

# ECE 546

## Lecture -27

# Co-Design for HI

Spring 2026

Jose E. Schutt-Aine  
Electrical & Computer Engineering  
University of Illinois  
jesa@illinois.edu

# AI Energy Requirements

- Data centers used 2.5% of US electricity in 2022
- Projected to increase to 20% by 2030

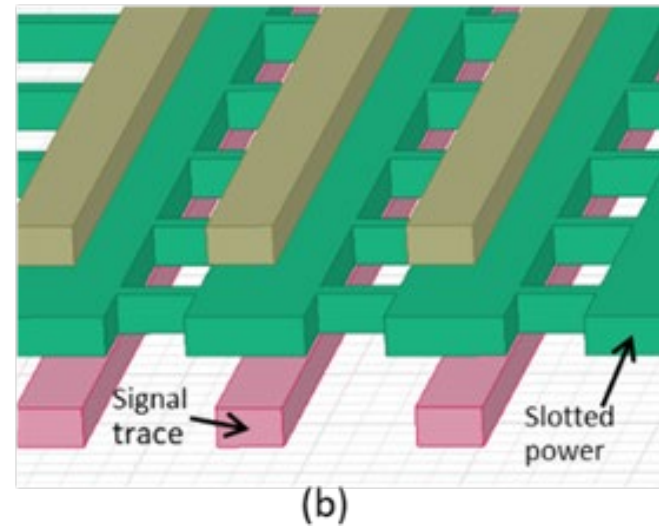
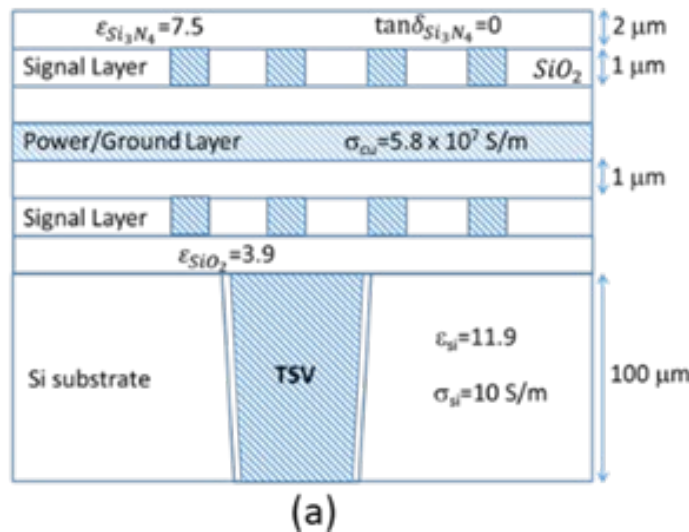
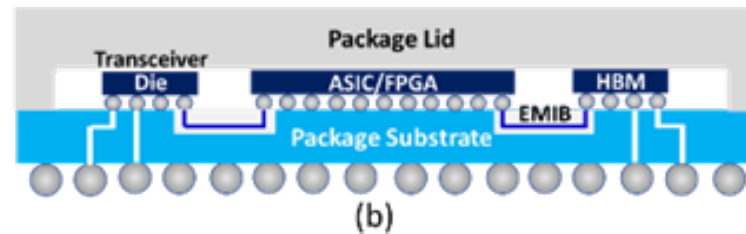
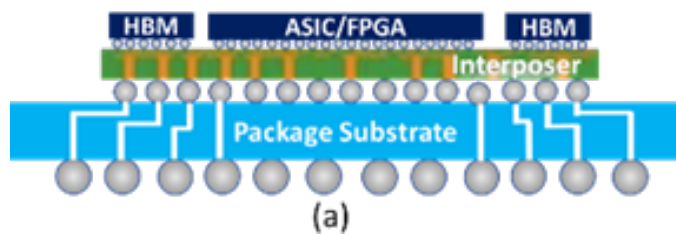
Training a large language model like **GPT-3** is estimated to use **1,300 megawatt hours (MWh)** of electricity.



*80% of power consumption is due to data movement through interconnects*

# HI – The Interconnect Challenge

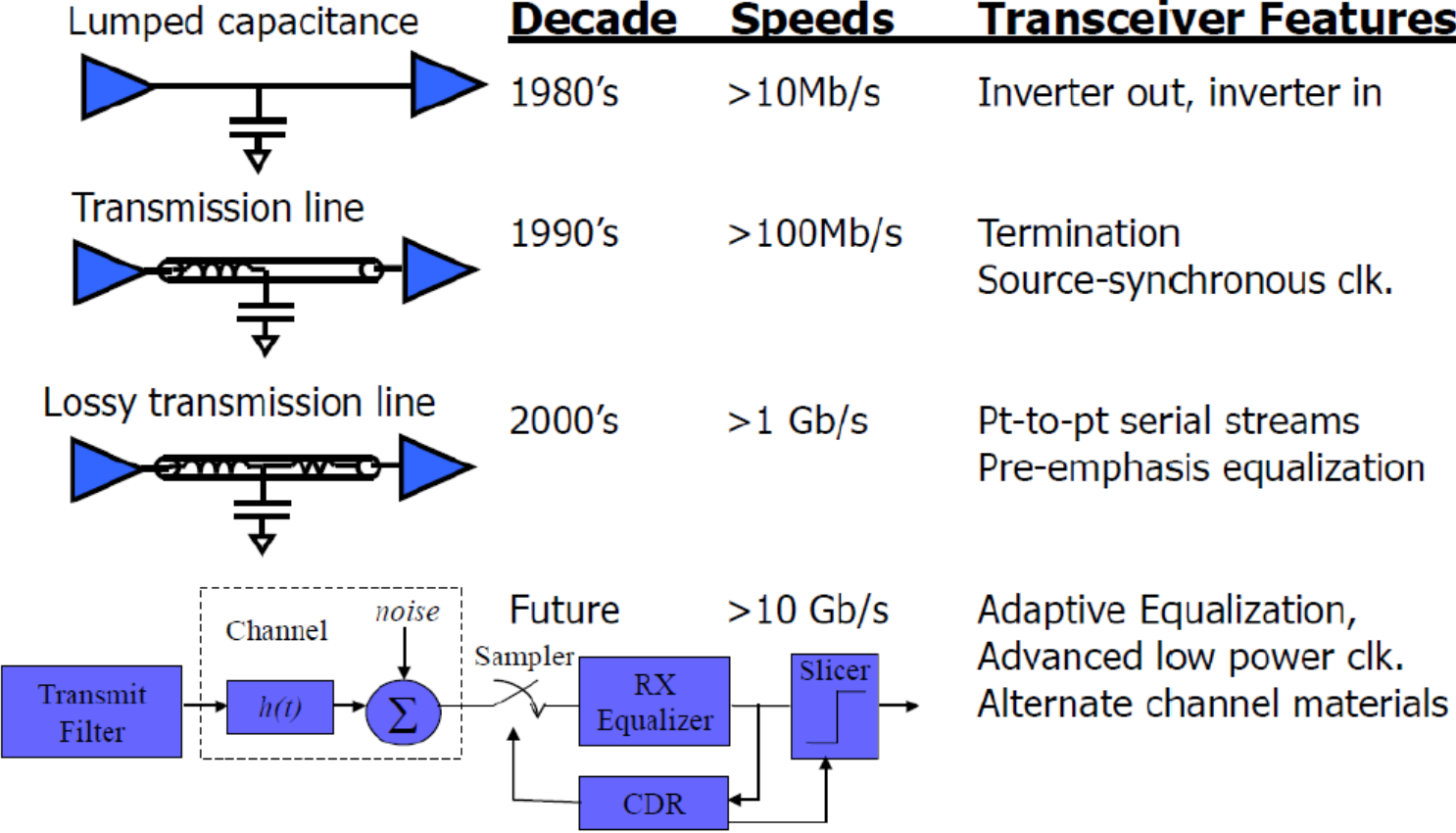
*“There are many solutions and techniques for improving transistors; options for improving interconnects are very few...”*



## Challenges

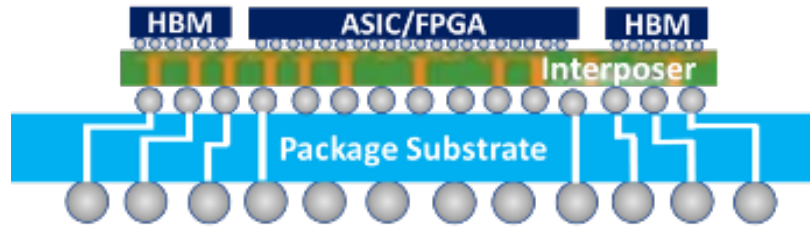
- Signal/power integrity
- Thermal effects
- IR Drop
- Power dissipation
- High I/O count
- Reliability
- Security
- Environment
- Architecture

# Example - Interconnect Evolution



Slide Courtesy of Frank O'Mahony & Brian Casper, Intel

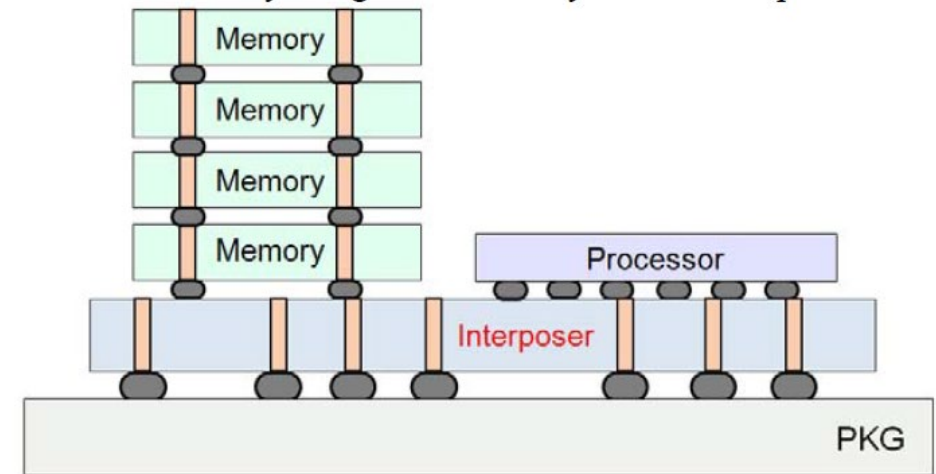
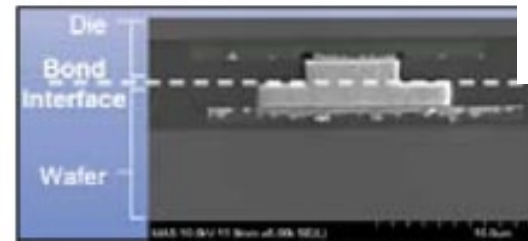
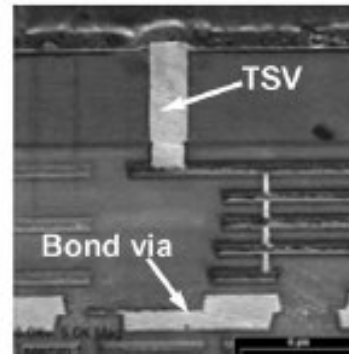
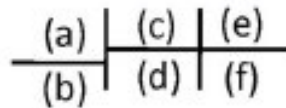
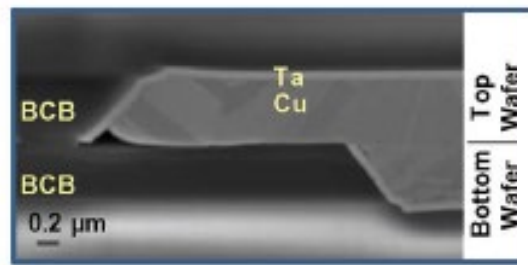
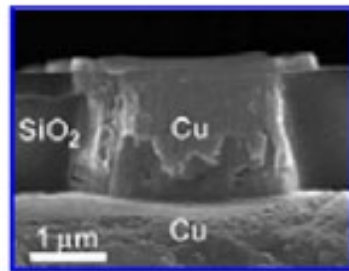
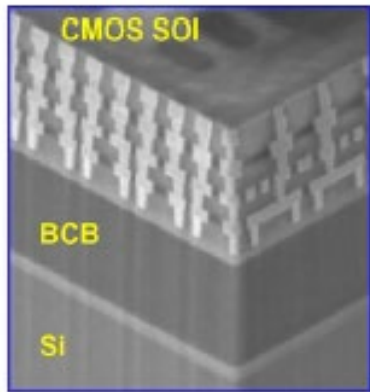
# Advanced Packaging



(a)



(b)



***In order to meet the demands of AI, HPC, and next-generation data centers, advanced packaging will be key to performance, power efficiency, and scalability.***

# Hybrid Bonding

*simultaneous bonding of dielectric and metal bond pads in one bonding step*

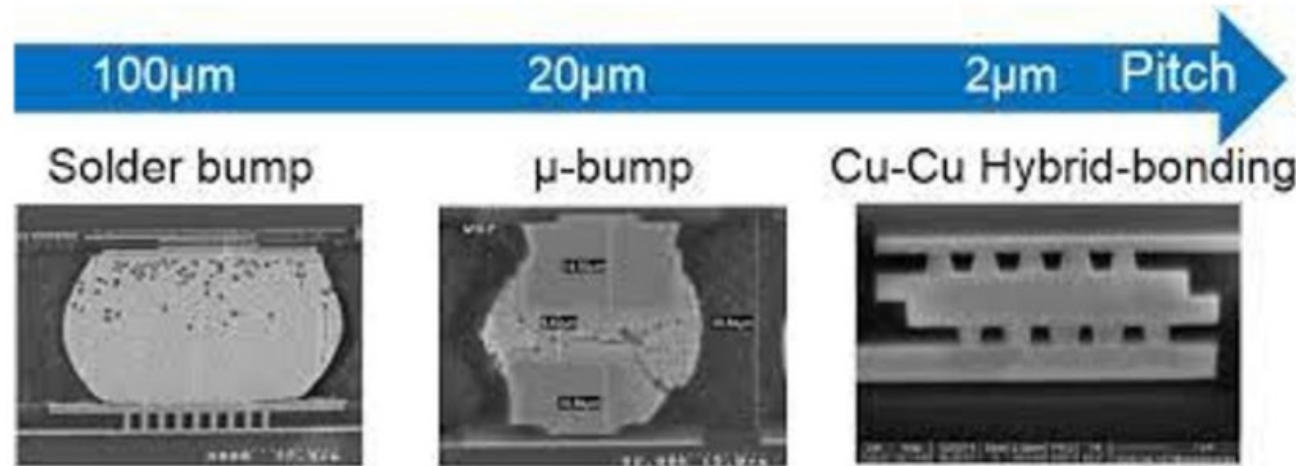
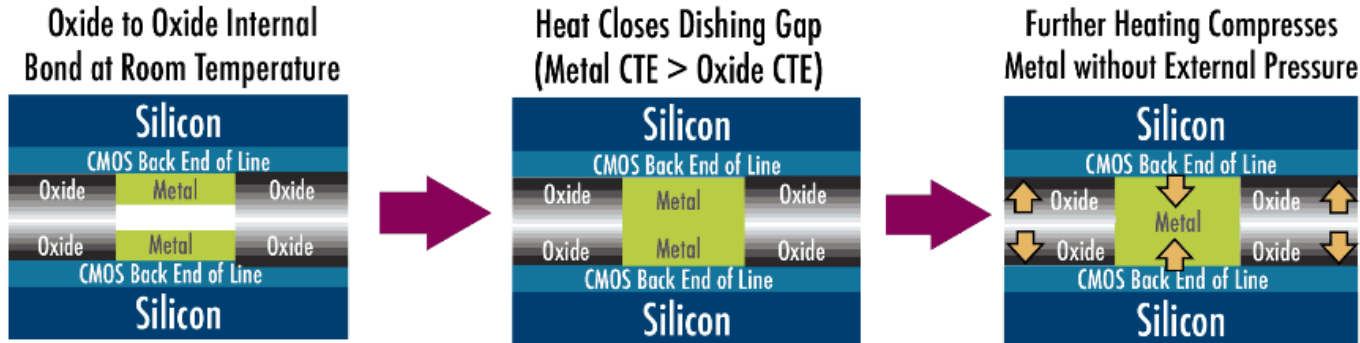
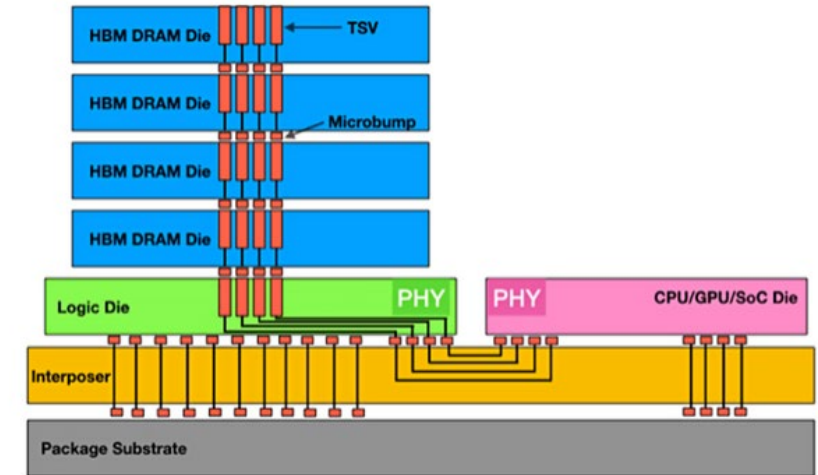


Image Credit: Imed Jani. Test and characterization of 3D high-density interconnects. Micro and Nanotechnologies/Microelectronics. Université Grenoble Alpes, 2019. English. NNT : 2019GREAT094 . tel- 02634259



HBM stack for maximum data throughput. Source: Rambus

- Allows advanced 3D device stacking
- Highest I/O
- Enables sub-10-µm bonding pitch
- Higher memory density
- Expanded bandwidth
- Increased power
- Improved speed efficiency
- Eliminates the need for bumps, improving performance with no power or signal penalties

# Heterogeneous Integration

- **Heterogeneous Integration is defined as the integration of separately manufactured components into a higher-level assembly (Chiplets, SiPs, Modules) that, in the aggregate, provides enhanced functionality and improved operating characteristics**
- **The size, complexity and stochasticity associated with future advanced packaging platforms will render their design intractable using the traditional approaches. In addition, such design will need to be performed with consideration for thermal management, reliability, architecture and floorplanning/placement and routing constraints**

# Objectives for HI

- Focus on minimizing **energy** and **delay**
- Move away from Von Neumann computing model
- Identify and address conflicting requirements,
- Make use of new physics and emerging technologies
- Future applications are drivers
- End of Moore's Law necessitates new paradigm

*Universal co-design implies collaboration between many disciplines with the common goal of improved performance at lower cost*

# Need for: Heterogeneous Integration

- **Motivation: Ubiquitous materials and technologies**
  - Die, interposer, substrate
  - Scaling is with wavelength rather than technology
- **Opportunity: New interconnect technologies**
  - Increased density of I/Os, finer pitches
  - 2.5D/3D integration
- **Performance objective: minimize energy and delay**
- **Approach: Use chiplets**

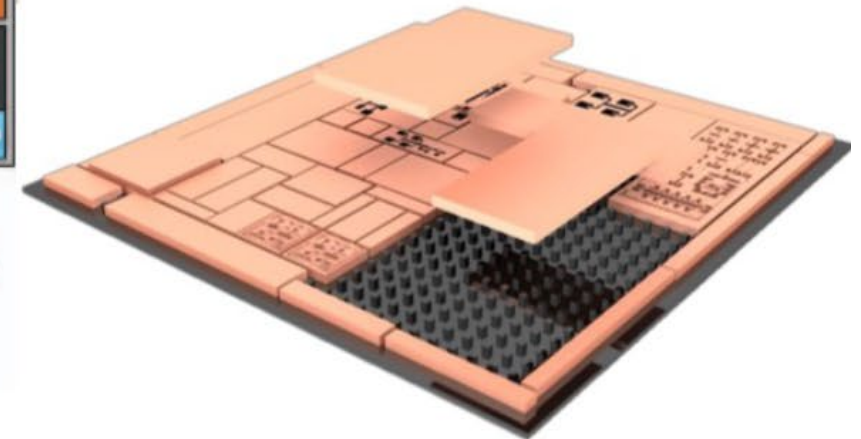
# Chiplet-Based Design: AMD EPYC



Monolithic 32-core Chip  
777mm<sup>2</sup> total area  
1.0x Cost

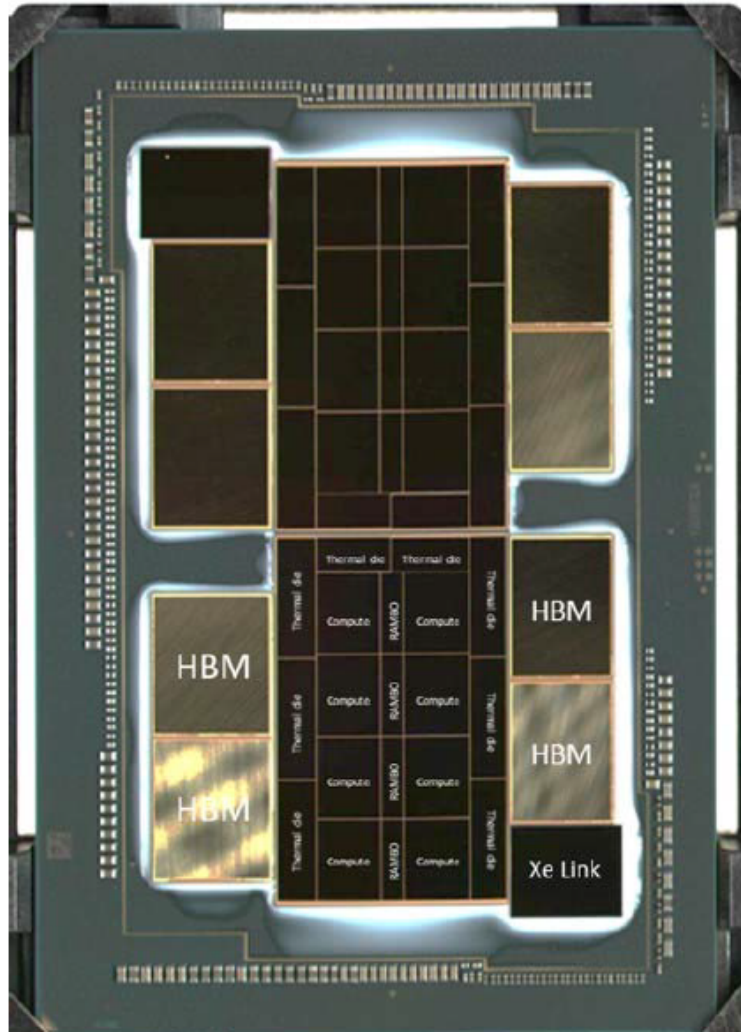


4 x 8-core Chiplet, 213mm<sup>2</sup> per chiplet  
852mm<sup>2</sup> total area (+9.7%)  
0.59x Cost



Naffziger, S., Noah Beck, T. Burd, K. Lepak, Gabriel H. Loh, M. Subramony and Sean White.  
"Pioneering Chiplet Technology and Design for the AMD EPYC™ and Ryzen™ Processor Families :  
Industrial Product." 2021 ACM/IEEE 48th Annual International Symposium on Computer  
Architecture (ISCA) (2021): 57-70.

# Chiplet-Based Design: Intel's Ponte Vecchio



<b>Integration</b>	Foveros + EMIB
<b>Power Envelope</b>	600W
<b>Transistor count</b>	> 100B
<b>Total Tiles</b>	63 ( 47 functional + 16 thermal Tiles)
<b>HBM count</b>	8
<b>Package Form factor</b>	77.5 x 62.5 mm (4844 mm <sup>2</sup> )
<b>Platforms</b>	3 platforms
<b>IO</b>	4x16 90G SERDES, 1x16 PCIe Gen5
<b>Total Silicon</b>	3100 mm <sup>2</sup> Si
<b>Silicon footprint</b>	2330 mm <sup>2</sup> Si footprint
<b>Package layers</b>	11-2-11 (24 layers)
<b>2.5D Count</b>	11 2.5D connections
<b>Resistance</b>	0.15 mΩ R <sub>path</sub> /tile
<b>Package pins</b>	4468 pins
<b>Package Cavity</b>	186 mm <sup>2</sup> x4 cavities

Wilfred Gomes, Altug Koker, Pat Stover, Doug Ingerly, Scott Siers, Srikrishnan Venkataraman, Chris Pelto, Tejas Shah, Amreesh Rao, Frank O'Mahony, Eric Karl, Lance Cheney, Iqbal Rajwani, Hemant Jain, Ryan Cortez, Arun Chandrasekhar, Basavaraj Kanthi, Raja Koduri, "Ponte Vecchio: A Multi-Tile 3D Stacked Processor for Exascale Computing", ISSC 2022

# Need for: Chiplets

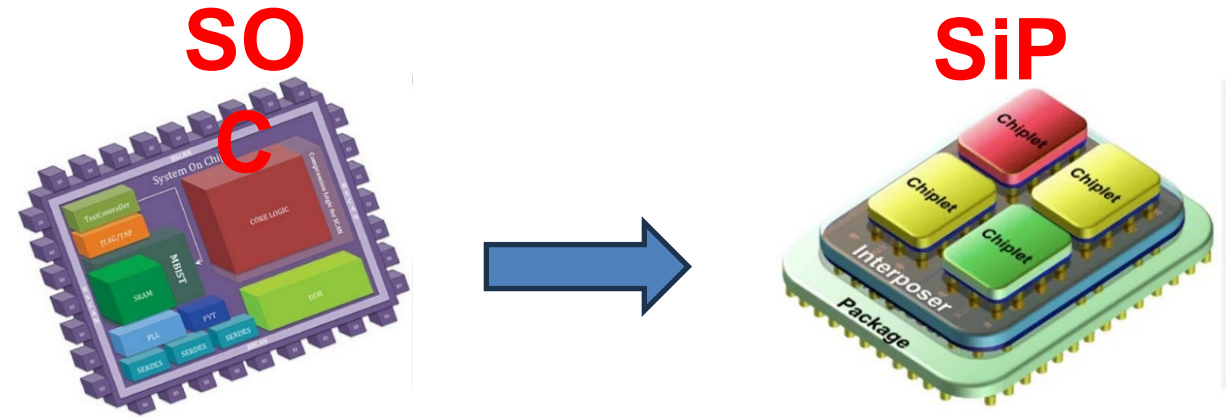
- **Advantages**

- Better cost per die, higher yield, mix and match, reusability

- **Challenges**

- Disaggregation is ad-hoc
- Multi-level, multi-scale, increased complexity

**Solution: Co-Design**



# What is Co-Design?

- **Level Co-design**
  - Chip
  - Package
  - Board
- **Physics-based Co-design**
  - Thermal
  - Electrical
  - Mechanical
  - Optical
- **Function-based Co-design**
  - Thermal aware
  - Signal integrity aware
  - Testability
  - Security aware
- **Domain-based Co-design**
  - Hardware
  - Software
  - Architecture

# Design Flow

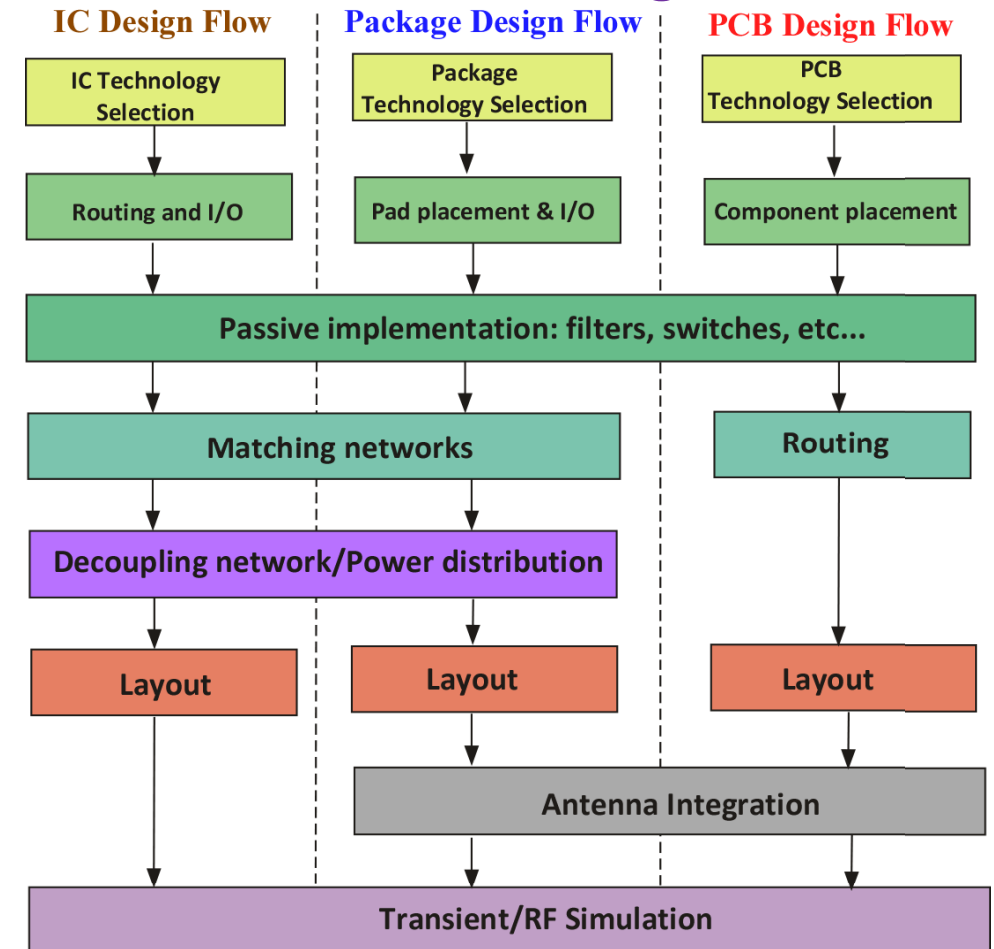
## VERTICAL

- Bridge IC, Package, Board
- Bridge System, Architecture, Layout
- Bridge Synthesis, Analysis, Verification
- Bridge Hardware, Software, Firmware

## HORIZONTAL

- Energy aware
- Signal/power integrity aware
- Stress/thermal aware
- Security aware
- Testing aware

## Traditional Co-Design Flow



# Co-Design

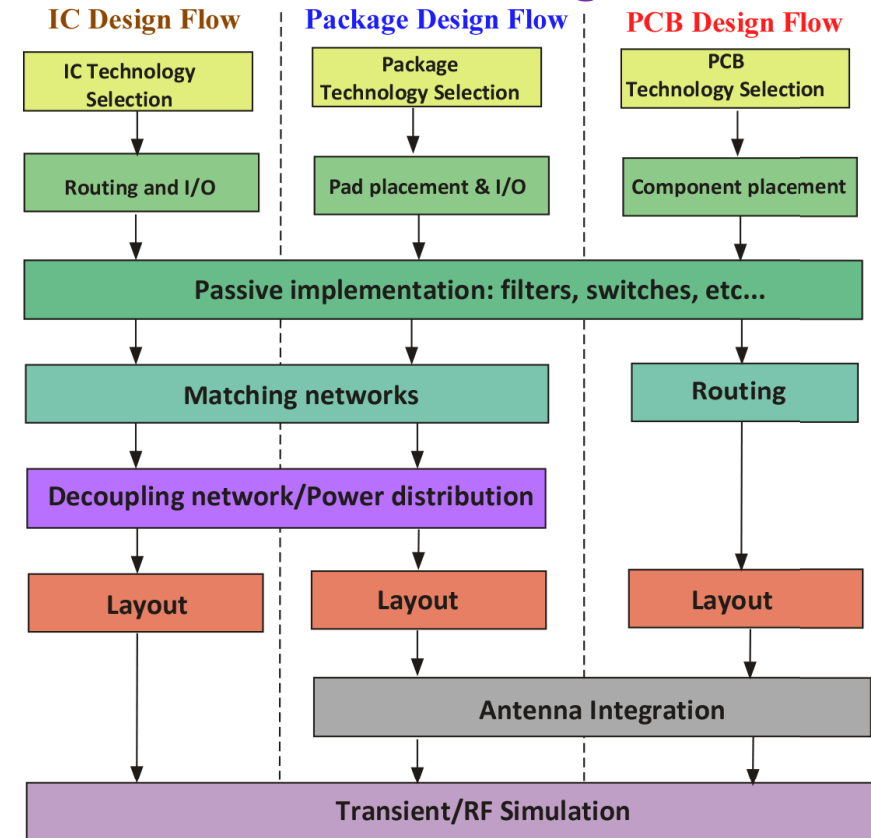
## VERTICAL

- Bridge IC, Package, Board
- Bridge System, Architecture, Layout
- Bridge Synthesis, Analysis, Verification
- Bridge Hardware, Software, Firmware

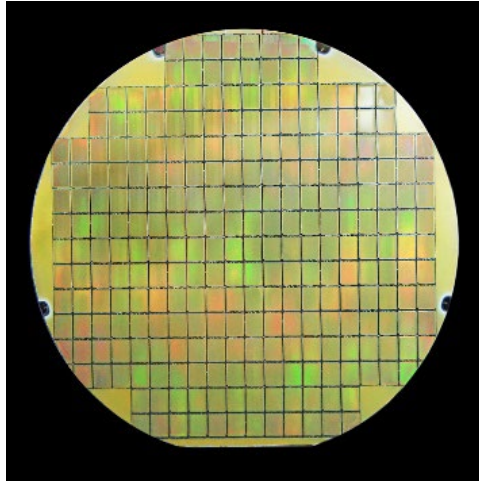
## HORIZONTAL

- Energy aware
- Signal/power integrity aware
- Stress/thermal aware
- Security aware
- Testing aware

## Traditional Co-Design Flow

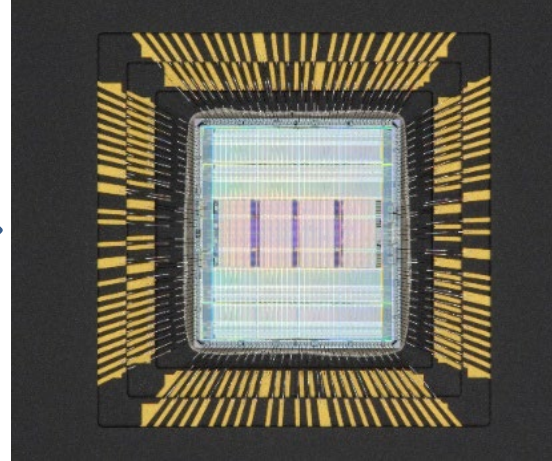
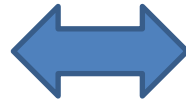


# Traditional Co-Design



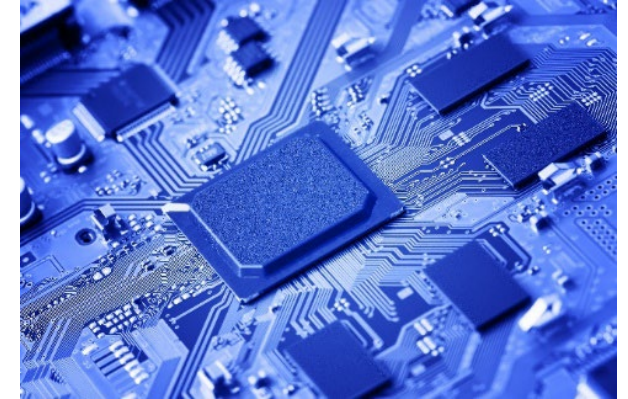
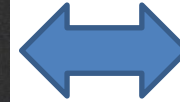
CHIP

- Transistors
- Nonlinear
- SPICE
- Scaling with tech



PACKAGE

- Interconnects
- Linear
- EM Tools
- Scaling with  $\lambda$



BOARD

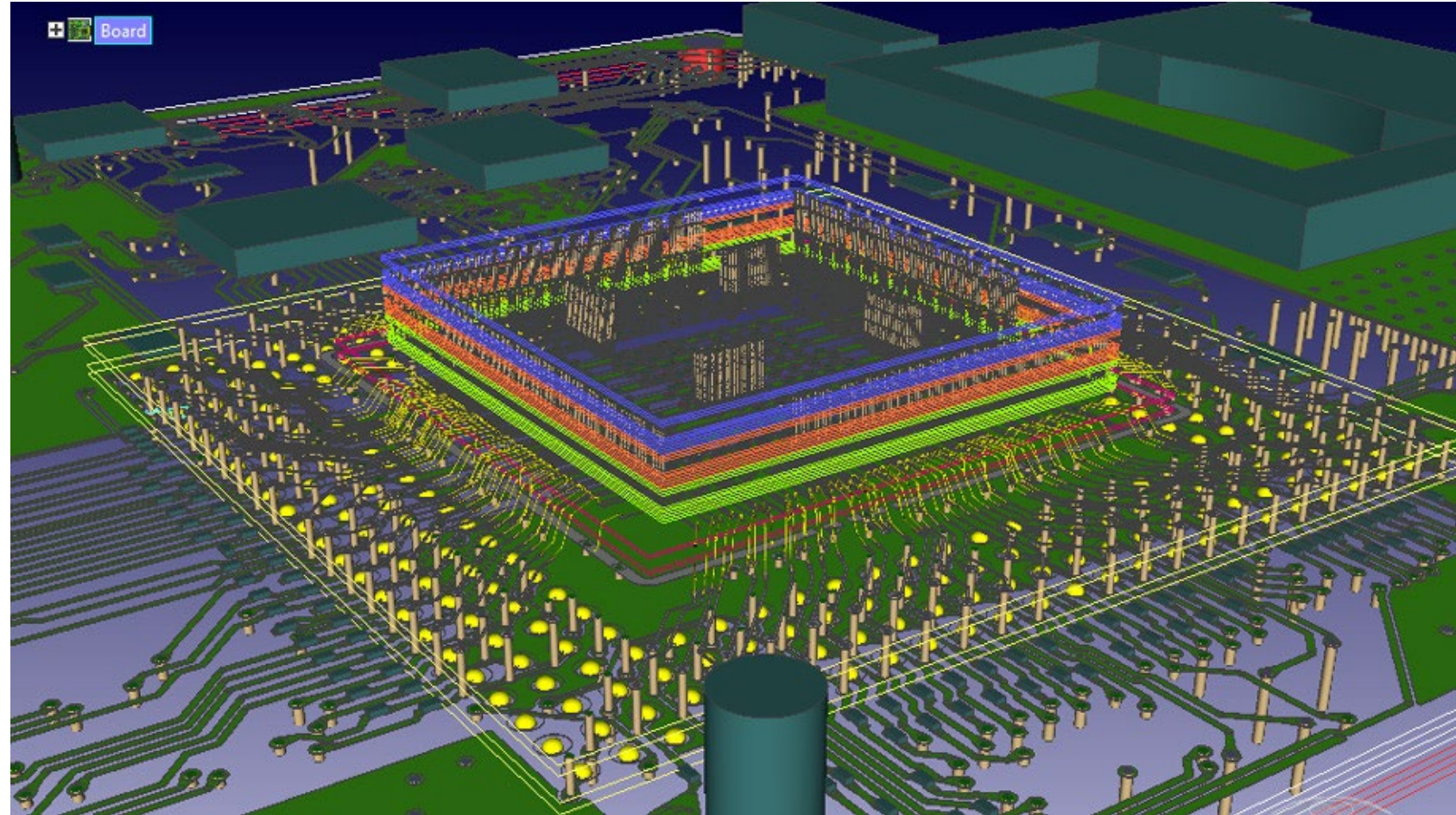
- Transmission lines, sensors
- Linear+Nonlinear
- EM Extraction, SPICE, IBIS,...
- Scaling with  $\lambda$

# Key Questions

- **What is the state of the art in co-design?**
- **What system products will drive creation of new co-design tools/methods?**
- **What are the key challenges that need to be overcome?**
- **What added value will new co-design tools/methods bring to those products?**
- **What needs to happen for these challenges to be overcome?**
- **When will answers to the above become apparent?**

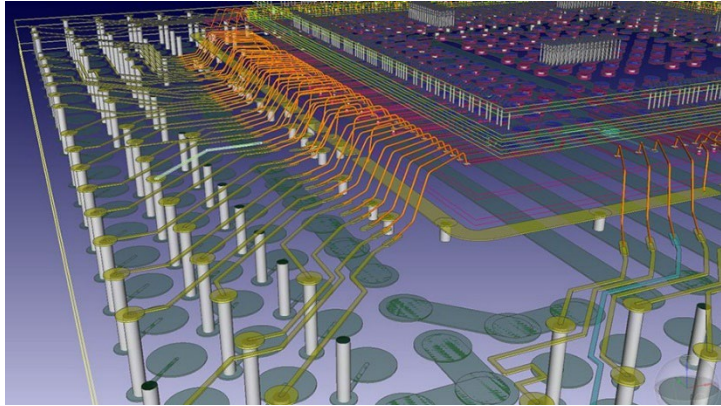
# Co-Design Requirements

- Tradeoffs in advance
- Translation and domains
- Propagate information
- Manage connectivity
- Database formats

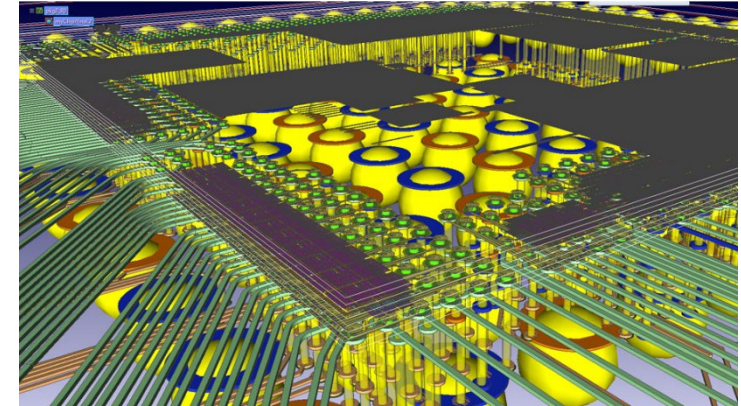


Courtesy of Zuken

# Pathfinding Methodologies



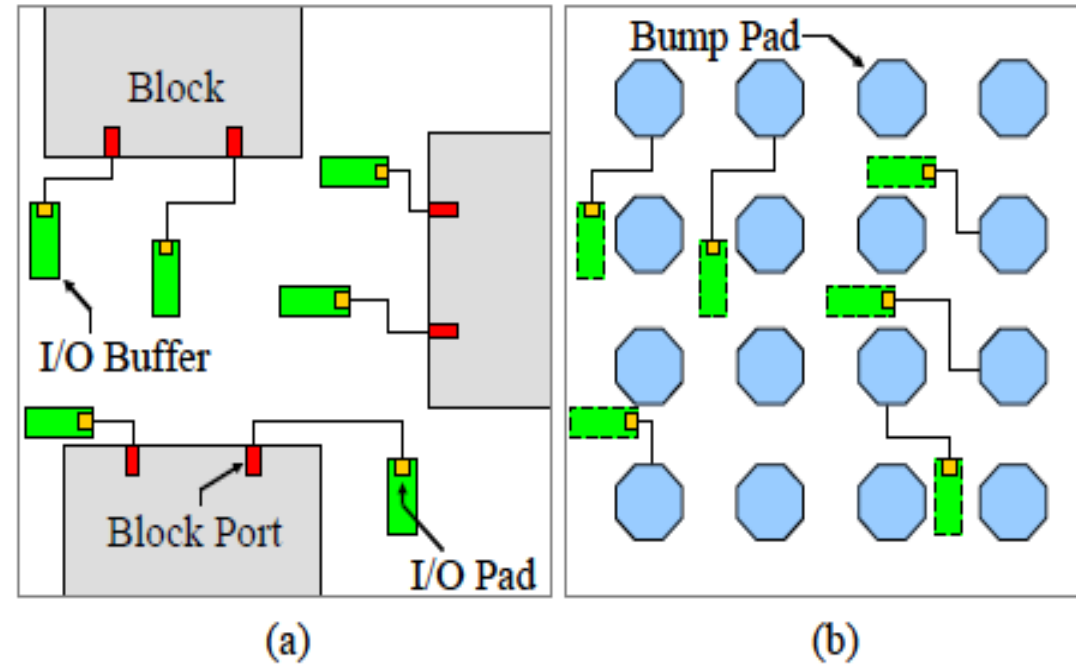
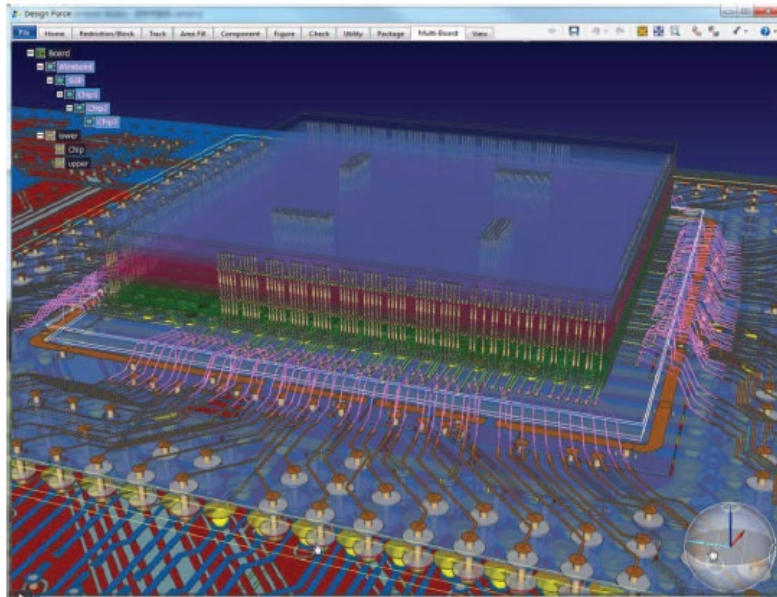
Courtesy of Zuken



Courtesy of Zuken

- Unified workflow, including partitioning, floor-planning, design of system-level interconnects, route pathway exploration and feasibility analysis. Capability to create abstract package models and virtual die models from multiple sources
- Ability to visualize and modify component placement scenarios and make connectivity changes in a preliminary floorplan. Provision of dynamic manipulation of pin arrays within the abstract models
- Preserving signal assignments and rules while making adjustments to the physical pin array. Support for multiple package variables and PCB form factors to verify and compare different system configurations
- Enabling the interaction of design tools from different EDA vendors

# Challenges for Placement and Routing

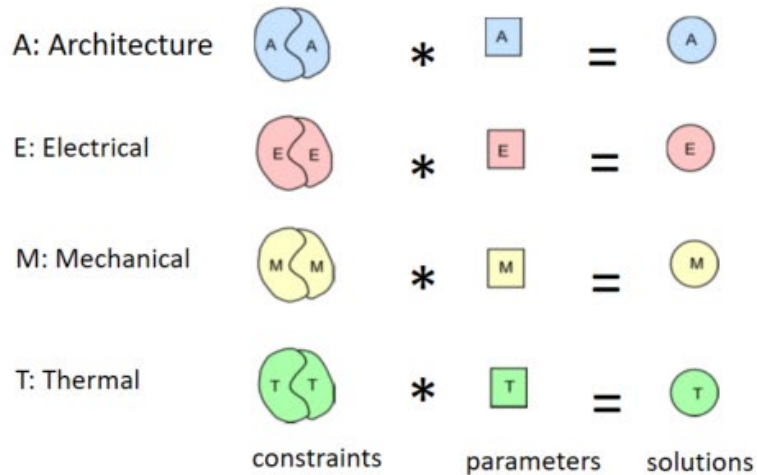


*Need to concurrently manage connectivity between intra- and inter-levels*

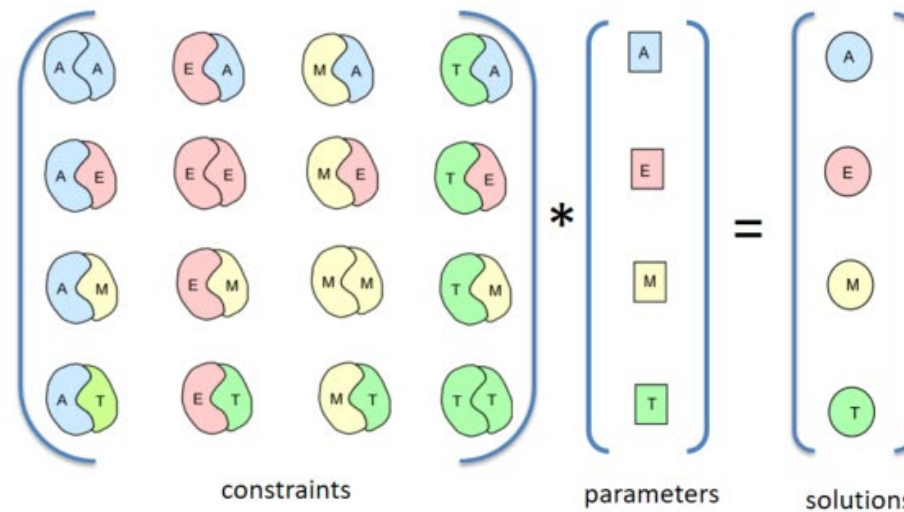
# Co-Design Challenges

Current co-design methods are simply a cascade or combination of independent solutions. Real co-design requires concurrent solution from a formulation that accounts for all multi-physics interactions while embodying conflicting requirements

## Traditional Approach



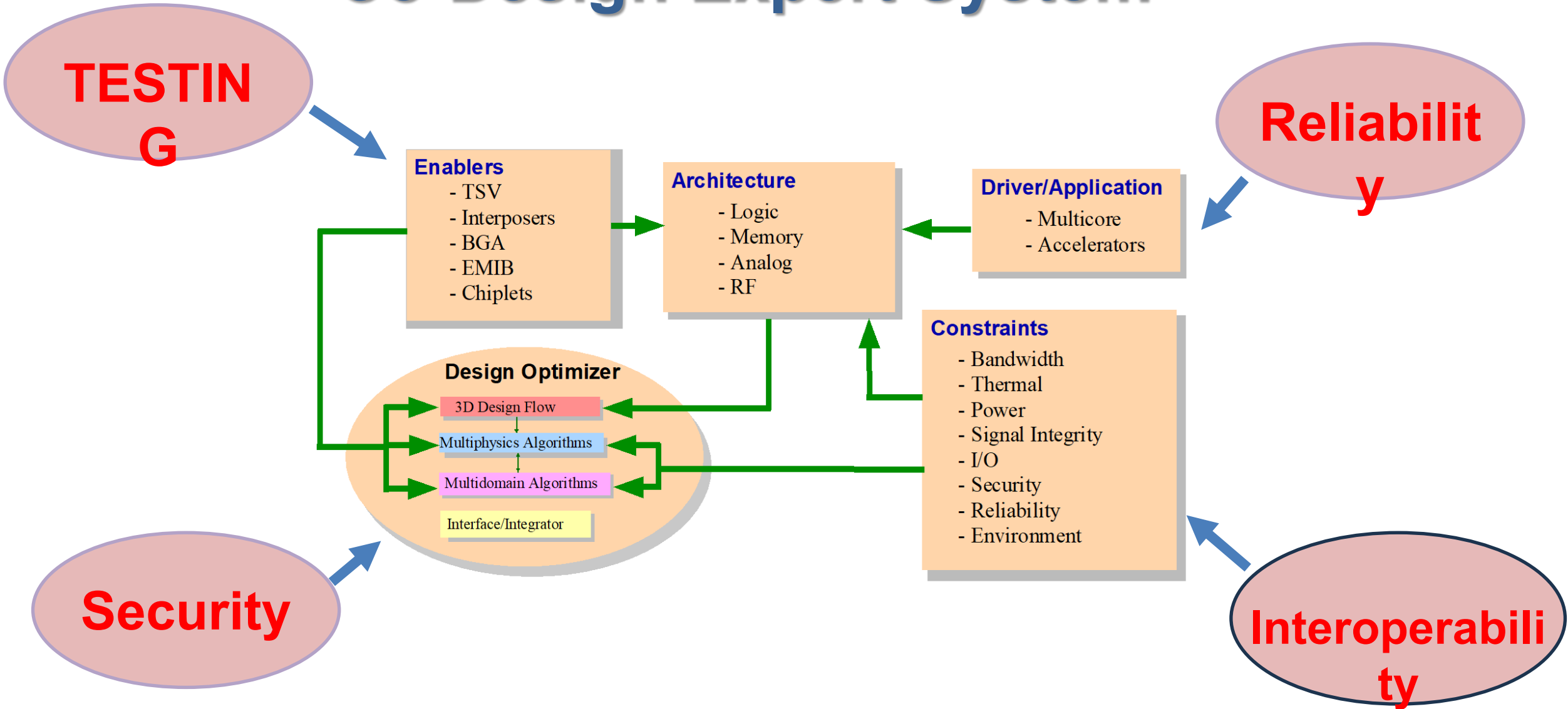
## Co-Design Approach



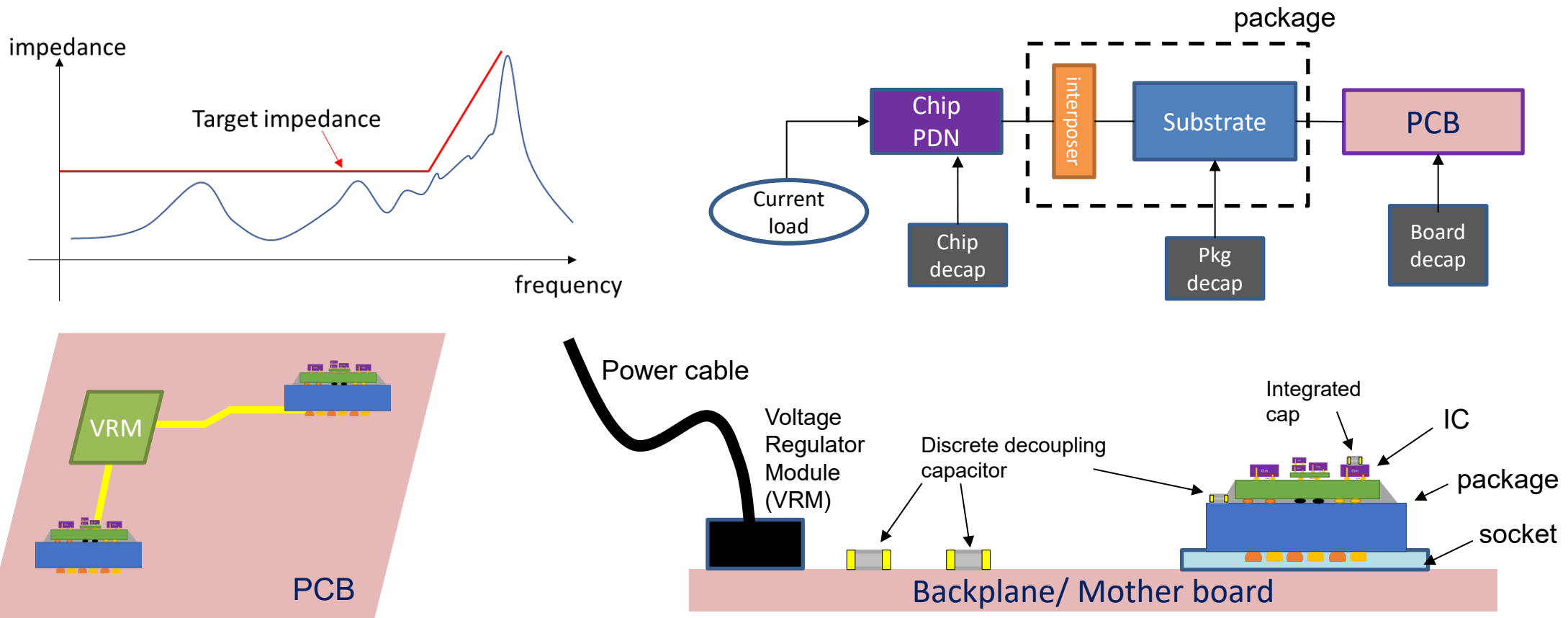
## Constraints

- Bandwidth
- Thermal
- Power
- Signal Integrity
- I/O
- Security
- Reliability
- Environment
- Size
- Cost

# Co-Design Expert System



# Power Distribution Network (PDN)

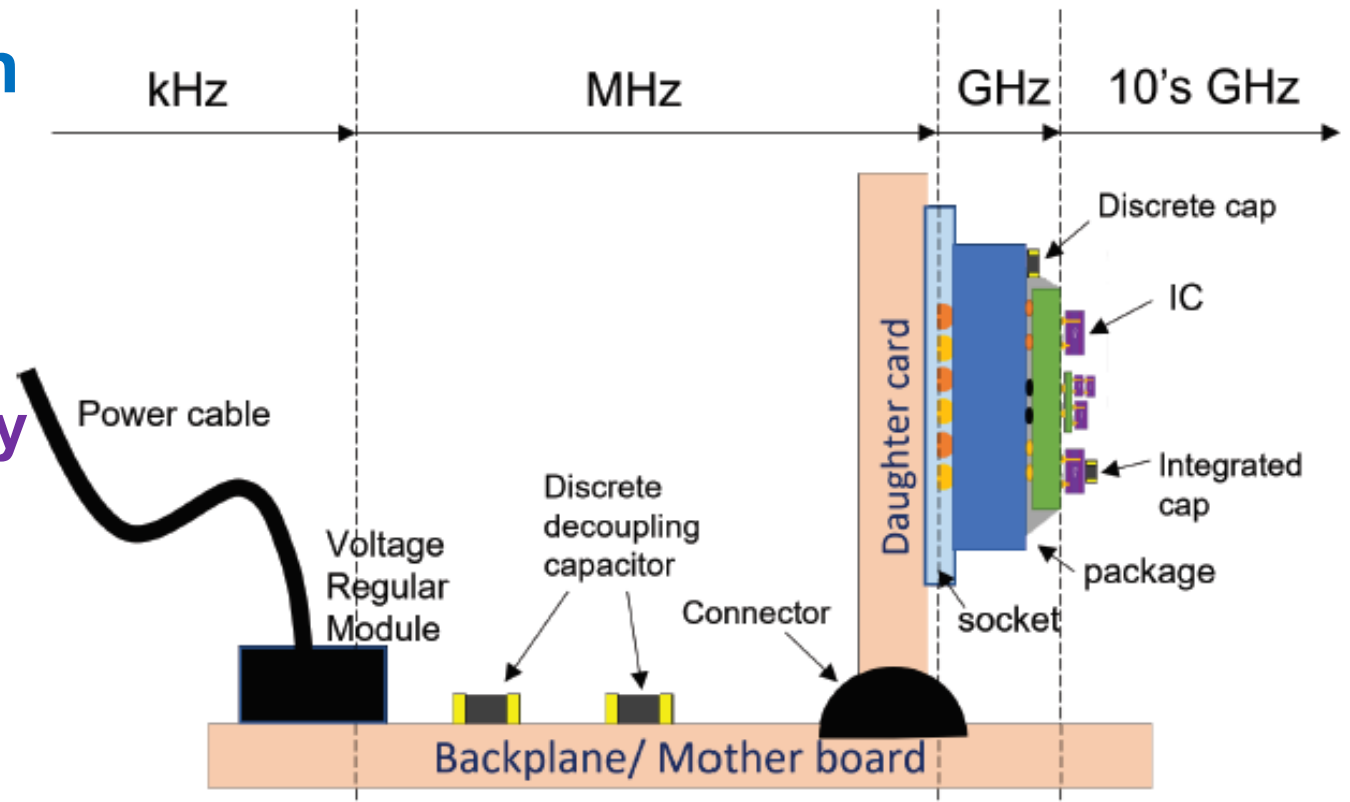


**Goal: Achieve target impedance through use of decoupling capacitors at chip, package and board levels**

# Chipllets and Heterogeneous Integration

- **Multilevel Power Distribution**

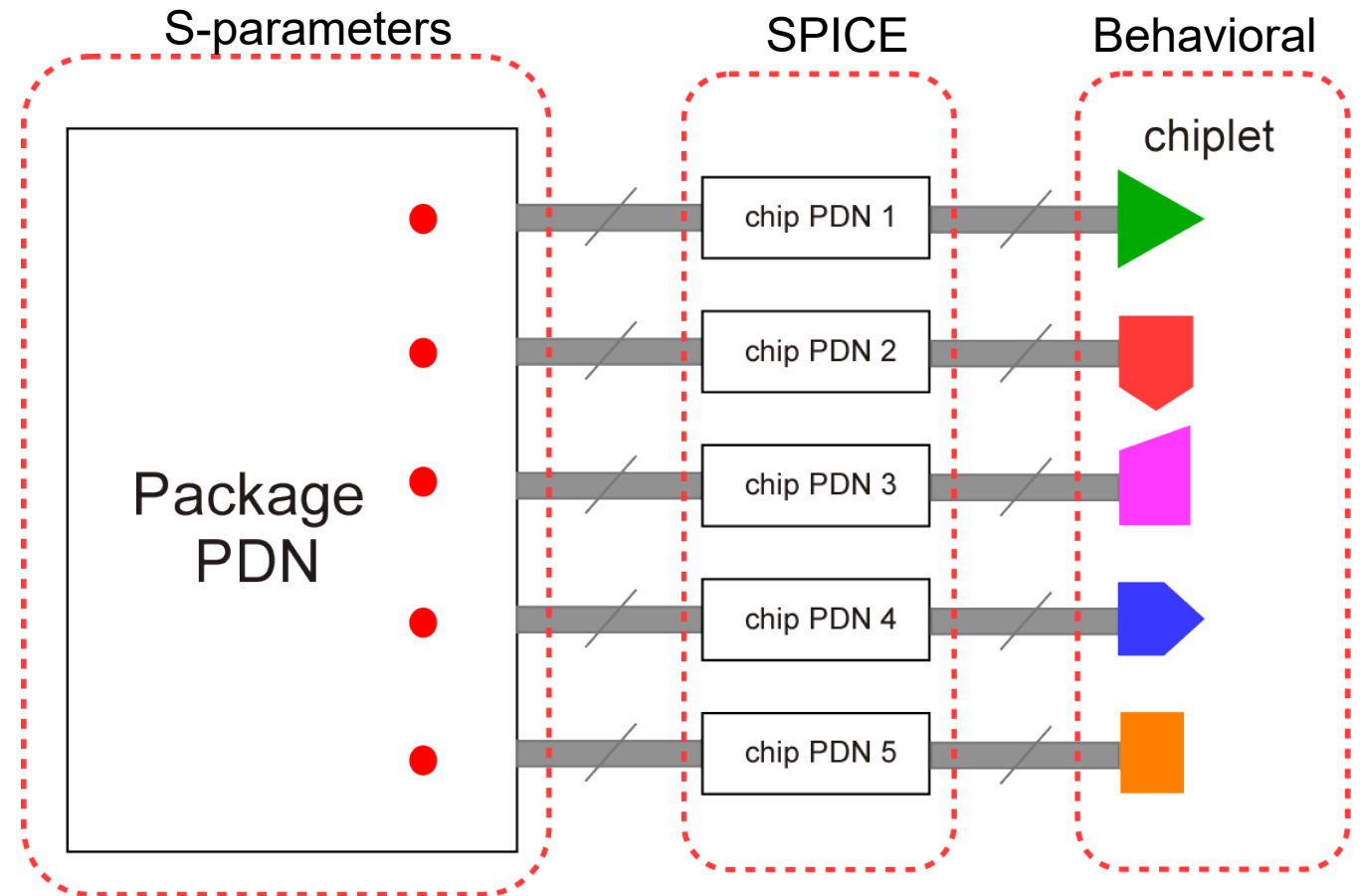
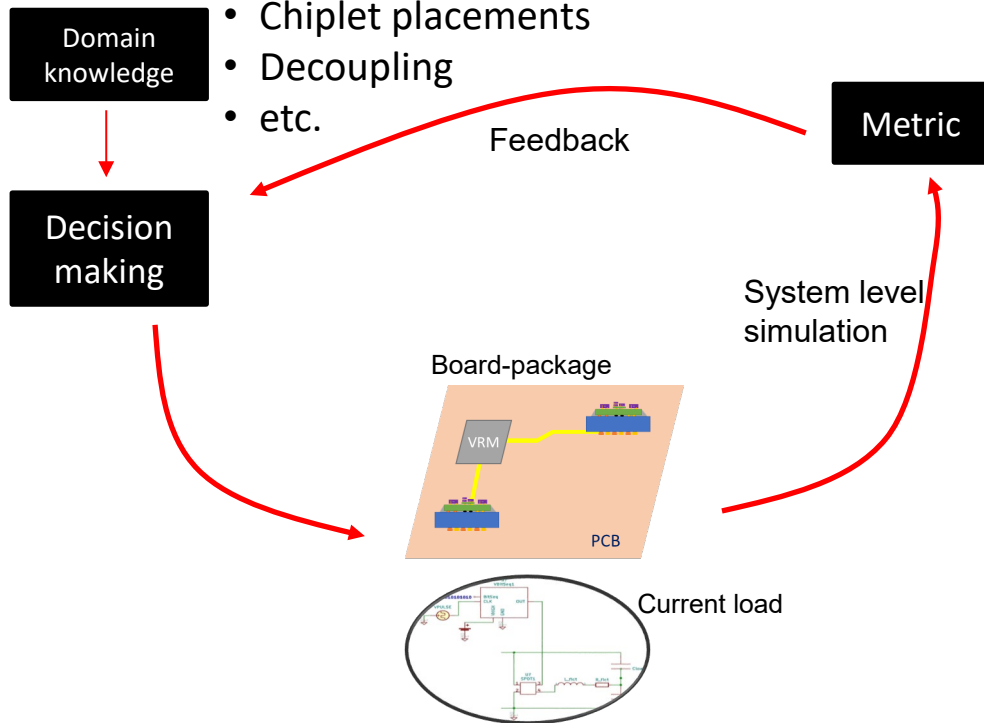
- Chiplets advantageous
- Boundary between chip and package is blurred
- Concurrent design is necessity
- SI/PI, multilevel frequencies
- System-level verification challenging



# Chiplet PDN

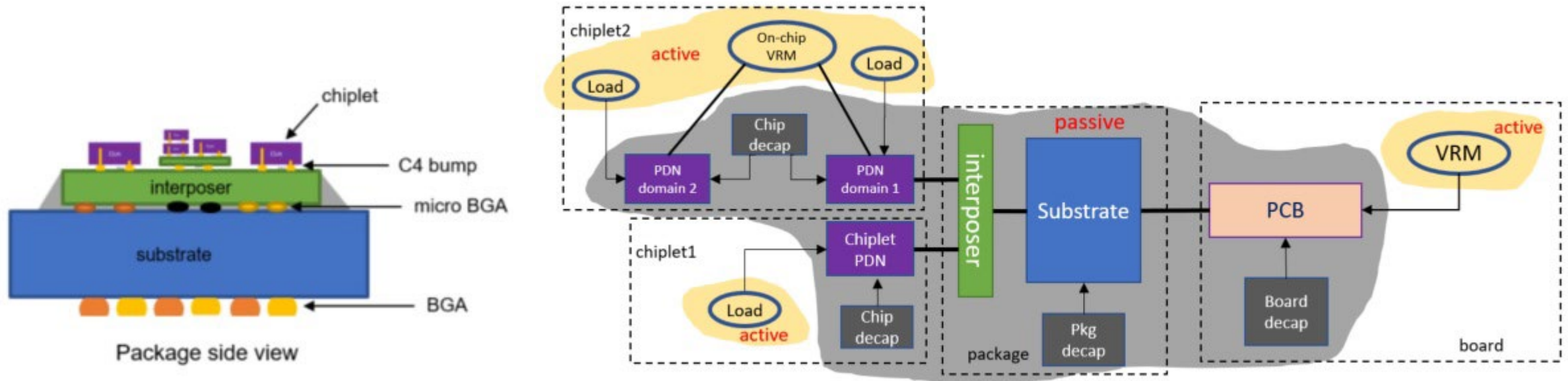
Update configuration:

- Package design (geometry, bumps placement etc.)
- Power domain allocation
- Chiplet placements
- Decoupling
- etc.



# Power Delivery for Chiplet-Based SiP

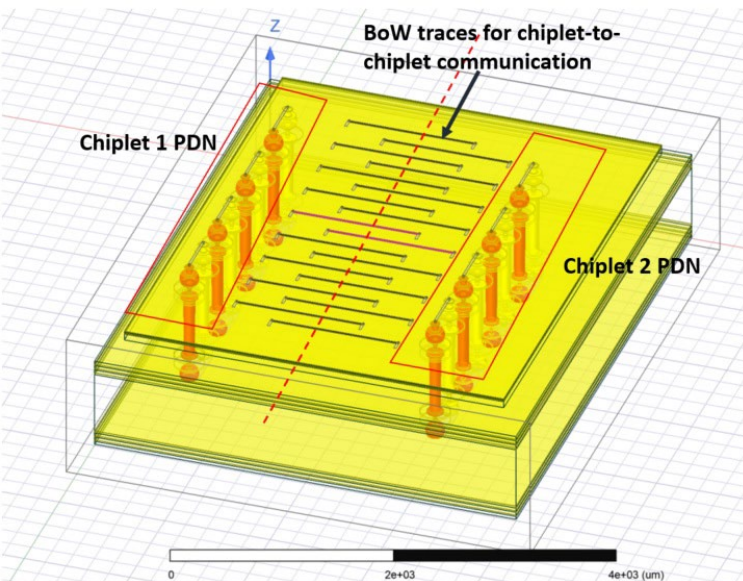
Chiplets must be positioned in their respective locations and provided with paths through the power distribution network. There can be multiple power domains. Package builder tools are used to automate the placement of vias, ball grid arrays and traces to the PDN. To mitigate supply voltage fluctuations decoupling capacitors (“decaps”) are placed on the PCB, package substrate, package interposer, and the silicon dies.



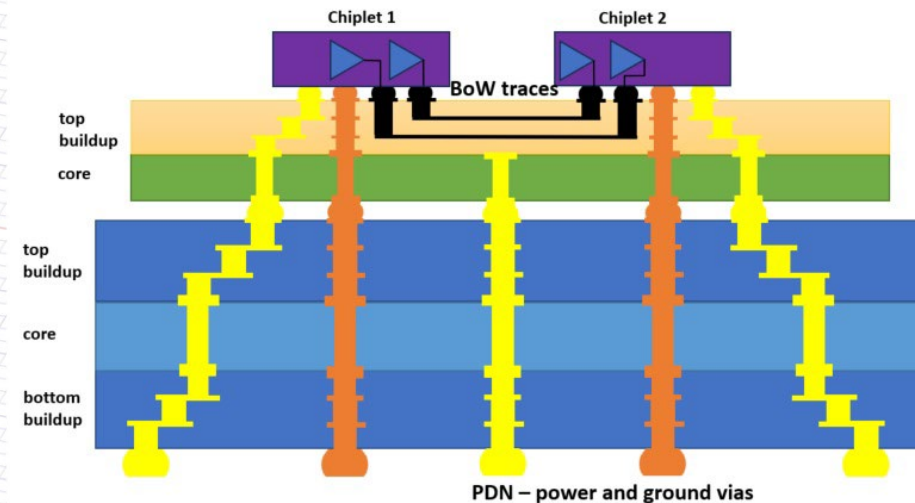
# Bump Mapping for Chiplets

Many problems take too long to simulate and thus make it difficult to predict the “what if” scenarios. This slows down the design process. This gets exacerbated when handling multiphysics situations (e.g. electrical + thermal) One approach is to first start with simplest possible models that will capture some physics and produce meaningful outputs

Top View



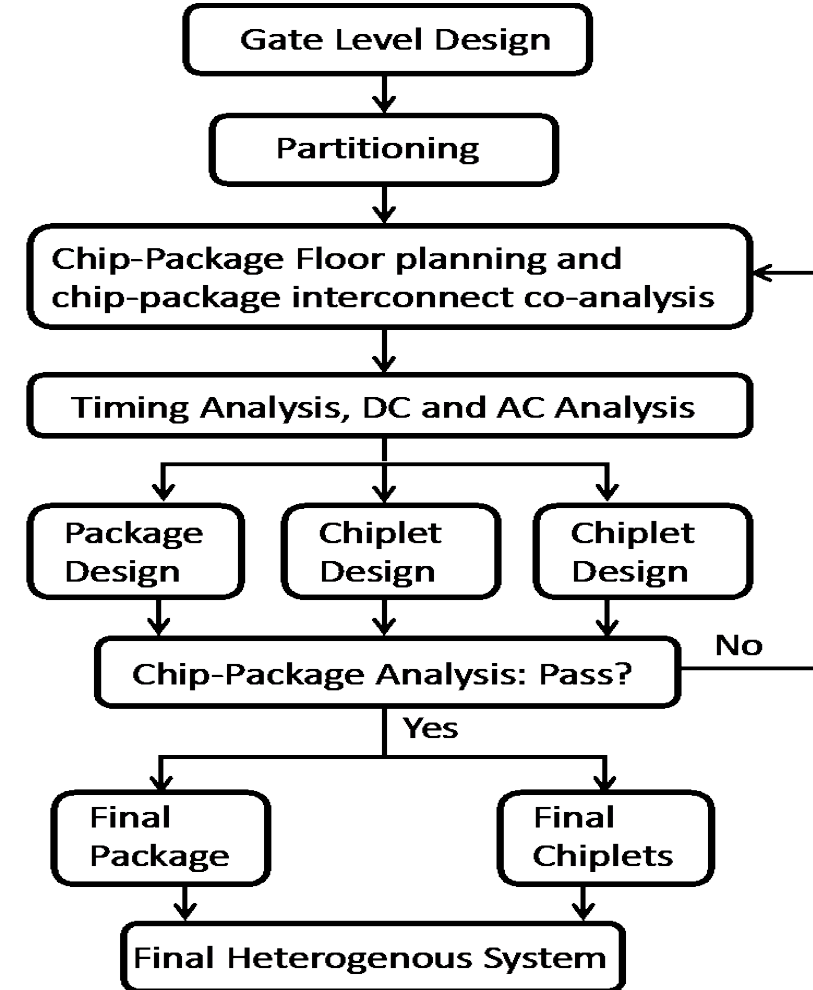
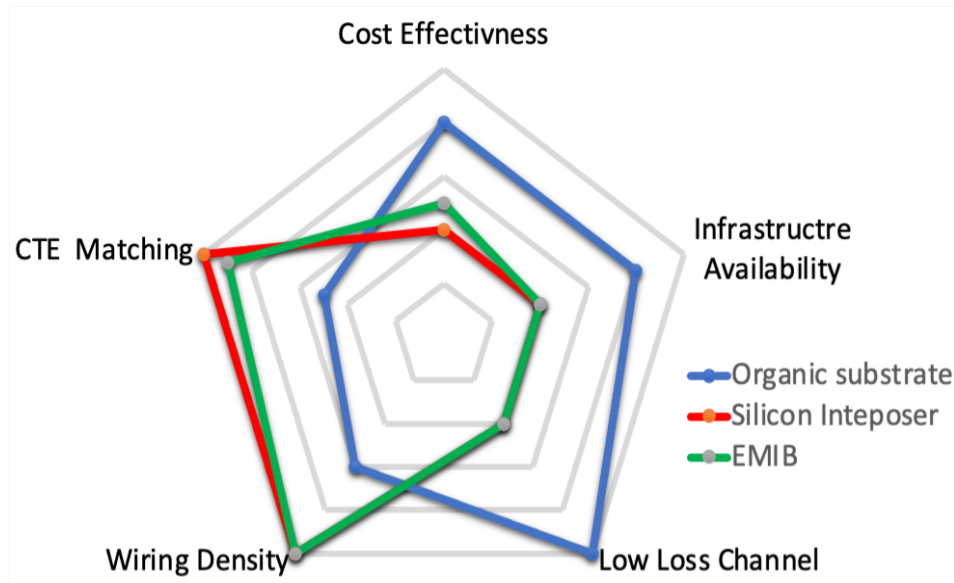
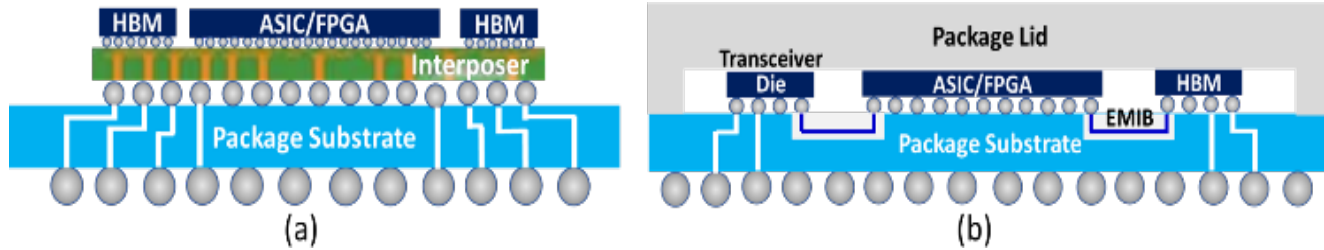
Cross-Section



**Chiplet-based design with PDN and D2D network**

**For 50 X 50 μBGA structure, extraction of network parameters takes more than a week!**

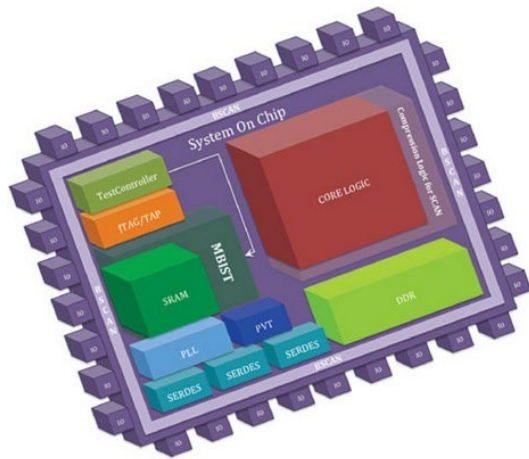
# Codesign for Chiplets



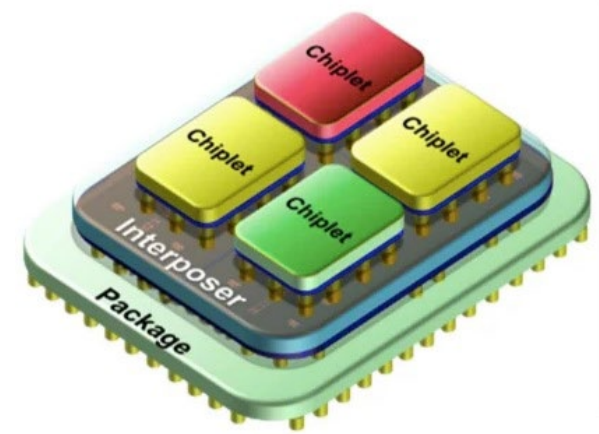
# The Problem of Disaggregation

Architecting an IC as a chiplet-based SiP rather than a SoC is referred to as disaggregation of function. Today, it is performed *ad hoc*; there is no established methodology to optimize the disaggregation, i.e., to determine how many separate chiplets should be used in order to meet specifications. Once chiplet design is democratized, there will be more choices from different vendors which could make the process more chaotic.

**SOC**

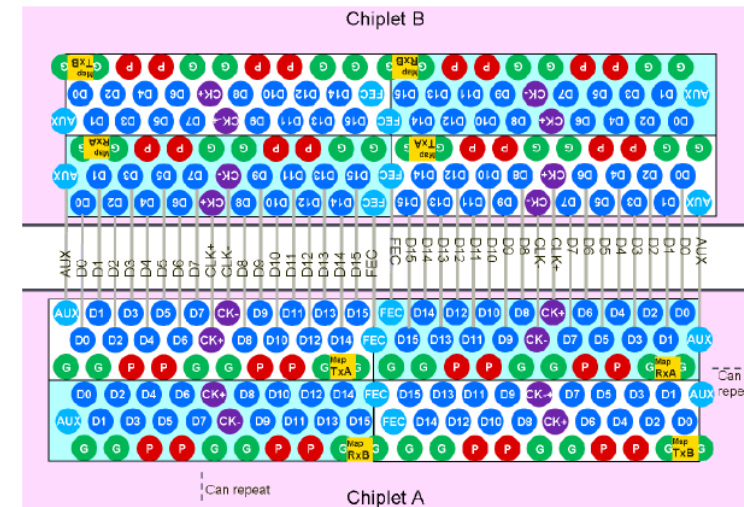
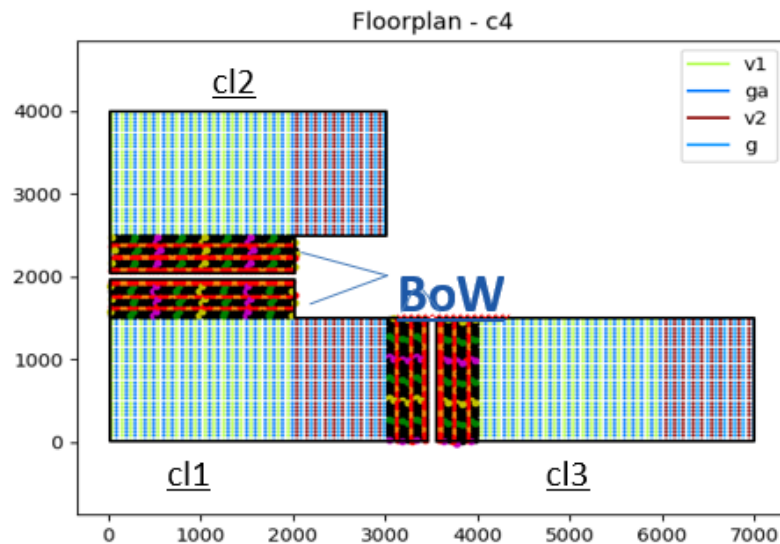


**SiP**



# Automating Disaggregation

- Implement co-design into EDA tools
- Assist with AI/ML
- Integrate interconnection standards (e.g. UciE, BoW)



[1] Shahab Ardalan, Ramin Farjadrad, Mark Kuemerle, Ken Poulton, Suresh Subramaniam, Bapiraju Vinnakota, "Chiplet Communication Link: Bunch of Wires (BoW)", IEEE Micro, Jan/Feb 2021



# Barriers to HI Co-Design

- **Technical Challenges**

- problem size and complexity → computational bottleneck
- design tools need more intelligent interfaces

- **Socio-Academic Challenges**

- - lack of information sharing within entities (data)
- - lack of common standardized formats
- - lack of access to low-volume manufacturing for validation
- - lack of access to design tools for exploration purposes
- - lack of multidisciplinary curriculum in advanced packaging

# Possible Solutions for Co-Design

- **Quantum Computing**
  - Several decades away
  - Solutions exist for specialized problems
- **Artificial Intelligence (AI) or Machine Learning (ML)**
  - Data must be available
  - Training can be very tedious
  - ➔ ***Need FASTER Electromagnetic Solvers (>10X)***

# Research Needs

- **Methodical abstractization**

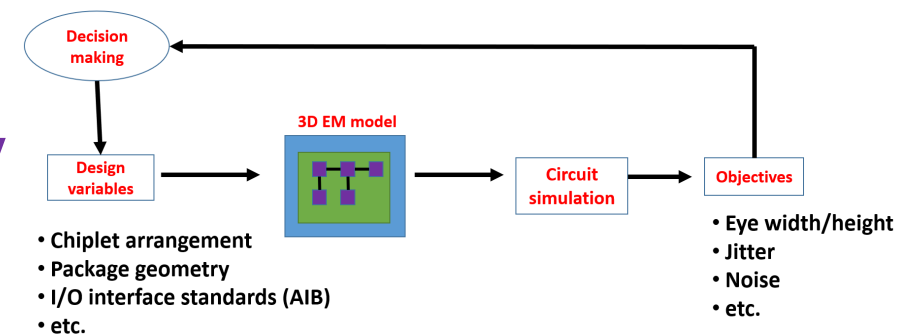
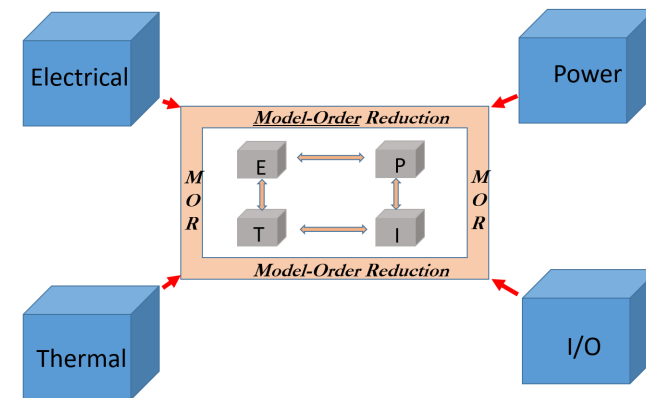
- Compact models
- Reduced-order models
- Behavioral models

- **Faster verification platforms (>10X)**

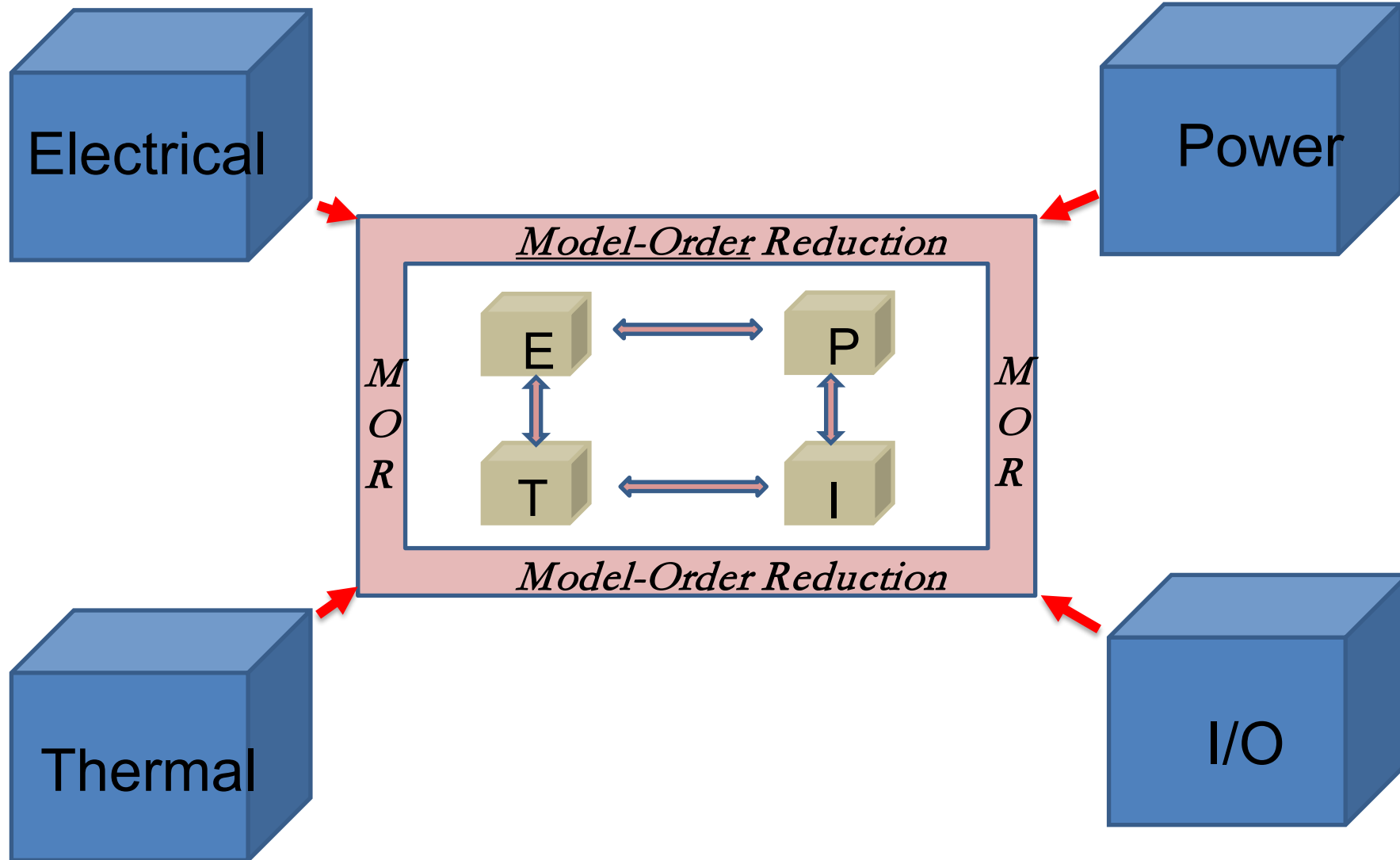
- Multi-physics solvers (EM , thermal management, materials)
- Transistor-level circuit simulation

- **AI/ML assisted solutions**

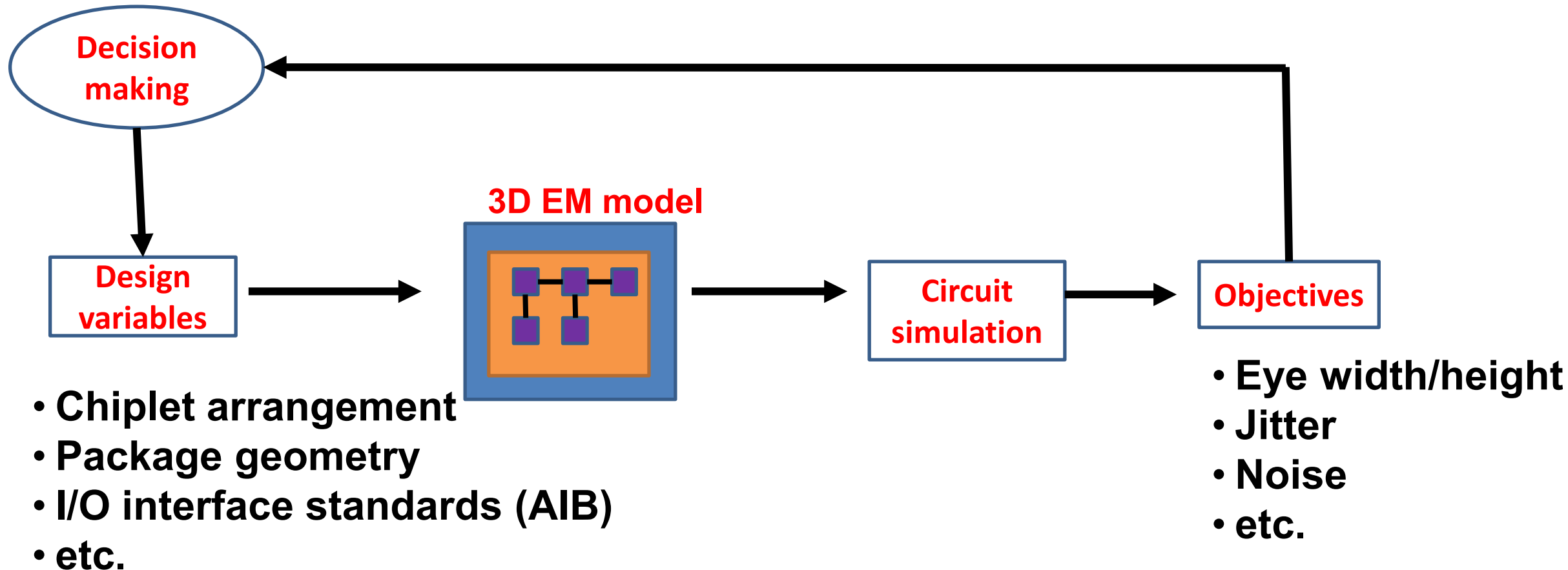
- Digital twins
- Optimization, mitigation of uncertainty



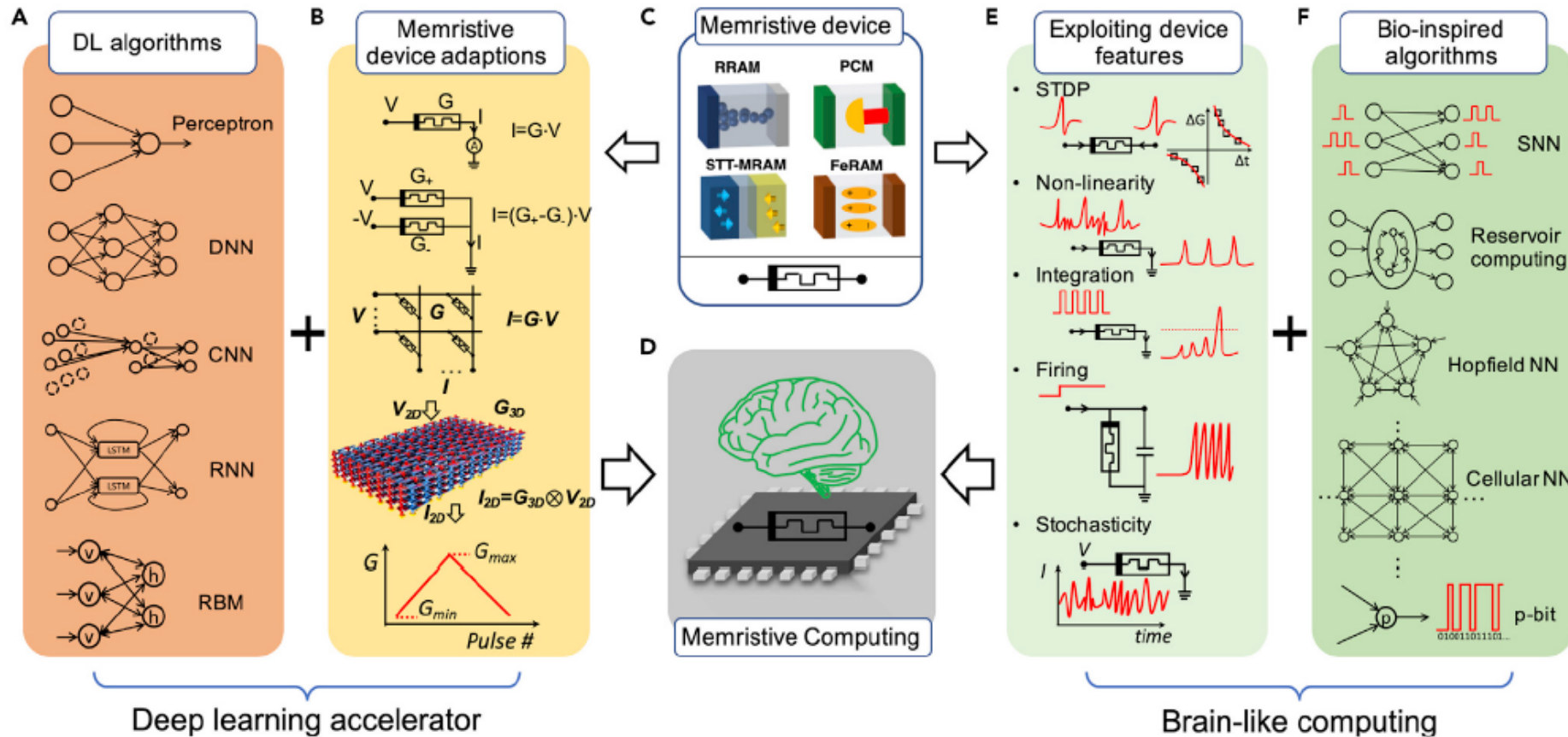
# Multiphysics & MOR



# AI/ML for System-Level



# Driver: Brain-Inspired AI Chipllets



SOURCE: Wei Wang, Wenhao Song, Peng Yao, Yang Li, Joseph Van Nostrand, Qinru Qiu, Daniele Ielmini and J. Joshua Yang, "Integration and Co-design of Memristive Devices and Algorithms for Artificial Intelligence", *iScience* 23, 101809, December 18, 2020

# Neurobiology-Driven Algorithms

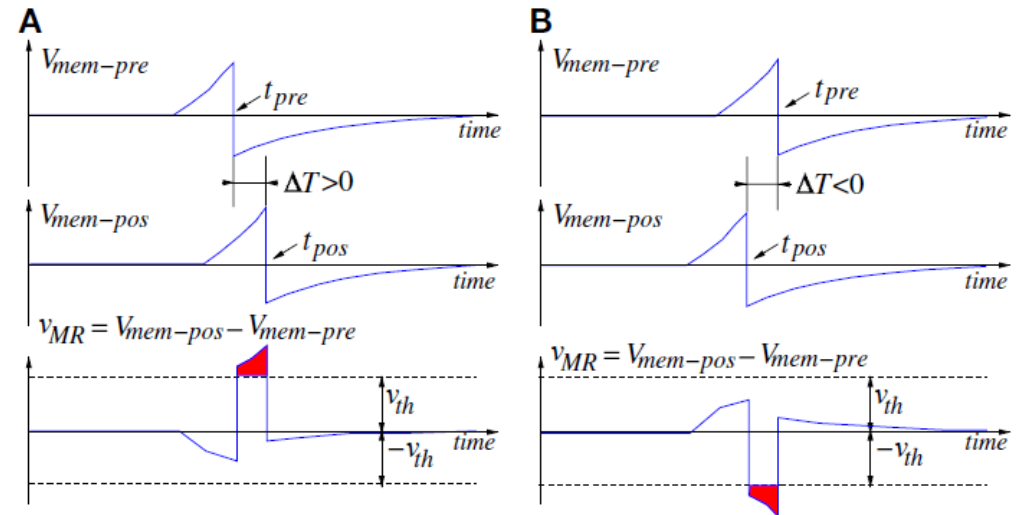
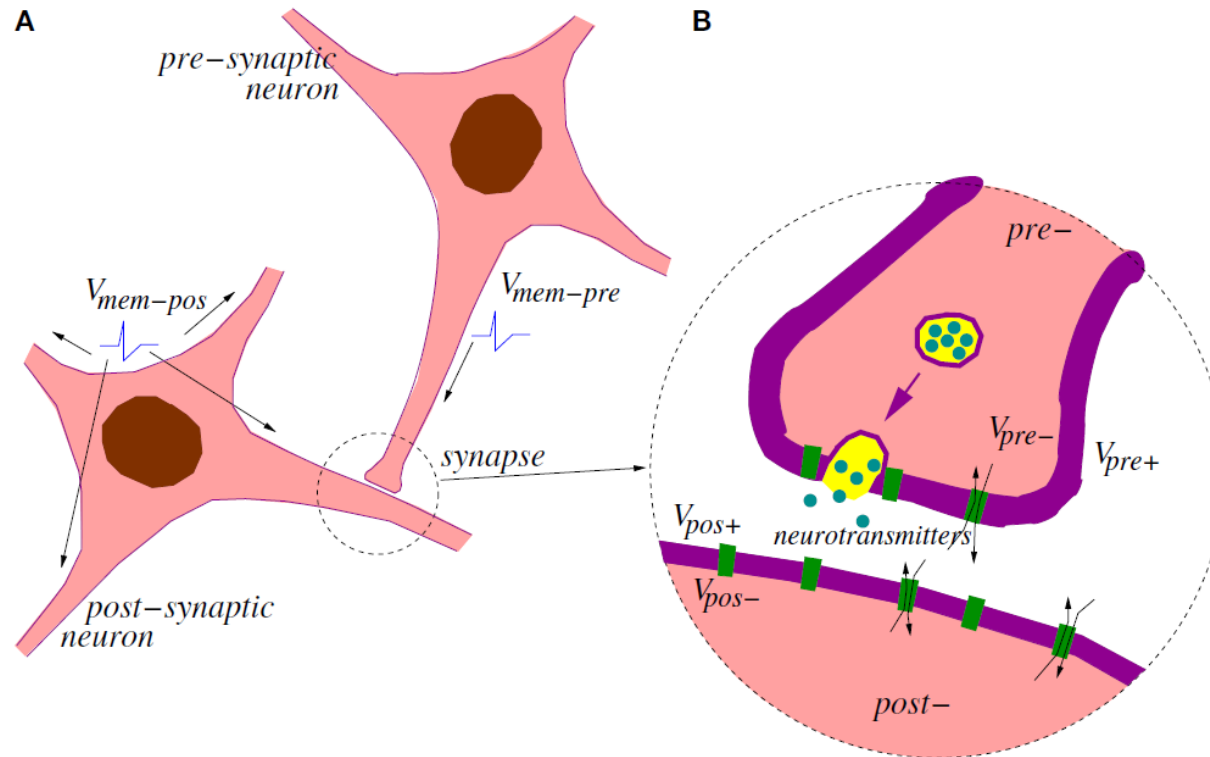
- **Similarities with biological synapses**
  - Stochastic
  - pulse-pair facilitation
  - short-term plasticity
- **Approach**
  - Modify devices and materials to emulate algorithms
  - Modify algorithms to exploit special properties of memristive devices

# STDP-Based Memristor

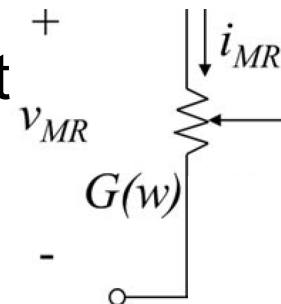


## On spike-timing-dependent-plasticity, memristive devices, and building a self-learning visual cortex

Carlos Zamarreño-Ramos<sup>1</sup>, Luis A. Camuñas-Mesa<sup>1</sup>, Jose A. Pérez-Carrasco<sup>1</sup>, Timothée Masquelier<sup>2</sup>, Teresa Serrano-Gotarredona<sup>1</sup> and Bernabé Linares-Barranco<sup>1\*</sup>



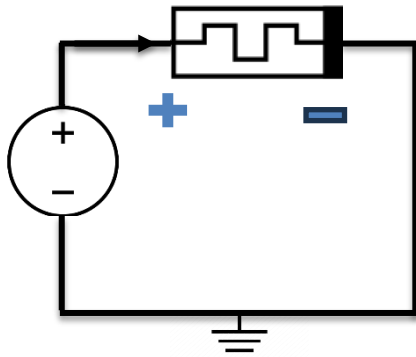
Equivalent Circuit



$$\Delta G(v_{MR}) = \begin{cases} G_0 \left[ e^{v_{MR}/v_o^+} - e^{v_{th}^+/v_o^+} \right], & v_{MR} > v_{th}^+ \\ -G_0 \left[ e^{-v_{MR}/v_o^-} - e^{-v_{th}^-/v_o^-} \right], & v_{MR} < v_{th}^- \\ 0, & \text{otherwise} \end{cases}$$

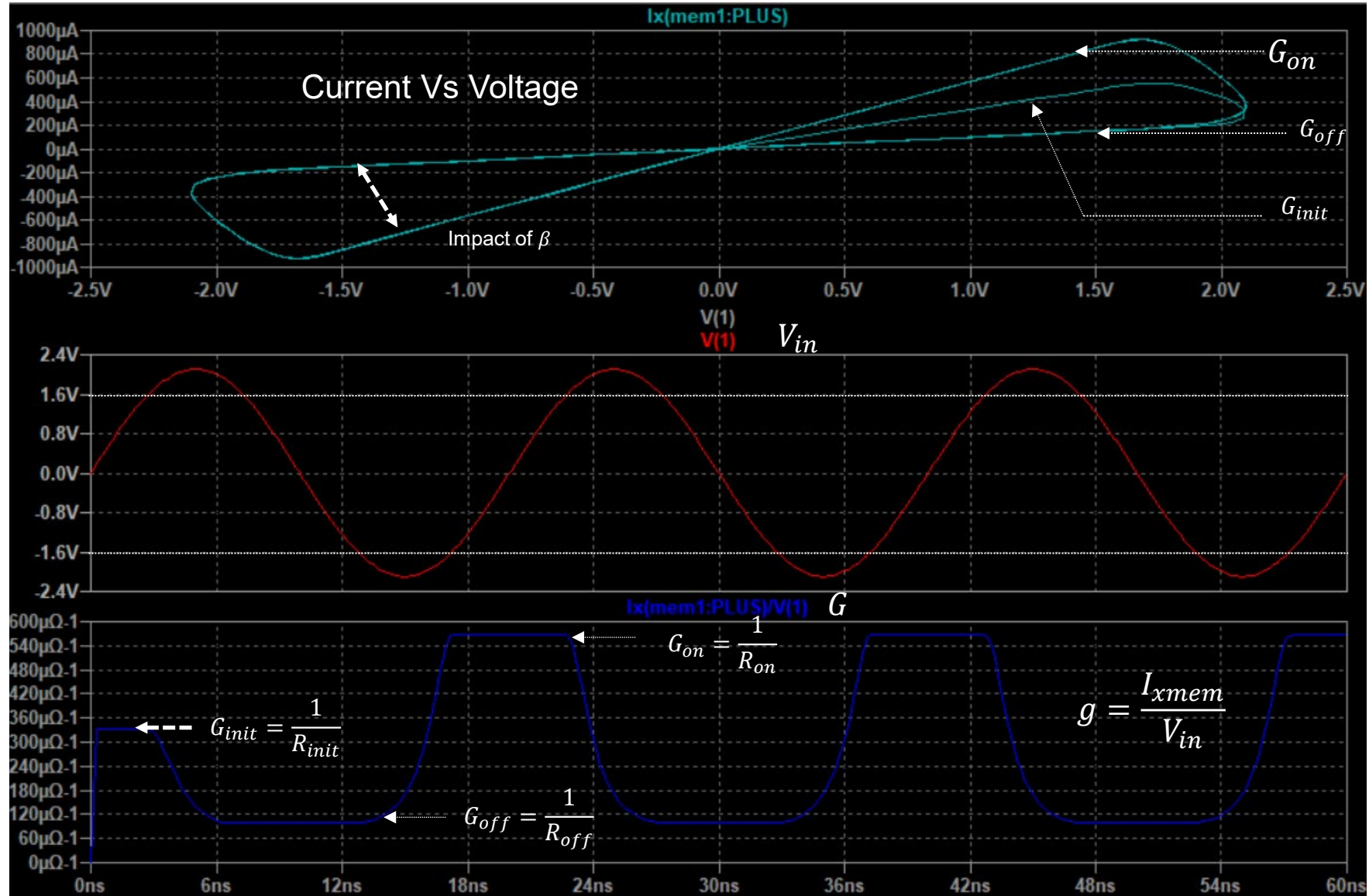
# Memristor Model: Bipolar memristor with threshold

Reliable SPICE Simulations of Memristors, Memcapacitors and Meminductors, 2013



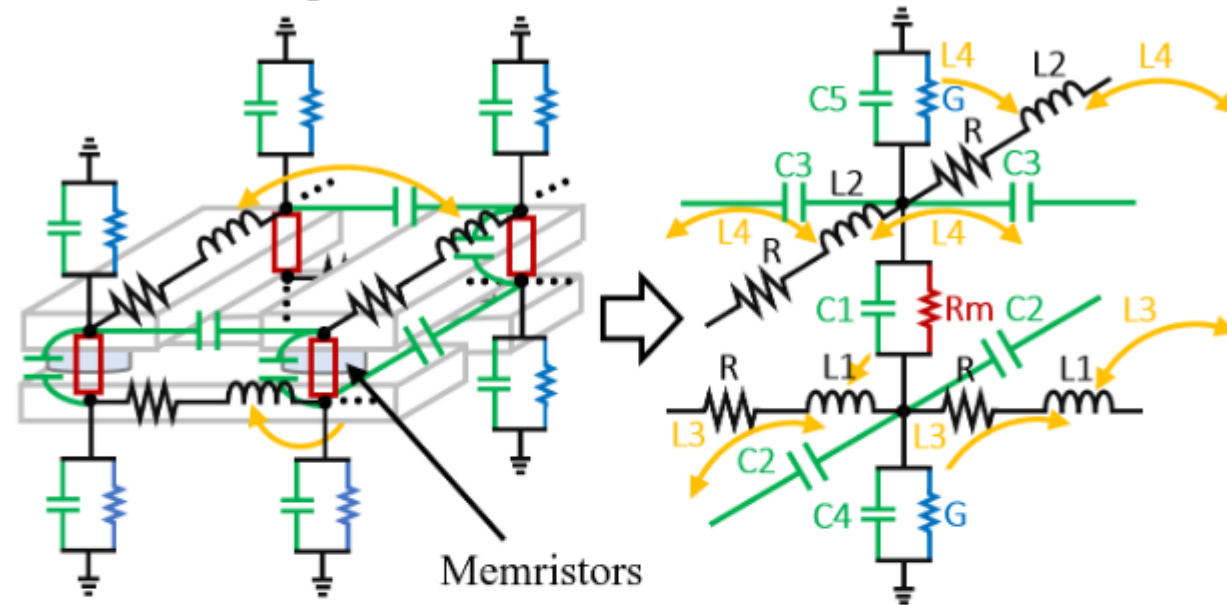
- Behavioral Model is used
- Old and New Parameters:

- $V_{th} = 4.6V \rightarrow 1.6V$
- $R_{on} = 1k$
- $R_{off} = 10k$
- $R_{init} = 5k \rightarrow 3K$
- $\beta = gain = 5e12$

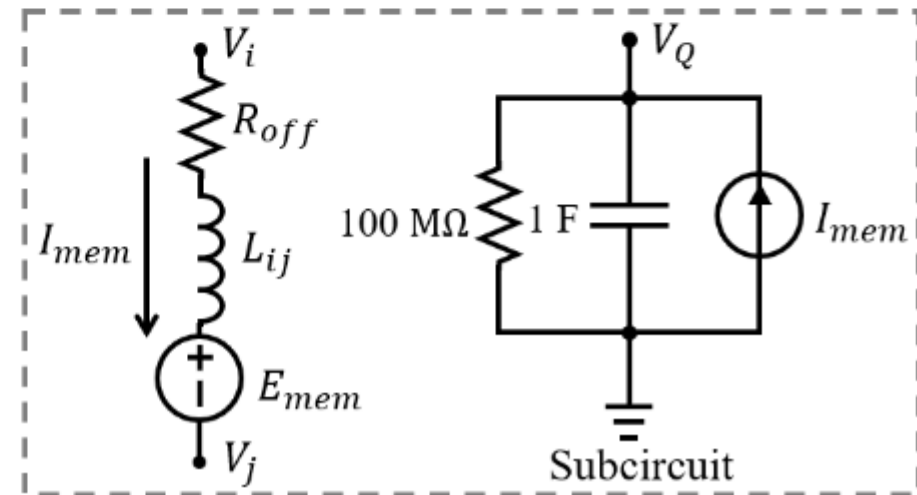
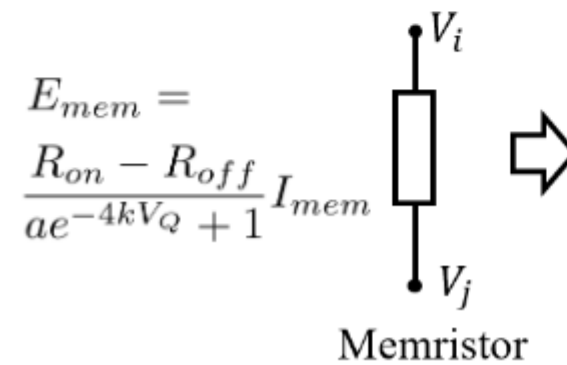


# Memristor

## Crossbar array:

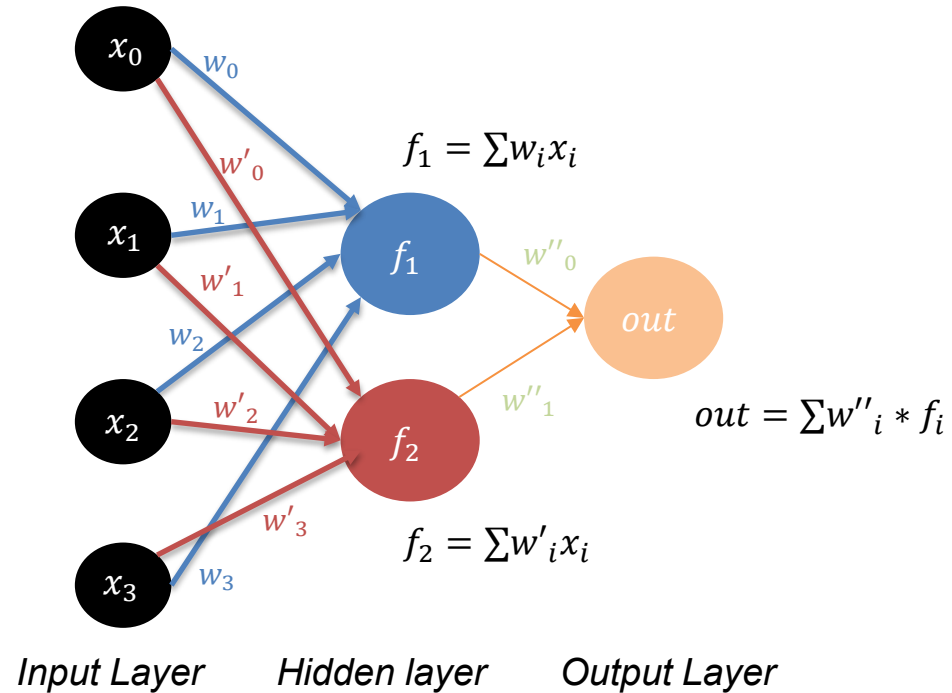


## Memristors:

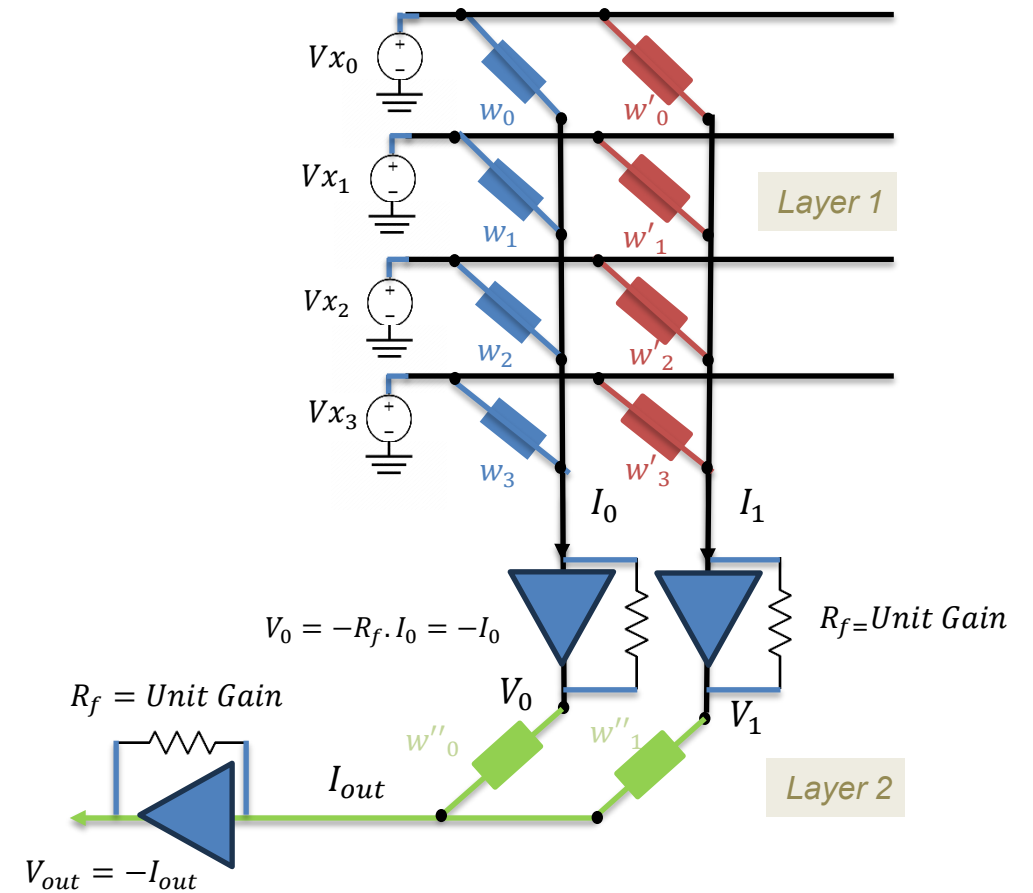


# Fully Connected Layer for NN

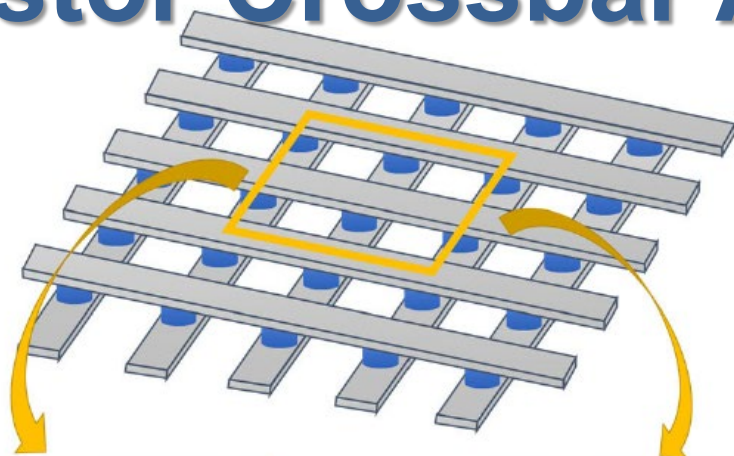
Linear Neural Network  
(Without any activation Function)



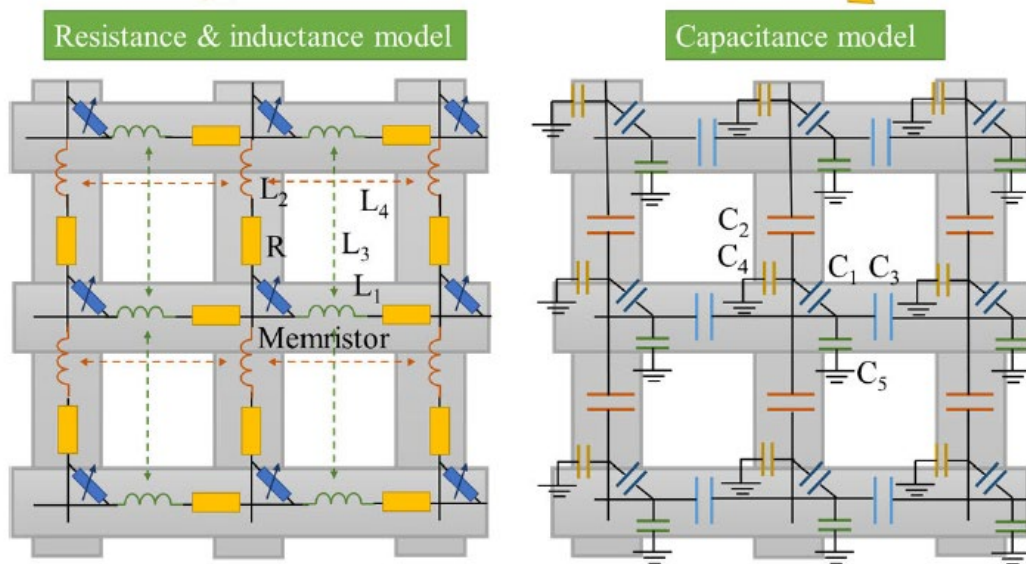
Conductance Mapping of NN in memristor array



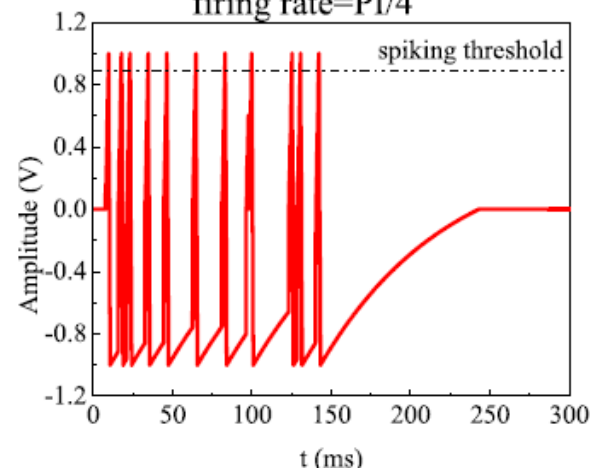
# Memristor Crossbar Array



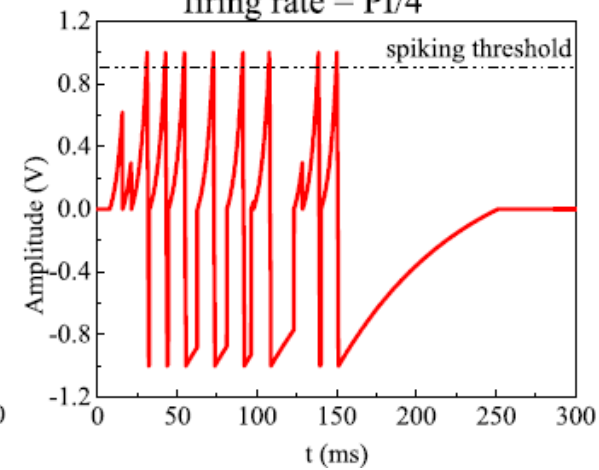
Will perform VMMA in analog domain



Input spikes when  $\tau^+ = 5$  ms,  $\text{tail}^+ = 2$  ms, firing rate =  $\pi/4$



Input spikes when  $\tau^+ = 5$  ms,  $\text{tail}^+ = 10$  ms, firing rate =  $\pi/4$



# Conclusion

- **Chipelets: HI solution for future SiP**
- **Can only be addressed through co-design**
- **Ecosystem still in nascent stage**
- **Size, complexity and stochasticity are increasing**
- **Brain-inspired computing will be a driver**
- **Great opportunities for MOR**