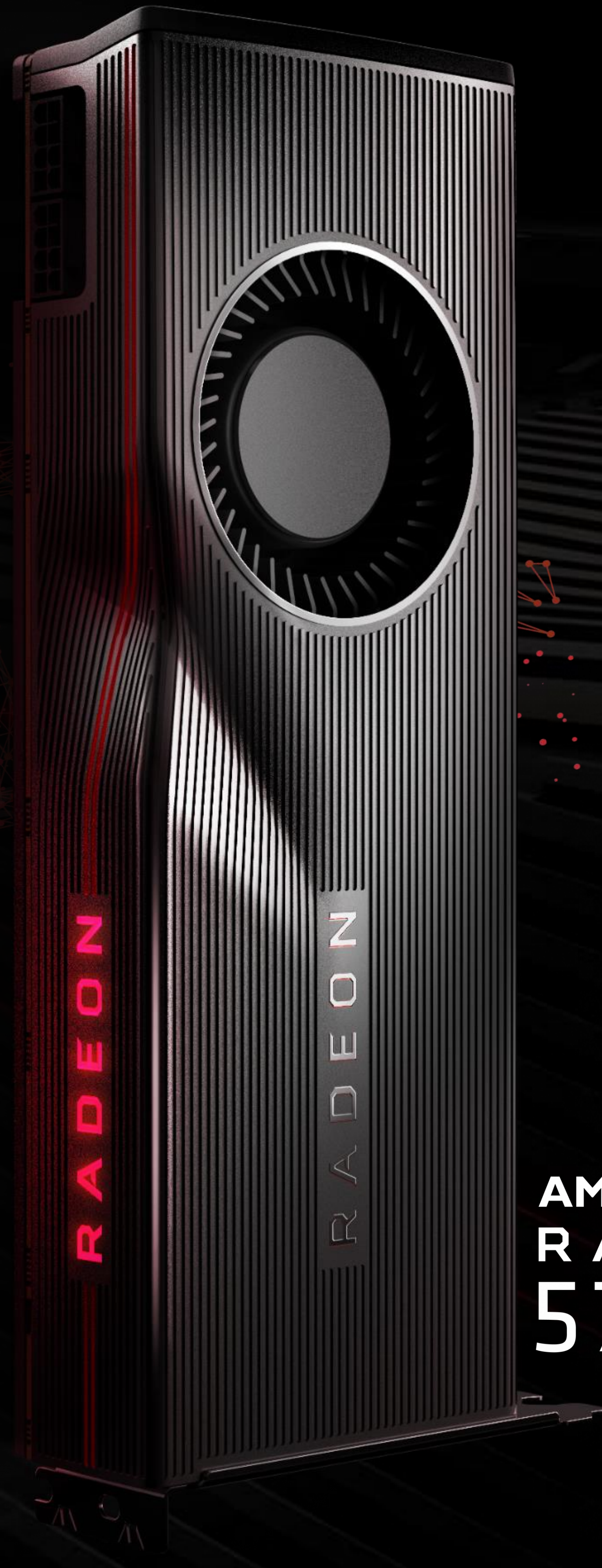


# RDNA ARCHITECTURE



7

nm

251

sqmm

10.3

Billion Transistors

GDDR

6

PCIe®

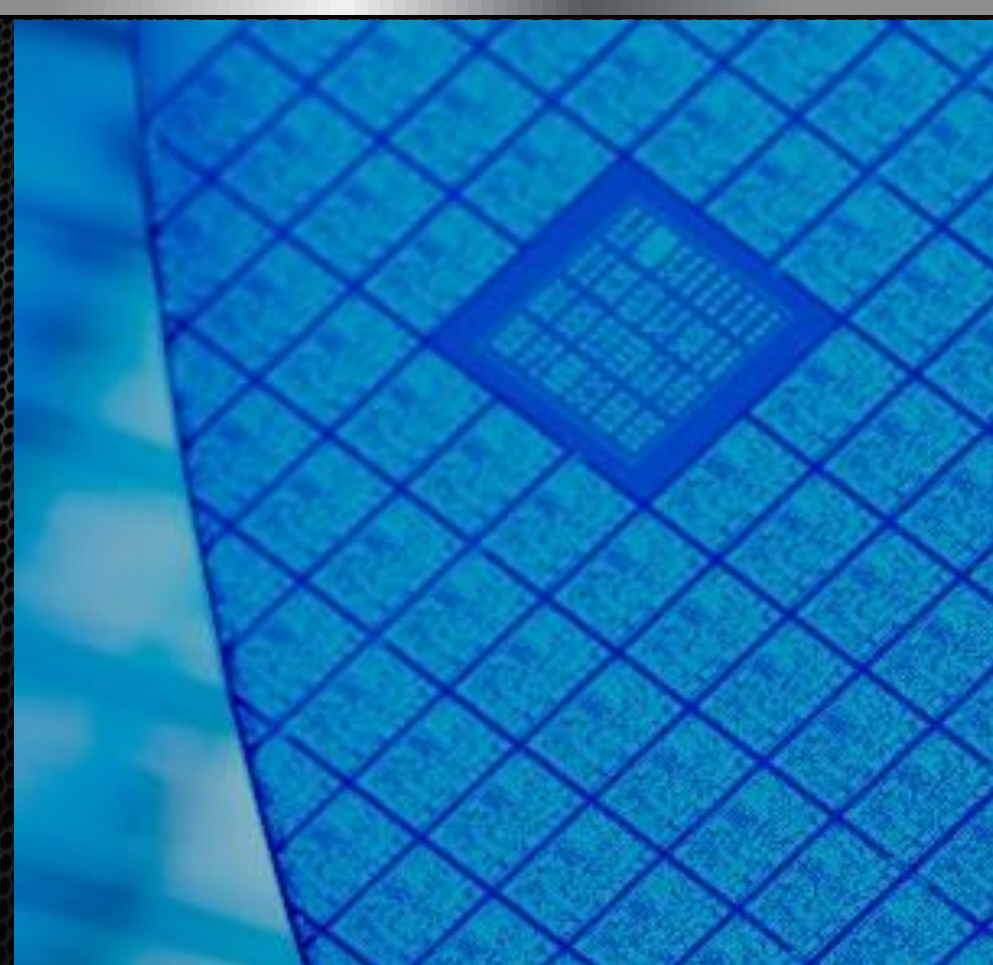
4.0 Support

AMD  
RADEON **RX**  
5700 XT



# “NAVI”

## KEY TECHNOLOGY INFLECTIONS



**PROCESS**

**7NM**

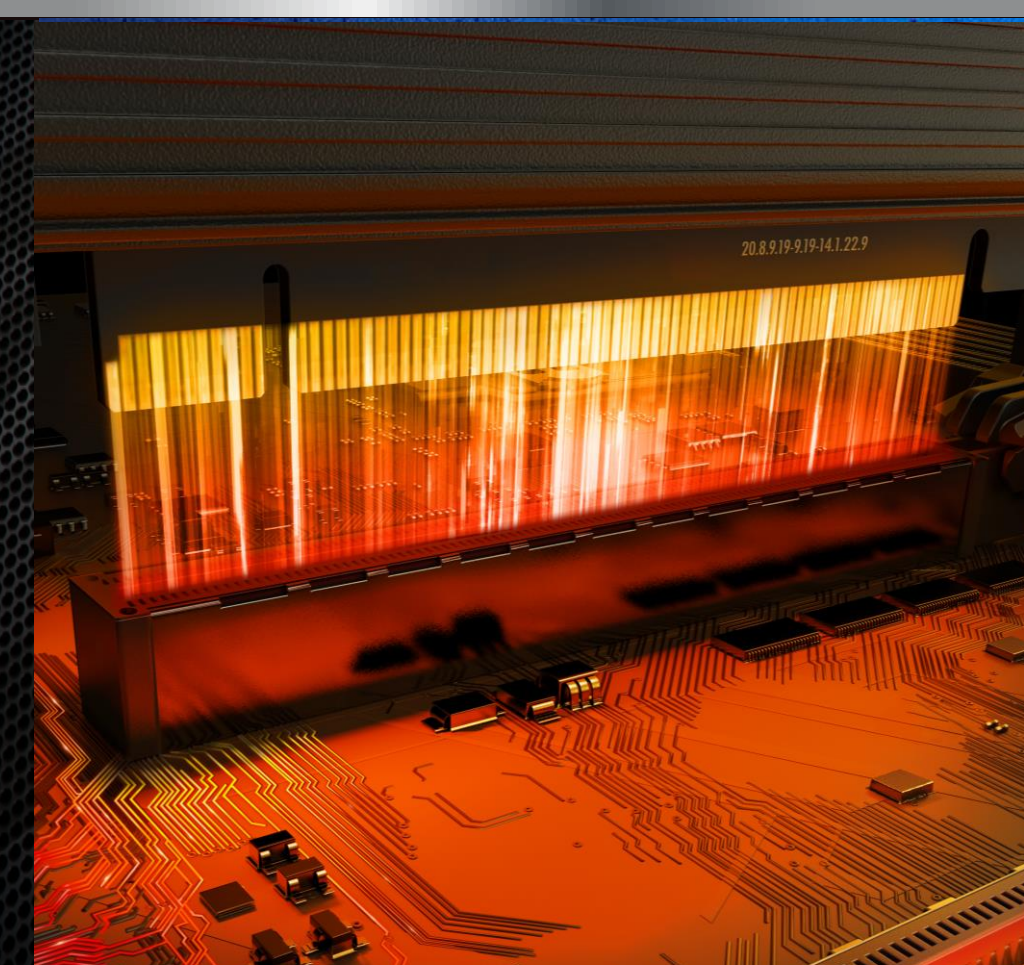
Faster, Smaller, Lower Power Transistors



**DRAM**

**GDDR6**

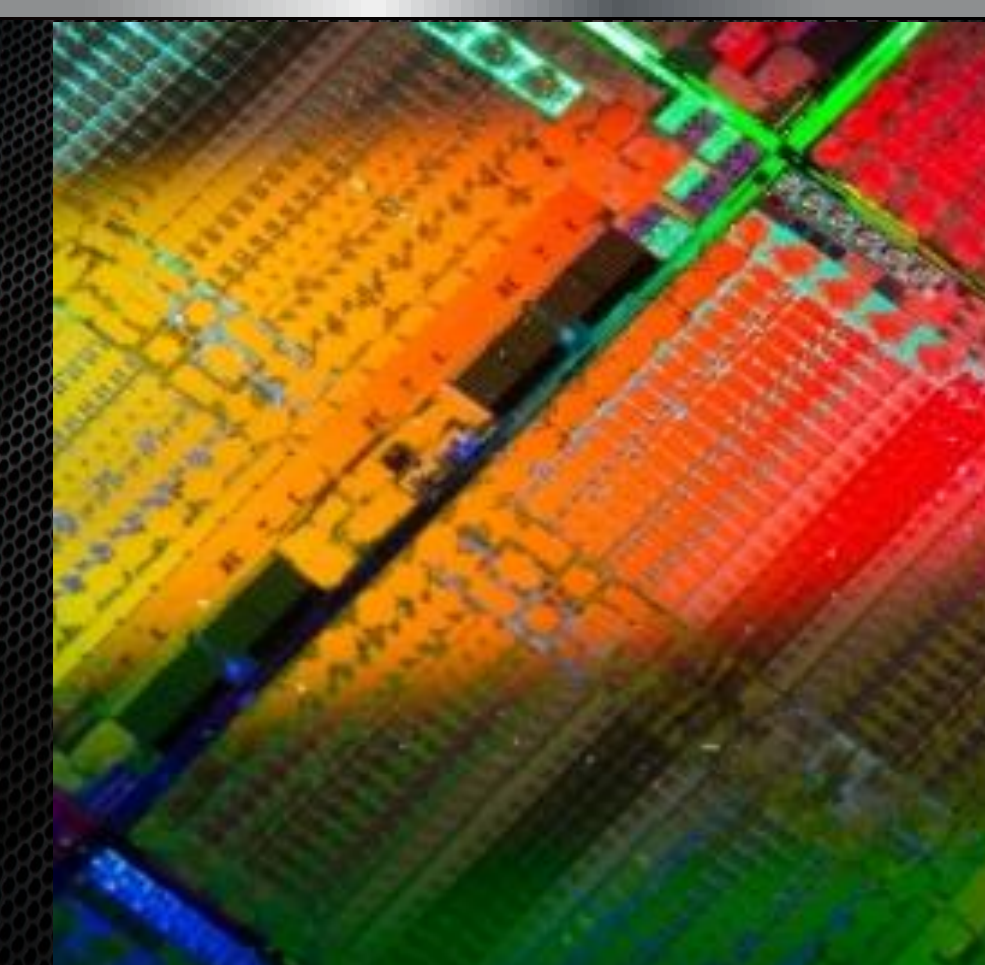
Cost Effective 448 GB/S of Memory Bandwidth



**INTERCONNECT BANDWIDTH**

**PCIe® 4.0 Support**

Up to 2X Interconnect Bandwidth Of PCIe® Gen3



**ARCHITECTURE**

**New GFX RDNA**

Designed For Gaming Performance & Efficiency



# “NAVI”

## Radeon™ Display Engine

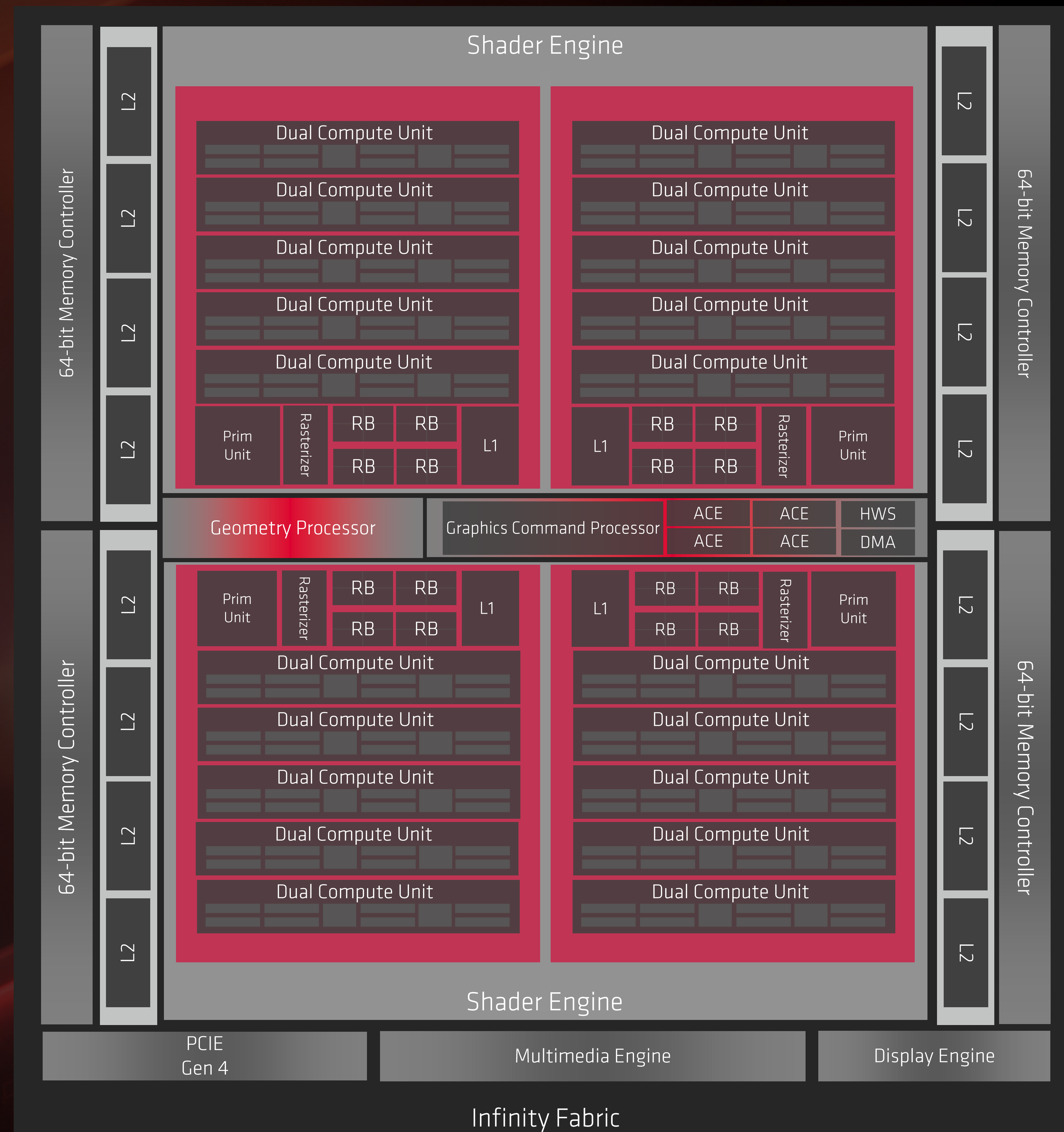
New High Resolution HDR Displays  
New Levels of Compression

## Radeon™ Multi-Media Engine

Seamless Streaming  
Improved Encoding

## New Graphics RDNA Architecture

New Compute Units  
Multilevel Cache  
Streamlined Graphics Engine





# RADEON™ DISPLAY ENGINE

FEATURING  
AMD RADEON FREESYNC™ TECHNOLOGY

**HDMI 2.0 & DisplayPort 1.4 HDR**

**Display Stream Compression 1.2a**

**Direct Read of DCC Compressed Surfaces**

Optimized for High Resolution HDR displays

4K 240Hz | SINGLE CABLE | 8K 60Hz

Optimized for Head Mounted Displays

Single IO connectivity

High Fidelity Internal Color Depth

30 bpp color

Better Power Efficiency

Multi Plane Overlay Protocol with Low voltage mode



# RADEON™ MULTIMEDIA ENGINE

SEAMLESS STREAMING

## IMPROVED ENCODING

NEW HDR/WCG ENCODE (HEVC)

8K ENCODE (HEVC & VP)

40% ENCODER SPEEDUPS

VP9

YouTube

DECODE

4K90

8K24

H.264  
MPEG-4

twitch



DECODE

1080p600

4K150

ENCODE

1080p360

4K90

H.265  
HEVC



NEXT  
GEN

DECODE

1080p360

4K90

8K24

ENCODE

1080p360

4K60



# RDNA

## NEW GRAPHICS ARCHITECTURE

DESIGNED FOR  
THE FUTURE  
OF GAMING

NEW COMPUTE UNIT  
DESIGN

Great Efficiency for Diverse  
Modern Gaming

MULTILEVEL CACHE  
HIERARCHY

Low Latency, High Bandwidth,  
Low Power

STREAMLINED  
GRAPHICS PIPELINE

Excellent Performance per  
Higher Clock Frequencies



# “NAVI” STATS

## RDNA, ALL NEW ARCHITECTURE

### 40 RDNA Compute Units

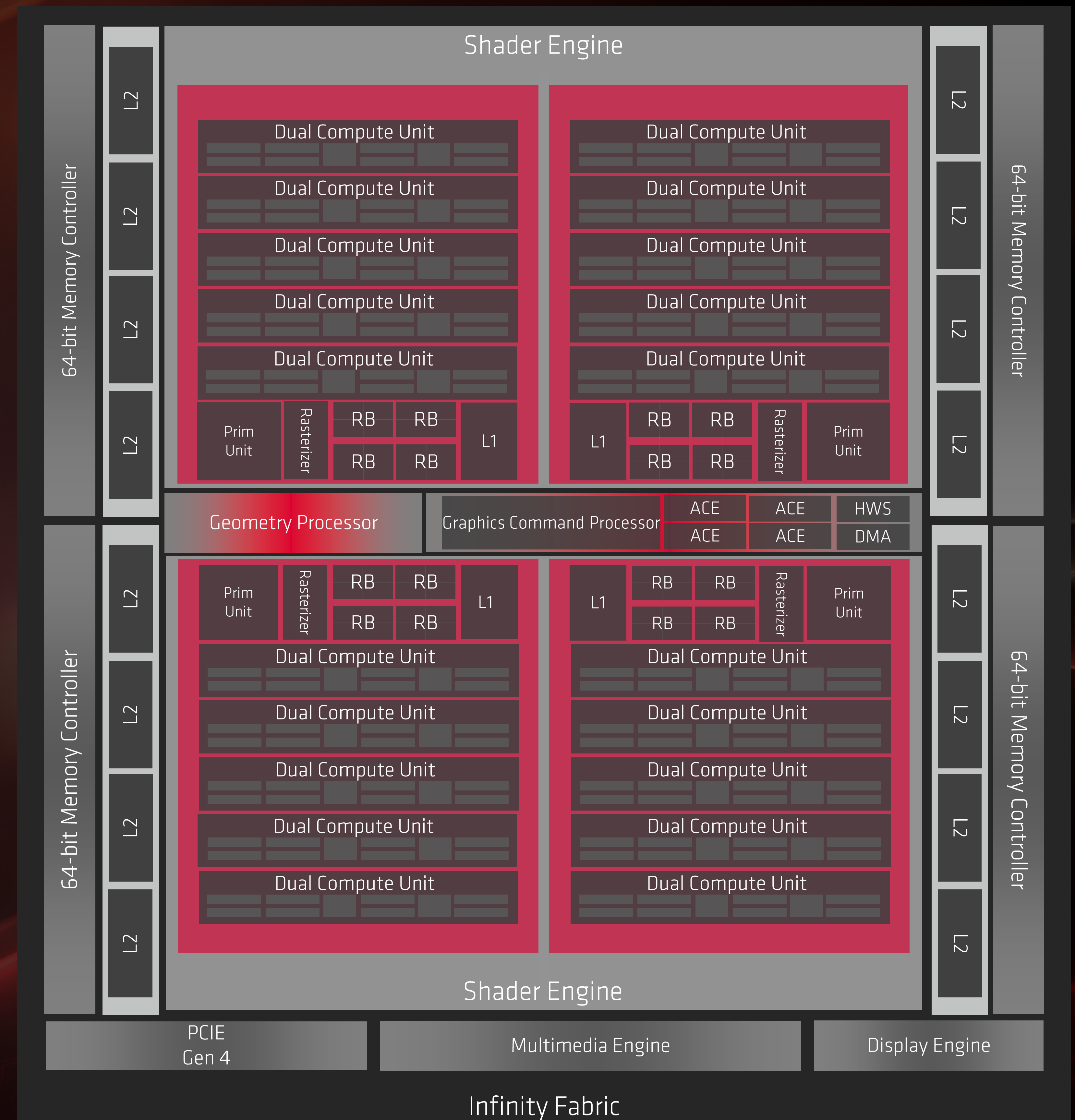
- 80 Scalar Processors
- 2560 Stream Processors
- 160 64b Bilinear Filter units

### Multilevel Cache

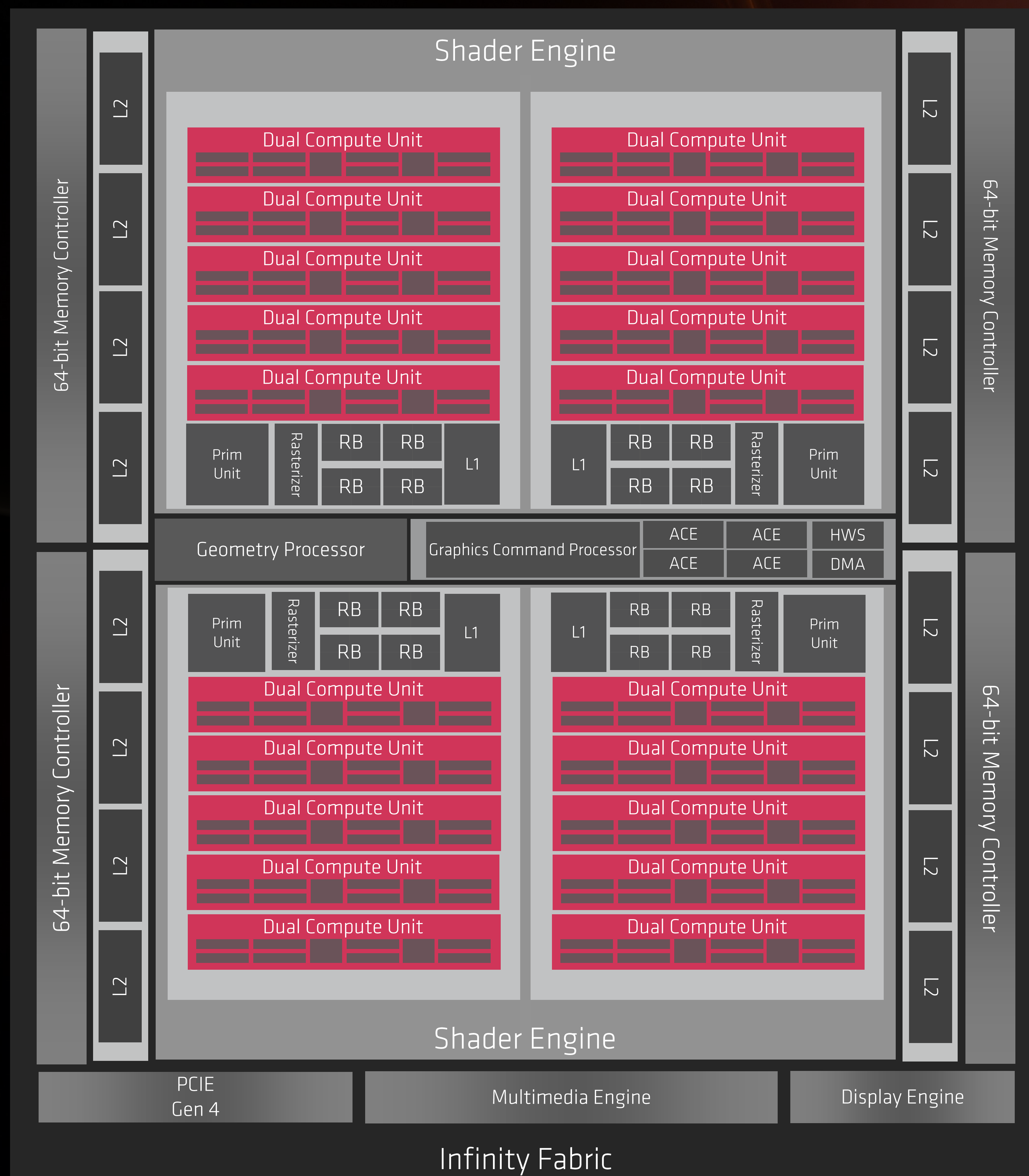
- 4MB L2, 512KB L1, (V\$, I\$, K\$) L0
- 2x V\$L0 Load Bandwidth
- DCC Everywhere

### Streamlined Graphics Engine

- Geometry Engine (4 Prim shader out, 8 Prim shader in)
- 64 Pixel Units
- 4 Asynchronous Compute Engines
- Balanced Work Distribution & Redistribution
- Designed for higher frequencies at lower power







# RDNA COMPUTE UNIT

GREAT EFFICIENCY FOR DIVERSE  
WORKLOADS

**2x Vector & Scalar Instruction Rate**

**Single Cycle Instruction Issue**

**Dual Wave Length (32\64) Modes**

**Adjacent CU Resource Pooling**



# RADEON™ ARCHITECTURAL ADVANCES OF PROGRAMMABLE GRAPHICS

## PRE 2000 1st ERA R100

Fixed Function

3D GEOMETRY TRANSFORMATION

$$V_{eye} \begin{bmatrix} x \\ y \\ z \\ w \end{bmatrix} = MVP \begin{bmatrix} m_0 & m_4 & m_8 & m_{12} \\ m_1 & m_5 & m_9 & m_{13} \\ m_2 & m_6 & m_{10} & m_{14} \\ m_3 & m_7 & m_{11} & m_{15} \end{bmatrix} \cdot V_{obj} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

$$V_{tex} \begin{bmatrix} s \\ t \\ r \\ q \end{bmatrix} = M_{proj} \cdot V_{in}$$

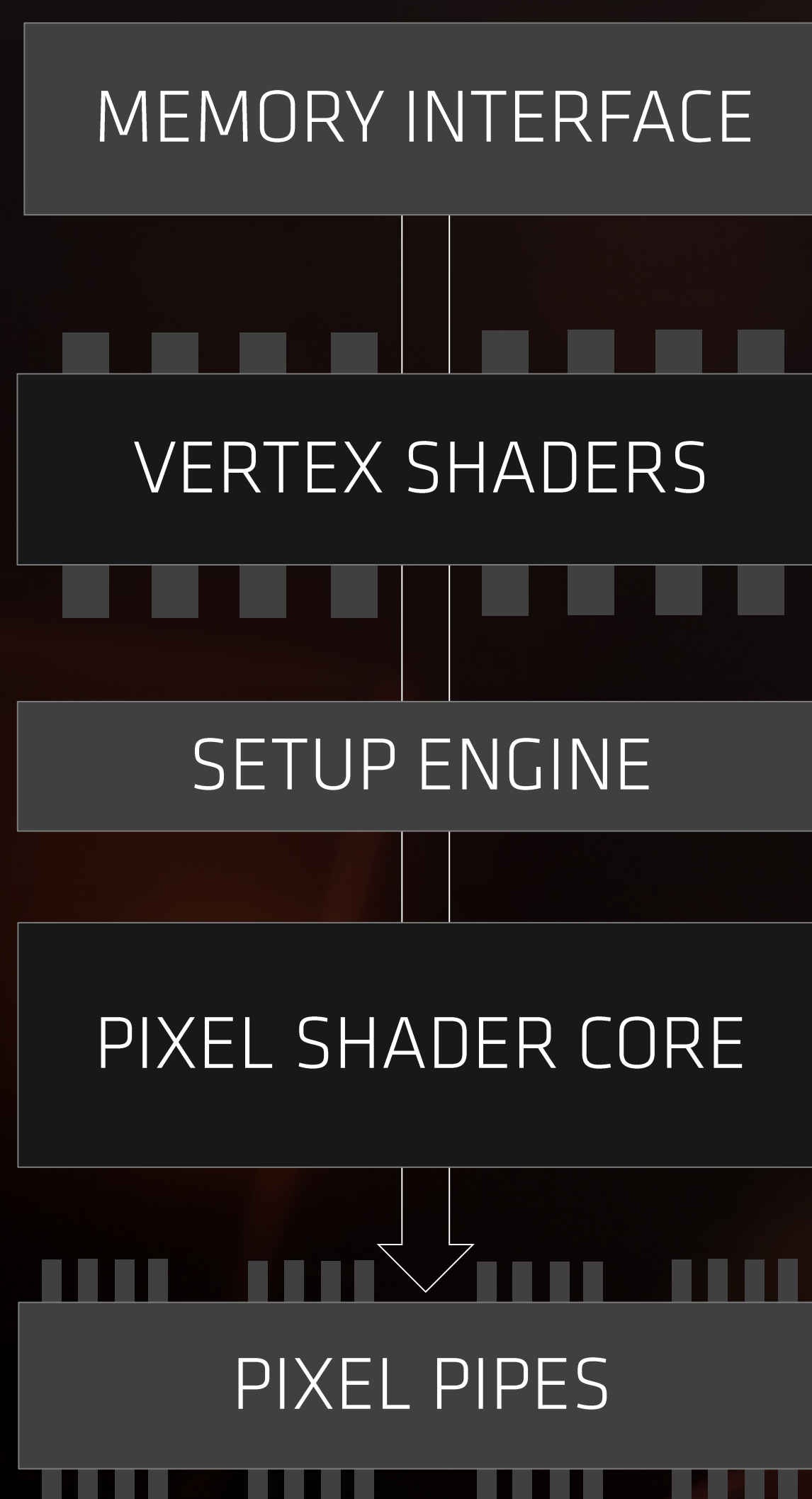
LIGHTING

$$C_p = k_a L_a + \sum_{n-lights} Att_n (k_d (\hat{L}_n \cdot \hat{N}) + k_s (\hat{R}_n \cdot \hat{V})^\alpha)$$



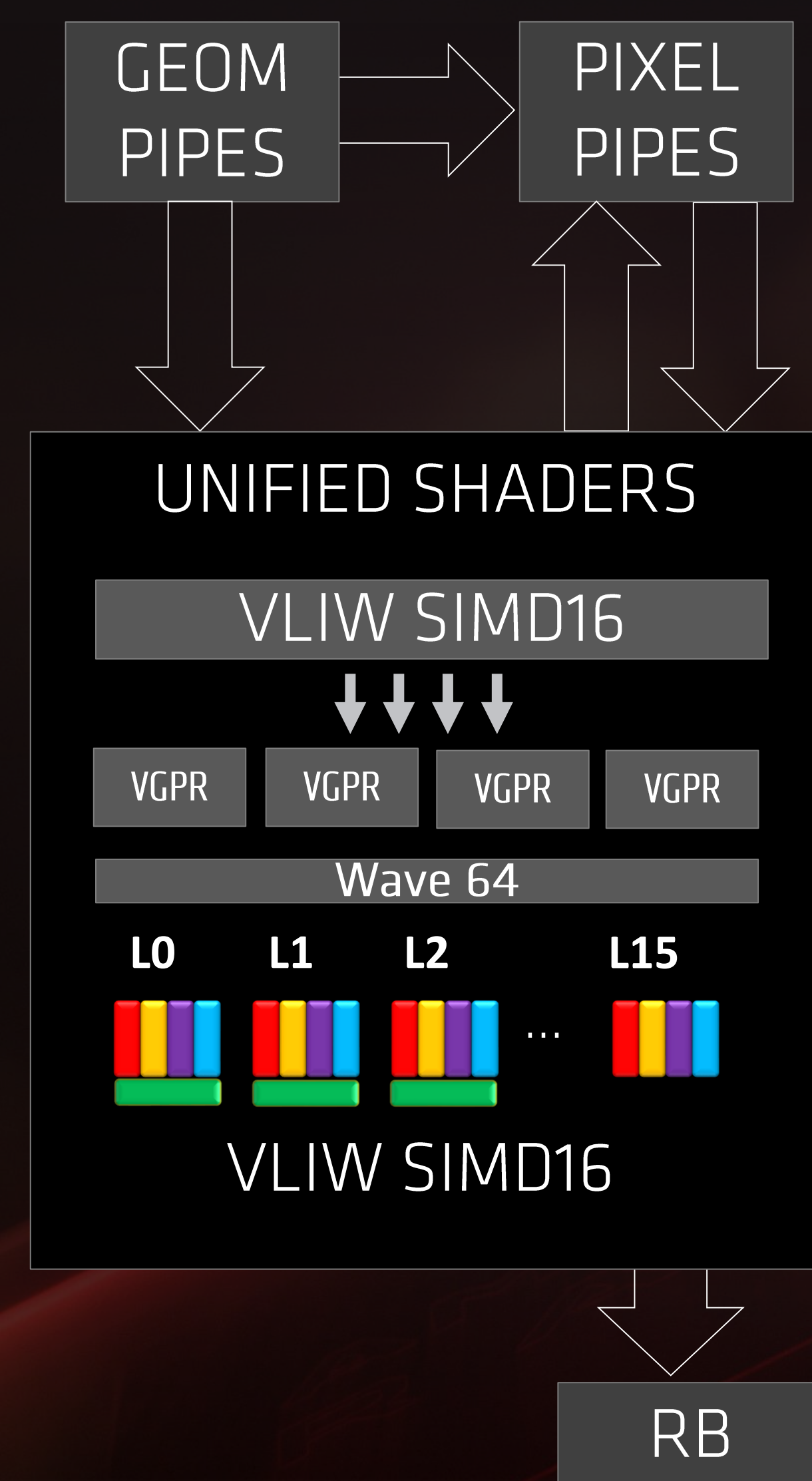
## 2001-2007 2nd ERA R200-R500

Simple VS/PS Shaders



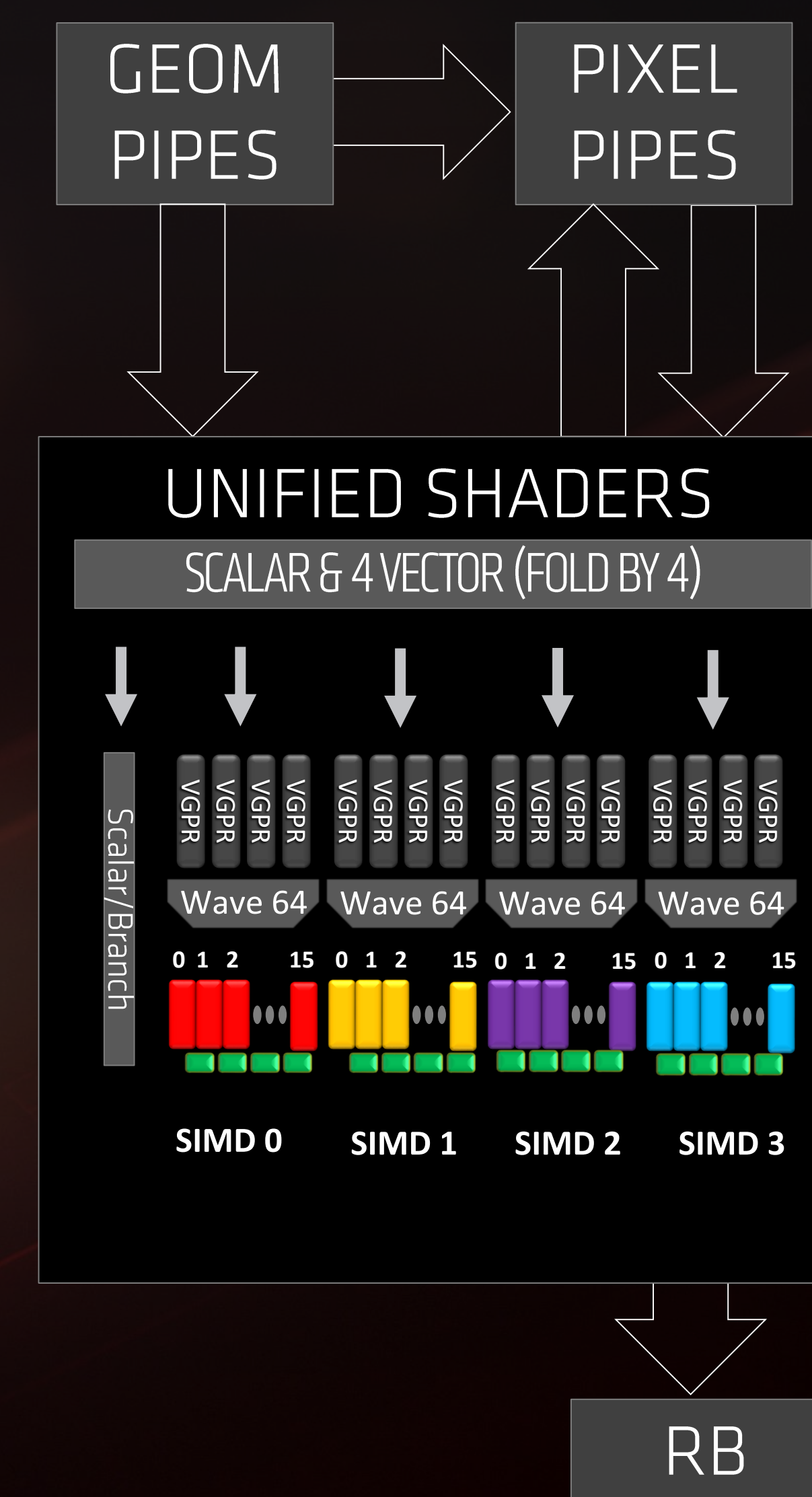
## 2008-2011 3rd ERA R600

**TeraScale**  
Unified Shaders with VLIW



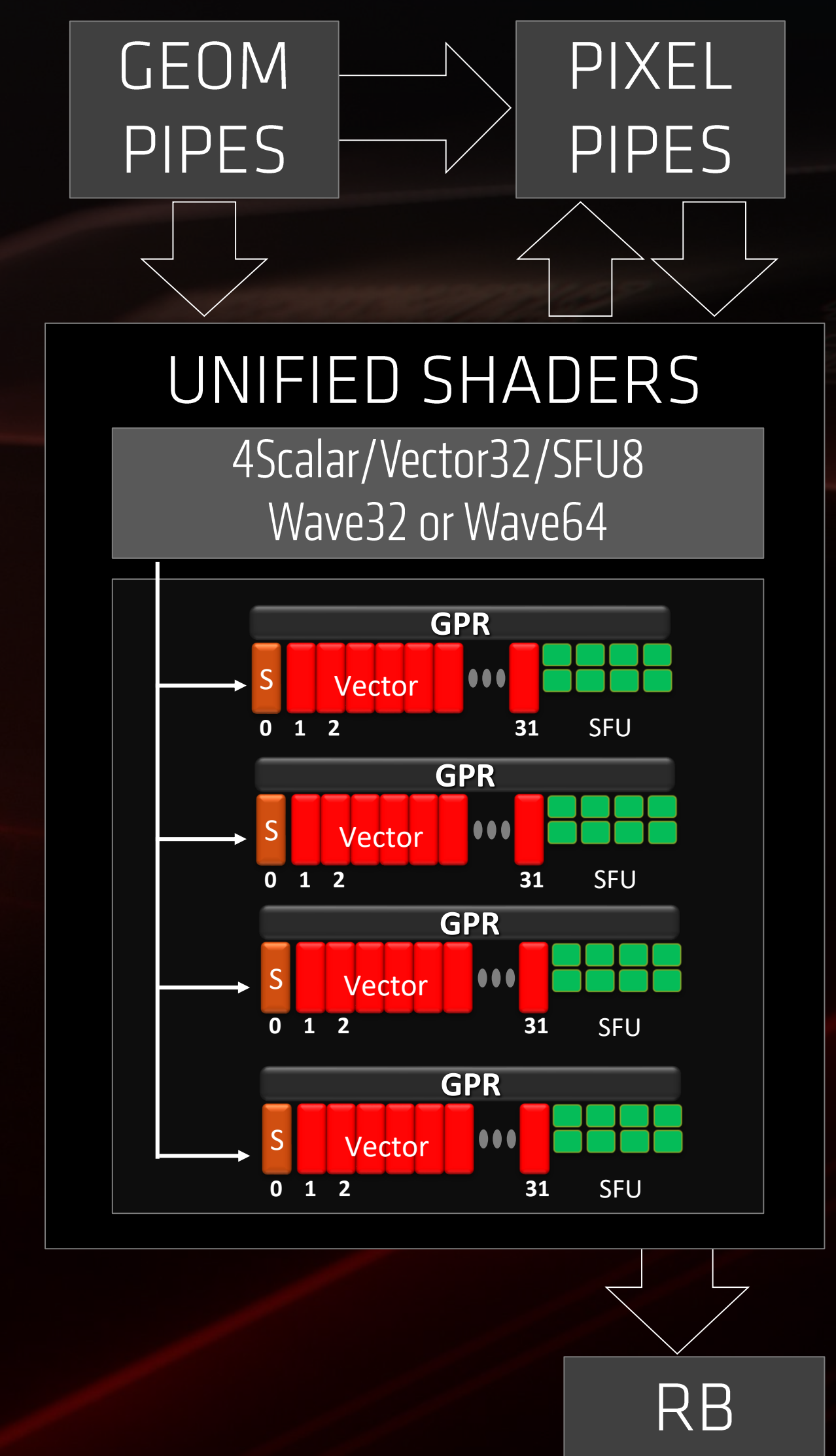
## 2012-2018 4th ERA Southern Islands GCN

**GCN**  
Unified Shader with Scalar & Vector  
(Fold by 4)



## 5th ERA "NAVI" RDNA

**RDNA**  
Unified Shader with Scalar Vector  
(SIMT ILP Capable)



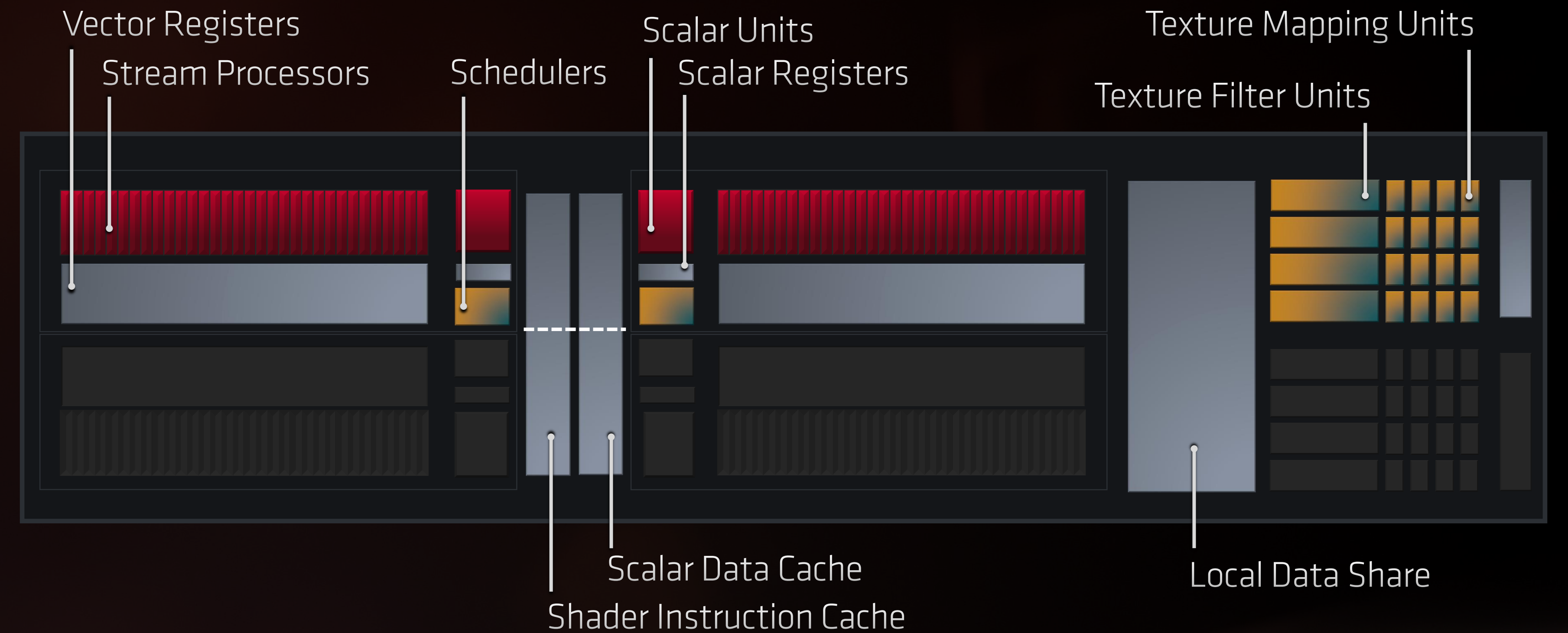


# RDNA COMPUTE UNIT

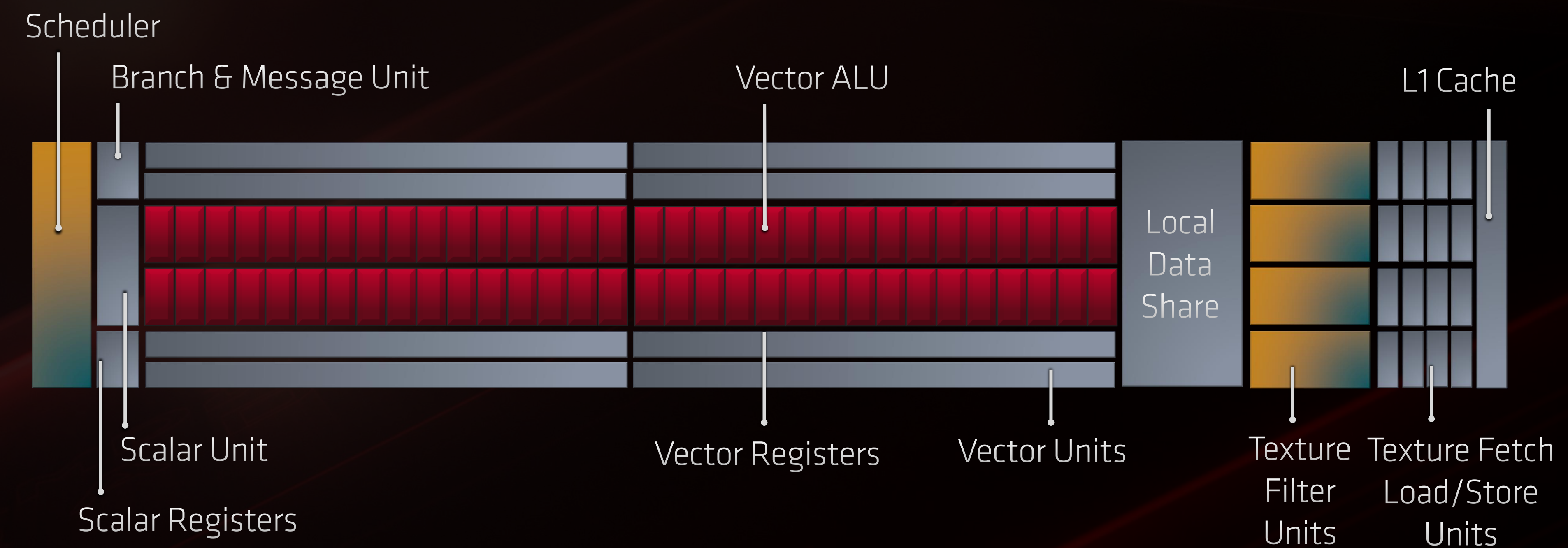
REDESIGNED FOR SINGLE  
THREADED PERFORMANCE

NEW COMPUTE UNIT

## RDNA COMPUTE UNIT

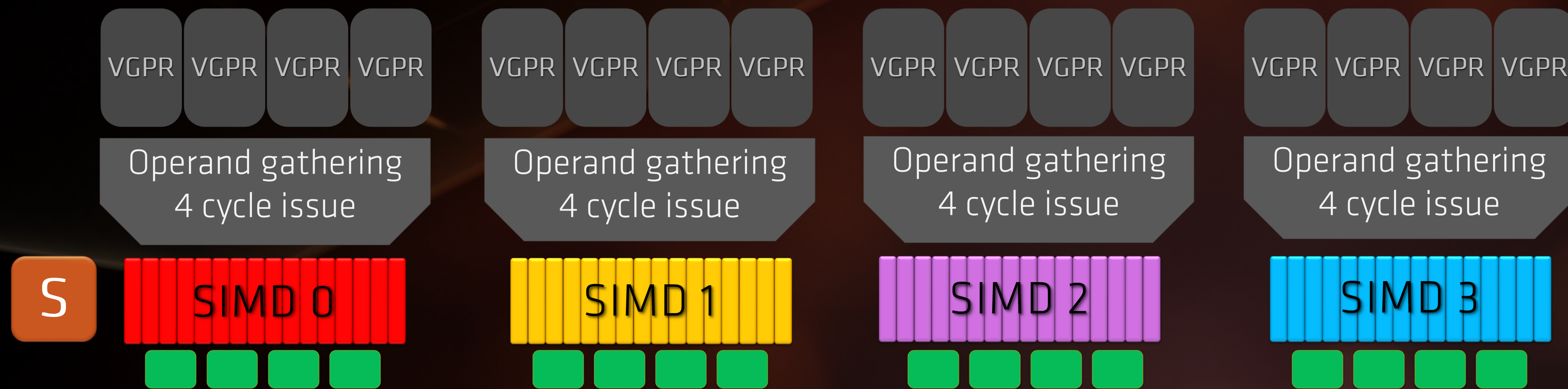


## GCN COMPUTE UNIT



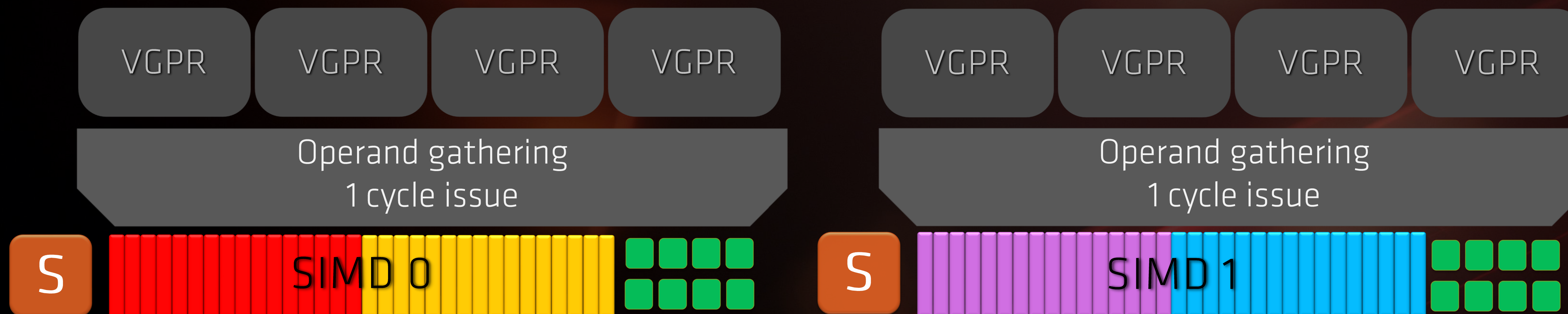


# EXECUTION UNITS



## GCN

- 4 SIMD16 Vector Unit
- 4 SIMD4 Special Function Unit
- 1 Shared Scalar Decode & Issue Unit
- 1 Shared Vector Decode & Issue Unit
- 256 KB VGPR

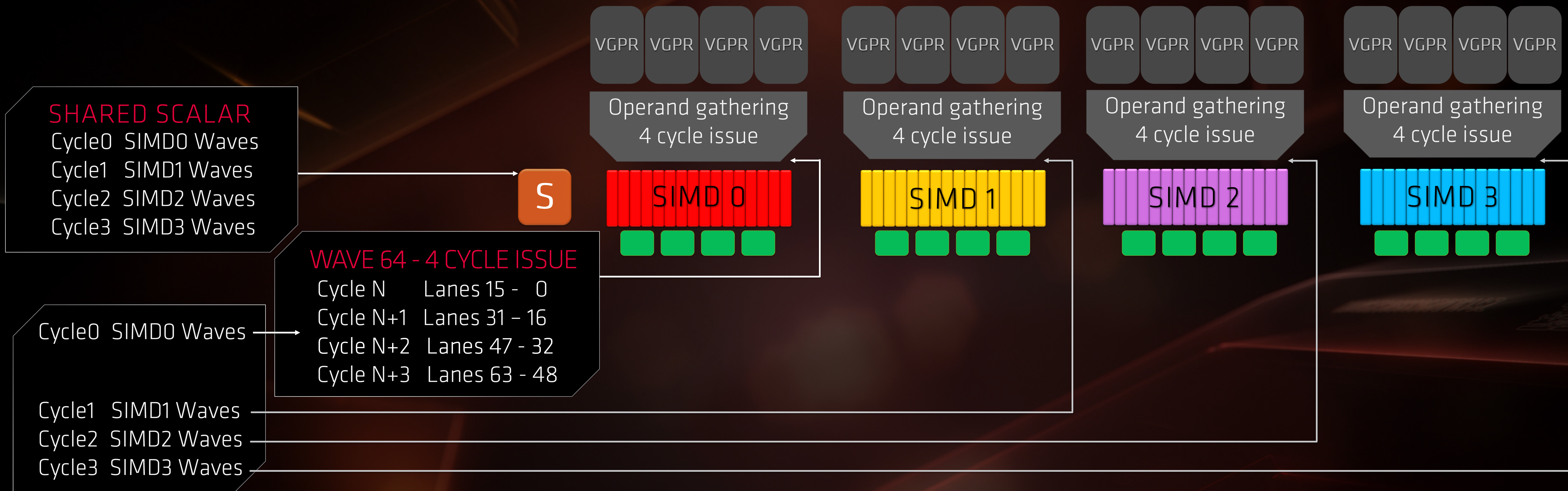


## RDNA

- 2 SIMD32
- 2 SIMD 8 Special Function Unit
- 2 Scalar Decode and Issue Units
- 2 Vector Decode and Issue Units
- 256 KB VGPR



# GCN INSTRUCTION ISSUE

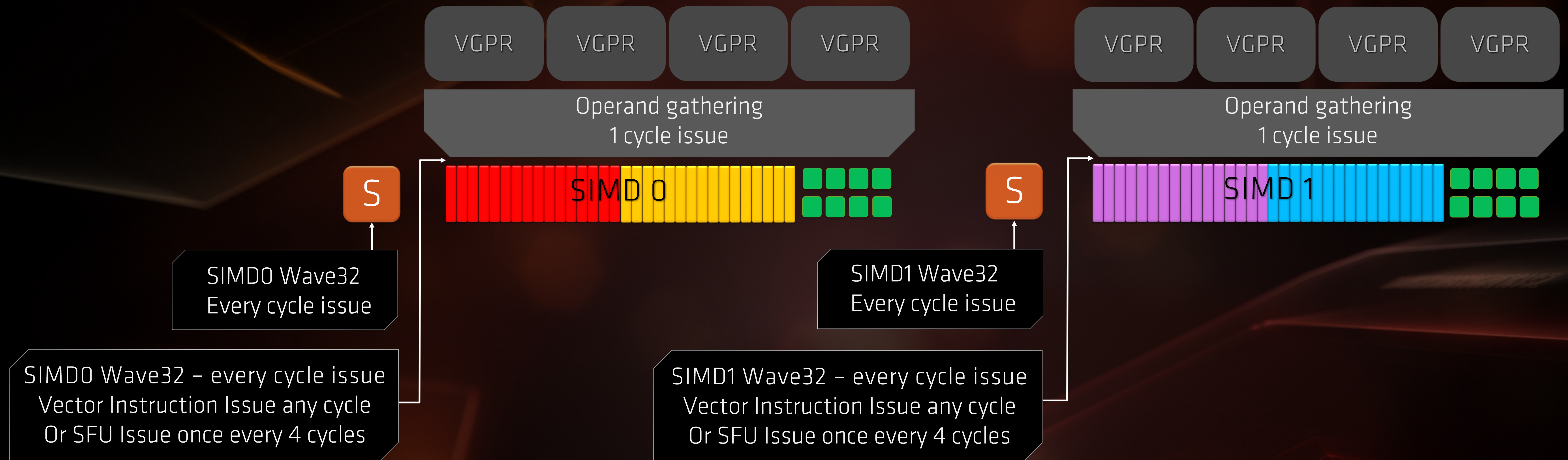


All work-items of a wave64 have an opportunity to do work once every 4 clocks due to hardware interleaving  
 Special Function Unit alternate execution unit running at 1/4 rate

A wave from a SIMD has an opportunity to accomplish a scalar instruction once every 4 clocks



# RDNA INSTRUCTION ISSUE



Vector Units - All work-items of one wave32 have an opportunity to do work every clock  
Special Function Unit uses 1 issue cycle and then executes in parallel  
Each SIMD equipped with a scalar unit for an instruction execution every cycle



# INSTRUCTION ISSUE EXAMPLE

## EXAMPLE SHADER

```
s_add_i32 s0, s1, s2
v_mul_f32 v0, v1, s0
v_add_f32 v5, v4, v3
v_sub_f32 v6, v7, v0
```

## CYCLE

## “VEGA” EXECUTION

```
0 s_add_i32 s0, s1, s2
1 ...
2 ...
3 ...
4 v_mul_f32 v0, v1, s0
5 ... (simd busy 4 cycles)
6 ...
7 ...
8 v_add_f32 v5, v4, v3
9 ...
10 ...
11 ...
12 v_sub_f32 v6, v7, v0
13 ...
14 ...
15 ...
```

## “NAVI” WAVE32

```
s_add_i32 s0, s1, s2
... (salu dependency stall on S0)
v_mul_f32 v0, v1, s0
v_add_f32 v5, v4, v3
... (valu dependency stall on V0)
...
v_sub_f32 v6, v7, v0
```

SHORTEST  
WAVE ISSUE  
LATENCY

## “NAVI” WAVE64

```
s_add_i32 s0, s1, s2
... (salu dependency stall on S0)
v_mul_f32 v0, v1, s0 (lo)
v_mul_f32 v0, v1, s0 (hi)
v_add_f32 v5, v4, v3 (lo)
v_add_f32 v5, v4, v3 (hi)
... (valu dependency stall on V0 lo)
v_sub_f32 v6, v7, v0 (lo)
v_sub_f32 v6, v7, v0 (hi)
```

44%  
REDUCTION IN  
ISSUE CYCLES



# BENEFITS OF NEW WORK DISTRIBUTION & INSTRUCTION EXECUTION CHANGES

## SINGLE THREADED PERFORMANCE IMPROVEMENT

WORK LOAD EXAMPLE: 64 WORK-ITEMS ALU INTENSIVE CODE

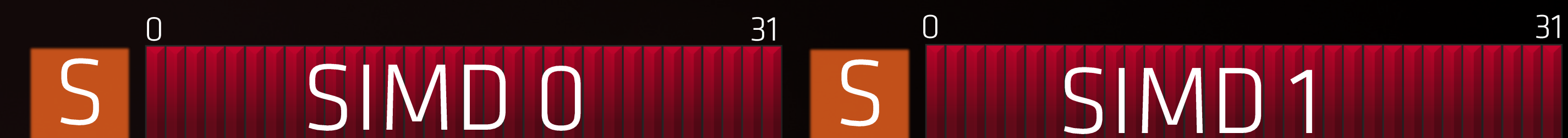
### GCN



1 Wave64 → SIMD16  
Instruction Issue → 4 clock  
CU ALU → 25% utilized

EFFECTIVE THROUGHPUT

### RDNA



2 Wave32 → 2 SIMD32  
Instruction Issue → 1 clock  
CU ALU → 100% utilized

ILP UNLOCKS UP TO 4X FASTER EXECUTION

## RDNA MORE EFFECTIVELY UTILIZES THE MACHINE

- Engage machine quicker by uniformly distributing work to all ALUs
- Optimize efficiency and latency by preferring highest priority/oldest work
- Extract program ILP and scheduling to benefit from data locality
- Utilize multi-threading of waves to hide remaining latencies for throughput



# RDNA SIMD UNIT

REDESIGNED FOR  
**SINGLE THREADED PERFORMANCE**  
AND  
**EFFECTIVE IPC**

UP TO  
**20 Wave  
Controllers**  
WAVE 32 or WAVE64

**Dedicated  
Instruction Issue Unit**

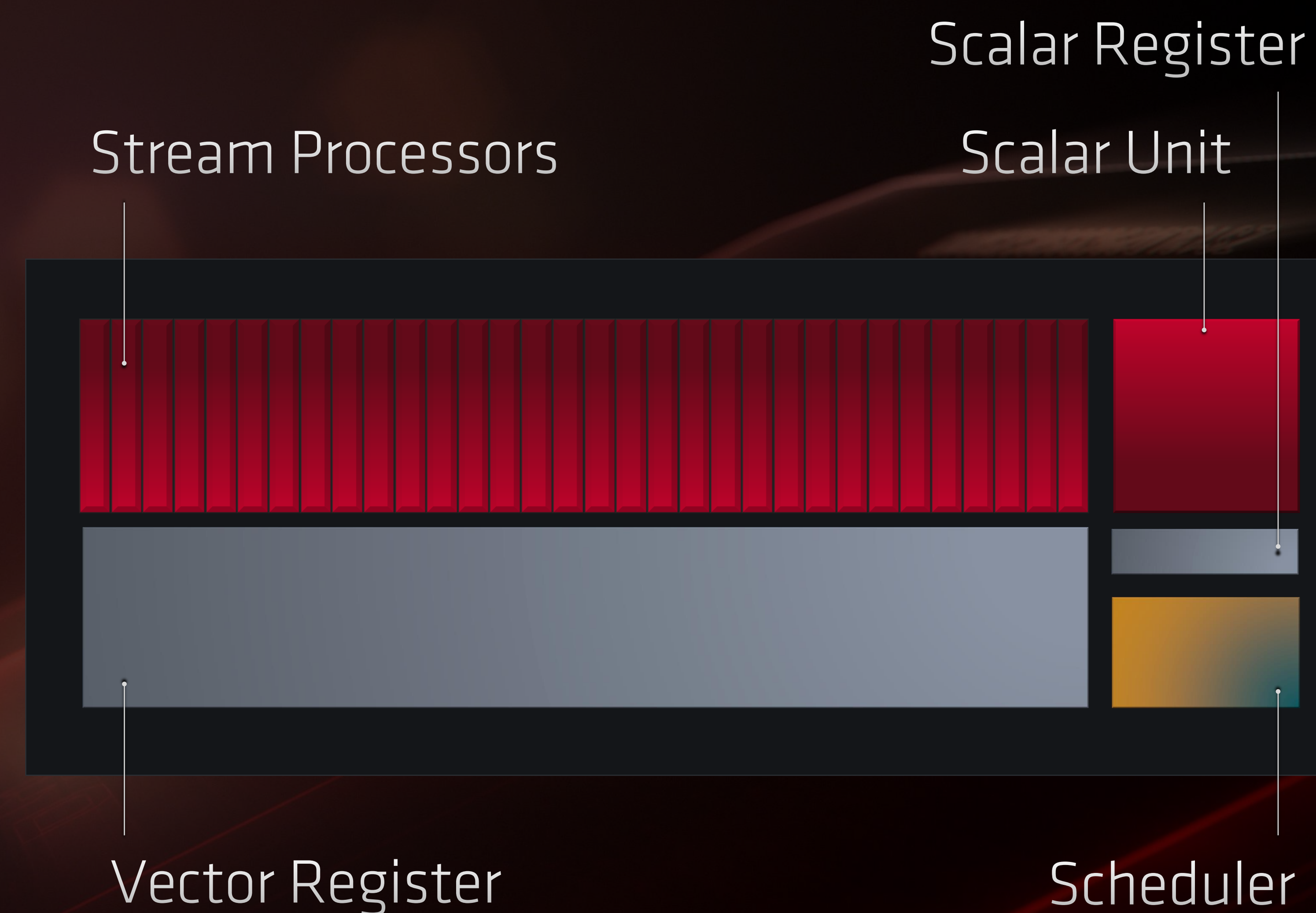
**Dedicated  
Scalar Unit**

**Increased SIMD Width**  
SIMD 16 → SIMD 32

**Improved Instruction  
Arbitration & Prefetch**

**Full Rate 32b  
& Dual 16b ALU**

RESULTS  
**REDUCED WAVE-LIFETIME & IMPROVED EFFICIENCY**





# WAVE EXECUTION

ENABLING DETERMINISTIC GAMEPLAY

MORE TOLERANT TO DIVERGENT CODES  
COMPILER DRIVEN CACHE AWARE SCHEDULING

## WAVE 32 EXECUTION

### NATIVE ONE CLOCK EXECUTION

Can operate in a smaller  
cache footprint

Less work to  
hide total machine latency

Fewer resources (GPRS)  
required for wave launch

More SIMDs engaged with  
small workload

## WAVE 64 EXECUTION

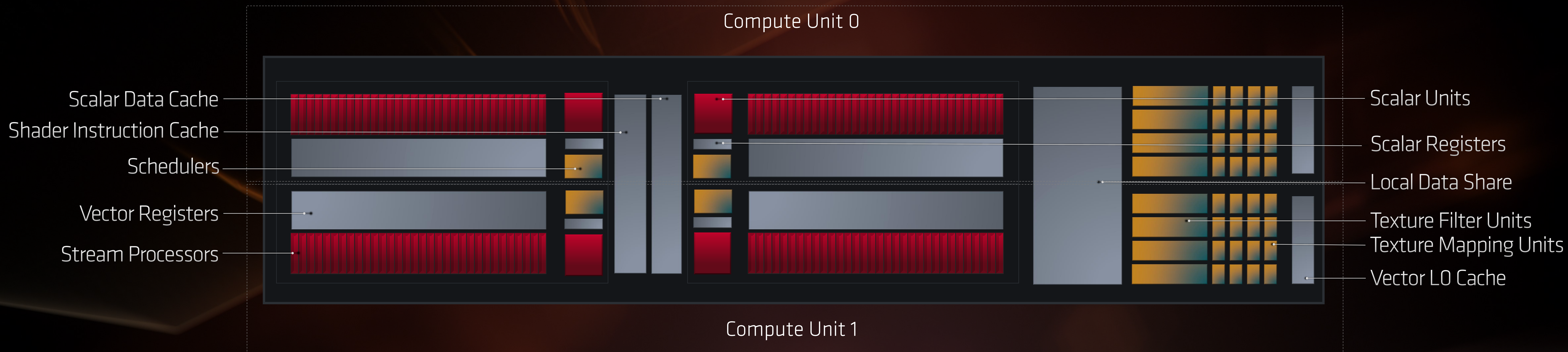
### NATIVE TWO CLOCK EXECUTION

Half wave execution  
improves efficiency

Reduces exposure to  
ALU pipe stalls

Enable double work item occupancy for latency hiding





# COMPUTE UNIT COOPERATION

WORKGROUP PROCESSOR

UP TO  
**2X**  
ALUs

ACCESS TO  
**2X**  
Registers

ACCESS TO  
**4X**  
Cache Bandwidth



# “NAVI” RDNA ARCHITECTURE

## FUNDAMENTAL CHANGES IN PROGRAMMABLE CORE

### PRE GCN

- VLIW5/VLIW4
- Hard To Program for Performance
- Complex Compiler Technology
- Per Work-Item IPC = 1.25 potential

### GCN

- Wave64 on SIMD16 (4clk issue)
- Easy To Program For Performance
- Standard Compiler Techniques
- Per Work-Item IPC = 0.25

### RDNA

- Wave32 on SIMD32 (1clk issue)
- Easier Achieved Performance
- Enables New Compiler Techniques
- Per Work-Item IPC = 1 potential



# RDNA CACHE HIERARCHY

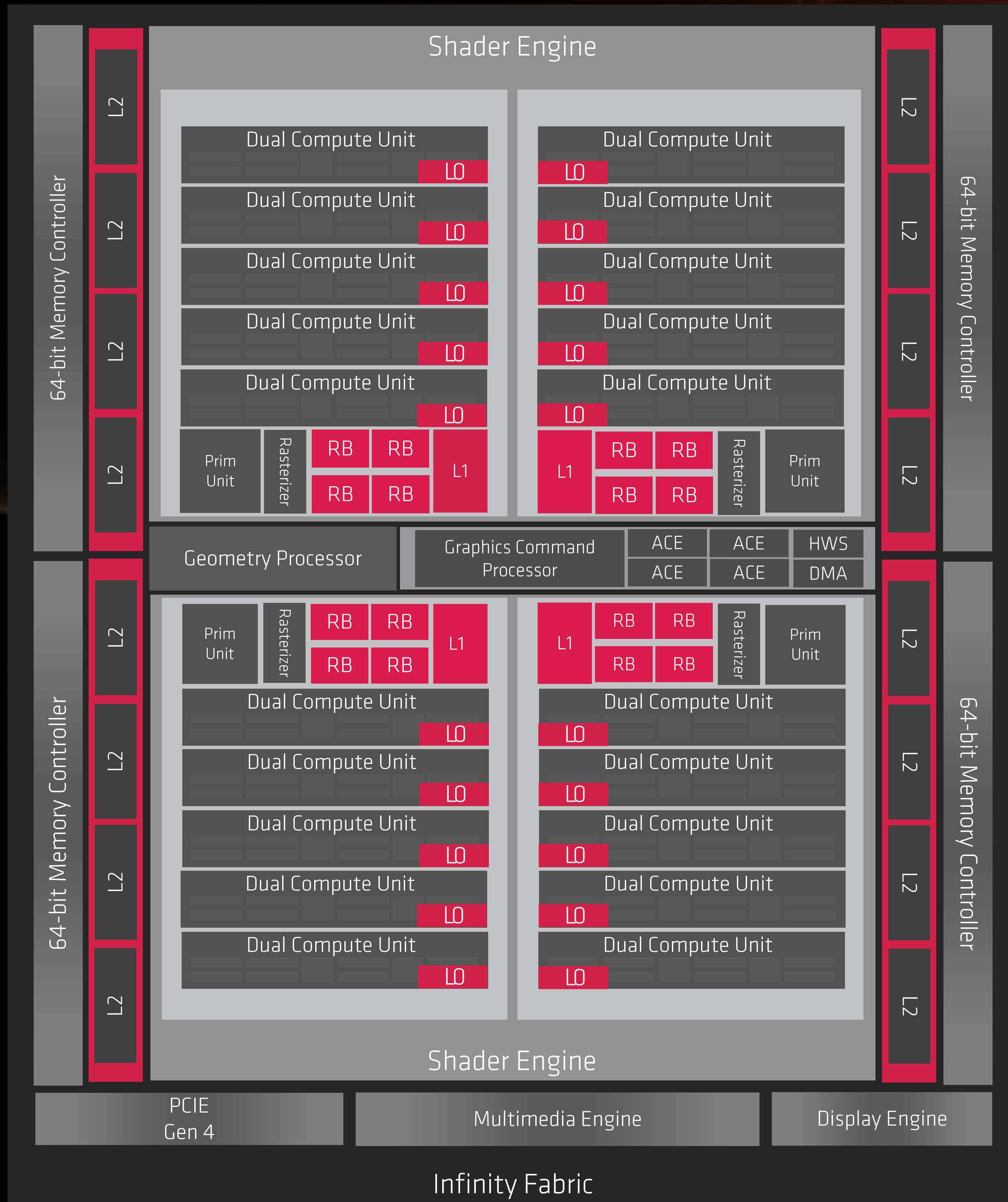
## NEW MULTI LEVEL SYSTEM

**New L1 Level Cache**

**Improved Bandwidth Amplification**

**Reduced Latency and Power**

**Reduced Congestion at L2 Level**





# MULTILEVEL CACHE HIERARCHY

LOW LATENCY, HIGH BANDWIDTH, LOW POWER

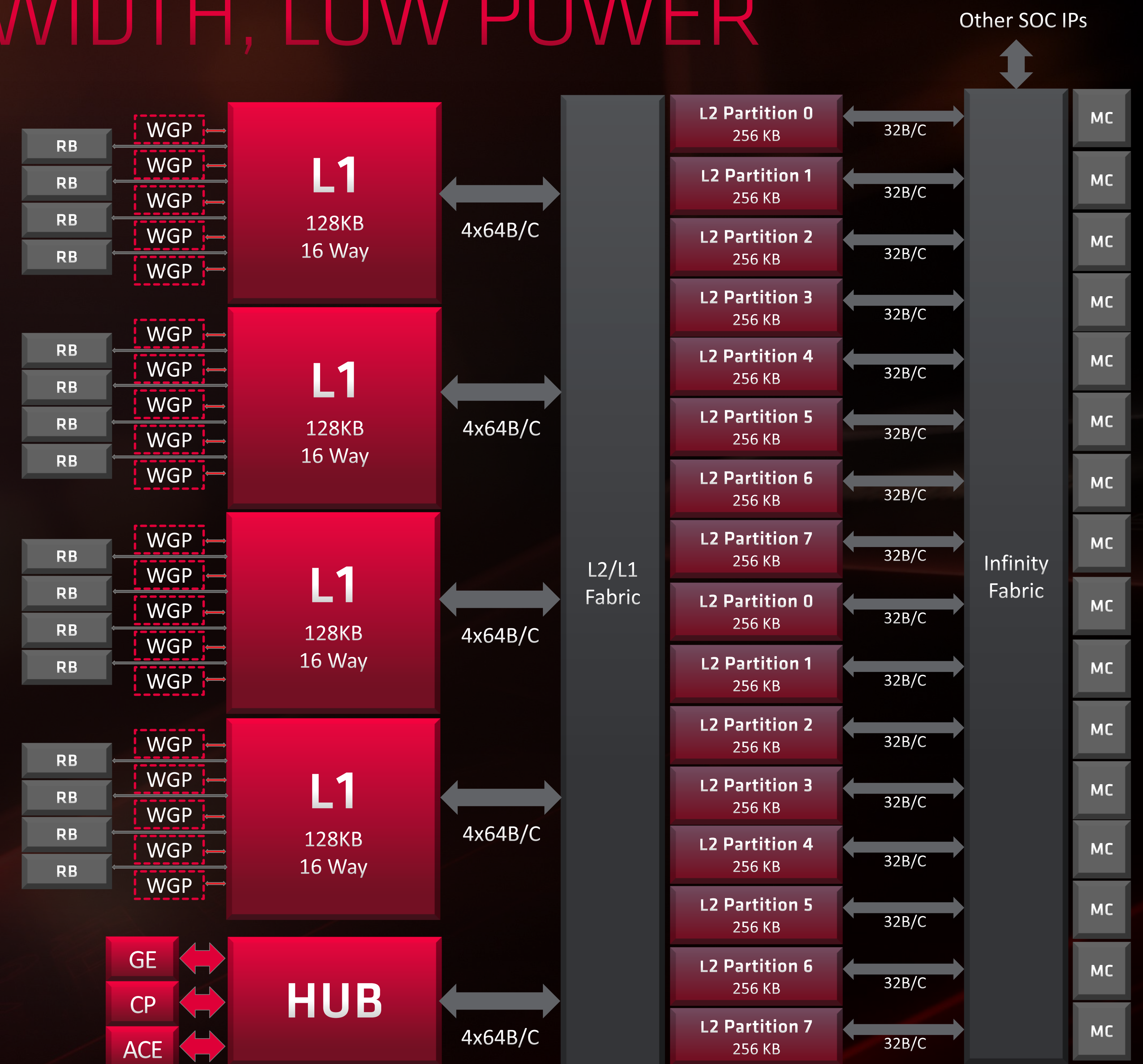
Unified LLC for GFX/ACE Pipes

Instruction Range Based Actions

OOO between R/W, L0, L1, L2, Mem

Reduced Latency and Power

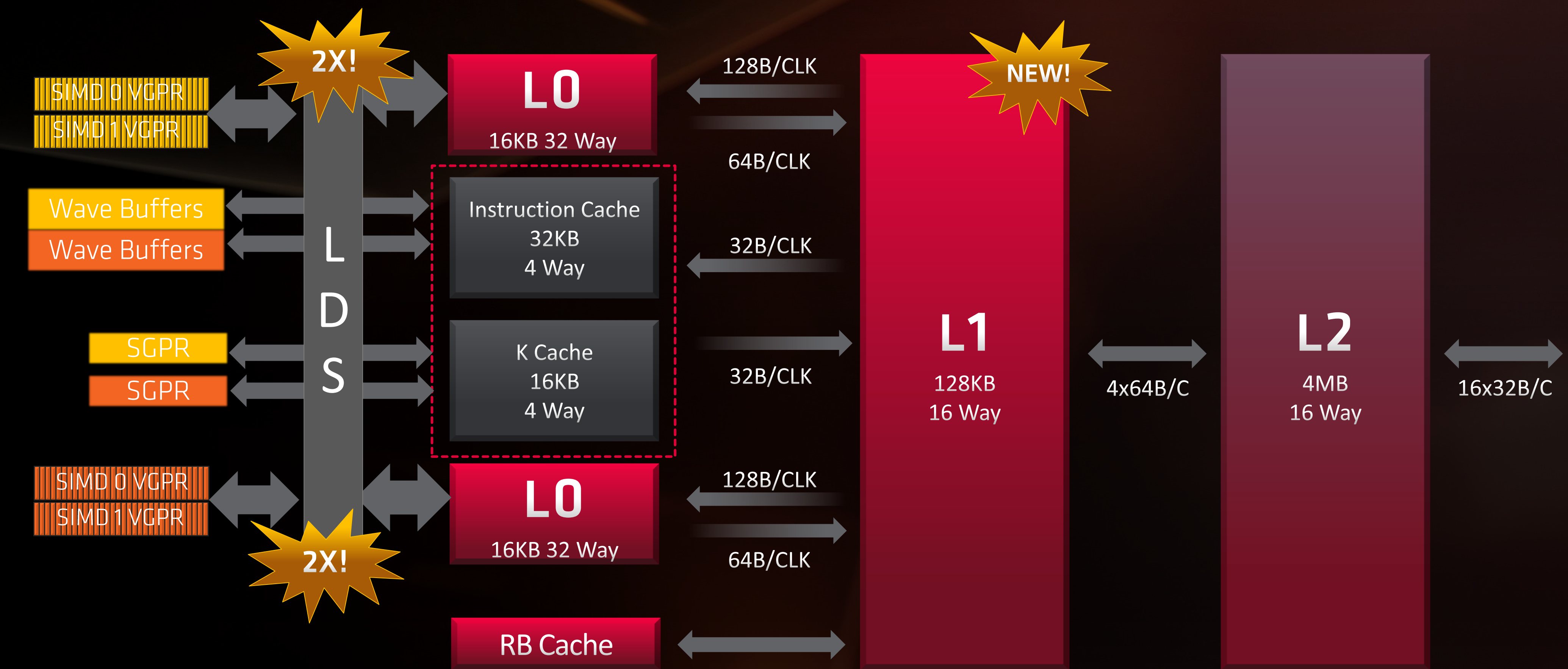
Reduced Data Movement





# MULTILEVEL CACHE HIERARCHY

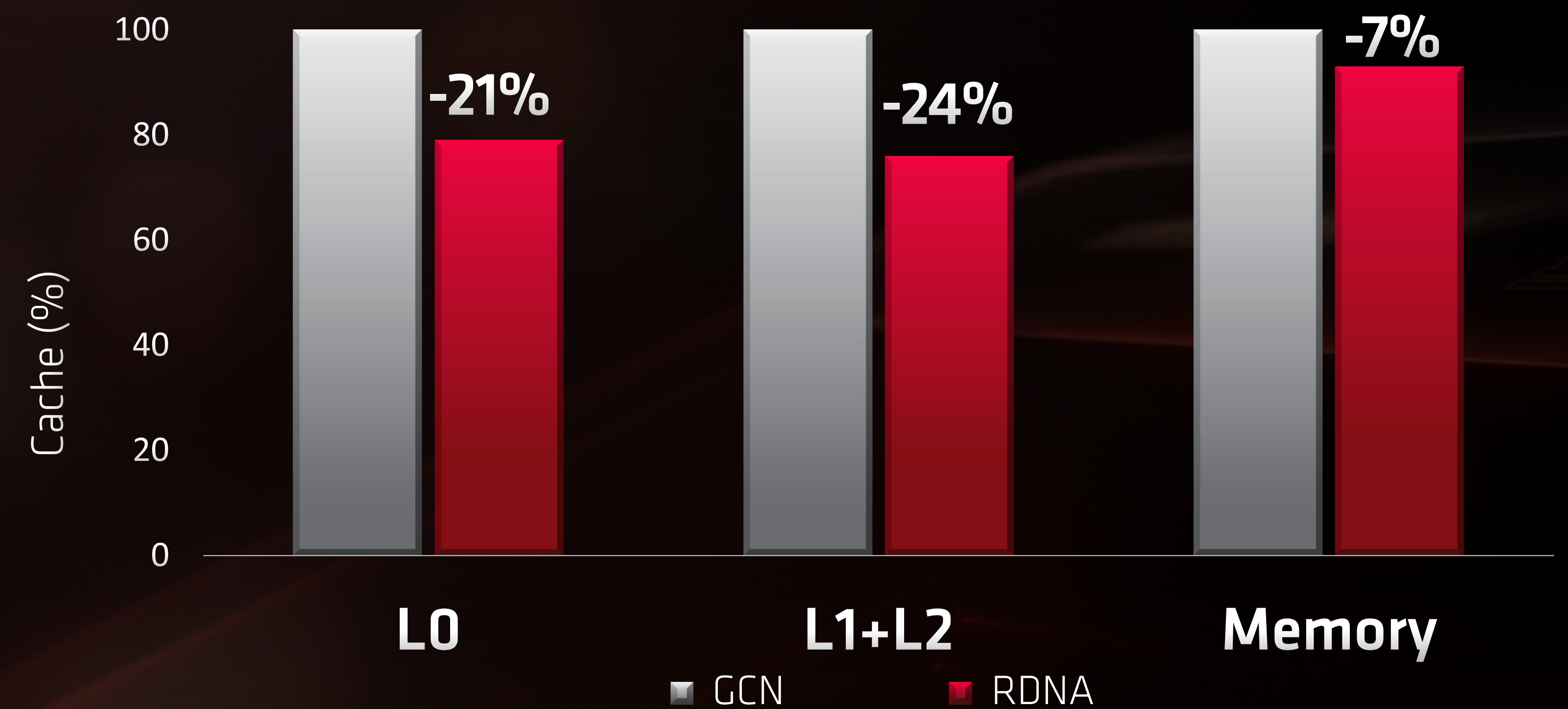
LOW LATENCY, HIGH BANDWIDTH, LOW POWER



Introduce L1 Cache Hierarchy

Double the Load Bandwidth from L0 to ALU

## Relative Cache Latency



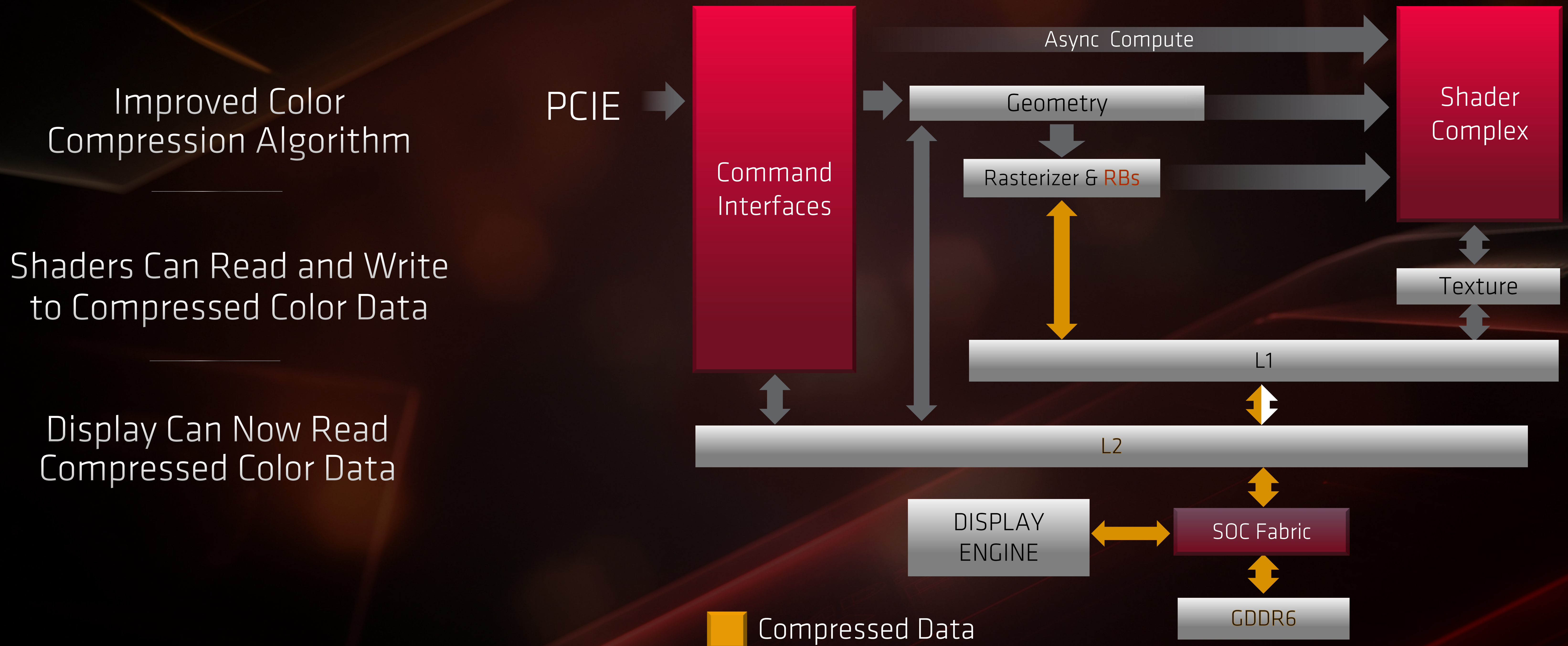
Reduce Latency at Each Level

Improve Effective Bandwidth



# MULTILEVEL CACHE HIERARCHY

## DELTA COLOR COMPRESSION EVERYWHERE





# STREAMLINED GRAPHICS ENGINE

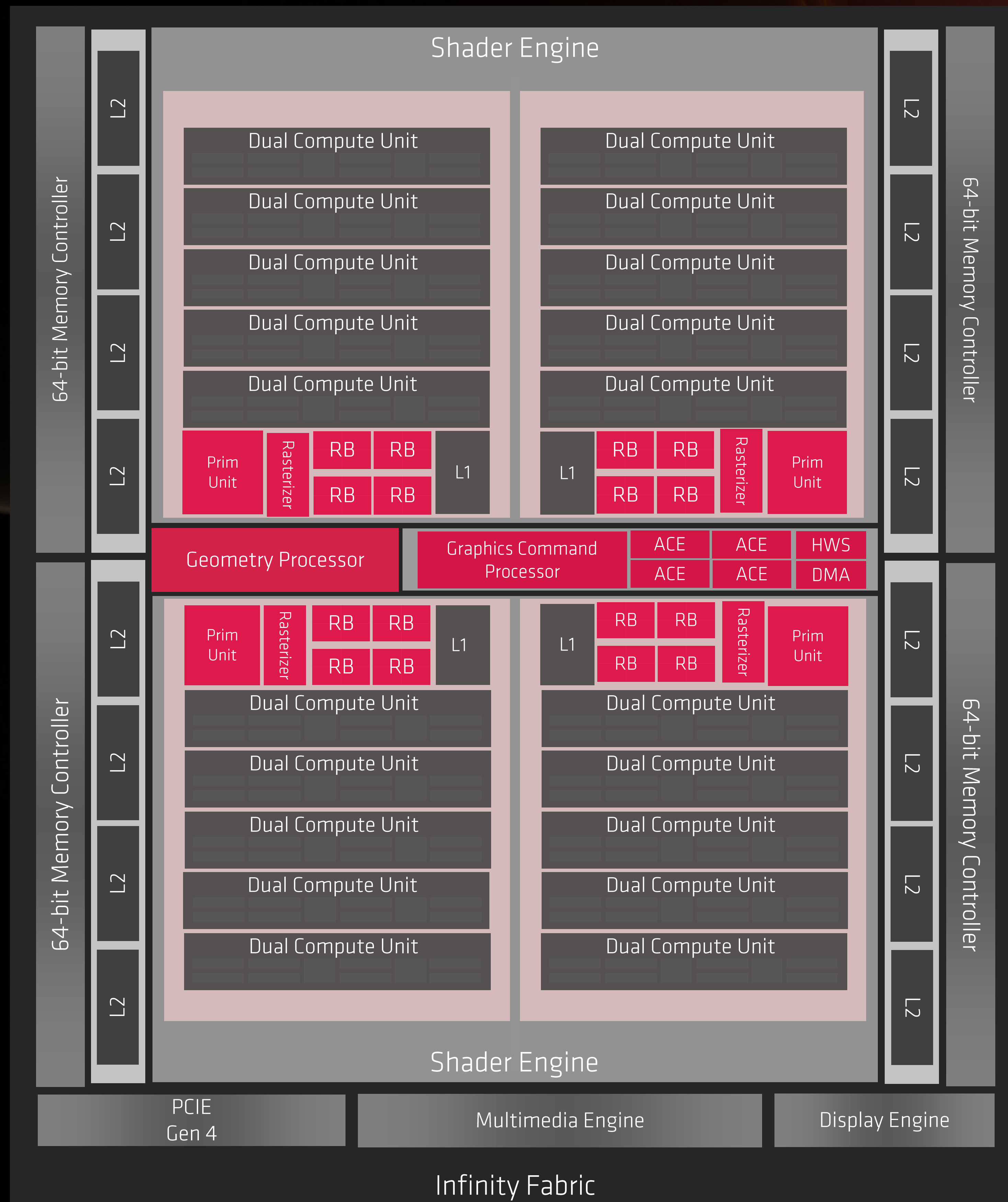
IMPROVED PERFORMANCE PER CLOCK

**4 Enhanced Asynchronous Compute Engines  
Priority Tunneling**

**Centralized Geometry Processor with 4 Prim Units  
Uniformly handle:**

**Vertex reuse, primitive assembly, reset index  
Uniformly distribute pre/post tessellation work  
Shader Culling - 4 Prim out, 8 Prim in**

**64 Pixel Units  
Cache aware pixel wave packing**





# ENHANCED ASYNC COMPUTE

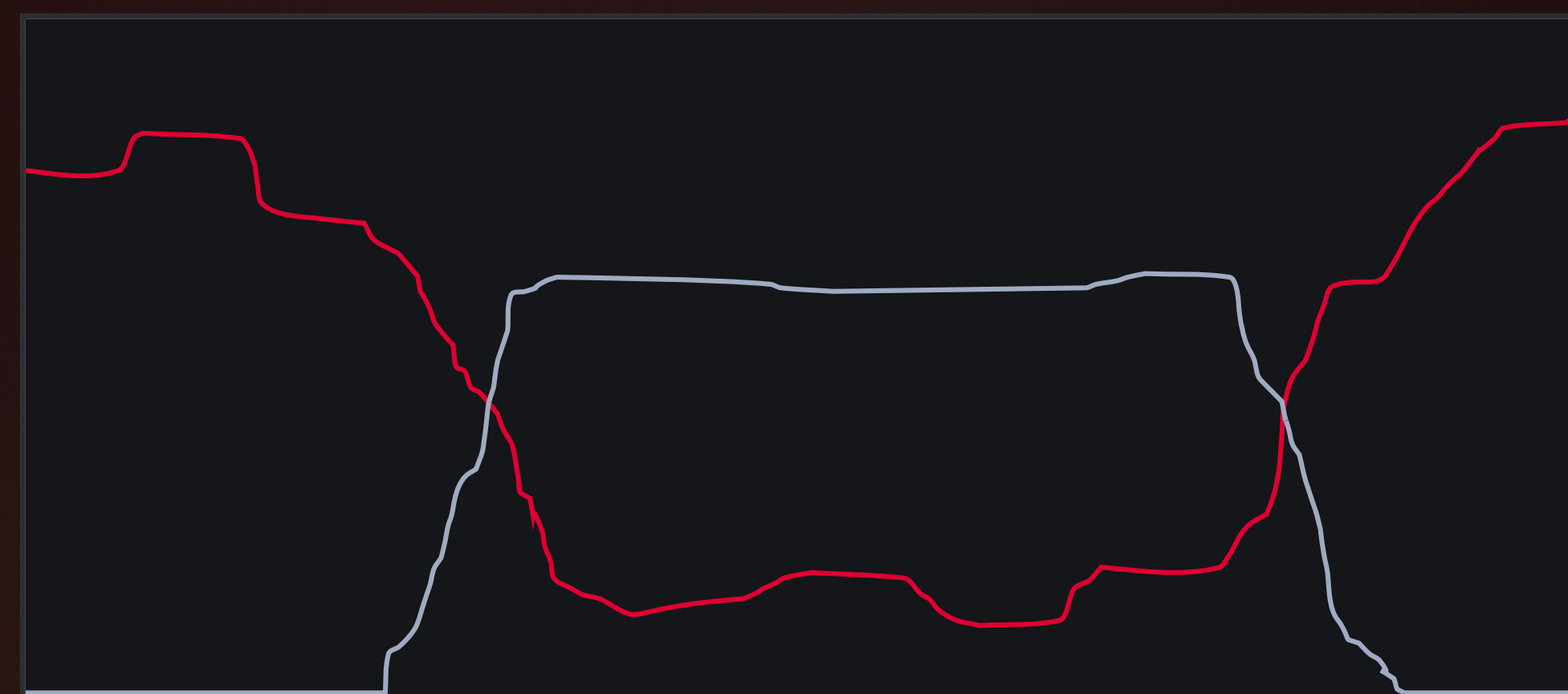
ACCELERATED PERFORMANCE  
HIGH PRIORITY COMPUTE



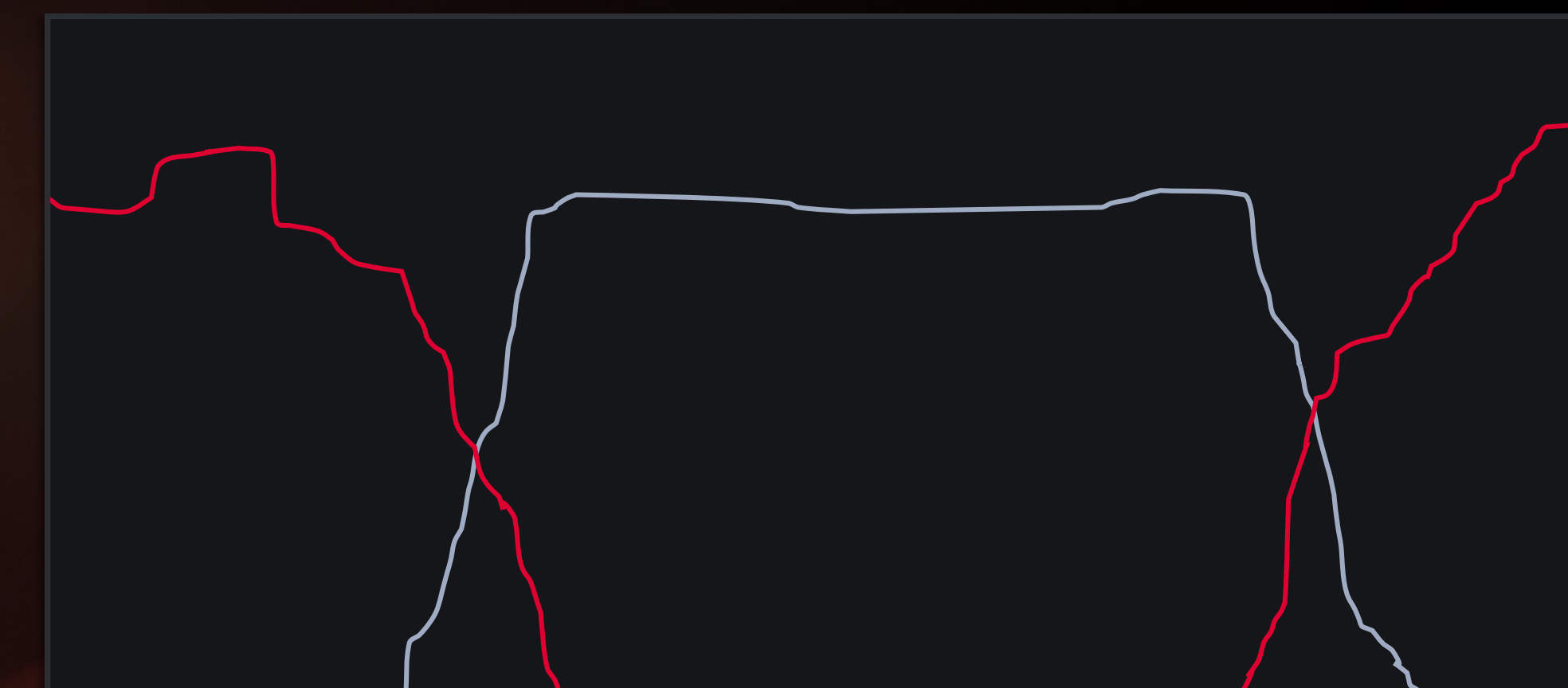
ENABLING SMOOTH VR

## ASYNC COMPUTE TUNNELING

PRECISE CONTROL OF OTHER WORK IN FLIGHT



*Current Solution*



*With Tunneling*

Stall other pipeline launch and complete draining of other pipeline shader waves

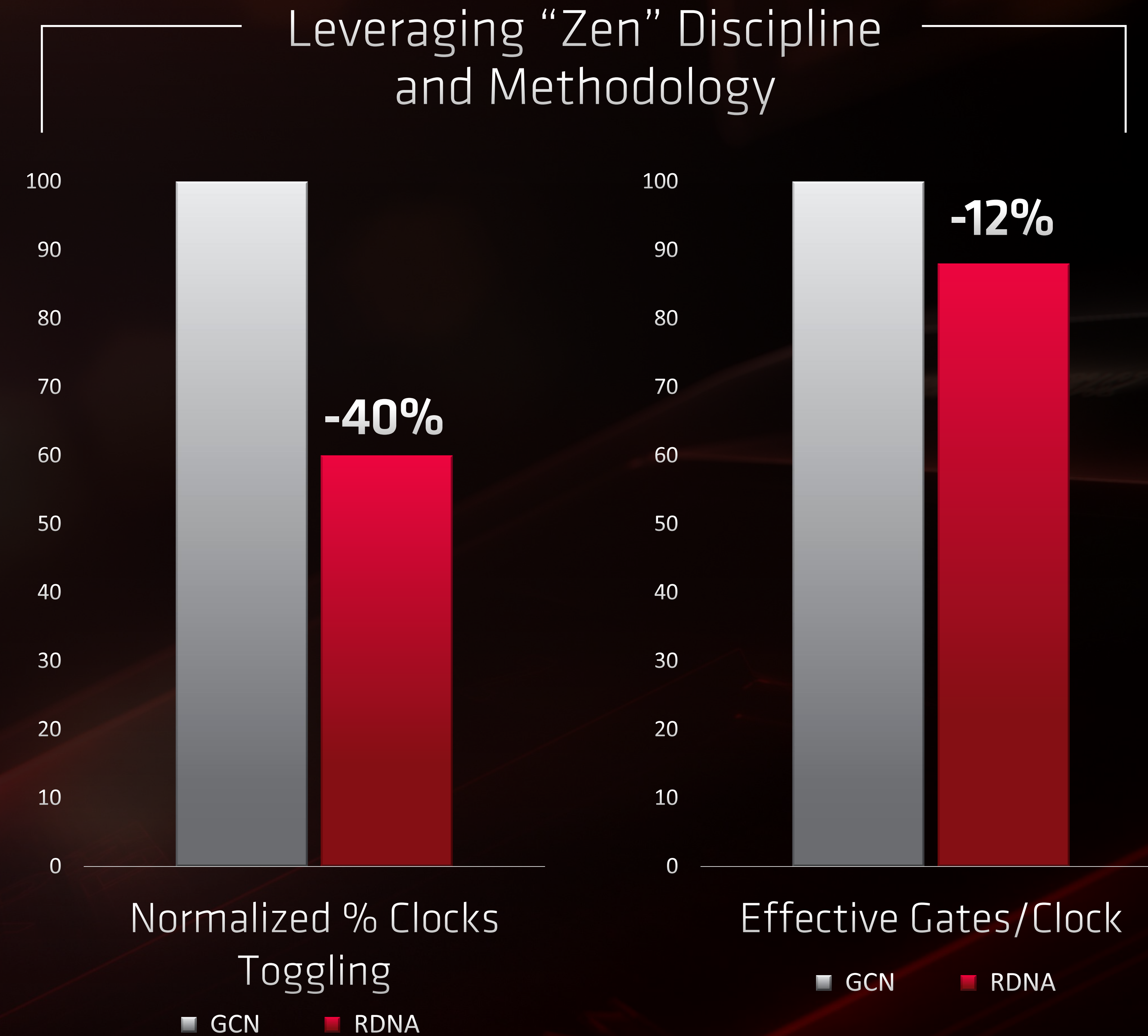
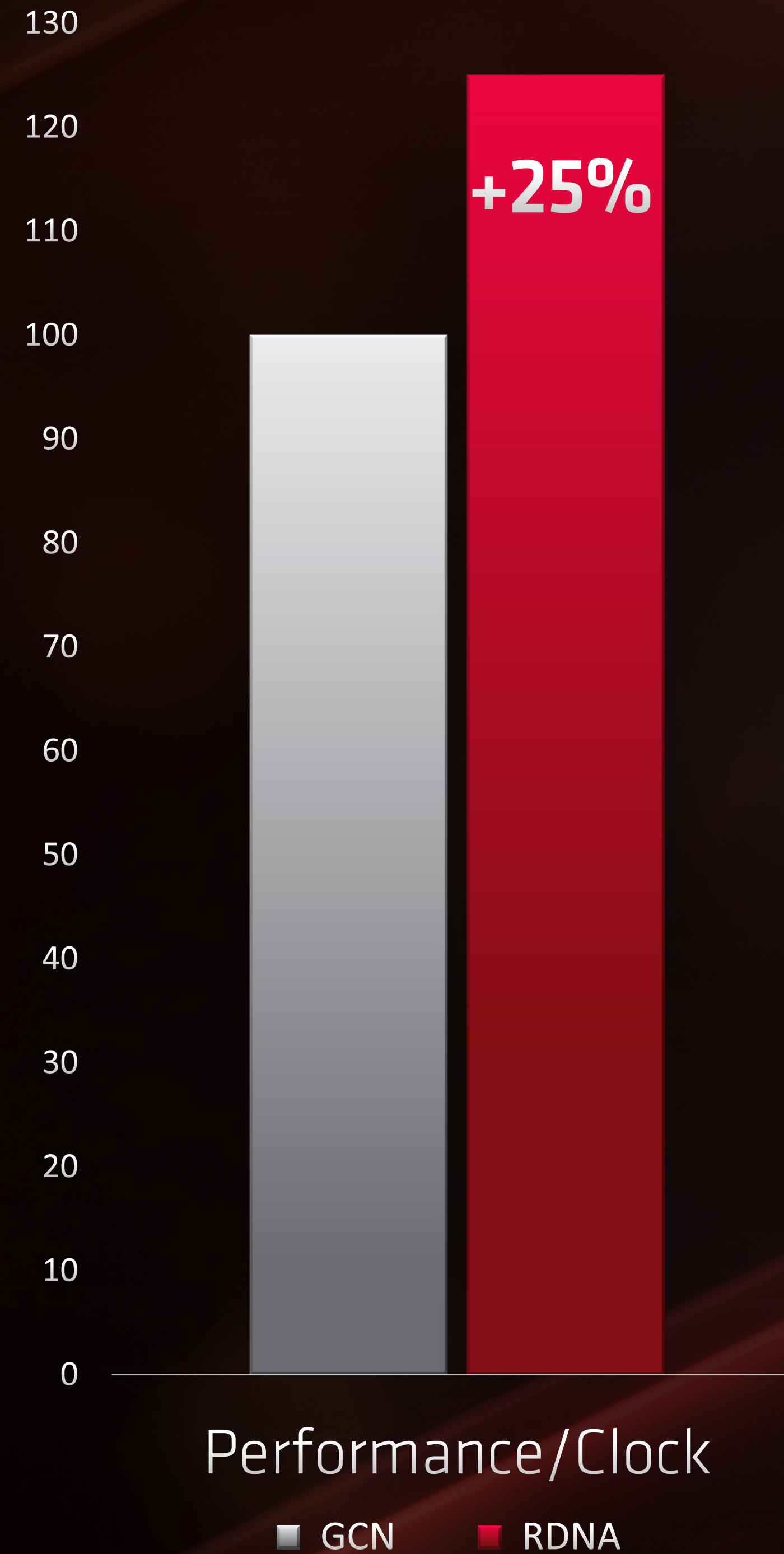


# STREAMLINED GRAPHICS PIPELINE

Improved Architectural Efficiency for Performance

Hyper-Effective Clock Gating for Power Efficiency

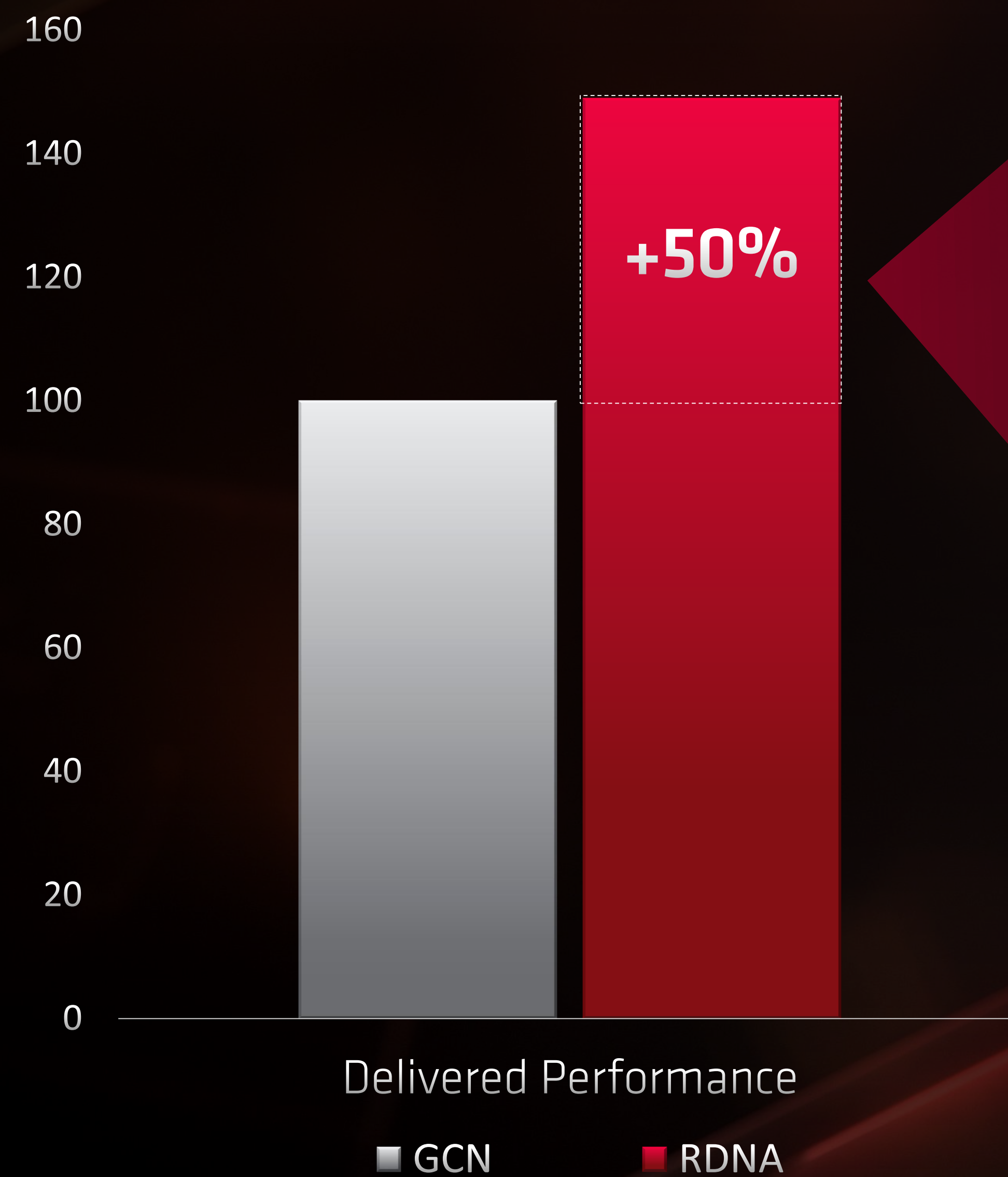
Reduced Levels of Logic For Higher Frequency



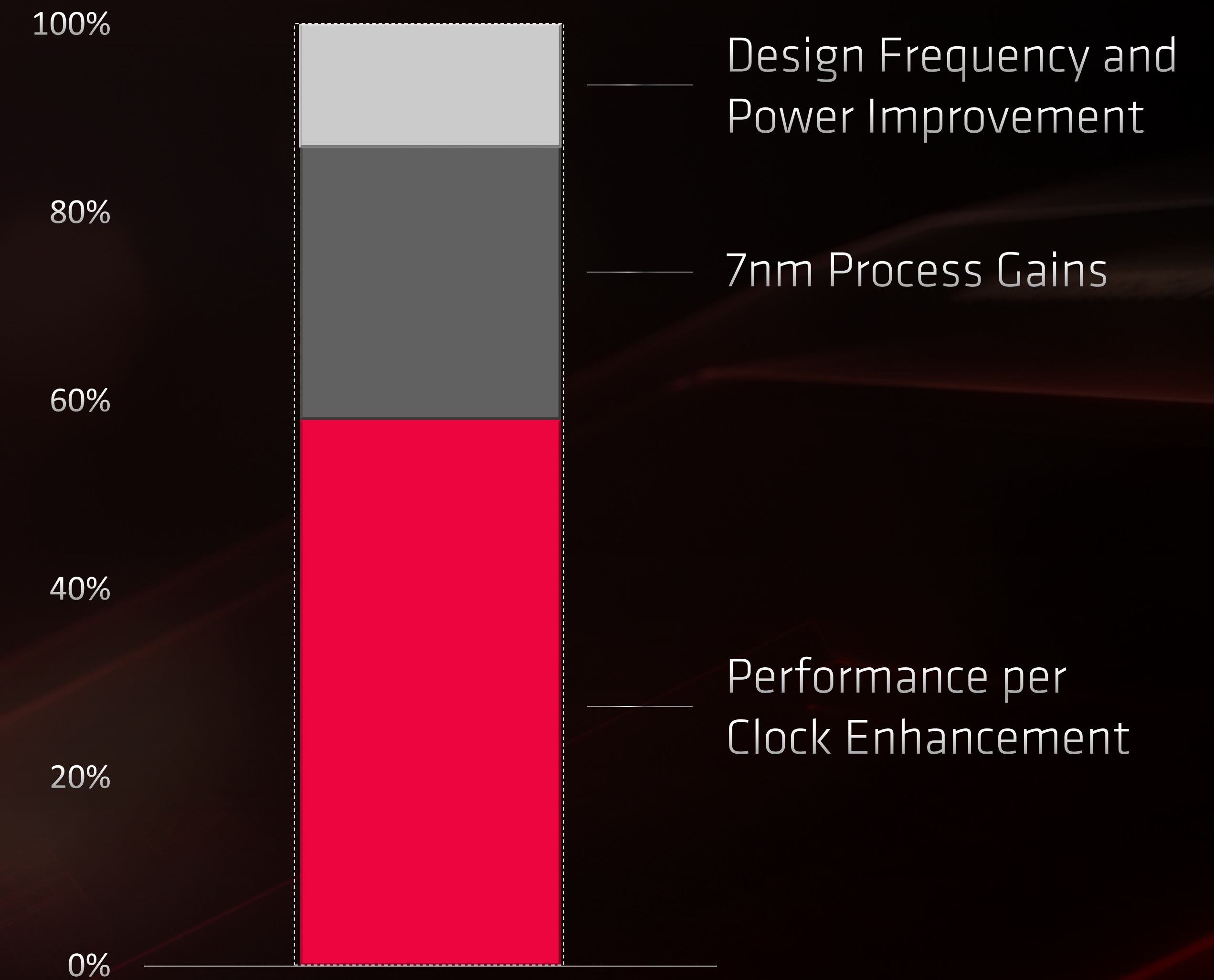


# RDNA DELIVERS

## SAME POWER, GREATER PERFORMANCE



## RDNA PERFORMANCE CONTRIBUTORS



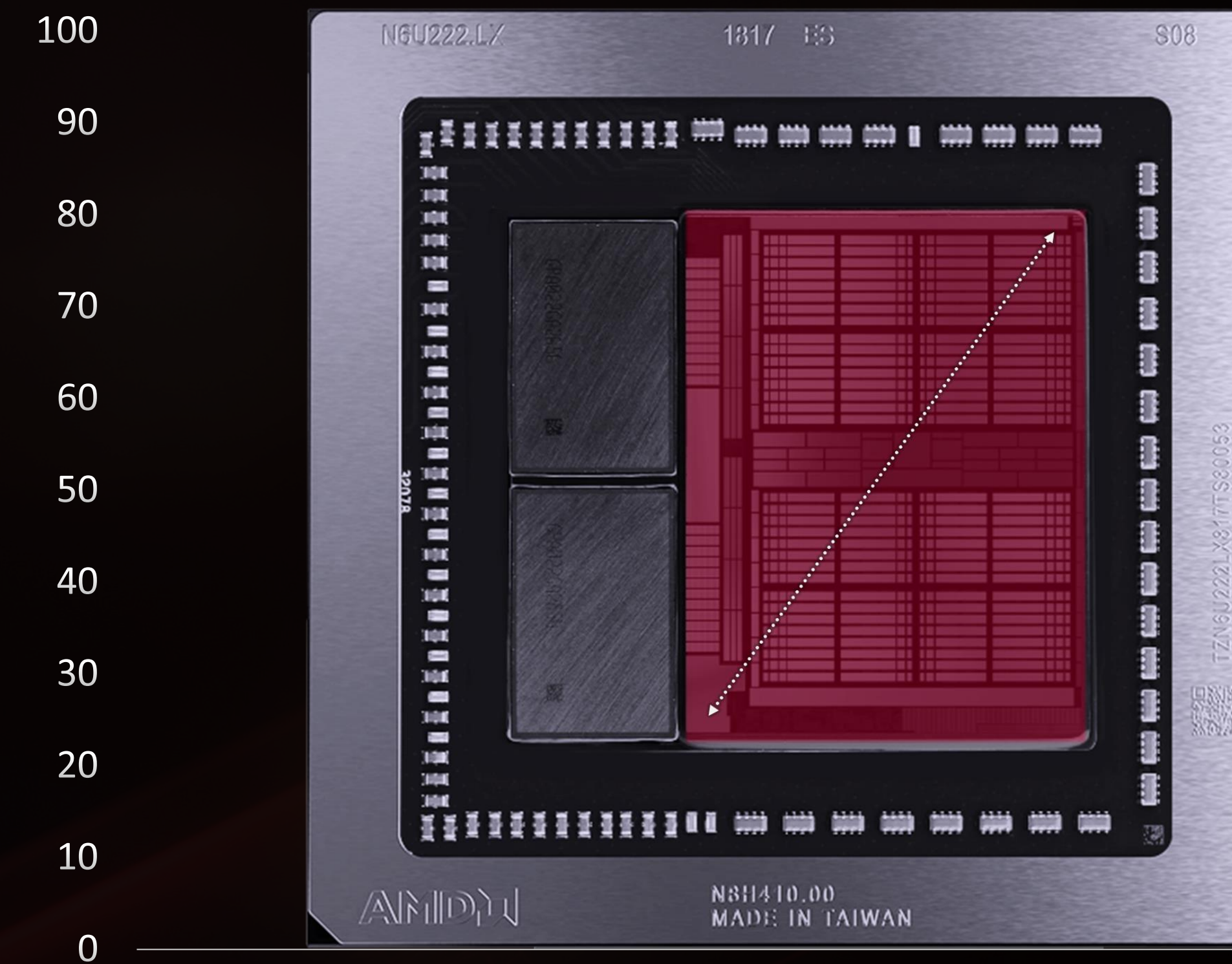
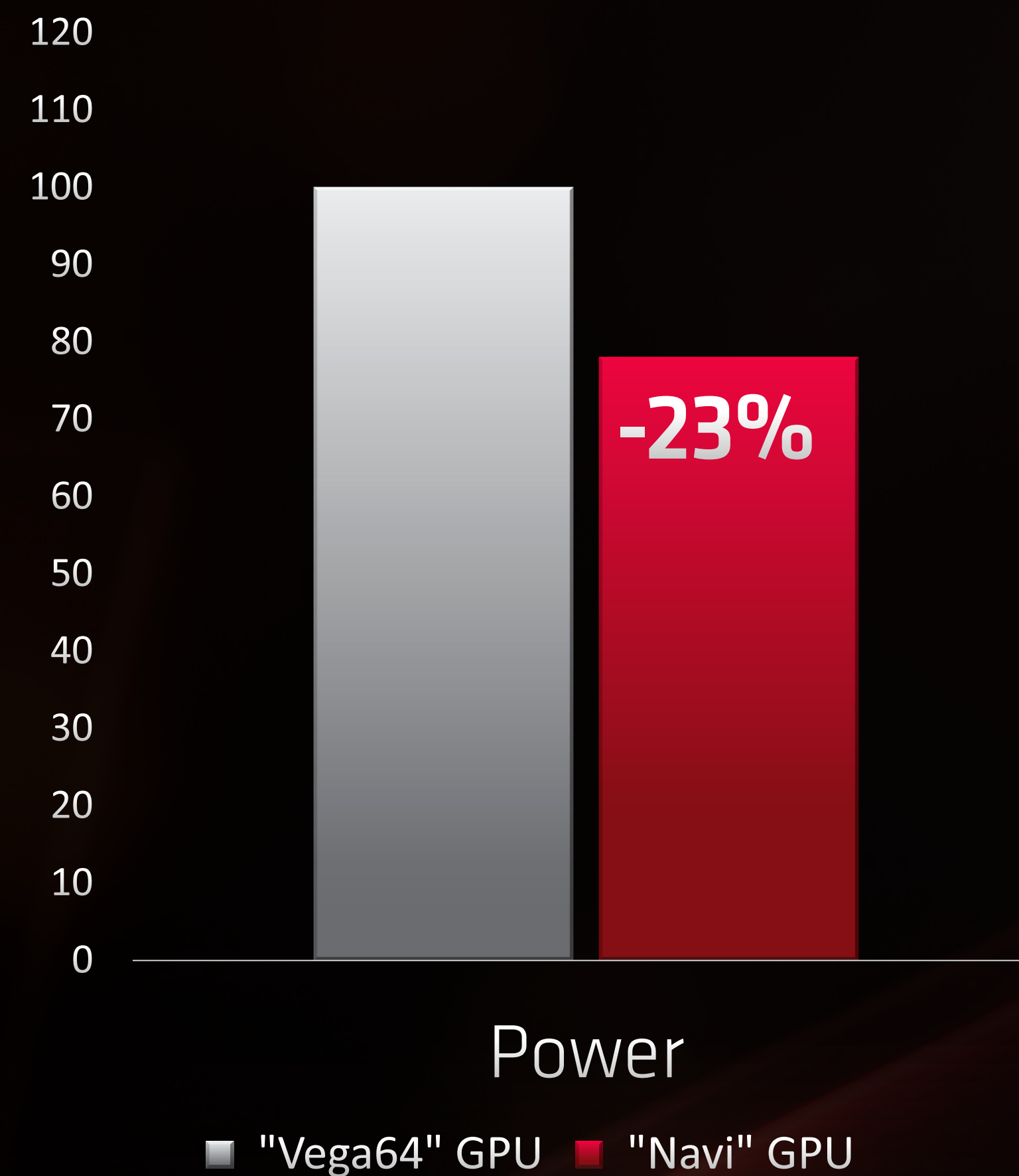
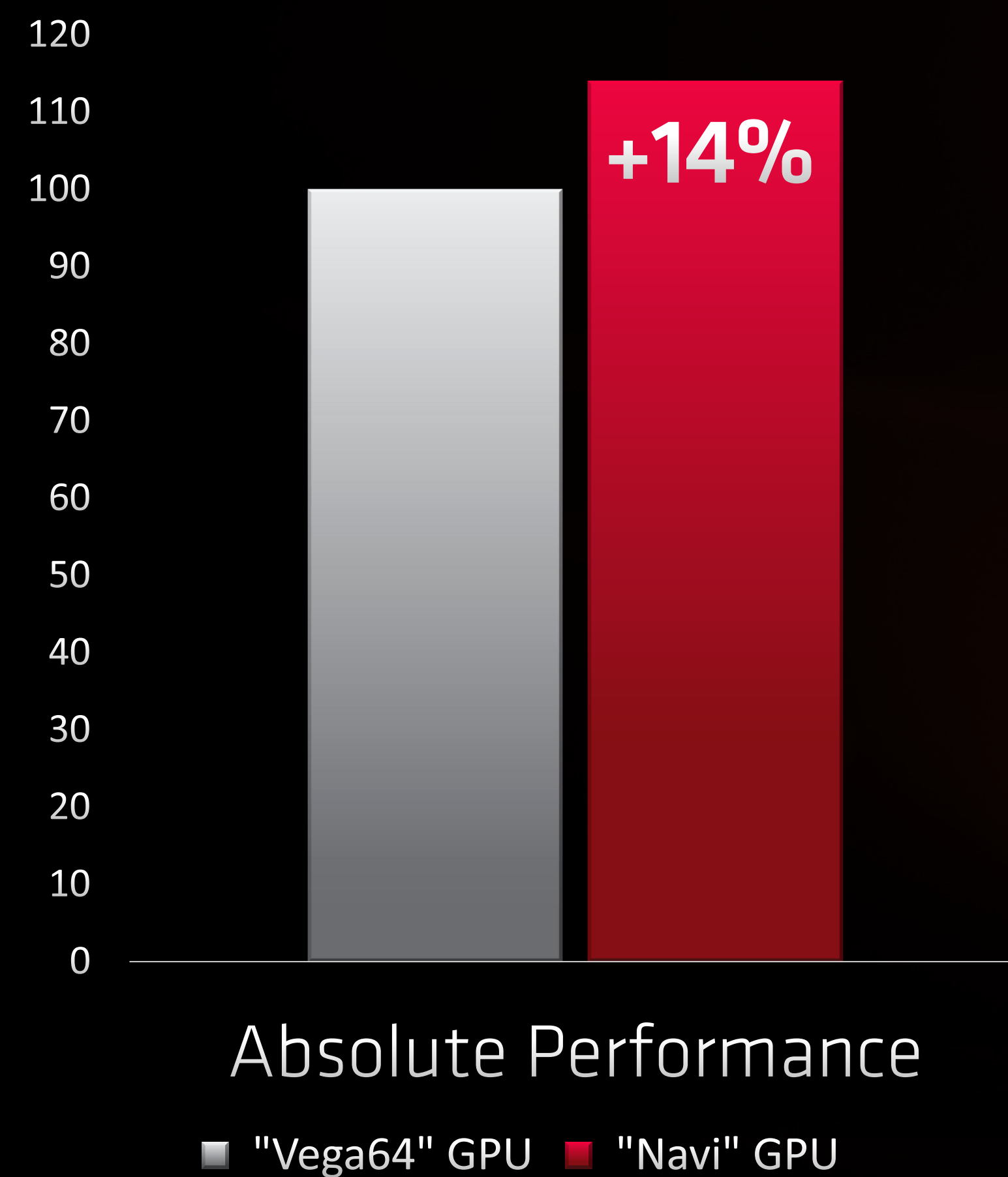


# RDNA

THE FOUNDATION FOR GREAT PRODUCTS

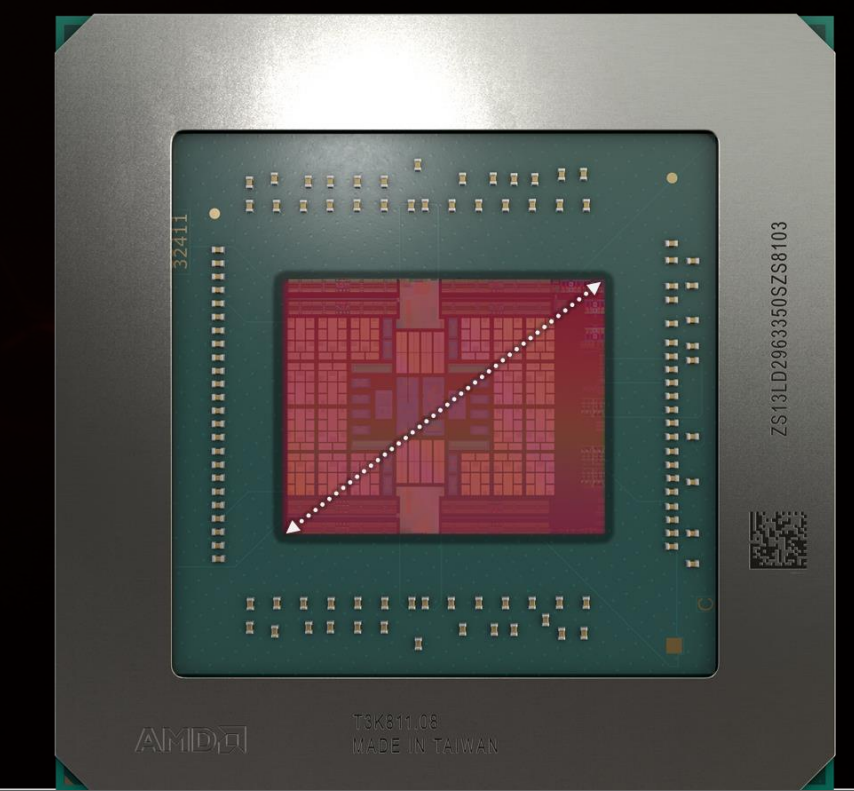
**1.5X** PERFORMANCE PER WATT

**2.3X** PERFORMANCE PER AREA



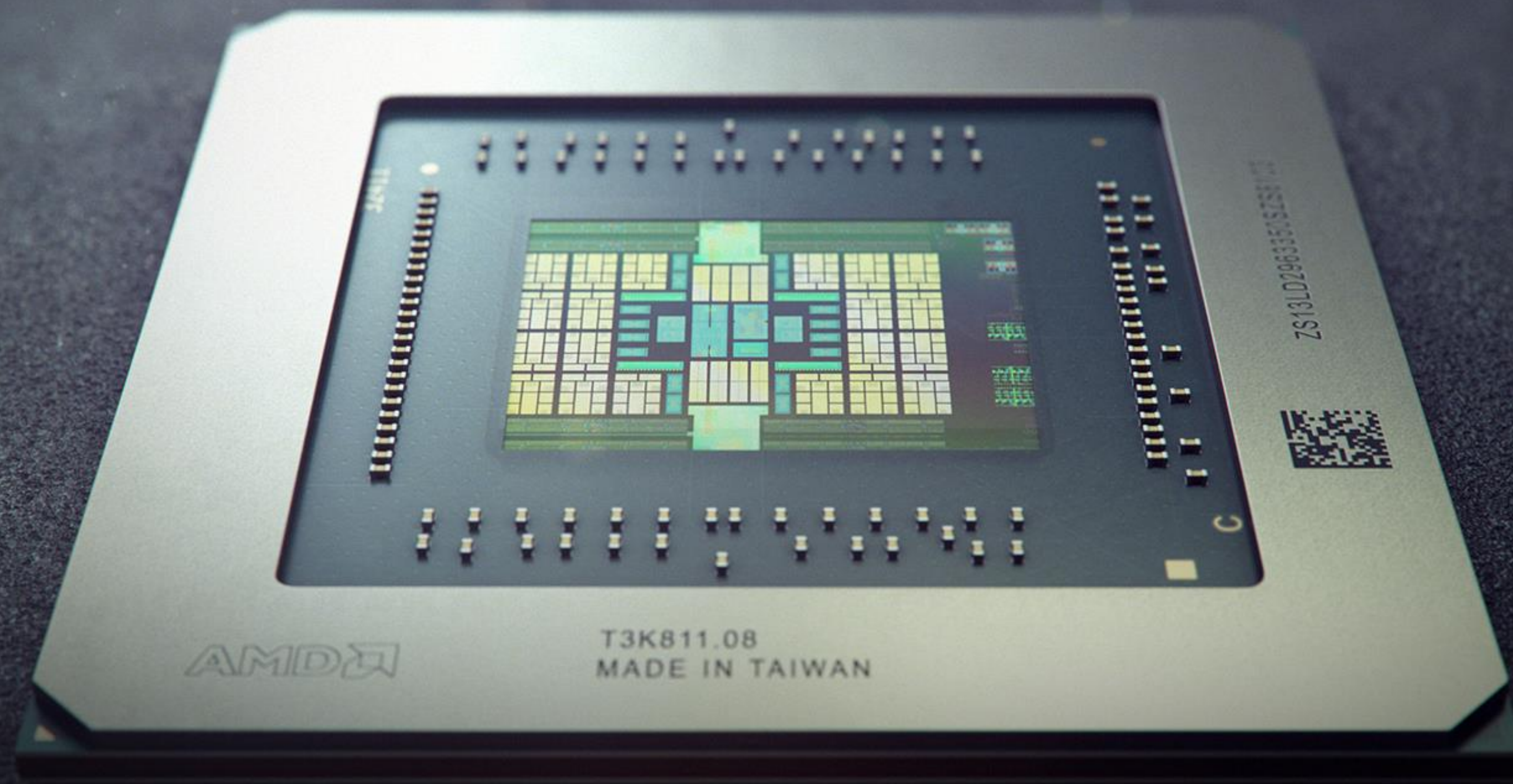
**495mm<sup>2</sup>**  
**14nm "VEGA 10"**

Die Size



**251mm<sup>2</sup>**  
**7nm "NAVI"**





THE ALL NEW  
“NAVI” GPU FAMILY FEATURES

**RDNA**  
Architecture

**7nm**  
Process

**GDDR6**  
Memory

**PCIe<sup>®</sup> 4.0**  
Support

**RADEON**  
Media Engine

**RADEON**  
Display Engine



# AMD RADEON™ GAMING IS EVERYWHERE

## RDNA: EXPANDING THE RADEON UNIVERSE



PCs



Macs



Consoles



Cloud



Mobile



# ENDNOTES

---

## **RX-325**

Testing done by AMD performance labs 5/23/19, using the Division 2 @ 25x14 Ultra settings. Performance may vary based on use of latest drivers. RX-325

## **RX-327**

Testing done by AMD performance labs 5/23/19, showing a geomean of 1.25x per/clock across 30 different games @ 4K Ultra, 4xAA settings. Performance may vary based on use of latest drivers. RX-327

## **RX-329**

Testing conducted by AMD Performance Labs as of 05/30/2019 on Radeon RX 5700XT with AMD Driver 19.10 (1902270946) on Intel i7-6900k, and on Radeon Vega Frontier Edition with AMD Driver 19.30 (1904231814) on Intel i7-5960k. Both systems used 2x8GB DDR4 2133Mhz RAM, Asus ROG Rampage V Edition Motherboard, and Windows 10 Enterprise. Performance may vary. RX-329.

## **RX-358**

Testing done by AMD performance labs on June 4 2019. Systems were tested with: Intel(R) Core(TM) i7-5930K CPU @ 3.50GHz (6 core) with 16GB DDR4 @ 2133 MHz using a Asus X99-E Motherboard running Windows 10 Enterprise 64-bit (Ver. 1809, build 17763.053). Using the following graphics cards: Navi (Driver 19.30\_1905161434 (CL# 1784070)) with 40 compute units, versus a Vega 64 (Driver 19.4.1) with 40 compute units enabled. Running 3D Mark 11 GT1 (1280 x 720), 3D Mark 11 GT2 (1280 x 720), 3d Mark Firestrike GT1 (2560 x 1440), 3d Mark Firestrike GT2 (2560 x 1440), UnigineHeaven (1920 x 1080) , the Navi (with a die size of 251mm^2) achieved an average FPS score of 140 , 136, 49, 37, and 84 respectively. Compared to the Vega 56 (with a die size of 486mm^2) which achieved 103, 113, 41, 32, and 72 respectively. RX-358

## **GD-81**

HEVC (H.265), H.264, and VP9 acceleration are subject to and not operable without inclusion/installation of compatible HEVC players. GD-81

## **GD-127**

Radeon FreeSync technology requires a monitor and AMD Radeon™ graphics, both with FreeSync support. See [www.amd.com/freesync](http://www.amd.com/freesync) for complete details. Confirm capability with your system manufacturer before purchase. GD-127

## **DISCLAIMER**

The information contained herein is for informational purposes only, and is subject to change without notice. While every precaution has been taken in the preparation of this document, it may contain technical inaccuracies, omissions and typographical errors, and AMD is under no obligation to update or otherwise correct this information. Advanced Micro Devices, Inc. makes no representations or warranties with respect to the accuracy or completeness of the contents of this document, and assumes no liability of any kind, including the implied warranties of noninfringement, merchantability or fitness for particular purposes, with respect to the operation or use of AMD hardware, software or other products described herein. No license, including implied or arising by estoppel, to any intellectual property rights is granted by this document. Terms and limitations applicable to the purchase or use of AMD's products are as set forth in a signed agreement between the parties or in AMD's Standard Terms and Conditions of Sale.

## **ATTRIBUTION**

© 2019 Advanced Micro Devices, Inc. All rights reserved. AMD, the AMD Arrow logo, Radeon, Ryzen, FreeSync, and combinations thereof are trademarks of Advanced Micro Devices, Inc. in the United States and/or other jurisdictions. Other names are for informational purposes only and may be trademarks of their respective owners.