# HBM Package Integration:  Technology Trends, Challenges and Applications

**Suresh Ramalingam**

**Hot Chips**
**Aug 21, 2016**

# Agenda

- ▶ **Motivation**
- ▶ **HBM Packaging Options**
- ▶ **Interposer Design**
- ▶ **Supply Chain**
- ▶ **Application and Challenges**
- ▶ **Summary**

*S*tacked *S*ilicon *I*nterconnect *T*echnology Refers to Xilinx 3D solutions

© Copyright 2016 Xilinx

# FPGA's in the Data Center Today

- ❯ **The Accelerator (FPGA or GPU) is used to offload only certain tasks**
  - These tasks are called "Workloads", and FPGA's are well suited for many workloads
  - Note: Accelerators <u>don't</u> replace the CPU!

- ❯ **Hard and Soft is the basic approach**
  - **Hard** is the IO, Memory and PCIe interfaces
    - Does not change
  - **Soft** is the workload being accelerated
    - Is configured on the fly using P.R.

- ❯ **API's are run on the CPU to reprogram the FPGA to accelerate the workload as needed.**
  - Average P.R. happens every 15 minutes!

- ❯ **Acceleration requires lots of memory BW**

**Convey-Xilinx Accelerator**

**IBM Power8 Processor**

Hard

Soft

€ XILINX ❯ ALL PROGRAMMABLE™

# Multi-Die Technology for HBM (Side-by-Side)

MCM: Multi-chip Module
FO-MCM: Fan-out MCM
FL-MCM: Fine-line MCM
NTI: No TSV Interconnection
SLIM: Silicon-Less Integrated Module
EMIB: Embedded Multi-die Interconnect Bridge

I/O density (continuous interface)

$>10^4$

$10^3$

$10^2$

Organic Interposer/MCM
W/S 5/5um, ML ~12

**Cisco, ECTC 2016**

50 mm

38 mm

30 mm

50 mm

HBM-M
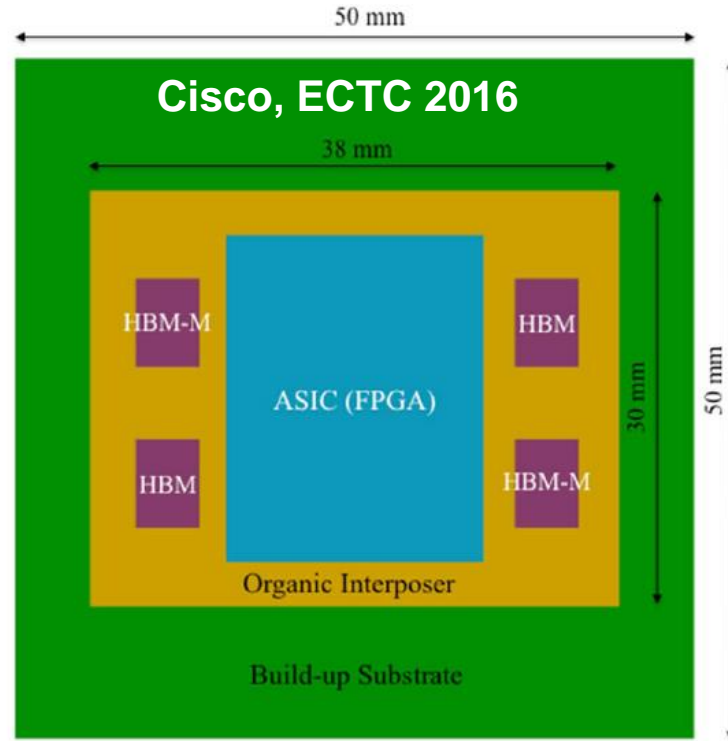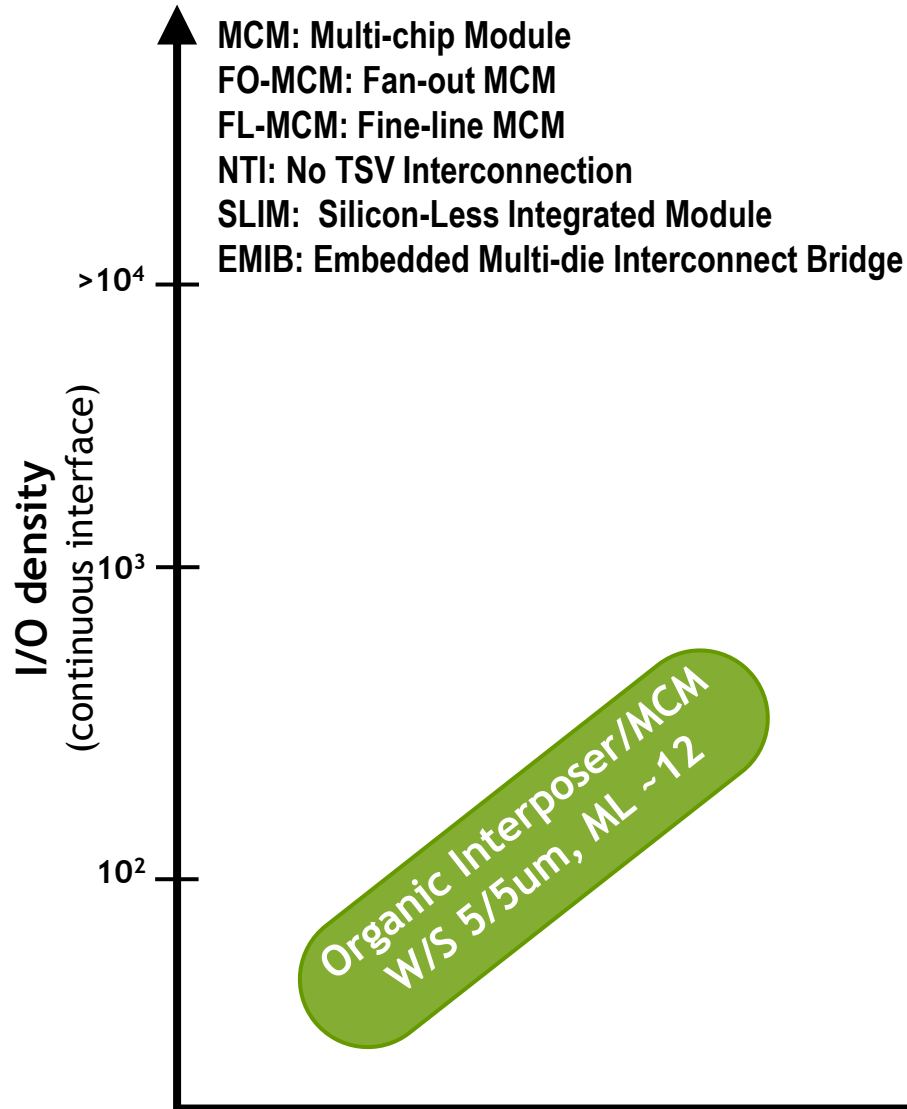
HBM

HBM

HBM-M

ASIC (FPGA)

Organic Interposer

Build-up Substrate

Figure 2. A schematic top view of the 3D SiP designed.

HBM_Functional

u Micro-pillar

HBM_Mechanical

Organic Interposer

C4 bumps
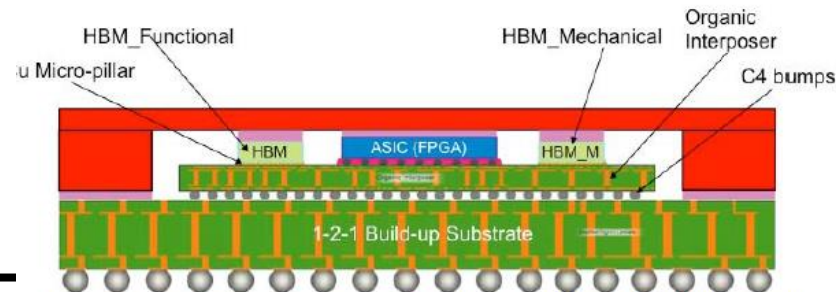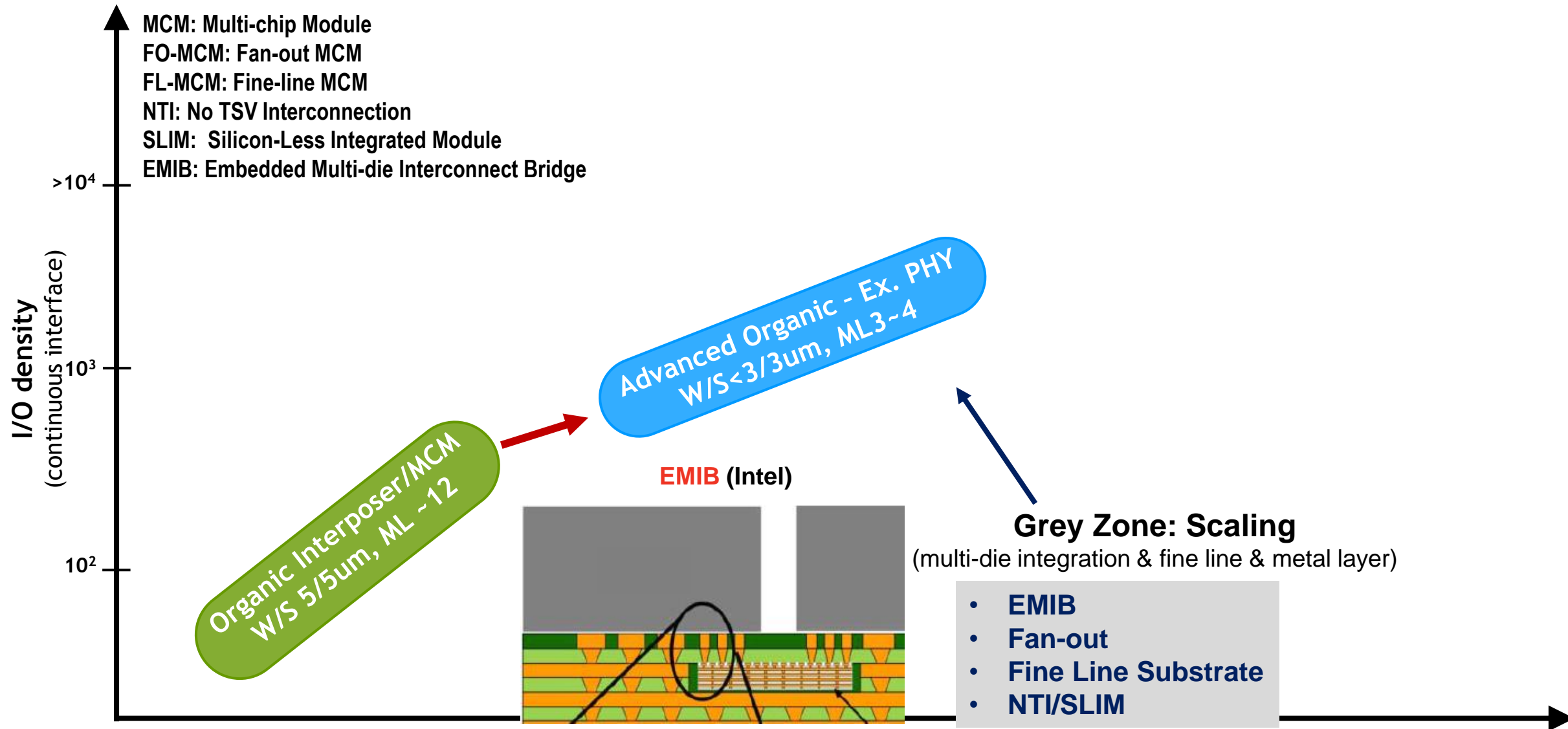
HBM

ASIC (FPGA)

HBM_M

1-2-1 Build-up Substrate

Figure 3. A schematic cross-sectional view of the 3D SiP designed.

**£ XILINX ➤ ALL PROGRAMMABLE.™**

# Multi-Die Technology for HBM (Side-by-Side)

**MCM: Multi-chip Module**
**FO-MCM: Fan-out MCM**
**FL-MCM: Fine-line MCM**
**NTI: No TSV Interconnection**
**SLIM:  Silicon-Less Integrated Module**
**EMIB: Embedded Multi-die Interconnect Bridge**

**I/O density**
(continuous interface)

$>10^4$

$10^3$

$10^2$

**Advanced Organic - Ex. PHY**
**W/S<3/3um, ML3~4**

**Organic Interposer/MCM**
**W/S 5/5um, ML ~12**

**EMIB** (Intel)

**Grey Zone: Scaling**
(multi-die integration & fine line & metal layer)

- **EMIB**
- **Fan-out**
- **Fine Line Substrate**
- **NTI/SLIM**

XILINX ➤ ALL PROGRAMMABLE.

# Multi-Die Technology for HBM (Side-by-Side)

**MCM: Multi-chip Module**
**FO-MCM: Fan-out MCM**
**FL-MCM: Fine-line MCM**
**NTI: No TSV Interconnection**
**SLIM:  Silicon-Less Integrated Module**
**EMIB: Embedded Multi-die Interconnect Bridge**

I/O density (continuous interface)

$>10^4$

$10^3$

$10^2$

Advanced Organic - Ex. PHY
W/S<3/3um, ML3~4

Organic Interposer/MCM
W/S 5/5um, ML ~12

**FL-MCM**
**(Shinko/Hitachi/UMTC/Ibiden)**



micro bump
(45mm pitch)

Die

Fine line layer

**Organic**/
LTCC layer

BGA bump

**XILINX** ➤ ALL PROGRAMMABLE.

# Multi-Die Technology for HBM (Side-by-Side)

MCM: Multi-chip Module
FO-MCM: Fan-out MCM
FL-MCM: Fine-line MCM
NTI: No TSV Interconnection
SLIM: Silicon-Less Integrated Module
EMIB: Embedded Multi-die Interconnect Bridge

I/O density (continuous interface)

$>10^4$

$10^3$

$10^2$

Advanced Organic - Ex. PHY W/S<3/3um, ML3~4

Organic Interposer/MCM W/S 5/5um, ML ~12

**NTI** (SPIL)



| Layer | Scheme (thk) | Dielectric |
|---|---|---|
| M1: L/S 2/2 um | FS- RDL1 | SiO2/SiNx |
| M2: L/S 5/5 um | BS- RDL2 | PBO |
| M3: L/S 10/10um Contact to BGA Ball | BS- RDL3 | PBO |

Figure 16. Hybrid Integration Scheme

**FO-MCM** (TSMC/ASE/SPIL/Amkor)

XILINX ➤ ALL PROGRAMMABLE.

# Multi-Die Technology for HBM (Side-by-Side)

MCM: Multi-chip Module
FO-MCM: Fan-out MCM
FL-MCM: Fine-line MCM
NTI: No TSV Interconnection
SLIM:  Silicon-Less Integrated Module
EMIB: Embedded Multi-die Interconnect Bridge

I/O density
(continuous interface)

$>10^4$

$10^3$

$10^2$

2.5D Si interposer
W/S<1/1um, ML ≤ 3

Advanced Organic - Ex. PHY
W/S<3/3um, ML3~4

Organic Interposer/MCM
W/S 5/5um, ML ~12

SSIT (Xilinx)

Fiji (AMD)



28nm FPGA Slice  28nm FPGA Slice  28nm FPGA Slice  28nm FPGA Slice
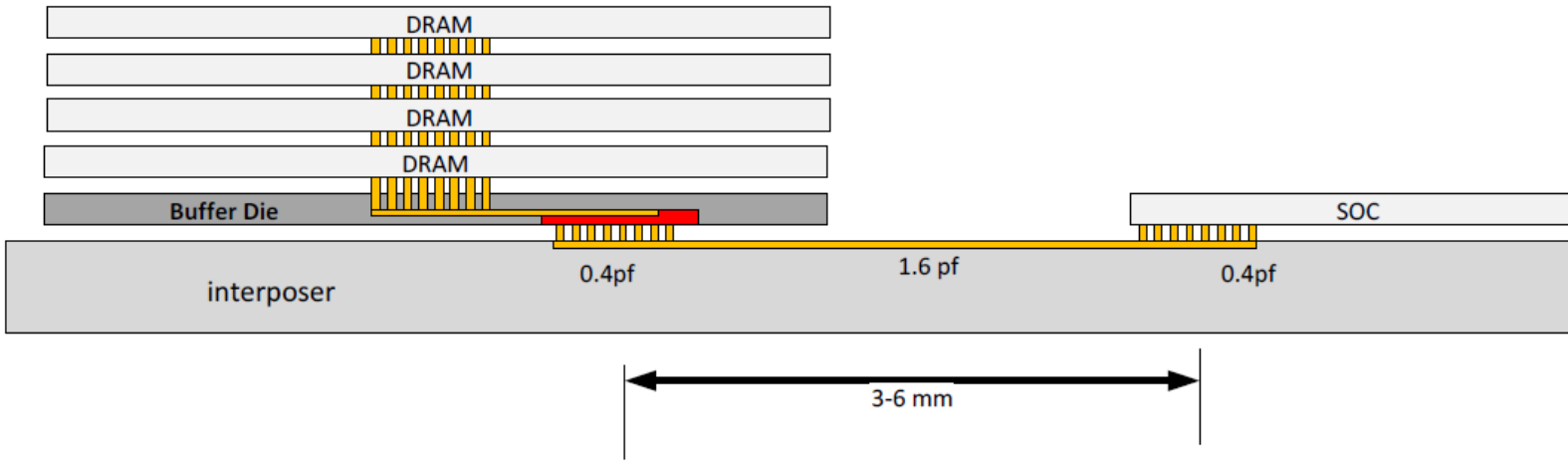
Package Substrate

XILINX ➤ ALL PROGRAMMABLE.

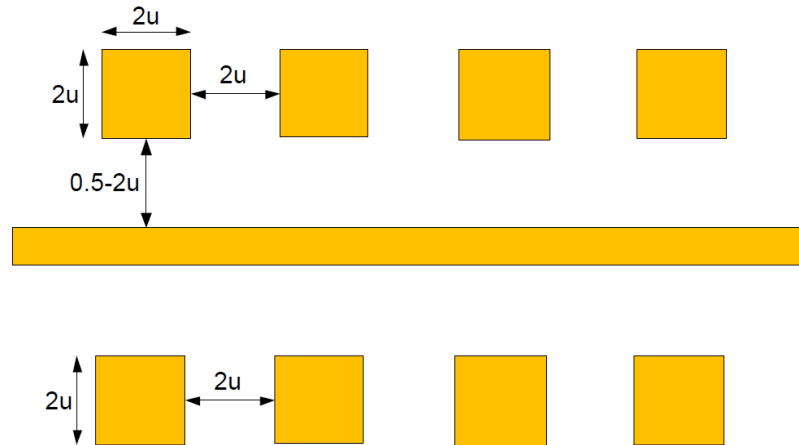# Multi-Die Package Design Rule Comparison

| Design Rules for Die to Die interconnection | MCM (Substrate) **Integrated Fine Layers** | EMIB (Embedded Multi-die Interconnect bridge) | Silicon Interposer (65 nm BEOL) | WLFO (Wafer Level Fan-out) |
|---|---|---|---|---|
| Minimum Bump pitch (um) | 130 (C4) 40 (u-bump d2d interface) | 130 (C4) 40 (u-bump) bridge | < 40 (u-bump) | 40 um RDL pad pitch |
| Via size / pad size (um) | 10 / 25 | 0.4 / 0.7 | 0.4 / 0.7 | 10/25 |
| Minimum Line & Space (um) | 2 / 2 | 0.4 / 0.4 | 0.4 / 0.4 | 2 / 2 |
| Metal thickness (um) | 2-5 | 1 | 1 | 2-5 |
| Dielectric thickness (um) | ~5 | 1 | 1 | < 5 |
| # of die-to-die connections per layer + GND shield layer (2L) | 1000's | 1000's (bridge interface length limited) | 10,000's | 1000's |
| Minimum die to die spacing (um) | < 500 | <2500 | <100 | < 250 |
| # of High density layers feasible | Not a limitation 1-3L | Not a limitation | Not a limitation | 1-3L layers |
| Die Sizes for assembly and # of assemblies | Not a concern d2d interconnect only | Size & # limitation? | Not a concern | Size limitation? |
| In Production | **No** (2018) | **No** (2017) | **Yes** | **No** not for 2/2um L/S (2018) |

**ΣXILINX ➤ ALL PROGRAMMABLE.**
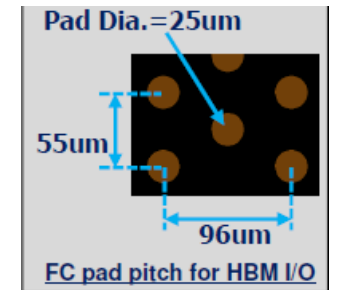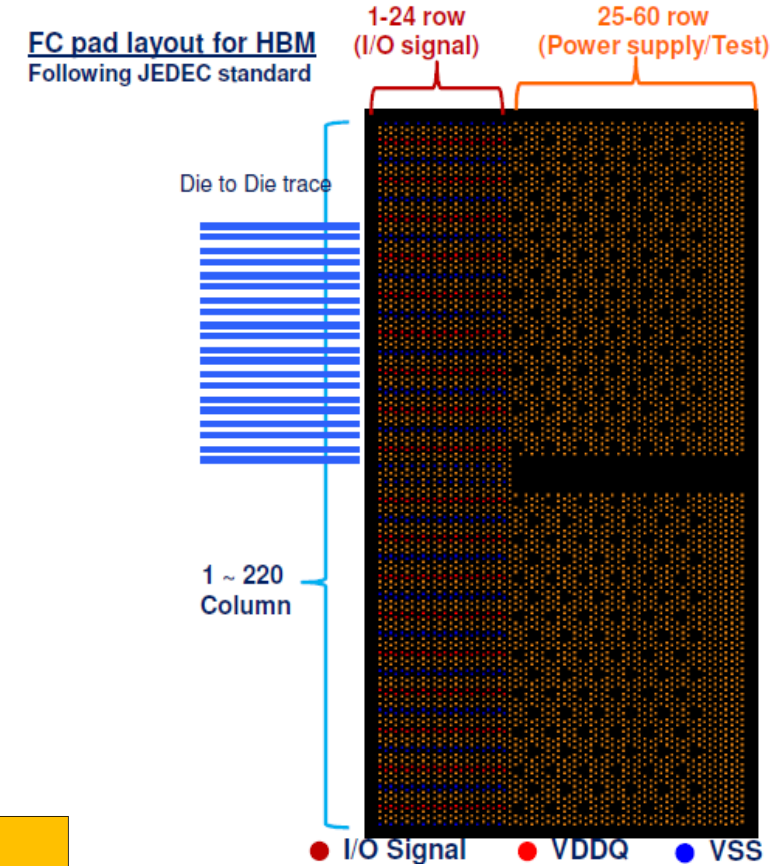
# HBM2 System Overview (Jedec)

❯ **HBM2 system with SOC/DRAM on interposer with 3-6mm length**

❯ **24 signals across 55um u-bump pitch across interface**

❯ **Supports 2Gb/s PHY (1Tb/sec bandwidth for 4-Hi)**

DRAM
DRAM
DRAM
DRAM

**Buffer Die**

SOC

interposer

0.4pf   1.6 pf   0.4pf

3-6 mm

| (Jedec) Interposer Parameters * | | | | |
|---|---|---|---|---|
| Width | Space | Thickness | Length | Resistance |
| 2u | 2u | 0.5-2u | 6mm | 36 ohms |

2u

2u   2u

0.5-2u

2u   2u

FC pad layout for HBM
Following JEDEC standard

1-24 row (I/O signal)
25-60 row (Power supply/Test)

Die to Die trace

1 ~ 220 Column

● I/O Signal   ● VDDQ   ● VSS

Pad Dia.=25um

55um

96um

FC pad pitch for HBM I/O

**XILINX** ❯ ALL PROGRAMMABLE.

# Interposer Design Tools & Methodology

- ## Vertical Routing
  - a model from ubump to package pin is generated and used by high frequency designs (e.g. GT and IO)
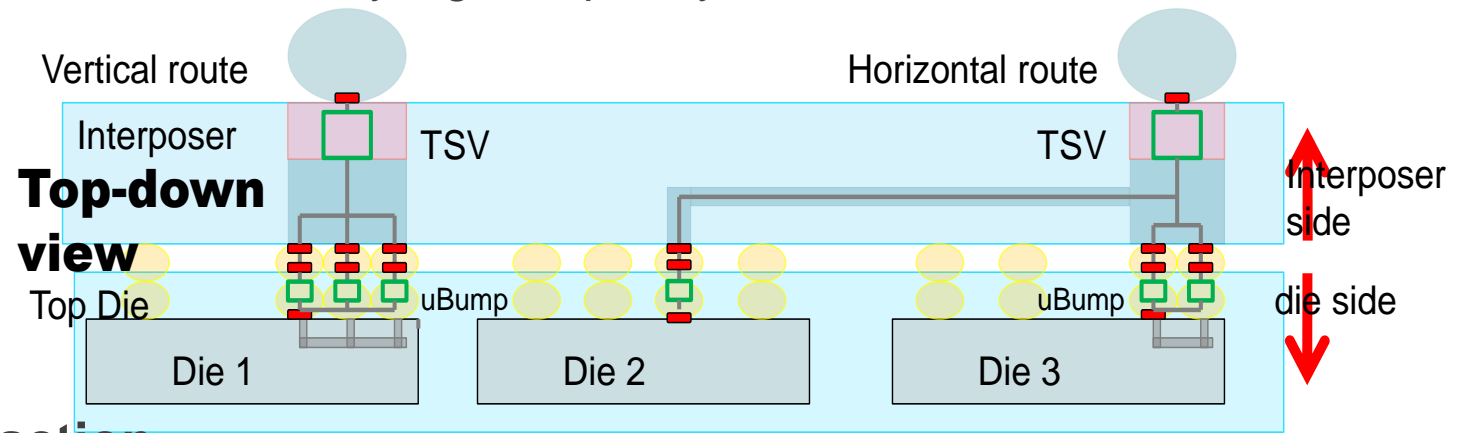- ## Horizontal Routing
  - ### Die LVS and extraction
    - Standard extraction
    - ubump is extracted as a subcircuit
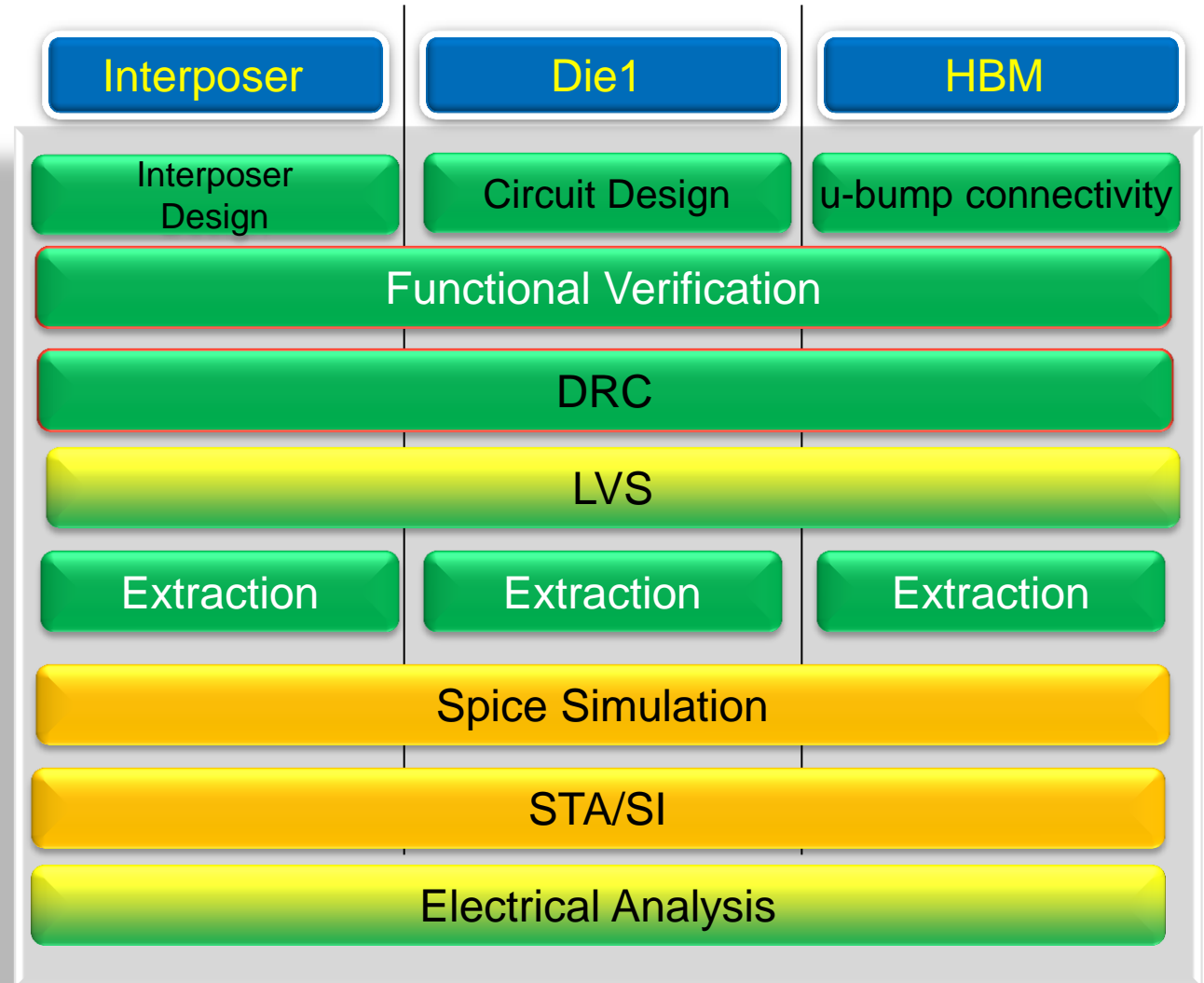  - ### Interposer with box die LVS and extraction
    - Interposer metal extraction
    - TSV is extracted as a subcircuit
- ## Combine the extracted netlists from die and interposer for simulation



Vertical route

Horizontal route

**Top-down view**

Interposer

TSV

TSV

Interposer side

Top Die

uBump

uBump

die side

Die 1

Die 2
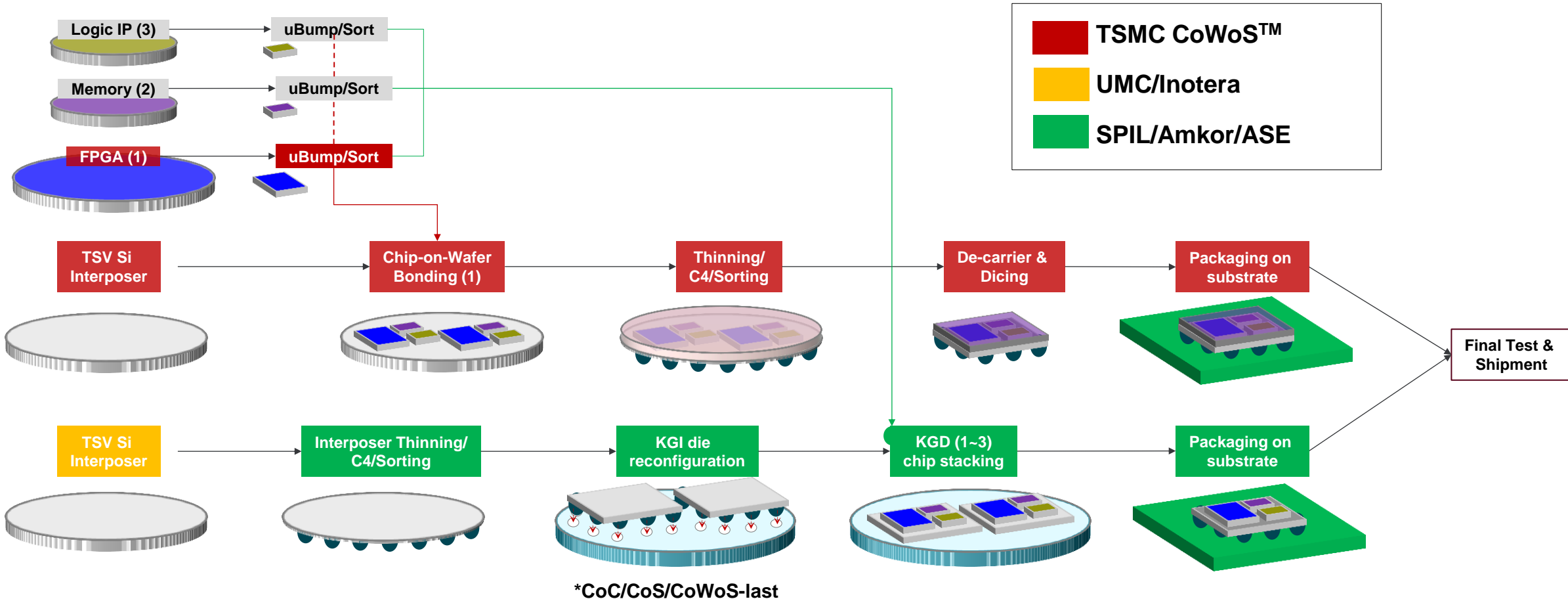
Die 3

XILINX ➤ ALL PROGRAMMABLE.

# Interposer Design Tools & Methodology

- **2.5D - uses the same tool sets as single die design with customized interposer / top die PDK**
- **EDA vendor Tools are validated by TSMC design reference flow**
- **PKG - uses same tool sets as Flip chip (C4-to-BGA)**
  - TSV budget is handled in the Silicon design environment
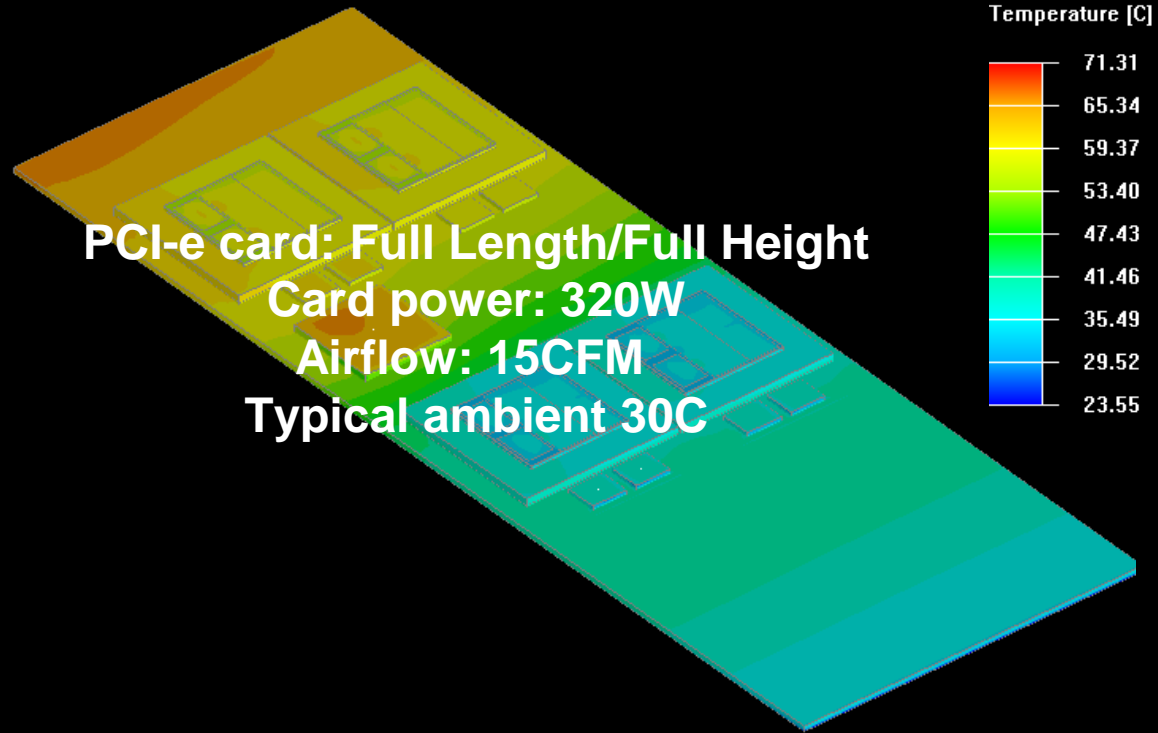  - Layout and PI tools must be capable to handle large data sets

| Interposer | Die1 | HBM |
|---|---|---|
| Interposer Design | Circuit Design | u-bump connectivity |
| Functional Verification | | |
| DRC | | |
| LVS | | |
| Extraction | Extraction | Extraction |
| Spice Simulation | | |
| STA/SI | | |
| Electrical Analysis | | |

**XILINX** ➤ ALL PROGRAMMABLE.

# Supply Chain – Silicon Interposer Approach

➤ **Xilinx in production with 2nd generation of products with TSMC CoWoS**



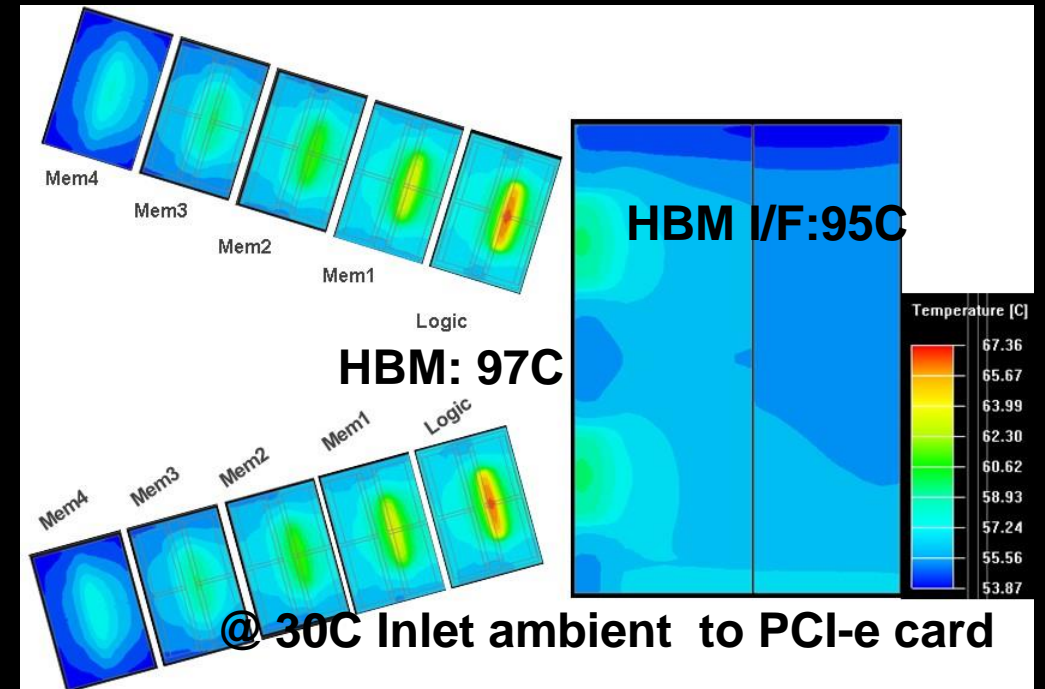| | TSMC CoWoS™ |
| --- | --- |
| | UMC/Inotera |
| | SPIL/Amkor/ASE |

Logic IP (3) → uBump/Sort

Memory (2) → uBump/Sort

FPGA (1) → uBump/Sort

**TSV Si Interposer** → **Chip-on-Wafer Bonding (1)** → **Thinning/ C4/Sorting** → **De-carrier & Dicing** → **Packaging on substrate** → **Final Test & Shipment**

**TSV Si Interposer** → **Interposer Thinning/ C4/Sorting** → **KGI die reconfiguration** → **KGD (1~3) chip stacking** → **Packaging on substrate**

*CoC/CoS/CoWoS-last

© Copyright 2016 Xilinx

**XILINX** ➤ ALL PROGRAMMABLE.

# HBM Integration – HPC Application



Temperature [C]

71.31
65.34
59.37
53.40
47.43
41.46
35.49
29.52
23.55

**PCI-e card: Full Length/Full Height**
**Card power: 320W**
**Airflow: 15CFM**
**Typical ambient 30C**

- **HBM Power map provided by vendors**
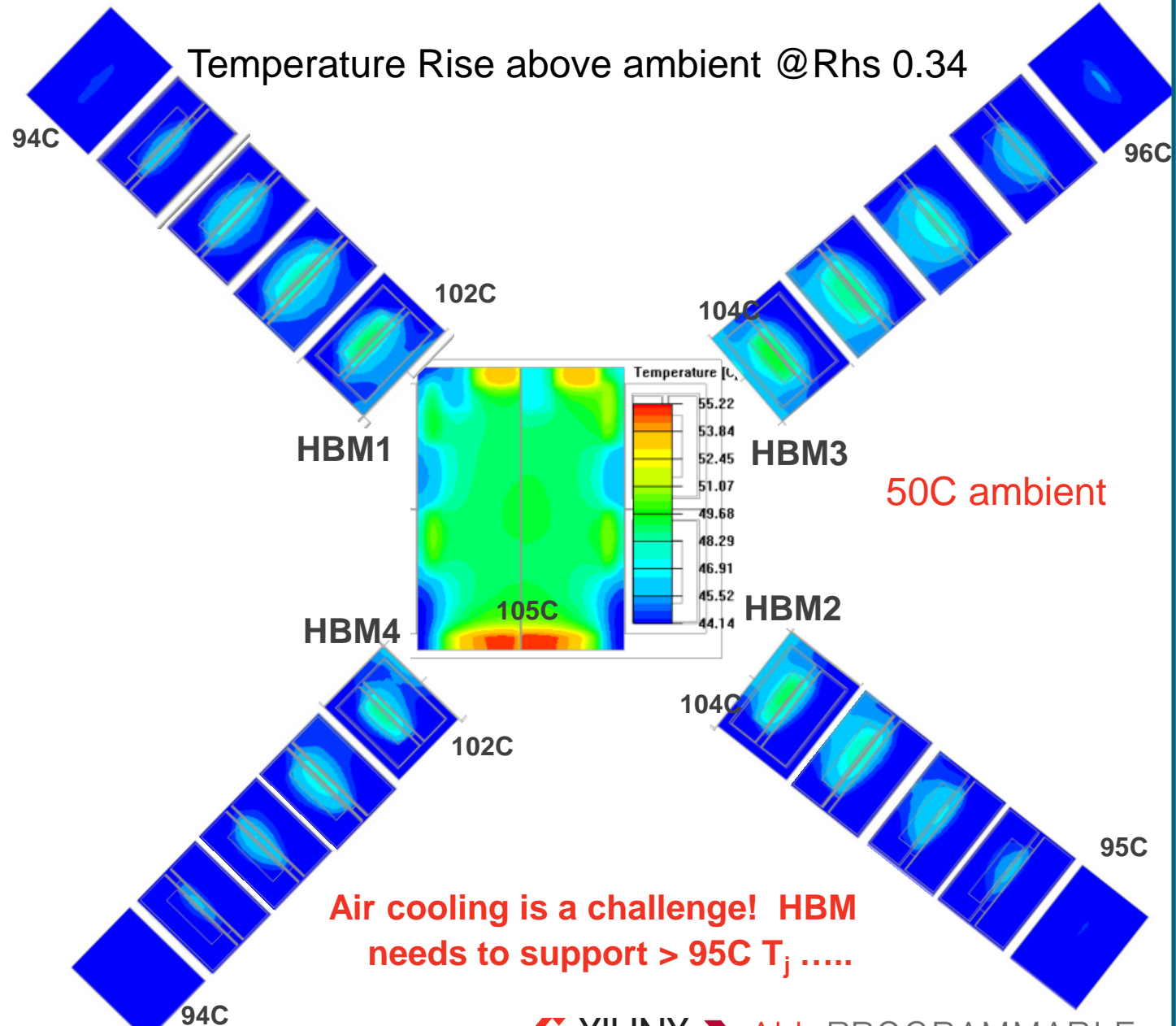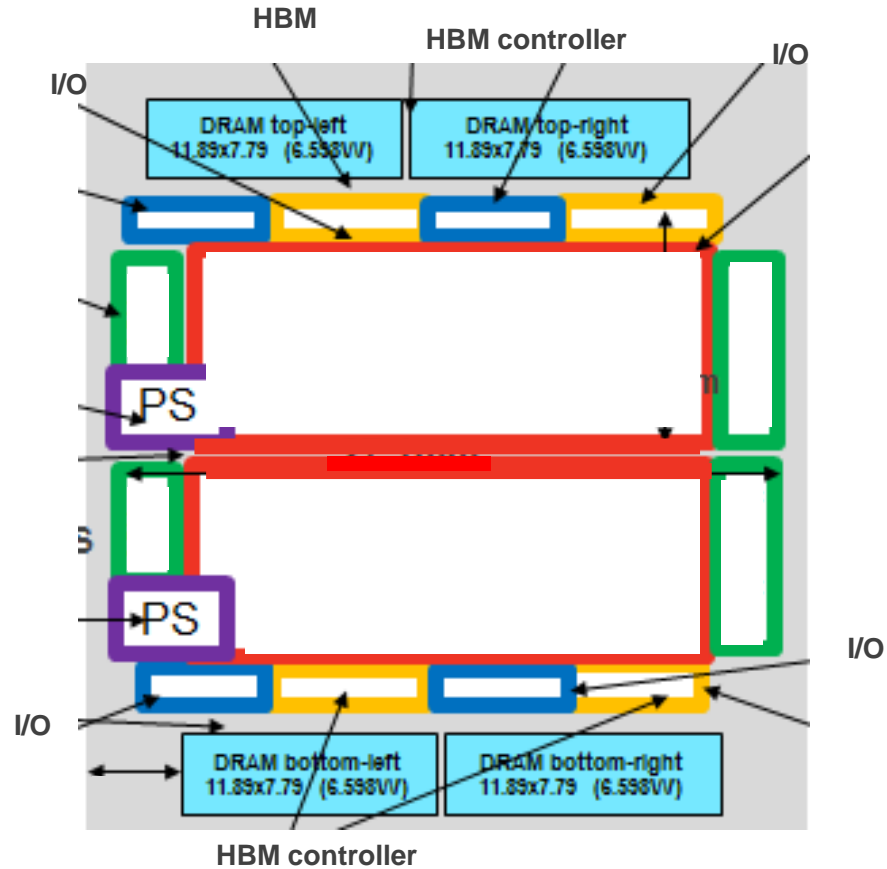- **Thermal model can be done in Flotherm or IcePak environments for example**

HBM can **be 97C** and HBM I/F 96C @30C

HBM gradient ~14C (~2.5C/Layer)

**Air cooling can be a challenge!  HBM 8-Hi needs to support > 95C $T_j$ .....**



HBM I/F:95C

HBM: 97C

@ 30C Inlet ambient  to PCI-e card

Temperature [C]

67.36
65.67
63.99
62.30
60.62
58.93
57.24
55.56
53.87

XILINX ➤ ALL PROGRAMMABLE.

# Telecom Application with HBM

**PKG size:52.5x52.5 mm**
**Total Heat Dissipation is 116.3 W**



Temperature Rise above ambient @Rhs 0.34

94C

96C

102C

104C

**HBM1**

**HBM3**

50C ambient

**HBM4**

105C

**HBM2**

104C

102C

95C

94C

**Air cooling is a challenge!  HBM**
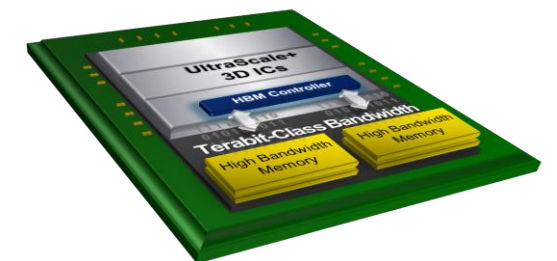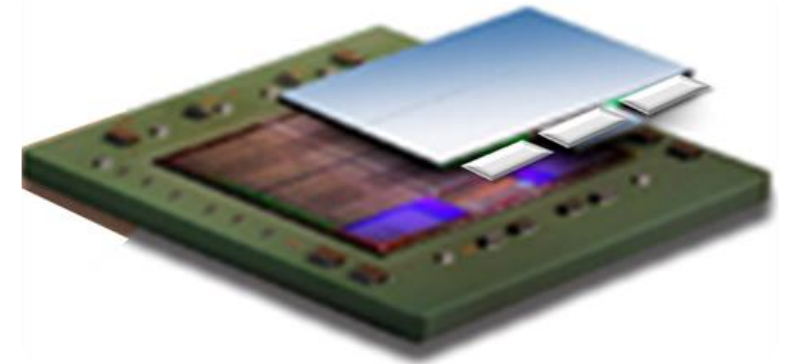**needs to support > 95C $T_j$ …..**

XILINX ALL PROGRAMMABLE

# Summary

> **Tb/s low latency bandwidth and lower system power is driving the need for HBM adoption**

> **Silicon Interposer (2.5D) is the incumbent technology of choice. Potentially lower cost, fine pitch interconnect wafer-level and substrate based technologies are emerging**

> **To drive broader adoption of HBM applications (cooling limited) and higher performance stacks (8-Hi), higher HBM junction temperature (>95C) needs to be supported**

# Follow Xilinx

facebook.com/XilinxInc

twitter.com/XilinxInc

youtube.com/XilinxInc

linkedin.com/company/Xilinx

plus.google.com/+Xilinx