## **Re-thinking CNN Frameworks for Time-Sensitive Autonomous-Driving Applications:** Addressing an Industrial Challenge



NORTH CAROLINA at CHAPEL HILL

- Ming Yang<sup>1</sup>, Shige Wang<sup>2</sup>, Joshua Bakita<sup>1</sup>, Thanh Vu<sup>1</sup>, F. Donelson Smith<sup>1</sup>, James H. Anderson<sup>1</sup>, and Jan-Michael Frahm<sup>1</sup>
  - <sup>1</sup>The University of North Carolina at Chapel Hill <sup>2</sup>General Motors Research



## Re-thinking CNN Frameworks for Time-Sensitive Autonomous-Driving Applications: Addressing an Industrial Challenge

Ming Yang<sup>1</sup>, Shige Wang<sup>2</sup>, Joshua Bakita<sup>1</sup>, Thanh Vu<sup>1</sup>, F. Donelson Smith<sup>1</sup>, James H. Anderson<sup>1</sup>, and Jan-Michael Frahm<sup>1</sup>

<sup>1</sup>The University of North Carolina at Chapel Hill <sup>2</sup>General Motors Research



THE UNIVERSITY of NORTH CAROLINA at CHAPEL HILL



З



З



З



З



https://blogs.nvidia.com/blog/2016/01/04/ automotive-nvidia-drive-px-2/







https://blogs.nvidia.com/blog/2016/01/04/ automotive-nvidia-drive-px-2/







Icons made by <u>Freepik</u> from <u>Flaticon</u> is licensed by <u>CC 3.0 BY</u>





https://blogs.nvidia.com/blog/2016/01/04/





Icons made by <u>Freepik</u> from <u>Flaticon</u> is licensed by <u>CC 3.0 BY</u>





https://blogs.nvidia.com/blog/2016/01/04/





# 2. Accuracy

Icons made by <u>Freepik</u> from Flaticon is licensed by CC 3.0 BY





https://blogs.nvidia.com/blog/2016/01/04/





- 2. Accuracy
- 3. Throughput

Icons made by <u>Freepik</u> from Flaticon is licensed by CC 3.0 BY





https://blogs.nvidia.com/blog/2016/01/04/ automotive-nvidia-drive-px-2/

- 1. Response time
- 2. Accuracy
- 3. Throughput





Ming Yang - RTAS 2019

# Our focus



https://blogs.nvidia.com/blog/2016/01/04/ automotive-nvidia-drive-px-2/

- 1. Response time
- 2. Accuracy
- 3. Throughput





Ming Yang - RTAS 2019

# Our focus



**NVIDIA** TEGRA

https://blogs.nvidia.com/blog/2016/01/04/ automotive-nvidia-drive-px-2/

- 1. Response time
- 2. Accuracy
- 3. Throughput

**CNN software** underutilizes the hardware.

# Our focus





















**Issues:** 



### **Issues**:

1. Memory requirements multiply, limiting the number of instances.



### **Issues**:

- 1. Memory requirements multiply, limiting the number of instances. 2. Context switches on GPU cause overheads.



### **Issues:**

- 1. Memory requirements multiply, limiting the number of instances.
- 2. Context switches on GPU cause overheads.
- 3. Fast synchronization between cameras becomes harder.



- 2. Context switches on GPU cause overheads.
- 3. Fast synchronization between cameras becomes harder.

### **Part I:**



### **Part II:**



Ming Yang - RTAS 2019



## **Parallel Execution** for CNN frameworks

## Multi-camera Composite Images to provide high throughput for multiple cameras.



## **Proposed Solutions**

### **Part I:**











Ming Yang - RTAS 2019

## **Parallel Execution** for CNN frameworks

## Multi-camera Composite Images to provide high throughput for multiple cameras.



### Part I:



Parallel **Execution** 







### Part I:



Parallel **Execution**   CNN models are graphs of layers.









### Part I:



Parallel Execution



- CNN models are graphs of layers.
- Processing of images can be independent, e.g., object detection.





### Part I:



Parallel Execution



Ming Yang - RTAS 2019

- CNN models are graphs of layers.
- Processing of images can be independent, e.g., object detection.

# shared CNN for multiple cameras.



We enable *parallel execution* for CNN frameworks and











- Generalize concept of layers into stages
- Communicate data between stages using **PGMRT** (a processing graph management tool)









- Generalize concept of layers into stages
- Communicate data between stages using **PGMRT** (a processing graph management tool)
- Share CNN among multiple cameras





## **Different Execution Methods**

### Part I:



Parallel Execution


## **Different Execution Methods**

SERIAL

#### Part I:



Parallel Execution



Ming Yang - RTAS 2019

## private CNN in one process

# **Different Execution Methods**

#### SERIAL

#### Part I:



Parallel **Execution** 

**Part II:** 

Multicamera Composite Images

#### Ming Yang - RTAS 2019

#### PIPELINE

shared CNN that has one thread per stage

## private CNN in one process

# **Different Execution Methods**

#### SERIAL

#### Part I:



Parallel Execution

**Part II:** 

Multicamera Composite Images

#### PIPELINE

PARALLEL

shared CNN that has *multiple* threads per stage

Ming Yang - RTAS 2019

## private CNN in one process

## shared CNN that has one thread per stage



#### Part I:



Parallel **Execution** 

#### PIPELINE

Part II:

Multicamera Composite Images

PARALLEL



#### Part I:



Parallel **Execution** 

#### PIPELINE

Part II:

Multicamera Composite Images

PARALLEL



#### Part I:



Parallel **Execution** 

#### **PIPELINE**

Part II:

Multicamera Composite Images

PARALLEL



#### Part I:



Parallel **Execution** 

#### PIPELINE

Part II:

Multicamera Composite Images

PARALLEL



#### Part I:



Parallel **Execution** 

#### PIPELINE

Part II:

Multicamera Composite Images

PARALLEL



#### Part I:



Parallel **Execution** 

#### PIPELINE

Part II:

Multicamera Composite Images

PARALLEL



#### Part I:



Parallel **Execution** 

#### PIPELINE

Part II:

Multicamera Composite Images

PARALLEL



#### Part I:



Parallel Execution

#### PIPELINE

Part II:

Multicamera Composite Images

PARALLEL



#### Part I:



Parallel Execution

#### PIPELINE

Part II:

Multicamera Composite Images

PARALLEL



#### Part I:



Parallel Execution

#### PIPELINE

Part II:

Multicamera Composite Images

PARALLEL



#### Part I:



Parallel Execution

#### PIPELINE

Part II:

Multicamera Composite Images

PARALLEL

Ming Yang - RTAS 2019



14



Parallel **Execution** 



- - SERIAL  $\bullet$
  - SERIAL x6
  - PIPELINE
  - PARALLEL
- Darknet

Ming Yang - RTAS 2019

# Evaluation

## We compared latency and throughput between

## • With CNN model **Tiny YOLOv2** on CNN framework

## • On hardware platform: NVIDIA Drive PX 2.

## **Evaluation (Hardware)**



#### Part I:



Parallel Execution



Multicamera Composite Images





Parallel **Execution** 



Multicamera Composite Images

Ming Yang - RTAS 2019

## **Evaluation (Hardware)**







Parallel Execution

Part II: Multicamera Composite Images

#### Ming Yang - RTAS 2019

# **Evaluation Results**

Camera frame rate (frames per second)







Part II: Multicamera Composite Images

#### Ming Yang - RTAS 2019

# **Evaluation Results**



# **Evaluation Results**



Part I:



Parallel Execution

Part II:Multi-<br/>camera<br/>Composite<br/>Images









Parallel Execution

Part II: Multicamera Composite Images

#### Ming Yang - RTAS 2019

# **Evaluation Results**

Camera frame rate (frames per second)







Parallel Execution

Part II: Multicamera Composite Images

Ming Yang - RTAS 2019

# **Evaluation Results**







Parallel Execution

Part II: Multicamera Composite Images

#### Ming Yang - RTAS 2019

# **Evaluation Results**

Camera frame rate (frames per second)







Parallel Execution

Part II: Multicamera Composite Images

#### Ming Yang - RTAS 2019

# **Evaluation Results**

Camera frame rate (frames per second)

# **Evaluation Results (cont.)**

#### Part I:



Parallel Execution

Part II:



Multicamera Composite Images

Ming Yang - RTAS 2019

#### SERIAL

## SERIAL X6

## PIPELINE

(Single thread per stage

PARALLEL

(10 threads per stage)

	CPUs (%)	Memory (MB)
	92	774
	536	4,644
e)	219	1,132
	239	1,136



# Evaluation Results (cont.)CPUs (%)Memory (MB)

#### Part I:



Parallel Execution

Part II:



Multicamera Composite Images

Ming Yang - RTAS 2019

#### SERIAL

## SERIAL X6

## PIPELINE

(Single thread per stage)

PARALLEL

(10 threads per stage)





# **Evaluation Results (cont.)** CPUs (%) Memory (MB)

92

#### Part I:



Parallel Execution SERIAL

#### PIPELINE

(Single thread per stage)

### PARALLEL

(10 threads per stage)

**Part II:** Multicamera Composite Images

Ming Yang - RTAS 2019







774

# **Evaluation Results (cont.)**

#### Part I:



Parallel Execution

## Enabling intra-stage parallelism takes slight overheads.



Multicamera Composite Images

#### PIPELINE

(Single thread per stage)

PARALLEL

(10 threads per stage)

Ming Yang - RTAS 2019

CPUs (%)

Memory (MB)







Parallel **Execution** 



## Part I: Parallel Execution

- to 71 FPS
- No accuracy loss

Ming Yang - RTAS 2019

Pipeline and Parallel improve throughput from 28 FPS

#### With acceptable overheads and



#### Part I:



## **Execution**

#### Part II:



Multicamera Composite Images

#### Ming Yang - RTAS 2019



Icons made by **Butterflytronics** from **Flaticon** is licensed by <u>CC 3.0 BY</u>



#### Part I:



## **Execution**

#### Part II:



Multicamera Composite Images

#### Ming Yang - RTAS 2019



Icons made by **Butterflytronics** from **Flaticon** is licensed by <u>CC 3.0 BY</u>



#### Part I:



## **Execution**

#### Part II:



Multicamera Composite Images

#### Ming Yang - RTAS 2019



Icons made by **Butterflytronics** from **Flaticon** is licensed by <u>CC 3.0 BY</u>



#### Part I:



Parallel Execution

#### Part II:



Multicamera Composite Images



#### Part I:



**Parallel Execution** 

#### Part II:



Multicamera Composite Images

#### **Virtual Camera**



#### Part I:



#### Part II: Multicamera Composite Images



#### **Virtual Camera**





Ming Yang - RTAS 2019

Images are from PASCAL dataset



#### Part I:



#### Part II: Multicamera Composite Images



#### **Virtual Camera**





Ming Yang - RTAS 2019

## Shared CNN

Images are from PASCAL dataset


## Multi-camera Composite Images



### Part I:



**Parallel Execution** 

### Part II:



Multicamera Composite Images



Ming Yang - RTAS 2019



## Multi-camera Composite Images



### Part I:



Parallel **Execution** 

### Part II:



Multicamera Composite Images



Ming Yang - RTAS 2019







### Part I:



### Part II:



Multicamera Composite Images

- between
  - Full-size images

Ming Yang - RTAS 2019

## Evaluation

### We compared latency, throughput and accuracy

• Four-camera composite images



Ming Yang - RTAS 2019

### Part I:



### Part II:



Multicamera Composite Images

Ming Yang - RTAS 2019



### Part I:



## **Part II:**



Multicamera Composite Images

Classes: bicycle, bus, car, motorbike, train, bird, person, cat, cow, dog, horse, sheep

### Ming Yang - RTAS 2019

Table 1: accuracy (mAPs) of object classes relevant to autonomous driving





### **Original YOLO**

**Part II:** Multicamera Composite Images

Parallel

Execution

Classes: bicycle, bus, car, motorbike, train, bird, person, cat, cow, dog, horse, sheep

### Ming Yang - RTAS 2019

Part I:

Table 1: accuracy (mAPs) of object classes relevant to autonomous driving





### Part I:



Parallel Execution

### **Part II:**



Multicamera Composite Images

**Original YOLO** 

### **Retrained YOLO**

Classes: bicycle, bus, car, motorbike, train, bird, person, cat, cow, dog, horse, sheep

Ming Yang - RTAS 2019

Table 1: accuracy (mAPs) of object classes relevant to autonomous driving





# Conclusions

• We presented an industrial study that addresses the challenge of supporting multiple cameras.



### **Parallel execution**

**Multi-camera composite image** 

- Evaluation results showed  $\bullet$ 
  - Significant throughput improvements



- No accuracy loss with parallel execution
- ~7.4% accuracy drop with multi-camera composite image (but 4-fold throughput improvement!)



### Table 1: mAPs of object classes relevant to autonomous driving

	Full-size Test	Composite Test
Original YOLO	63.66	44.91
<b>Retrained YOLO</b>	66.20	56.21



# Conclusions

٠ We presented an industrial study that addresses the challenge of supporting multiple cameras.



### **Parallel execution**

**Multi-camera composite image** 

- **Evaluation results** showed  $\bullet$ 
  - Significant throughput improvements



- No accuracy loss with parallel execution
- ~7.4% accuracy drop with multi-camera composite image (but 4-fold throughput improvement!)

Other considerations in the paper:

- Configurable stages
- Multi-GPU execution



Table 1: mAPs of object classes relevant to autonomous driving

	Full-size Test	Composite Test
Original YOLO	63.66	44.91
<b>Retrained YOLO</b>	66.20	56.21





## Future Work

 Dynamically apply composite-image technique with criticality change

• Finer granularity of stages

• Dynamically share CNN among multiple models

Ming Yang - RTAS 2019







# Thank you!

## Re-thinking CNN Frameworks for Time-Sensitive Autonomous-Driving Applications: Addressing an Industrial Challenge

Ming Yang<sup>1</sup>, Shige Wang<sup>2</sup>, Joshua Bakita<sup>1</sup>, Thanh Vu<sup>1</sup>, F. Donelson Smith<sup>1</sup>, James H. Anderson<sup>1</sup>, and Jan-Michael Frahm<sup>1</sup>

<sup>1</sup>The University of North Carolina at Chapel Hill <sup>2</sup>General Motors Research



THE UNIVERSITY of NORTH CAROLINA at CHAPEL HILL

Ming Yang - RTAS 2019

