# AN EVOLUTION OF MOBILE GRAPHICS

**Michael C. Shebanow**

**Vice President, Advanced Processor Lab**
**Samsung Electronics**

July 20, 2013

# DISCLAIMER

- The views herein are my own

- They do not represent Samsung's vision nor product plans

- The Mobile Market

- Review of GPU Tech

- GPU Efficiency
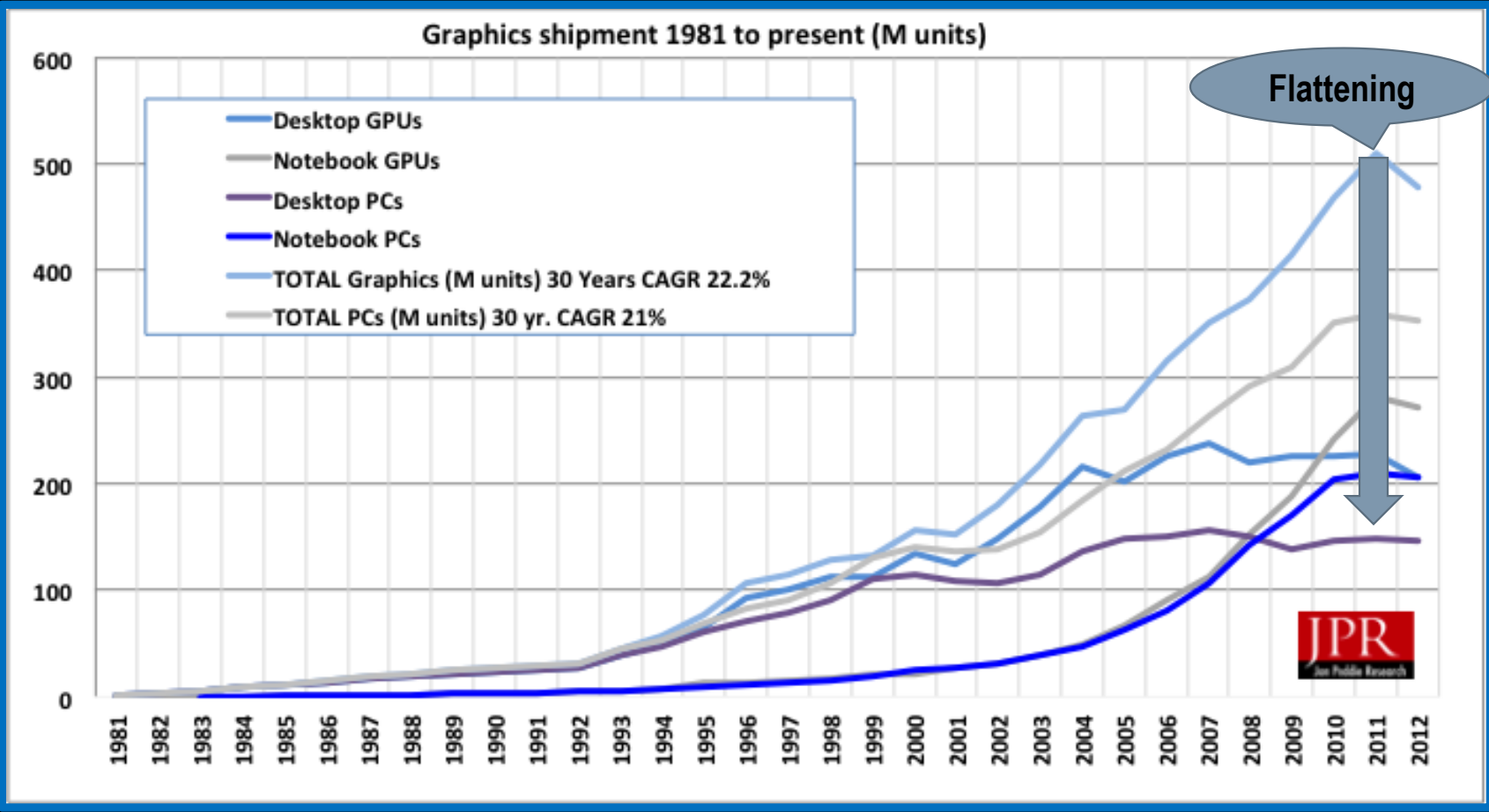
- User Experience

- Tech Challenges

- Summary

AGENDA
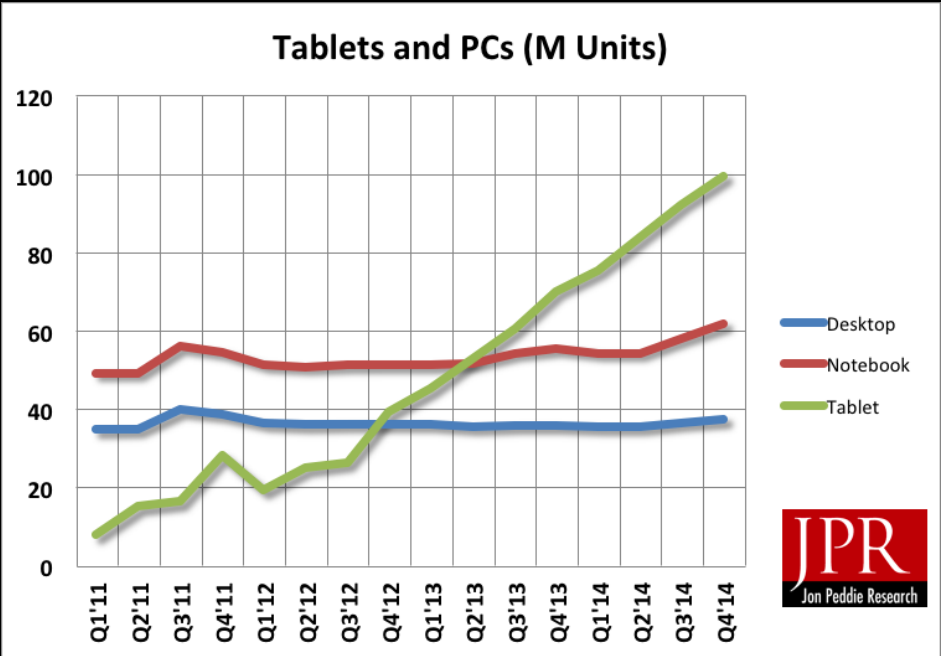
The Rise of the Mobile GPU & Connectivity

# A NEW WORLD COMING?

# DISCRETE GPU MARKET



Graphics shipment 1981 to present (M units)

Legend:
- Desktop GPUs
- Notebook GPUs
- Desktop PCs
- Notebook PCs
- TOTAL Graphics (M units) 30 Years CAGR 22.2%
- TOTAL PCs (M units) 30 yr. CAGR 21%

Flattening

JPR
Jon Peddie Research

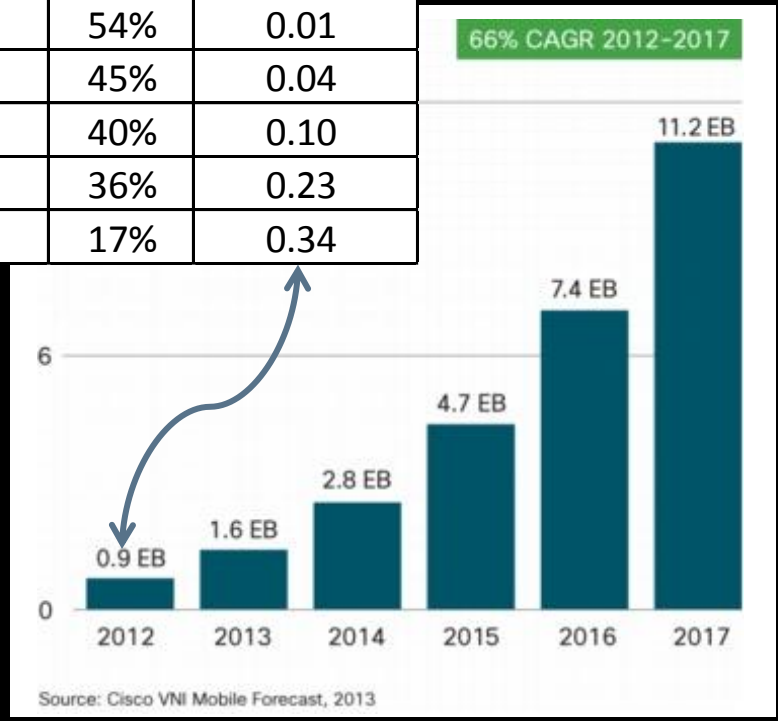# MOBILE GPU MARKET

- In 2012, an estimated 800+ million mobile GPUs shipped
    - ~123M tablets
    - ~712M smart phones
- Will easily exceed 1B in the coming years

- **Trend:**
    - **Discrete GPU relatively flat**
    - **Mobile is growing rapidly**

"Phablets"

Smart Phones

Tablets

**Tablets and PCs (M Units)**

Desktop
Notebook
Tablet

JPR
Jon Peddie Research

# WW INTERNET TRAFFIC

- Source: Cisco VNI

- Internet traffic growth rate is staggering

  - **2012** total traffic is **13.7 GB per person per month**

  - **2012** smart phone traffic at **0.342 GB per person per month**

  - **2017** smart phone traffic expected at **2.7 GB per person per month**

| Year | IP Traffic (TB/sec) | growth rate | Mobile INET Traffic (TB/sec) |
|---|---|---|---|
| 2005 | 0.9 | | 0.00 |
| 2006 | 1.5 | 65% | 0.00 |
| 2007 | 2.5 | 61% | 0.01 |
| 2008 | 3.8 | 54% | 0.01 |
| 2009 | 5.6 | 45% | 0.04 |
| 2010 | 7.8 | 40% | 0.10 |
| 2011 | 10.6 | 36% | 0.23 |
| 2012 | 12.4 | 17% | 0.34 |

66% CAGR 2012-2017

11.2 EB — 2017
7.4 EB — 2016
4.7 EB — 2015
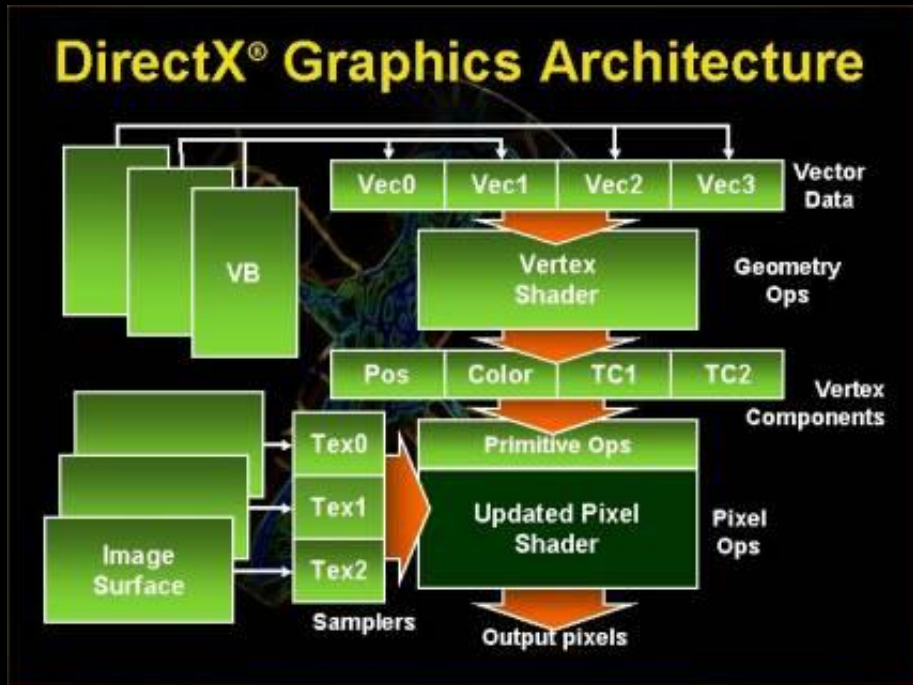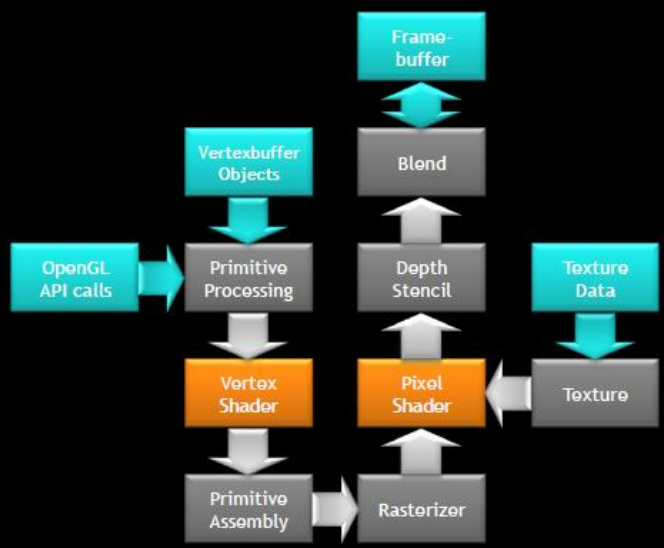2.8 EB — 2014
1.6 EB — 2013
0.9 EB — 2012

Source: Cisco VNI Mobile Forecast, 2013

# WHERE ARE WE HEADED?…

- Enormous quantity of GPUs

- Large amount of interconnectivity

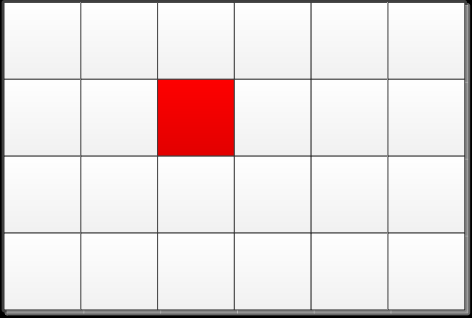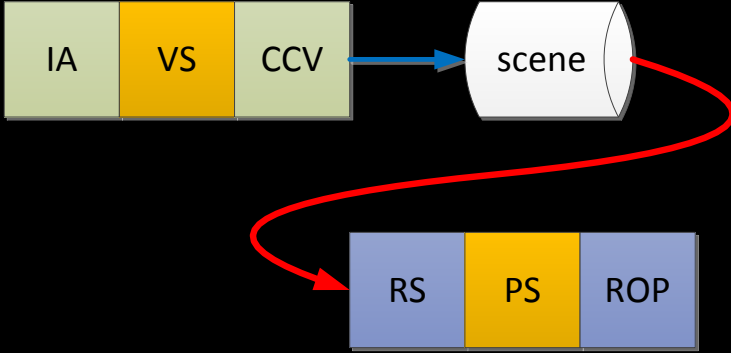- Better I/O

*GPU Pipelines*

# A BRIEF REVIEW OF GPU TECH

# MOBILE GPU PIPELINE ARCHITECTURES
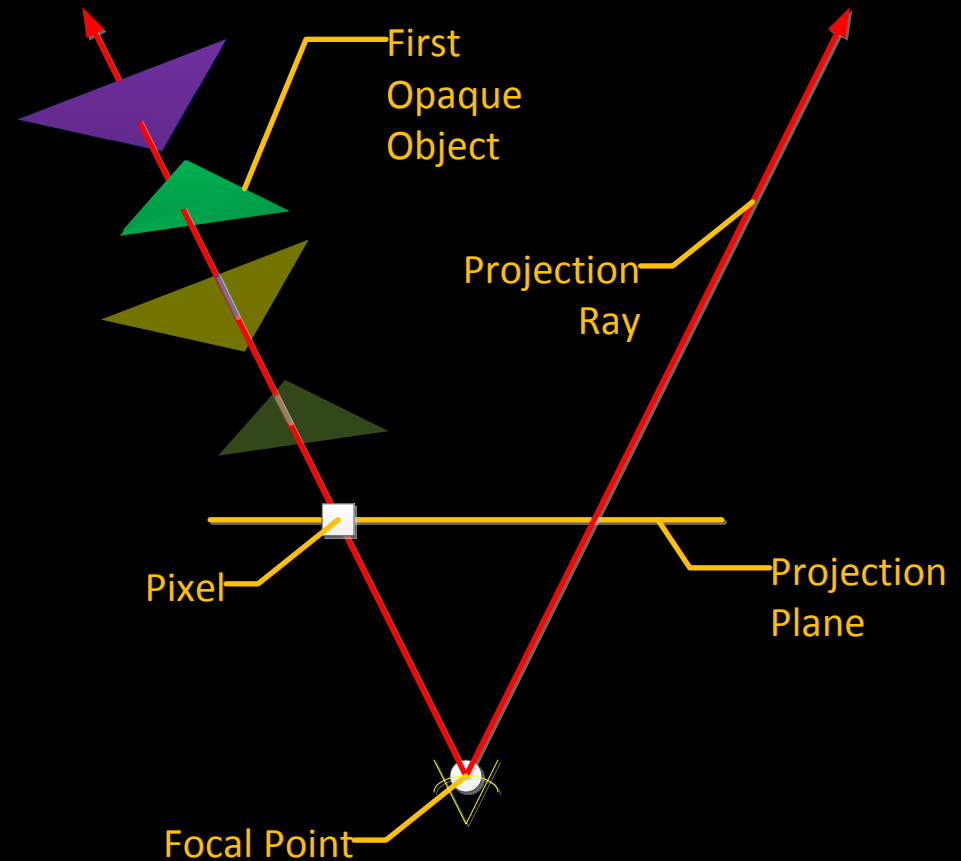
*Tile-based immediate mode rendering (TBIMR)*
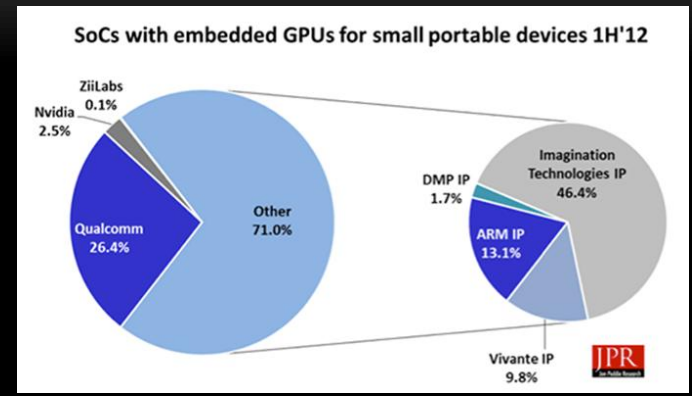
*Tile-based deferred rendering (TBDR)*

| IA | VS | CCV | RS | PS | ROP |
|----|----|----|----|----|----|

| IA | VS | CCV | → scene |
|----|----|----|----|

| RS | PS | ROP |
|----|----|----|

*IA = input assembler*
*VS = vertex shader*
*CCV = cull, clip, viewport transform*
*RS = rasterization, setup*
*PS = pixel shader*
*ROP = raster operations (blend)*

# TBDR W/ HSR

- ## HSR = *hidden surface removal*

  - ### Sort all objects across each projection ray

    - #### Use tiling to reduce data set size

  - ### Only nearest opaque and closer transparent objects need to be drawn

  - ### Remaining fragments can be killed => not drawn

First Opaque Object

Projection Ray

Pixel

Projection Plane

Focal Point

# MOBILE GPU LANDSCAPE

SoCs with embedded GPUs for small portable devices 1H'12

| Company | Product | Pipeline | Notes |
|---------|---------|----------|-------|
| ARM | Mali | TBIMR | Unified shader, 2-4 math pipes per core |
| Imagination | PowerVR | TBDR/HSR | Latest is Rogue (S6). Unified shader. DX11 support |
| Qualcomm | Adreno | FlexRender | Unified shader. "FlexRender" = automatic switching between direct render (IMR) and tile-based deferred rendering (TBDR). |
| NVIDIA | Tegra | TBDR & TBIMR | Evolution:<br>• Tegra 1/2/3/4: non-unified TBDR architecture<br>• Logan: Kepler-based GPU, TBIMR<br>• Parker: Maxwell-based GPU, TBIMR |
| Vivante | ScalarMorphic | IMR | Unified Shader. |
| Intel | Gen \| Atom | IMR \| PowerVR | Market leader in integrated graphics.<br>Atom-based devices using Imagination PowerVR |
| AMD | Radeon | IMR | Hondo/Temash pipes. |

# A PATH TO A BETTER MOBILE GPU?
# [PART 1]

# WHAT IS IMPORTANT?

- More with less

- Better user experience

# POWER EFFICIENCY

- Performance = power efficiency

- Two types of efficiency:

  - "perf@watts":

    *The ability to deliver maximum performance*
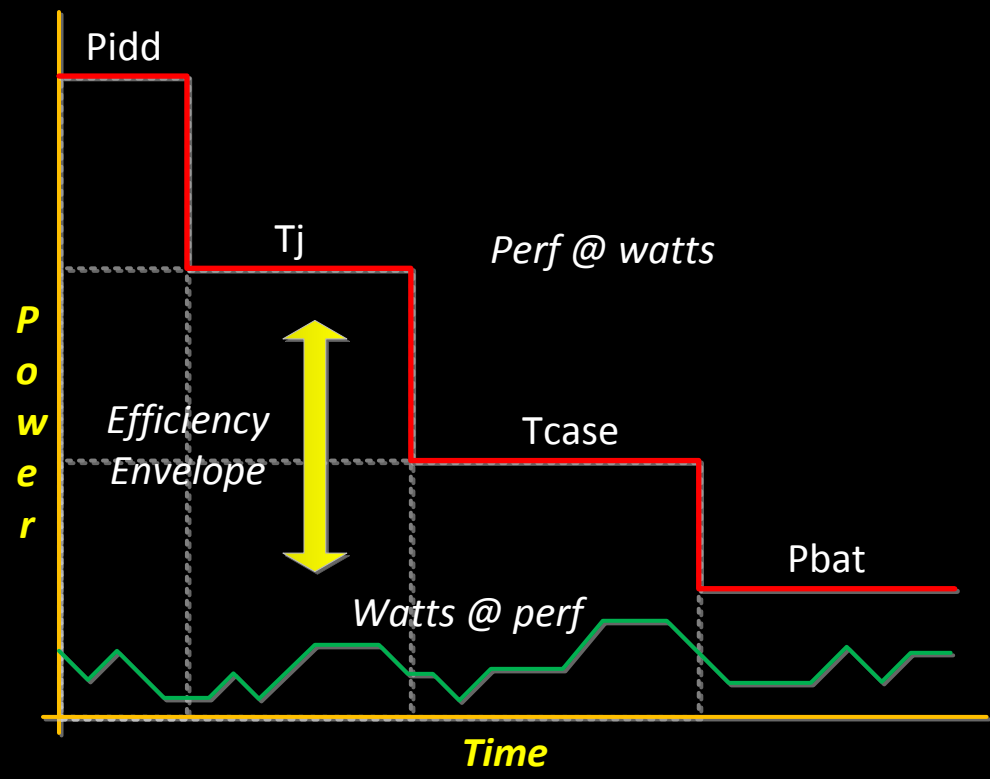
  - "watts@perf":

    *The ability to deliver maximum battery life at a minimum required performance*
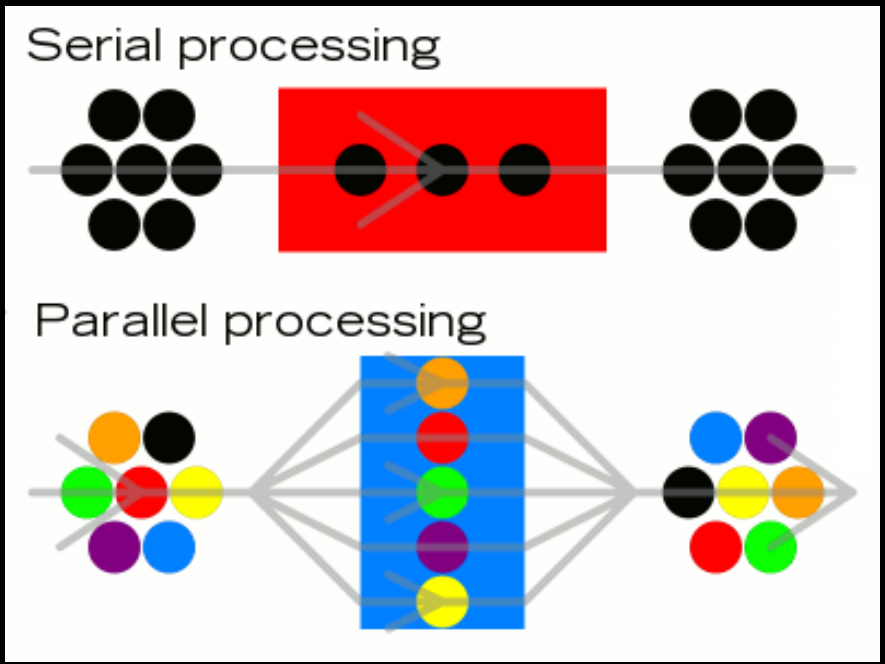
# WHAT IS EFFICIENCY?

- ## Perf @ Watts

  - *Maximum performance at some power limit*

  - Limits:

    - electrical (Pidd)

    - die temp (Tj)

    - skin temp (Tcase)

    - battery life (Pbat)

- ## Watts @ Perf

  - *Minimum power at constant performance*

  - Example: deliver 60 frames/sec at lowest power

# PARALLELISM

- Parallel vs. Sequential
  - Parallel → independence
  - Sequential → dependence
- Three fundamental forms of parallelism
  - Spatial: executing operations between threads at the same time
  - Temporal: executing operations between threads at the same place
  - ILP: executing operations from within the same thread in parallel
- Fundamental differences between ILP-only machines and massive TLP-ILP machines
  - CPUs vs. GPUs



Serial processing

Parallel processing

# THROUGHPUT VS. LATENCY

- Throughput = rate at which operations complete

- Latency = time it takes to complete an operation or set of operations

- CPUs versus GPUs

    - In CPUs, the primary objective is low latency

    - In GPUs, the primary objective is high throughput

- CPUs versus GPUs

    - In an application suitable for CPUs, we assume a low degree of TLP

    - In an application suitable for GPUs, we assume a high degree of TLP
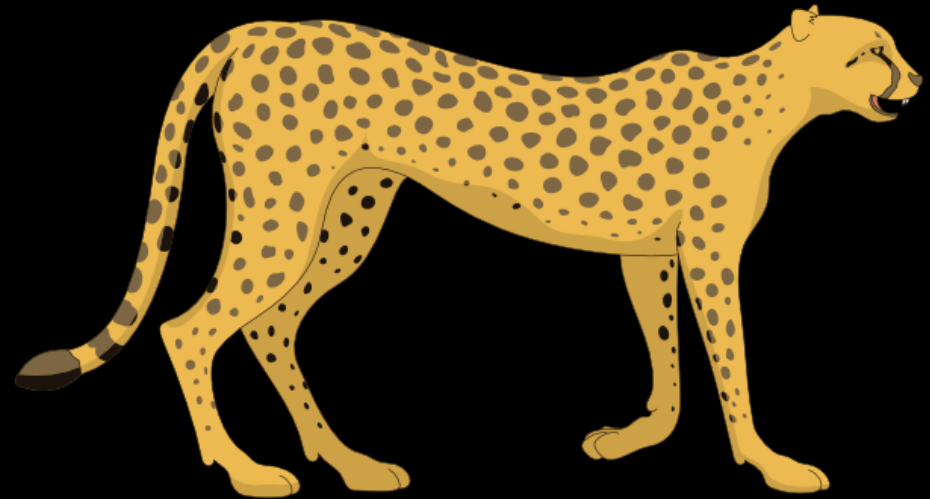
# GPU PERFORMANCE

- Supply and demand:

$$\vec{S} \geq \lambda \vec{D}$$

("limiter equation")

- Lambda ($\lambda$) is throughput
- Supply examples:
  - FP BW (flops/clock)
  - Texture BW (quads/clock)
  - Memory BW (bytes/clock)
- Demand density examples:
  - FP ops per shader
  - Sample ops per shader

# ENERGY REDUCTION TECHNIQUES

- Work Reduction

- Memory Avoidance

- Memory BW Reduction

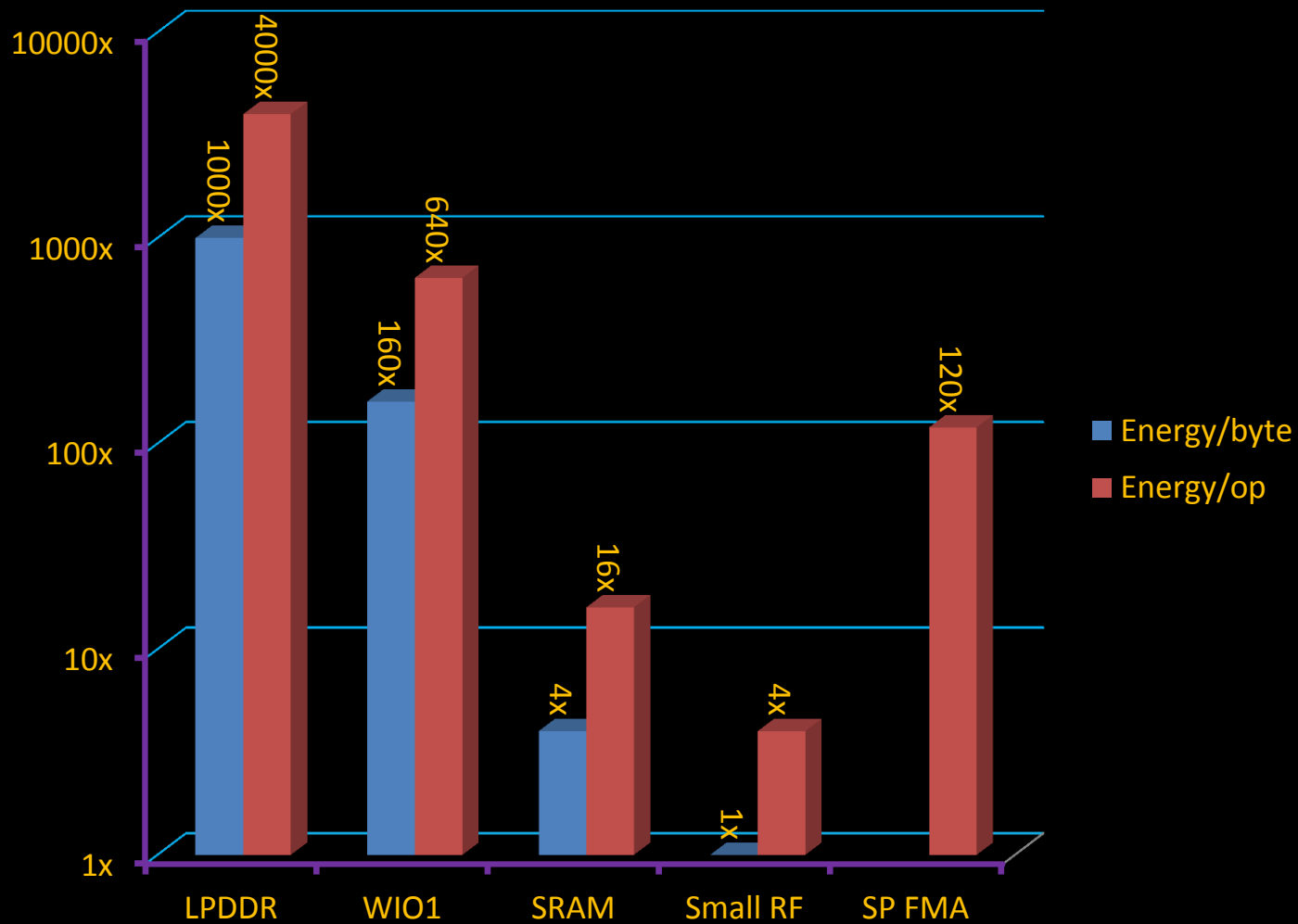- Memory Access Management

# WORK REDUCTION

- Pixel shaders in ES games ~95% of the shader load

    - A pixel shader killed is raw power savings

    - HSR can kill 30-50% of the shader threads

- Geometry in DX11 a problem

    - Unigine Heaven ~10M Tri/frame

- Inter-frame work reduction?
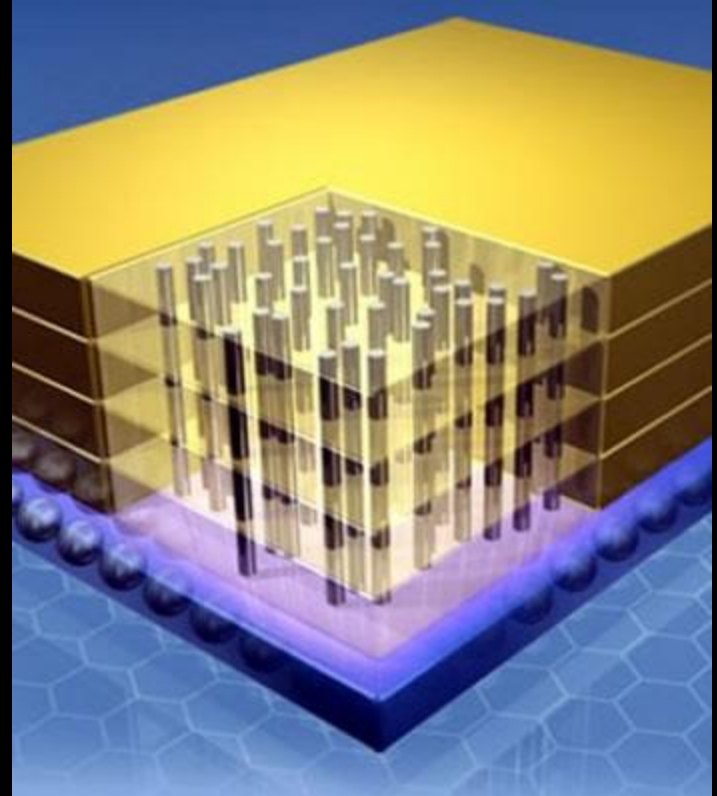
# RELATIVE ACCESS ENERGY COSTS

# MEMORY AVOIDANCE

- Memory power a problem

  - LPDDR ~150 pJ/byte
    *(150 mW @ 1 GB/sec)*

  - WIO1 ~24 pJ/byte
    *(24 mW @ 1 GB/sec)*

  - On-chip SRAM ~0.6 pJ/byte
    *(0.6 mW @ 1 GB/sec)*

- Reduction in working set for non-essential traffic (i.e., not texture, attribute, command, or render target)

  - Rematerialize? (computation vs. BW)

  - Scheduling to reduce lifetimes?

# MEMORY BW REDUCTION

- Texture compression (RD)
    - Better compression?
    - Tessellation use of textures?

- Tile compression (WT)
    - TB-based signature checking
    - Lossless compression

- Attribute compression (RD)
    - Reduce stream BW

# MEMORY ACCESS MANAGEMENT

- SOC memory architecture

  - Blood rivals (antagonists)

  - Effect of CPU/GPU traffic on Memory Controller (MC)

    - Intelligent page open/close management

    - Balance latency vs. BW

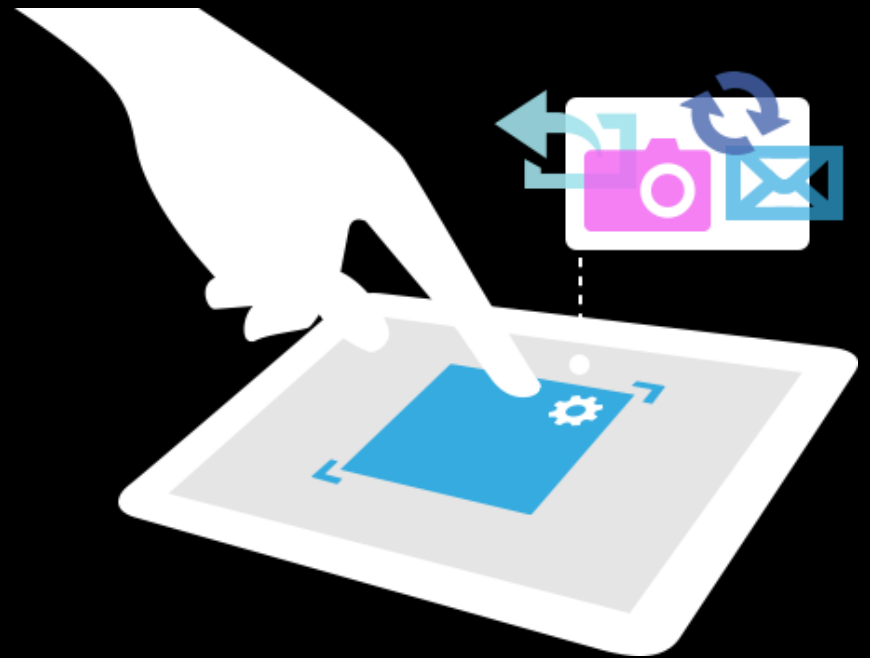- Mismanaging DRAM results in both performance loss AND extra energy – double whammy

4 cores

240 cores

A better user experience…

# A PATH TO A BETTER MOBILE GPU?
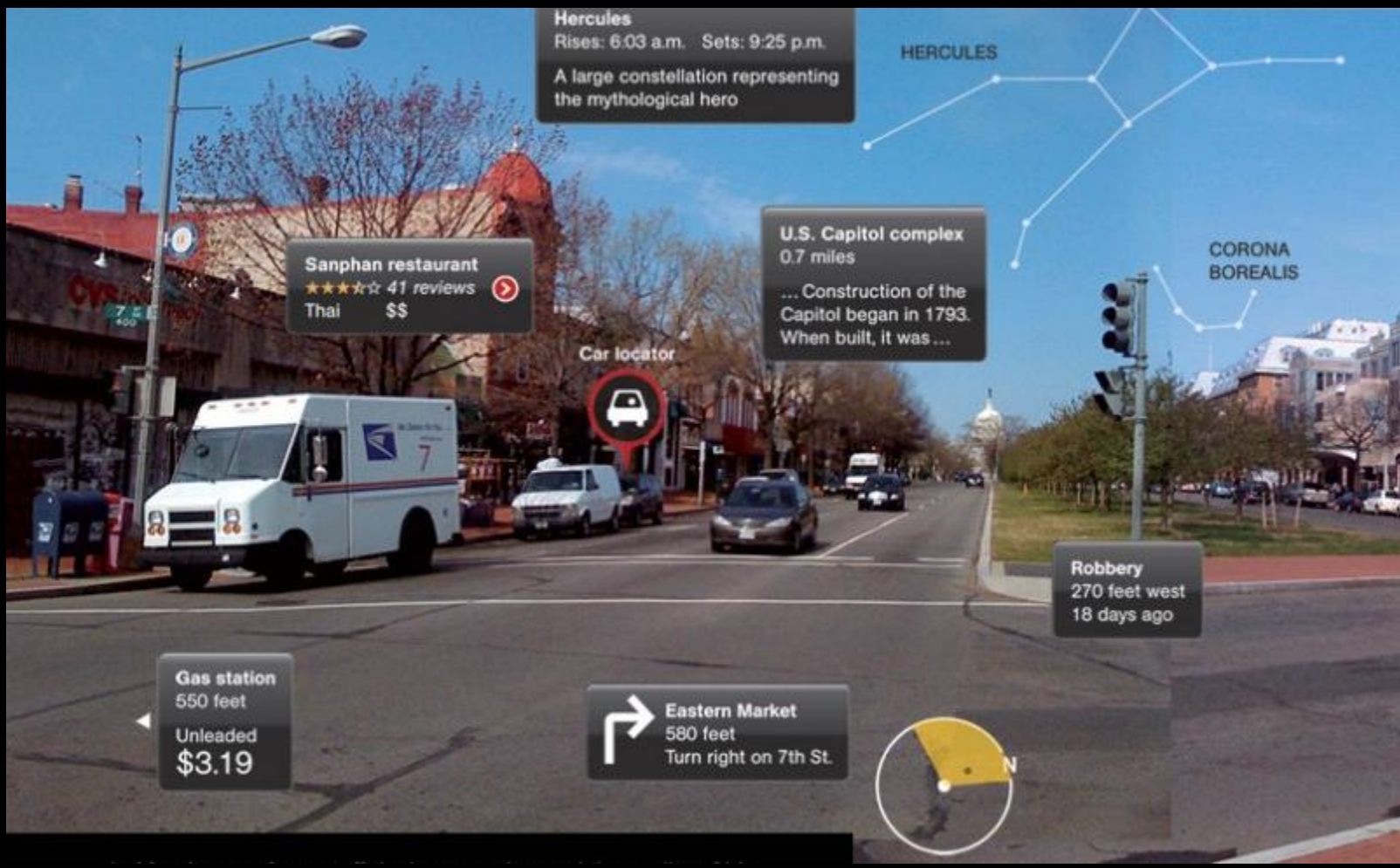# [PART 2]

# USER EXPERIENCE (UX)

- ## User Experience = perception of device:

  - ### Functionality

  - ### Integration into every day life

  - ### Ease of use (intuitive)

*ISO 9241-210[1] defines user experience as "a person's perceptions and responses that result from the use or anticipated use of a product, system or service". - Wikipedia*

# APPLICATION: NAVIGATION

# APPLICATION: FACE RECOGNITION

# APPLICATION: TELEPRESENCE

"General-Purpose Telepresence with Head-Worn Optical See-Through Displays and Projector-Based Lighting." by Maimone A., Yang, X., Dierk, N., State, A., Dou, M., and Fuchs, H. , IEEE Virtual Reality 2013

# APPLICATION: VIRTUAL COMPUTER

# THE UX OPPORTUNITY

- Killer apps will be integration of:
    - AR/MR technology
    - Big Data operations

- Subject to:
    - Real-time constraints
    - Parallelization on a massive scale

*Making a better UX*

# FUTURE MOBILE TECH CHALLENGES?

# KEY CHALLENGES

- I/O:

  - AR Headsets

  - Environment Imaging

- Computational:

  - API Improvements

  - Cloud-device integration

# AR HEADSETS

- Google Glass is pretty cool, but…

- Better imaging

  - Stereo/Light field

  - HD → UHD

  - Speed

- More sensors

- Wireless power?

- Fashion/ubiquity
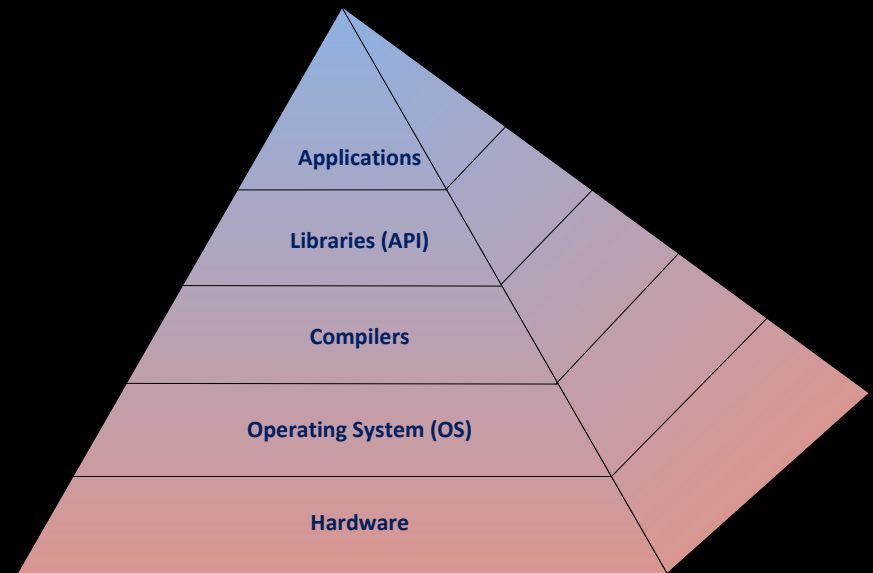
# ENVIRONMENT IMAGING

- For telepresence, headset camera is insufficient

- Need "environment cameras"

- Lots of privacy concerns

- Localizing environment to a client?

# API IMPROVEMENTS

- Today's APIs are power inefficient

- Needed:
  - Hints
  - State-less rendering
    - API commands supply state with action
  - Frame-less rendering
    - Compositing deferred and on-demand
  - Hierarchical geometry
    - Deferred detail

# CLOUD-DEVICE INTEGRATION

- SW Challenge:

  - Making cloud queries easier

  - Utilizing the parallelism of the cloud

- Ultimate challenge:

  - The "network GPU"

  - Analogously extend the GPU model to network scale

  - $10^9$ GPUs $\rightarrow$ $10^{21}$ FLOPs?

# SUMMARY

- Mobile computing, in particular graphics, is growing rapidly and becoming ubiquitous

- Tomorrow's machines:
    - Ever improving efficiency
    - Integrated visual UX
    - Tied to the cloud

- Challenges remain to make this a reality

- Exciting prospects…