



# Harnessing The Power Of Multiple GPUs

Games Developer Conference 2008

Holger Gruen, [holger.gruen@amd.com](mailto:holger.gruen@amd.com)

Material:

Holger Gruen,

Ignacio Llamas, [illamas@nvidia.com](mailto:illamas@nvidia.com)

Jon Story, [jon.story@amd.com](mailto:jon.story@amd.com)

[WWW.GDCONF.COM](http://WWW.GDCONF.COM)



# Agenda

- ④ Why MGPU?
- ④ Driver Considerations
- ④ Programming for MGPU
- ④ Common Pitfalls & Solutions



Game Developers  
Conference

08

# Why MGPU?



CMP

United Business Media

[WWW.GDCONF.COM](http://WWW.GDCONF.COM)



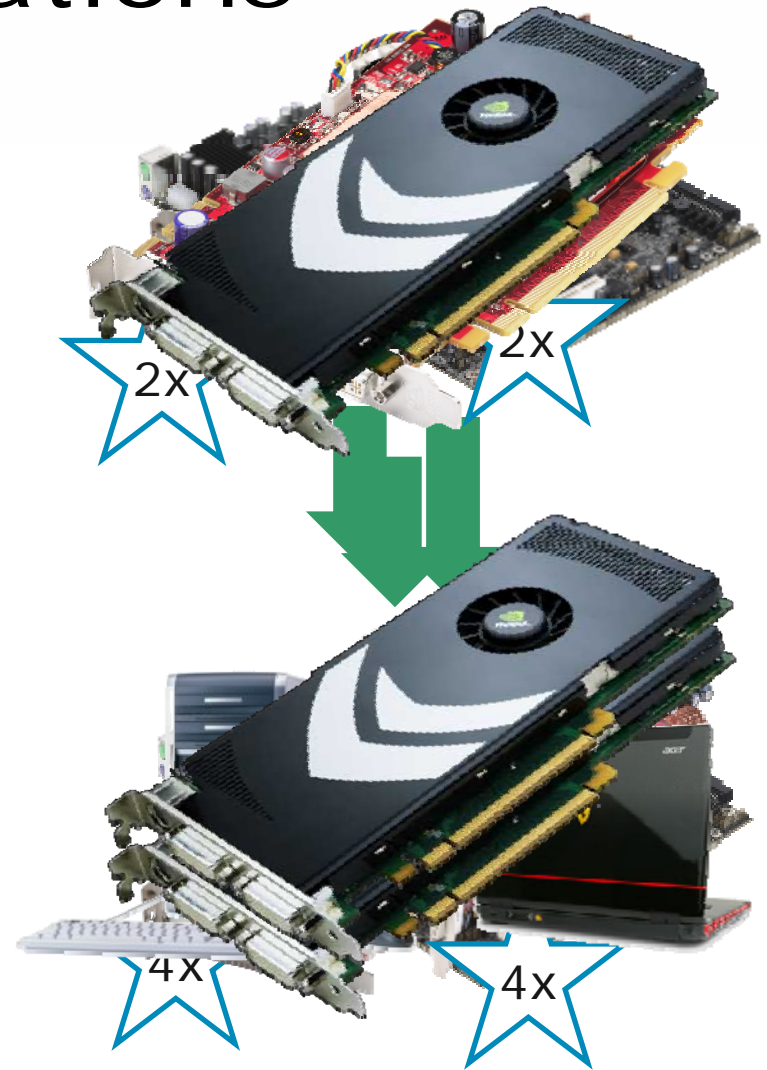
# Why MGPU?

- ③ Many apps are GPU limited at high resolutions
  - ③ Generally CPU limited at lower resolution
- ③ High res monitors have become affordable
  - ③ Consumer expectations have risen
- ③ MGPUs can dramatically increase performance
  - ③ Especially at higher screen resolutions
- ③ Next gen performance on today's HW
  - ③ Prototype your next engine



# MGPU configurations

- Multiple Boards
- Multiple GPUs per Board
- Hybrid MGPU

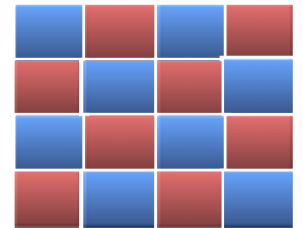




# MGPU Rendering Modes

## ⊗ SuperTiling (AMD)

- ⊗ Screen divided in a grid
- ⊗ GPUs take alternate tiles



## ⊗ Split Frame Rendering / Scissor

- ⊗ Screen is divided between GPUs
- ⊗ Dynamic load balancing



## ⊗ Alternate Frame Rendering

- ⊗ GPUs take alternate frames
- ⊗ Highest performing mode





# Driver Considerations



GameDevelopers  
Conference



# Driver modes for MGPU

- ④ AFR compatible mode
  - ④ Default – driver works around problems
  
- ④ App Profile mode
  - ④ Profile fully defines driver behaviour
  
- ④ Forced AFR - speed test mode
  - ④ AFR-FriendlyD3D.exe – no work arounds



CMP

United Business Media

[WWW.GDCONF.COM](http://WWW.GDCONF.COM)





Game Developers  
Conference

08

# Programming for MGPU



CMP

United Business Media

[WWW.GDCONF.COM](http://WWW.GDCONF.COM)



Game Developers  
Conference



# Programming for MGPU

- ⊕ MGPU no shared memory architecture  
Apps need to behave well to scale
- ⊕ Use AMD / NVIDIA libraries  
Allow to query MGPU topology
- ⊕ Know what GPU/s are rendering the  
current frame  
Critical to adapt application behavior



CMP  
United Business Media

[WWW.GDCONF.COM](http://WWW.GDCONF.COM)



# Common Pitfalls & Solutions

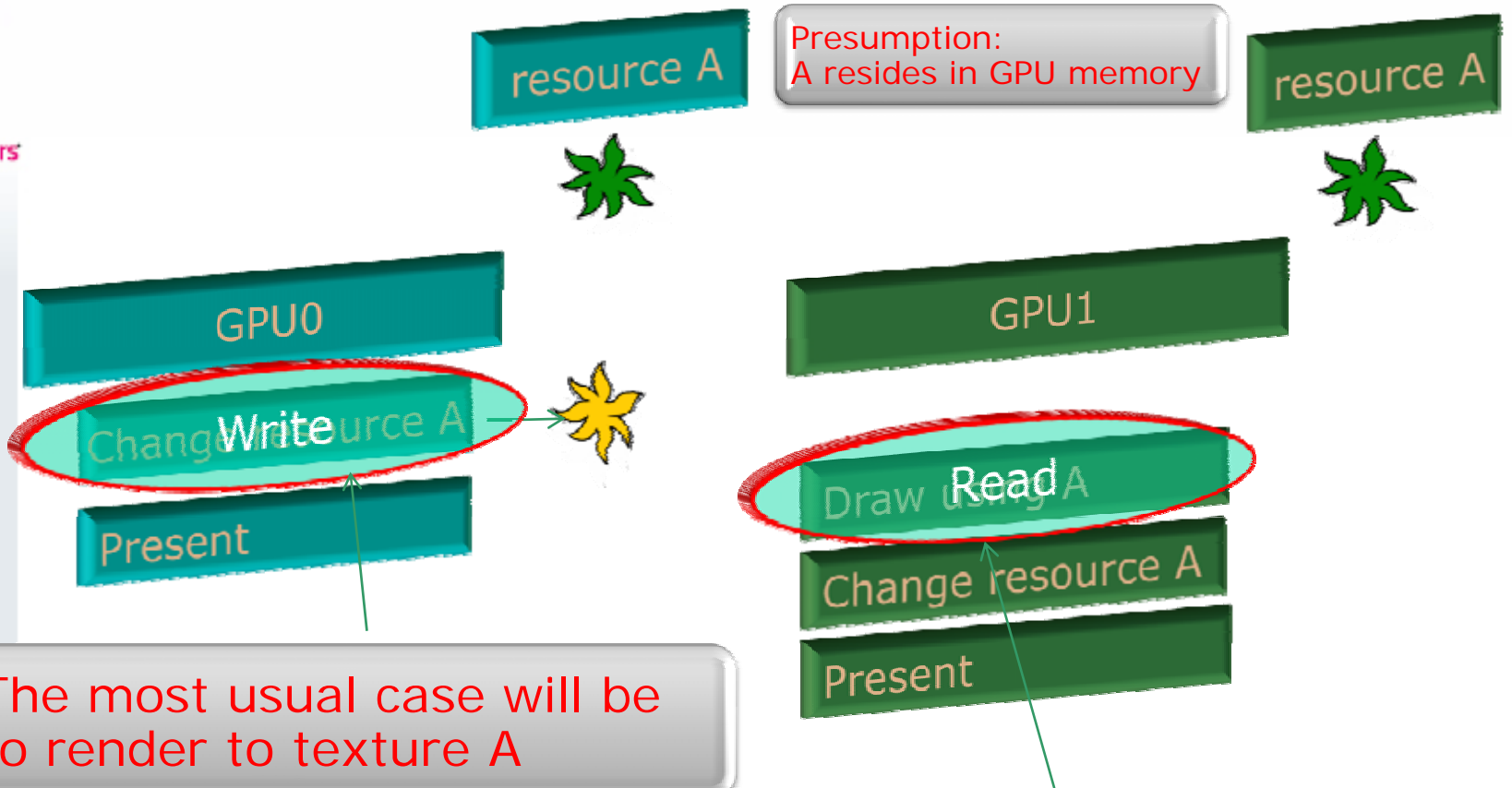


[WWW.GDCONF.COM](http://WWW.GDCONF.COM)



# Pitfalls: Dependencies between frames

Game Developers Conference 08

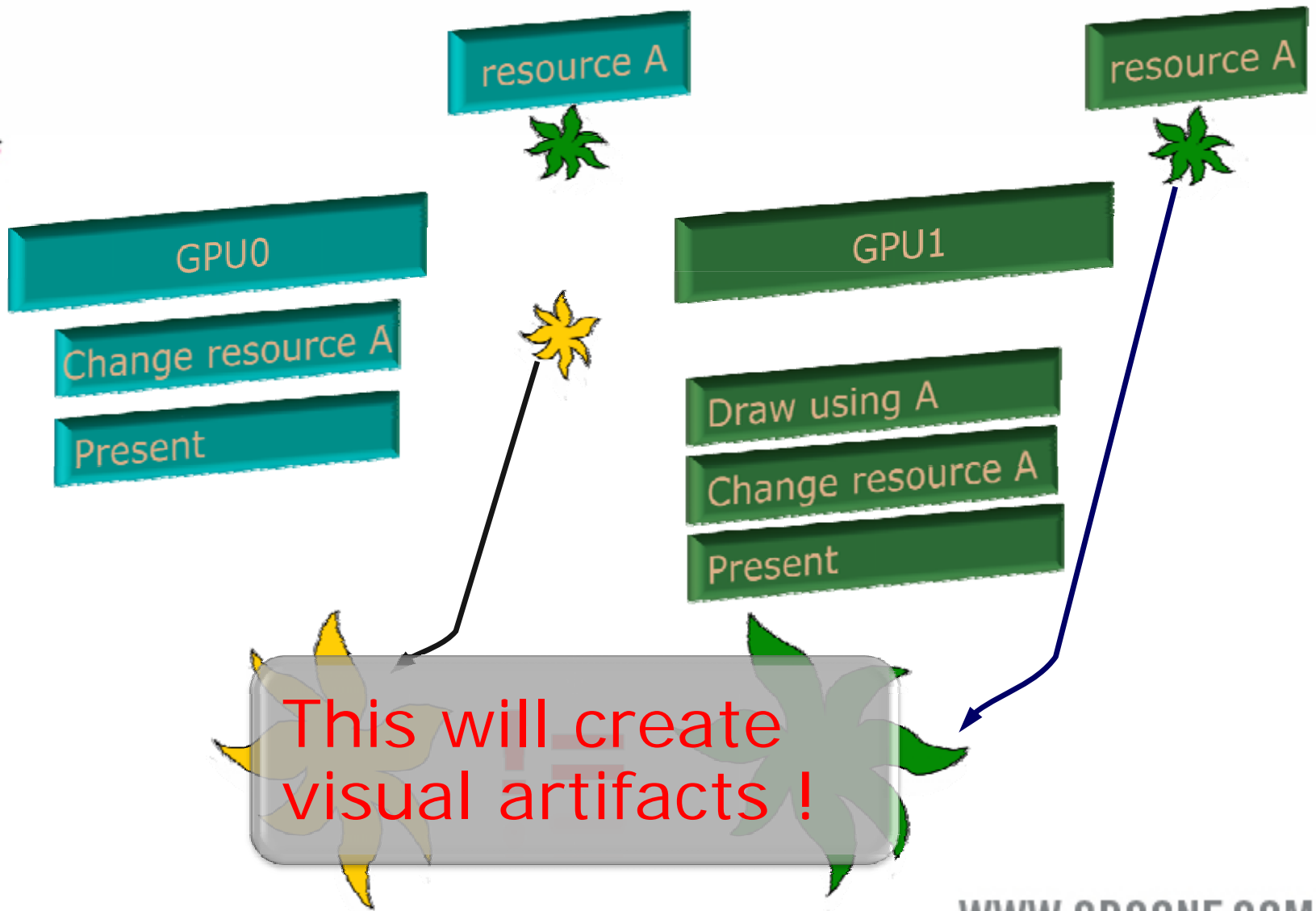


The most usual case will be to render to texture A

This is not A as it resides on GPU0



# Pitfalls: Dependencies between frames







# Solutions: Change resource every frame before using it

Good for data that changes every frame

Game Developer Conference

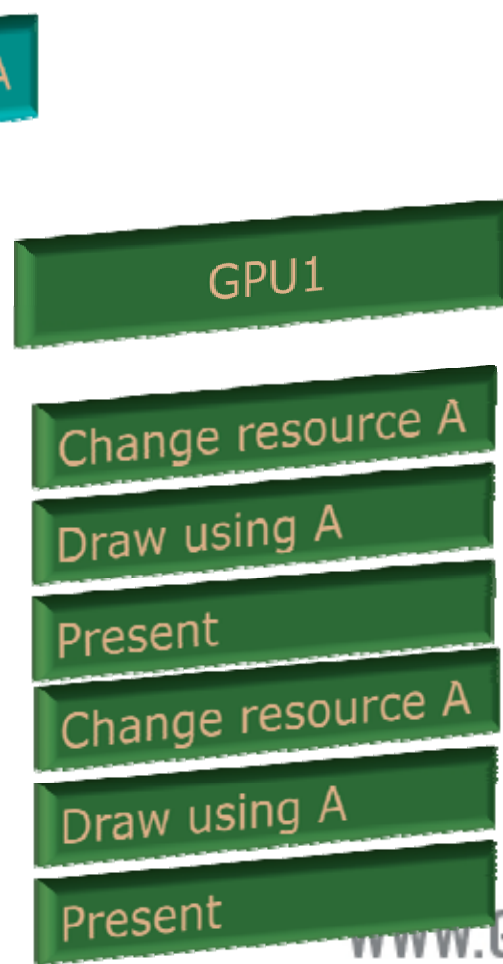
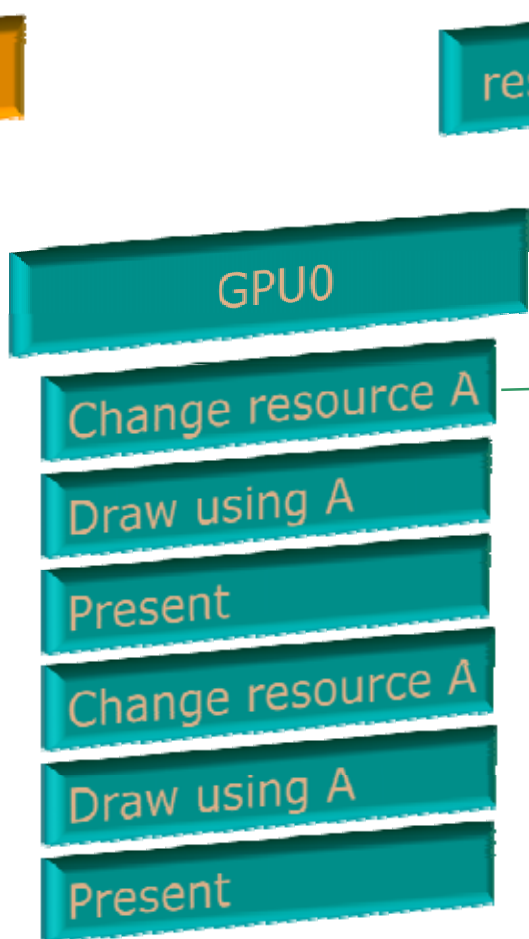
08

Driver state

Resources in sync

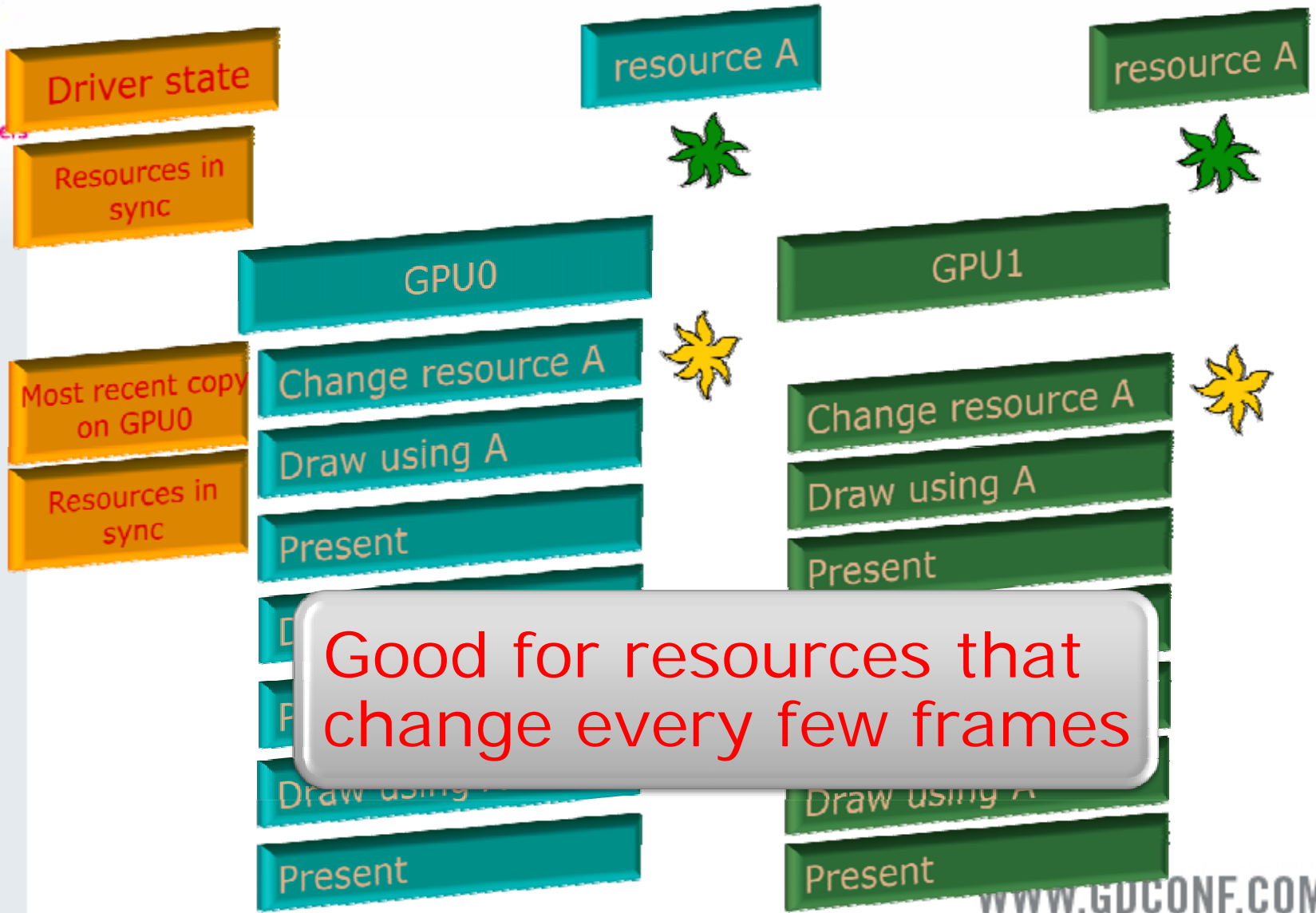
Most recent copy on GPU

Resources in sync





# Solutions: Repeat change on each GPU



Good for resources that change every few frames



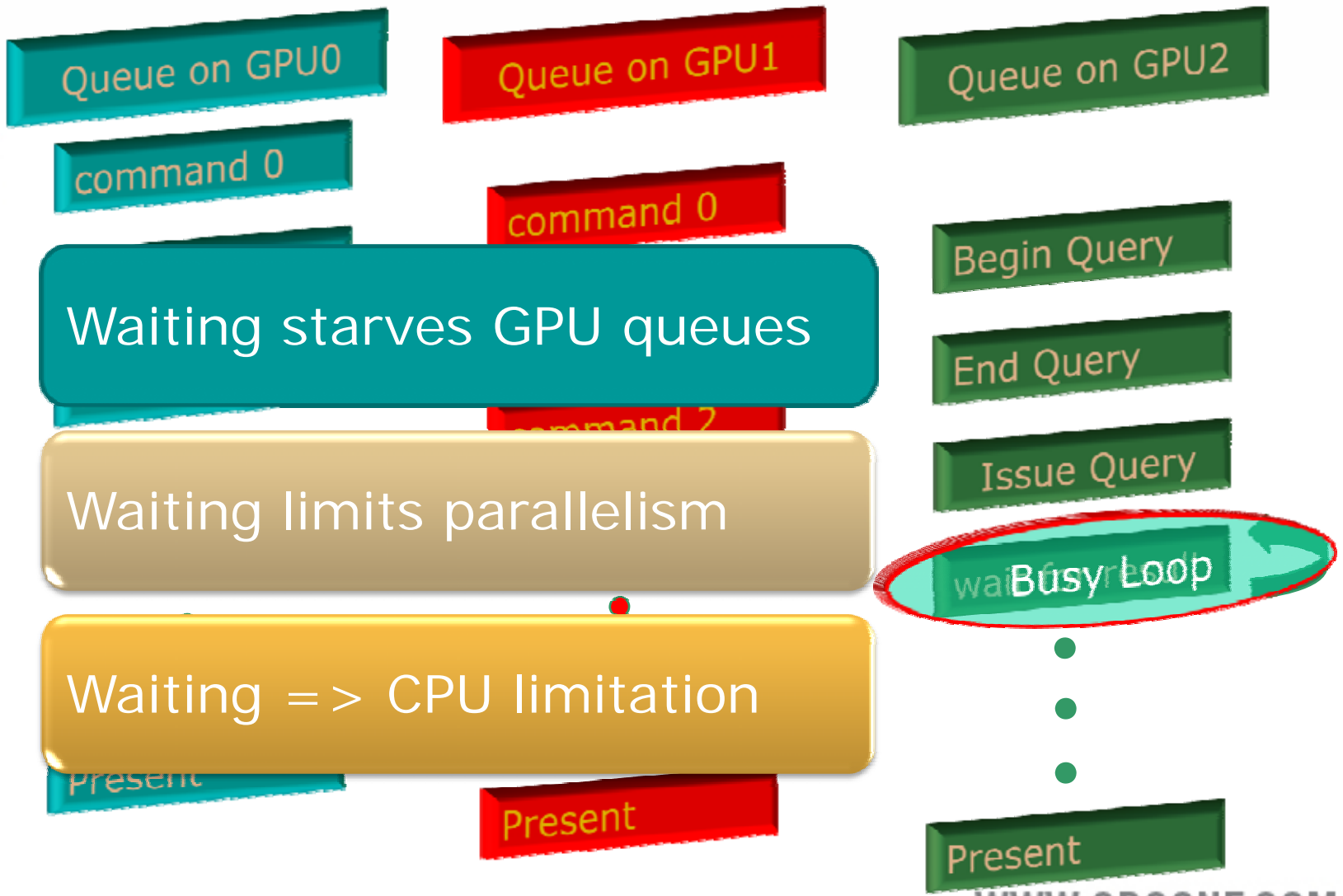


# Pitfalls: Things to watch out for under DX10

- ⊗ Drawing to vertex/index buffers
- ⊗ Stream output buffers
- ⊗ CopyResource calls
- ⊗ CopySubresourceRegion calls
- ⊗ GenerateMips calls
- ⊗ ResolveSubresource calls
- ⊗ Do not use same resource as destination of both Map(WRITE\_DISCARD) and CopyResource/CopySubresourceRegion calls



# Pitfalls: Busy waits on Queries





Game Developers  
Conference



# Solutions for Queries

- ④ Begin/End queries in same frame
- ④ Use N-GPU queries if used every frame
- ④ Expect results *starting* N-GPU frames after ending the query



CMP

United Business Media

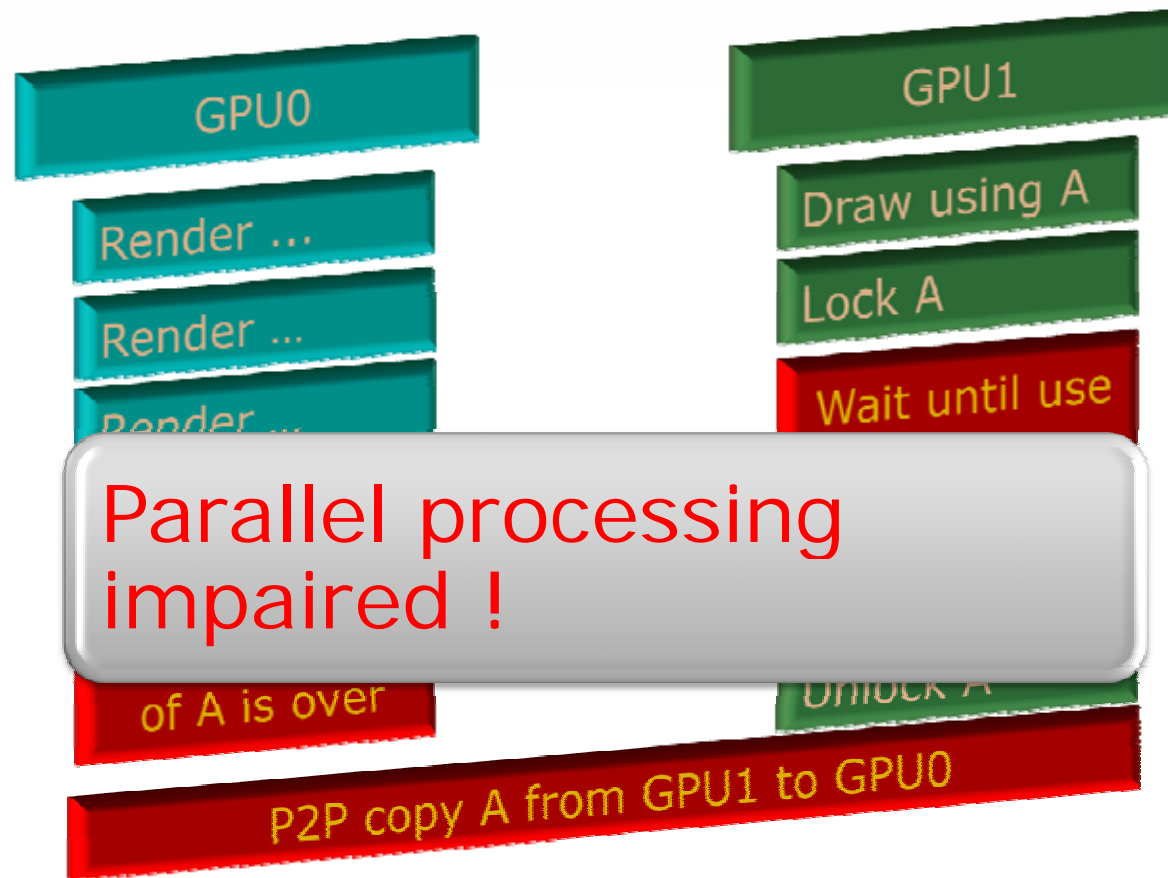
[WWW.GDCONF.COM](http://WWW.GDCONF.COM)



GameDevelopers  
Conference

08

# Pitfalls: Locks/Maps on renderable resources



Parallel processing  
impaired !

of A is over

UNLOCK A

P2P copy A from GPU1 to GPU0



CMP

United Business Media

WWW.GDCONF.COM



# Solutions/Pitfalls: Locks/Maps

## Lock/Map flip-chain or render-able resources

- ③ On DX10 call UpdateXX() (copy from STAGING resources)
- ③ On DX9 blits are always better than locks
- ③ Use dynamic hint/CPU-flags at creation
- ③ Hint discard/no-overwrite during Lock/Map
  - ③ DX10: don't use discard a lot !!

## Lock (DX9) static vertex/index buffers

- ③ Change happens only on one card => P2P copies



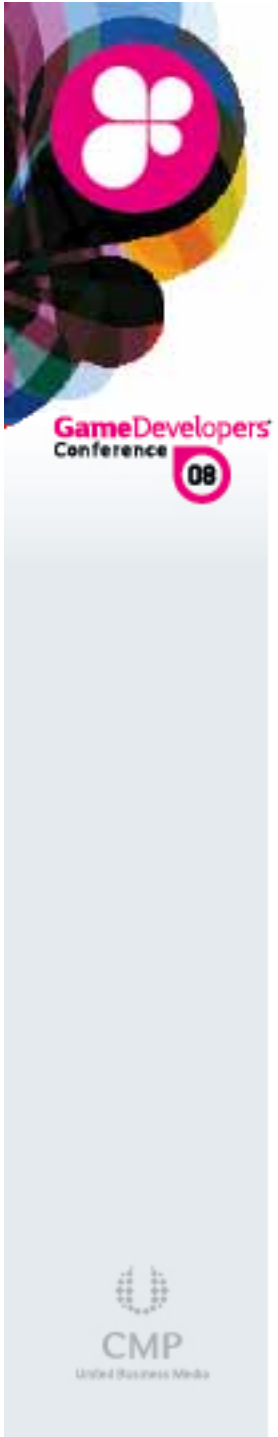
# Concluding Pitfalls & Solutions

- ⌚ Not all AFR unfriendliness causes artifacts
- ⌚ One frame old data may be OK
- ⌚ Compatible AFR mode can't detect this
- ⌚ Fixing invisible problems sacrifices perf.
- ⌚ Never use shared resources on DX10
  - ⌚ No way to detect update by other app



# Call to Action

- ④ Use AMD / NVIDIA libraries to detect MGPU topology
- ④ Write AFR friendly rendering code
- ④ Find out about your scaling
- ④ Talk to us if you do encounter problems



Questions?

[WWW.GDCONF.COM](http://WWW.GDCONF.COM)