

**GPU** TECHNOLOGY  
CONFERENCE

# POWER EFFICIENT VISUAL COMPUTING ON MOBILE PLATFORMS

BRANT ZHAO, NVIDIA

MAX LV, NVIDIA

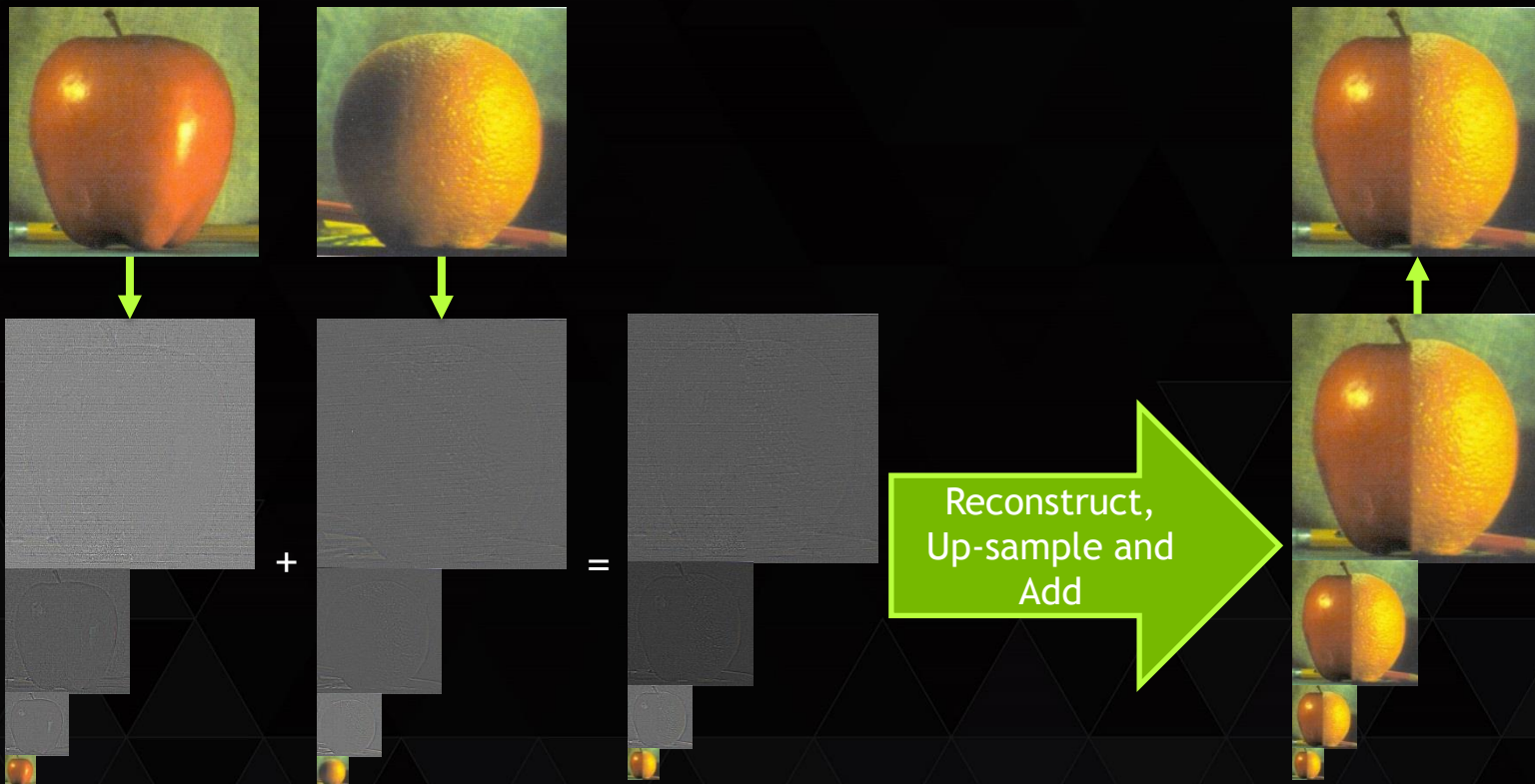
- Performance
- Energy Efficiency



# Power Efficient GPU Programming - Case Studies & Findings

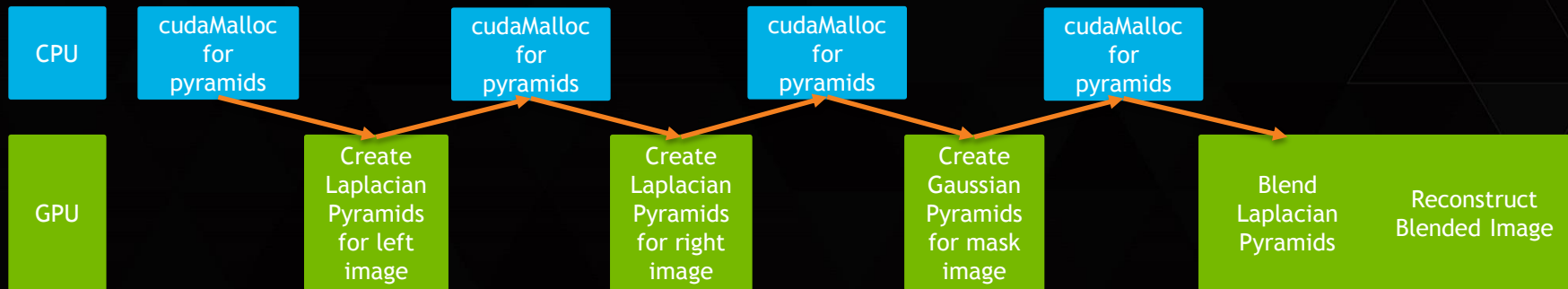
# Case study #1: Image Pyramid Blending

# Image Pyramid Blending



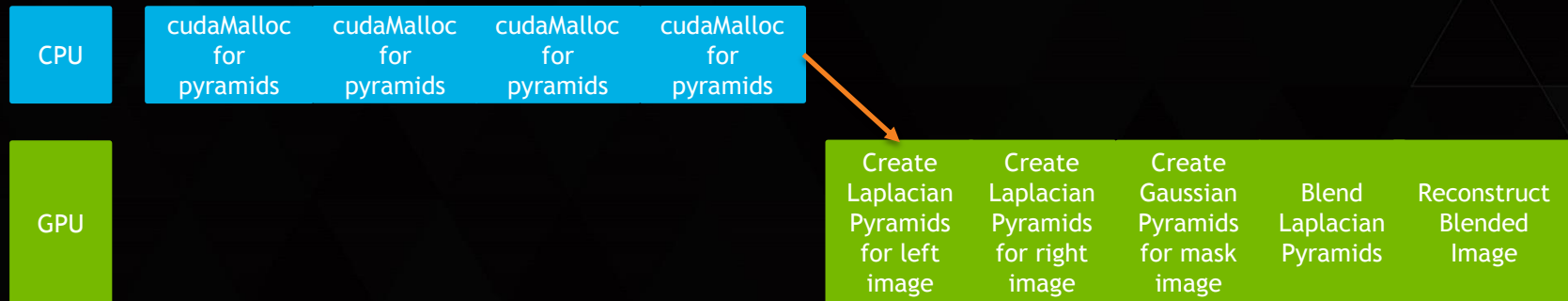
# Image Pyramid Blending

## - A naïve CUDA implementation



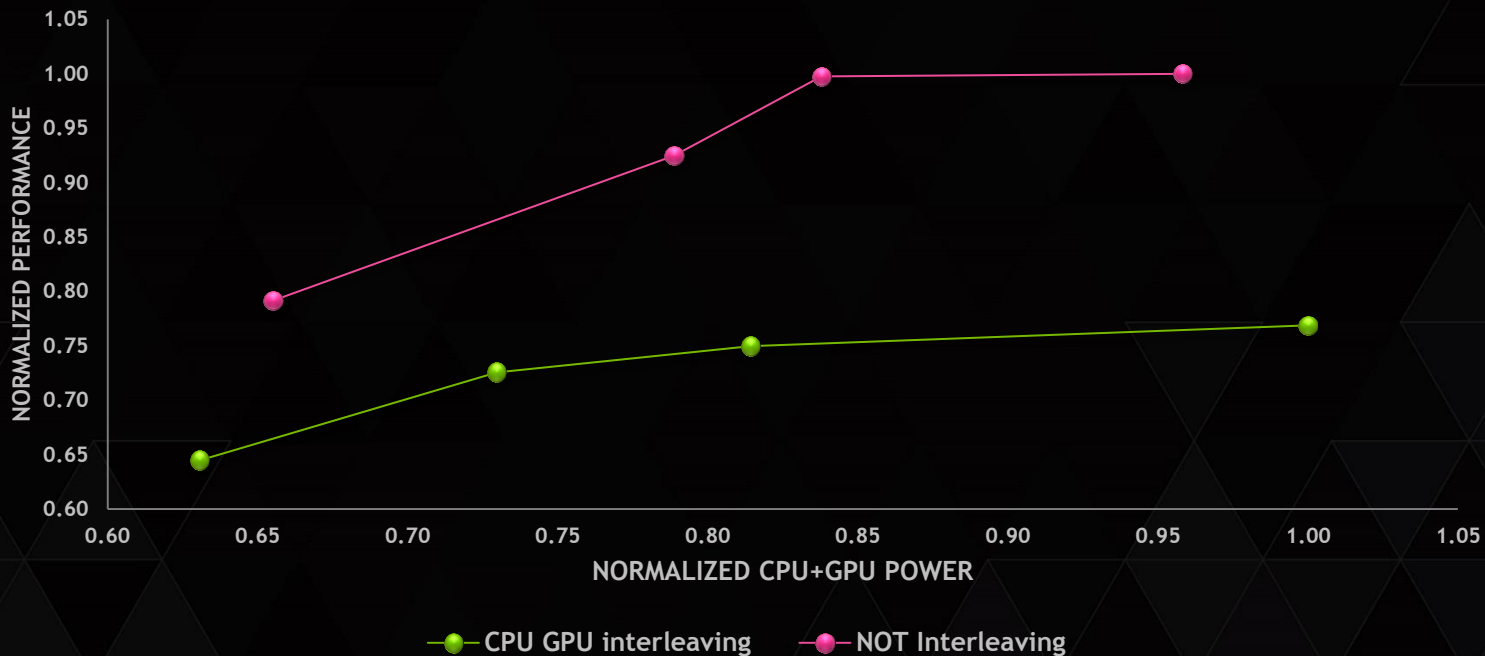
# Image Pyramid Blending

- Power optimized: Avoid CPU $\leftrightarrow$ GPU interleaving



# Image Pyramid Blending

## - Perf/Watt comparison





# Case study #2: 2D Convolution

# 2D Convolution

0.25	0	0
0	0	0
0	0	0.75

+

1	2	1	2	
0	0	2	3	
2	0	1	10	

=

	1			

# 2D Convolution

0.25	0	0
0	0	0
0	0	0.75

+

1	2	1	2	
0	0	2	3	
2	0	1	10	

=

	1	8		

# 2D Convolution

- 3x3 2D convolution with FP16

0.25	0	0
0	0	0
0	0	0.75

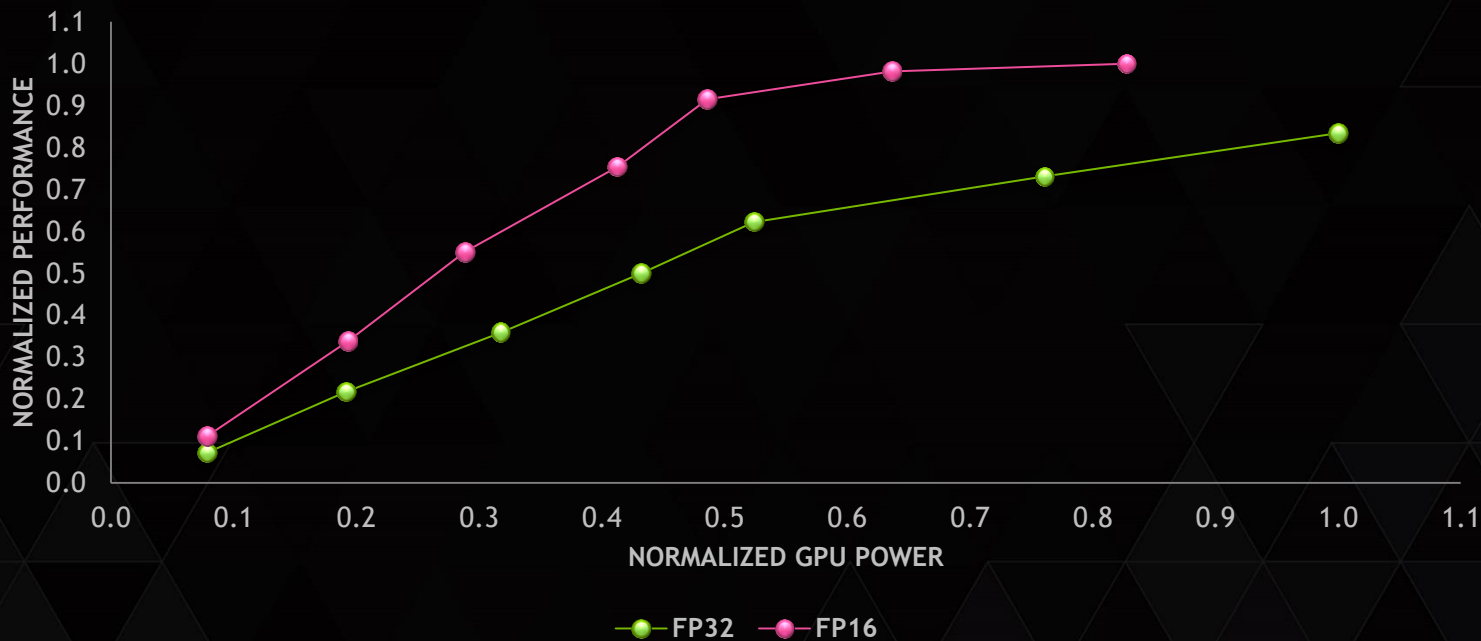
1 pack 0	2 pack 1	1 pack 2	2
0 pack 3	0 pack 4	2 pack 5	3
2 pack 6	0 pack 7	1 pack 8	10

1	8
---	---

- Basic operations for 2 output pixels
  - 9 packed FP16 MAD

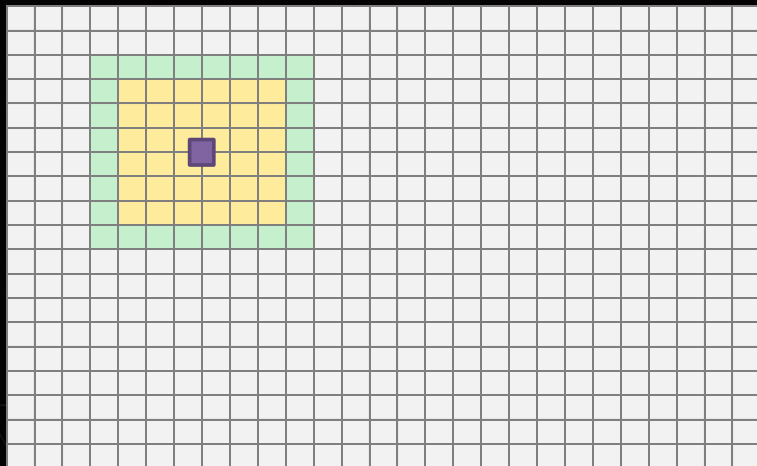
# 2D Convolution

## - Perf/Watt comparison

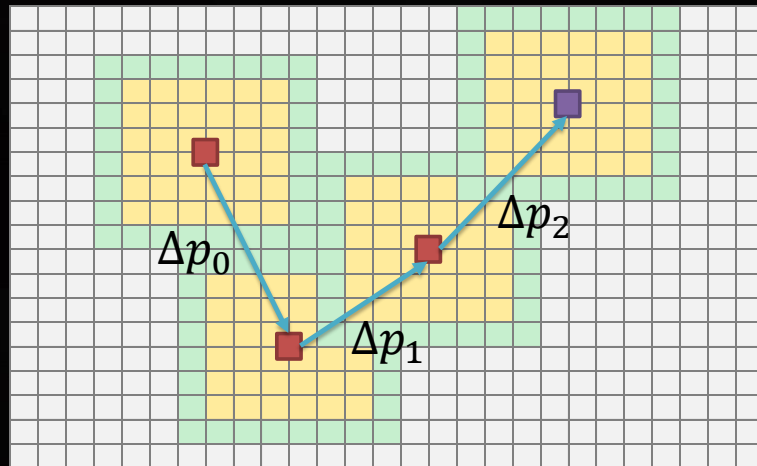


# Case study #3: Sparse Lucas- Kanade Optical Flow (SparseLK)

# SparseLK



First Frame  $I$



Second Frame  $I_{next}$

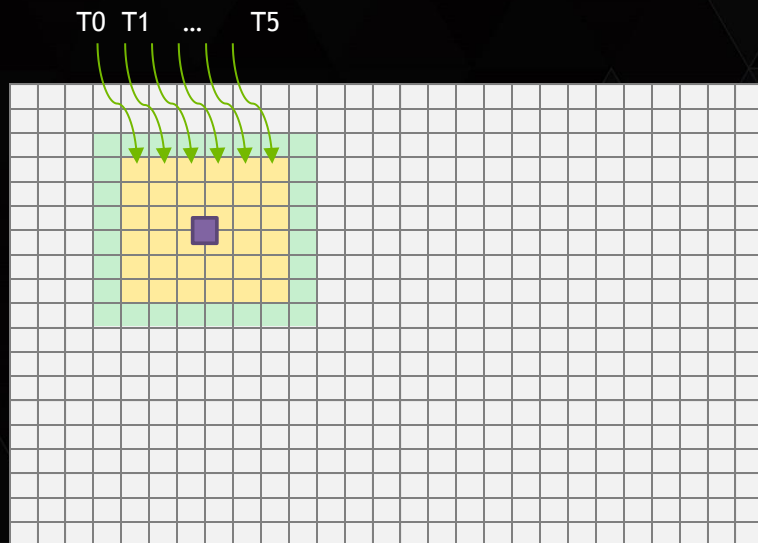
$$\Delta p = \left[ \sum_{x \in \Omega} \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix} \right]^{-1} \sum_{x \in \Omega} (I(x) - I_{next}(x + \Delta p_{prev})) \begin{bmatrix} I_x \\ I_y \end{bmatrix}$$

# SparseLK

## - Solution#1

$$\Delta p = \left[ \sum_{x \in \Omega} \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix} \right]^{-1} \sum_{x \in \Omega} (I(x) - I_{next}(x + \Delta p_{prev})) \begin{bmatrix} I_x \\ I_y \end{bmatrix}$$

- Multiple threads for a feature point
- Share data via shared memory or shuffle
- Reduction needed to get final results
- High thread level parallelism(TLP) but more instructions needed



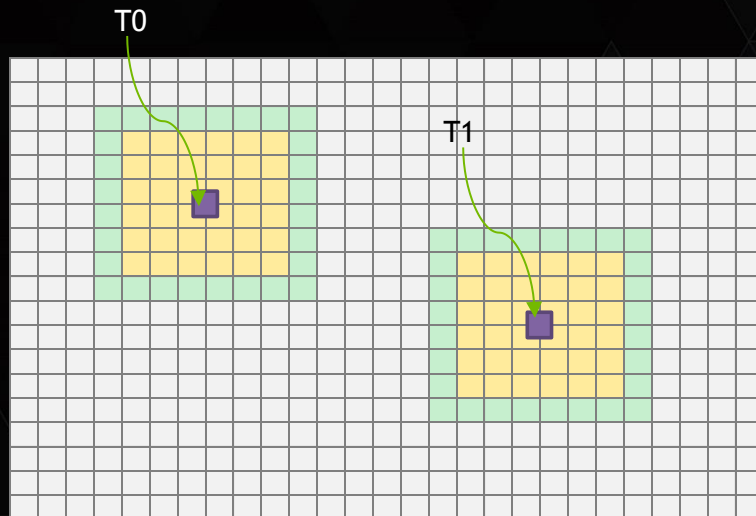


# SparseLK

## - Solution#2

$$\Delta p = \left[ \sum_{x \in \Omega} \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix} \right]^{-1} \sum_{x \in \Omega} (I(x) - I_{next}(x + \Delta p_{prev})) \begin{bmatrix} I_x \\ I_y \end{bmatrix}$$

- Each thread handles a feature point
- No need to shuffle data
- No need to do reduction
- Need more registers to hold data
- High instruction level parallelism(ILP) but low occupancy



# SparseLK

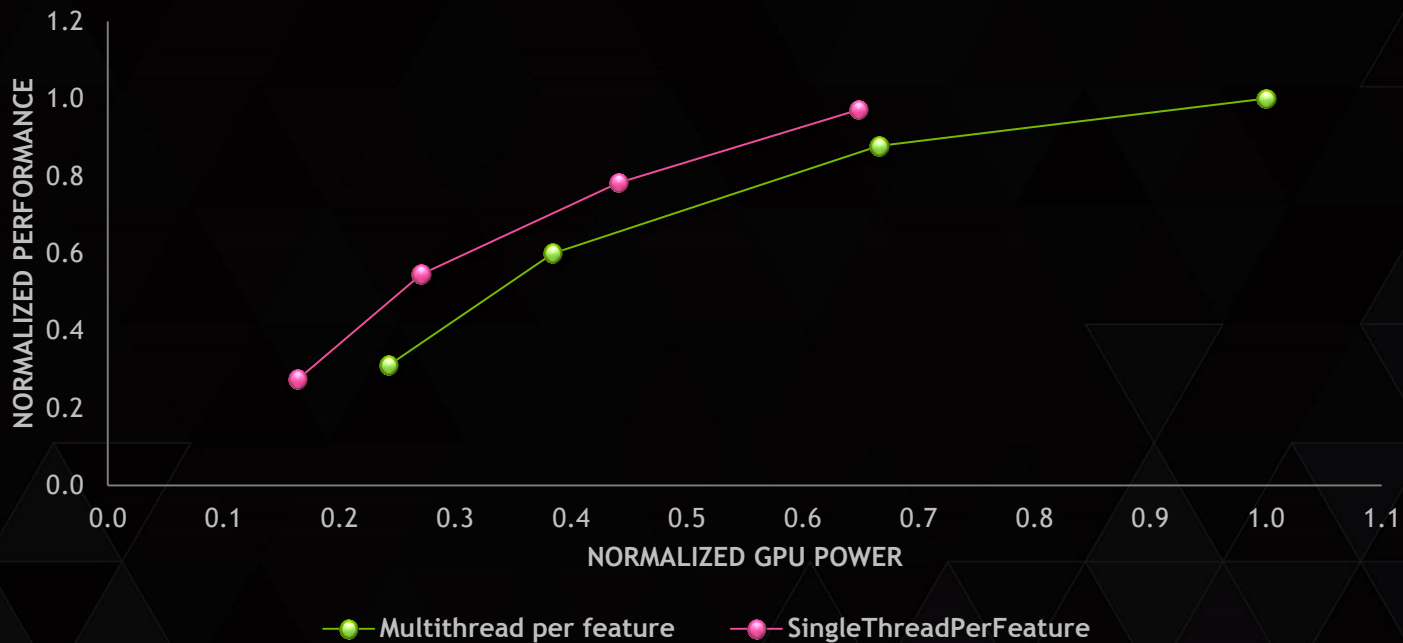
- Instruction# and Perf/Watt

$$\frac{\text{Perf}}{\text{Watt}} = \frac{\frac{\text{Workload}}{\text{Sec}}}{\frac{\text{Energy}}{\text{Sec}}} = \frac{\text{Workload}}{\text{Energy}} \quad \begin{matrix} \uparrow \\ \downarrow \end{matrix}$$

$$\begin{aligned} \text{Energy} = & \text{EnergyPerInst}_{\text{shuffle}} * \text{Instruction}_{\text{shuffle}} \\ & + \text{EnergyPerInst}_{\text{reduction}} * \text{Instruction}_{\text{reduction}} \\ & + \text{EnergyPerInst}_{\text{other}} * \text{Instruction}_{\text{other}} \\ & + \text{Power}_{\text{wasted}} * \text{Time} \end{aligned} \quad \downarrow$$

# SparseLK

## - Perf/Watt comparison



# Summary

- Analyze the whole pipeline at the system level
- Use energy efficient features on the target platform
- Balance between TLP and ILP

**GPU** TECHNOLOGY  
CONFERENCE

**THANK YOU**

**brantz@nvidia.com**

JOIN THE CONVERSATION

**#GTC15**   