



# **AMALTHEA-based GPU Response Time Analysis for NVidia's Jetson TX2**

Daniel Paredes



## Agenda:

1. Motivation
2. Nvidia GPU Model
3. AMALTHEA Model
4. Response Time Analysis for Jetson TX2's GPU
5. Results
- 6. Conclusions**

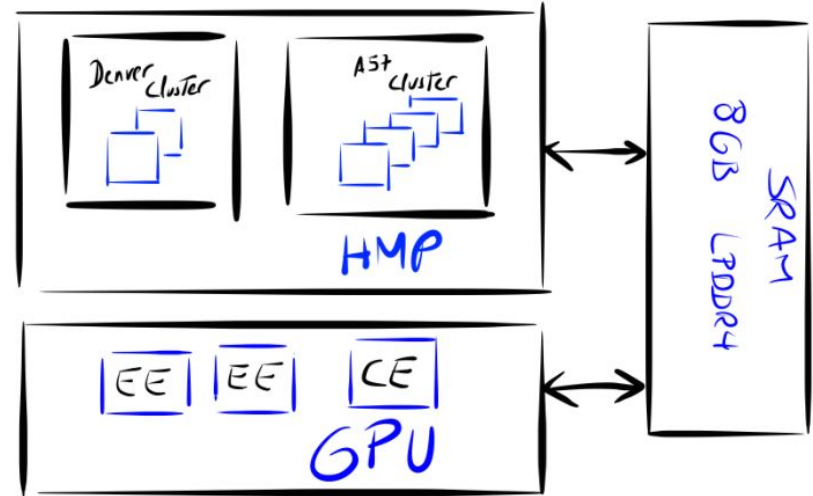
---

# Motivation

## Motivation

WATERS Industrial Challenge:

- Heterogeneous systems
- Centralized end-to-end architectures.
- Chosen platform: NVIDIA Jetson™ TX2



---

# NVidia GPU Model



# Nvidia GPU Model

Software Model:

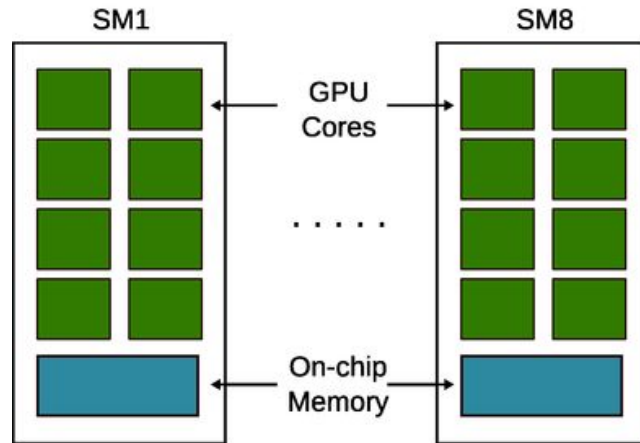
- Threads
- Blocks
- Grids

Hardware Model:

- Streaming Multiprocessors
- Memory

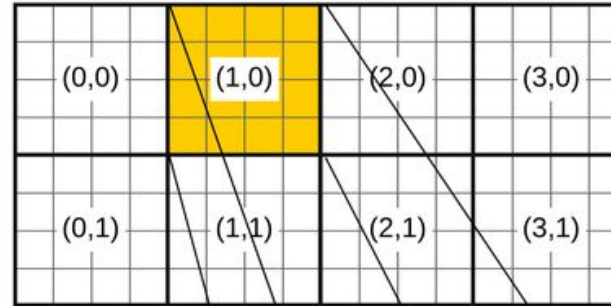
# Nvidia GPU Model

Streaming Multiprocessors



(a)

Grid of blocks



Threads within a block

(0,0)	(0,1)	(0,2)	(0,3)
(1,0)	(1,1)	(1,2)	(1,3)
(2,0)	(2,1)	(2,2)	(2,3)
(3,0)	(3,1)	(3,2)	(3,3)

scheduled to any  
of the available GPU cores

(b)

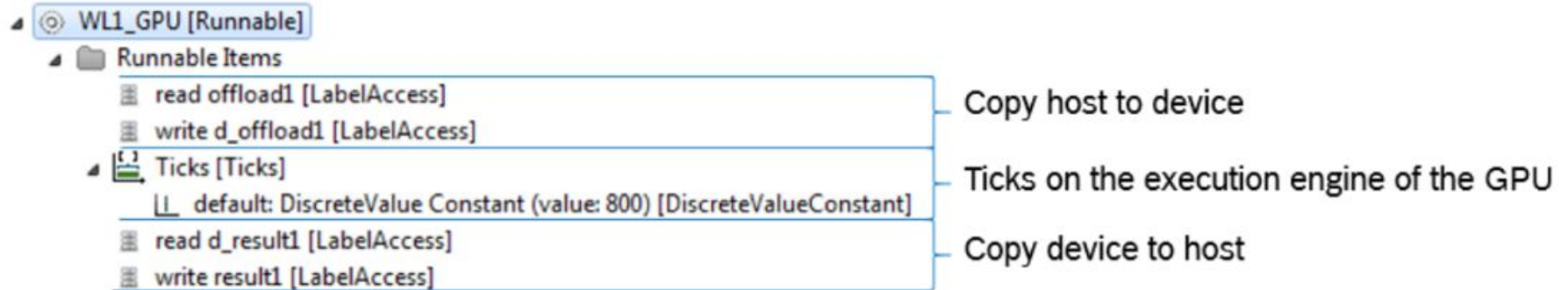
[1] **Computing Efficiently Spectral-Spatial Classification of Hyperspectral Images on Commodity GPUs** - Pablo Quesada-Barriuso

---

# AMALTHEA Model



# AMALTHEA Model



---

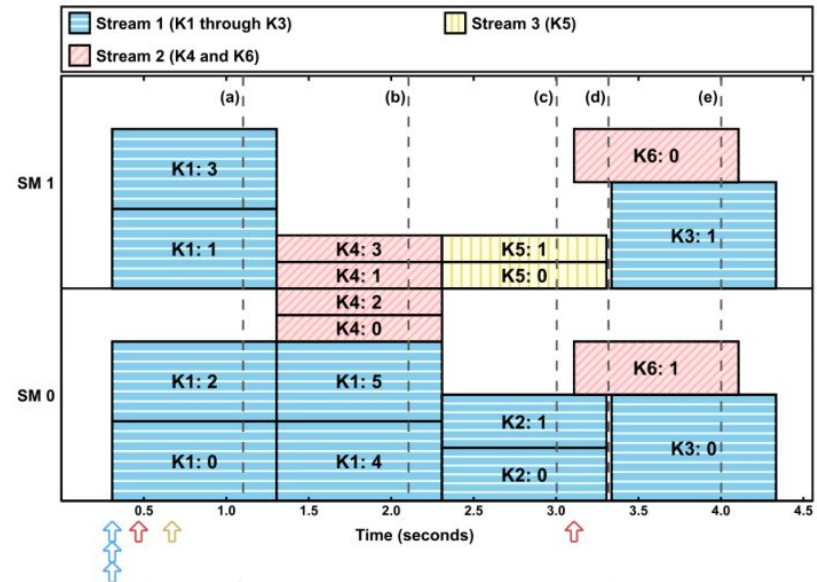
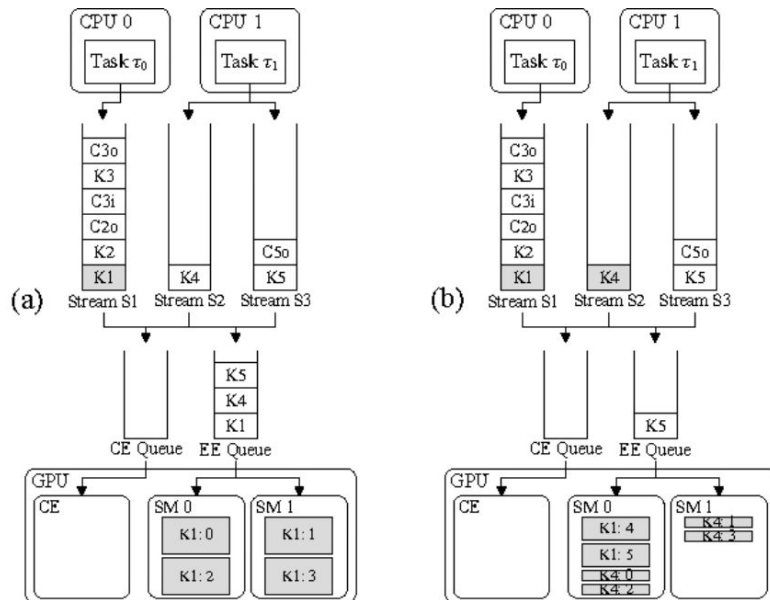
# Response Time Analysis



## Jetson TX2's GPU

- 2 Streaming Multiprocessors
- FIFO Queues
- Max. threads per Block: 1024
- Max. threads per SM: 2048

# GPU Scheduler



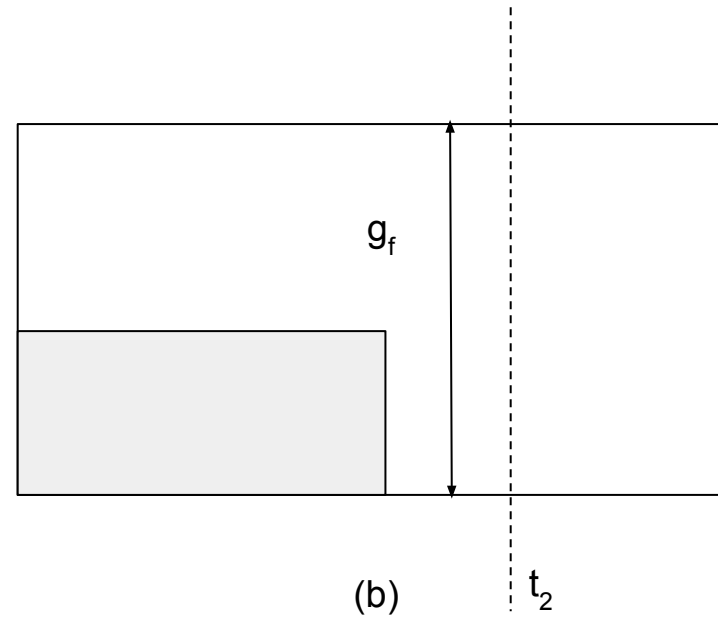
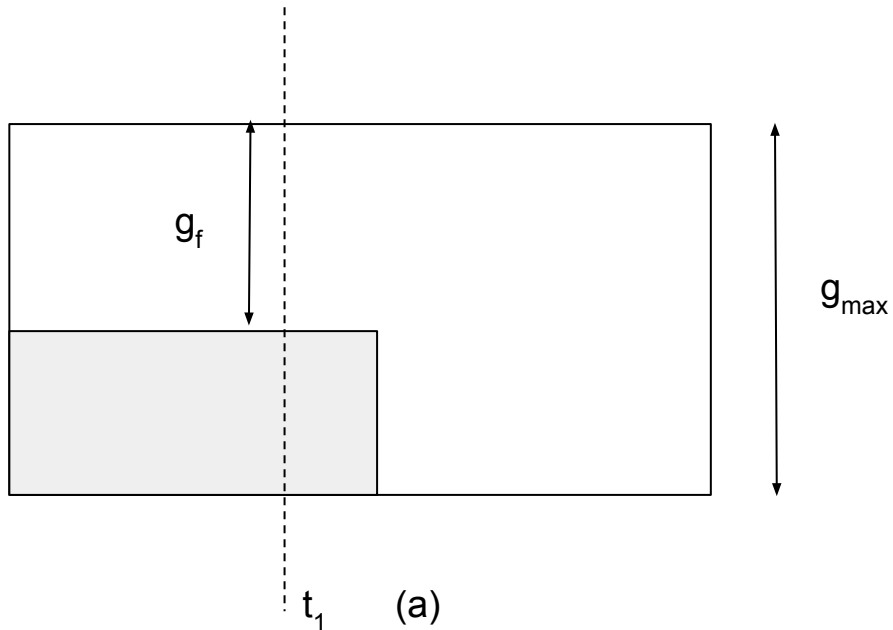
[2] GPU Scheduling on the NVIDIA TX2: Hidden Details Revealed - Tanya Amert et al.



## Assumptions

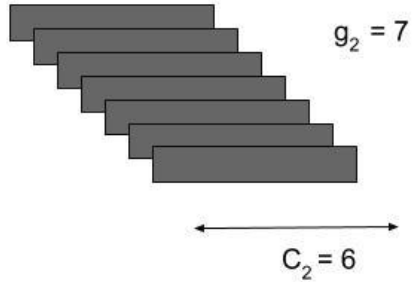
- All blocks have the same amount of threads.
  - Best practices on CUDA programming
- One big Streaming Multiprocessor.
  - $\text{\#Threads per SM} / \text{\#threads per block} = N$

# Response Time Calculation



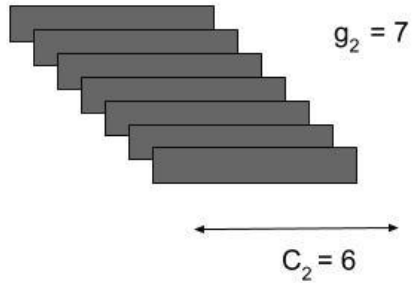


# Response Time Calculation

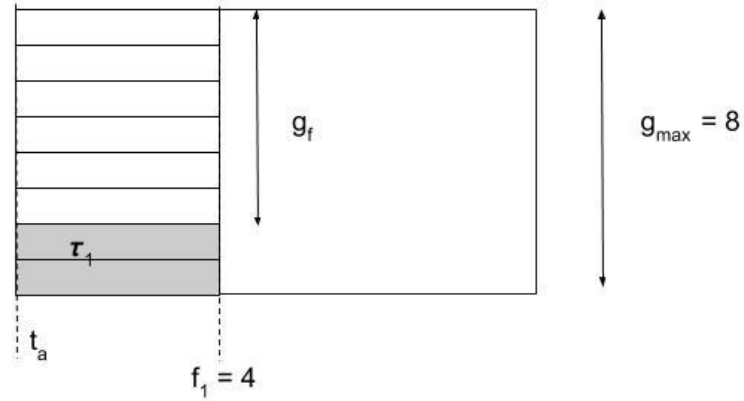


(a)

# Response Time Calculation



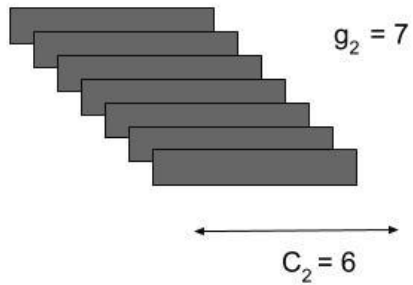
(a)



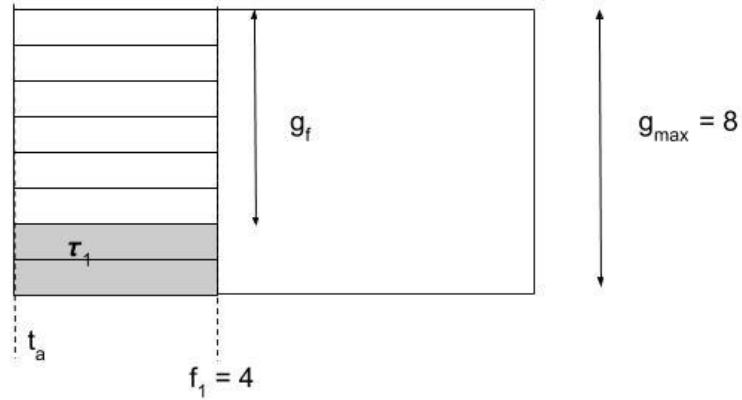
(b)



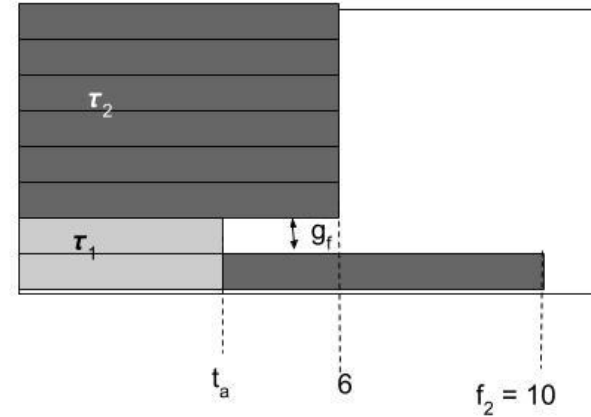
# Response Time Calculation



(a)



(b)



(c)



# Algorithm

$\tau$ : set of kernels

$f_i$ : completion time of kernel  $i$

$C_i$ : thread execution time

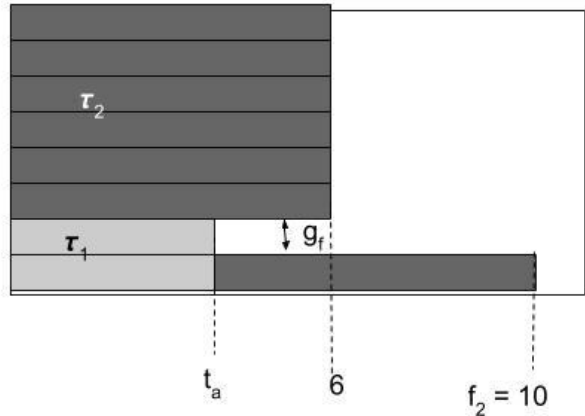
Input :  $\tau$

Output:  $f_1, \dots, f_n$

Initialization:  $t_a = 0, g_f = g_{max}, i = 1$

```
while  $i \leq n$  do
    if  $g_f \geq g_i$  then
         $f_i = t_a + C_i$ ;
        Update  $g_f$  and  $t_a$ ;
         $i++$  ; // Next kernel
    else
         $g_i = g_i - g_f$ ;
        Update  $g_f$  and  $t_a$ ;
    end
end
```

# Algorithm



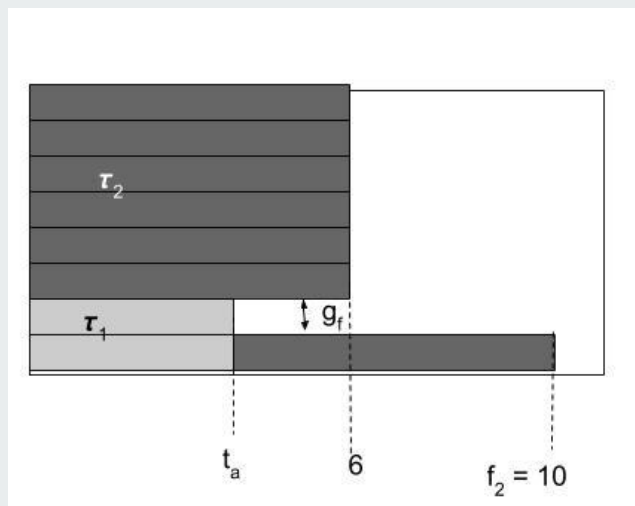
Input :  $\tau$

Output:  $f_1, \dots, f_n$

Initialization:  $t_a = 0$ ,  $g_f = g_{max}$ ,  $i = 1$

```
while  $i \leq n$  do
  if  $g_f \geq g_i$  then
     $f_i = t_a + C_i$ ;
    Update  $g_f$  and  $t_a$ ;
     $i++$  ; // Next kernel
  else
     $g_i = g_i - g_f$ ;
    Update  $g_f$  and  $t_a$ ;
  end
end
```

# Algorithm



Input :  $\tau$

Output:  $f_1, \dots, f_n$

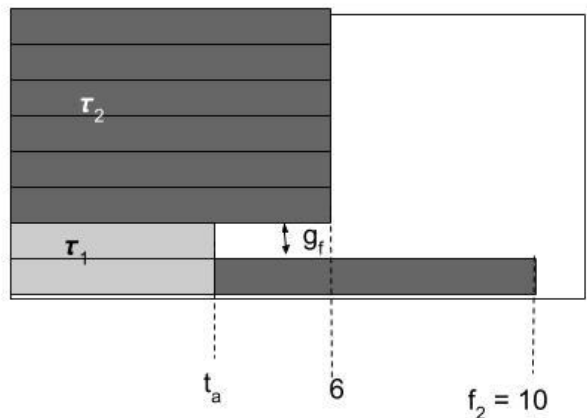
Initialization:  $t_a = 0, g_f = g_{max}, i = 1, h = \{\}$

```

while  $i \leq n$  do
    if  $g_f \geq g_i$  then
         $f_i = t_a + C_i$ ;
         $h = \{h; (f_i, g_i)\}$ ;
         $t_a = t_a$ ;
         $g_f = g_f - g_i$ ;
         $i++$  ; // Next kernel
    else
         $g_i = g_i - g_f$ ;
         $h = \{h; (t_a + C_i, g_f)\}$ ;
         $[t_a, \text{index}] = \min(h[:, 1])$ ;
         $g_f = h[\text{index}, 2]$ ;
        Update  $h$ ;
    end
end
end

```

# Algorithm



Input :  $\tau$

Output:  $f_1, \dots, f_n$

Initialization:  $t_a = 0$ ,  $g_f = g_{max}$ ,  $i = 1$ ,  $h = \{\}$

while  $i \leq n$  do

    if  $g_f \geq g_i$  then

$f_i = t_a + C_i$ ;

$h = \{h; (f_i, g_i)\}$ ;

$t_a = t_a$ ;

$g_f = g_f - g_i$ ;

$i++$  ; // Next kernel

    else

$g_i = g_i - g_f$ ;

$h = \{h; (t_a + C_i, g_f)\}$ ;

$[t_a, \text{index}] = \min(h[:, 1])$ ;

$g_f = h[\text{index}, 2]$ ;

        Update  $h$ ;

    end

end

# AMALTHEA implementation

```
while  $i \leq n$  do
  if  $g_f \geq g_i$  then
     $f_i = t_a + C_i$ ;
     $h = \{h; (f_i, g_i)\}$ ;
     $t_a = t_a$ ;
     $g_f = g_f - g_i$ ;
     $i++$  ; // Next kernel
  else
     $g_i = g_i - g_f$ ;
     $h = \{h; (t_a + C_i, g_f)\}$ ;
     $[t_a, \text{index}] = \min(h[:, 1])$ ;
     $g_f = h[\text{index}, 2]$ ;
    Update  $h$ ;
  end
end
```

```
// Main loop
while ( current_kernel < rList.size() ) {
  if (  $g_f \geq g_i.get(\text{current\_kernel})$  ) {
     $f.add(\text{current\_kernel}, t_a + c_i.get(\text{current\_kernel}) )$  ;

     $h = \text{updateH}(h, f.get(\text{current\_kernel}), g_i.get(\text{current\_kernel}) )$  ;

     $g_f = g_f - g_i.get(\text{current\_kernel})$ ;
     $\text{current\_kernel}++$ ;
  }
  else {
     $g_i.set(\text{current\_kernel}, g_i.get(\text{current\_kernel}) - g_f)$ ;

     $h = \text{updateH}(h, t_a + c_i.get(\text{current\_kernel}), g_f)$ ;
     $\text{minimumRegisteredTicks} = \text{findIndexOffMinValue}(h)$ ;

     $g_f = h.get(\text{minimumRegisteredTicks})$ ;
     $t_a = \text{minimumRegisteredTicks}$ ;

     $h.remove(\text{minimumRegisteredTicks})$ ;
  }
}
```

# AMALTHEA implementation

```
while  $i \leq n$  do
  if  $g_f \geq g_i$  then
     $f_i = t_a + C_i$ ;
     $h = \{h; (f_i, g_i)\}$ ;
     $t_a = t_a$ ;
     $g_f = g_f - g_i$ ;
     $i++$  ; // Next kernel
  else
     $g_i = g_i - g_f$ ;
     $h = \{h; (t_a + C_i, g_f)\}$ ;
     $[t_a, \text{index}] = \min(h[:, 1])$ ;
     $g_f = h[\text{index}, 2]$ ;
    Update  $h$ ;
  end
end
```

```
// Main loop
while ( current_kernel < rList.size() ) {
  if (  $g_f \geq g_i.get(\text{current\_kernel})$  ) {
     $f.add(\text{current\_kernel}, t_a + c_i.get(\text{current\_kernel}) )$  ;

     $h = \text{updateH}(h, f.get(\text{current\_kernel}), g_i.get(\text{current\_kernel}) )$  ;

     $g_f = g_f - g_i.get(\text{current\_kernel})$ ;
     $\text{current\_kernel}++$ ;
  }
  else {
     $g_i.set(\text{current\_kernel}, g_i.get(\text{current\_kernel}) - g_f)$ ;

     $h = \text{updateH}(h, t_a + c_i.get(\text{current\_kernel}), g_f)$ ;
     $\text{minimumRegisteredTicks} = \text{findIndexOffMinValue}(h)$ ;

     $g_f = h.get(\text{minimumRegisteredTicks})$ ;
     $t_a = \text{minimumRegisteredTicks}$ ;

     $h.remove(\text{minimumRegisteredTicks})$ ;
  }
}
```

# AMALTHEA implementation

```
while  $i \leq n$  do
  if  $g_f \geq g_i$  then
     $f_i = t_a + C_i$ ;
     $h = \{h; (f_i, g_i)\}$ ;
     $t_a = t_a$ ;
     $g_f = g_f - g_i$ ;
     $i++$  ; // Next kernel
  else
     $g_i = g_i - g_f$ ;
     $h = \{h; (t_a + C_i, g_f)\}$ ;
     $[t_a, \text{index}] = \min(h[:, 1])$ ;
     $g_f = h[\text{index}, 2]$ ;
    Update  $h$ ;
  end
end
```

```
// Main loop
while ( current_kernel < rList.size() ) {
  if (  $g_f \geq g_i.get(\text{current\_kernel})$  ) {
     $f.add(\text{current\_kernel}, t_a + c_i.get(\text{current\_kernel}) )$  ;

     $h = \text{updateH}(h, f.get(\text{current\_kernel}), g_i.get(\text{current\_kernel}) )$  ;

     $g_f = g_f - g_i.get(\text{current\_kernel})$ ;
     $\text{current\_kernel}++$ ;
  }
  else {
     $g_i.set(\text{current\_kernel}, g_i.get(\text{current\_kernel}) - g_f)$ ;

     $h = \text{updateH}(h, t_a + c_i.get(\text{current\_kernel}), g_f)$ ;
     $\text{minimumRegisteredTicks} = \text{findIndexOffMinValue}(h)$ ;

     $g_f = h.get(\text{minimumRegisteredTicks})$ ;
     $t_a = \text{minimumRegisteredTicks}$ ;

     $h.remove(\text{minimumRegisteredTicks})$ ;
  }
}
```



# AMALTHEA implementation

```
while  $i \leq n$  do
  if  $g_f \geq g_i$  then
     $f_i = t_a + C_i$ ;
     $h = \{h; (f_i, g_i)\}$ ;
     $t_a = t_a$ ;
     $g_f = g_f - g_i$ ;
     $i++$  ; // Next kernel
  else
     $g_i = g_i - g_f$ ;
     $h = \{h; (t_a + C_i, g_f)\}$ ;
     $[t_a, \text{index}] = \min(h[:, 1])$ ;
     $g_f = h[\text{index}, 2]$ ;
    Update  $h$ ;
  end
end
```

```
// Main loop
while ( current_kernel < rList.size() ) {
  if (  $g_f \geq g_i.get(\text{current\_kernel})$  ) {
     $f.add(\text{current\_kernel}, t_a + c_i.get(\text{current\_kernel}) )$  ;

     $h = \text{updateH}(h, f.get(\text{current\_kernel}), g_i.get(\text{current\_kernel}) )$  ;

     $g_f = g_f - g_i.get(\text{current\_kernel})$ ;
     $\text{current\_kernel}++$ ;
  }
  else {
     $g_i.set(\text{current\_kernel}, g_i.get(\text{current\_kernel}) - g_f)$ ;

     $h = \text{updateH}(h, t_a + c_i.get(\text{current\_kernel}), g_f)$ ;
     $\text{minimumRegisteredTicks} = \text{findIndexOffMinValue}(h)$ ;

     $g_f = h.get(\text{minimumRegisteredTicks})$ ;
     $t_a = \text{minimumRegisteredTicks}$ ;

     $h.remove(\text{minimumRegisteredTicks})$ ;
  }
}
```

# AMALTHEA implementation

```
while  $i \leq n$  do
  if  $g_f \geq g_i$  then
     $f_i = t_a + C_i$ ;
     $h = \{h; (f_i, g_i)\}$ ;
     $t_a = t_a$ ;
     $g_f = g_f - g_i$ ;
     $i++$  ; // Next kernel
  else
     $g_i = g_i - g_f$ ;
     $h = \{h; (t_a + C_i, g_f)\}$ ;
     $[t_a, \text{index}] = \min(h[:, 1])$ ;
     $g_f = h[\text{index}, 2]$ ;
    Update  $h$ ;
  end
end
```

```
// Main loop
while ( current_kernel < rList.size() ) {
  if (  $g_f \geq g_i.get(\text{current\_kernel})$  ) {
     $f.add(\text{current\_kernel}, t_a + c_i.get(\text{current\_kernel}) )$  ;

     $h = \text{updateH}(h, f.get(\text{current\_kernel}), g_i.get(\text{current\_kernel}) )$  ;

     $g_f = g_f - g_i.get(\text{current\_kernel})$ ;
     $\text{current\_kernel}++$ ;
  }
  else {
     $g_i.set(\text{current\_kernel}, g_i.get(\text{current\_kernel}) - g_f)$ ;

     $h = \text{updateH}(h, t_a + c_i.get(\text{current\_kernel}), g_f)$ ;
     $\text{minimumRegisteredTicks} = \text{findIndexOffMinValue}(h)$ ;

     $g_f = h.get(\text{minimumRegisteredTicks})$ ;
     $t_a = \text{minimumRegisteredTicks}$ ;

     $h.remove(\text{minimumRegisteredTicks})$ ;
  }
}
```

---

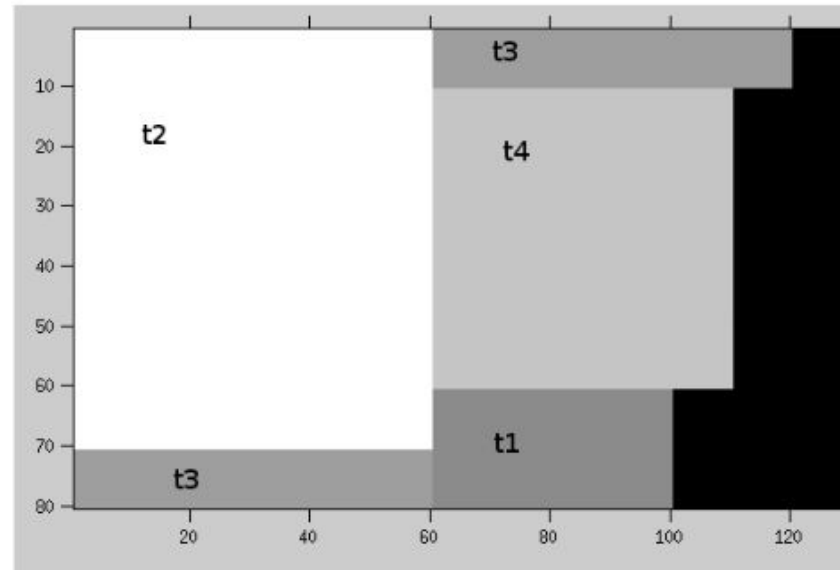
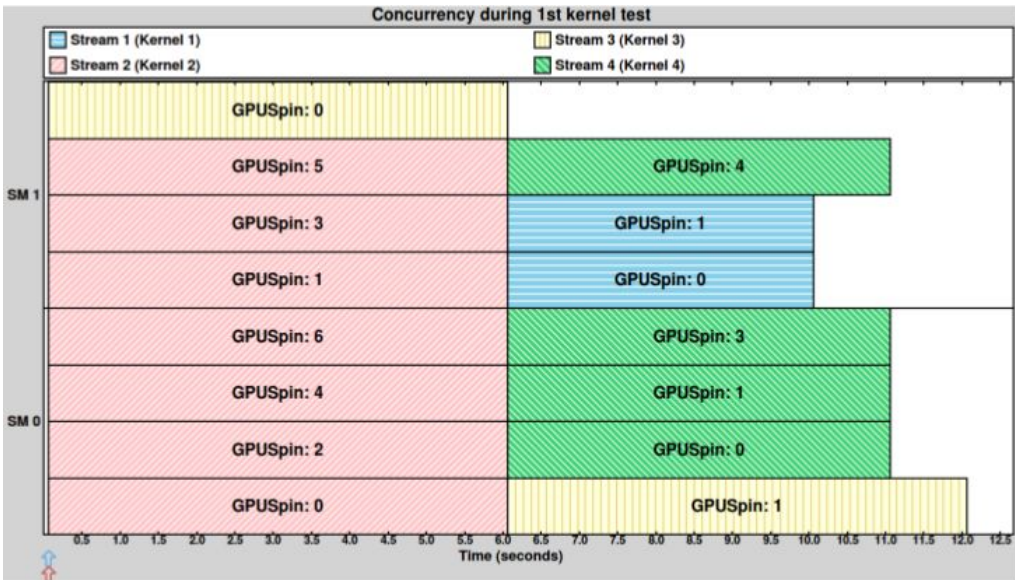
# Results



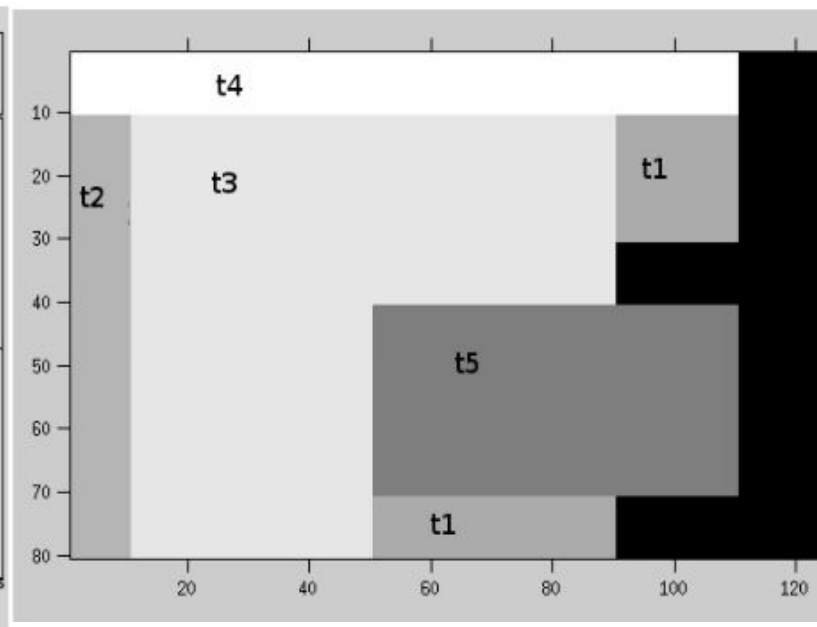
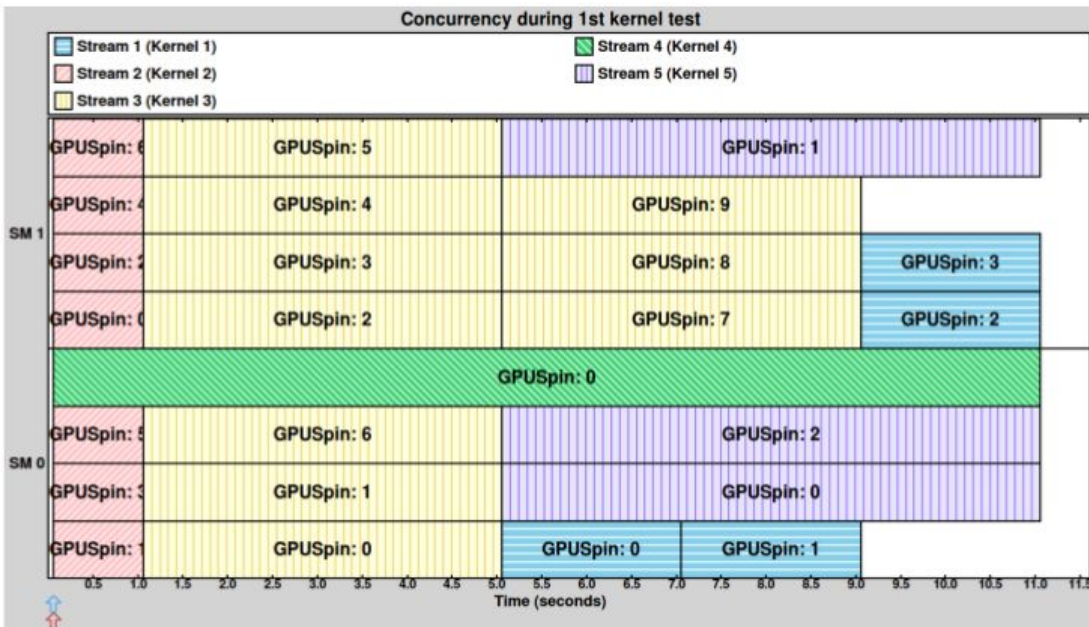
## Set-up

- Dummy kernels on Jetson TX2
- AMALTHEA implementation
- Timing measurement
- Block allocation diagrams

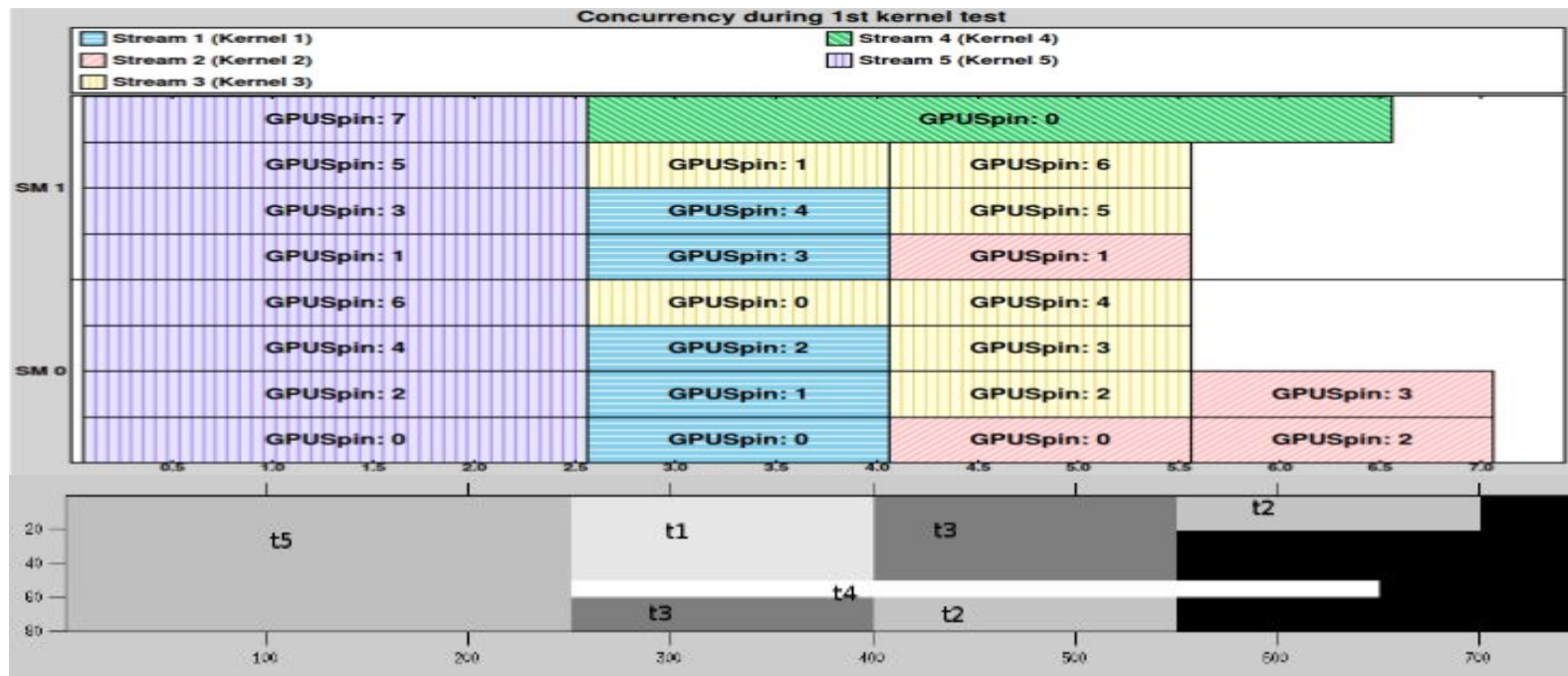
# Results



# Results



# Results



---

# Conclusions and Future work





## Conclusions

- Simple and easy to implement algorithm
- Accuracy in calculation completion times
- Integration with AMALTHEA models
- Simulation of complex use cases on AMALTHEA models



## Future Work

- Automatic CUDA code generation
- Consideration of memory transactions
- Not fully understand of Jetson TX2's GPU behaviour



## References:

1. Quesada-Barriuso, Pablo & Argüello, Francisco & B. Heras, Dora. (2013). Computing Efficiently Spectral-Spatial Classification of Hyperspectral Images on Commodity GPUs. 10.1007/978-3-319-01649-8\_2.
2. T. Amert, N. Otterness, M. Yang, J. H. Anderson, and F. D. Smith, "GPU Scheduling on the Nvidia TX2: Hidden Details Revealed," in 2017 IEEE Real-Time Systems Symposium (RTSS), 2017, pp. 104–115.



# **AMALTHEA-based GPU Response Time Analysis for NVidia's Jetson TX2**

Daniel Paredes