

Porting **Maxwell** to the GPU

Top Challenges

Juan Cañada
Head of Visualization
Next Limit Technologies

Agenda

- Maxwell overview
- Why porting to the GPU was challenging
- Performance considerations
- Using the CPU to improve the GPU engine
- Summary

Agenda

- Maxwell overview
- Why porting to the GPU was challenging
- Performance considerations
- Using the CPU to improve the GPU engine
- Summary

Maxwell Overview

GPU TECHNOLOGY
CONFERENCE



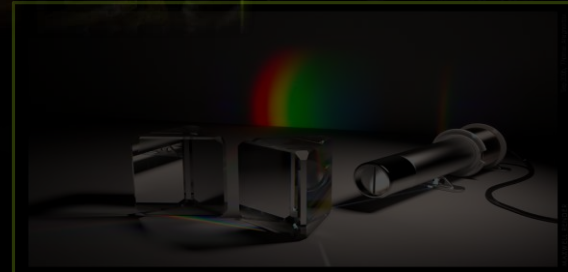
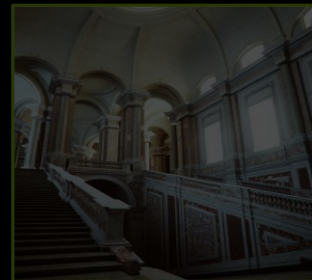
MAXWELL

- First physically based render in the market (2004)
- Ground-truth reference render
- Predictive rendering tool
- Light analysis tool



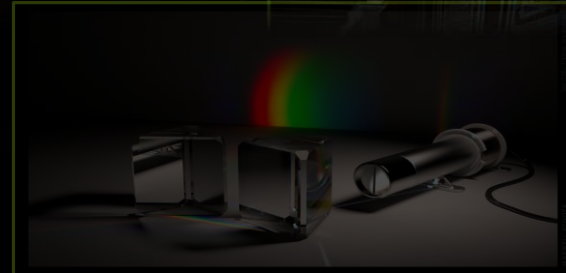
Maxwell in use

- Animation & VFX
- Architecture
- Industrial Design
- Science
- Others



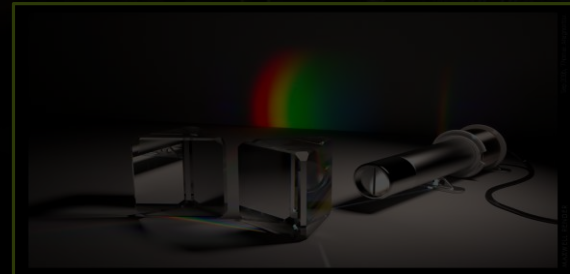
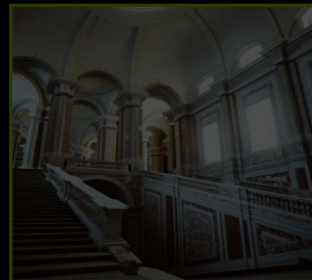
Maxwell in use

- Animation & VFX
- **Architecture**
- Industrial Design
- Science
- Others



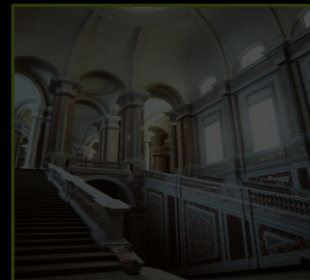
Maxwell in use

- Animation & VFX
- Architecture
- **Industrial Design**
- Science
- Others



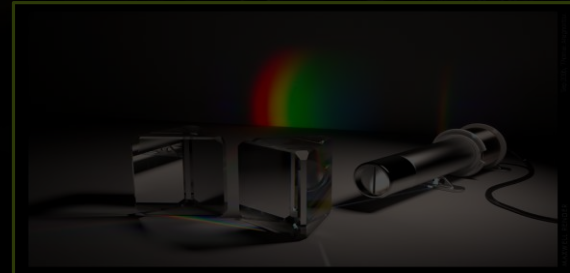
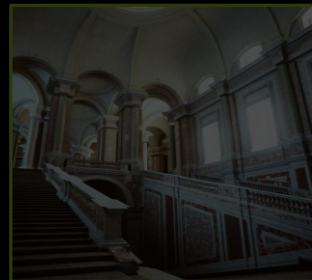
Maxwell in use

- Animation & VFX
- Architecture
- Industrial Design
- Science
- Others



Maxwell in use

- Animation & VFX
- Architecture
- Industrial Design
- Science
- Others



Agenda

- Maxwell Render overview
- Why porting to the GPU was challenging
- Performance considerations
- Using the CPU to improve the GPU engine
- Summary

Challenges

- Keep pixel accuracy
- Use GPU for predictive rendering
- Improve performance
- Spectral, unbiased, accurate PBR
- Support CPU & GPU resuming & merging
- ...



Predictive Rendering



Photo.

Nikon D70. Focal length 24mm. CCD sensor size 23.7x15.6mm. Manual settings
ISO 200 / fStop 8 / Shutter speed 90 / WB set with 6500K lamp in Cornell box
Cornell box dimensions 50x50x50cm
Compact fluorescent lamp Philips Master PL electronic 865 (6500k daylight white)
27W / 1700 lumen / efficacy 62,96 lumen/W



Maxwell render 1.0

Maxwell camera. Focal length 24mm. Film size 23.7x15.6mm.
ISO 200 / fStop 8 / Shutter speed 90 / Burn 1 / Gamma 2.2
Scene box dimensions 50x50x50cm
Maxwell emitter 6500k / 27W / 1700 lumen / efficacy 62,96 lumen/W

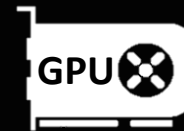
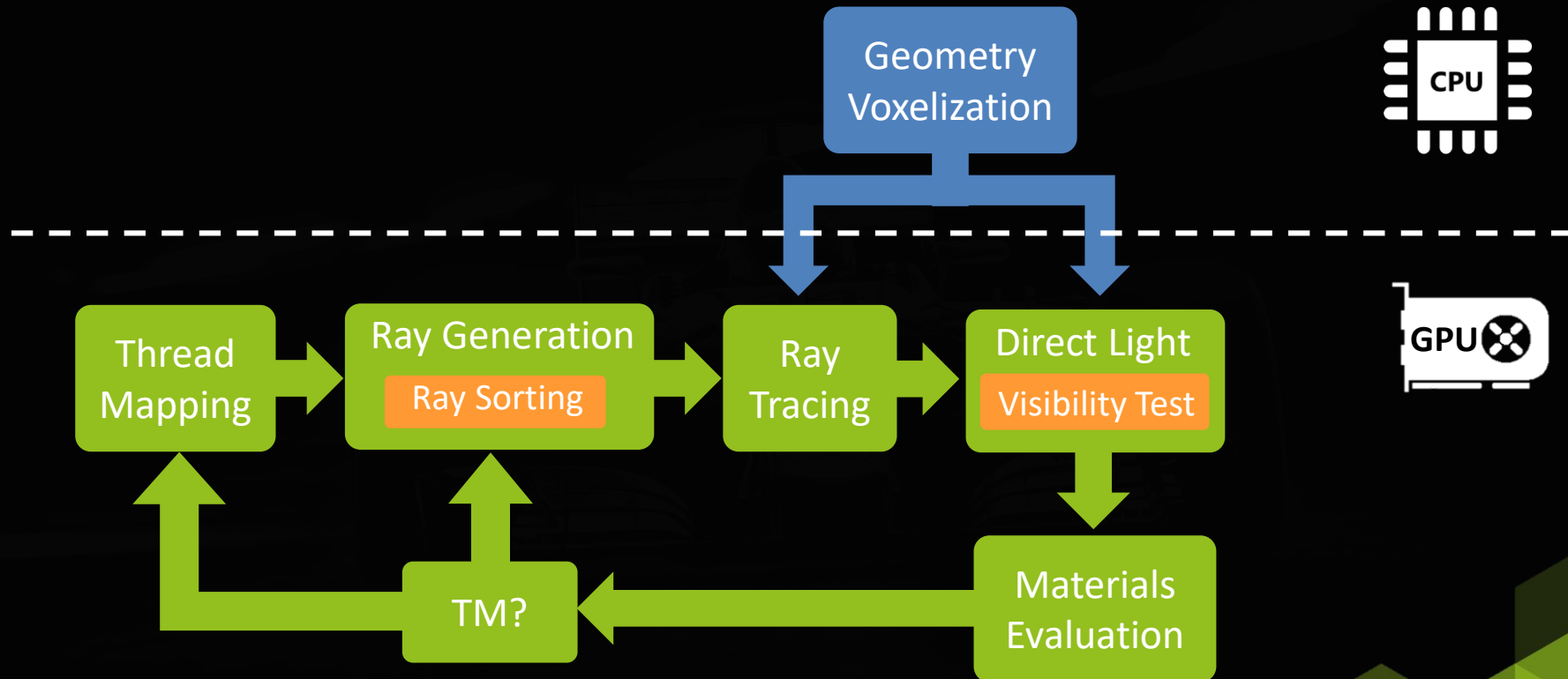
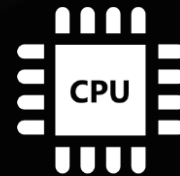
Correct → Fast 😊
Fast → Correct ☹️

Agenda

- Maxwell overview
- Why porting to the GPU was challenging
- Performance considerations
- Using the CPU to improve the GPU engine
- Summary

Maxwell GPU Architecture

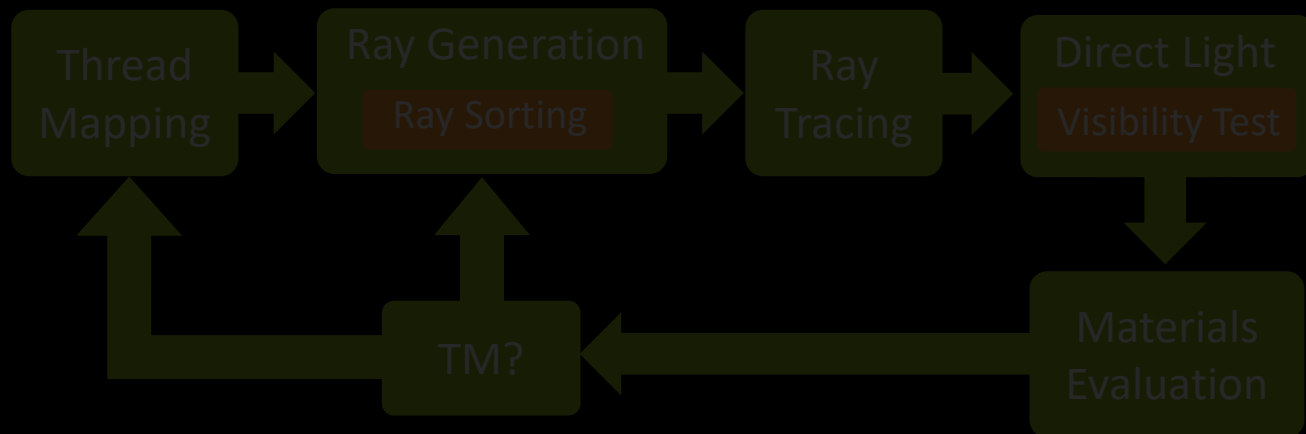
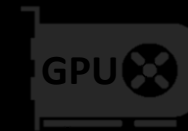
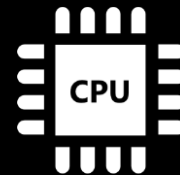
GPU TECHNOLOGY
CONFERENCE



GPU Maxwell

GPU TECHNOLOGY
CONFERENCE

Geometry
Voxelization



- **Voxelization**

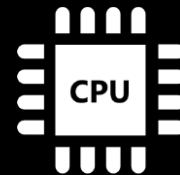
- Same Voxelization system as the CPU render
- Currently performed in CPU just once
- BVH
 - Binary tree (each node has 2 childs)
 - Coherent traversal



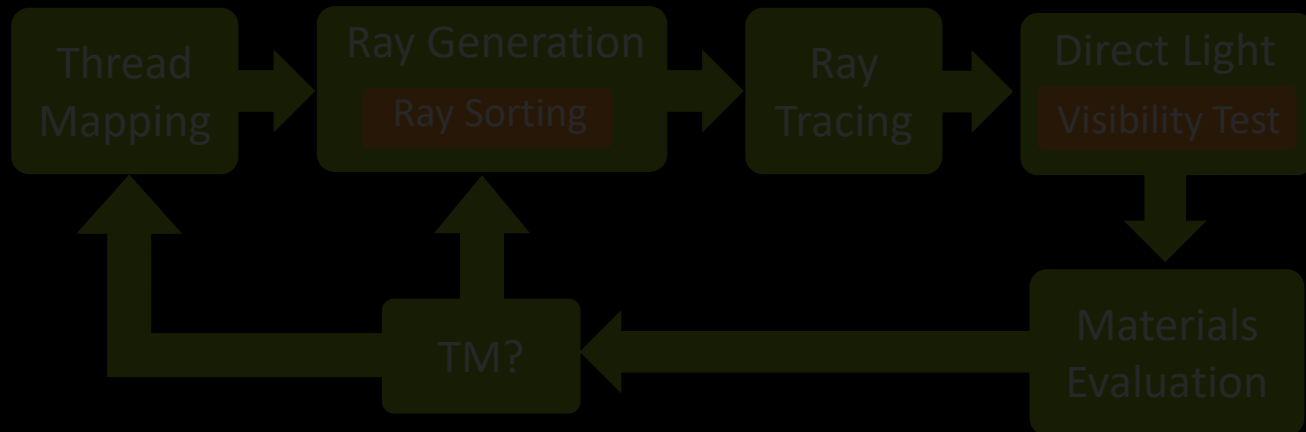
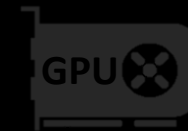
- + All threads fetch same amount of data / node
- + Increase coherence in performance
- Trees become bigger

GPU Maxwell

GPU TECHNOLOGY
CONFERENCE

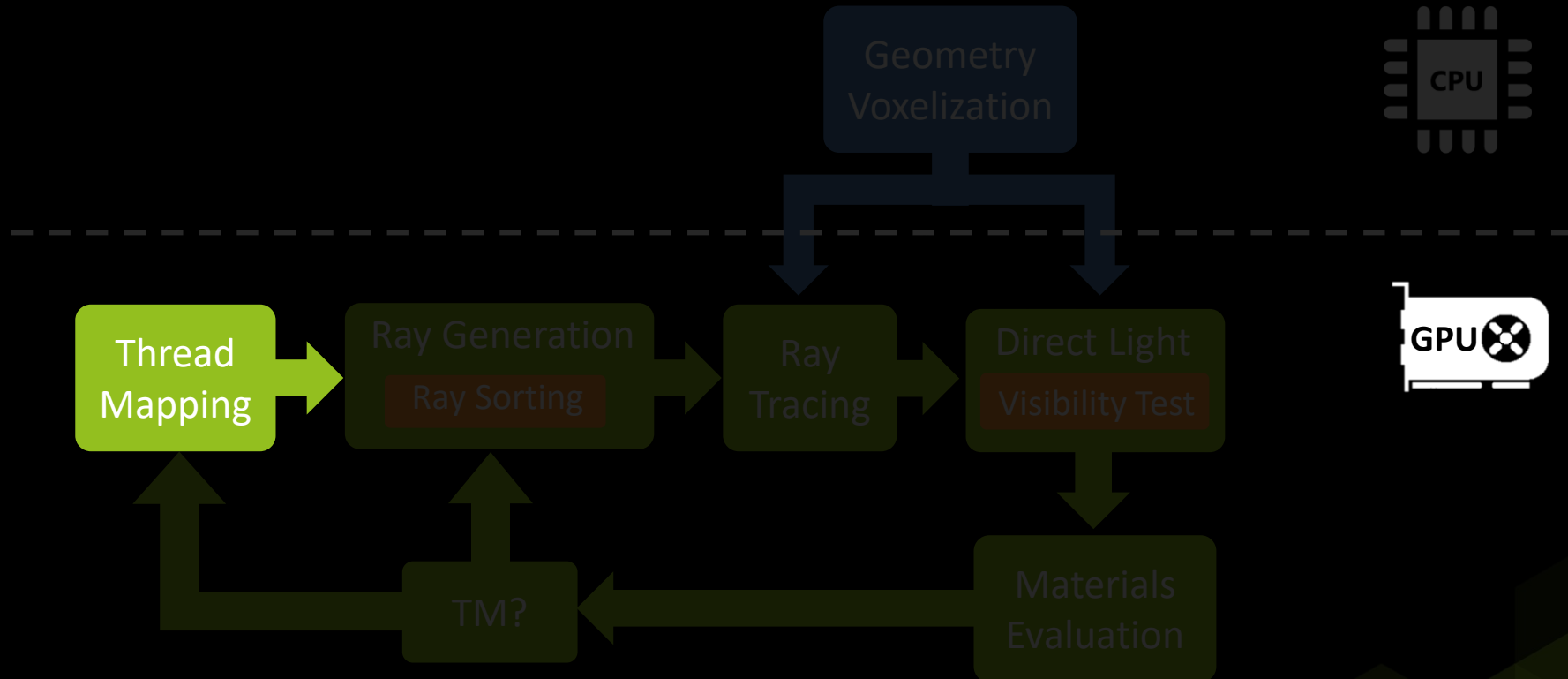
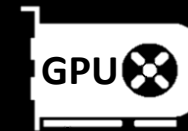
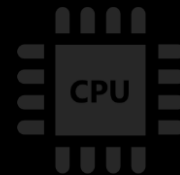


Geometry
Voxelization



GPU Maxwell

GPU TECHNOLOGY
CONFERENCE



GPU Maxwell

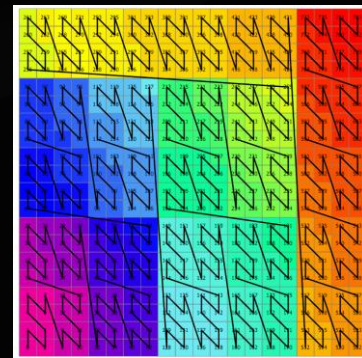
- Thread Mapping

- Module that manages THREAD / PIXEL mapping

- Sampling Level (SL)

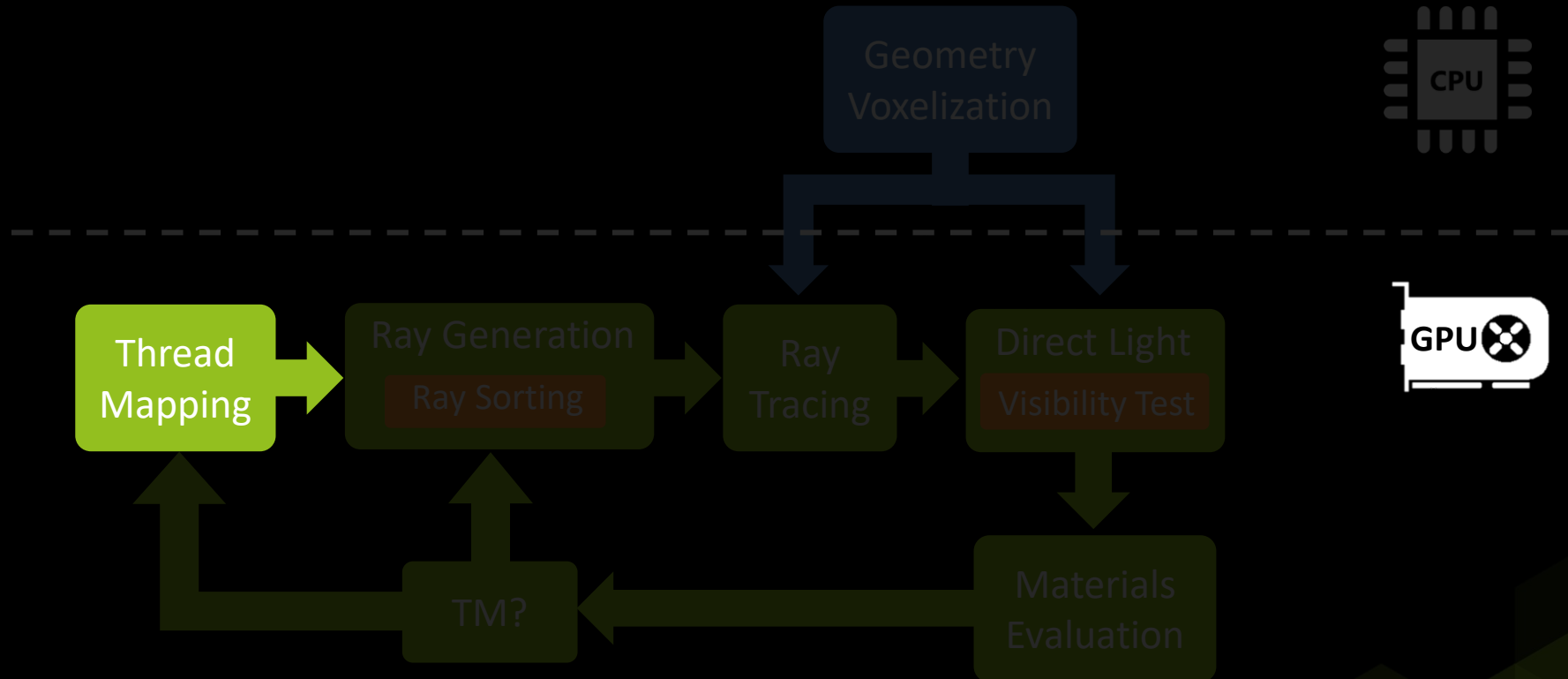
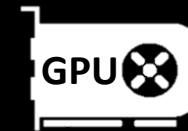
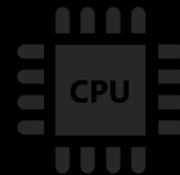
- Low Morton → Curve
 - Medium → Balances SPP
 - High → Uses Variance

Morton Curve



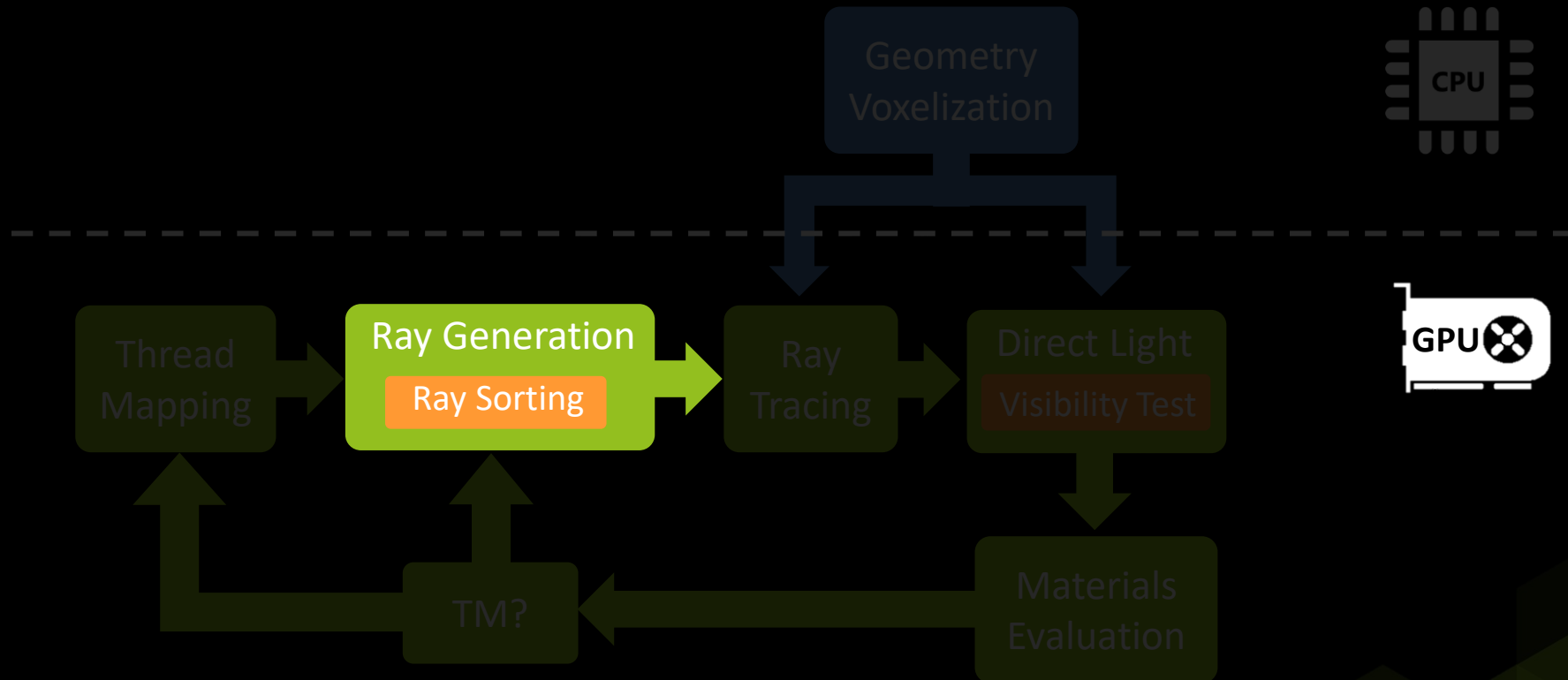
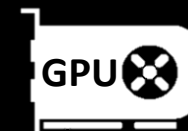
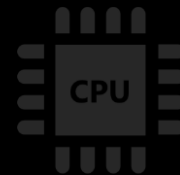
GPU Maxwell

GPU TECHNOLOGY
CONFERENCE



GPU Maxwell

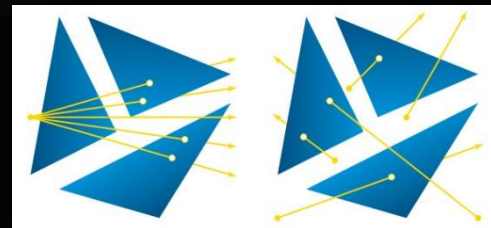
GPU TECHNOLOGY
CONFERENCE



- Ray Generation Module

- Primary Rays (PR)

- Rays shot from camera
- High degree of coherence
- Two neighboring rays will hit nearby similar objects



- Secondary Rays (SR)

- Rays shot from surfaces
- No coherence
- Two neighbouring rays might hit different objects

- Ray Generation Module
 - Thread blocks with just PR
 - High degree of coherence
 - Best performance situation
 - Thread blocks with just SR
 - All will take much more time than PR
 - The worst SR will drive the performance
 - Thread blocks with PR and SR
 - SR will hurt PR performance

- Ray Generation Module
 - How do we handle it?
 - GPU Ray sorting by Ray Type



- Ray Generation Module
 - How do we handle it?
 - GPU Ray sorting by Ray Type

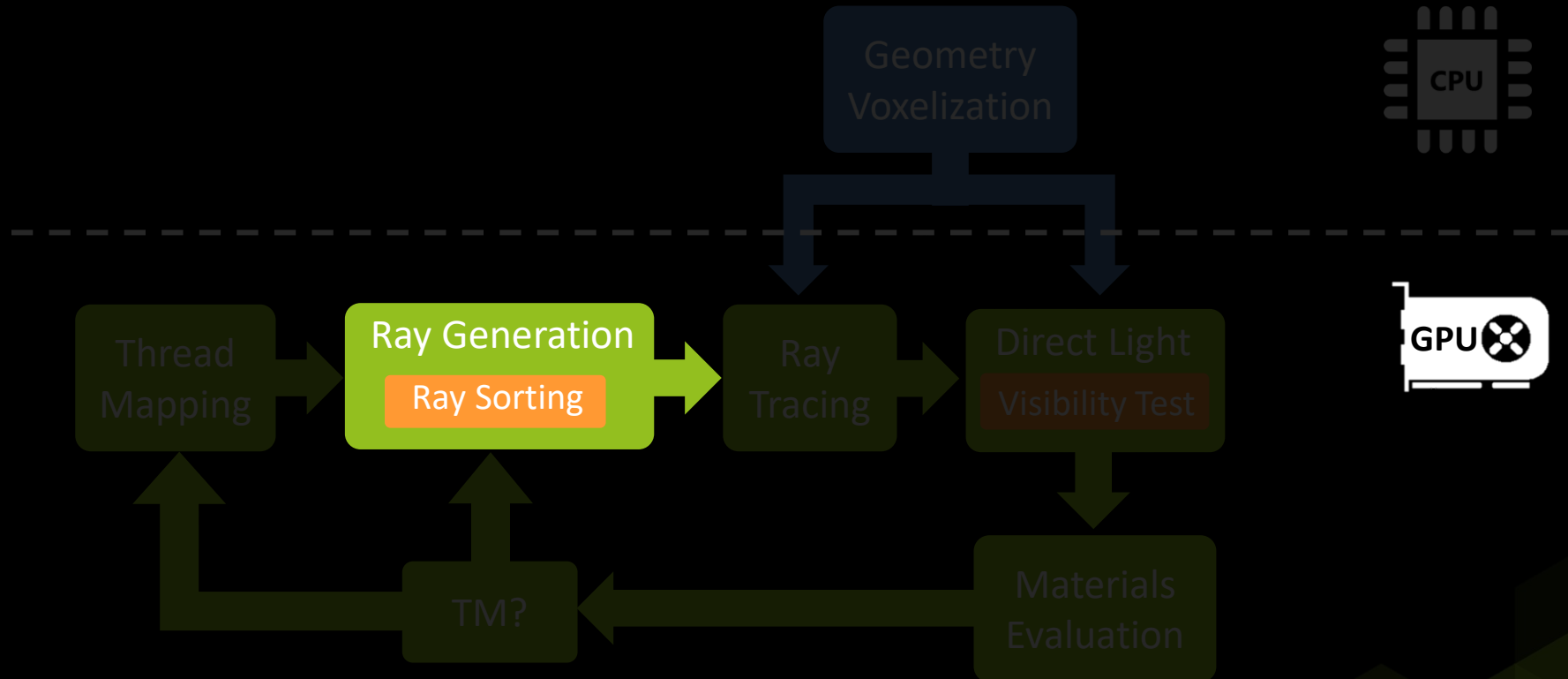
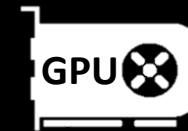
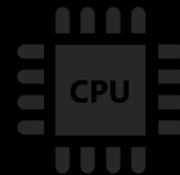


- Ray Generation Module
 - How do we handle it?
 - GPU Ray sorting by Ray Type
 - Sorting is really fast
 - Simple, yet powerful
 - Do it just after 2nd bounce
 - Not needed for PR
 - Performance boost is scene dependant

- Ray Generation Module
 - How do we handle it?
 - GPU Ray sorting by Ray Type
 - **Considerations**
 - Not useful for medium to small-res images
 - Use an indirection buffer
 - Cleaner code
 - Avoids moving global data
 - Much better performance

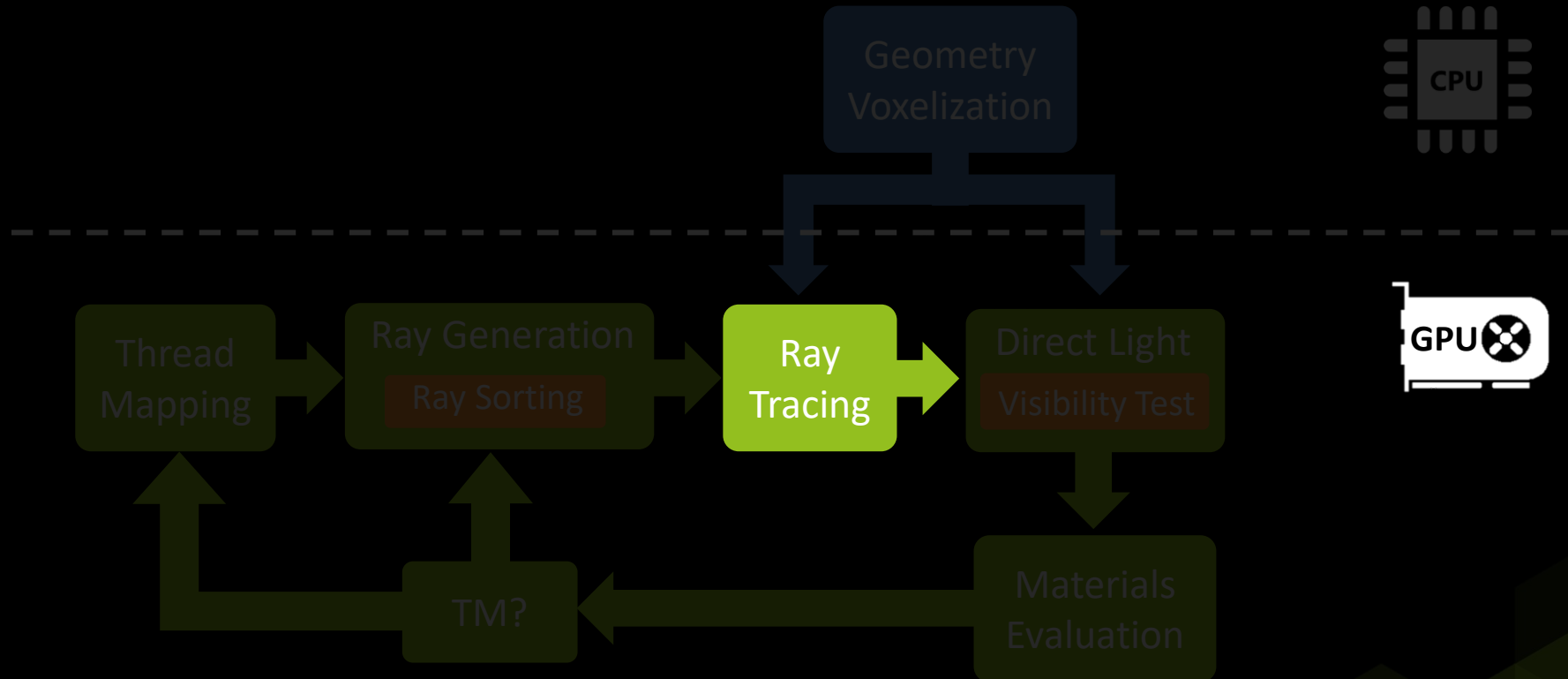
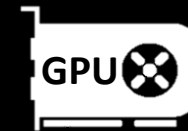
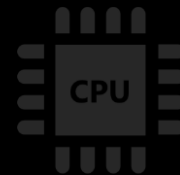
GPU Maxwell

GPU TECHNOLOGY
CONFERENCE



GPU Maxwell

GPU TECHNOLOGY
CONFERENCE

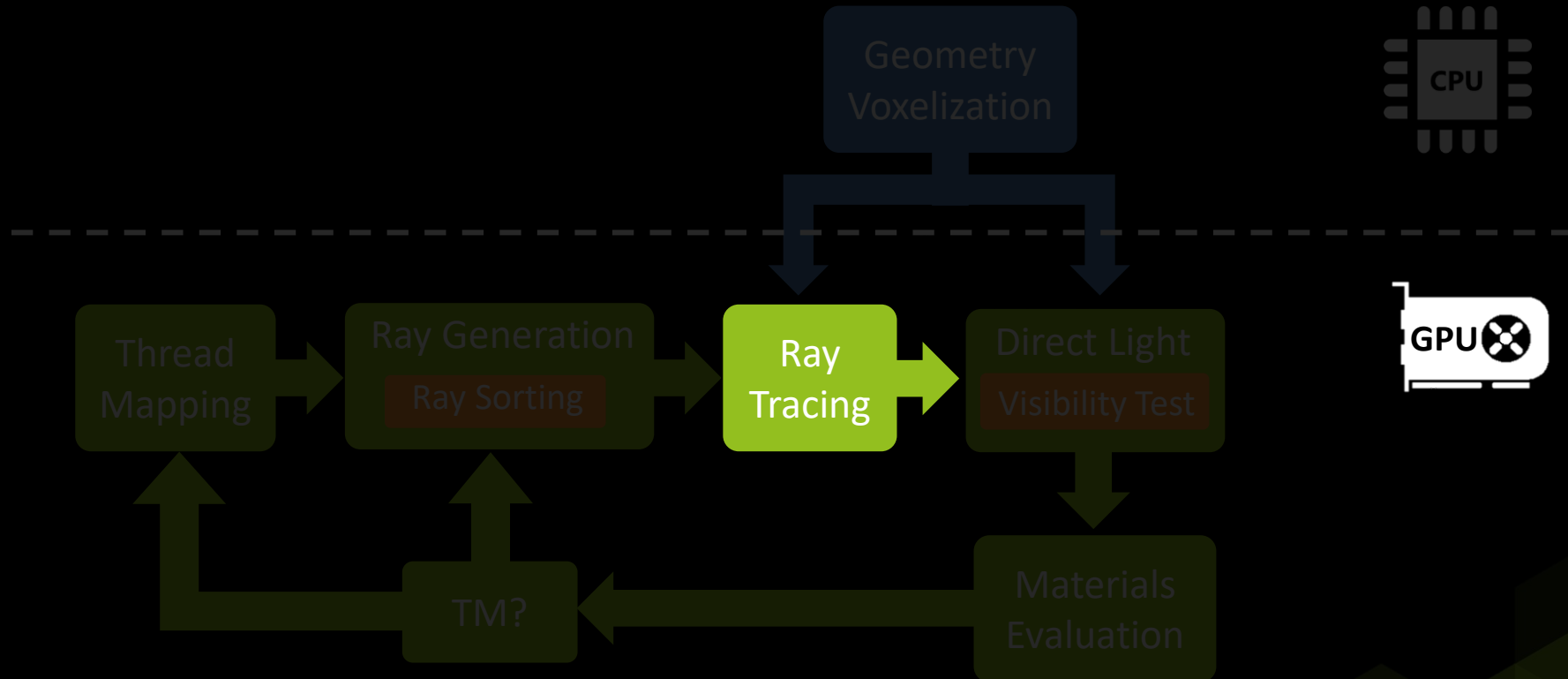
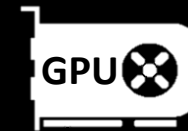
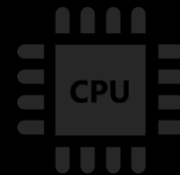


- Ray Tracing Module
 - GPU architecture dependent kernels
 - Fermi, Kepler, Maxwell
 - Use every architecture strengths



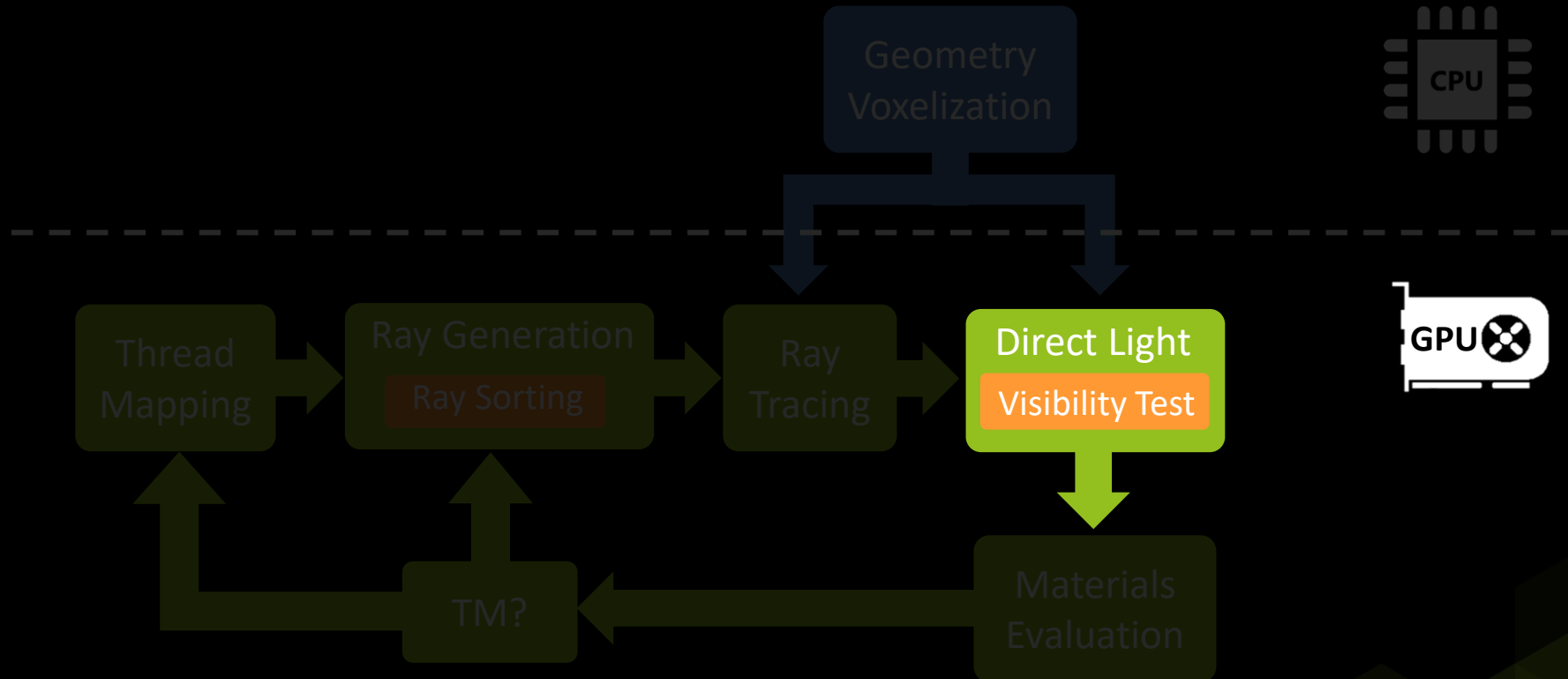
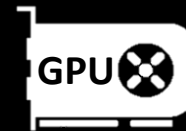
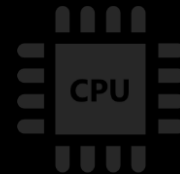
GPU Maxwell

GPU TECHNOLOGY
CONFERENCE



GPU Maxwell Render

GPU TECHNOLOGY
CONFERENCE

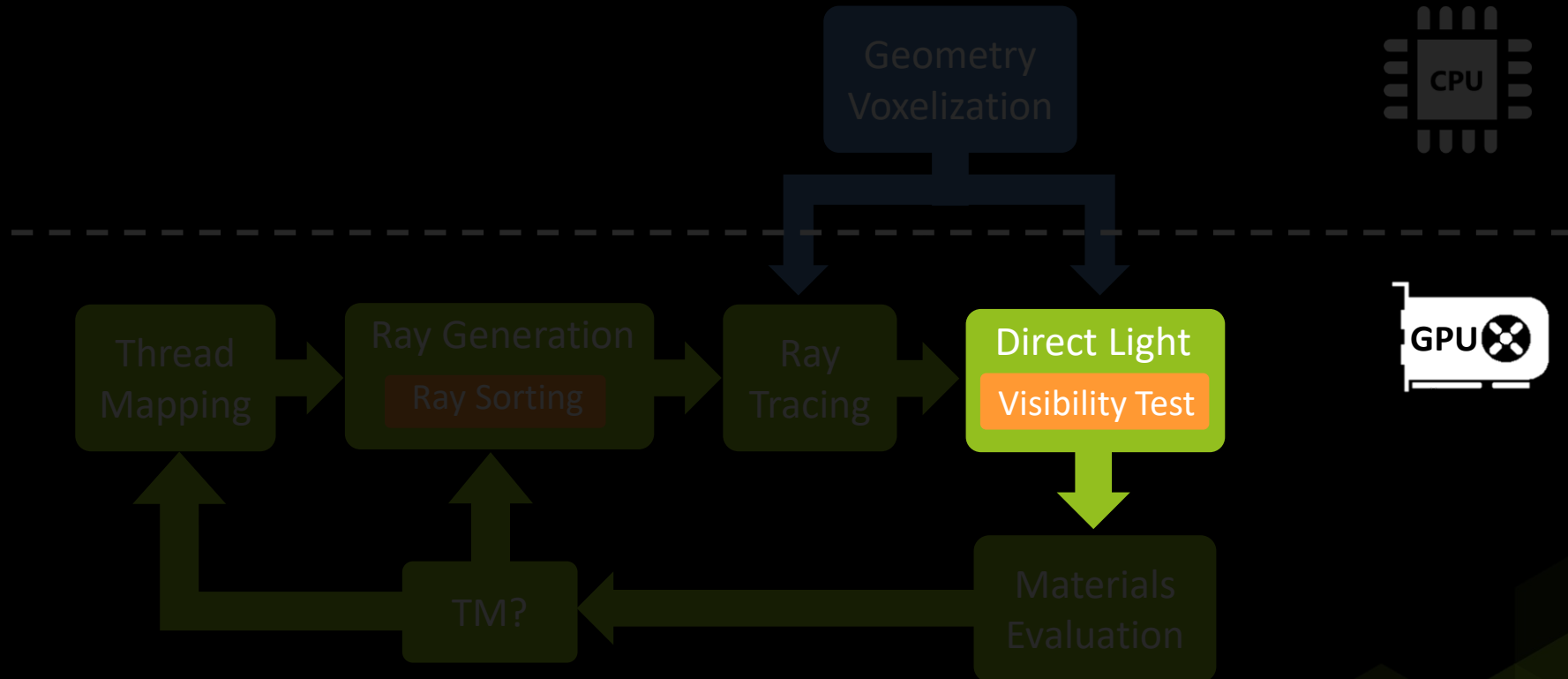
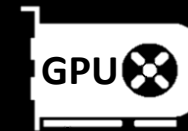
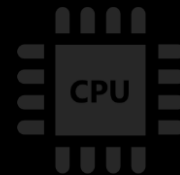


Direct Light Module

1. Sample scene emitters at each path node
 - Two strategies
 - Sample 1 random emitter / sample
 - Sample all emitters / sample
2. Visibility test
 - Trace shadow rays
 - Incoherent rays → Ray sorting does not help
3. Many other optimizations

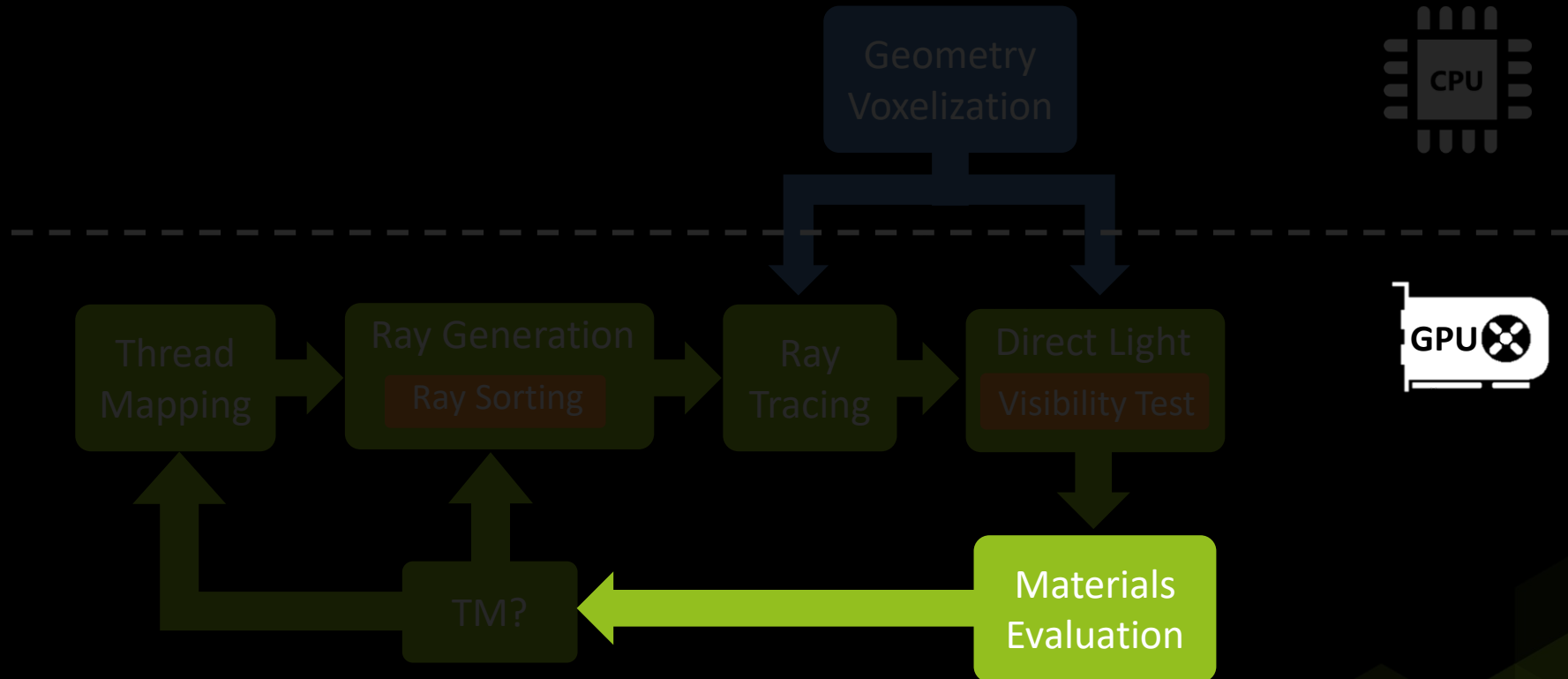
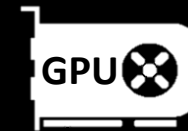
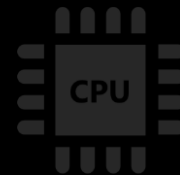
GPU Maxwell

GPU TECHNOLOGY
CONFERENCE



GPU Maxwell

GPU TECHNOLOGY
CONFERENCE



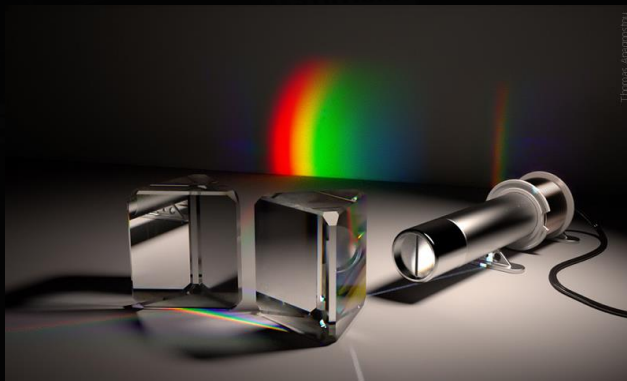
GPU Maxwell

- Materials Evaluation Module
 - Maxwell materials are complex
 - Many layers and many BSDFs / layer → very generic



Materials Evaluation Module

- Bbig kernels are harmful
- Samples evaluating different materials
 - Access different data
 - Execute different code



Thomas Wegmayer

- Materials Evaluation Module

- Materials Group Queue System (MGQS)

1. Every material is assigned a Material Group ID
2. Queue system for Material Groups (MG)
3. Every queue has specific kernels
 - + Avoid big kernels
4. Samples are queued to the corresponding MG Queue
5. All samples evaluating the same MG are executed together
 - + Increased coherence in execution time
 - + Increased coherence in data access

- Materials Evaluation Module
 - Materials Group Queue System (MGQS)
 1. Every material is assigned a Material Group ID
 2. Queue system for Material Groups (MG)
 3. Every queue has specific kernels
 - + Avoid big kernels
 4. Samples are queued to the corresponding MG Queue
 5. All samples evaluating the same MG are executed together
 - + Increased coherence in execution time
 - + Increased coherence in data access

- Materials Evaluation Module
 - Materials Group Queue System (MGQS)
 1. Every material is assigned a Material Group ID
 2. Queue system for Material Groups (MG)
 3. Every queue has specific kernels
 - + Avoid big kernels
 4. Samples are queued to the corresponding MG Queue
 5. All samples evaluating the same MG are executed together
 - + Increased coherence in execution time
 - + Increased coherence in data access

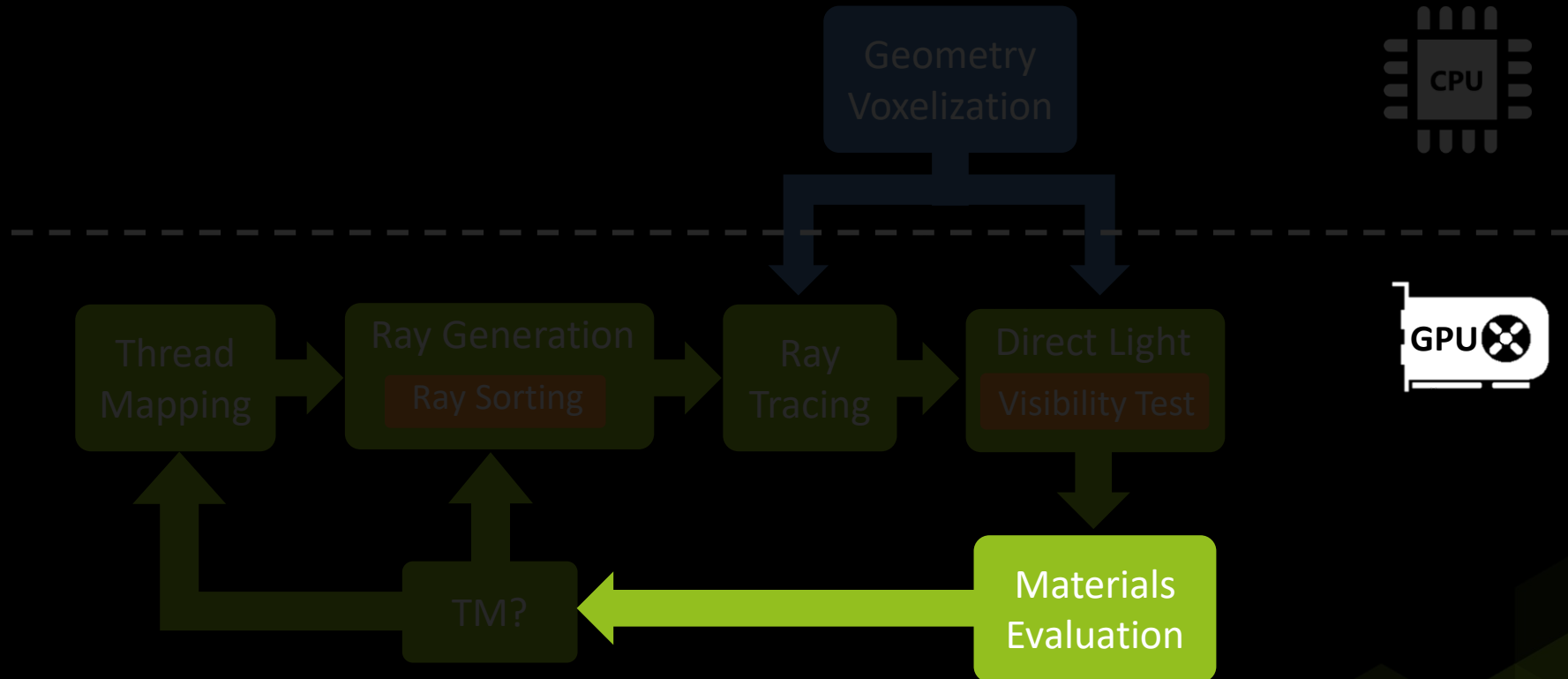
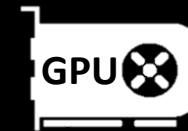
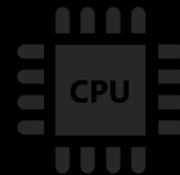
- Materials Evaluation Module
 - Materials Group Queue System (MGQS)
 1. Every material is assigned a Material Group ID
 2. Queue system for Material Groups (MG)
 3. Every queue has specific kernels (Avoid big kernels)
 4. Samples are queued to the corresponding MG Queue
 5. All samples evaluating the same MG are executed together
 - + Increased coherence in execution time
 - + Increased coherence in data access

- Materials Evaluation Module
 - Materials Group Queue System (MGQS)
 1. Every material is assigned a Material Group ID
 2. Queue system for Material Groups (MG)
 3. Every queue has specific kernels (Avoid big kernels)
 4. Samples are queued to the corresponding MG Queue
 5. All samples evaluating the same MG are executed together
 - Increased coherence in execution time
 - Increased coherence in data access

- Materials Evaluation Module
 - Materials Group Queue System (MGQS)
 1. Every material is assigned a Material Group ID
 2. Queue system for Material Groups (MG)
 3. Every queue has specific kernels (Avoid big kernels)
 4. Samples are queued to the corresponding MG Queue
 5. All samples evaluating the same MG are executed together
 - Increased coherence in execution time
 - Increased coherence in data access

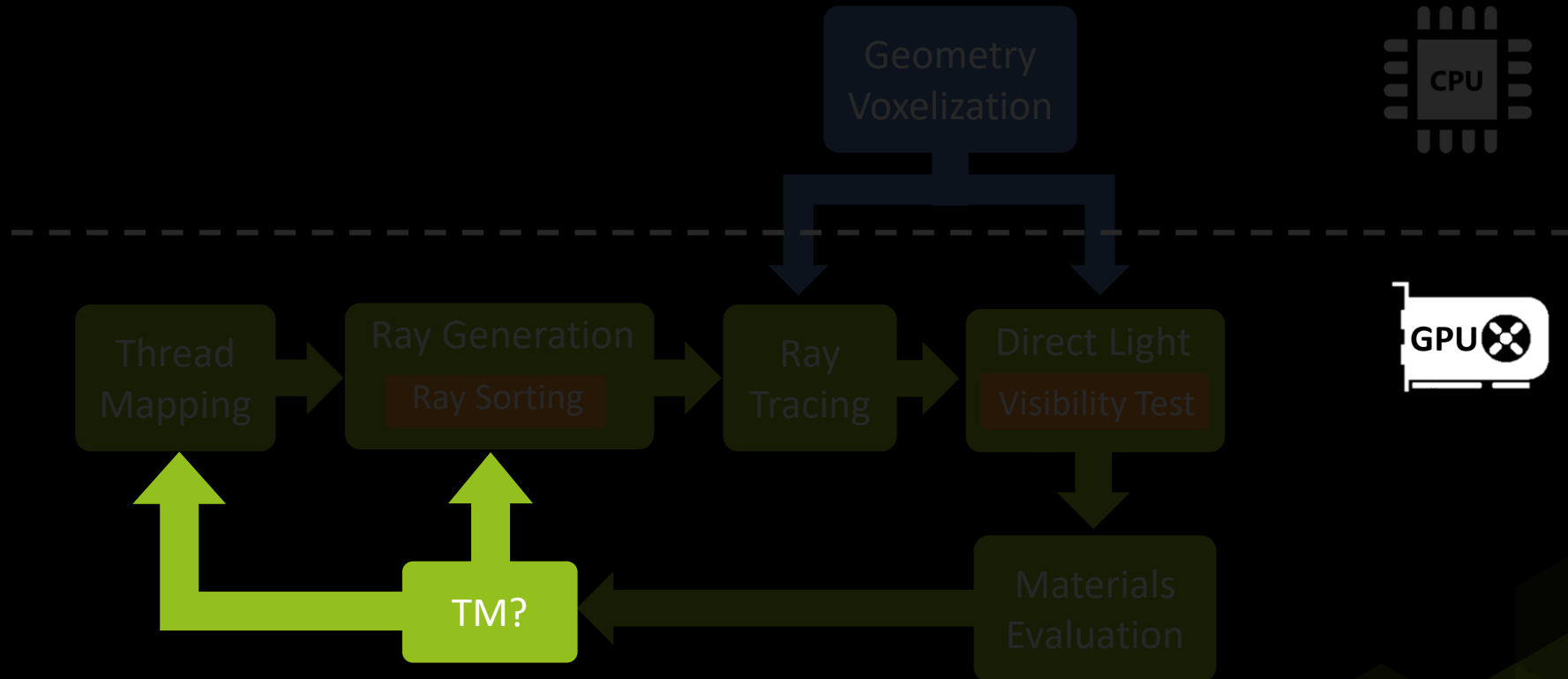
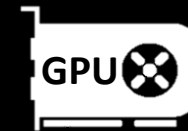
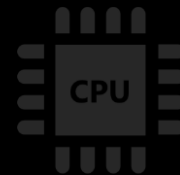
GPU Maxwell

GPU TECHNOLOGY
CONFERENCE



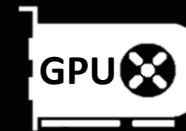
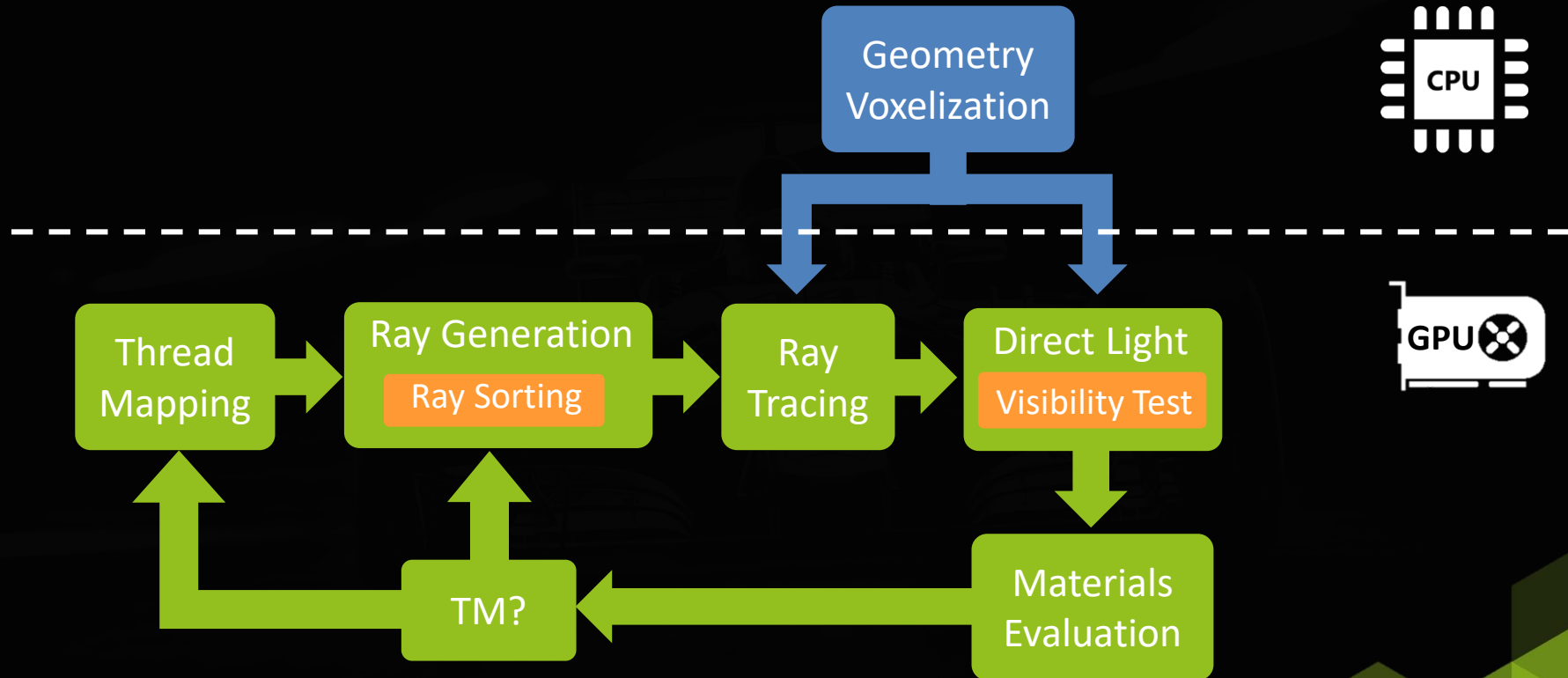
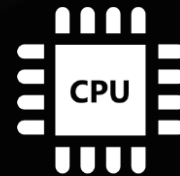
GPU Maxwell

GPU TECHNOLOGY
CONFERENCE



GPU Maxwell

GPU TECHNOLOGY
CONFERENCE



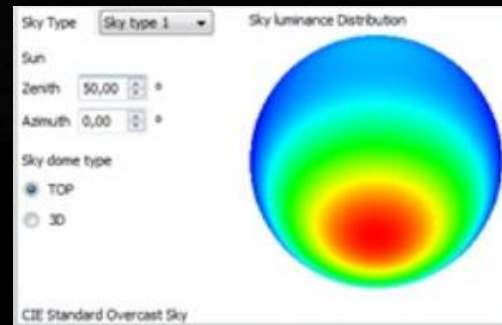
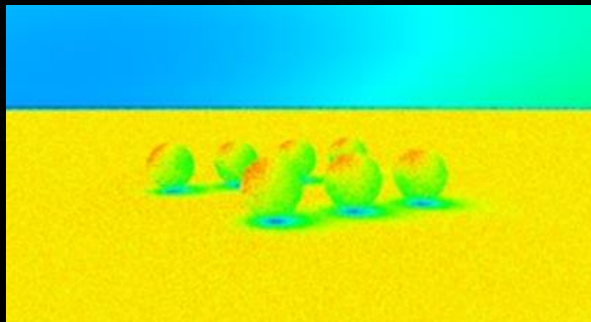
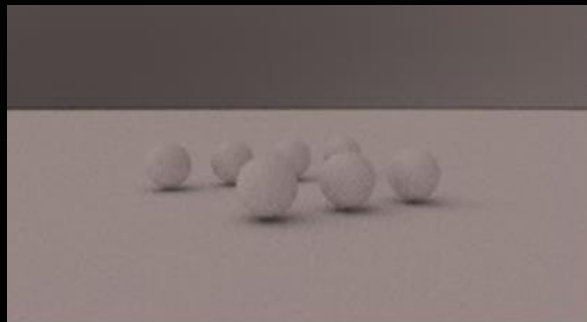
Agenda

- Maxwell overview
- Why porting to the GPU was challenging
- Performance considerations
- Using the CPU to improve the GPU engine
- Summary

Using the CPU to improve the GPU engine

Why using our CPU engine as ground truth?

- 12 years old → Stable & Robust
- Used many times for validation purposes



CPU vs GPU Case Studies

Guggenheim scene



Teapot scene



Guggenheim Scene



Guggenheim Scene



ISSUES

- Slight differences in intensity
- Noise in some areas
- Subtle changes in glossy surfaces

STRATEGY

- Simplifying & Isolating (surprise :P)
- Automated numerical comparisons
 - Raytracing text output
- Ray viewer



Guggenheim Scene

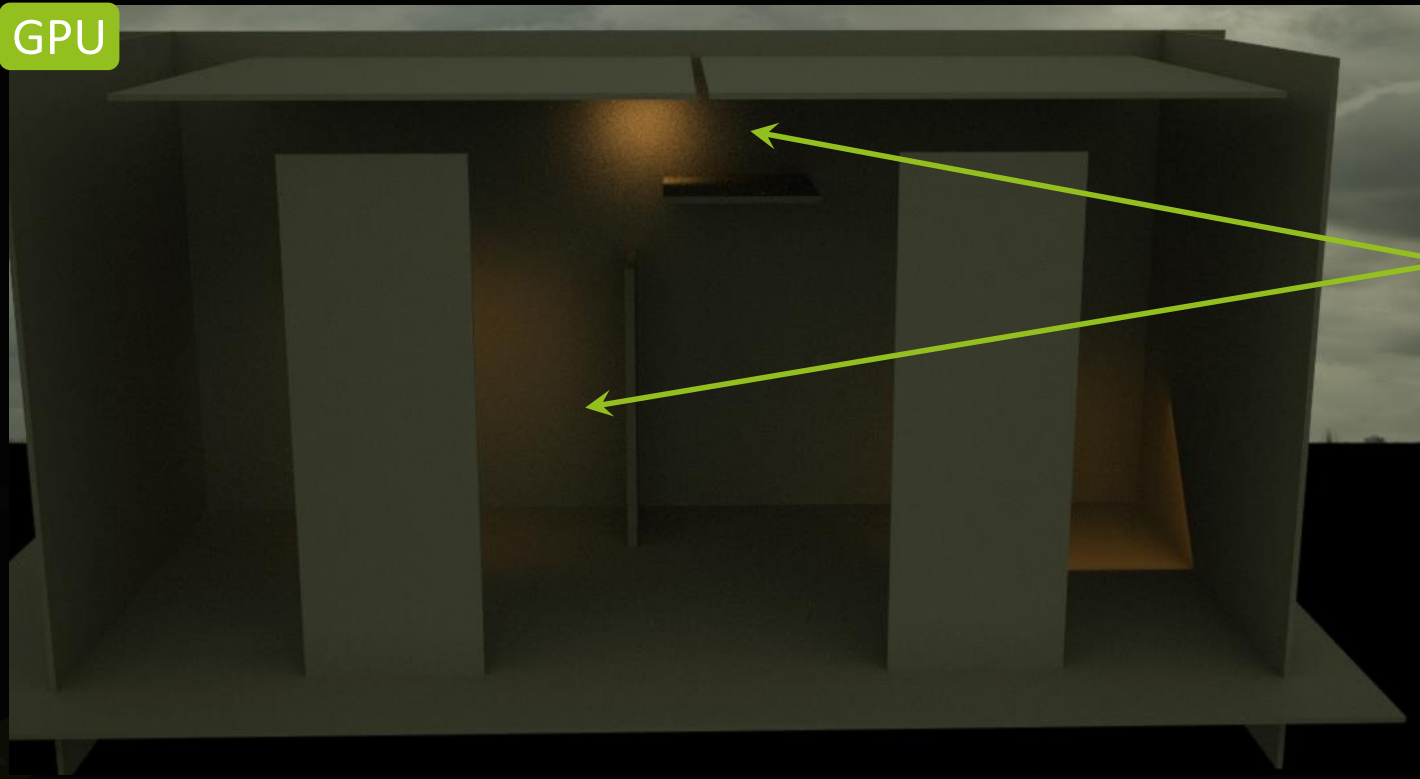
GPU TECHNOLOGY
CONFERENCE

CPU



Guggenheim Scene

GPU



Different Intensity
+
Noise Problems



Guggenheim Scene – Intensity & Noise

CPU



GPU



CPU



GPU



FINDINGS

- Emitters intensity
 - Hidden property of emitters was not working properly
 - Non-visible emitters were causing occlusions
 - Loss of energy
- Noise
 - QMC had some problems for higher dimensions

Guggenheim Scene

CPU

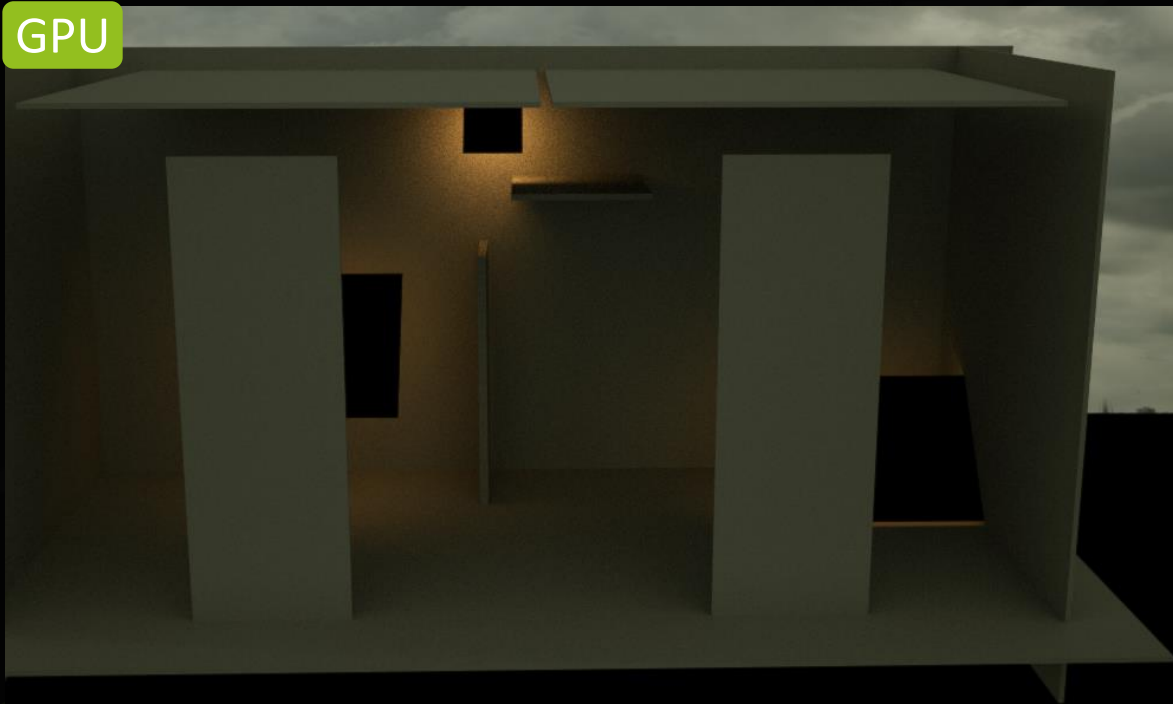


FIXED



Guggenheim Scene

GPU



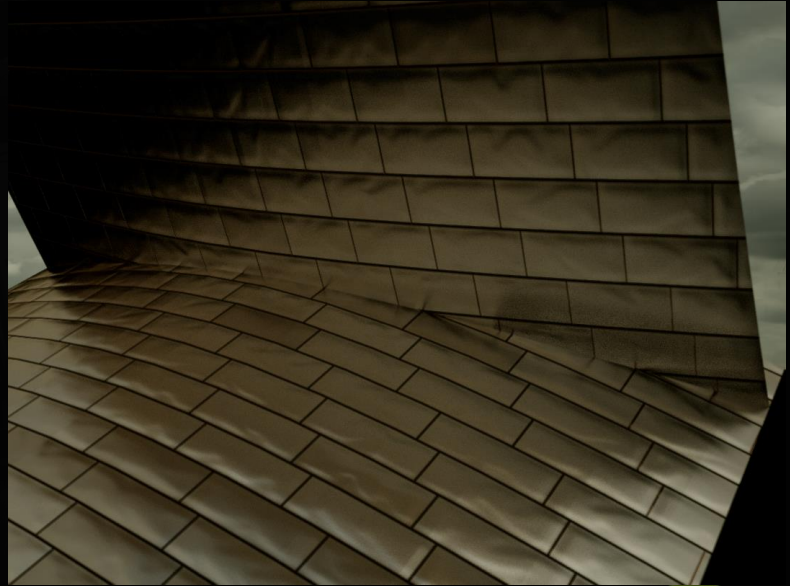
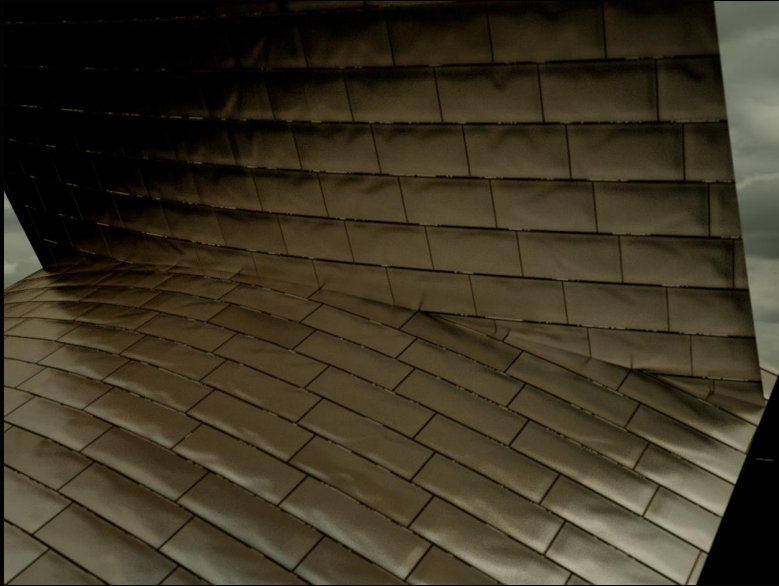
FIXED



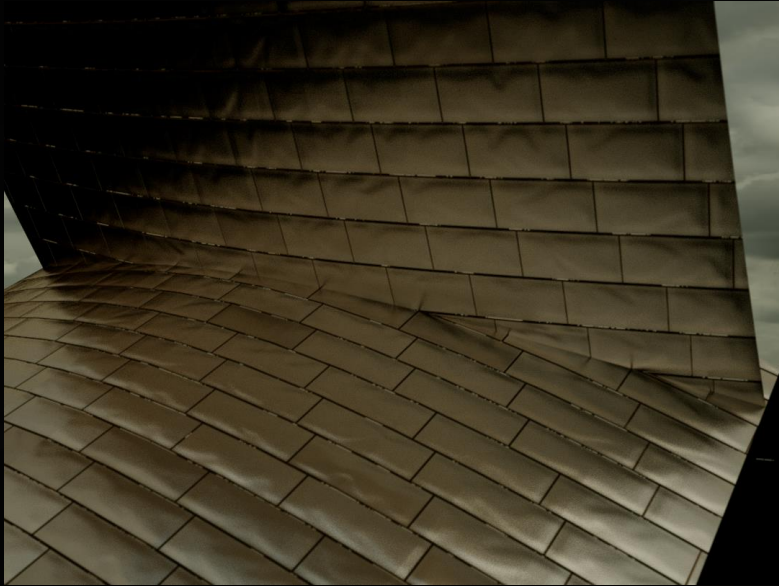
Guggenheim Scene



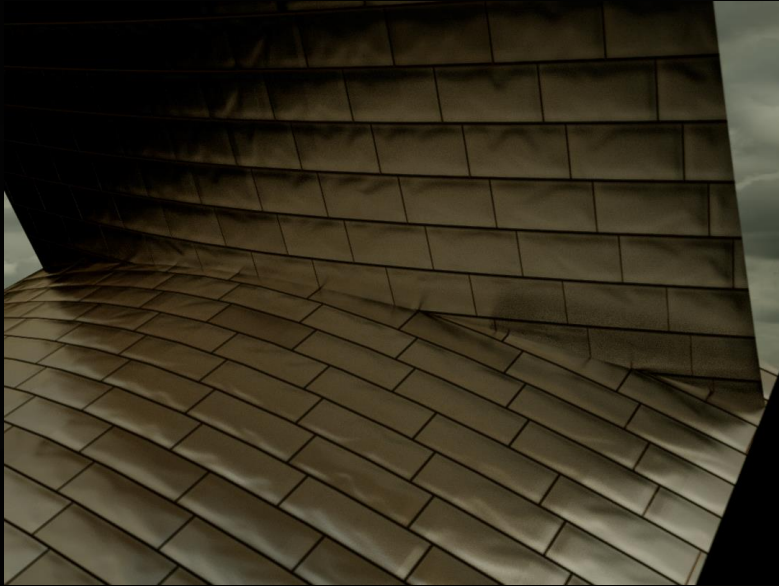
Guggenheim Scene – Differences in glossies



Guggenheim Scene – Differences in glossies



Guggenheim Scene – Differences in glossies

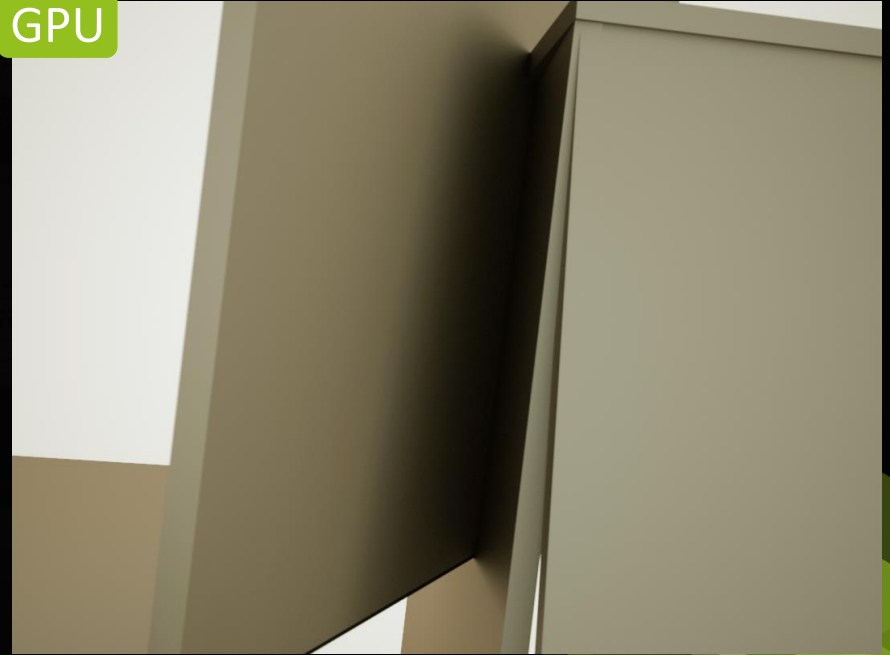


Guggenheim Scene – Differences in glossies

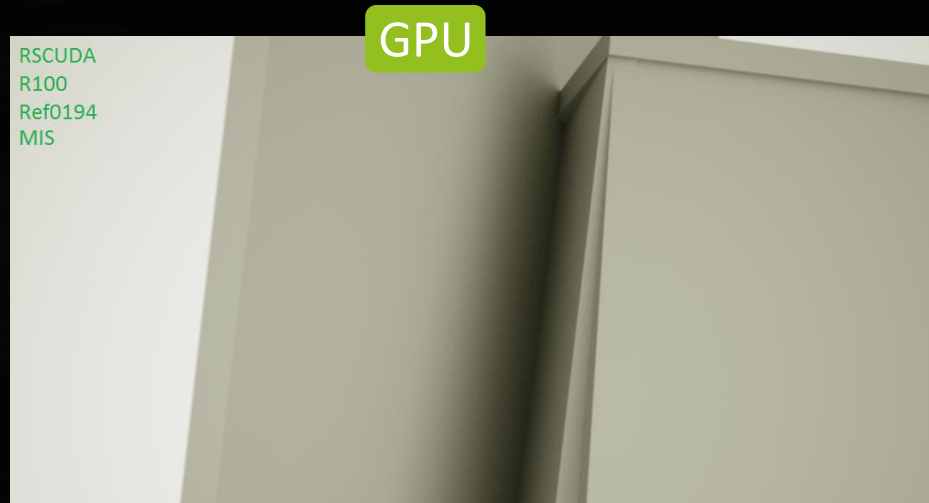
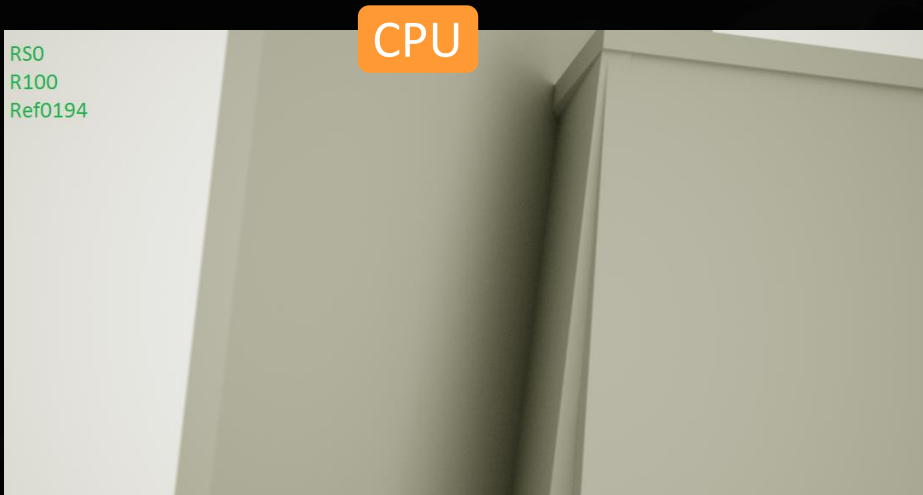
CPU



GPU

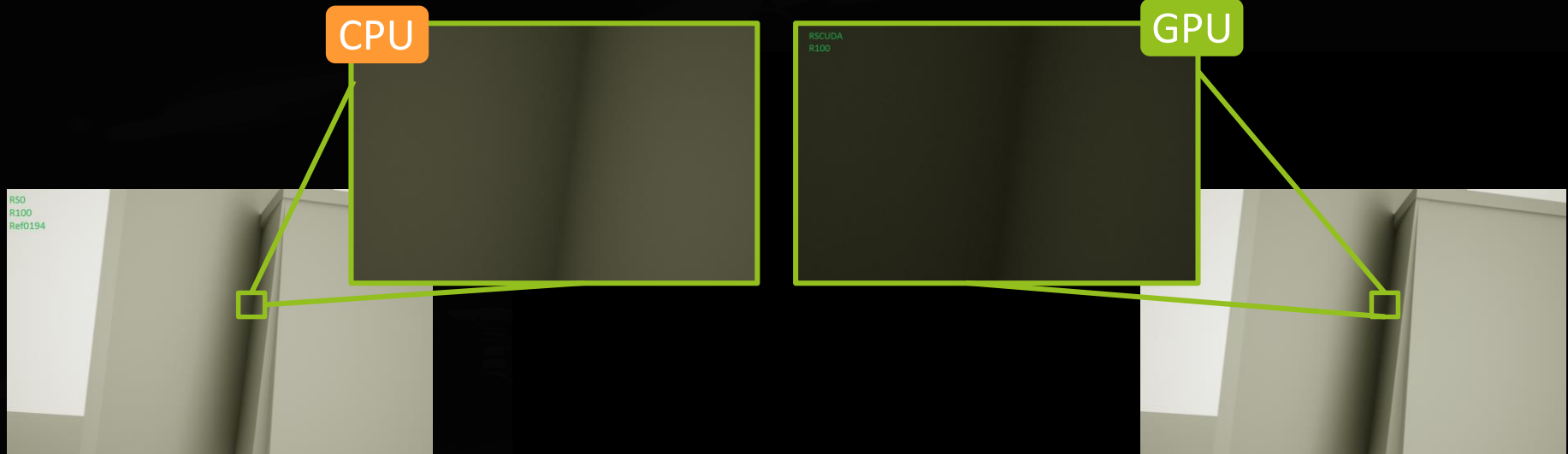


Guggenheim Scene – Differences in glossies

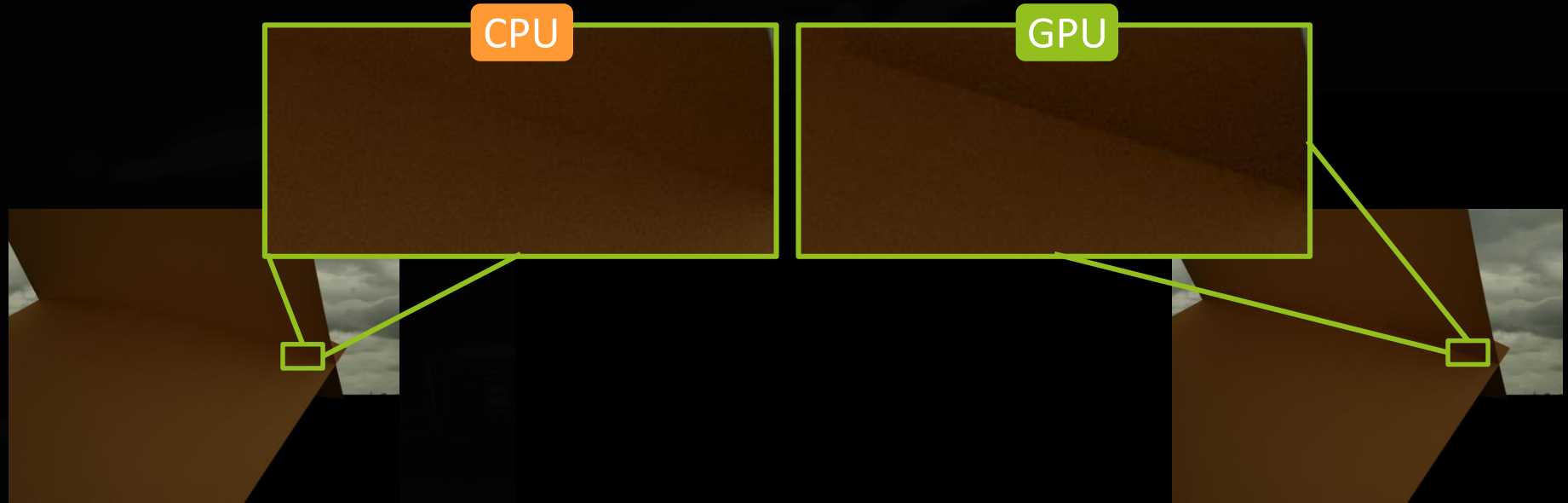


- Simplify the material → Lambert

Guggenheim Scene – Differences in glossies



Guggenheim Scene – Differences in glossies

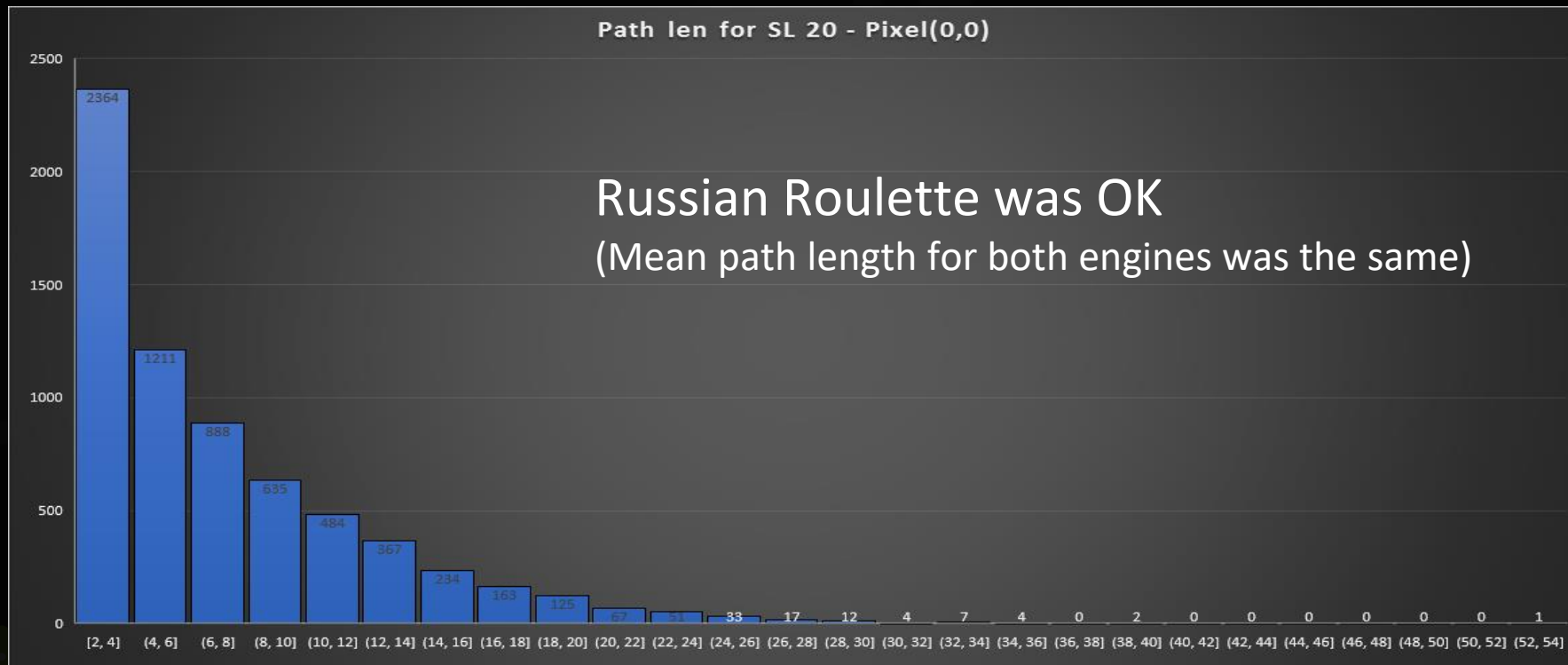


Guggenheim Scene – Differences in glossies

- It turned out it was not related to materials
 - Both glossy and lambert have the same problem
 - Difficult to isolate
- Possible problems
 - QMC numbers bug?
 - Russian Roulette bug?
 - Ray / triangle intersection issues with indirect bounces?
 - Energy accumulation problem?
 - Precision issues?
 - ...

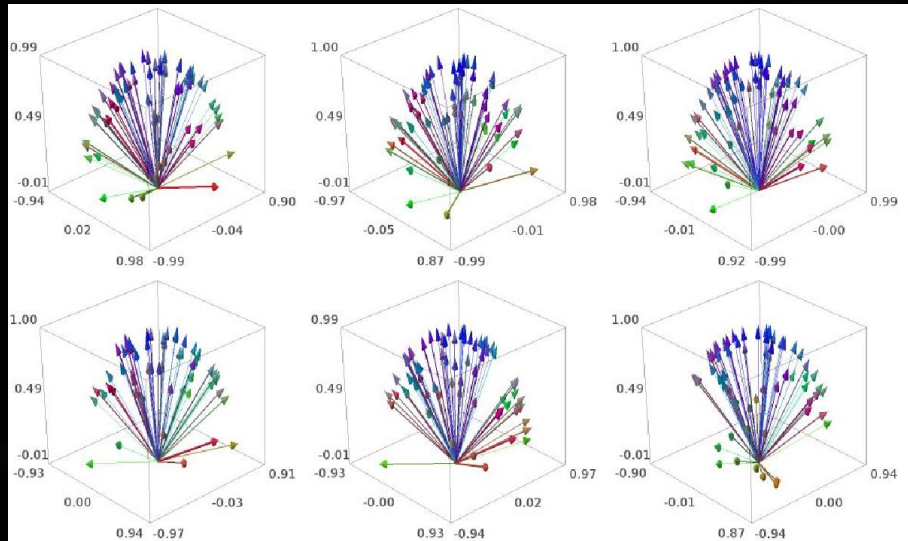


Guggenheim Scene – Differences in glossies

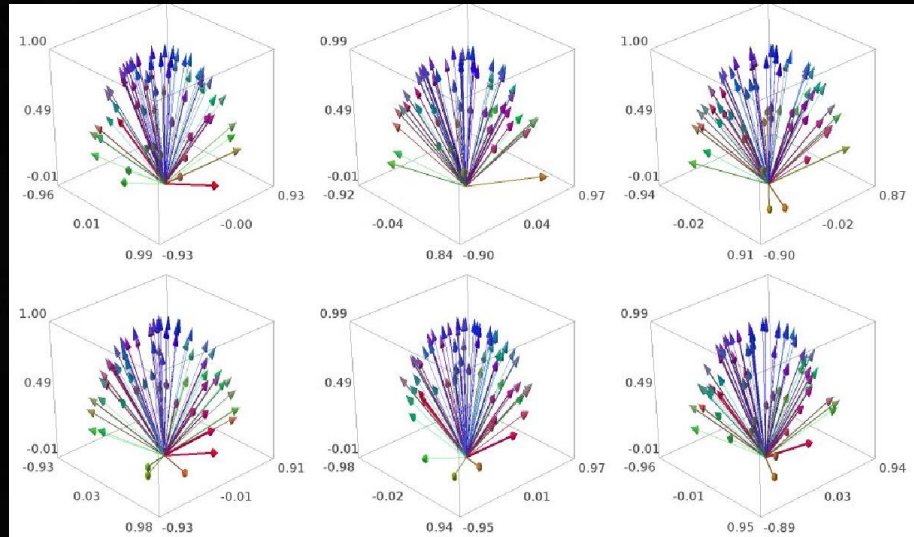


Guggenheim Scene – Differences in glossies

CPU – QMC Distributions

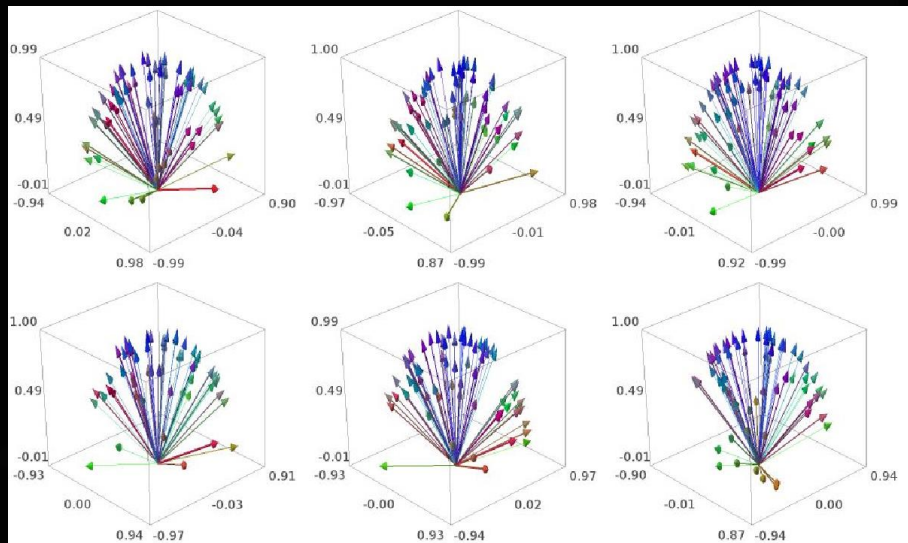


GPU – QMC Distributions

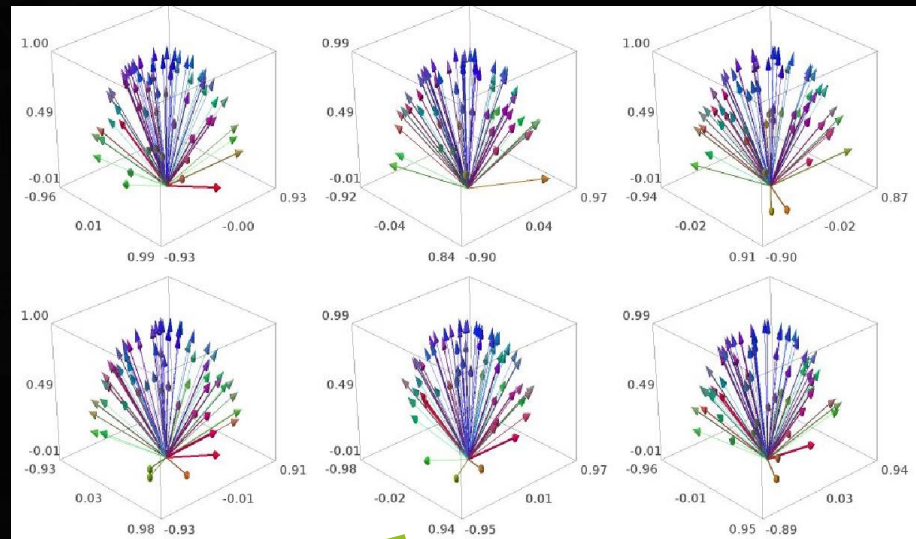


Guggenheim Scene – Differences in glossies

CPU – QMC Distributions



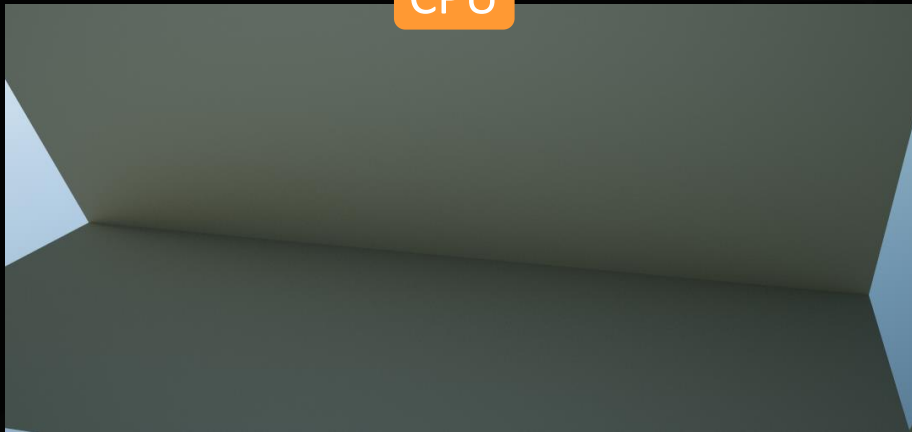
GPU – QMC Distributions



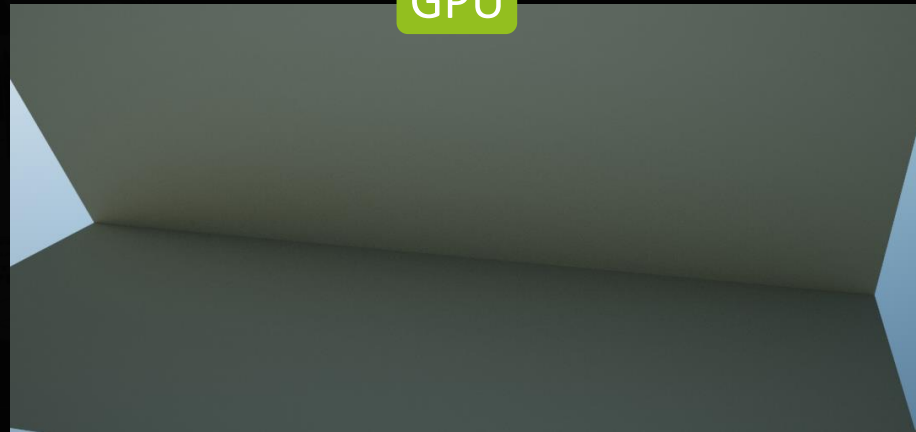
Automated tests detected differences!

Guggenheim Scene – Differences in glossies

CPU

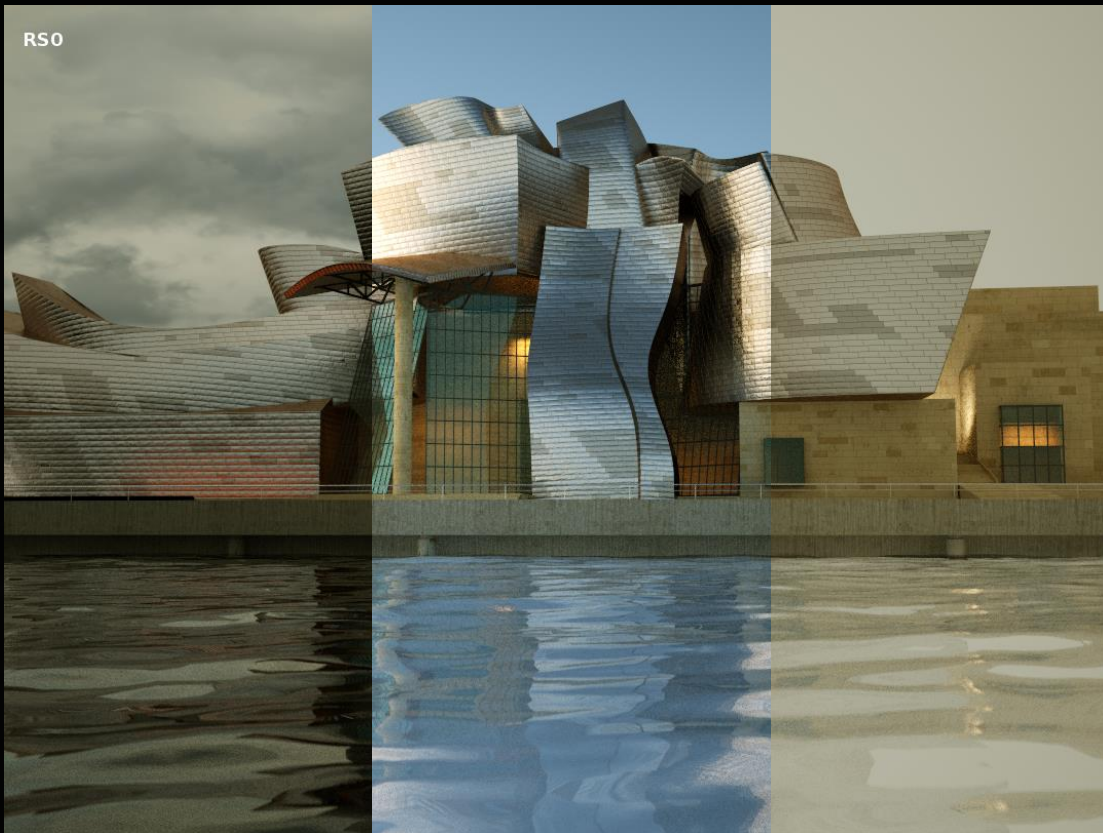


GPU



SOLVED

Guggenheim Scene



CPU == GPU 😊



Teapot Scene

CPU



GPU



Teapot Scene



Teapot Scene



ISSUES

- Subtle differences in bump/normal mapping
- Differences in materials with many layers/bsdfs
- Small changes in intensity



Use cases where CPU Maxwell helped... A LOT!!!



Test 1 : Lambert materials + Constant Sky → OK

Teapot Scene



Test 2 : Added textures + Normal maps → **WRONG**

Teapot Scene

CPU



GPU



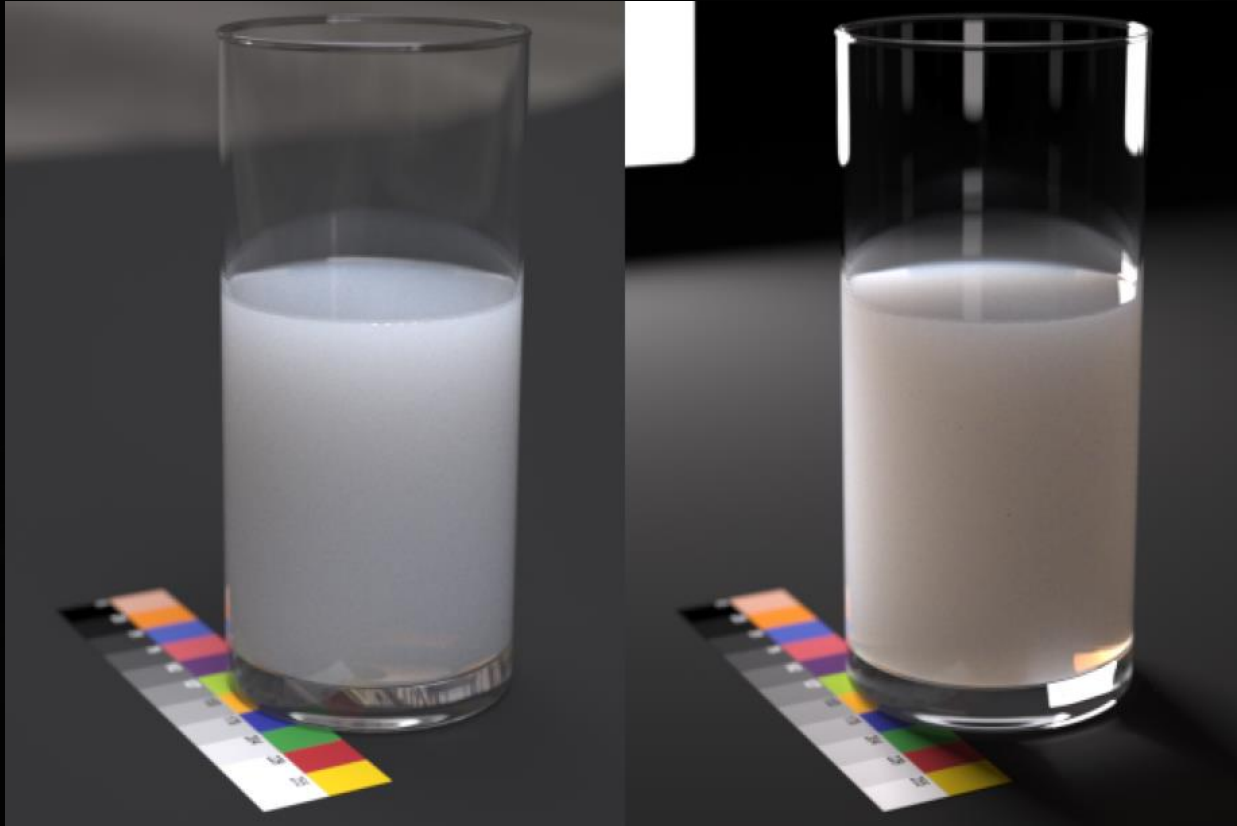
Test 3 : Added multilayered materials → WRONG

FINDINGS

- Automated CPU vs GPU numerical comparisons were **key**
 - Rays reaching IBL were not accumulating energy properly
 - Multilayered weights were not properly computed
 - Bug introduced when porting CPU **optimized** code
 - Precision issues creating TBN bases (Affected bump/normal mapping)



Next Steps Unbiased, GPU friendly SSS



Agenda

- Maxwell Render overview
- Why porting to the GPU was challenging
- Performance considerations
- Using the CPU to improve the GPU engine
- Summary

Main sources of bugs:

- CPU optimized code not easy to port
- Refactoring to make code GPU friendly
- Precision issues with some math operators



- **90% of the complexity of Maxwell already ported**
 - Very happy with the results: Speed boost: 5x-15x
 - CUDA made it possible
- **Validating using a ground truth renderer**
 - Was painful
 - 100% worth in the long run (quality first, speed second)





Thanks!

Juan Cañada
Head of Visualization
Next Limit Technologies