# Lunar Lake Architecture Session

HotChips 2024

**Arik Gihon**

Senior Principal Engineer, CCG
SoC Architecture

# Lunar Lake

## Flagship SoC for the next gen of AI PCs

| **Breakthrough x86 power efficiency** | **Exceptional core performance** | **Massive leap in graphics** | **Unmatched AI compute** |
|---|---|---|---|
| Up to 40% lower SoC power* | Similar ST perf at half the power* | Up to 1.5X better graphics* | Up to 120 platform TOPS |

*Versus Intel's previous gen. For more information, go to Intel.com/PerformanceIndex.

# Lunar Lake Construction

## Built with advanced packaging

Memory

Compute tile

Platform Controller tile

Package

Base tile

Foveros

Stiffener

Memory

Compute tile

Platform Controller tile

Package

Filler tile

Base tile

Foveros

# Memory on Package

First ever Intel integration onto package
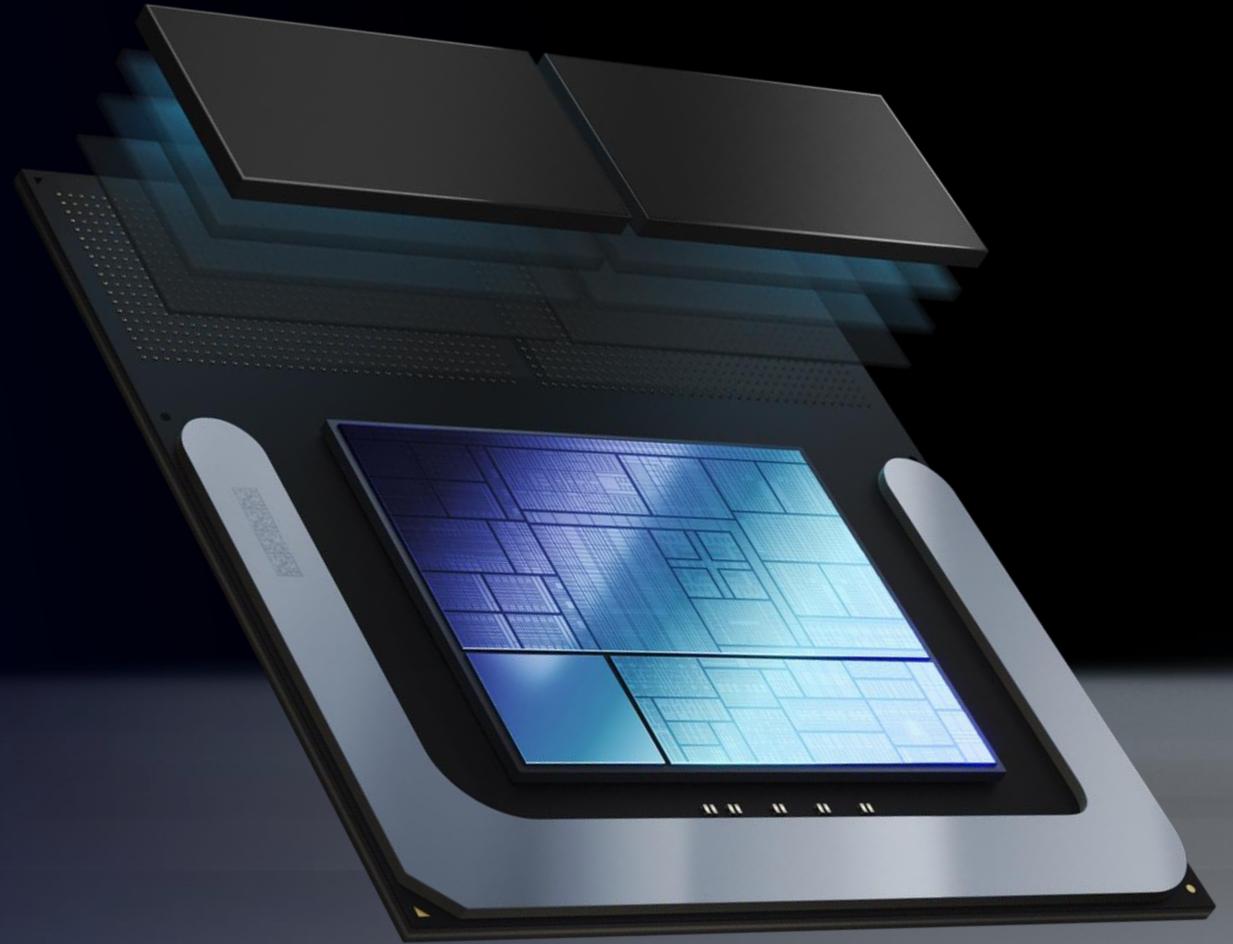
Up to
**32GB**
with 2 ranks

Support for
**LPDDR5x**
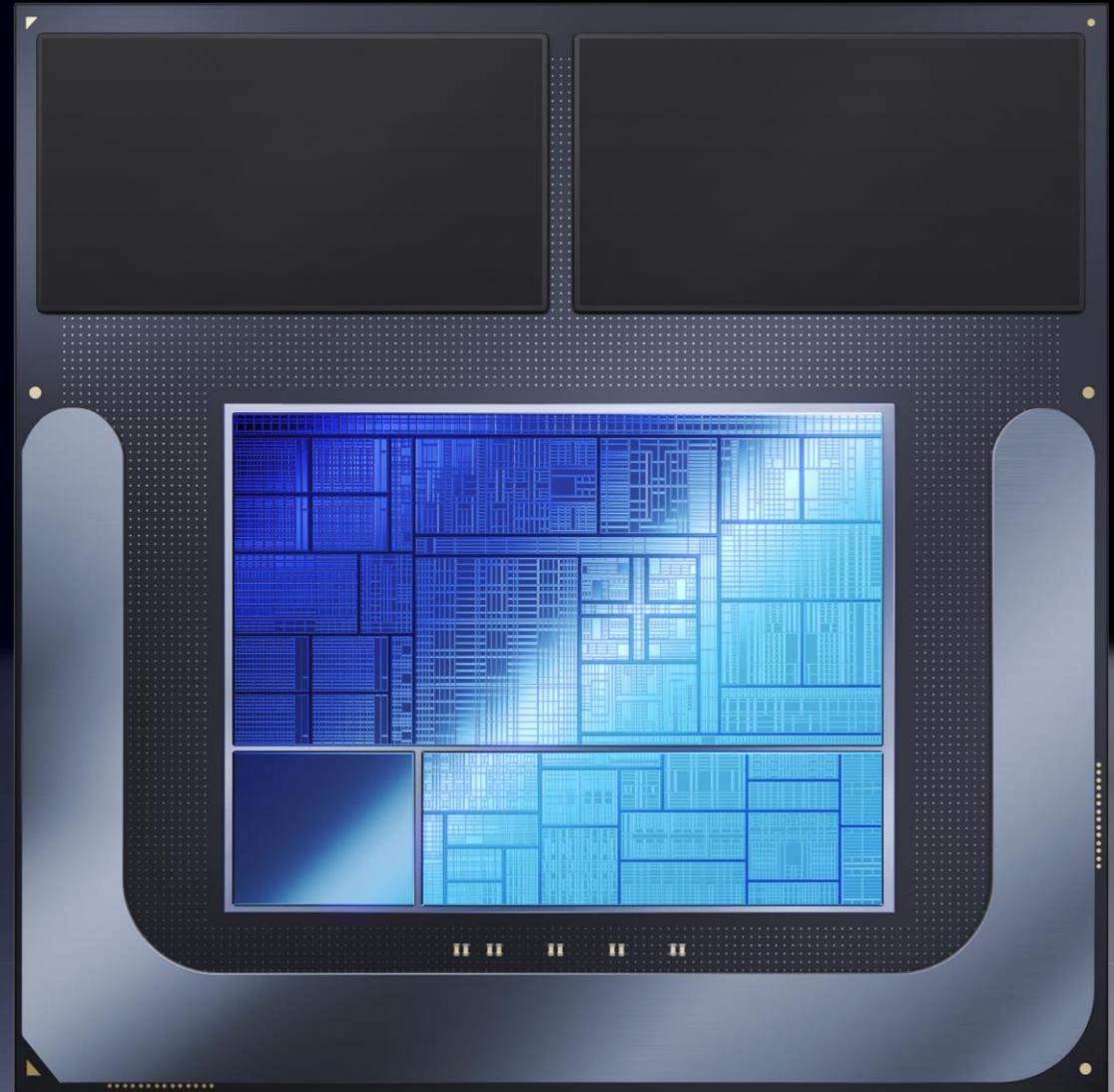DRAM

Up to
**8.5GT/s**
per chip

Support for
**16b x4**
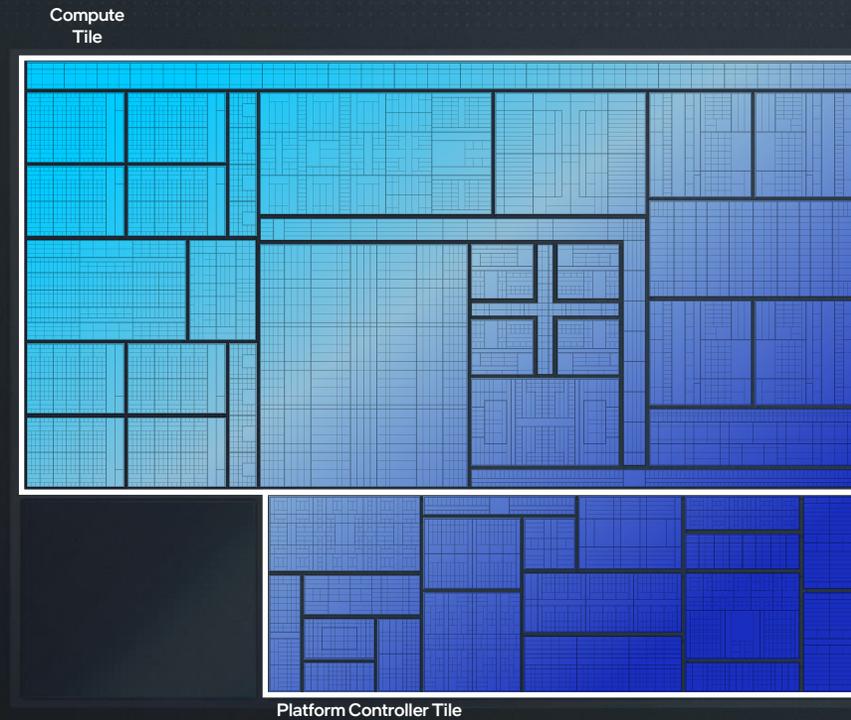channels

Up to
**40%**
lower PHY power

Up to
**250mm²**
area savings

# Lunar Lake

## Architecture overview

Lunar Lake

# Architectural Framework

Compute Tile

Platform Controller Tile

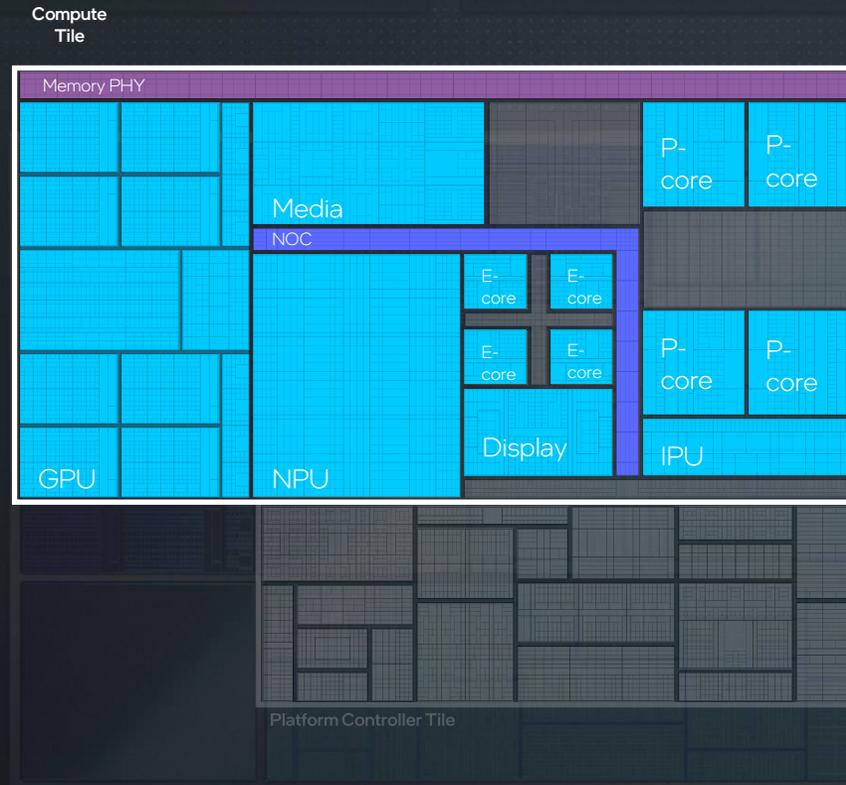- Compute tile structure
- Memory side cache
- Enhanced E-core cluster
- New power delivery & management

# Lunar Lake

# Enhanced SoC Structure

For better performance efficiency

**Compute Tile**

Memory PHY

Media

NOC

GPU

NPU

E-core

E-core

E-core

E-core

Display

IPU

P-core

P-core

P-core

P-core

Platform Controller Tile

**Compute tile structure**

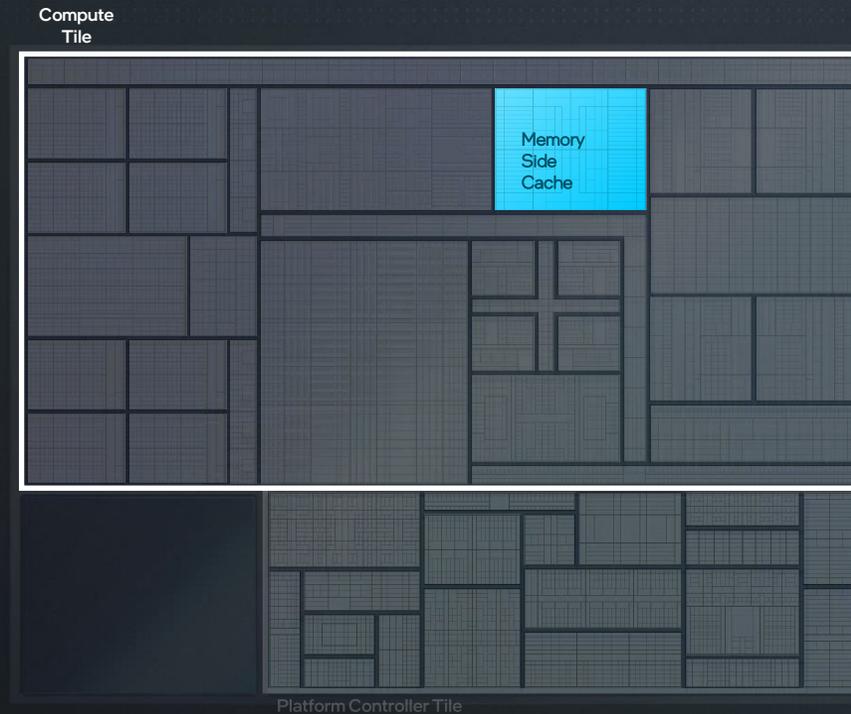Monolithic die, on a leading process

New NOC, with enhanced caching

Optimized memory latency

# Lunar Lake

# Memory Side Cache

For efficient performance & lower power consumption for other engines

**Compute Tile**

Memory Side Cache

**Platform Controller Tile**

**Memory side cache**

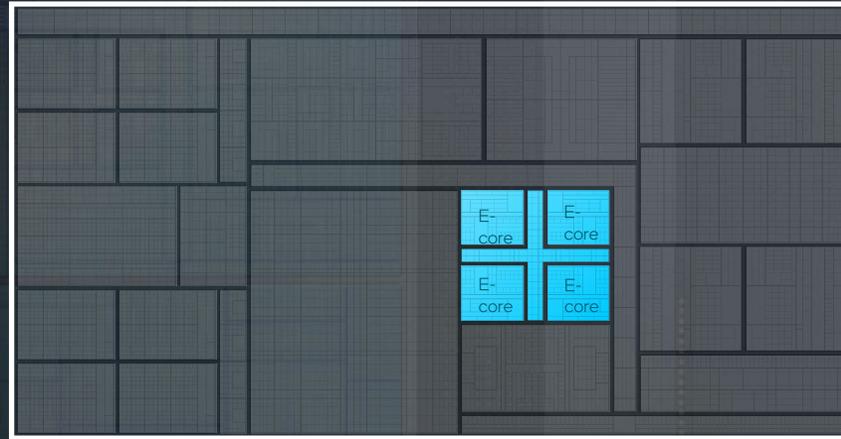8MB physical cache

Reduce DRAM traffic and power

Enhanced both latency and bandwidth

Caching for IO engines

Lunar Lake

# Enhanced E-core Cluster

For more efficient performance & lower power for almost every use[1]

[1]vs. previous gen low power island

Compute Tile

E-core  E-core
E-core  E-core

Platform Controller Tile

Enhanced E-core cluster

Doubled the core count

Increased 4MB L2 cache

Leading-edge process for frequency and power

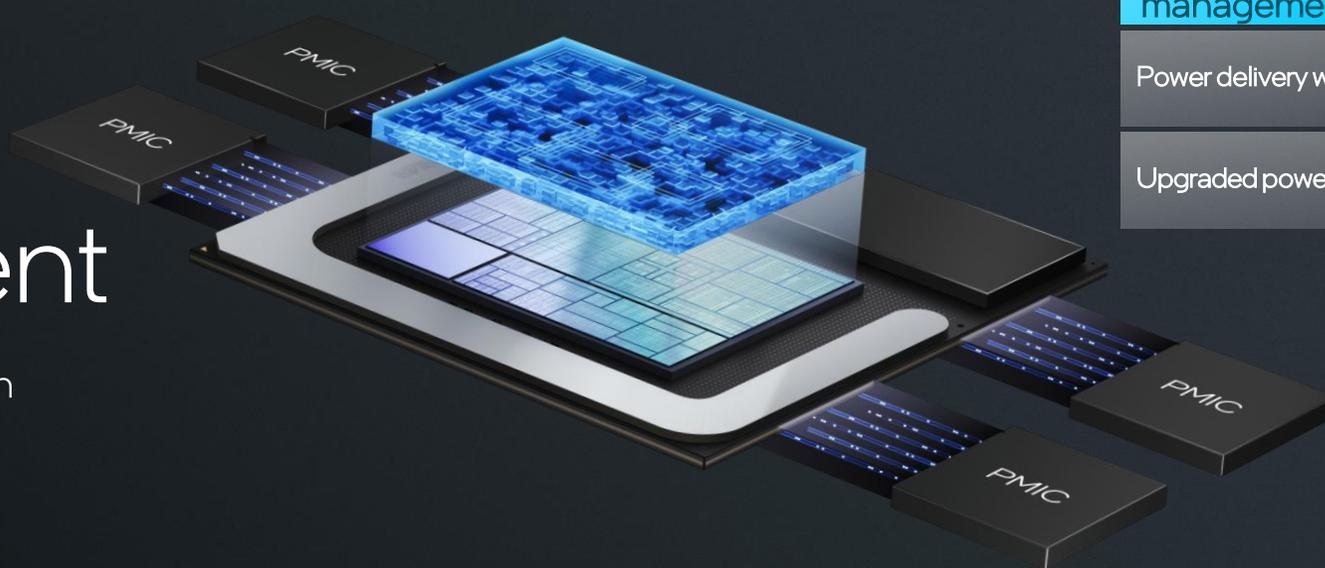Mem side cache for power and latency

Independent power delivery

Lunar
Lake

# Power Delivery & Management

For optimal SOC power utilization
& performance efficiency

New power delivery &
management

Power delivery with 4 PMICs

Upgraded power management

Lunar
Lake

# Power Delivery & Management
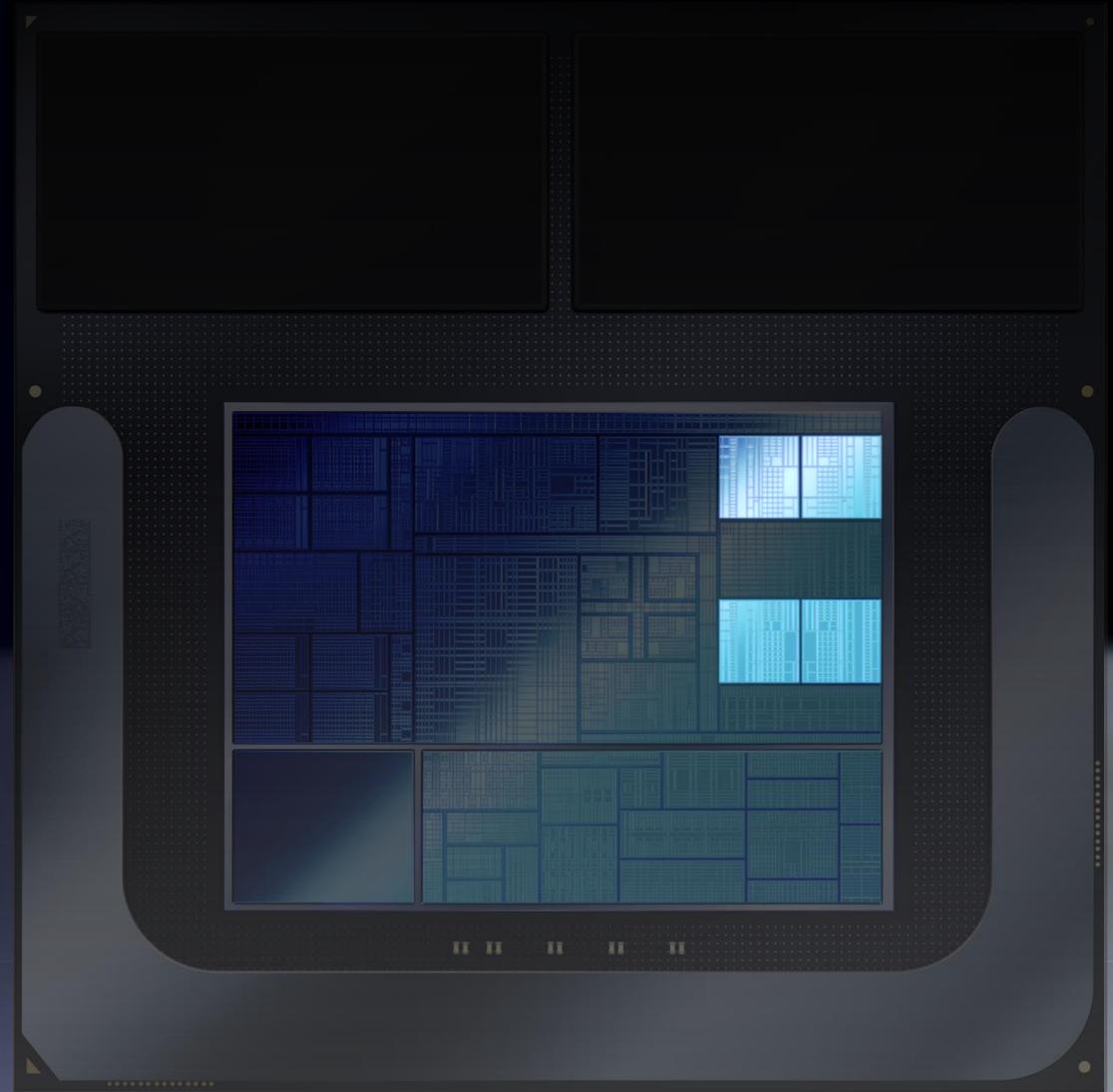
For optimal SOC power utilization & performance efficiency

**New power delivery & management**

Power delivery with 4 PMICs

- More power rails
- Dynamic voltage ID
- More telemetry (IMON)

Lunar
Lake

# Power
# Delivery &
# Management

For optimal SOC power utilization
& performance efficiency

**Power Management**

**Thread Director**

New power delivery & management

Power delivery with 4 PMICs

Upgraded power management

Intel Thread Director focus on efficiency

Optimized power balancer per load type

Enhanced "sleep" states power and latency

ML-based WL classification & frequency control
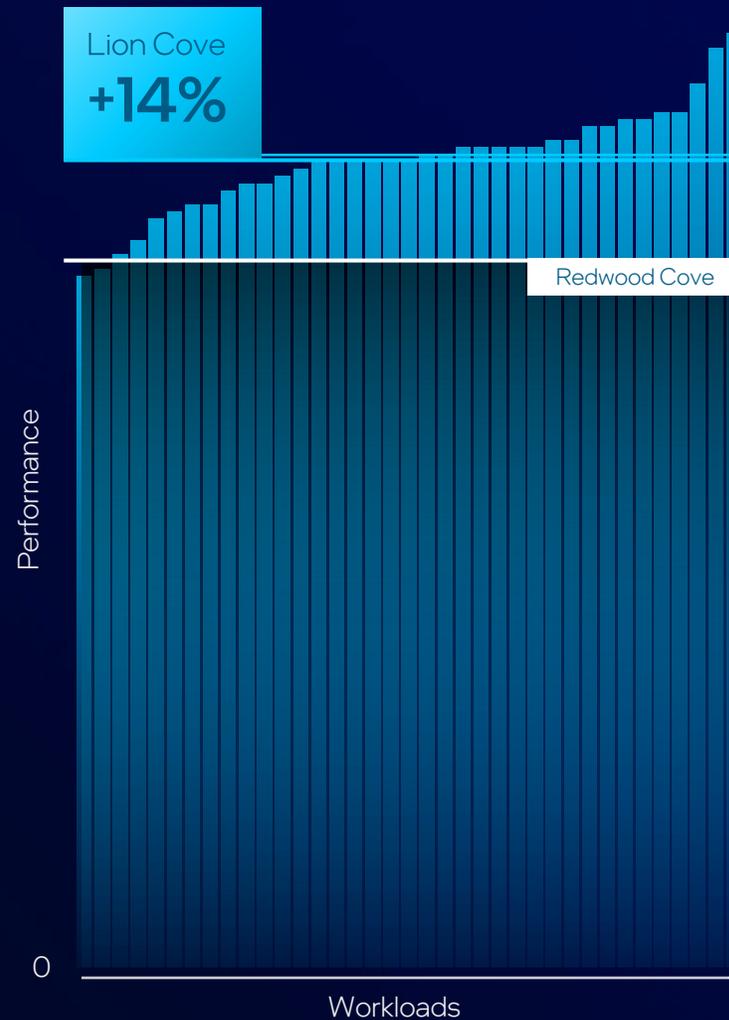
Lunar Lake
Performance cores

# Lion Cove
# P-core

Double-digit performance gains over prior generation

Lion Cove in Lunar Lake

Redwood Cove in Meteor Lake

**IPC**

Lion Cove **+14%**

Redwood Cove

Performance

0

Workloads

Iso frequency benefit estimate across: Components of SPECrate2017_int_base and SPECrate2017_fp_base (both estimated) running 1 copy, Cinebench R23 Single Core, Cinebench 2024 Single Core, Geekbench5.4.5 Single-Core, Geekbench6.2.1 Single-Core, WebXPRT 4, Speedometer
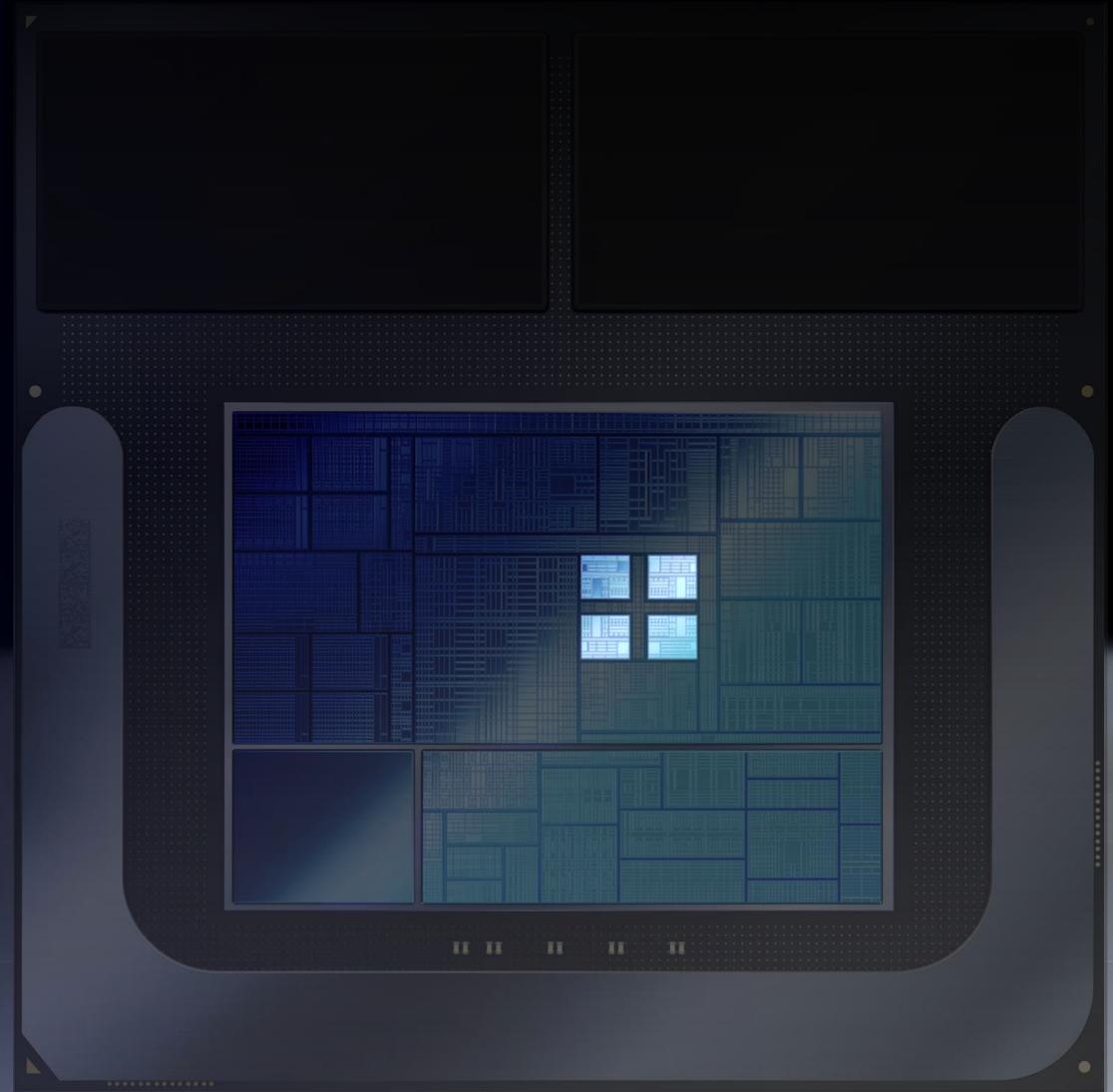
**Performance at power**

>10%

>12%

>15%

>18%

Performance

Power

Results are based on SPECrate2017_int_base (estimated) running n copies. Based on measurement on an Intel internal reference validation platforms at a fixed PL1 power

# Lunar Lake

Efficient cores

# Skymont IPC Gain

On Lunar Lake Low Power Island

vs. Meteor lake LP E-core

**Skymont** plus **Lunar Lake** fabric improvements combine to bring significant IPC improvements for many general workloads

SPEC CPU 2017 Rate
GCC12.1 -O2 Linux
Fixed Frequency (iso)

Results are based on Intel's internal projections/estimates (+/- 10% margin of error). See Intel.com/PerformanceIndex for additional details.

## ST integer improvement

SPECrate2017_int_base est / GCC

Skymont
**1.38x**

MTL LP E-core

Performance

Workloads

## ST FP improvement

SPECrate2017_fp_base est / GCC

Skymont
**1.68x**

MTL LP E-core

Performance

Workloads

# Skymont Power & Performance

On Lunar Lake Low Power Island

vs. Meteor lake LP E-core

**Significant increase in workload coverage** for the Lunar Lake Low Power Island

**More IPC, frequency & cores**

SPECrate2017_int_base est
GCC12.1-O2 Linux (iso)

Results are based on SPECrate2017_int_base (estimated) running n copies. Based on measurement on an Intel internal reference validation platforms at a fixed PL1 power

## Single threaded

Performance

Skymont
E-core cluster
(Lunar Lake)

**1/3** Power

**1.7x** Perf

**2x** Perf

**LP E-core** (Meteor Lake)

Power

# Latency Optimizations



Load to Use Latency in Memory

**Measured in out-of-the-box configuration (unfixed frequencies), using a pointer chasing latency test with 4KB Page size on Windows 11 at 15W-17W**

Latency Delta between MTL and LNL

# Efficient Cores BW & Cross Cluster latencies
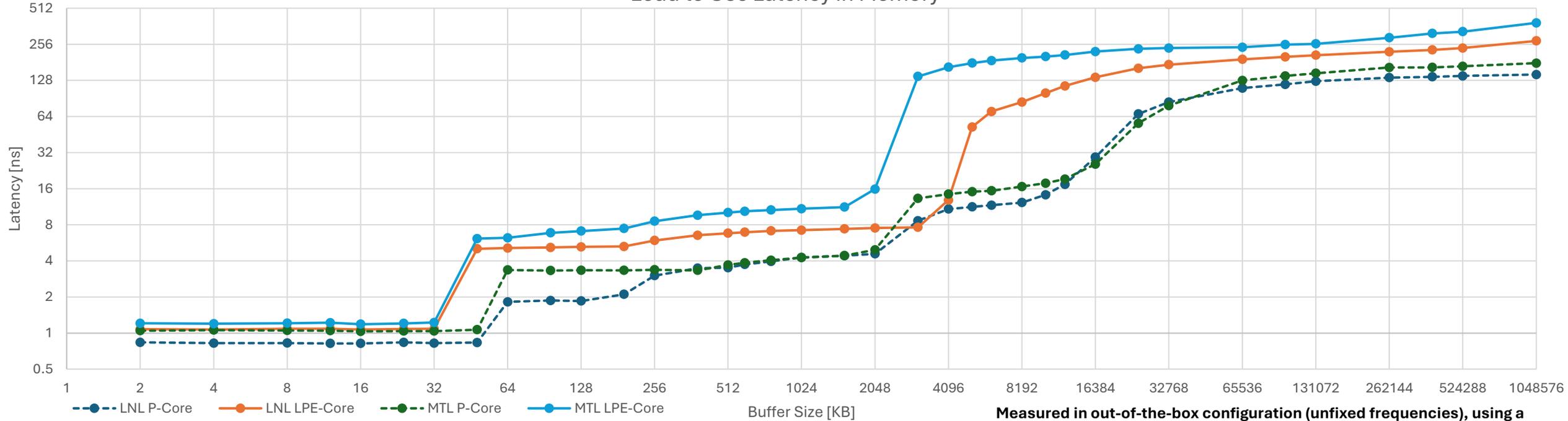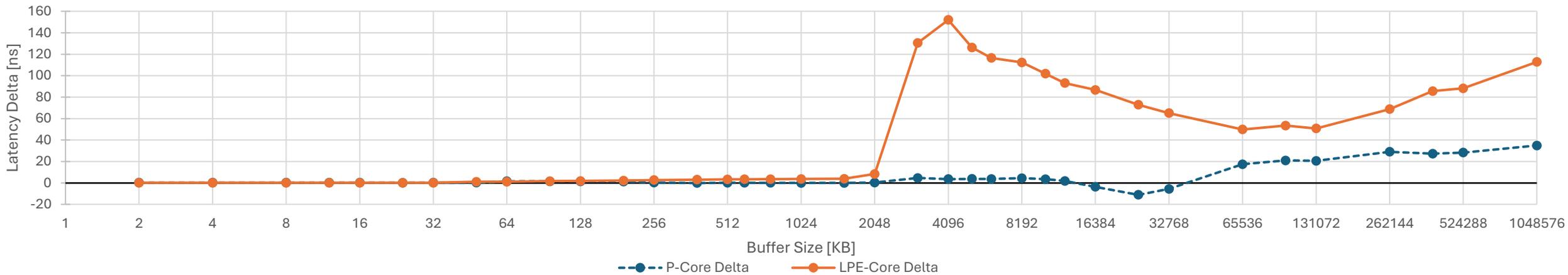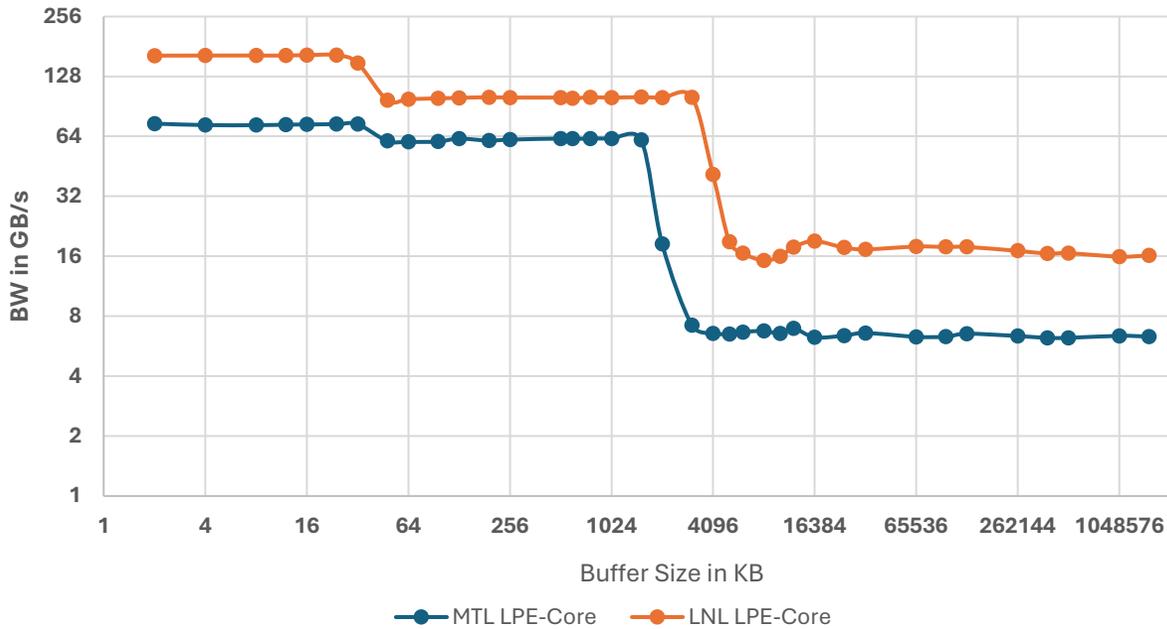
## 2.8x Memory Bandwidth to LP Efficient Cluster

### LP E-Core Memory Bandwidth



**Legend:** MTL LPE-Core, LNL LPE-Core

Y-axis: BW in GB/s (1, 2, 4, 8, 16, 32, 64, 128, 256)

X-axis: Buffer Size in KB (1, 4, 16, 64, 256, 1024, 4096, 16384, 65536, 262144, 1048576)

As part of the capable efficient Cores cluster that can keep the Compute Cluster powered off more often

## Low Latency Cross-Cluster Coherency

| | | Compute Cluster Cores | | | | Low Power Island Cores | | | |
|---|---|---|---|---|---|---|---|---|---|
| **Compute Cluster Cores** | | | 24.8 | 24.65 | 26 | 54.35 | 55.05 | 54.7 | 54.3 |
| | 23.85 | | | 27.25 | 27.45 | 54.65 | 55.2 | 55.6 | 54.1 |
| | 25.25 | 27.05 | | | 28.2 | 55.2 | 55.15 | 55.3 | 55.2 |
| | 25 | 27.5 | 27.9 | | | 56 | 55.15 | 55.3 | 56 |
| **Low Power Island Cores** | 55.05 | 54.45 | 55.4 | 55 | | | 23.7 | 23.6 | 22.7 |
| | 54.55 | 54.85 | 55.5 | 55.35 | 22.75 | | | 23.7 | 23.2 |
| | 55.1 | 55.05 | 55.55 | 55.15 | 23.4 | 23.6 | | | 23.5 |
| | 55.15 | 55.1 | 54.9 | 56.3 | 22.9 | 22.95 | 23.9 | | |

~55 ns on LNL , 3X better vs Prev Gen

**Low overhead for shared data in MT Workloads in both clusters enables efficient Cores cluster to also accelerate Multi-Threaded applications.**

# Breakthrough x86 Efficiency
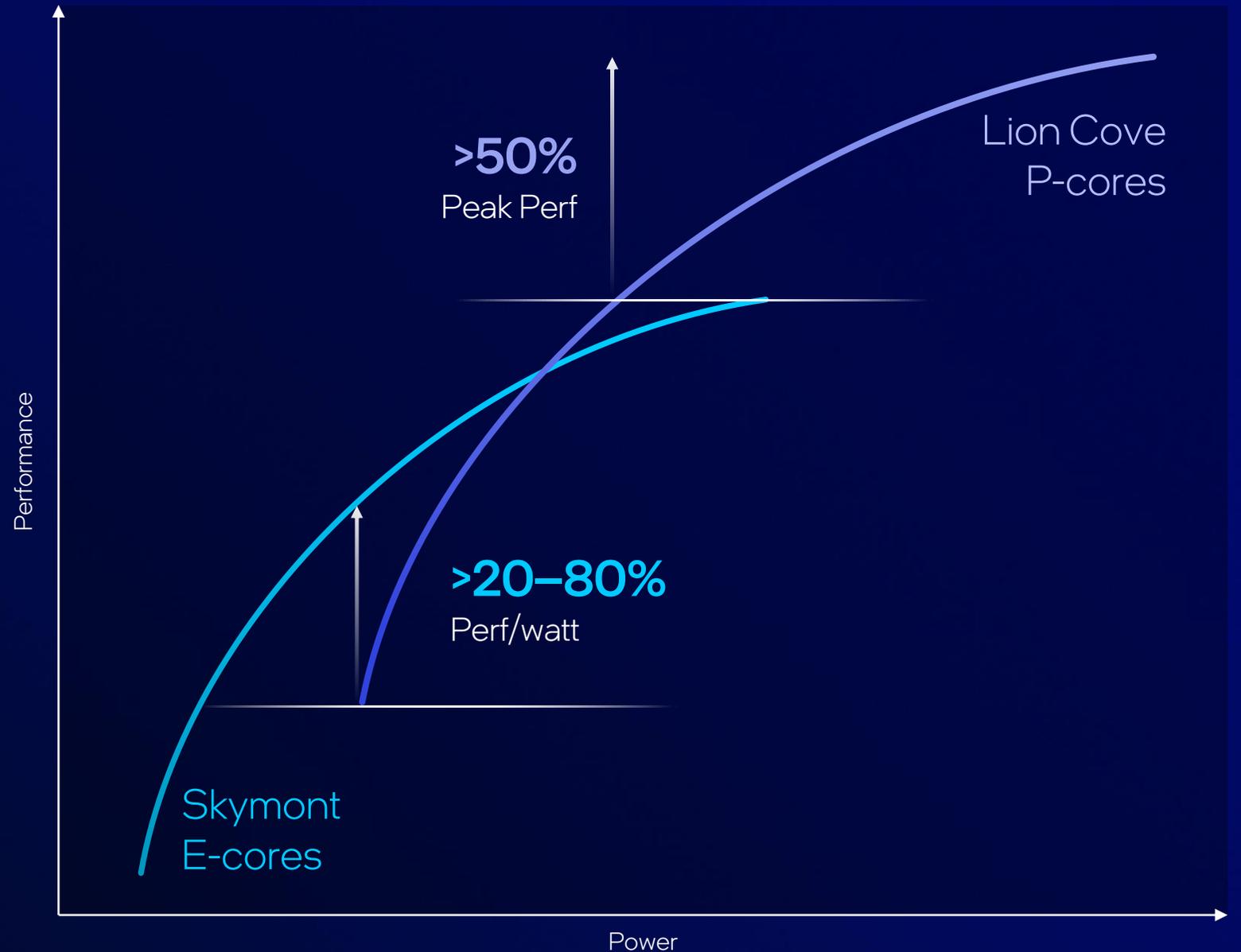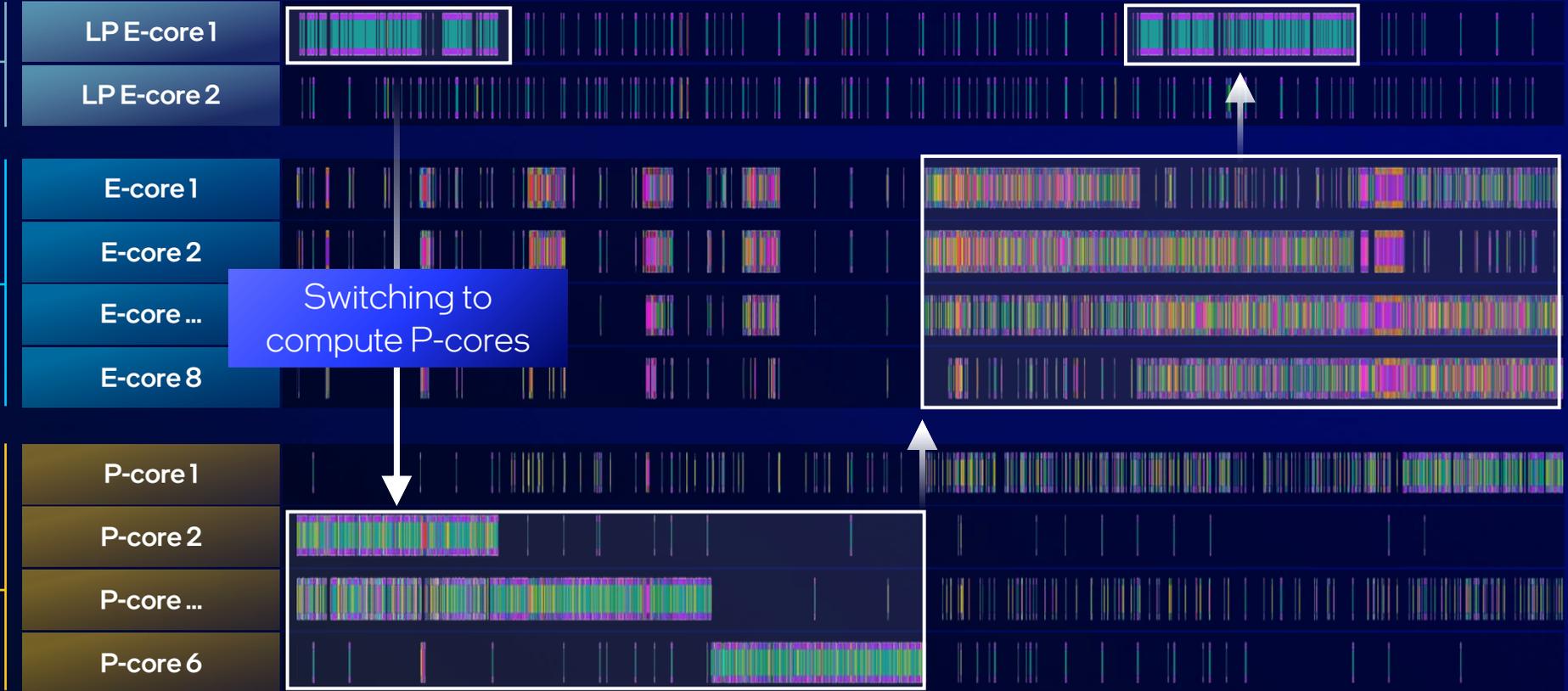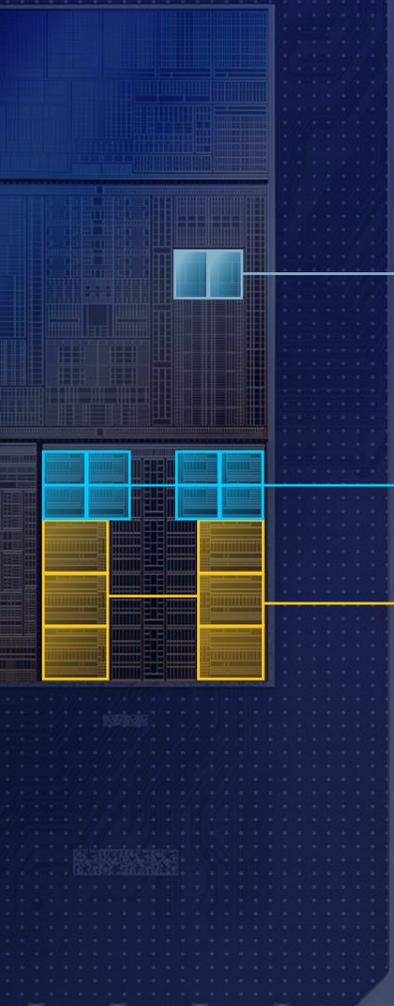
Whilst covering the full CPU performance range



**Performance** (vertical axis)

**Power** (horizontal axis)

**Lion Cove P-cores**

**>50%** Peak Perf

**>20–80%** Perf/watt

**Skymont E-cores**

Illustration of the relative Lunar Lake P-core and E-core performance across the SoC power range.
See Intel.com/PerformanceIndex for more details.

Lunar Lake

New X$^e$ 2 GPU

# Next Gen
# Xᵉ2 GPU

## Architecture goals

**Improved utilization**
of hardware functions

**Improved distribution**
of workload across architecture

**Improved integration**
Of hardware & software

# Next Gen
# Xᵉ2 GPU
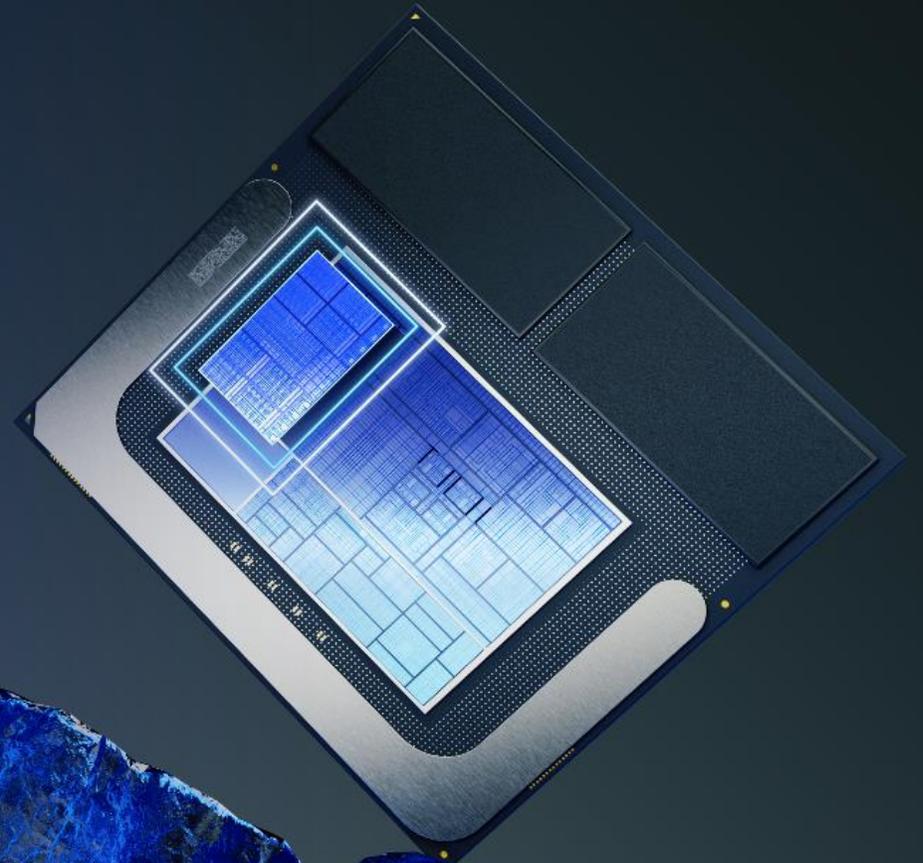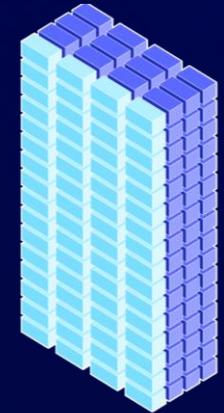
Major leap in graphics performance

up to
**67** TOPS

**New XMX engines**

**8** Larger ray tracing units

**Ray Tracing Unit**

BVH Cache

Traversal Pipeline | Traversal Pipeline | Traversal Pipeline

Box Int. | Box Int. | Box Int.

Xᵉcore

XVE | X M X | X M X
XVE | X M X | X M X
XVE | X M X | X M X
XVE | X M X | X M X

Load / Store

I$ | L1$ / SLM

**8** 2ⁿᵈ gen Xᵉ cores

Xᵉ2 vector engines

Xᵉ2 Vector Engine

Thread Control | Register File | FP | INT | EM | FP64 | XMX

Branch

Send

**1.5x** better vs. Meteor Lake GPU

intel ARC™
Software stack

eDP 1.5

Enhanced XᵉSS kernels

XᵉSS | Warp

Upscaling

**8** MB L2 cache

intel.

# Next Gen
# Xe2 GPU

Major leap in graphics performance

**~1.5x**

vs. previous gen

Xe2
Lunar Lake

Performance

**Higher Perf**
@iso power

**Lower Power**
@iso perf

Meteor Lake H

**Higher Perf**
@iso power

Meteor Lake U

Power

# Stable Diffusion



## Lunar Lake

```
C:\Windows\System32>call "C:\openvino\setupvars.bat"
Warning: Python is not installed. Please install one of Python 3.8 - 3.12 (64-bit) from https://www.python.org/
downloads/
[setupvars.bat] OpenVINO environment initialized
Press any key to continue . . .
OpenVINO version: OpenVINO Runtime
    Version : 2024.3.0
    Build   : 2024.3.0-15502-66093834e38

Loading and compiling text encoder: 2127.68 ms
Loading and compiling UNet: 6917.08 ms
Loading and compiling VAE decoder: 997.38 ms
Loading and compiling tokenizer: 82.82 ms
Started to run Stable Diffusion pipeline
image #1 : _
```
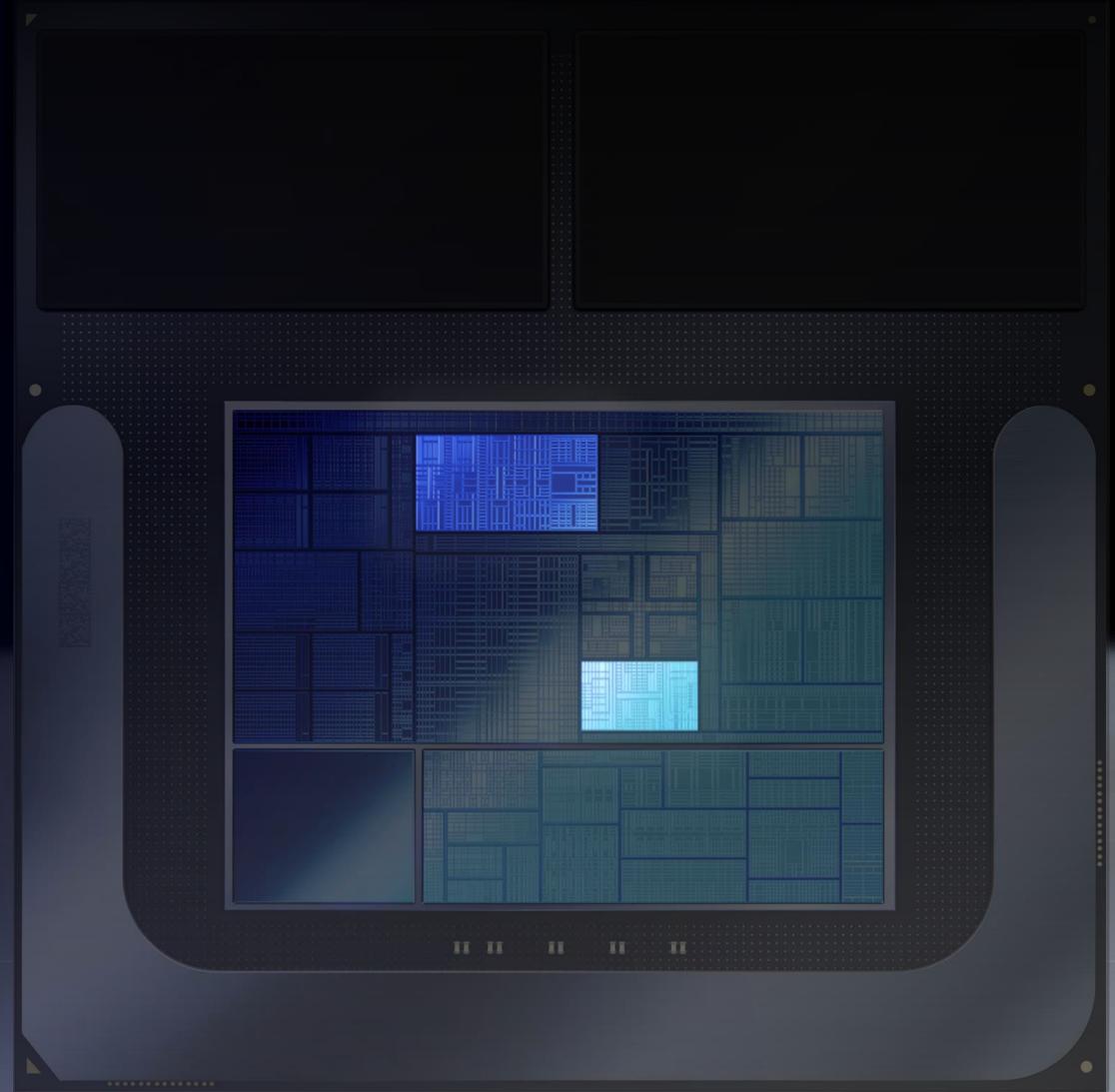
## Meteor Lake

```
C:\Windows\System32>call "C:\openvino\setupvars.bat"
Warning: Python is not installed. Please install one of Python 3.8 - 3.12 (64-bit) from https://www.pyt
downloads/
[setupvars.bat] OpenVINO environment initialized
Press any key to continue . . .
OpenVINO version: OpenVINO Runtime
    Version : 2024.3.0
    Build   : 2024.3.0-15502-66093834e38

Loading and compiling text encoder: 2459.51 ms
Loading and compiling UNet: 8908.48 ms
Loading and compiling VAE decoder:
```

"Older Intel Engineer wearing a blue hat with a blue jacket,
he is making a funny face with his mouth open, excited for new technology."

# Lunar Lake

## Media & display engines

# New Media & Display Engines

## Media engine

**AV1** — Encode & decode

**VVC** — Decode

## Display engine

**DisplayPort** — 1x eDP 1.5

**DisplayPort** — DisplayPort 2.1

**HDMI** HIGH-DEFINITION MULTIMEDIA INTERFACE™ — HDMI 2.1
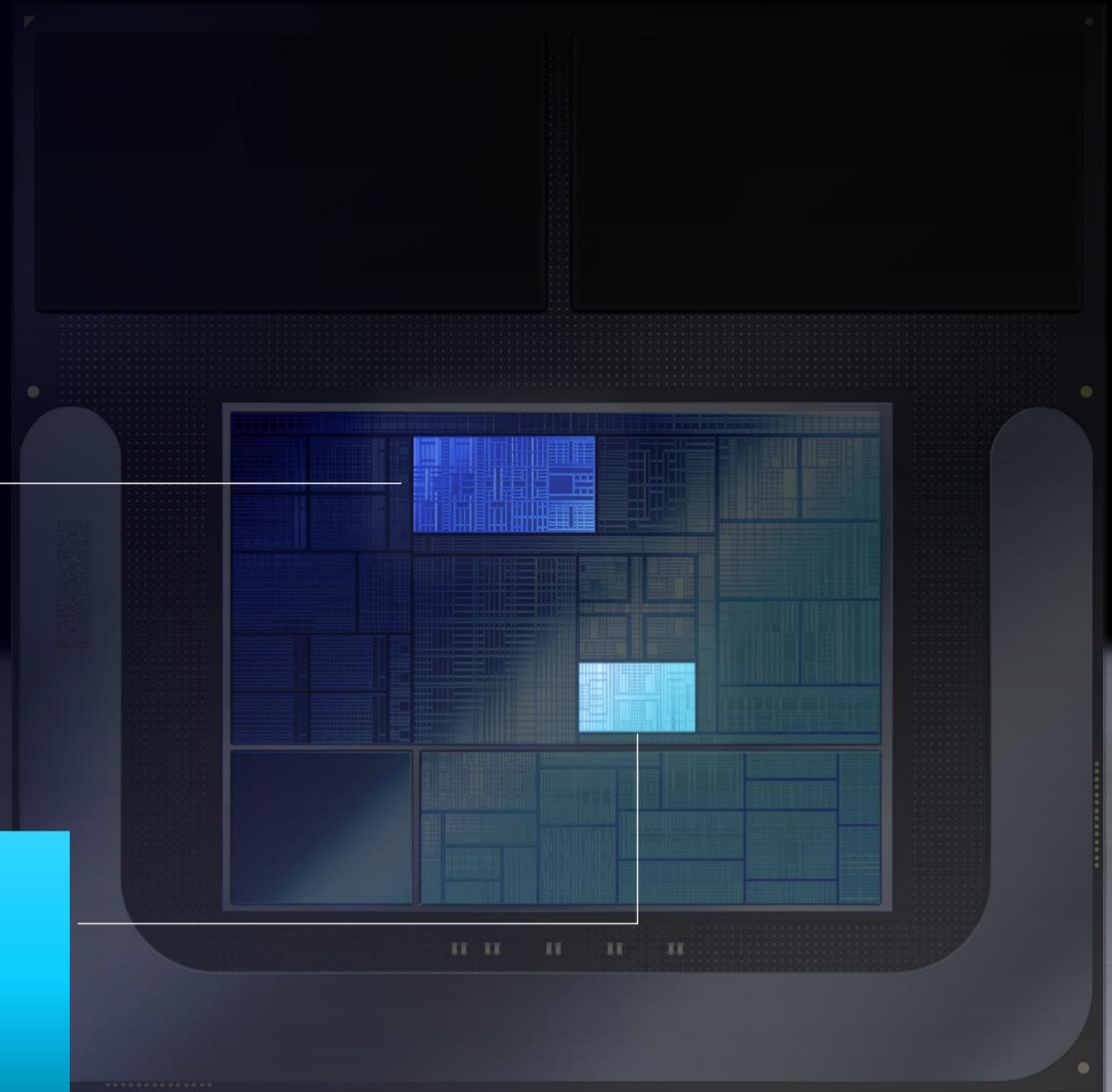
3 display pipes

# VVC Power reduction



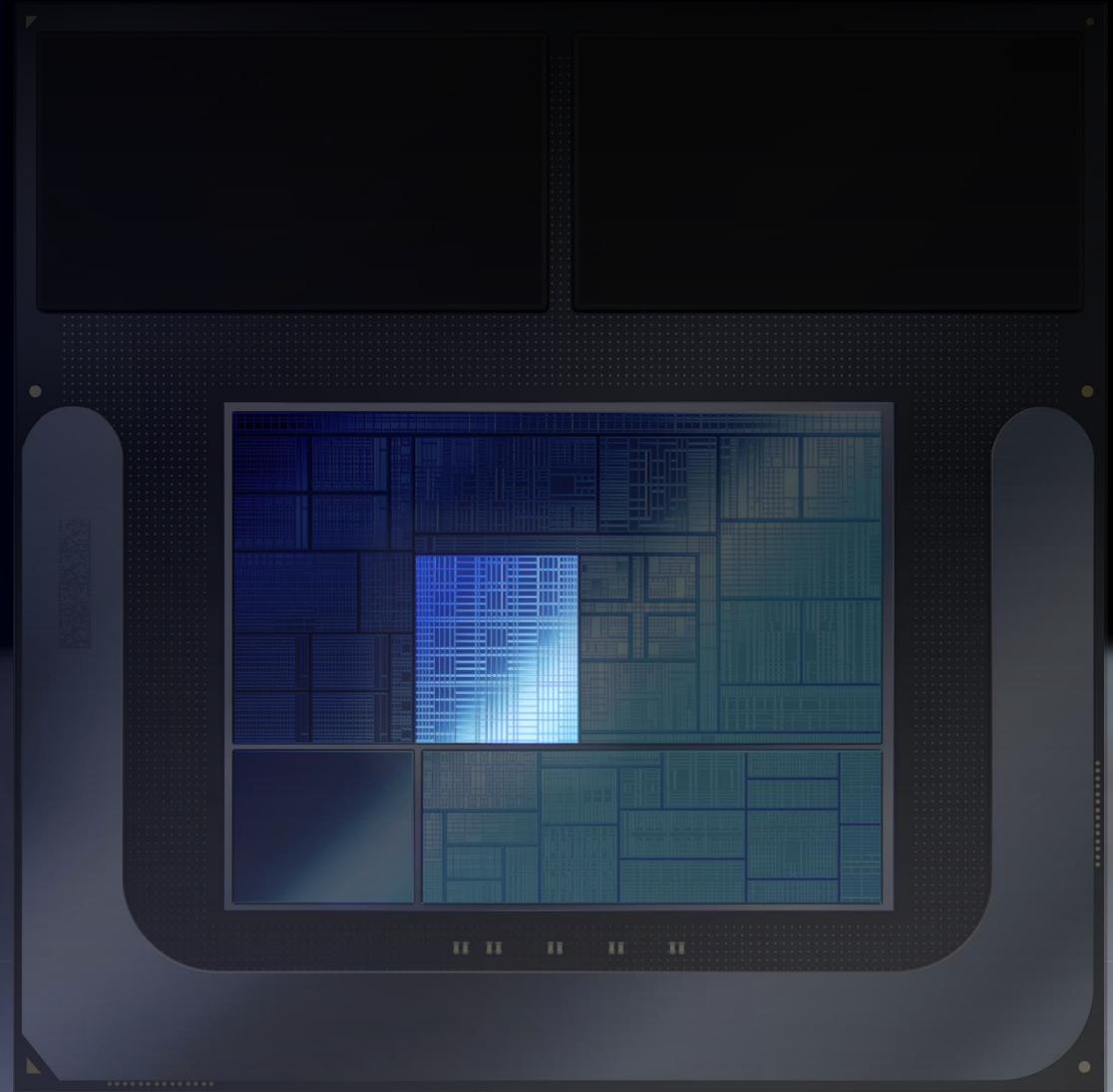CPU Package Power: 2.671 W
Lunar Lake (HW VVC)

CPU Package Power: 34.469 W
Meteor Lake (SW VVC)

# Lunar
# Lake

## New
## NPU 4.0

# Next Gen
# NPU 4
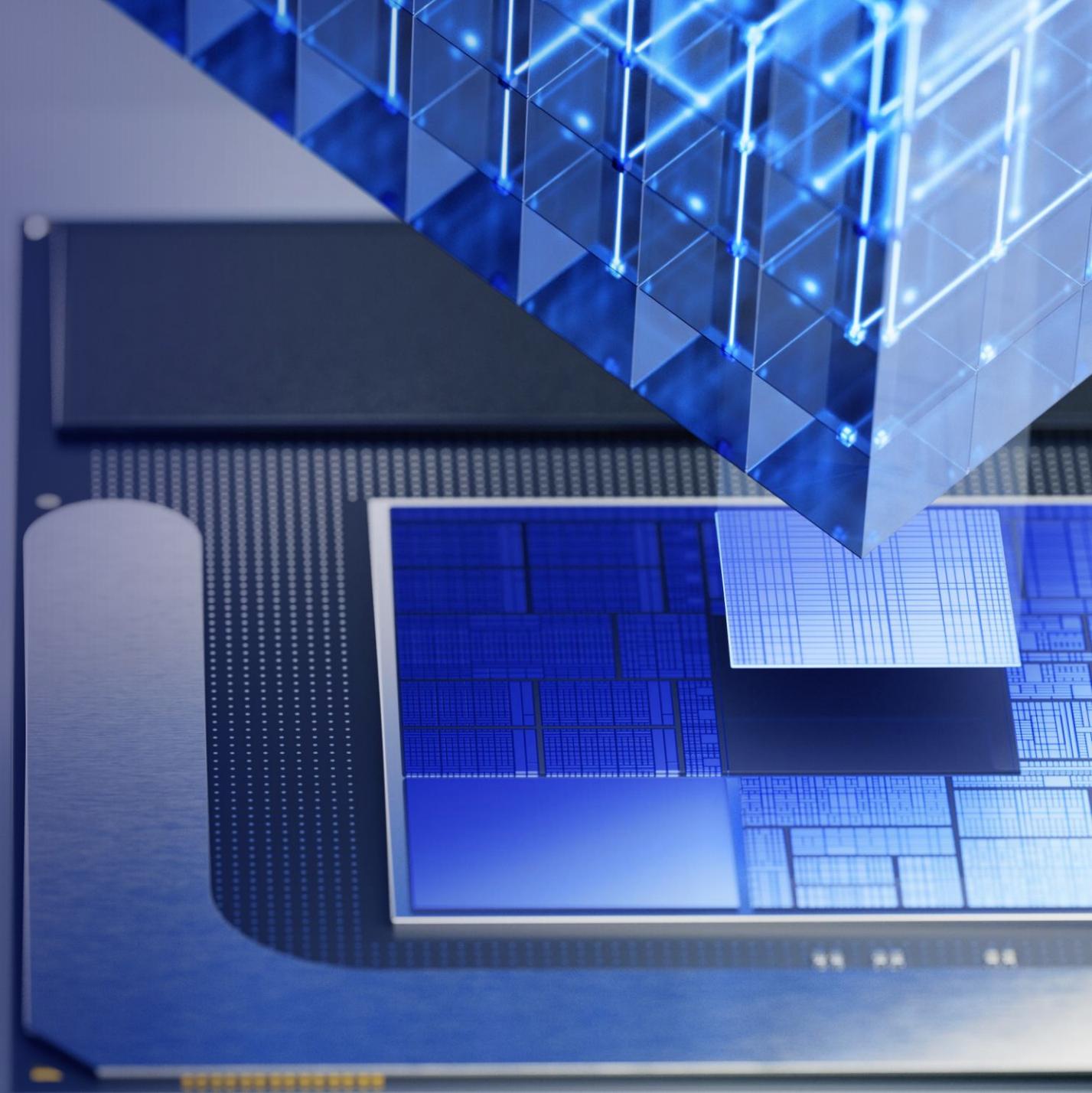## Architecture goals

**Increase NPU size**
To run next gen AI workloads

**Increase clock & efficiency**
To increase performance and battery life

**Optimize for modern AI**
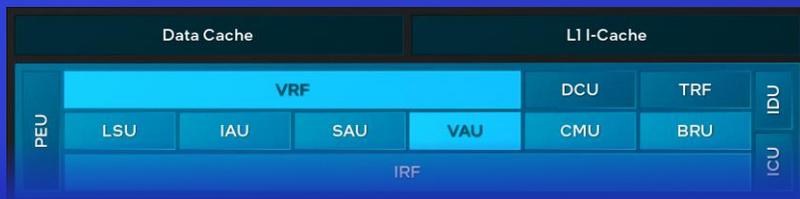For efficiently running LLMs and transformers

# Next Gen NPU 4

Largest integrated and dedicated AI accelerator for the AI PC

**12 Enhanced SHAVE DSPs**

Accelerating LLM & transformer operations

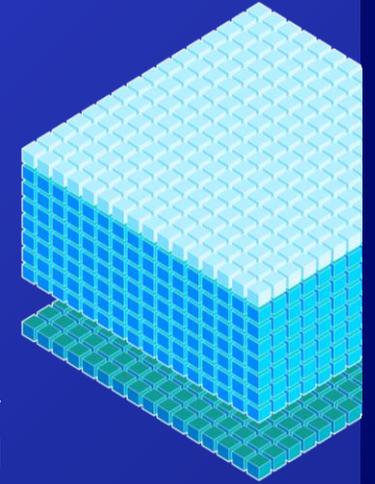Native activation function & data conversion support

up to **48 TOPS**

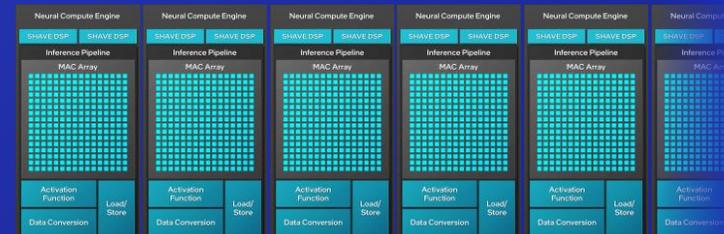**2x** Bandwidth

Efficiency optimized MAC array

Embedding tokenization used for LLMs

**6** Neural compute engines
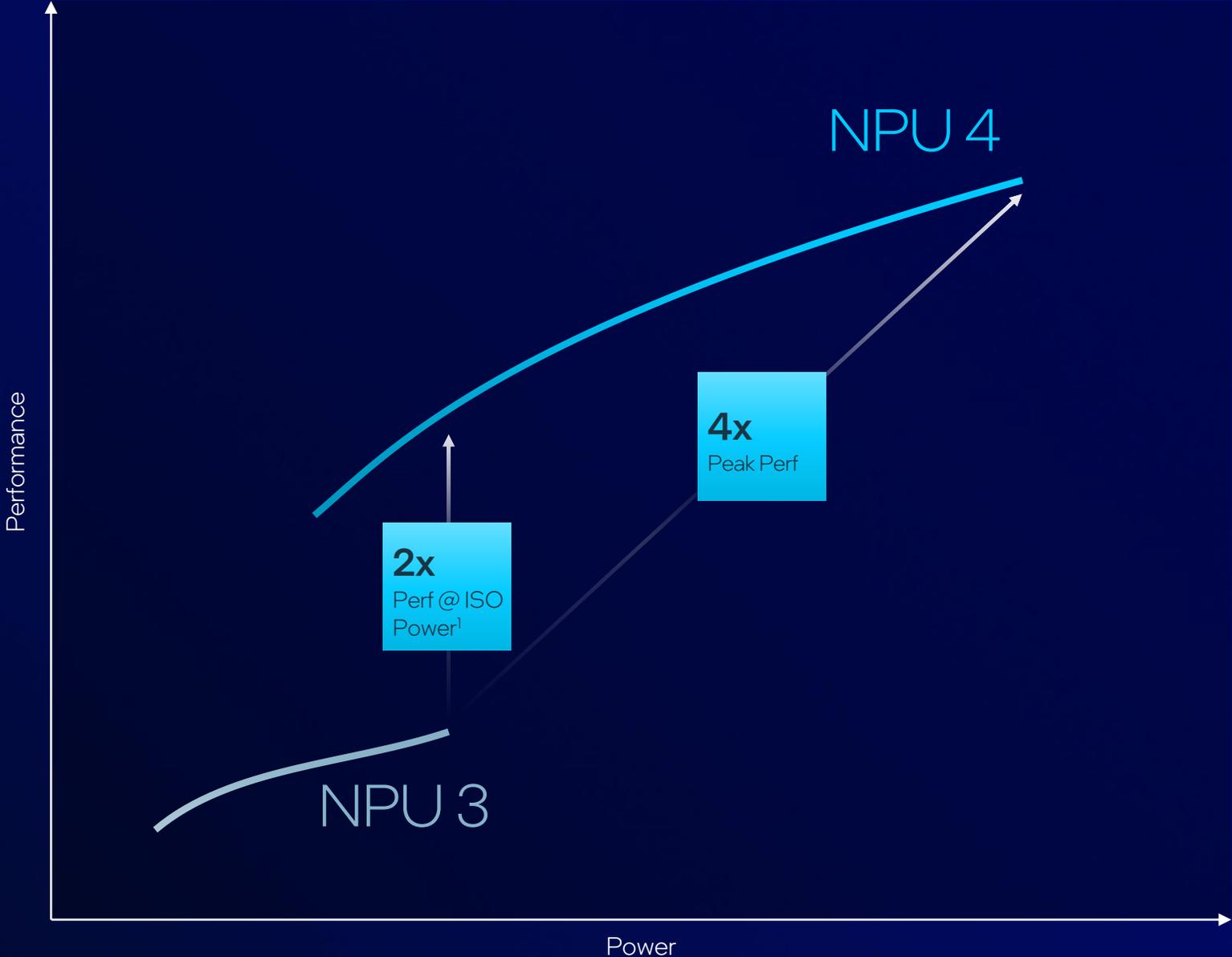
# Next Gen
# NPU 4

Scaling AI performance
and efficiency at AI pace

NPU 4 – in Lunar Lake

NPU 3 – in Meteor Lake

Performance

NPU 4

**4x**
Peak Perf

**2x**
Perf @ ISO
Power[1]

NPU 3

Power

[1] Based on pre-production simulation data of a real network. See backup for details.

# Lunar Lake

## Connectivity

Wi-Fi & Bluetooth

PCIe Gen 5.0

USB 3.0 & 2.0

Thunderbolt

PCIe Gen 4.0

# Leadership Connectivity

## Integrated right onto the package

**Intel® Bluetooth® 5.4**

For efficient & HD audio

**Integrated Intel® Wi-Fi 7 (5 Gig)**

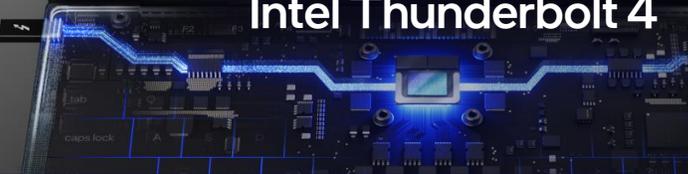**Intel Unison**

New multi-device experiences

Tablet control

Swift connect

Universal hotspot*

**5.8**Gb/s
Wi-Fi 7 speed

**40**Gb/s
TBT4 speed

**Up to 3x integrated Intel Thunderbolt 4**

**Thunderbolt Share**

Share between PCs at Thunderbolt speed

**PCI EXPRESS®**
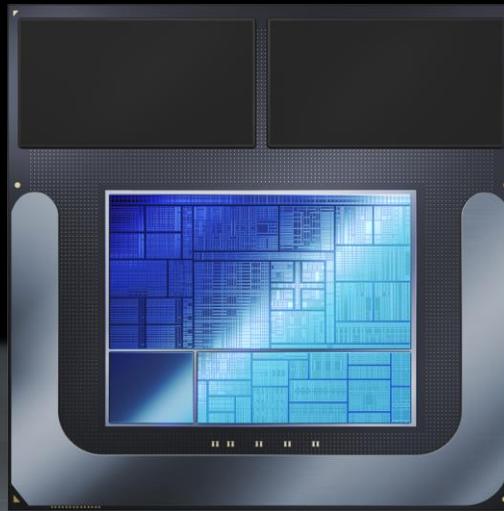
Up to
4x PCIe Gen 5.0
4x PCIe Gen 4.0

**Enhanced VR**

With Wi-Fi 7 & Intel Killer Wi-Fi

intel.

* Available post LNL launch

# Lunar Lake

Goals

| Multi-threaded | Single-threaded | AI compute | | Graphics performance | Core perf/watt | Energy efficiency |

- >20% MT perf with Lower Core Count

- > 20% Single Thread performance , but more important , delivers the same MTL Single thread performance at half the power.

- More than twice the performance per watt in productivity

- AI throughput of over 3 times the throughput across the NPU, GPU, and CPU

- For graphics apps, Lunar Lake delivers up to 50% uplift in performance,  much better user experiences for Gamers and creators

- For battery life, Lunar lake reduces SoC power by up to 40% , major step for mobile in which users will notice

Thank
You

# Notices & Disclaimers

# APPENDIX

| Claim # & Statement | Slide # & Title/Details |
|---|---|
| | SLIDE 2: Flagship SoC for the next gen of AI PCs |
| Up to 40% lower SoC power | Testing by Intel as of May 2024. Data based on Lunar Lake reference validation platform measurement vs Meteor Lake reference validation platform as measured by YouTube 4K 30 AV1. |
| Similar ST perf at half the power | Testing by Intel as of May 2024. Data based on Lunar Lake reference validation platform as measured CBR24 ST vs. prior generation. |
| Up to 1.5X better graphics | Testing by Intel as of May 2024. Data based on Lunar Lake reference validation platform measurement vs Meteor Lake reference validation platform as measured by 3DM Time Spy. 3DMark* |
| | SLIDE 9: Lion Cove P-core |
| Lion Cove core delivers 14% better IPC vs. Redwood Cove core | Iso frequency benefit estimate across: components of SPECrate2017_int_base and SPECrate2017_fp_base (estimated) running 1 copy |
| Lion Cove Performance at different power levels vs. Redwood Cove | Data based on measurements on an Intel internal reference validation platform. Performance: SPECrate2017_int_base (estimated) running 1 copy. Power: Fixed PL1 |

# APPENDIX

| Claim # & Statement | Slide # & Title/Details |
|---|---|
| | SLIDE 13: Skymont E-core |
| Skymont IPC on Lunar Lake Low Power Island: 1.38x integer and 1.68x floating point vs. Meteor Lake LP E-core (Crestmont) | Results are based on Intel's internal projections/estimates as of 5.13.2024(+/- 10% Margin of Error) on SPEC CPU 2017 Rate est, GCC12.1-O2 Linux at Fixed Frequency (ISO). |
| Skymont Power & Performance on Lunar Lake Low Power Island: up to 2x peak ST performance or 1/3 the power at similar ST performance and up to 4x peak MT performance or 1/3 of the power at similar MT performance | Results are based on Intel's internal projections/estimates as of 5.13.2024(+/- 10% Margin of Error) on SPEC CPU 2017_int_base est, GCC12.1-O2 Linux (ISO). Comparing a Skymont E-core cluster (4 Skymont cores) vs. Meteor Lake LP E-core cluster (2 Crestmont cores) to showcase workload coverage increase for the Lunar Lake Low Power Island. |
| | SLIDE 15: Breakthrough x86 Power Efficiency |
| >50% peak performance<br><br>>20-80% per/watt | Illustration of the relative Lunar Lake P-core and E-core performance across the SoC power range. |
| | SLIDE 21: Improved Experience |
| 35% power reduction when containment & power management optimization are enabled | As of May 2024, based performance estimated with measurements on Lunar Lake reference platform with power optimizations enabled vs. power optimizations disabled. |

intel. TECH tour.TW

# APPENDIX

| Claim # & Statement | Slide # & Title/Details |
|---|---|
| | SLIDE 26: Next Gen Xe2 GPU |
| 1.5x better vs. Meteor Lake GPU | Testing by Intel as of May2024. Data based on Lunar Lake reference validation platform measurement vs Meteor Lake reference validation platform as measured by 3DM Time Spy. 3DMark* |
| | SLIDE 28: Next Gen Xe2 GPU |
| 1.5x graphics performance over Meteor Lake | Testing by Intel as of May2024. Data based on Lunar Lake reference validation platform measurement vs Meteor Lake reference validation platform as measured by 3DM Time Spy. 3DMark* |
| | SLIDE 37: Next Gen NPU 4 |
| 2x performance at ISO power vs. Meteor Lake | Testing by Intel as of January 2024. Based on VPU-EM simulation. Power data is generated from the simulation tool based on power data that has been extracted from circuit simulation tools. This simulation, which is a ~100% utilization int8 network, is expected to correlate well with silicon. |
| 4x peak performance | 4x peak performance is based on TOPS increase from MTL (11 TOPS) to LNL (48 TOPS). |
| | SLIDE 41: Energy Efficiency |
| Up to 40% lower SoC power vs. Meteor Lake | Testing by Intel as of May 2024. Data based on Lunar Lake reference validation platform measurement vs Meteor Lake reference validation platform as measured by YouTube 4K 30 AV1. |

intel. TECH tour.TW