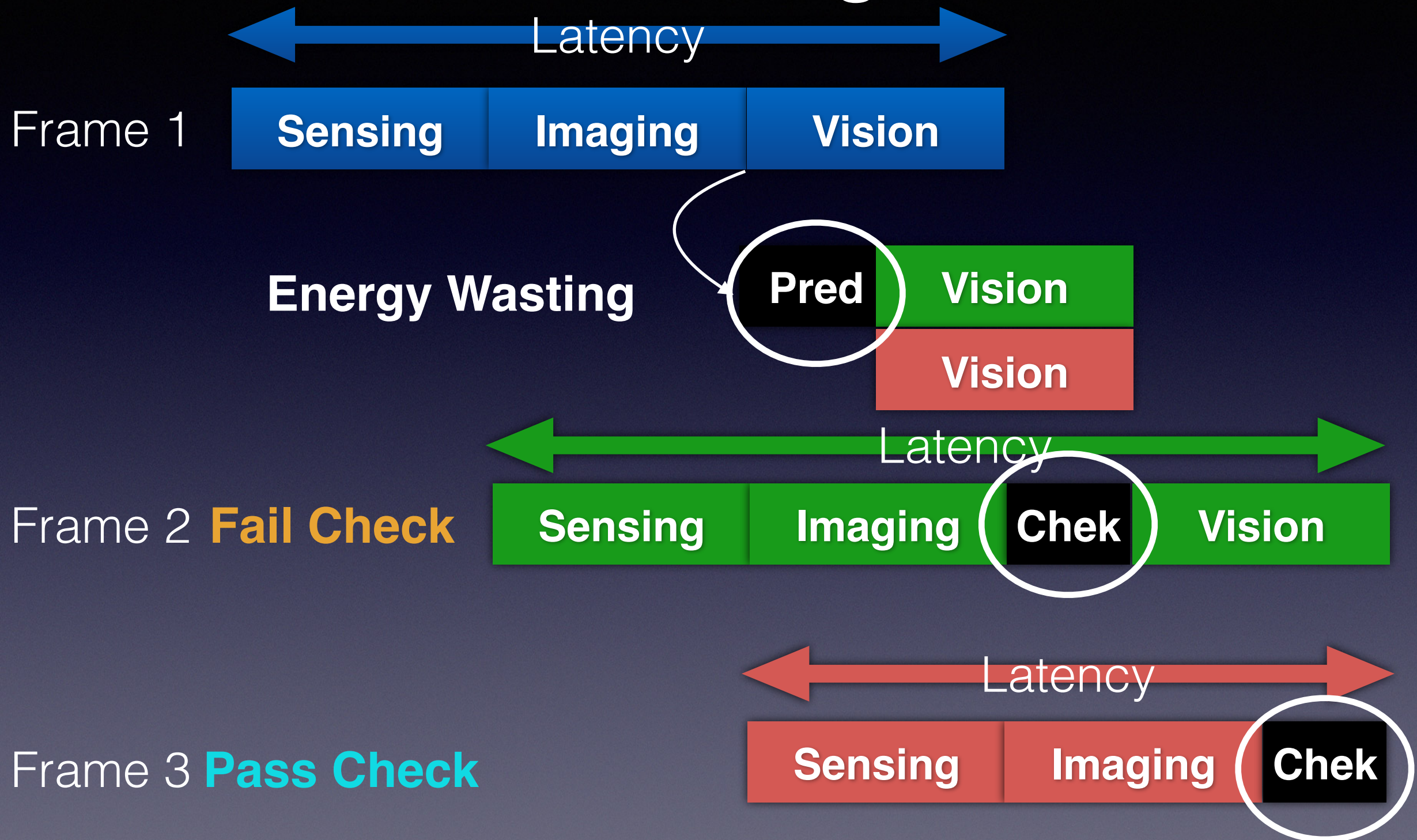
### Charging

Qualcomm® Quick  
Charge™ 4+ technology

Commercial devices expected in Q1 2019



# Challenges





# Solutions

- Relaxing Checking Criterion (Threshold  $T$ )

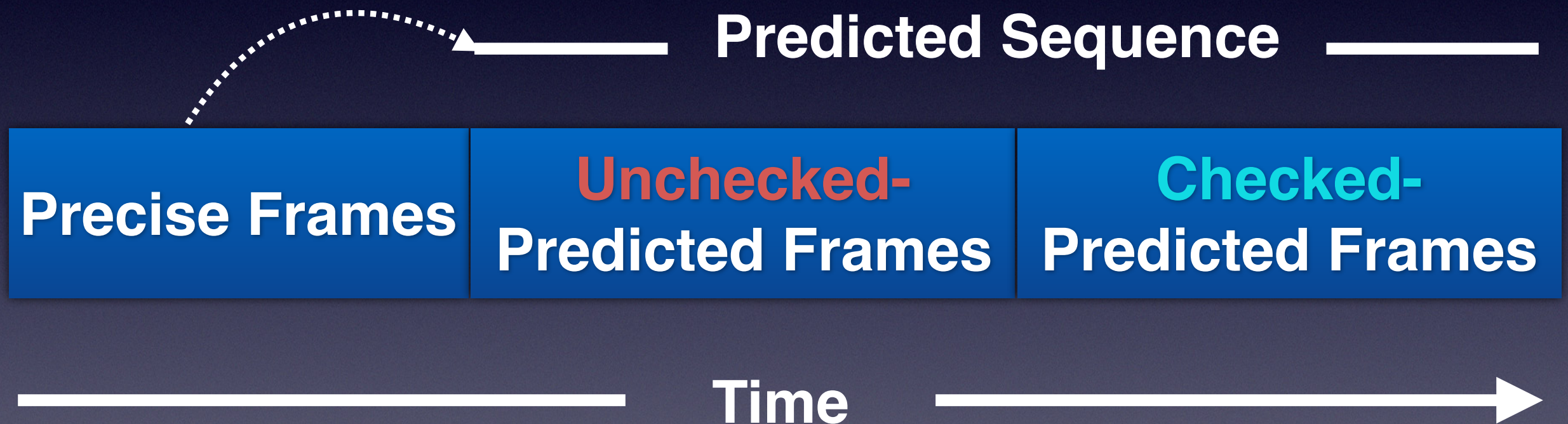


# Solutions

- Relaxing Checking Criterion (Threshold  $T$ )
- Relaxing Checking Frequency (Degree  $K$ )



# Frames Sequence

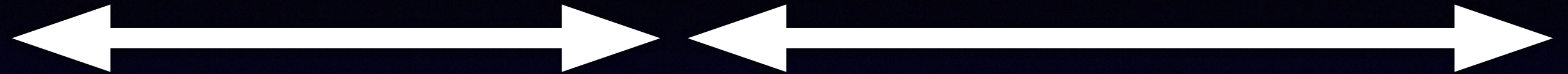




# PVF Framework

**Static**

**Dynamic**



**Vision  
Apps**

**SoC**

**CPU**

**NPU**

**DSP**

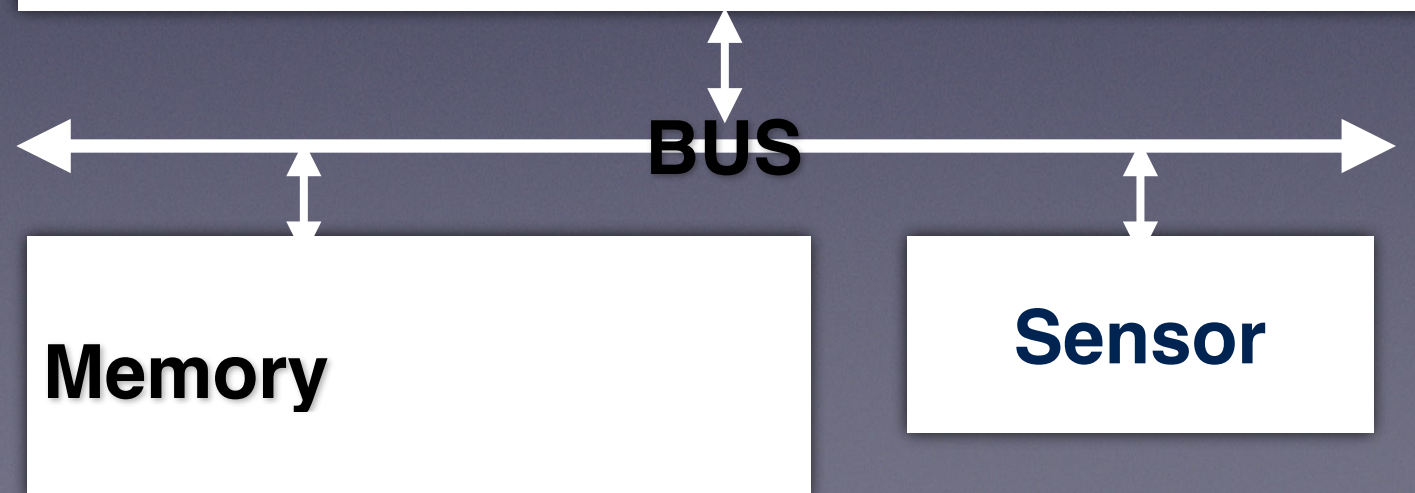
**GPU**

**ISP**

**BUS**

**Memory**

**Sensor**

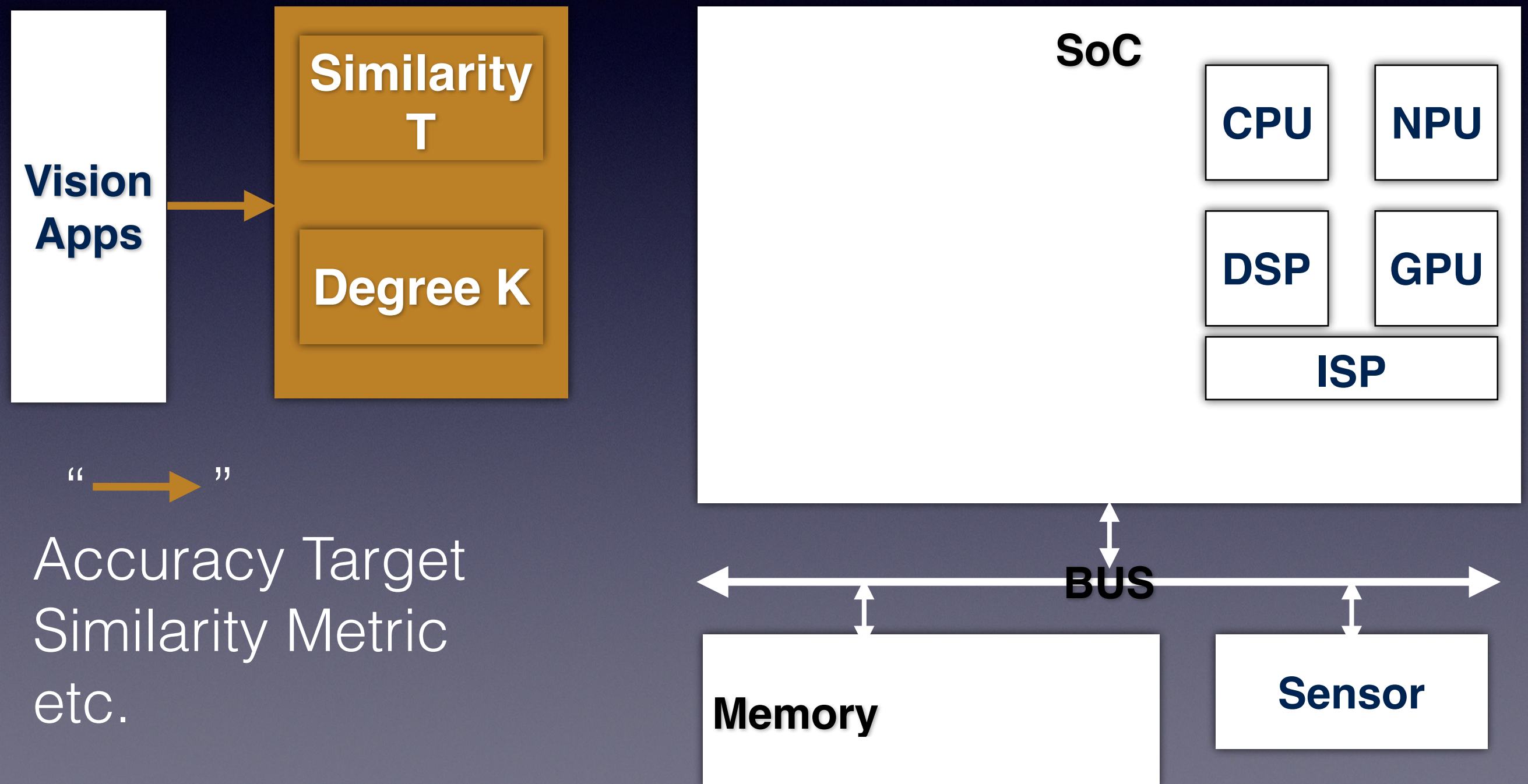




# PVF Framework

**Static**

**Dynamic**

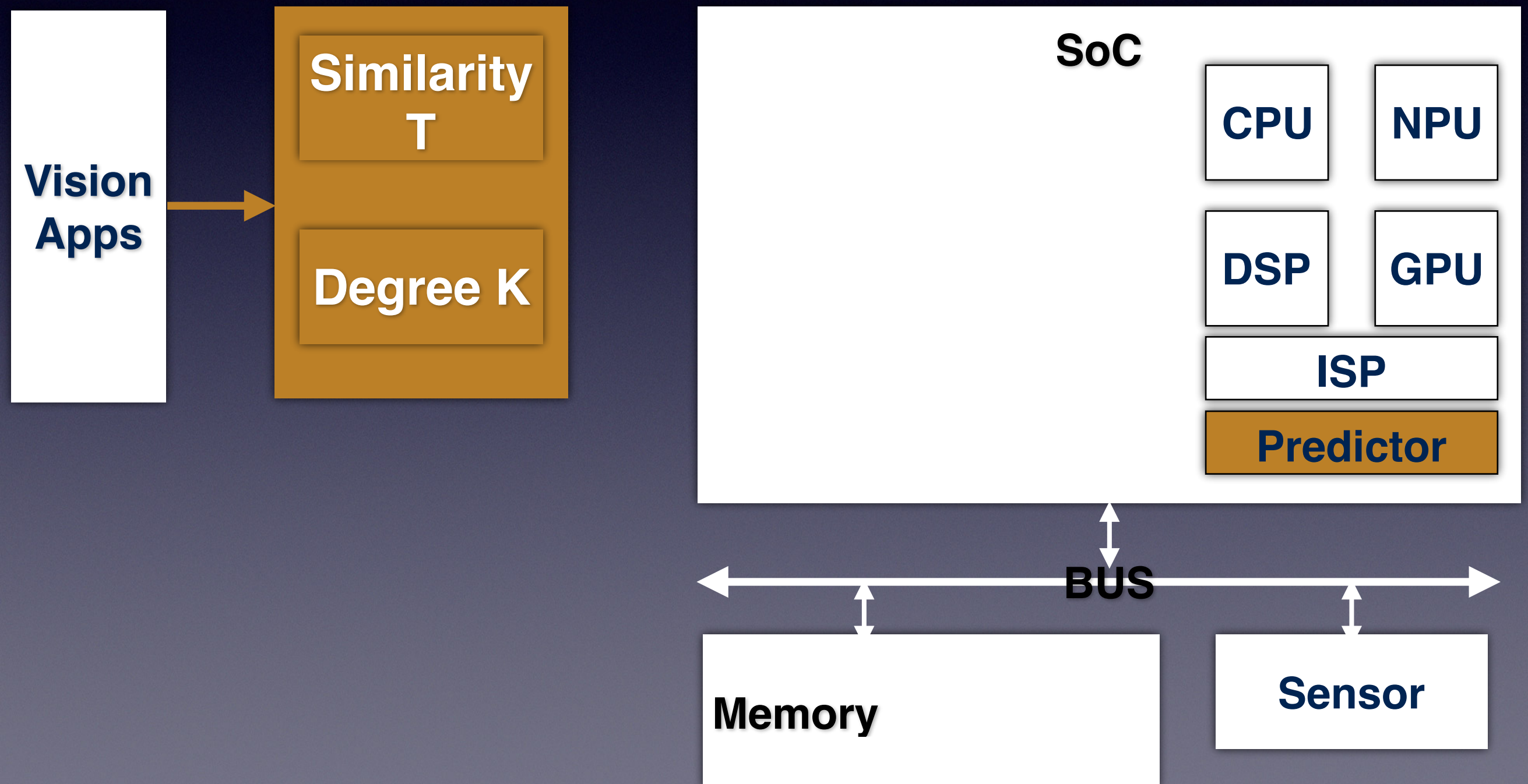




# PVF Framework

**Static**

**Dynamic**

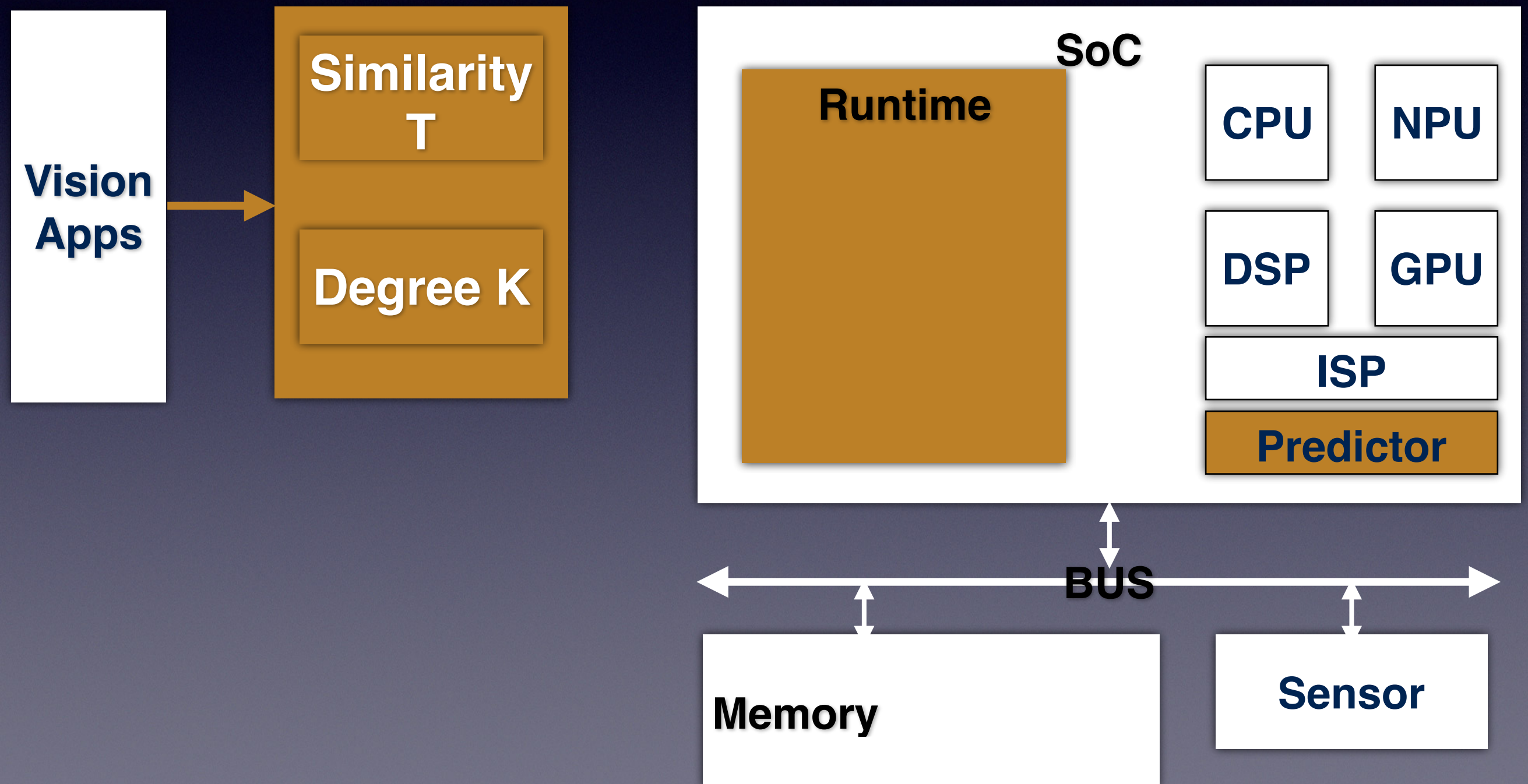




# PVF Framework

**Static**

**Dynamic**

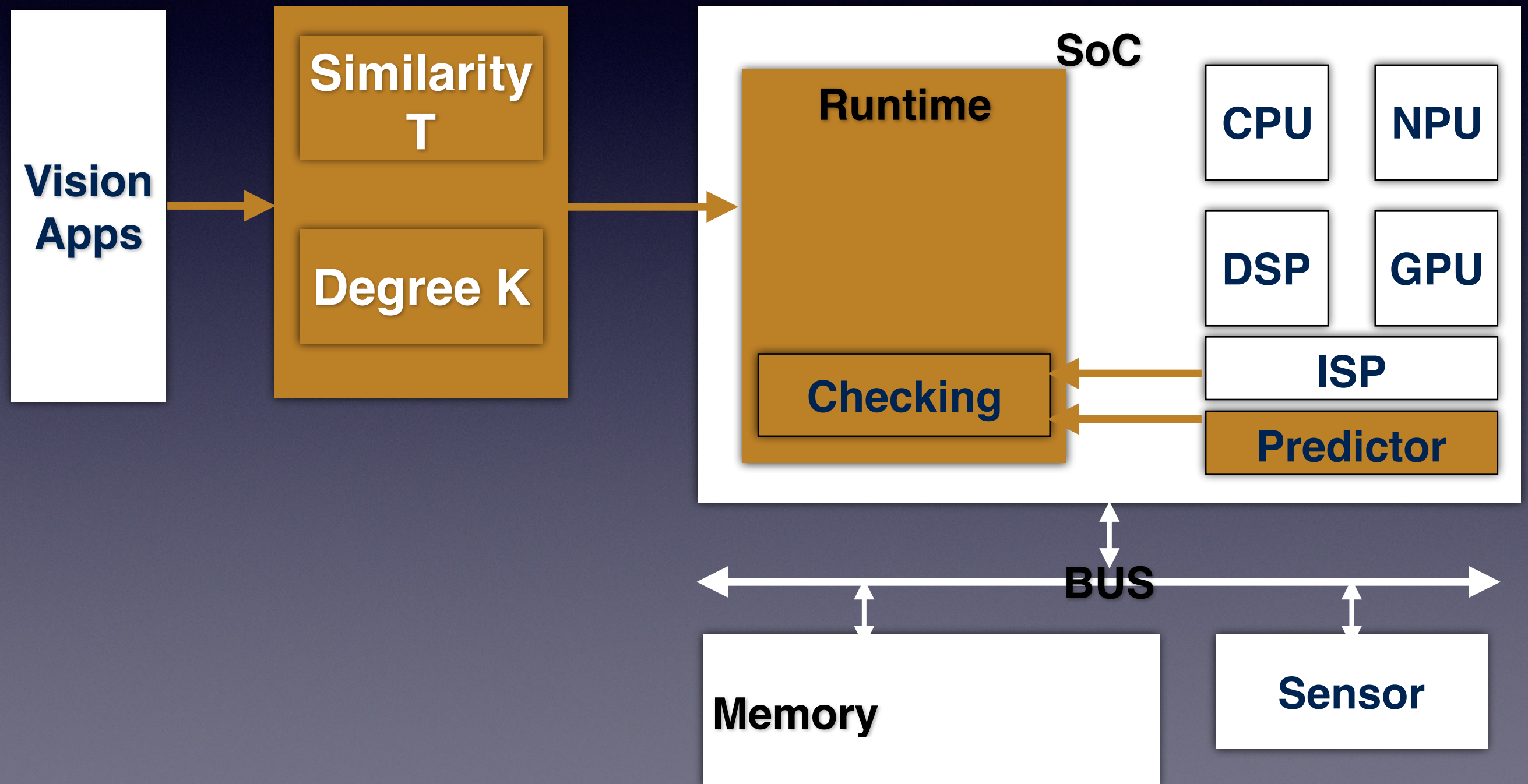




# PVF Framework

Static

Dynamic

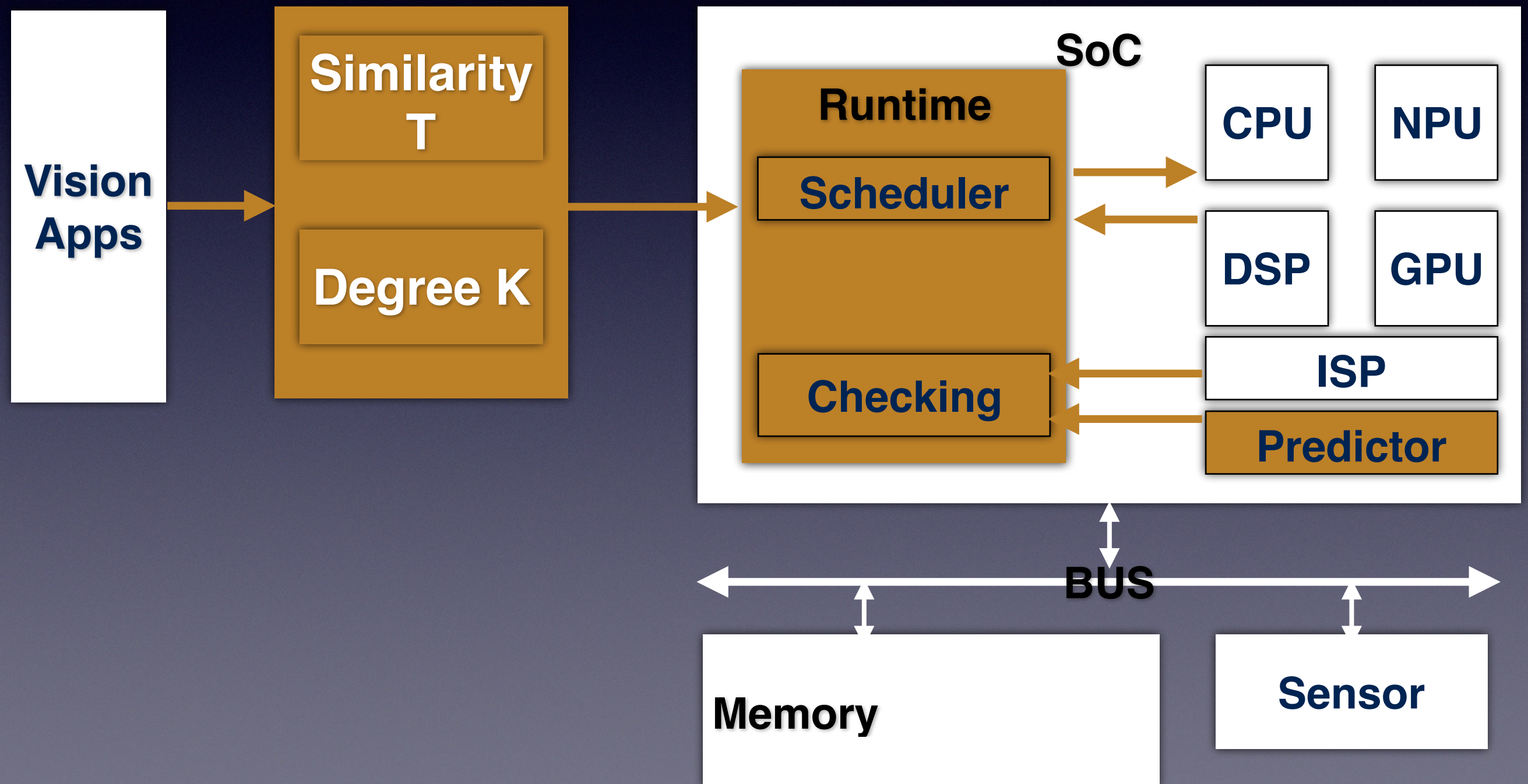




# PVF Framework

**Static**

**Dynamic**

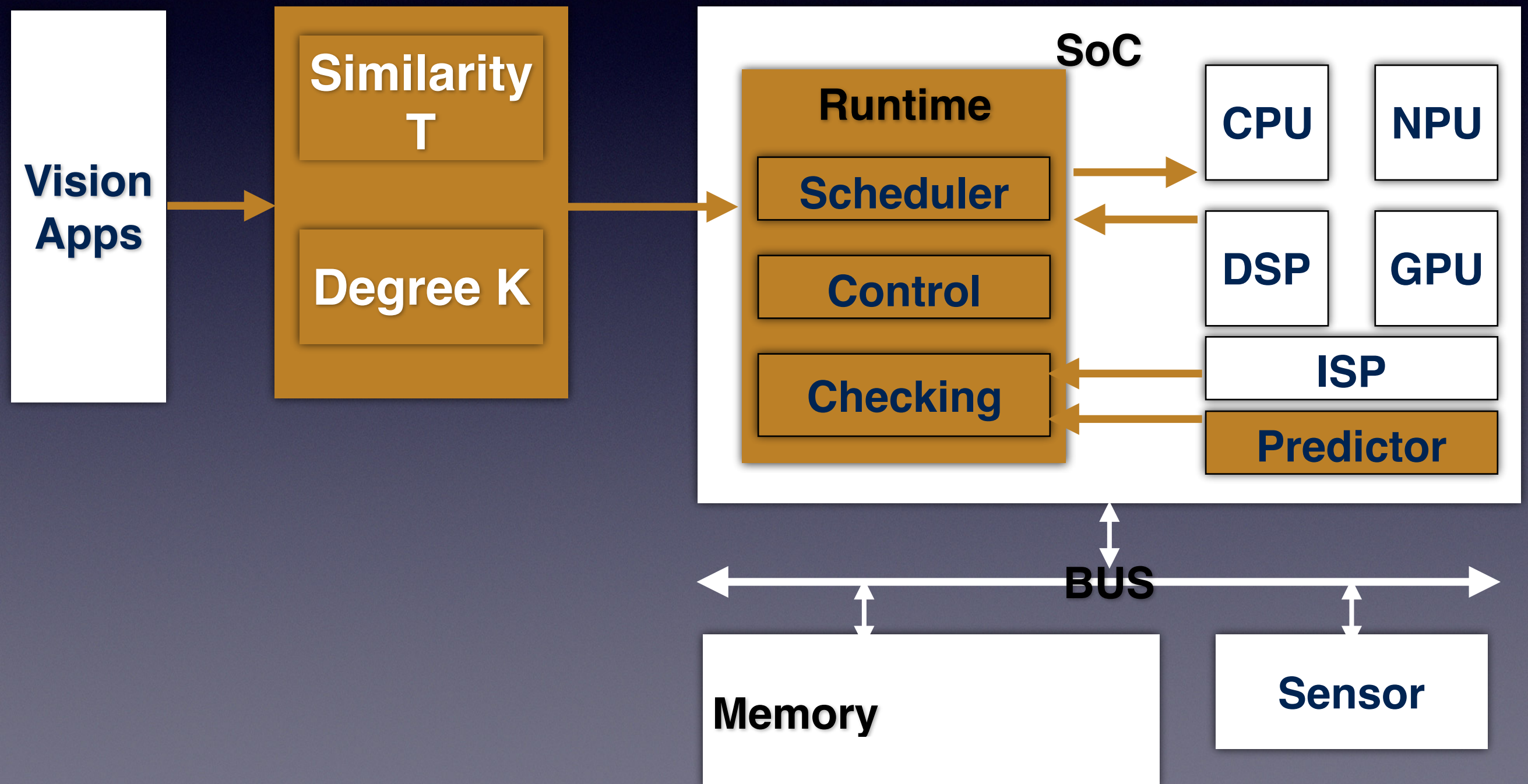




# PVF Framework

**Static**

**Dynamic**





# Experimental Setup

## **I. In house simulator modeling state-of-the art SoCs**

- Real measurement of latency and energy on different IPs.



# Experimental Setup

## **I. In house simulator modeling state-of-the art SoCs**

- Real measurement of latency and energy on different IPs.

## **II. RTL Implementations for NPU and Predictor**

- 20x20 Systolic Array for NPU, 10x10 Systolic Array for Predictor



# Experimental Setup

## **I. In house simulator modeling state-of-the art SoCs**

- Real measurement of latency and energy on different IPs.

## **II. RTL Implementations for NPU and Predictor**

- 20x20 Systolic Array for NPU, 10x10 Systolic Array for Predictor

## **III. Evaluate on Object Detection and Tracking**

- KITTI dataset for object detection, VOT-challenge for tracking.



# Experimental Setup

## **I. In house simulator modeling state-of-the art SoCs**

- Real measurement of latency and energy on different IPs.

## **II. RTL Implementations for NPU and Predictor**

- 20x20 Systolic Array for NPU, 10x10 Systolic Array for Predictor

## **III. Evaluate on Object Detection and Tracking**

- KITTI dataset for object detection, VOT-challenge for tracking.

## **IV. Different Input Resolutions**



# Baselines

## **I. Base**

- Baseline with traditional execution pipeline

## **II. BO**

- Baseline with optimized back-end

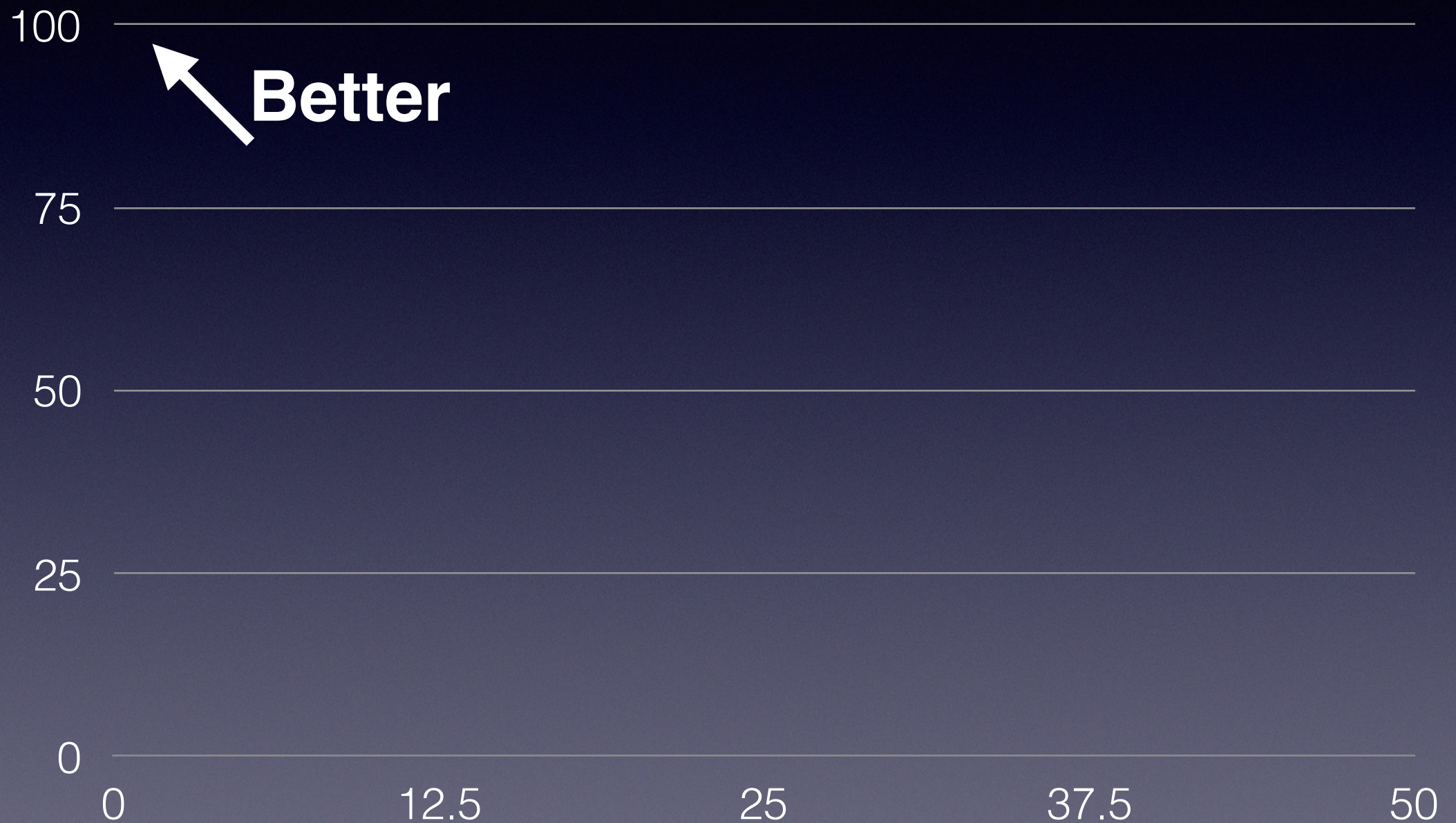
## **III. FCFS**

- Traditional pipeline with multiple hardware IPs



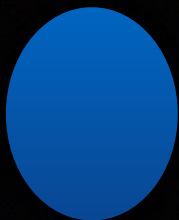
# Results

Latency Reduction (%)



Energy Budget (mJ)

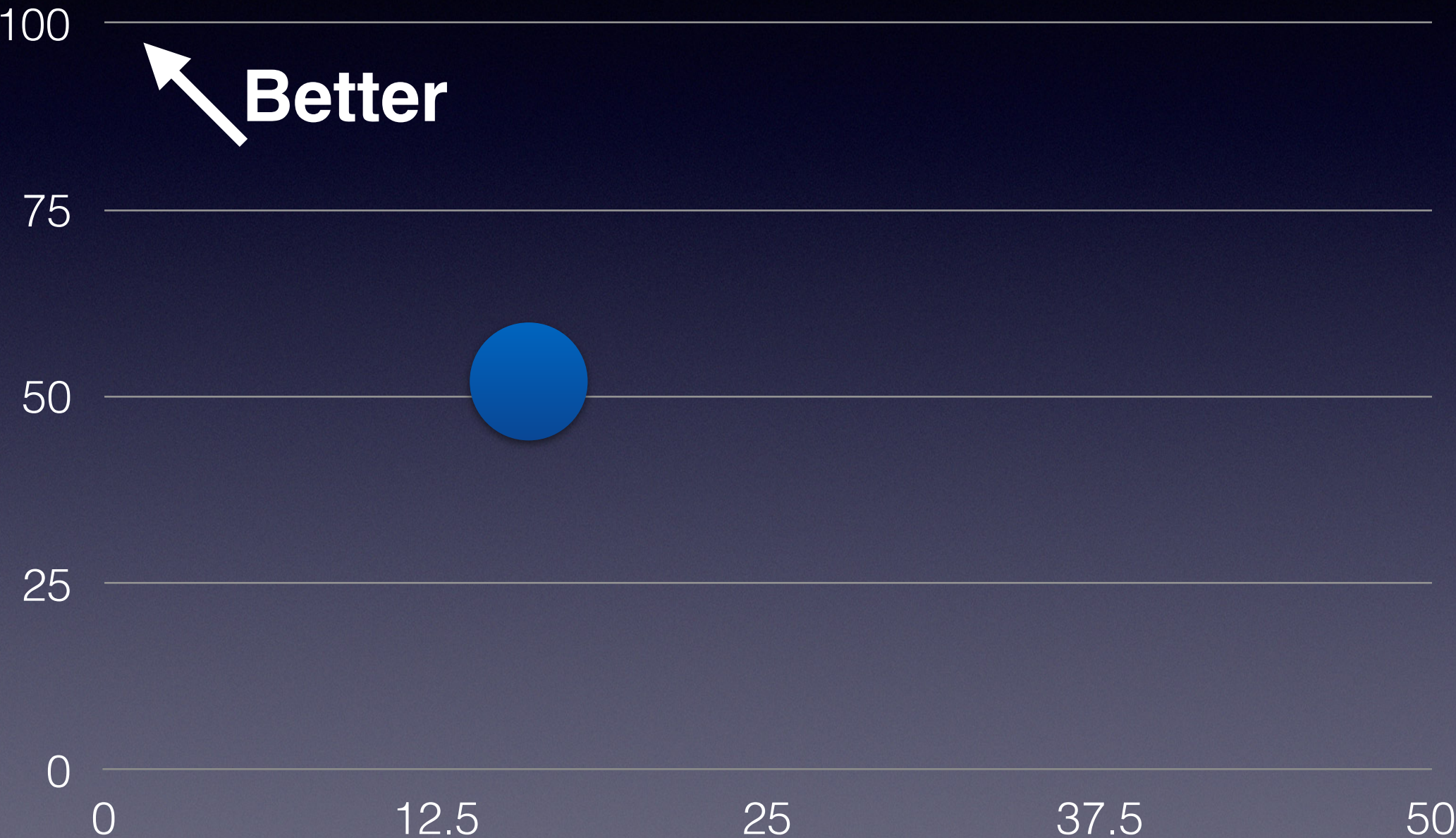




PVF

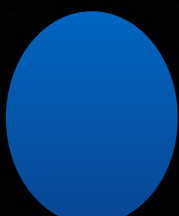
# Results

Latency Reduction (%)



Energy Budget (mJ)





PVF

# Results



Base



BO



FCFS

Latency Reduction (%)

100

75

50

25

0



Better

0

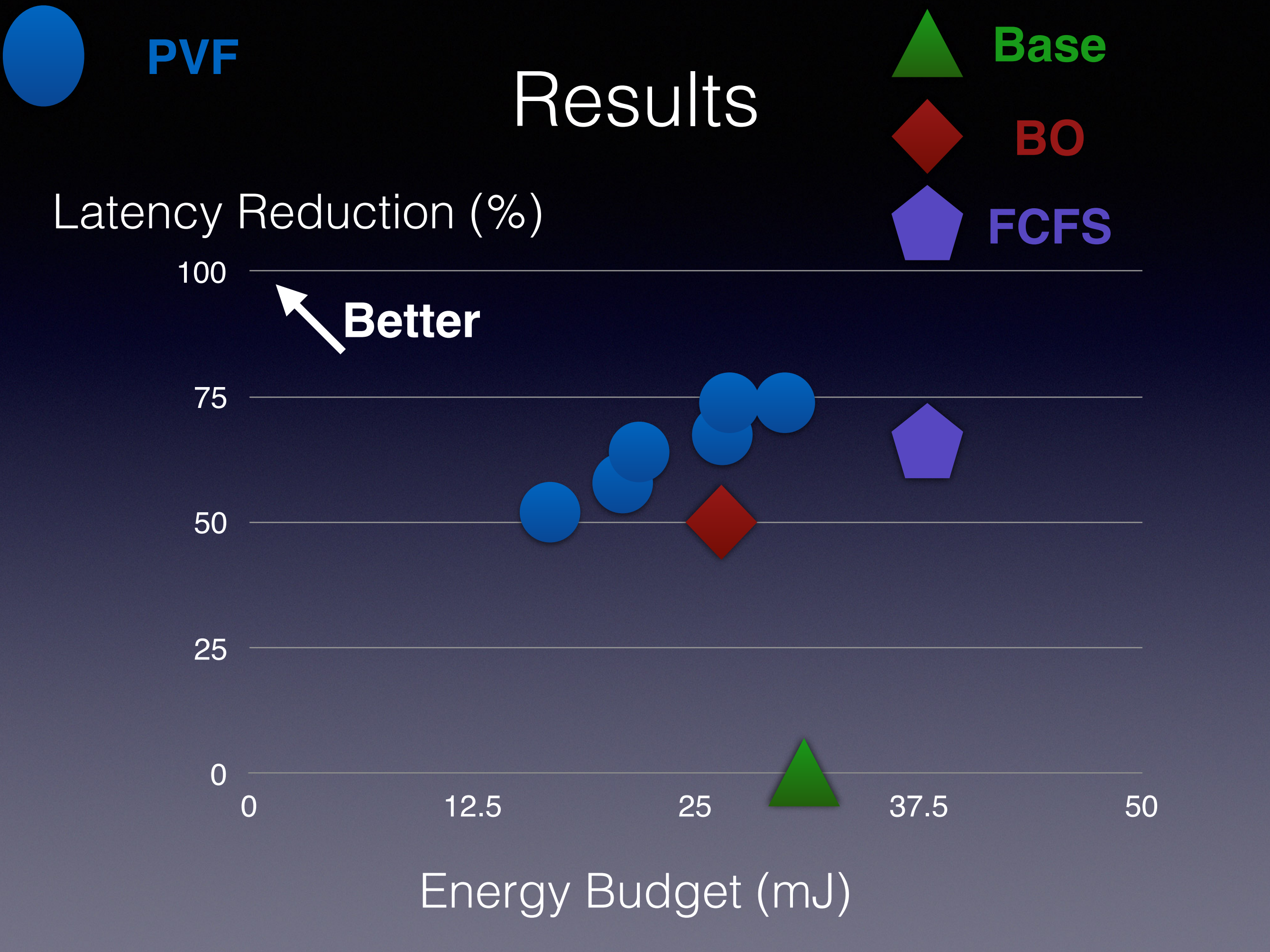
12.5

25

37.5

50

Energy Budget (mJ)





# Conclusion

## **I. Long Latency Bottleneck Continuous Vision**

## **II. Proactive Execution Pipeline**

- 1) Leveraging Heterogeneities in Mobile SoCs
- 2) Relaxed Checking

## **III. Non-mission-critical System**



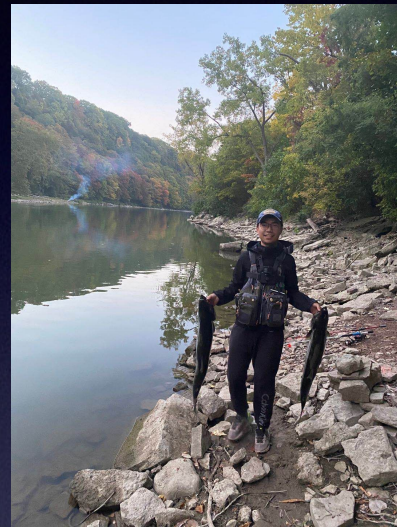
# Collaborators



Yuxian Qiu



Jingwen Leng



Lele Chen



Yuhao Zhu



# Questions