# NVIDIA SLI AND STUTTER AVOIDANCE:

A Recipe for Smooth Gaming and Perfect Scaling with Multiple GPUs

# NVIDIA SLI AND STUTTER AVOIDANCE:

Iain Cantlay (Developer Technology Engineer)

Lars Nordskog (Developer Technology Engineer)

# WHY SLI?

- "SLI" – Set of multi-GPU technologies

- Pixel counts increasing at a staggering rate (4K+)

- Emulating the "hardware of tomorrow"

- VR – 2 eyes, 2 GPUs

# SLI BASICS

- AFR – Alternate Frame Rendering
  - One frame per GPU in parallel

- Want linear performance improvements for each GPU added
  - "SLI scaling"

- AFR SLI abstracts all non-primary GPUs away from the runtime
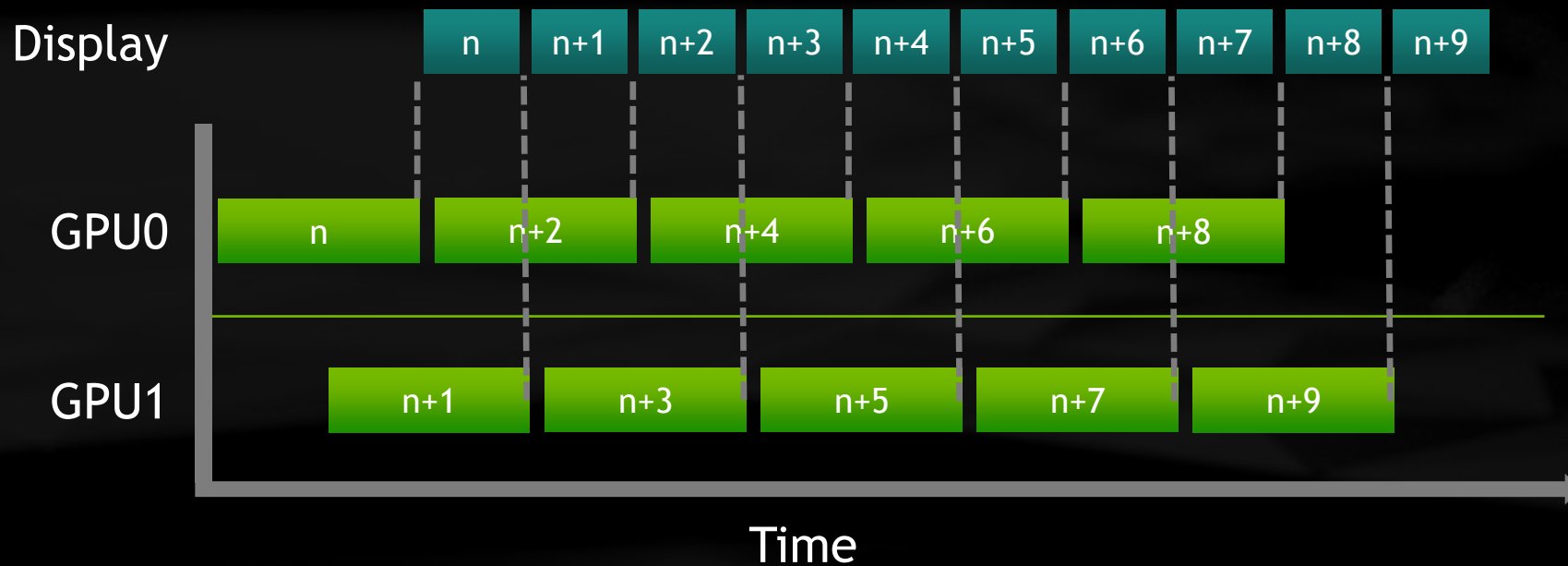  - Game sees one GPU
  - Driver does the "magic"

# SLI BASICS

- Single GPU frame rendering

# SLI BASICS

- 2-way Alternate Frame Rendering
  - Parallelism

# SLI BASICS

- Allocated resources replicated per AFR GPU

- Static GPU resource mirrored between GPUs
  - Reading from local memory is optimal
  - Static textures, IBs, VBs, etc.

- Dynamic GPU resources can diverge
  - RTs, UAVs

# SLI BASICS

## AFR Pros

- Up to linear performance scaling

- "Frame" provides natural data dependency boundary*

- Uniform workloads (frames similar)

## AFR Cons

- Non-uniform flip intervals (microstutter)

- *Interframe dependencies

- Input latency does not reduce with increased performance

# SLI BASICS

## AFR Pros

- Up to linear performance scaling

- "Frame" provides natural data dependency boundary*
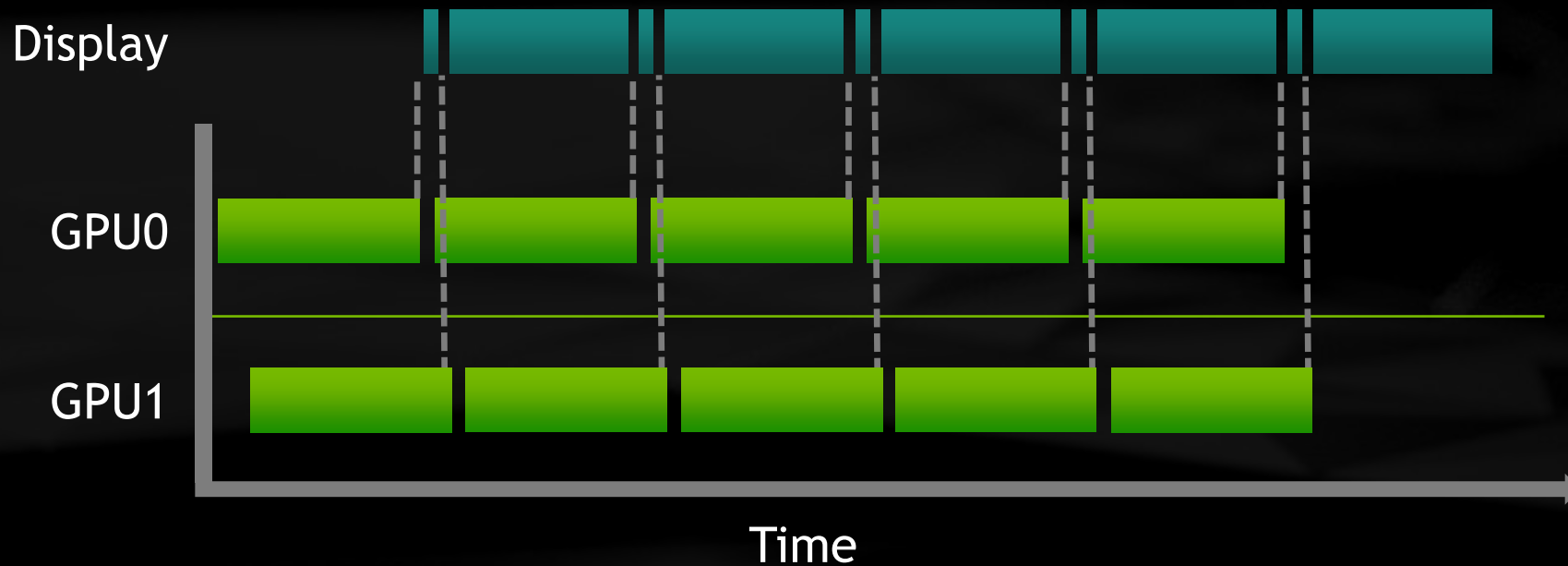
- Uniform workloads (frames similar)

## AFR Cons

- Non-uniform flip intervals (microstutter)

- *Interframe dependencies

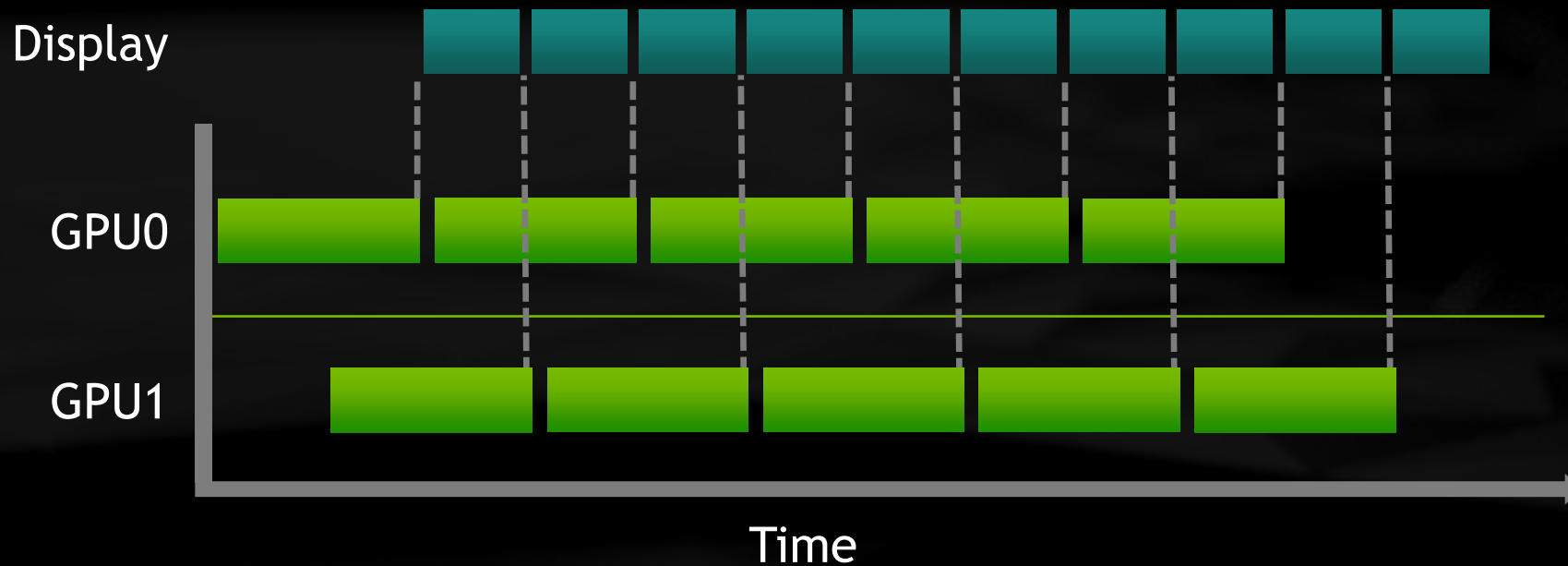- Input latency does not reduce with increased performance

# MICROSTUTTER

- Naïve parallelism -> non-uniform flip intervals...
  - Reported framerate 2x, but perceived framerate closer to single GPU

# FLIP METERING

- … but SLI driver handles frame flip metering, so you don't have to!
  - Back pressure to application avoids animation stutter

Display

GPU0

GPU1

Time

# INTERFRAME DEPENDENCIES

- We know static resources replicated to all GPUs
  - Never change, so no problem

- ...but some RTs/UAVs are modified by GPU
  - Correctness! Driver must keep RTs/UAVs in sync between GPUs
  - Sustain "illusion" of single GPU
  - Data transferred GPU->GPU when reference "dirty"

# INTERFRAME DEPENDENCIES

- Transferring data hurts SLI performance

- Some transfers not necessary
  - Game updates resource entirely each frame

- ...But other transfers are necessary
  - Techniques that need previous frame results as input
  - Temporal feedback (luminance adaptation, TXAA)
  - Compute (simulations)
  - Partial updates (tiled shadowmap, cubemap, atlas textures)

- Driver transfers entire mip slice/buffer

# INTERFRAME DEPENDENCIES

- SLI Profile skips transfers deemed unnecessary
  - Blunt instrument
  - Prioritize correctness

- NVIDIA tests, ships official SLI profile with driver
  - Profiles usually more complicated than AFR1/AFR2

# INTERFRAME DEPENDENCIES

- SLI-enabled in Control Panel without SLI profile = single GPU

- NVIDIA Control Panel AFR Modes
  - AFR1 transfers all dirty resources -> low scaling, but no corruption
  - AFR2 skips some transfers -> better scaling, but possible corruption

# SO WHERE'S THE PROBLEM?

Driver doesn't have all the information about game's intent ☹

Need game to behave well, provide hints to driver, and become "AFR-aware" for optimal performance!

# COMMON SCALING PITFALLS (CPU)

- Queries or APIs preventing queuing of frames (Bad!)
  - Solution: SLI input latency same as single, so allow n+1 frames in flight for n GPUs

- Readbacks to CPU (Dangerous!)
  - Solution: Avoid, or delay readback by n frames via buffering to avoid CPU stalling

- Stalling Map() writes (Bad!)
  - Solution: Use WRITE_DISCARD/WRITE_NO_OVERWRITE

# COMMON SCALING PITFALLS (GPU)

- Necessary transfer causing GPU->GPU serialization
  - Solution: Decouple GPUs, look back n frames on n-way config (input local)

- Mod of per-frame ping-pong buffer -> n-way dependent transfers
  - Solution: Always Discard/Clear dynamic resource before bind, QA SLI with 3-way config

- GPU-generated data not regenerated every frame
  - Solution: Regenerate data on each GPU, or hint to keep

# INITIAL STEPS

- Renaming EXE to AFR-FriendlyD3D.exe
  - Enables n-way AFR, skips **all** transfers
  - Corruption, but ideal for checking "speed of light"
  - No scaling with rename -> CPU-GPU serialization or CPU-boundedness


- Query NVAPI to detect SLI via number of GPUs
  - NvAPI_D3D_GetCurrentSLIState()
  - SLI profile "aware"… no profile returns "1 GPU"

# SLI COMPATIBILITY PROCESS

1. Think through your interframe dependent effects/systems

2. Run with exe renamed to "AFR-FriendlyD3D.exe" to skip all GPU->GPU transfers of dirty resources

3. Test thoroughly, looking for corruption

4. Address resources that are not "AFR Friendly"
   - Regenerate data for all GPUs, or hint to driver that data must persist
   - Hint what data can be discarded

# TAKING ACTION

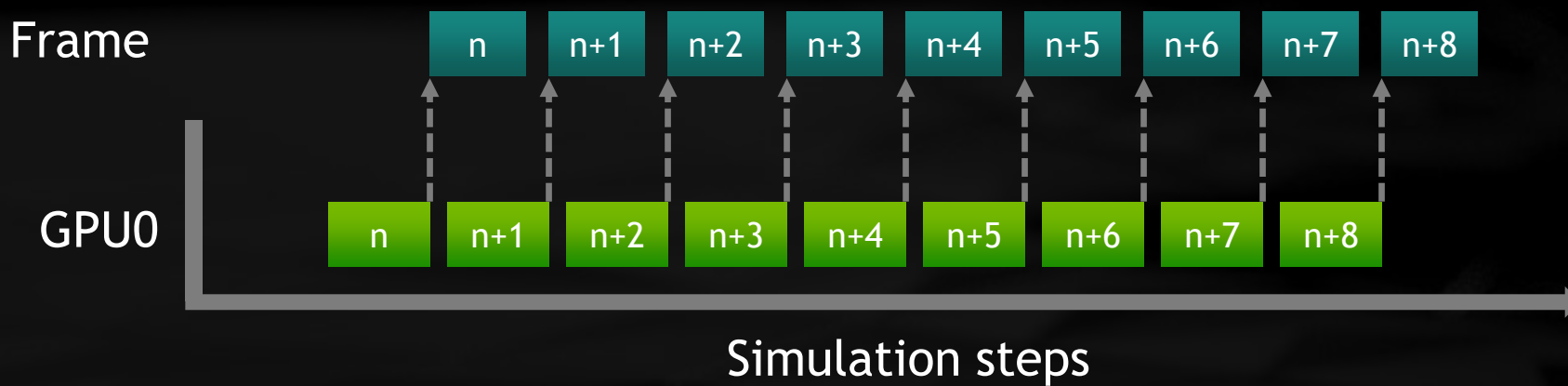*"How do I regenerate data for each GPU?"*

- "Explicit synchronization"
  - Keep track of which GPUs receive updates
  - Re-issue for each "dirty" GPU
  - Allows discarding of transfers + GPU coherency


- Regenerate work or hint to driver to keep?
  - Generally regenerating better, but case by case
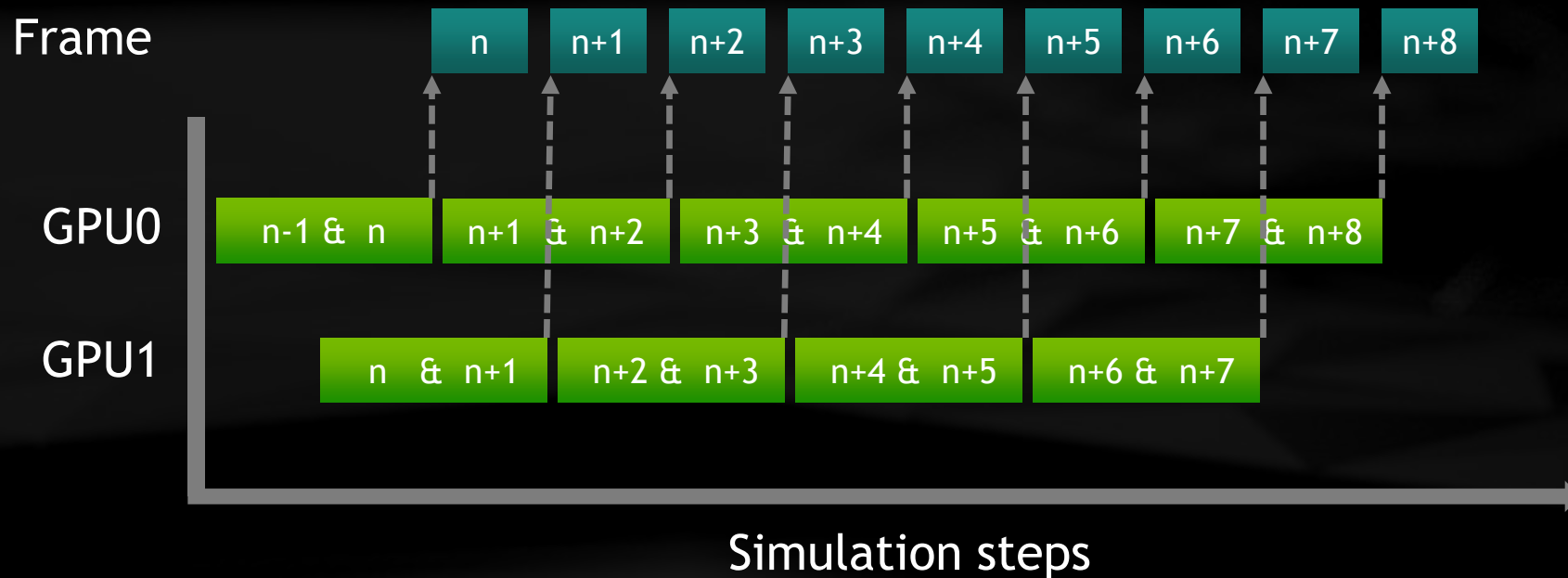  - Only so much data practically transferred per frame (performance)

# TAKING ACTION

*"How do I regenerate data for each GPU?" (cont)*



Frame

| n | n+1 | n+2 | n+3 | n+4 | n+5 | n+6 | n+7 | n+8 |

GPU0

| n | n+1 | n+2 | n+3 | n+4 | n+5 | n+6 | n+7 | n+8 |

Simulation steps

# TAKING ACTION

*"How do I regenerate data for each GPU?" (cont)*

- Simulation is duplicated for each GPU
- Still faster than doing a transfer!



| Frame | n | n+1 | n+2 | n+3 | n+4 | n+5 | n+6 | n+7 | n+8 |
|-------|---|-----|-----|-----|-----|-----|-----|-----|-----|

GPU0: n-1 & n | n+1 & n+2 | n+3 & n+4 | n+5 & n+6 | n+7 & n+8

GPU1: n & n+1 | n+2 & n+3 | n+4 & n+5 | n+6 & n+7

Simulation steps

# TAKING ACTION

*"How do I hint what I \*DO\* need to persist between frames?"*

- NvAPI_D3D_BeginResourceRendering()/NvAPI_D3D_EndResourceRendering()
  - Wrap update -> driver transfers early/efficiently
  - NVAPI_D3D_RR_FLAG_FORCE_DISCARD_CONTENT works as Discard/Clear for ping-pong case

- Begin/End assume only next GPU needs data
  - NVAPI_D3D_RR_FLAG_MULTI_FRAME if used for multiple frames

- USE WITH CAUTION!!!
  - Final update of resource in frame
  - Don't Begin/End > 1 time per frame, per resource
  - Begin/End takes precedence over profile
  - Entire resource transferred

# TAKING ACTION

*"How do I hint what I \*DON'T\* need to persist between frames?"*

- ID3D11DeviceContext1::DiscardView()/DiscardResource()
  - Ideal solution
  - Before bind in current frame
  - Only supported in DX11.1

- ID3D11DeviceContext::Clear*()
  - Before bind in current frame

- NvAPI_D3D_SetResourceHint()
  - Driver excludes resource from SLI "dirty" state tracking (never transfers)
  - Sticky through allocation lifetime

# TAKEAWAYS

- SLI excellent for substantially increasing GPU performance

- Ensure AFR friendly CPU behavior

- Use AFR-FriendlyD3D.exe

- Anticipate interframe dependent effects/systems
  - Design them to be AFR friendly

- At minimum, focus on regenerating data or hinting what to keep between frames
  - BeginResourceRendering/EndResourceRendering hints
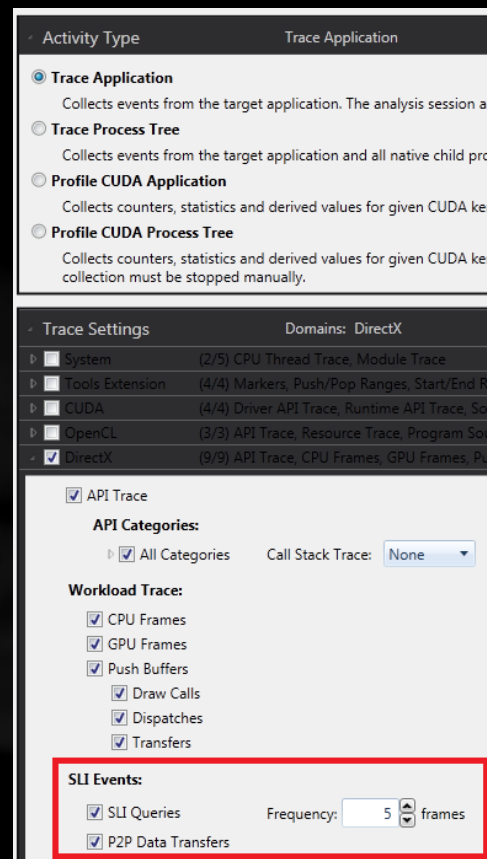  - NVIDIA can remove the rest with profiles, but Discard APIs better

# OH YEAH...

- Getting testing builds to NVIDIA early (Please ☺)
  - For SLI profiling, identification issues, advice


- QA SLI on 3-way configuration
  - Needs profile or AFR-FriendlyD3D.exe to scale

# OTHER RESOURCES

- Nsight SLILog
  - Plans to expand functionality and clarity

- GPUView

# QUESTIONS?

- Thank you!


- Iain Cantlay (icantlay@nvidia.com)
- Lars Nordskog (lnordskog@nvidia.com)

# GAMEWORKS

- Get the latest information for developers from NVIDIA and continue the discussion
- gameworks.nvidia.com