



January 24 - 26, 2023  
DoubleTree by Hilton San Jose  
ChipletSummit.com

# Open chipllets to enable a new era of silicon

Amber Huffman  
Google Cloud

With material & thanks to Partha Ranganathan, Martin Dixon,  
Peter Onufryk, and Rohit Mittal



# Challenge with Moore's Law



January 24 - 26, 2023  
DoubleTree by Hilton San Jose  
[ChipletSummit.com](http://ChipletSummit.com)



# Demand is increasing faster than ever ...

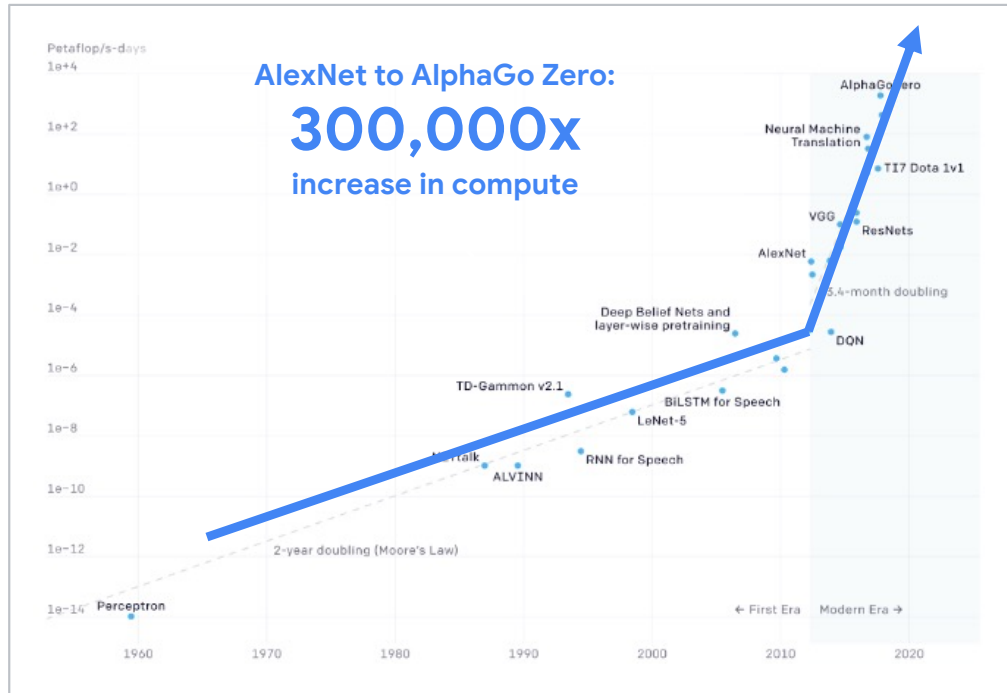


Image source: OpenAI

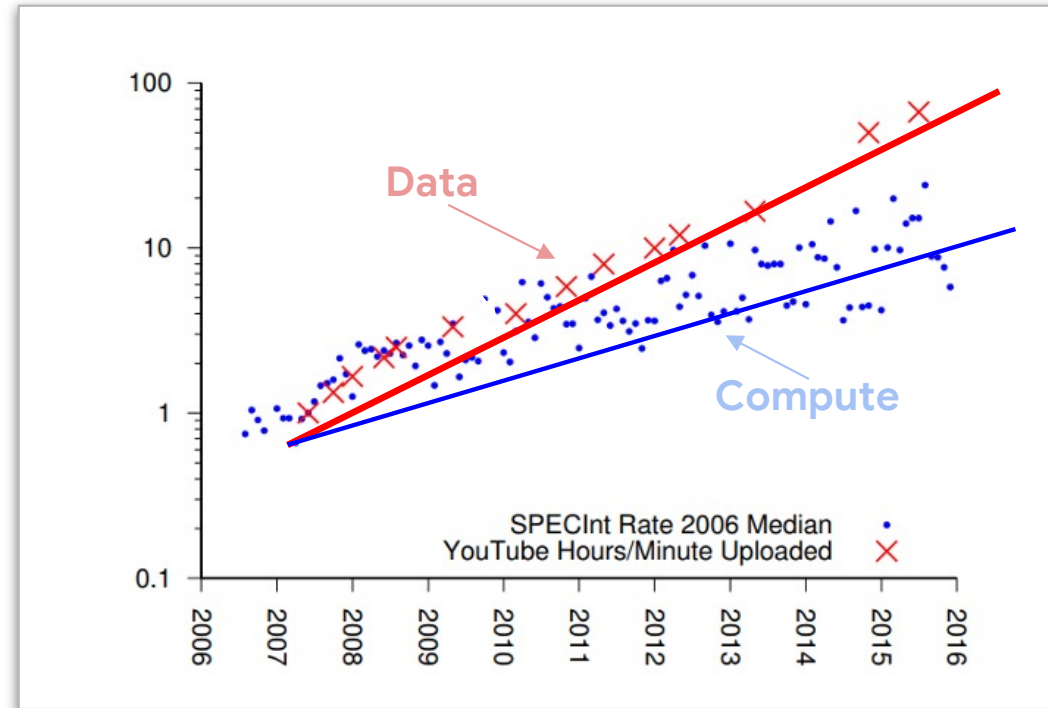
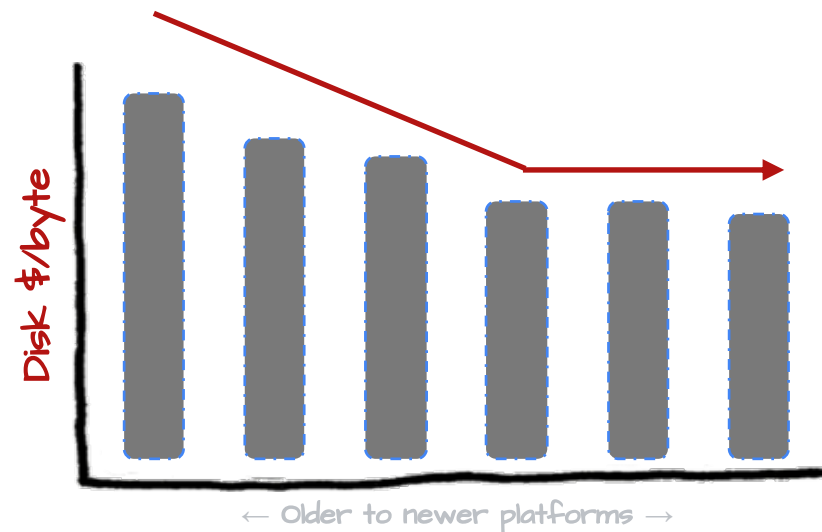
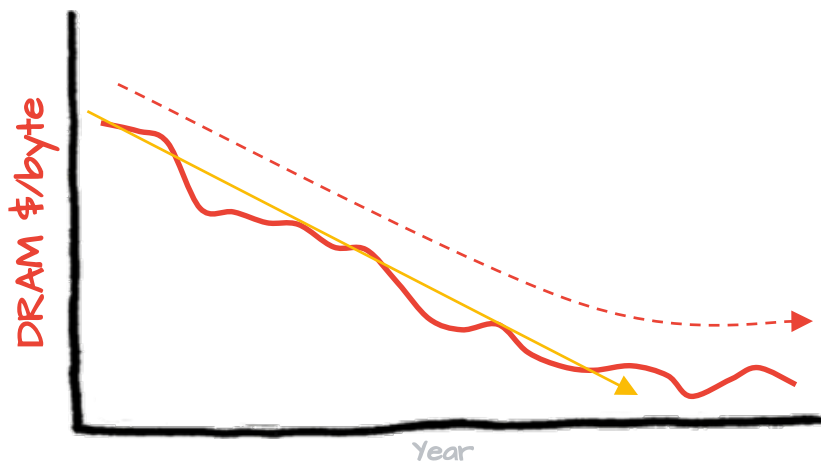
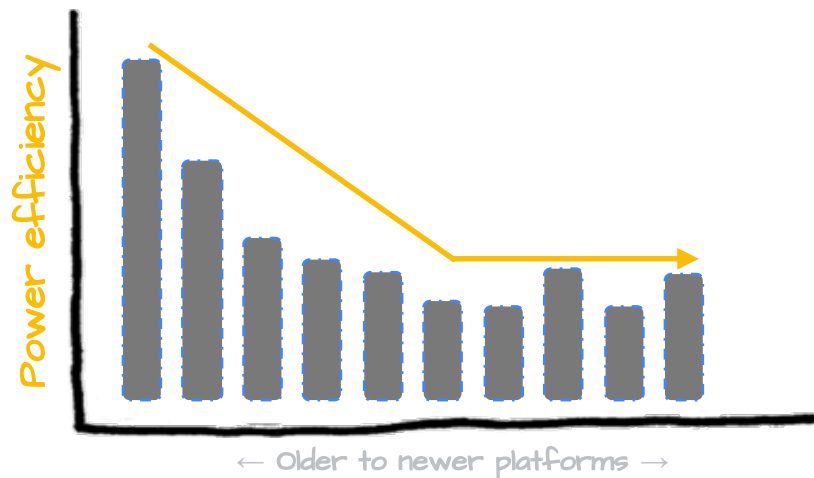
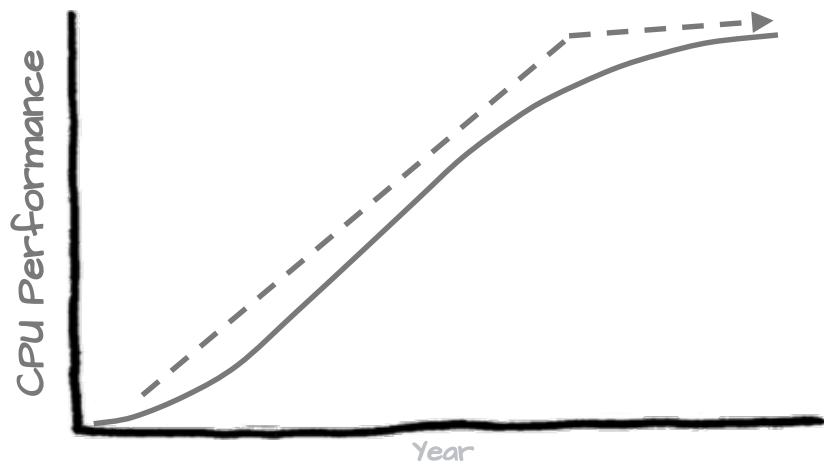


Image source: vbench, ASPLOS'18

# ... but **technology** plateaus ...



# Rise in Specialization



January 24 - 26, 2023  
DoubleTree by Hilton San Jose  
[ChipletSummit.com](http://ChipletSummit.com)



# Google VCU: Accelerating YouTube & More

Introduced to the world at ASPLOS 2021



Google

## Warehouse-scale Video Acceleration Co-design and Deployment in the Wild

ASPLOS 2021

[vcu@google.com](mailto:vcu@google.com)

Parthasarathy Ranganathan, Daniel Stodolsky, Jeff Calow, Jeremy Dorfman, Marisabel Guevara, Clinton Wills Smullen IV, Aki Kuusela, Raghu Balasubramanian, Sandeep Bhatia, Prakash Chauhan, Anna Cheung, In Suk Chong, Niranjani Dasharathi, Jia Feng, Brian Fosco, Samuel Foss, Ben Gelb, Sarah J. Gwin, Yoshiaki Hase, Da-ke He, C. Richard Ho, Roy W. Huffman Jr., Elisha Indupalli, Indira Jayaram, Poonacha Kongetira, Cho Mon Kyaw, Aaron Laursen, Yuan Li, Fong Lou, Kyle A. Lucke, JP Maaninen, Ramon Macias, Maire Mahony, David Alexander Munday, Srikanth Muroor, Narayana Penukonda, Eric Perkins-Argueta, Devin Persaud, Alex Ramirez, Ville-Mikko Rautio, Yolanda Ripley, Amir Salek, Sathish Sekar, Sergey N. Sokolov, Rob Springer, Don Stark, Mercedes Tan, Mark S. Wachsler, Andrew C. Walton, David A. Wickeraad, Alvin Wijaya, and Hon Kwan Wu.

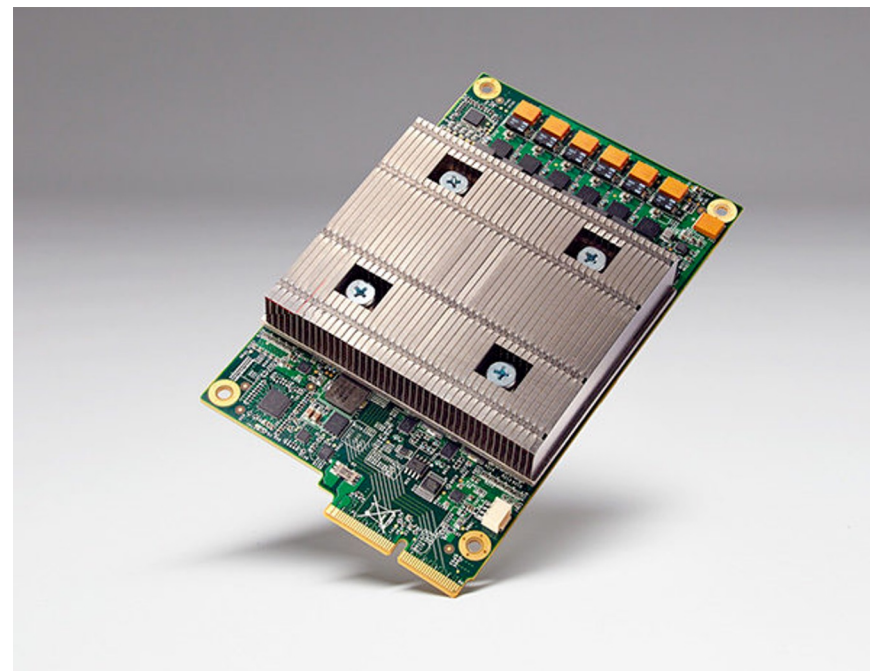
[goo.gle/vcu](https://goo.gle/vcu) for the paper and  
“[Warehouse Video Acceleration](#)” talk on YouTube



# Google TPU v1 (2015): Accelerating Inference

## First accelerator for ML Inference

- 92T Ops/sec (8-bit) @ 700MHz
- single-chip system
- built as coprocessor to a CPU



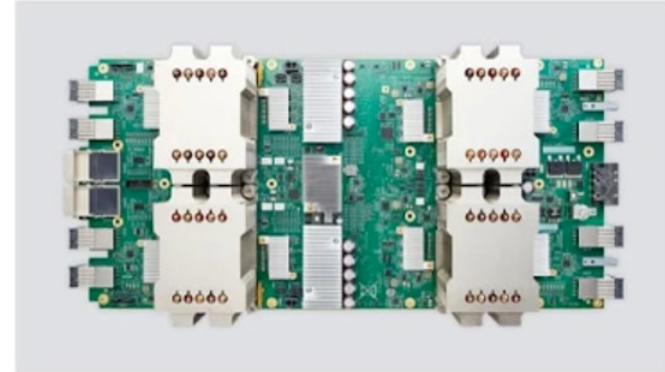
# Google TPU v2 (2017): Accelerating Training

## TPU v2 target ML training

- More flops and more memory
- High bandwidth **interconnect** enables multi-chip scaling

## TPU Pod -> **ML Supercomputer**

- Multi-chip Pod reduces ML training that takes 60-400 days to hours



Cloud TPU v2

180 teraflops

64 GB High Bandwidth Memory (HBM)



Cloud TPU v2 Pod

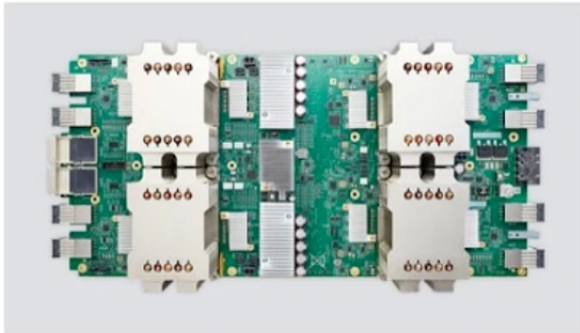
11.5 petaflops

4 TB HBM

2-D toroidal mesh network



# Google TPU v3 (2018): Accelerate Training More!



Cloud TPU v2

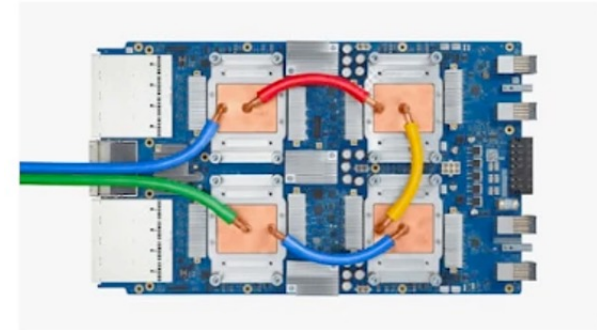
180 teraflops

64 GB High Bandwidth Memory (HBM)

**Increases v2 to v3**

2.3X teraflops

2X HBM



Cloud TPU v3

420 teraflops

128 GB HBM



Cloud TPU v2 Pod

11.5 petaflops

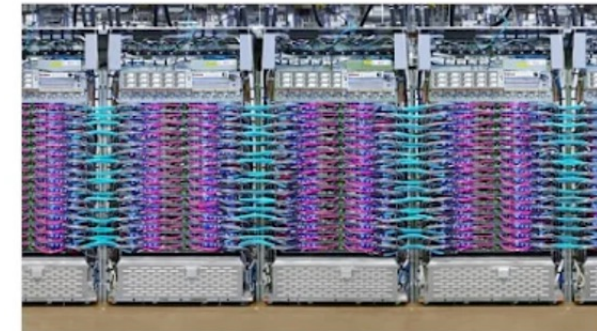
4 TB HBM

2-D toroidal mesh network

**Increases v2 to v3**

8X+ petaflops

8X HBM



Cloud TPU v3 Pod

100+ petaflops

32 TB HBM

2-D toroidal mesh network

# Custom Silicon is expensive



January 24 - 26, 2023  
DoubleTree by Hilton San Jose  
[ChipletSummit.com](http://ChipletSummit.com)



# Getting closer to the **Reticle Limit**

The mask/reticle is the 'glass' plate that has the exposure pattern for a modern semiconductor process

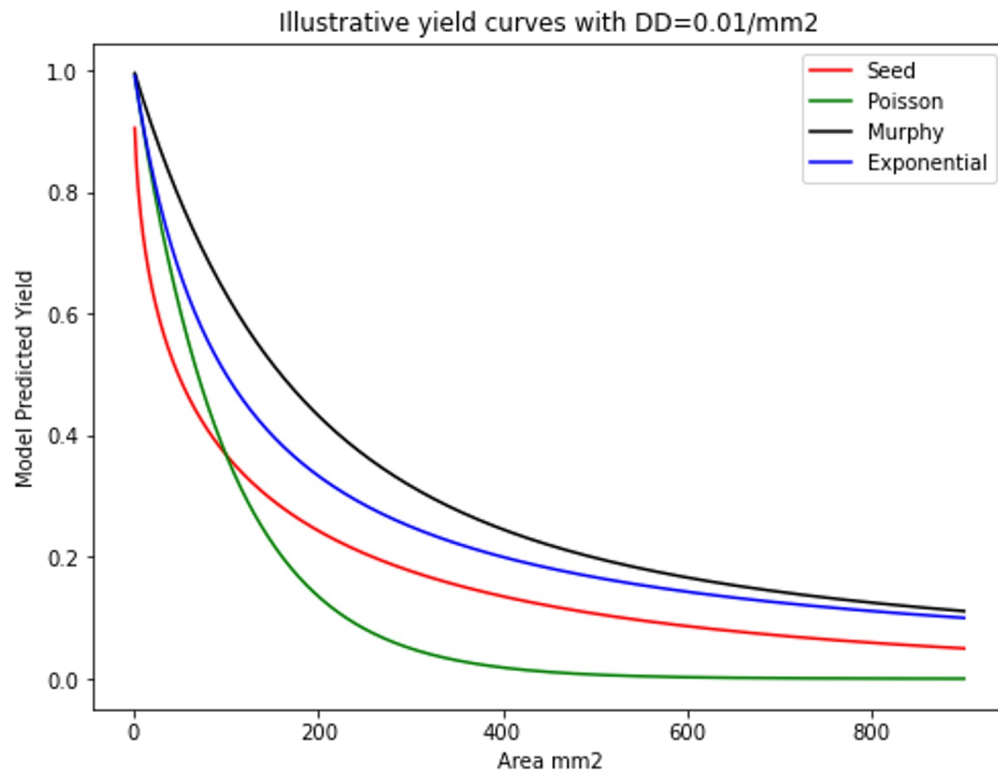
Current size limit for steppers is **850-900 mm<sup>2</sup>**

Yields at the largest sizes are **low**

<b>Chip feature</b>	<b>TPUv2</b>	<b>TPUv3</b>
Production deployment	Q3 2017	Q4 2018
Peak TOPS	46 (bf16)	123 (bf16)
Clock freq.	700 MHz	940 MHz
Tech. node, Die size	<b>16nm</b> <b>&lt; 625mm<sup>2</sup></b>	<b>16 nm</b> <b>&lt; 700 mm<sup>2</sup></b>
Transistor count	9 billion	10 billion

Source: [Jouppi, "Ten Lessons From Three Generations..."](#)

# Wafer Costs & Yield



not reflective of any specific process

## TSMC's Estimated Wafer Prices Revealed: 300mm Wafer at 5nm Is Nearly \$17,000

By Anton Shilov published September 18, 2020

High performance and high transistor density come at a cost

 Comments (3)



<https://www.tomshardware.com/news/tsmcs-wafer-prices-revealed-300mm-wafer-at-5nm-is-nearly-dollar17000>

# Chipllets arise



January 24 - 26, 2023  
DoubleTree by Hilton San Jose  
[ChipletSummit.com](http://ChipletSummit.com)





# What is a **Chiplet**?

Historically, die=package=chip.

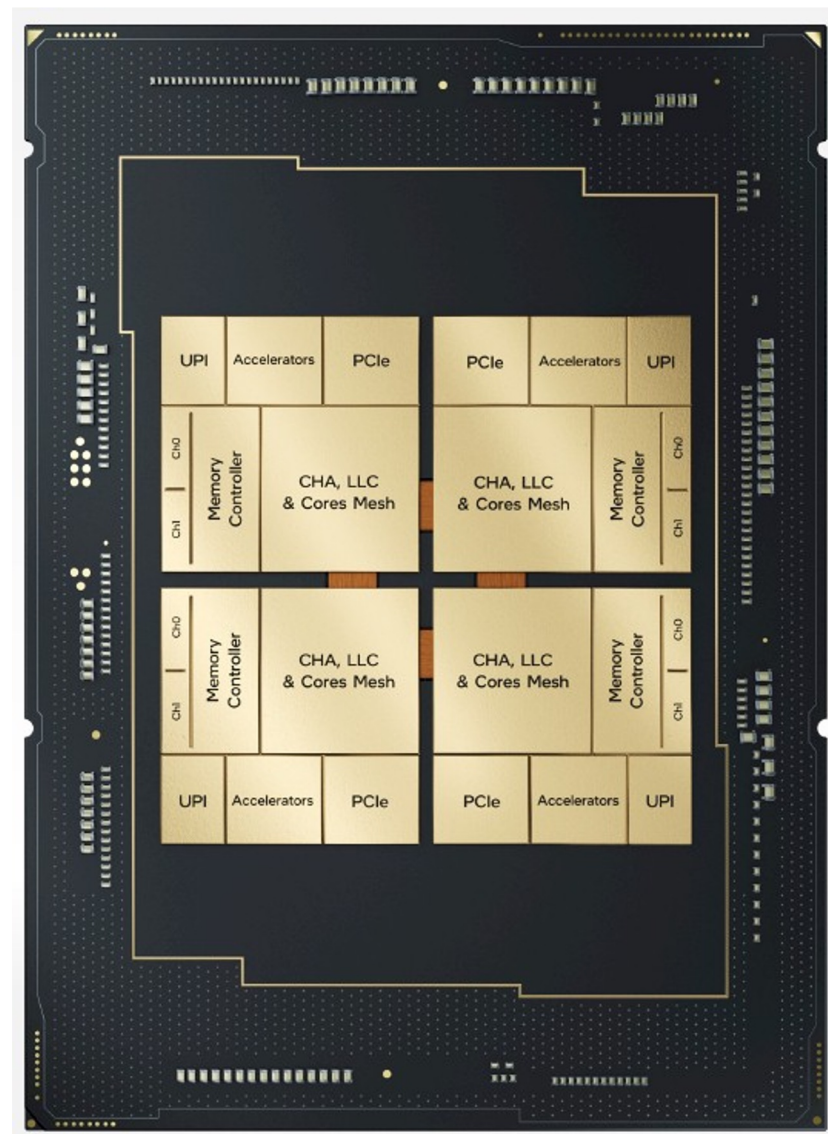
Recently, >1-2 die in a package/chip

die=chiplet

- GPUs have many dies
- AMD has 4-8 “compute” dies + 1 “I/O” die
- Intel showed images of four identical dies

Intel package shown on the right

- Four symmetric chiplets combined to build one x86-architecture chip





# The Rise & Promise of Chiplets

Devices that **exceed maximum reticle** size

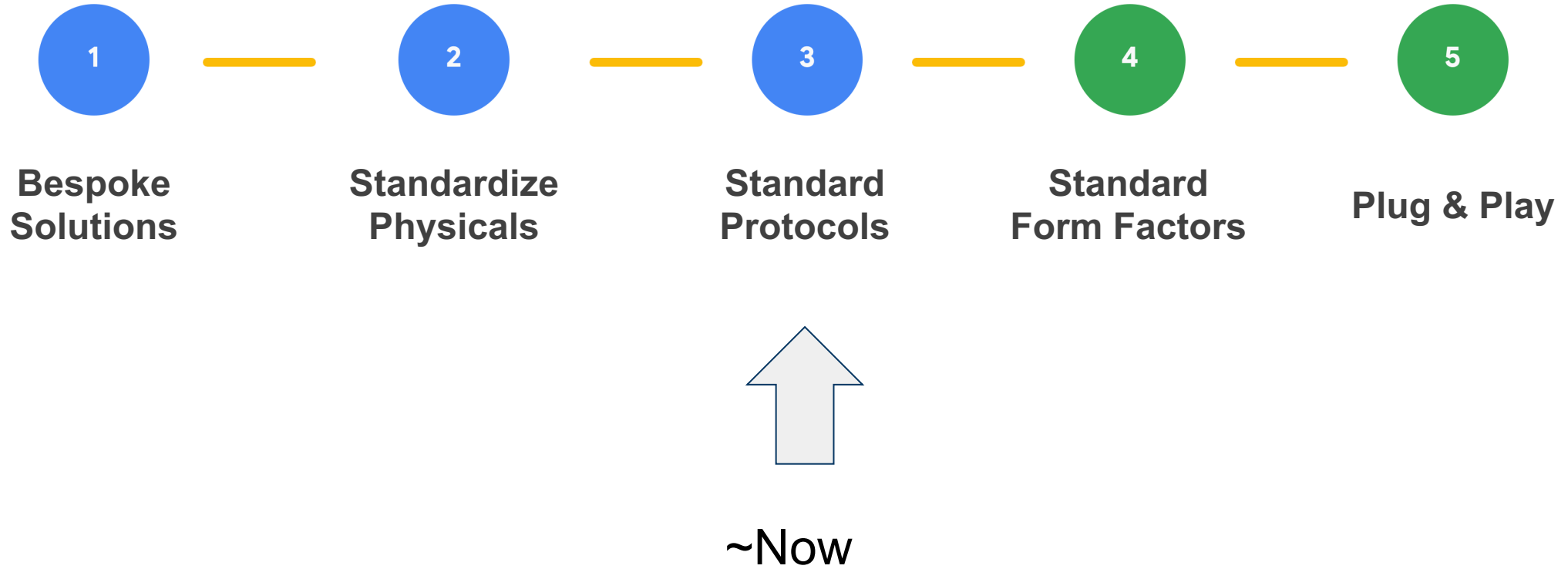
Reuse **reduces device design** time and cost

Smaller chiplets **improves silicon yield** and reduces cost

Different semiconductor processes reduces device cost and **enables greater levels of integration** and can lower risk

Different teams/companies may enable innovation & specialization

# Chiplet Journey



# Bespoke challenges

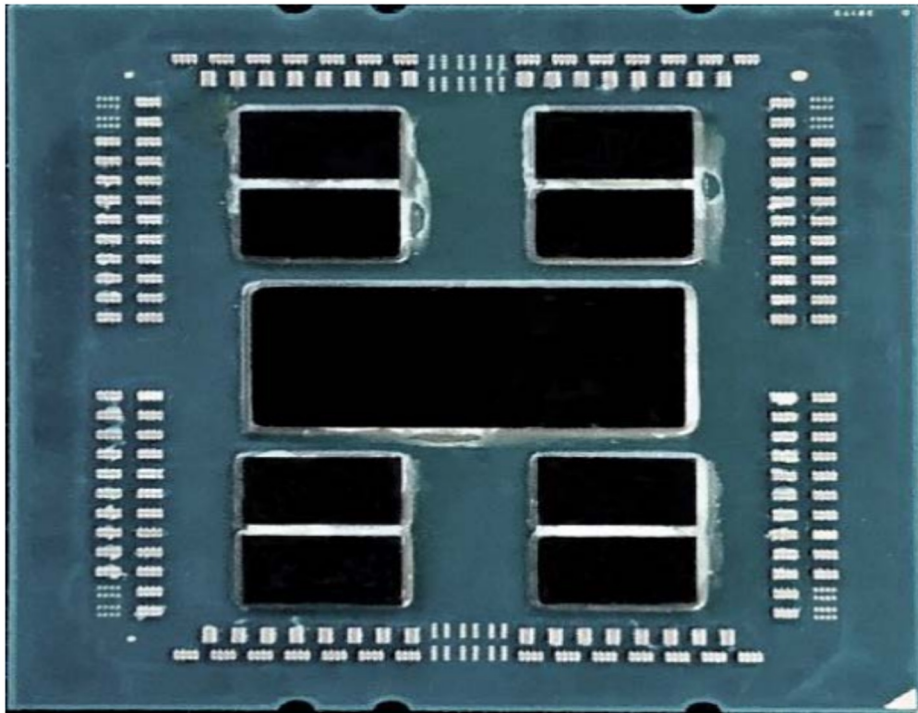


January 24 - 26, 2023  
DoubleTree by Hilton San Jose  
[ChipletSummit.com](http://ChipletSummit.com)



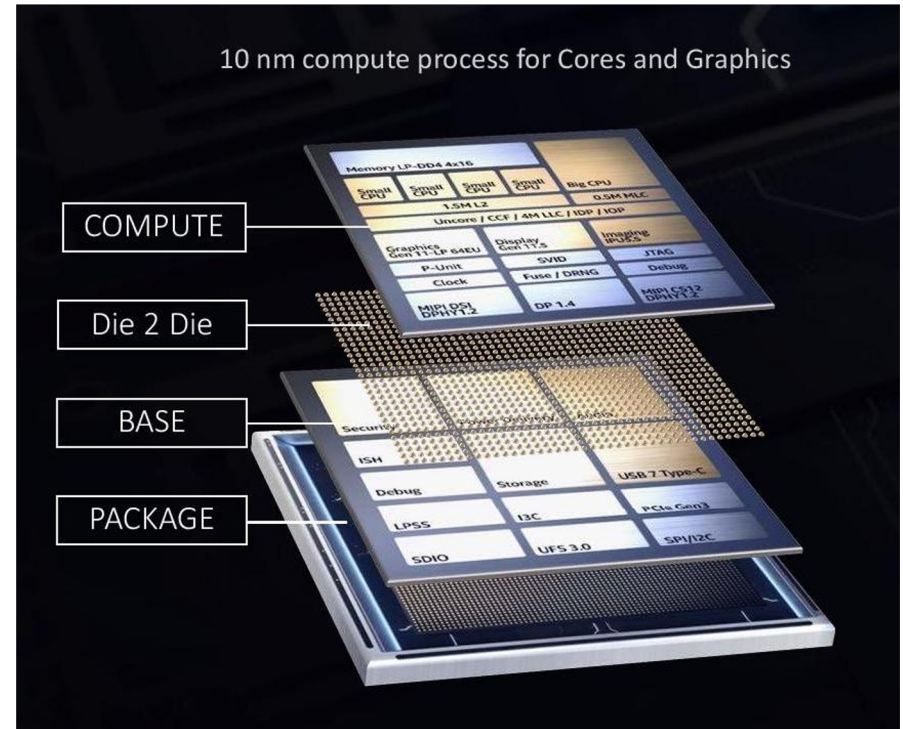
# Example: Bespoke Chiplet Solutions

## AMD Epyc 7002



<https://developer.amd.com/wp-content/resources/56827-1-0.pdf>

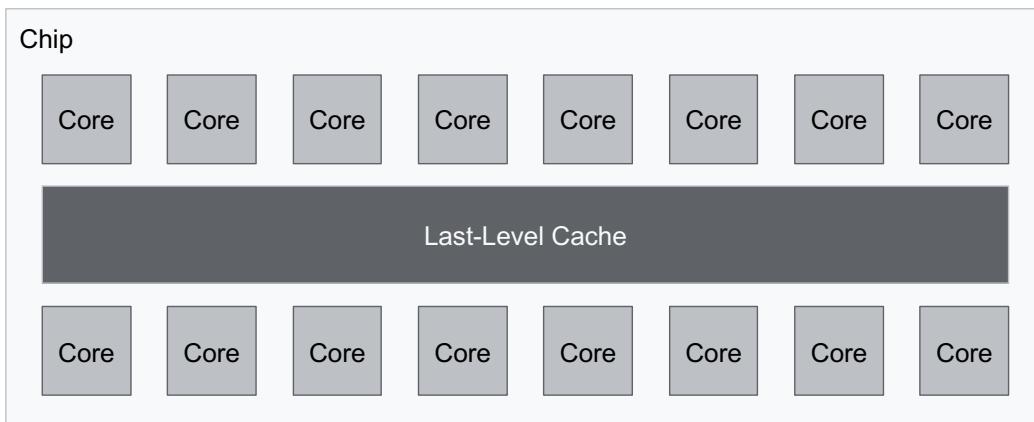
## Intel Lakefield



<https://newsroom.intel.com/press-kits/lakefield/>

# Chiplet Challenges: NUMA

Socket



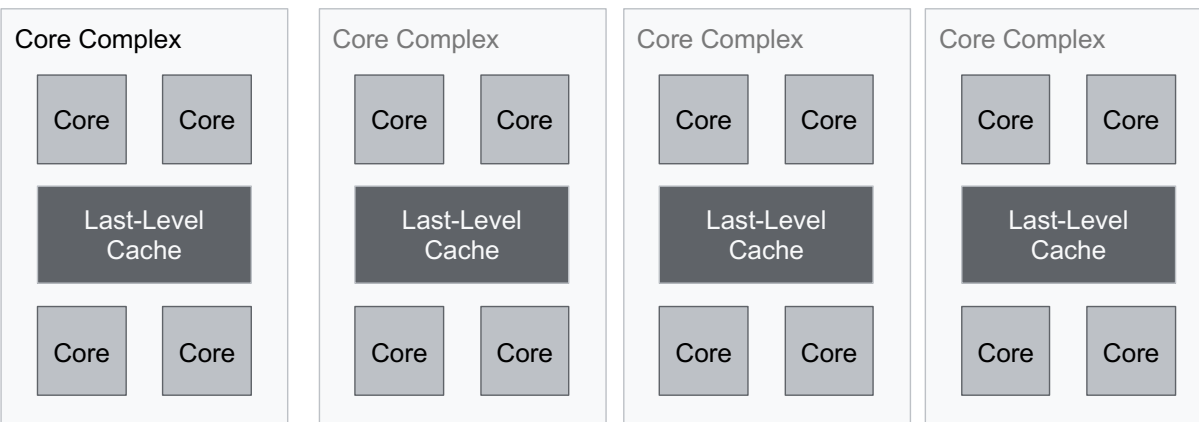
Monolithic

Shared LLC

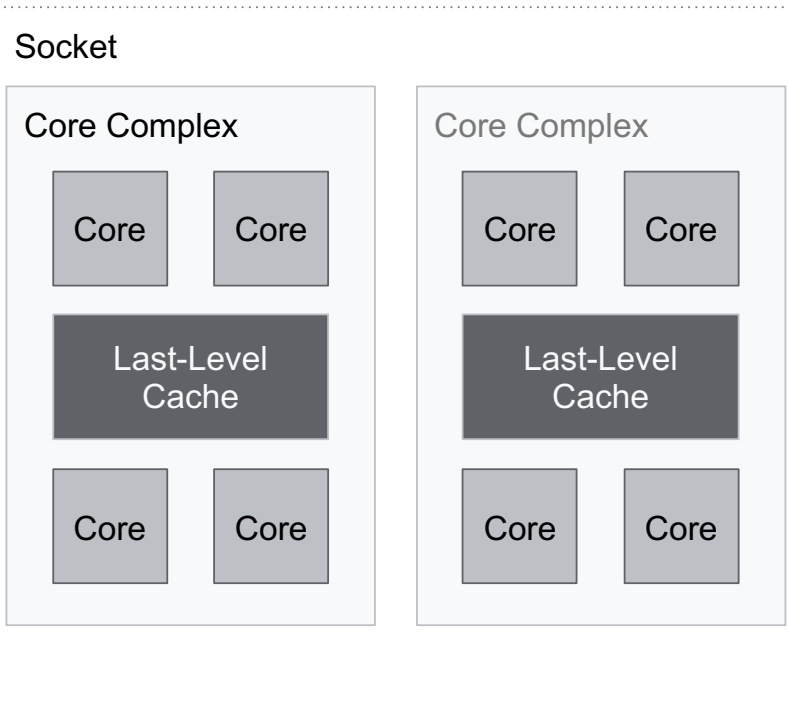
Chiplets

LLC **broken** up...

Socket

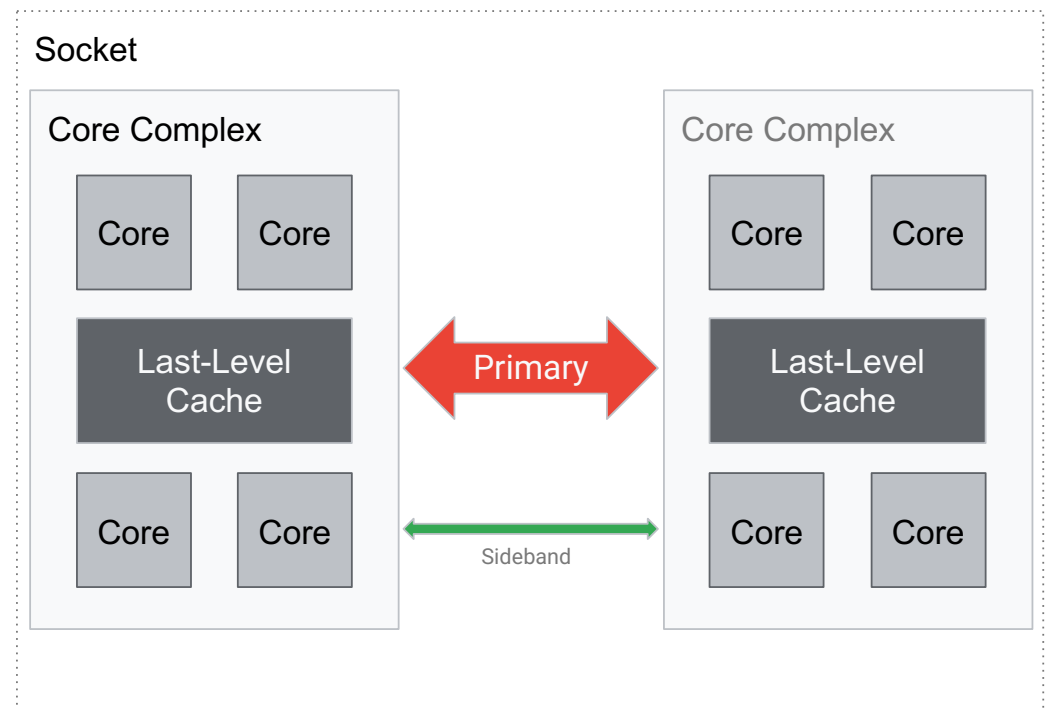


# Chiplet Challenges: Security & Manageability



## Monolithic

Communication across the chip is **easy** :-)



## Chiplets

Security and Manageability

... are **really** difficult!



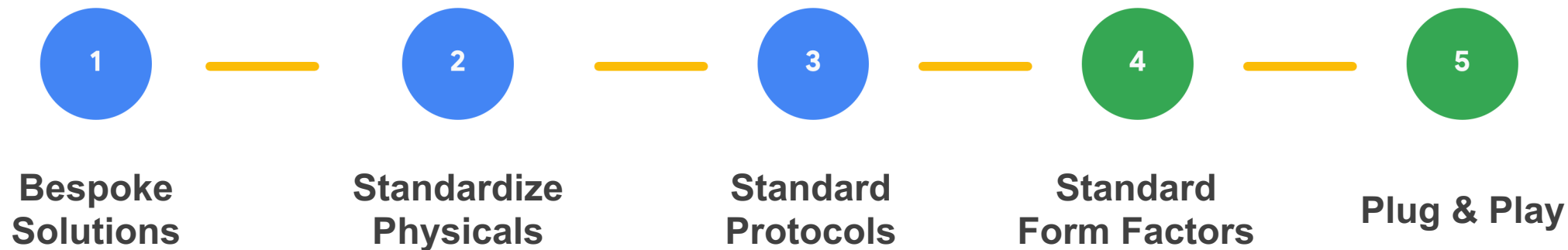
# Journey continues



January 24 - 26, 2023  
DoubleTree by Hilton San Jose  
[ChipletSummit.com](http://ChipletSummit.com)

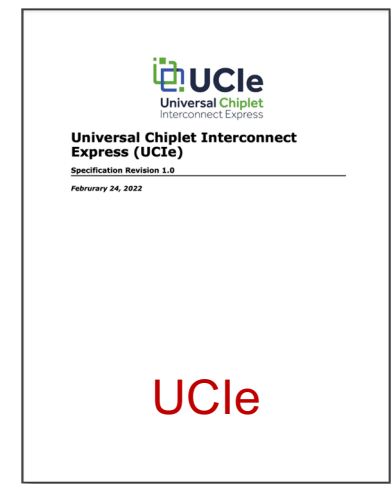
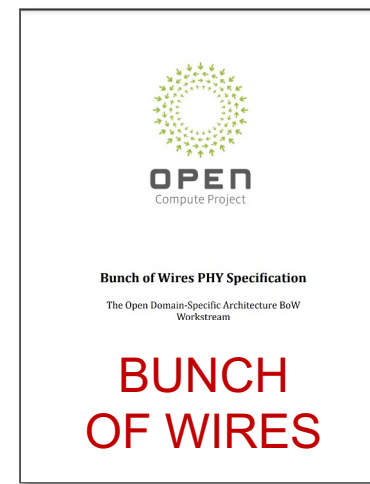
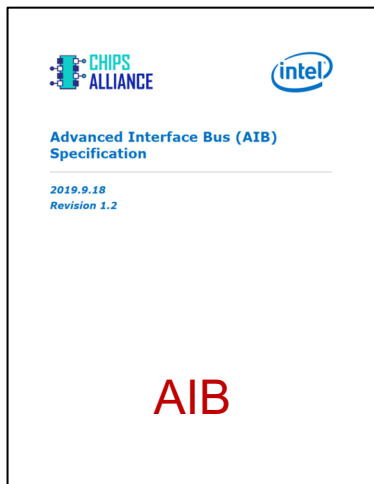


# Chiplet Journey



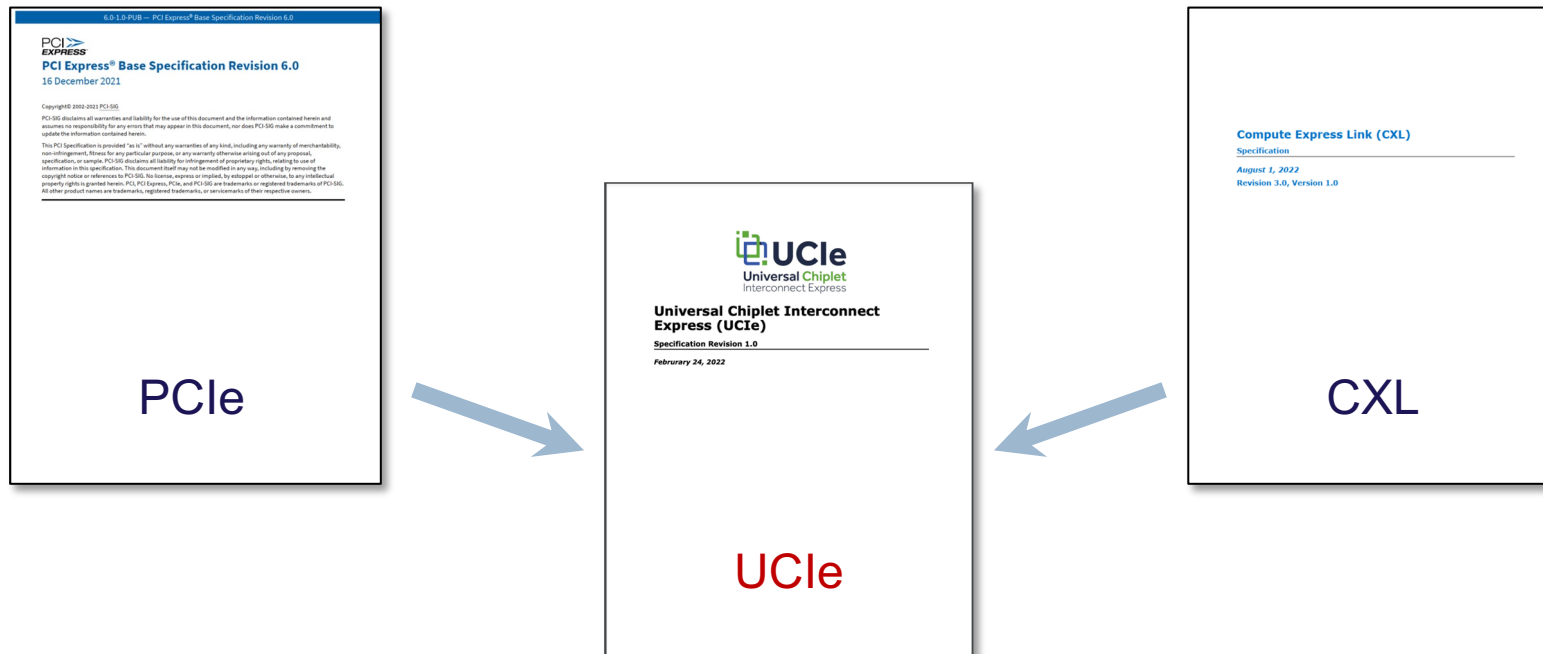
# Standard **Physical Interfaces**

Mechanical and electrical interfaces are defined enabling a die-to-die physical layer IP



# Standard Protocols

Standard protocols are defined enabling die-to-die fabric and controller IP and the beginning of an open chiplet ecosystem



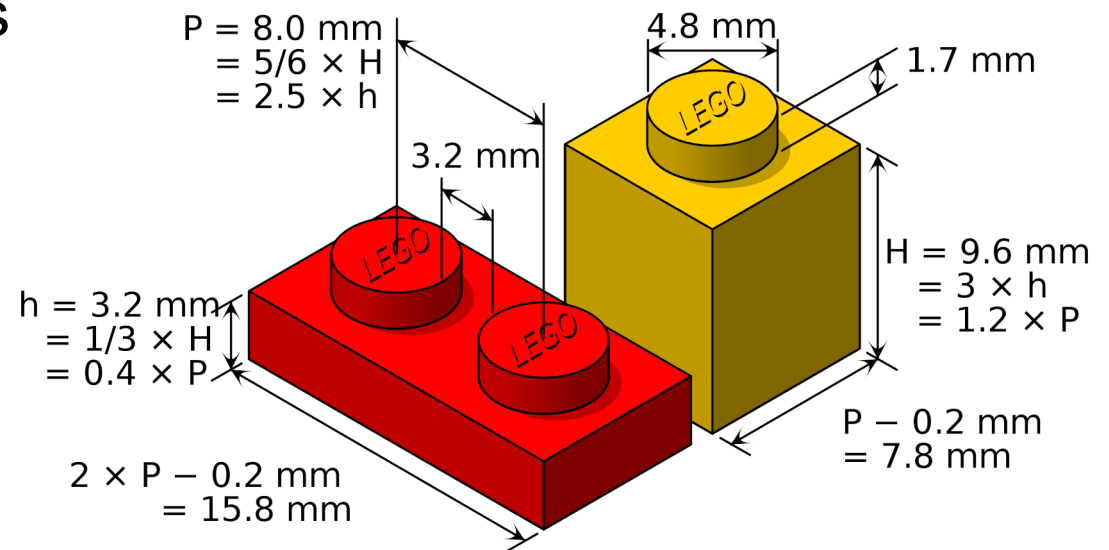
# Standard Form Factors, Plug & Play

Dreams of 'Legos'

Standardized form-factors are hard

What's the right size for a chiplet? Right shoreline?

Interesting journey ahead of us



# Learning More and **Getting Involved**

Specifications are emerging, but much more needs to be done

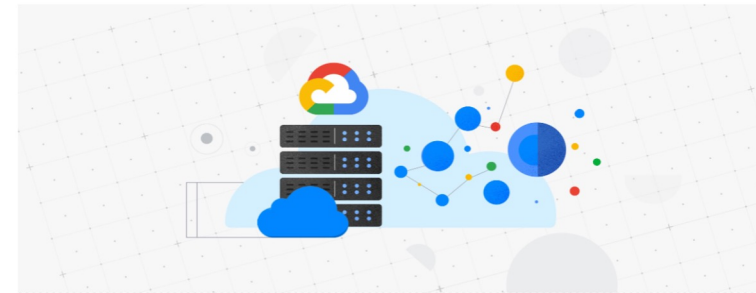
Biggest opportunities in standards near term are Security & Manageability, Form Factors, and Interoperability

## More Information

- Open HBI: Download [here](#)
- Bunch of Wires: Download [here](#)
- ODSA: Wiki [here](#)
- Universal Chiplet Interconnect Express: [www.uciexpress.org](http://www.uciexpress.org)
- Google blog available [here](#)

SYSTEMS

A chiplet innovation ecosystem for a new era of custom silicon



Parthasarathy Ranganathan  
VP, Technical Fellow  
March 2, 2022

As traditional Moore's law improvements slow down, we are now turning to custom chips to continue improving performance and efficiency. Innovations like Google's **Tensor Processing Units (TPUs)** and **Video Coding Units (VCUs)** have been incredibly valuable at **sustainably** meeting the growing demand for machine learning and video distribution services, and we expect to see additional custom chips that meet the emerging needs of our customers and users.

But building custom chips is a complex and costly endeavor. In particular, the semiconductor industry faces a key challenge. Each successive generation (technology

## Our Opportunity: **Make the SoC the New Motherboard!**



January 24 - 26, 2023  
DoubleTree by Hilton San Jose  
ChipletSummit.com







# Thank you.



January 24 - 26, 2023  
DoubleTree by Hilton San Jose  
[ChipletSummit.com](http://ChipletSummit.com)

