

Integrated Power Delivery for AI Computing Technology Gaps & Opportunities

Madhavan Swaminathan

Dept. Head Electrical Engineering

William E. Leonhard Endowed Chair

Director, CHIMES (an SRC JUMP 2.0 Center)

The Pennsylvania State University

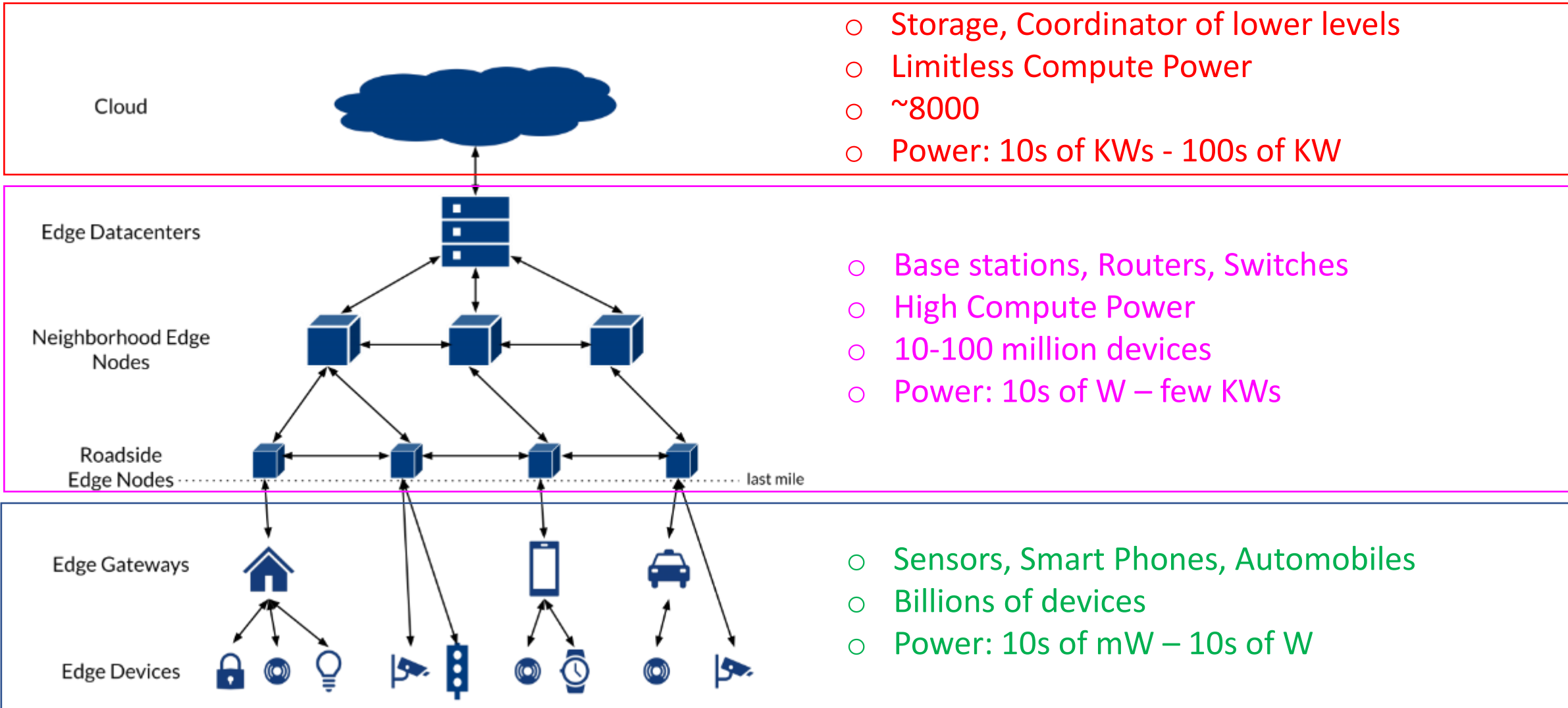
Emeritus Professor, ECE & MSE, Georgia Tech

Former Director, 3D Systems Packaging Research Center (NSF-ERC), Georgia Tech

Outline

- Deep Learning
- Challenges for Power Delivery
- Integrating Power Sources
- 3D Integration
- Center for Heterogeneous Integration of Micro Electronic Systems (CHIMES)
- Summary

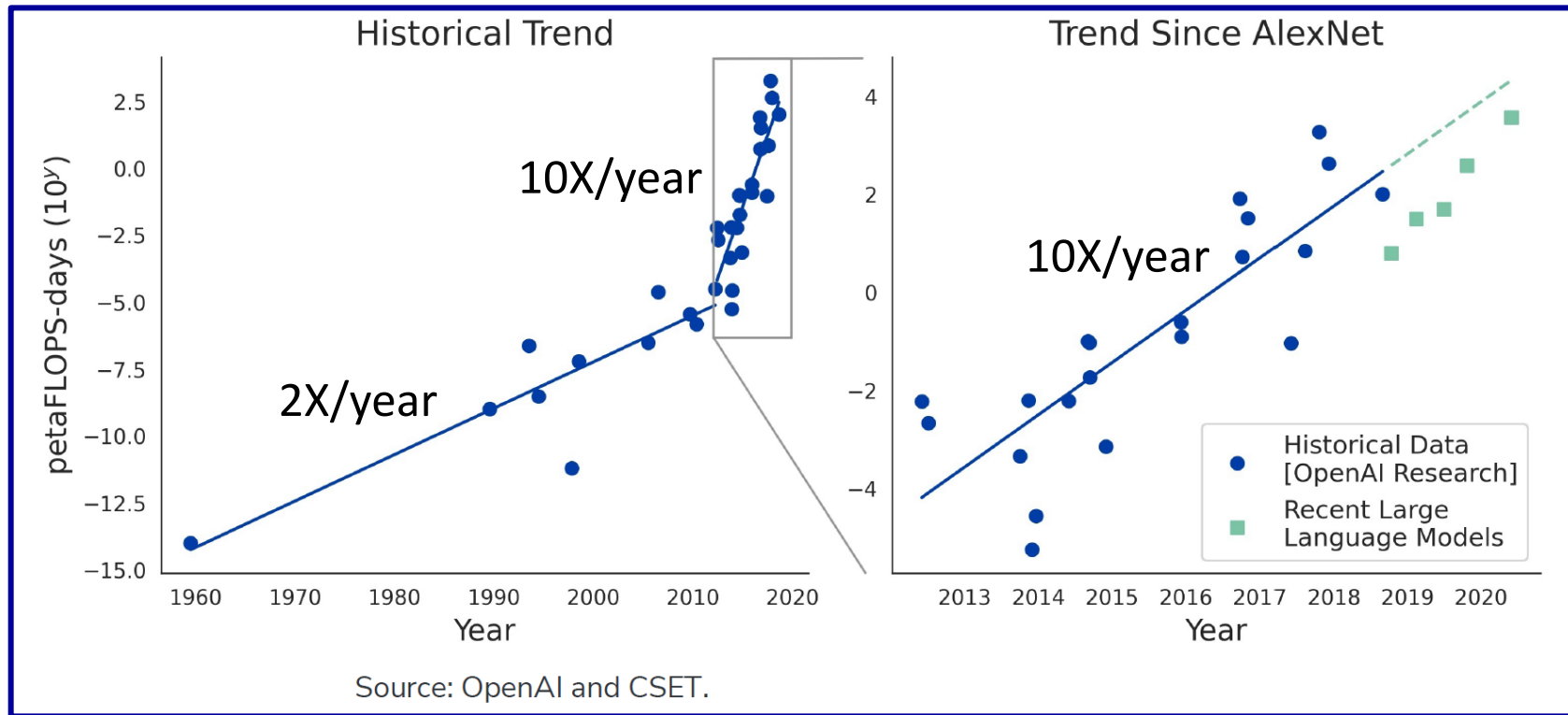
Emerging Distributed & Edge Computing



Ref: Jeff Burns, "Systems and Architectures for Distributed Compute", SRC Workshop, 2022.

<https://www.inovex.de/de/blog/edge-computing-introduction/>

Deep Learning & Computational Power (Large Data Centers)

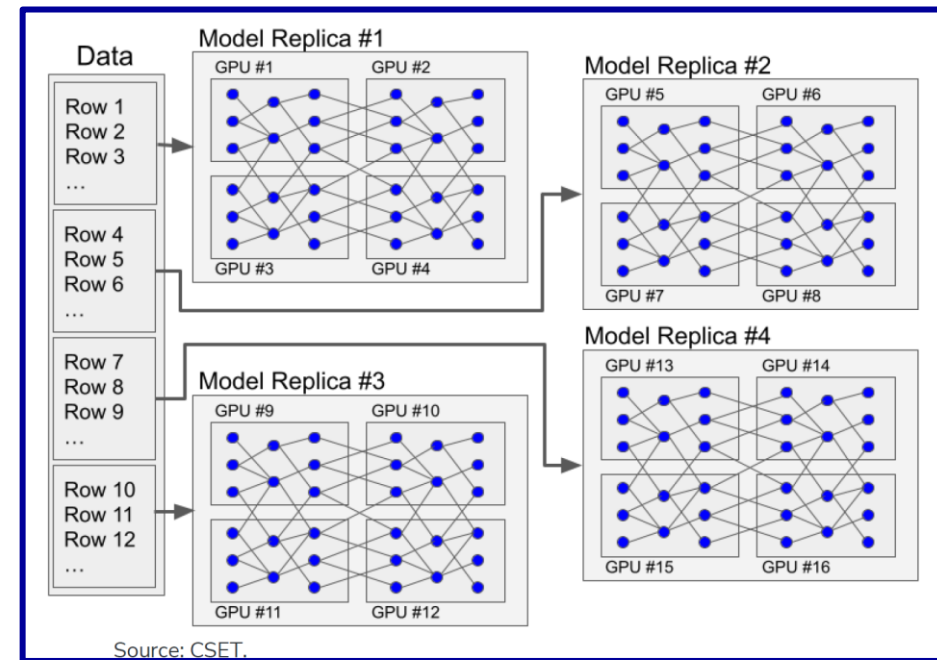
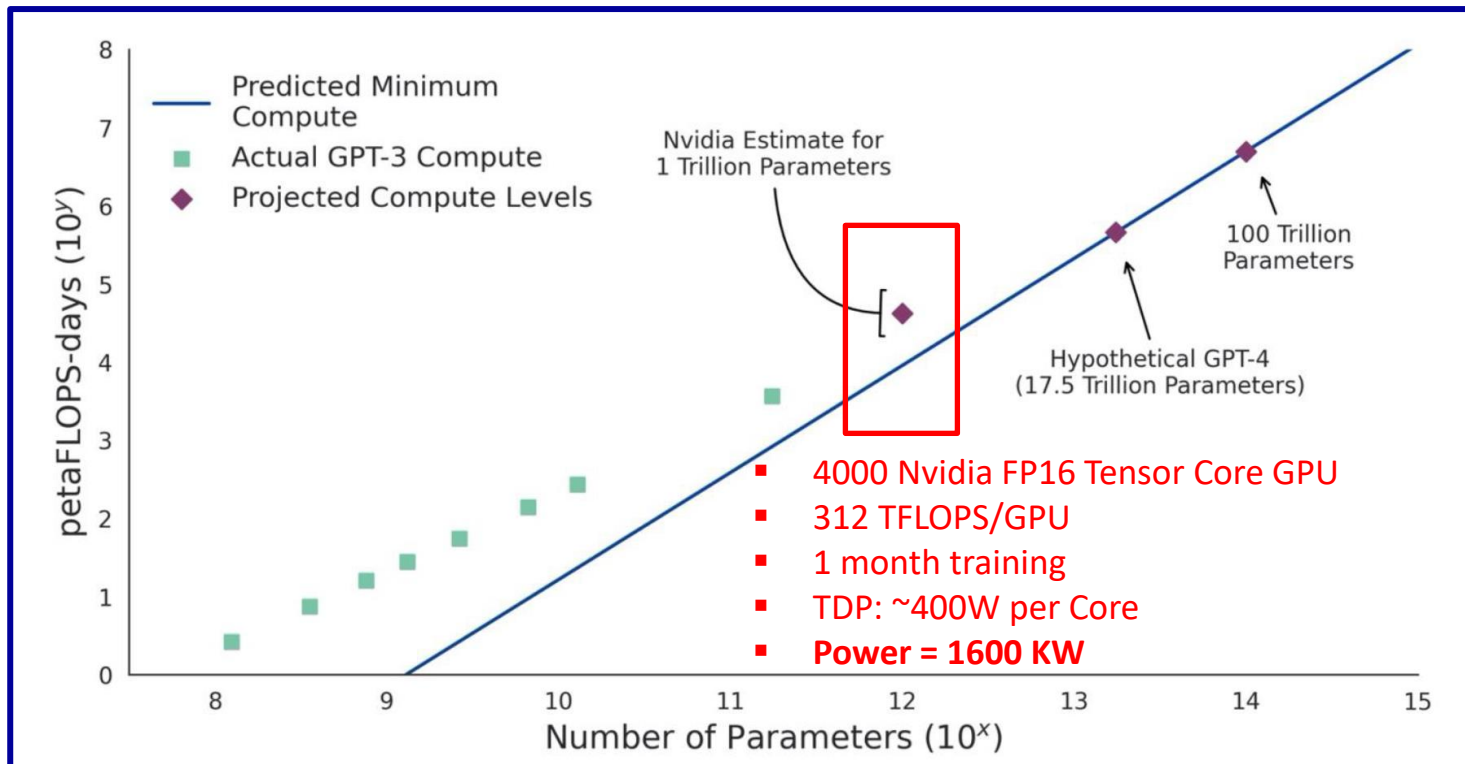


Observations

- Prior to 2012: Compute demand doubling every 2 years at the same rate as Moore's law
- Since 2012: Compute demand growing 10X/year due to Deep Learning

Ref: Andrew John and Micah Musser, "AI and Compute – How much longer can computing power drive artificial intelligence progress", CSET, 2022

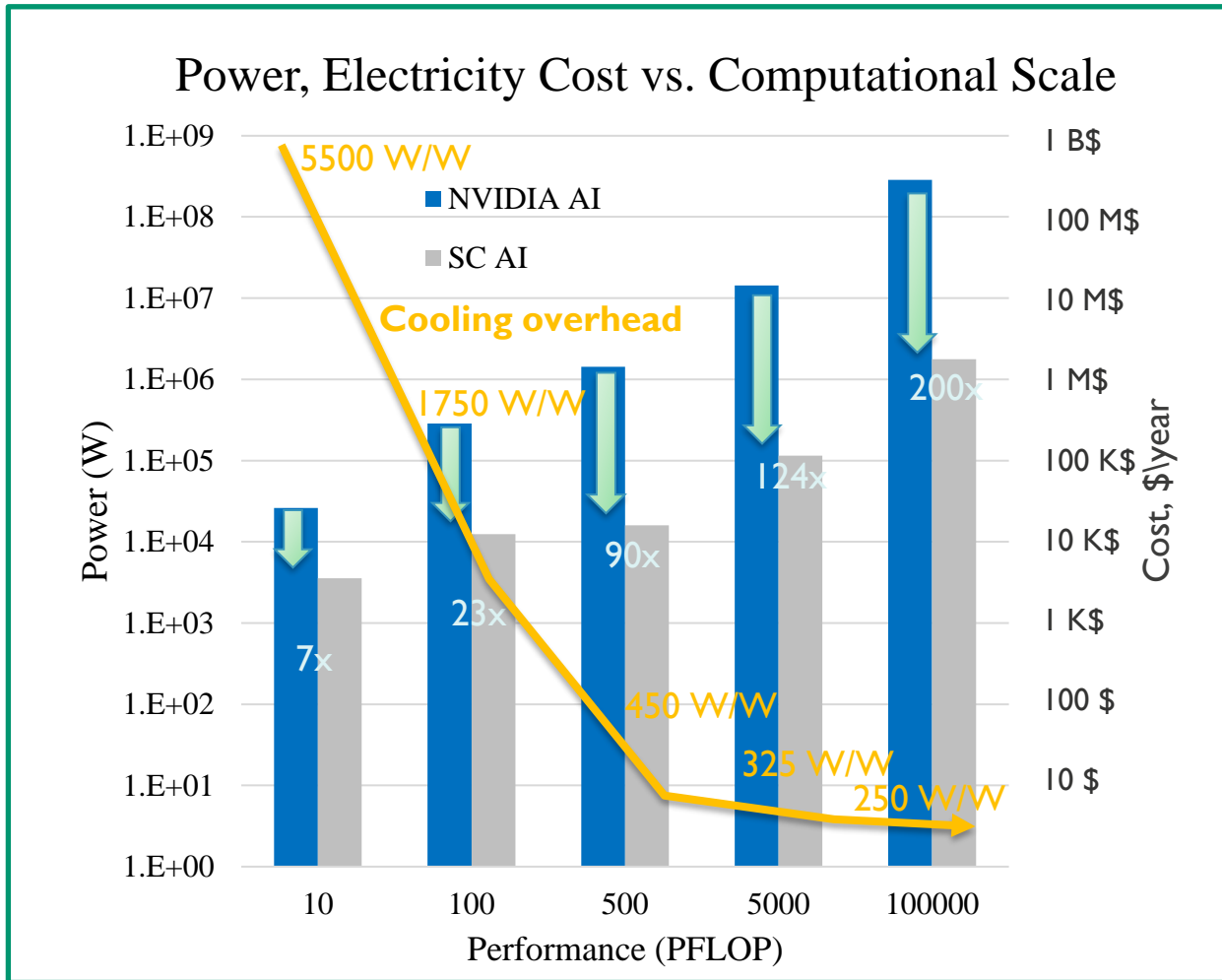
Compute Power & Scaling



- 2X Process Gains (8nm)
- 3X Architecture Gains (Tensor Flow)
- 300X System Scale out (Data Centers) – OUCH!
- Energy crisis if we continue this trend.....

Ref: Andrew John and Micah Musser, "AI and Compute – How much longer can computing power drive artificial intelligence progress", CSET, 2022
 Anna Herr and Quentin Herr, Superconducting AI & HPC, IMEC

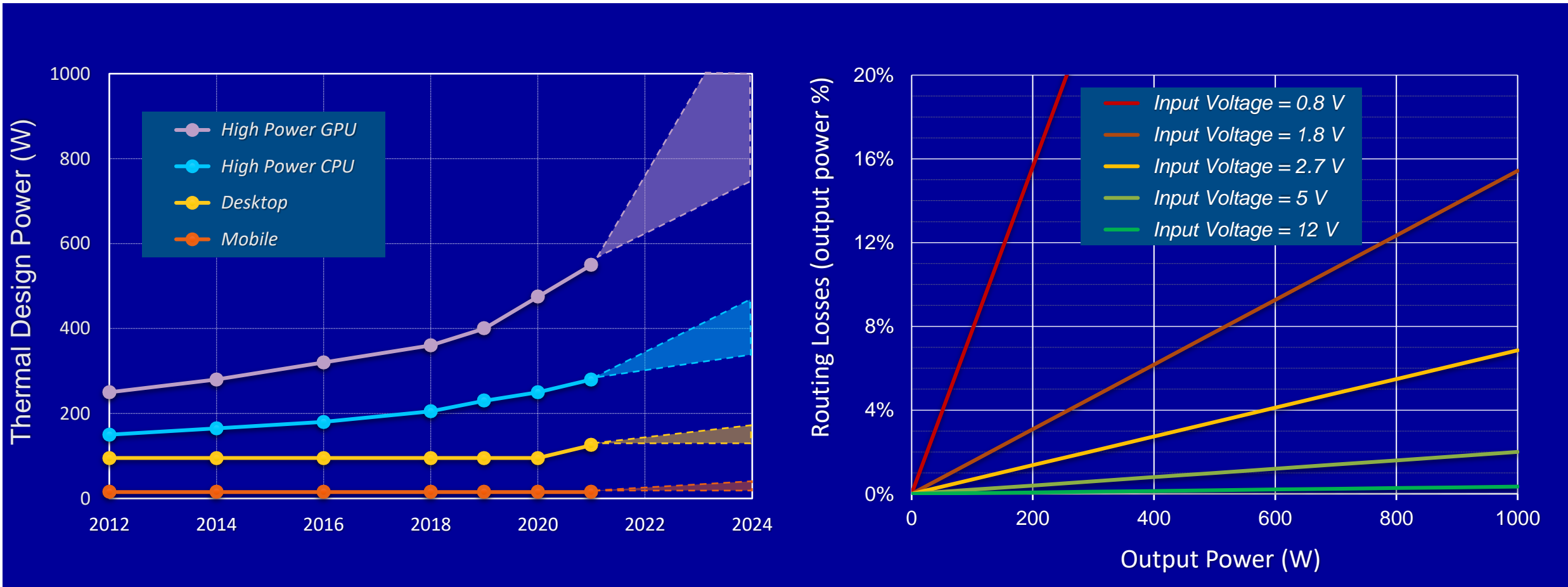
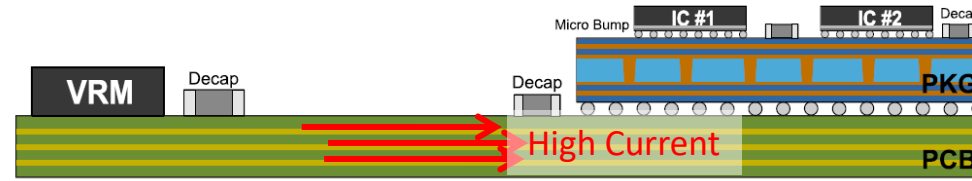
Achieving Energy & Cooling Efficiency at Scale



- Superconducting electronics breaks even at PFLOP scale
- Rapid increase in power efficiency with scale
- **100 M\$/year** savings in electricity

Courtesy: Anna Herr, IMEC

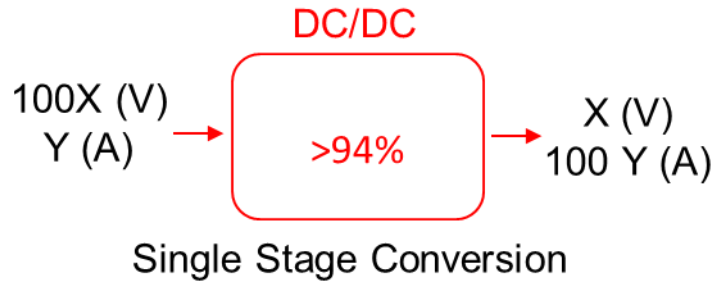
Let's start from the Fundamentals of Power Delivery



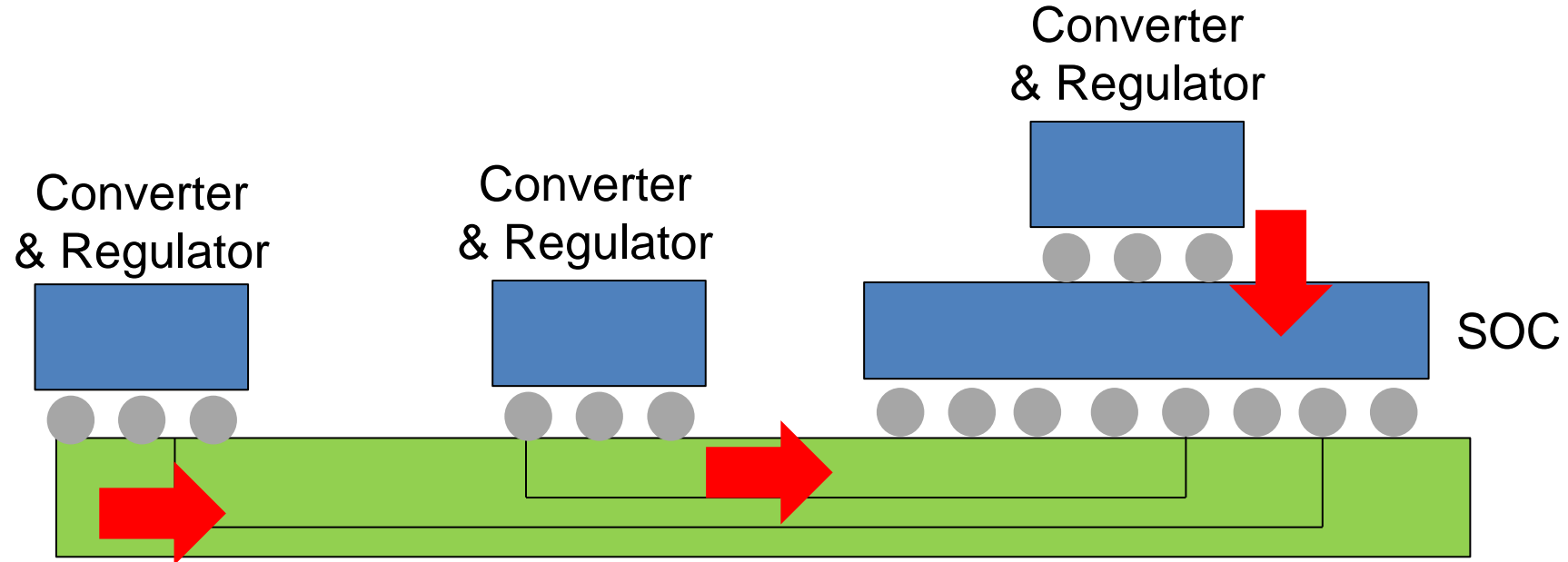
K. Radhakrishnan, EDAPS Keynote, 2021

K. Radhakrishnan, M. Swaminathan & B. Bhattacharyya, TCPMT, 2021

Objectives for Power Delivery

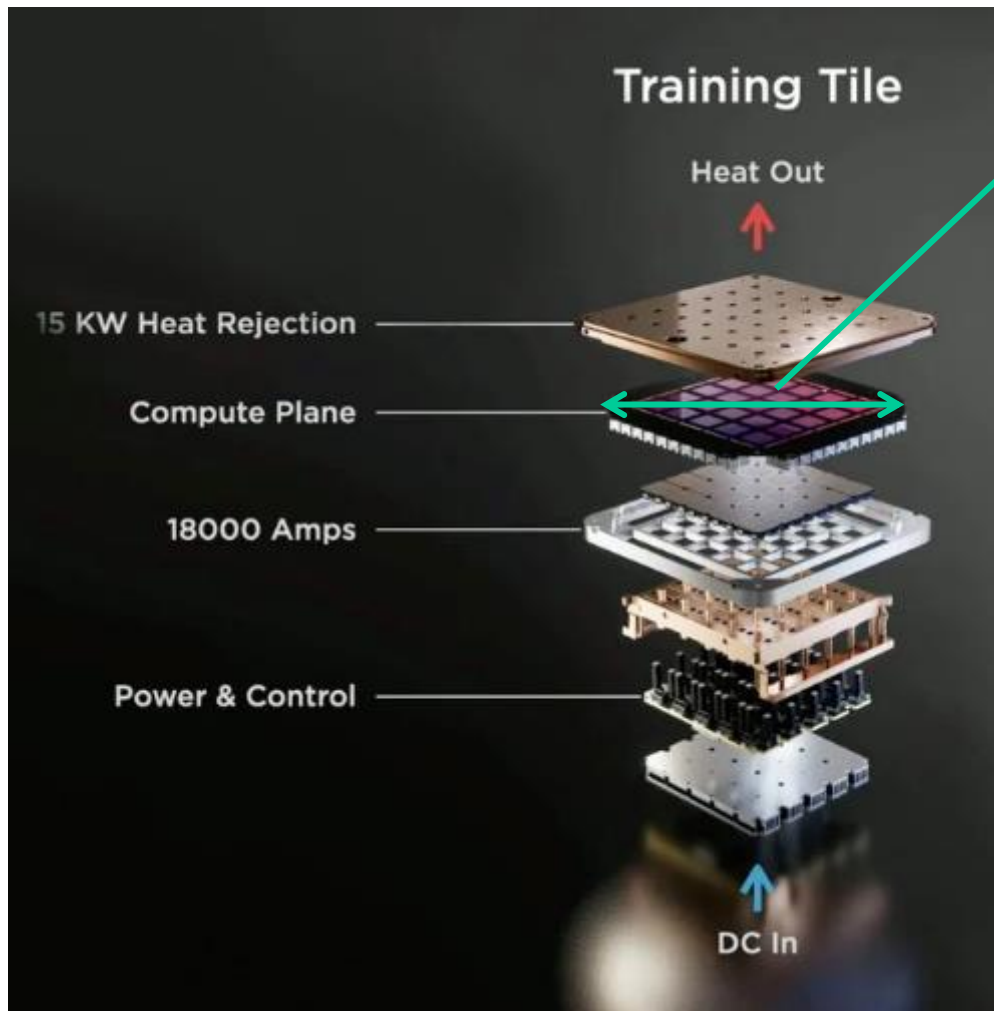


$$\eta = \frac{P_{out}}{P_{in}}$$



- ❑ Bring the converter (power source) near SOC
 - Significantly reduce Cu losses due to shorter current paths
- ❑ Integrate High Efficiency, Highly Integrated, Highly Miniaturized, High Conversion Ratio, Single Stage Converters
- ❑ Maximize power efficiency

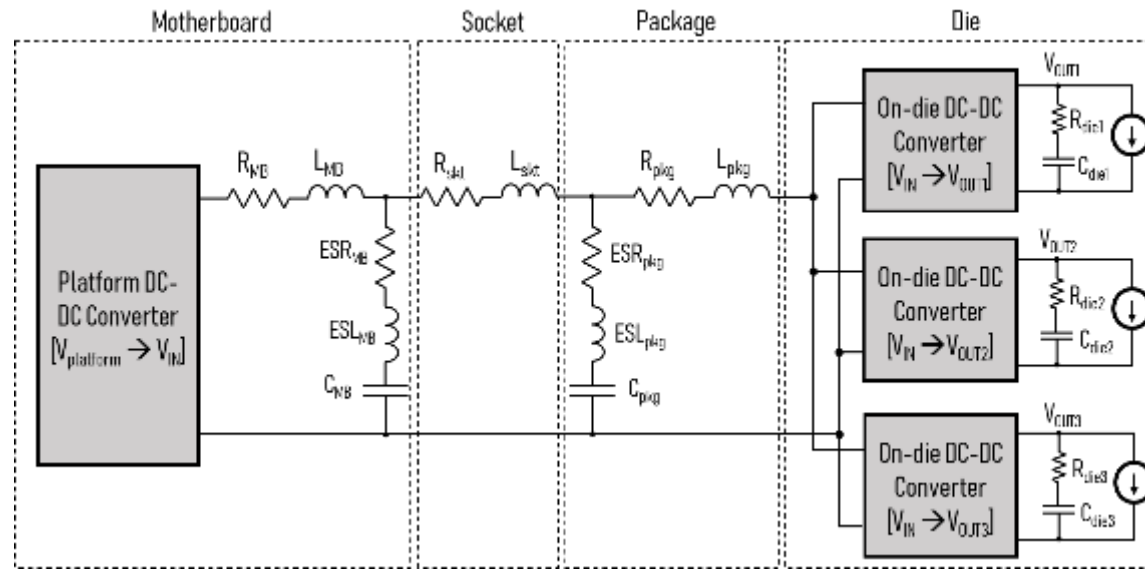
Tesla Dojo Power Delivery & Thermal - SOTA



- Signaling
 - Horizontal to maximize BW
 - 312 TFLOPS/die BF16
 - 22.6 TFLOPS/die FP32
 - TDP: 400W/die
- Power Delivery
 - Vertical to minimize losses
 - Power consumed: 10KW
 - Power Input: 15KW
 - Efficiency: ~67%
- Thermal
 - 15KW Heat Rejection
 - Liquid Cold Plate

<https://semianalysis.substack.com/p/tesla-dojo-unique-packaging-and-chip?s=r>

Need for Integrated Voltage Regulators (IVR)

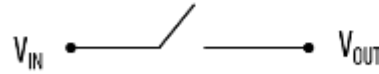


- ❑ Proliferation of on-die power domains requiring fine grain power management
 - Use small number of robust platform level voltage regulators to provide input power to the IVR
- ❑ With increasing power levels, routing losses in the Power Delivery Network (PDN) can have significant impact on the overall power system efficiency
 - Bring power at higher voltage to the processor

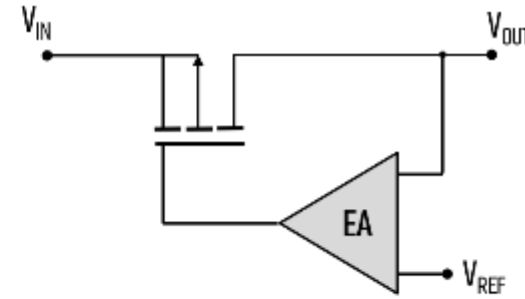
K. Radhakrishnan, M. Swaminathan & B. Bhattacharyya, TCPMT, 2021

Types of IVR

- Turn ON/OFF Power Domains
- Does not regulate voltage



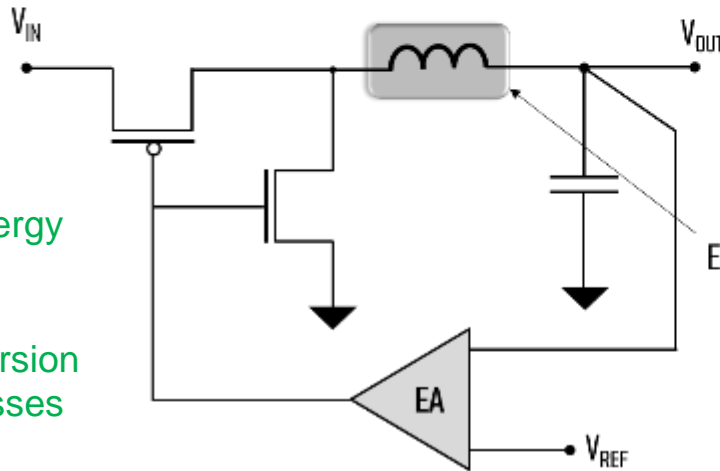
Simple Power Gate



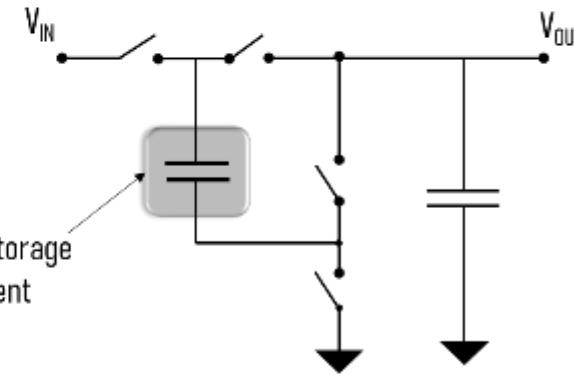
- Regulates Voltage
- No energy storage element
- Easy to integrate on-die
- Input close to output voltage
- Does not address routing losses

Low Drop-out (LDO) regulator

- Regulates Voltage
- Requires inductor energy storage element
- Higher input voltage
- High efficiency conversion
- Addresses routing losses



Switching Buck Regulator



- Regulates Voltage
- Requires capacitor energy storage element
- Higher input voltage
- Fixed ratio conversion from input to output
- Addresses routing losses

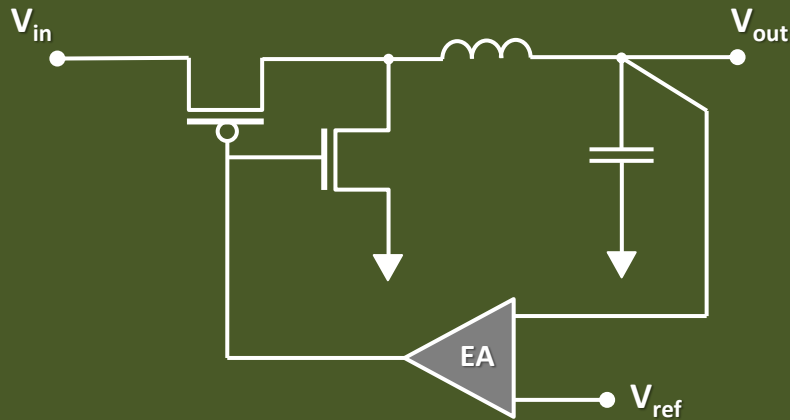
Switched capacitor voltage regulator

□ Hybrid: single-inductor-multiple-output (SIMO) regulators augmented with linear voltage regulators for transient management generally used.

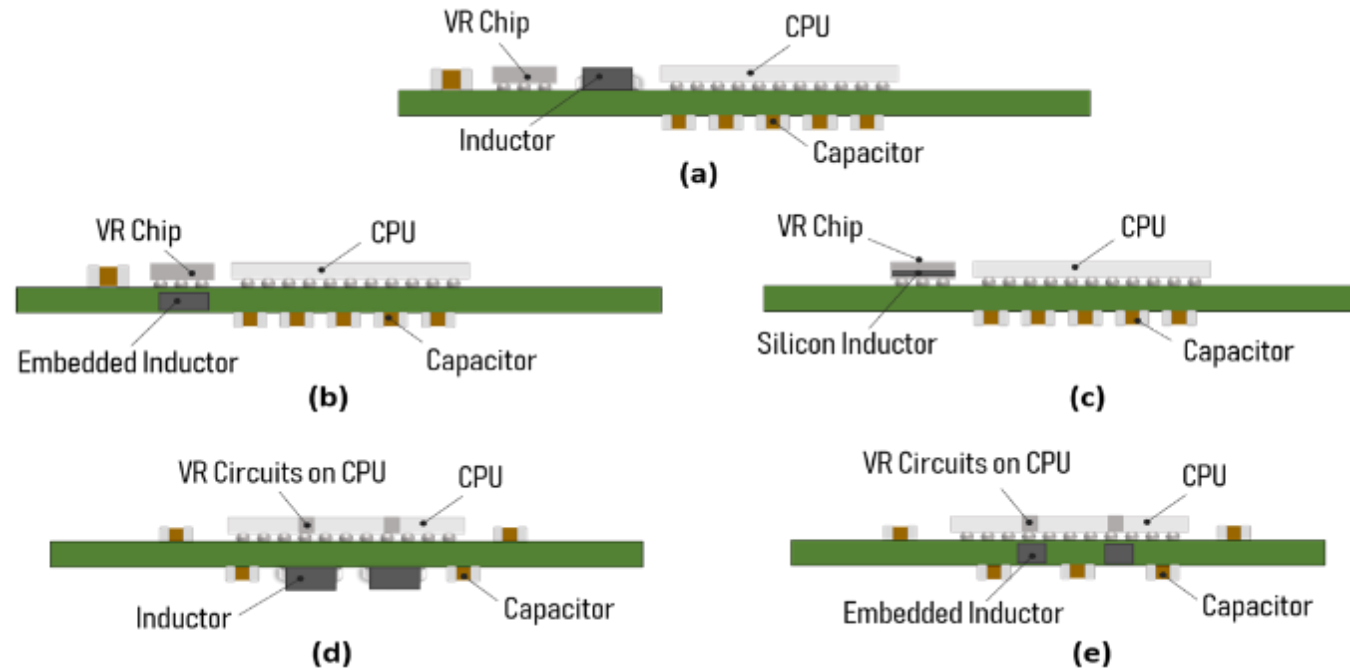
K. Radhakrishnan, M. Swaminathan & B. Bhattacharyya, TCPMT, 2021

IVR – Possible Embodiments

Buck Regulator



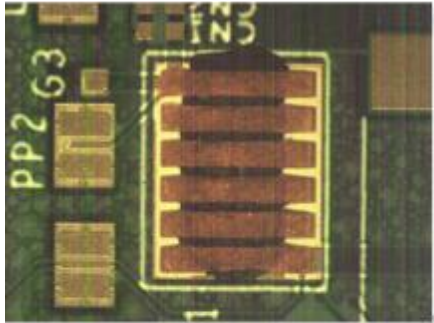
- Key Components
 - Power FETs
 - Inductor
 - Output Filter Capacitor
- More expensive & difficult to integrate
- Typically used for high power rails to take advantage of their high efficiency



- Separate VR Chip on Package or Integrated within CPU
- Discrete or Embedded Inductors in Package
- Air Core or Magnetic Core Inductors

K. Radhakrishnan, EDAPS Keynote, 2021

IVR with Solenoidal Inductors on package



Screen printed inductor on Organic Substrate

	[1]	[2]	[3]	[4]
Phases	4	Up to 16	8	8
Process node	130nm	22nm	45nm	45nm
Integration level	Package	Package	Silicon interposer	Die
Inductor type	Magnetic core	Air core	Magnetic core	Magnetic core
Inductance @100MHz	22.8nH	2.5nH	8nH	1.5nH
Peak eff. (1.7V:1V)	91.7%	90% (1.7:1.05)	82% (1.6:0.83)	84% (1.5:1.15)
Peak eff. (3V:1V)	89.6%	Not specified	Not specified	Not specified
Peak eff. (5V:1V)	81.3%	Not specified	Not specified	Not specified

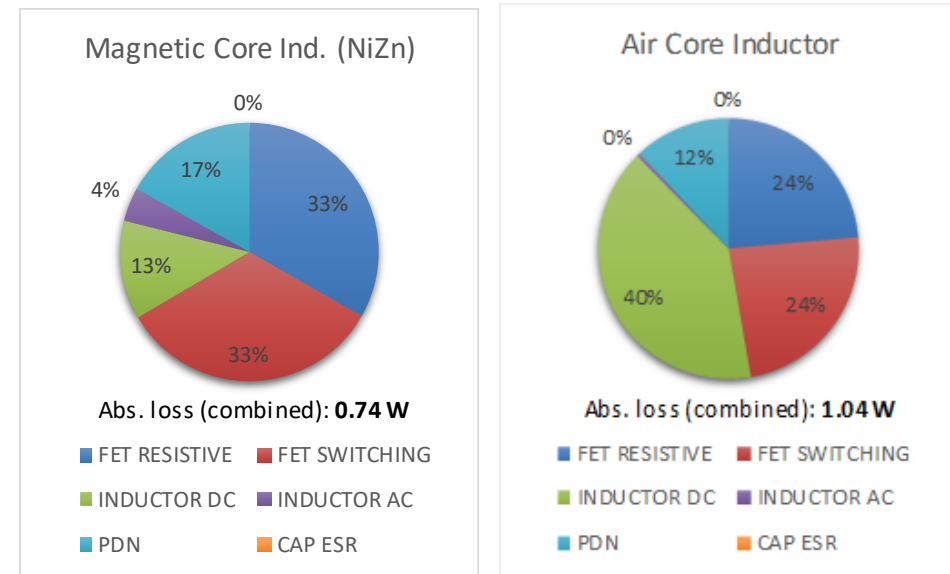
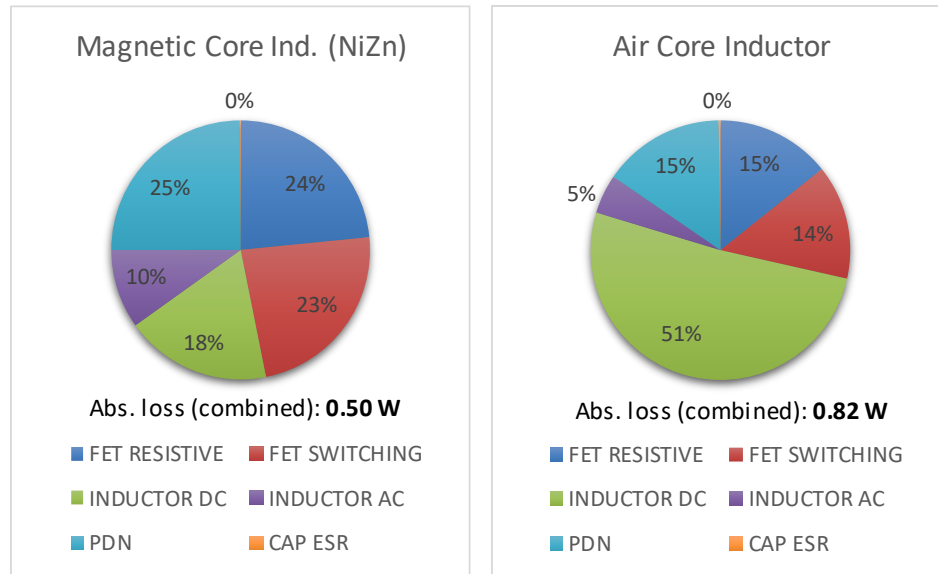
- [1] S. Mueller et al, *Design Exploration of Package-Embedded Inductors for High Efficiency Integrated Voltage Regulators*, IEEE Trans. CPMT, 2018.
- [2] E. A. Burton et al., “FIVR - Fully integrated voltage regulators on 4th generation Intel(R) Core™ SoCs,” in *2014 Twenty-Ninth Annual IEEE Applied Power Electronics Conference and Exposition (APEC)*, 2014.
- [3] K. Tien et al., “An 82%-efficient multiphase voltage-regulator 3D interposer with on-chip magnetic inductors,” in *2015 Symposium on VLSI Technology (VLSI Technology)*, 2015.
- [4] H. K. Krishnamurthy et al., “20.1 A digitally controlled fully integrated voltage regulator with on-die solenoid inductor with planar magnetic core in 14nm tri-gate CMOS,” in *2017 IEEE International Solid-State Circuits Conference (ISSCC)*, 2017.

Courtesy: PDES, GT & JUMP-SRC (ASCENT)

What contributes to Power Losses?

22 MHz

100 MHz



❑ 1.7/1V 5A Load Current (50% Load)

❑ 25nH inductance/phase

❑ 22MHz

- Ferrite Core: 47% FET; 28% Inductor (DC: 18% AC: 10%); 25% PDN
- Air Core: 29% FET; 56% Inductor (DC: 51% AC: 5%); 15% PDN

Magnetic core increases inductance density and lowers DC resistance

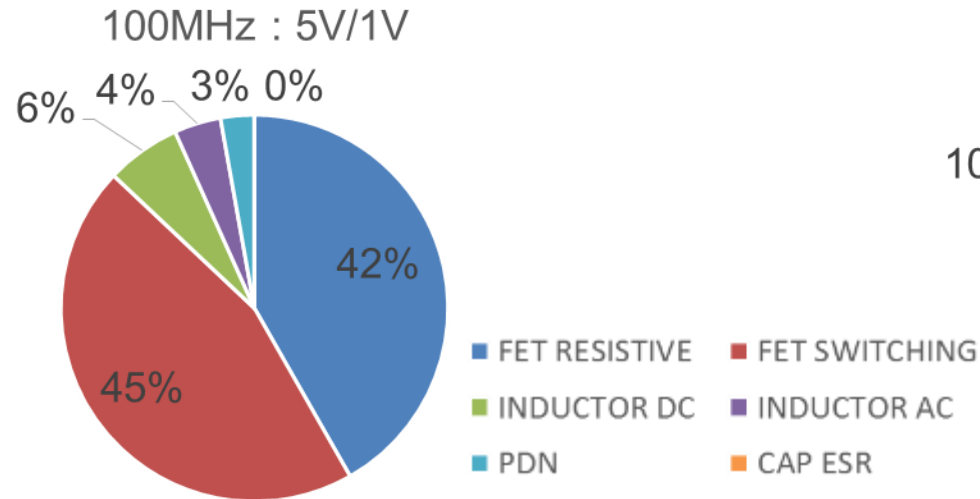
S. Mueller, Member, M.L. F. Bellaredj, A. K. Davis, P. A. Kohl, and M. Swaminathan, *Design Exploration of Package-Embedded Inductors for High Efficiency Integrated Voltage Regulators*, IEEE Trans. CPMT, Volume: 9, Issue: 1, pp: 96 – 106, 2019.

Courtesy: PDES Consortium

Impact of Heating on DC Loss

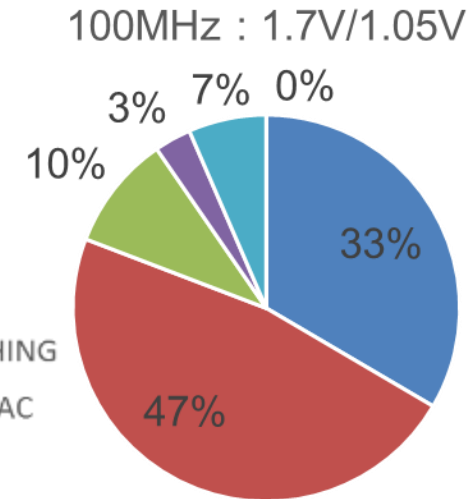
□ Impact of inductor heating on overall efficiency:

5V:1V conversion



- 9% of overall loss related to resistance (“inductor DC” and “PDN”)
- **Inductor heating changes overall efficiency at 100 MHz from 70.5% to 70.2% (small)**

1.7V:1V conversion

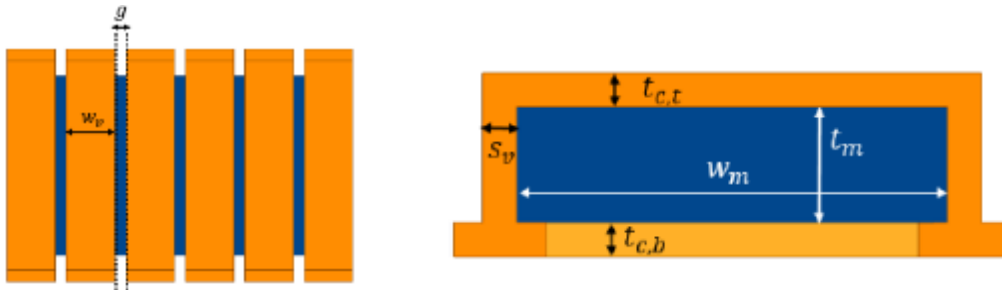
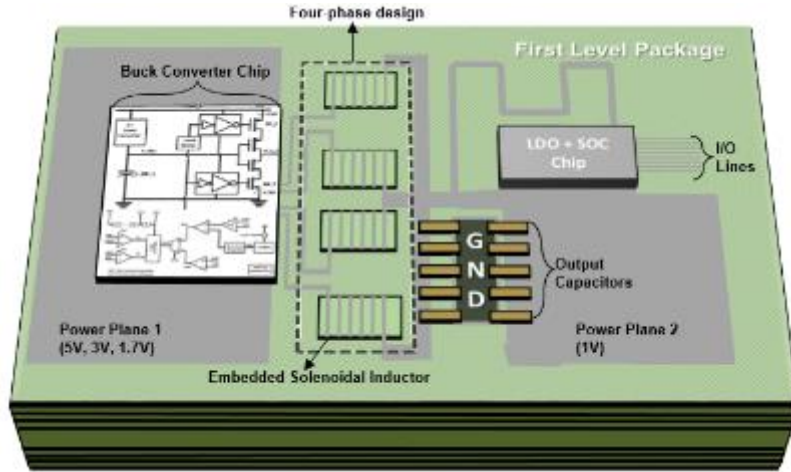


- 17% of overall loss related to resistance (“inductor DC” and “PDN”)
- **Inductor heating changes overall efficiency at 100 MHz from 86.7% to 86.4% (small)**

Courtesy: PDES Consortium

Modeling and Design of System-in-Package Integrated Voltage Regulator with Thermal Effects, S. Mueller, A. K. Davis, M. L. F. Bellaredj, A. Singh, K. Z. Ahmed, S. Mukhopadhyay, P. A. Kohl, M. Swaminathan, Y. Wang, J. Wong, S. Bharathi, Y. Mano, A. Beece, B. Fasano, H. Fathi Moghadam and D. Draper, To be presented at EPEPS, CA, Oct. 23-26, 2016

IVR Optimization (Inductors)



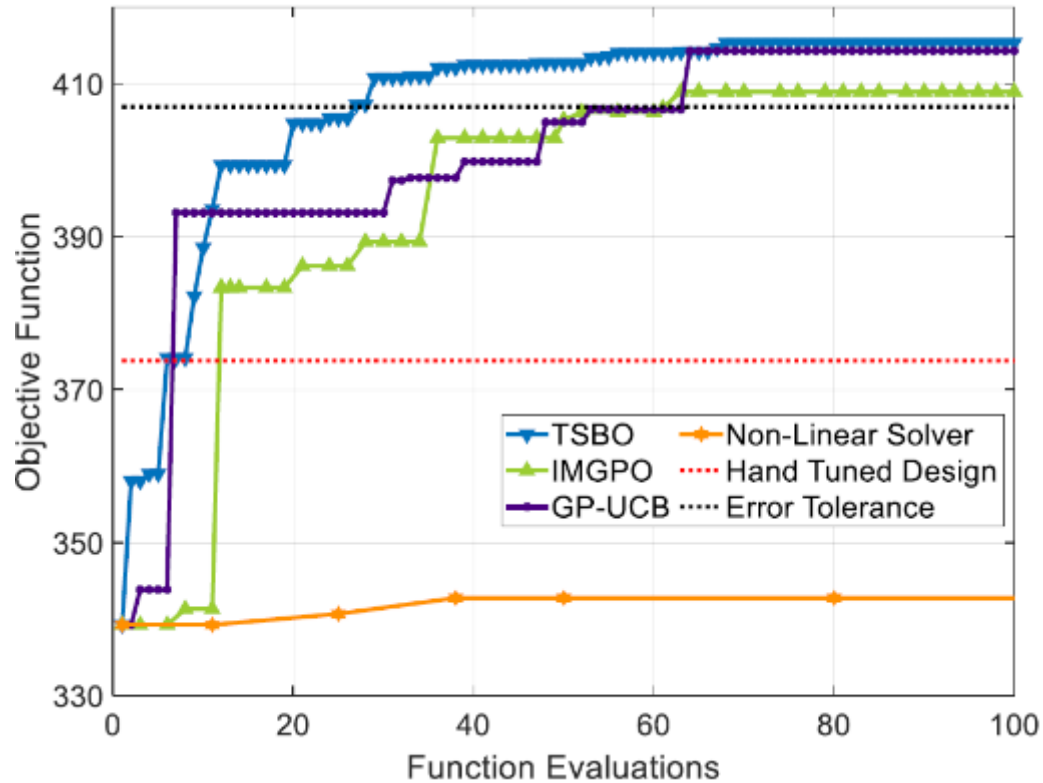
Design Space (10 Parameters)

Parameter	Unit	Min	Max
Gap between windings	g	mil	2 20
Number of windings	N		3 13
Size of via	s_v	μm	50 103
Copper Trace Width	w_c	mil	2 20
Copper Thickness Bottom	$t_{c,b}$	μm	35 170
Copper Thickness Top	$t_{c,t}$	μm	35 170
Dielectric Thickness	t_d	μm	50 650
Dielectric Width	w_d	μm	50 350
Magnetic Core Thickness Ratio	t_m		0.1 1
Magnetic Core Width offset	Δw_m	mil	0 100

- Integrated Voltage Regulators (IVR) are used to increase efficiency and conserve power in microprocessors (Ex: Intel Gen 4)
- Objective is to maximize IVR efficiency while minimizing inductor area
- Inductor is simulated using 3D EM, efficiency is then calculated analytically.
- Multiple trade-offs: ESR, DC resistance, inductance, lateral area
- Tune inductor control parameters to maximize efficiency (10 parameters)



Inductor Optimization Results



	Hand Tuned Design	Non-Linear Solver	GP-UCB	IMGPO	TSBO
Inductor Area	11.3 mm ² (+56.1%)	25.19 mm ² (+79.6%)	5.18 mm ² (%0.4)	6.64 mm ² (%28.1)	5.16 mm²
Peak Efficiency	79.4%	78.6%	84.9%	84.4%	85.1%
CPU Time	N/A	>185 min (+72.9%)	117.33 min (+57.4 %)	115.6 min (+56.7 %)	50.1 min

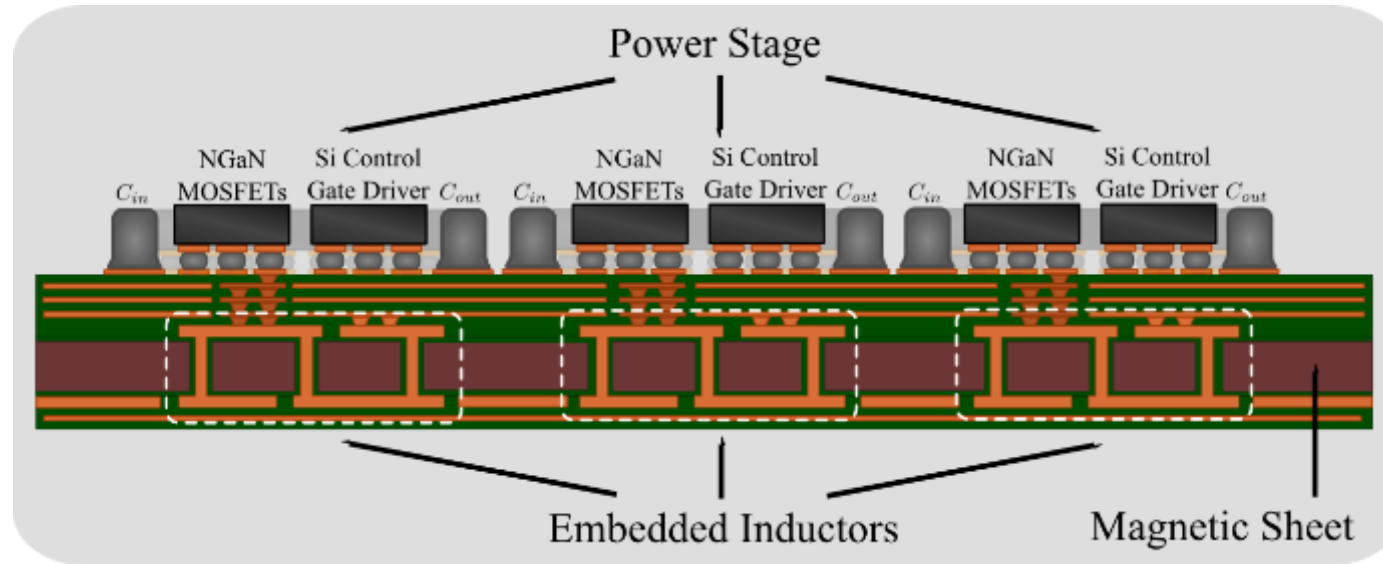


- Optimized IVR has 5.7% higher efficiency and 56.1% reduced area than the hand tuned design.
- TSBO converged ~2.3X faster compared to IMGPO and GP-UCB.

Hand Tuned: S. Mueller et al., "Design of High Efficiency Integrated Voltage Regulators with Embedded Magnetic Core Inductors," ECTC'16.

Optimized: H. M. Torun et al. "A Global Bayesian Optimization Algorithm and Its Application to Integrated System Design". TVLSI'18.

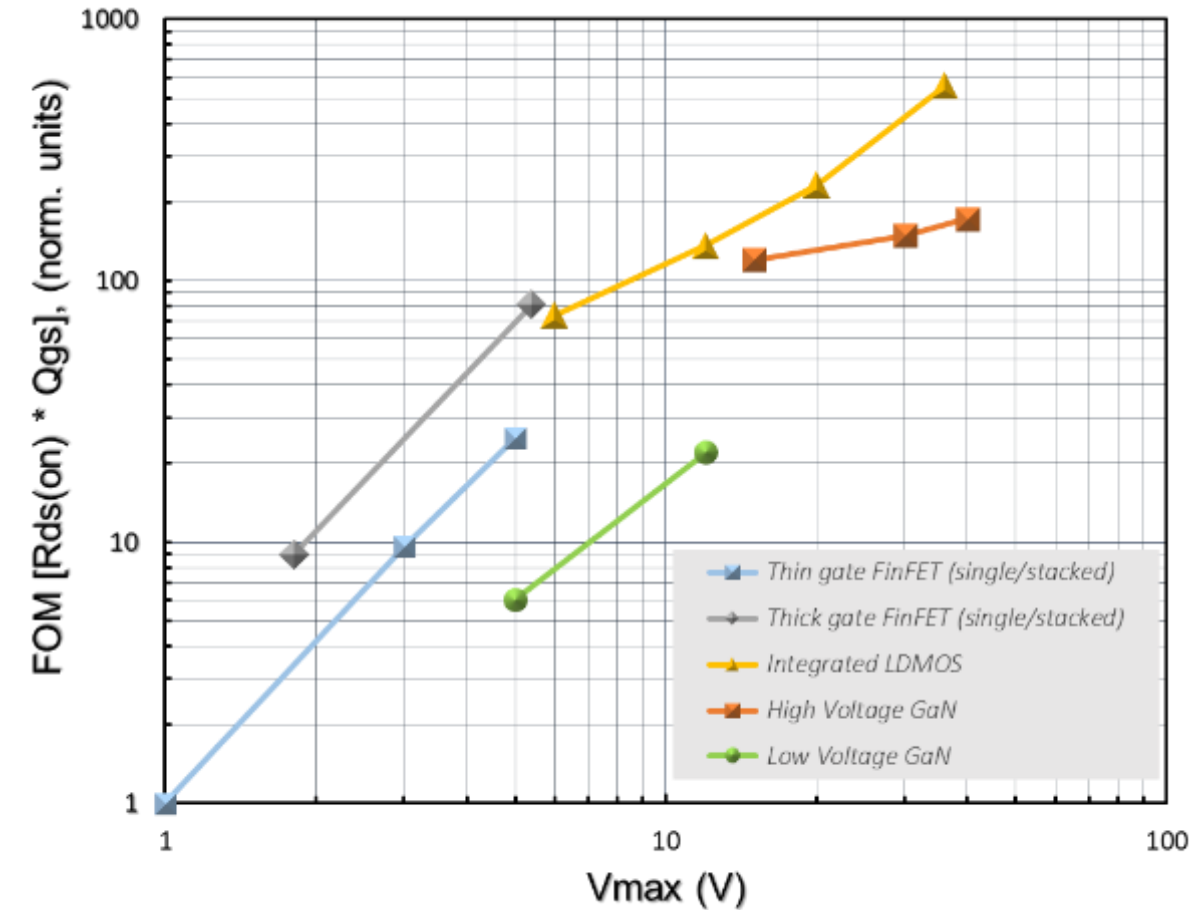
Higher Conversion Ratio IVR to Improve Power Efficiencies



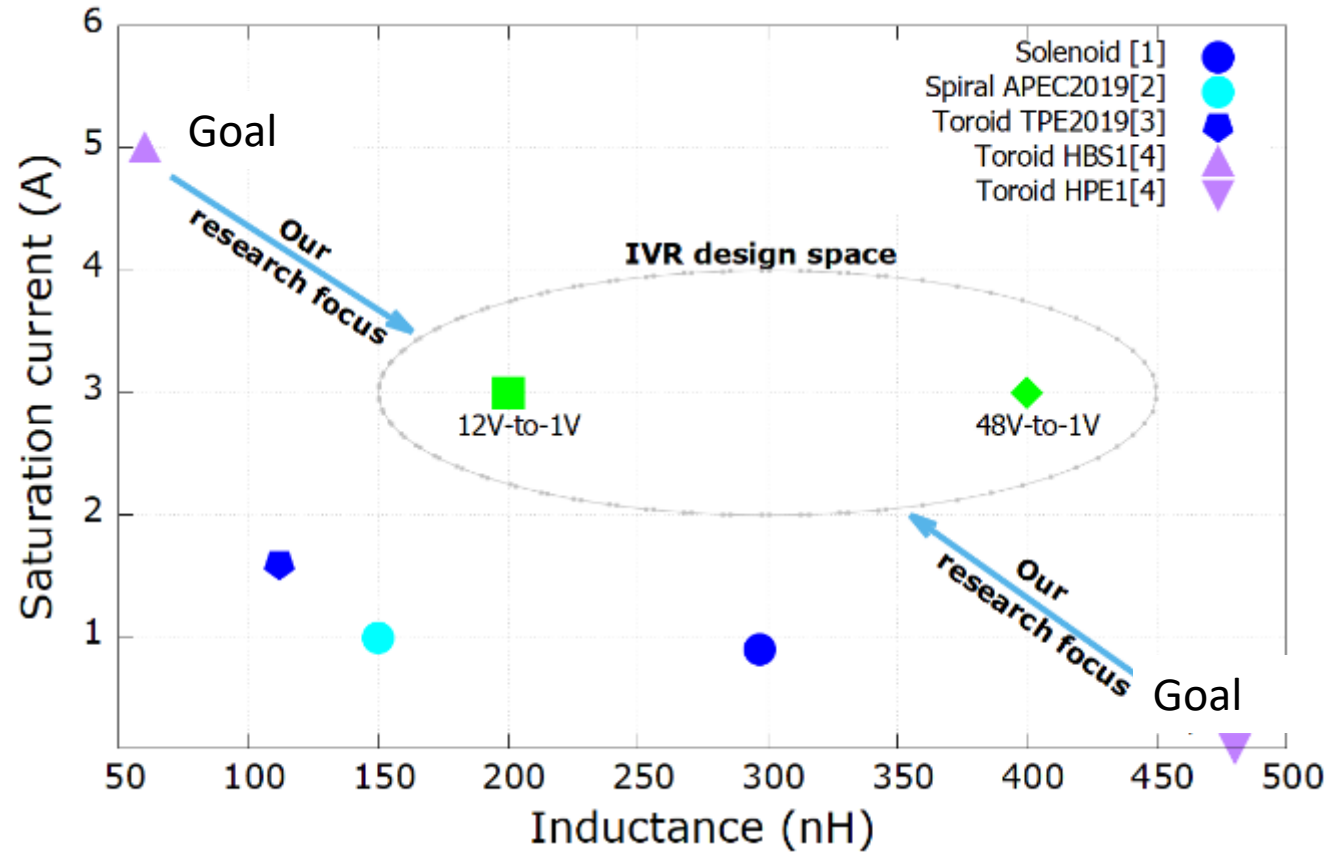
- ❑ 12V-1V and 48V-1V IVR can improve efficiencies.
- ❑ However, Si FETs will lead to low efficiencies.
- ❑ Embedded Inductors are desired to reduce losses.
- ❑ Are efficiencies $\eta > 90\%$ possible (Inductor $\eta_L > 95\%$) ?

K. Radhakrishnan, M. Swaminathan & B. Bhattacharyya, TCPMT, 2021

Device & Embedded Inductor Requirements



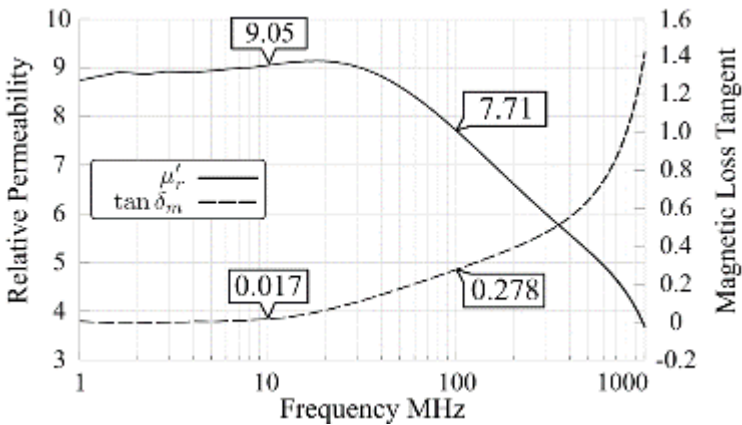
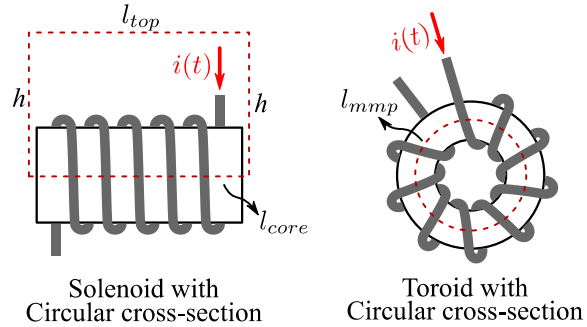
$$FoM = R_{ds(ON)} \times Q_{gs}$$



- [1] Y. Hsieh, S. Lin, C. Kung, P. Lee, and C. Wang, IEEE EPTC, 2019.
- [2] T. Fukuoka et al., 2019 IEEE APEC), 2019.
- [3] H. T. Le et al., IEEE Transactions on Power Electronics2019.
- [4] Claudio, PhD thesis, 2021, Georgia Tech

Challenges with Embedded Inductors

1. Magnetic Material

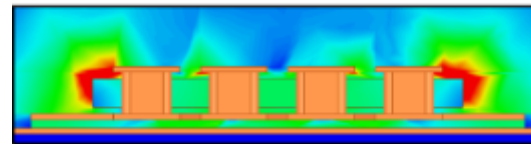


$$L \approx \mu_0 \mu_e \frac{N^2 A}{l}$$

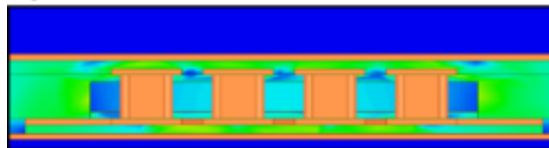
2. Closed Loop High Q Inductors



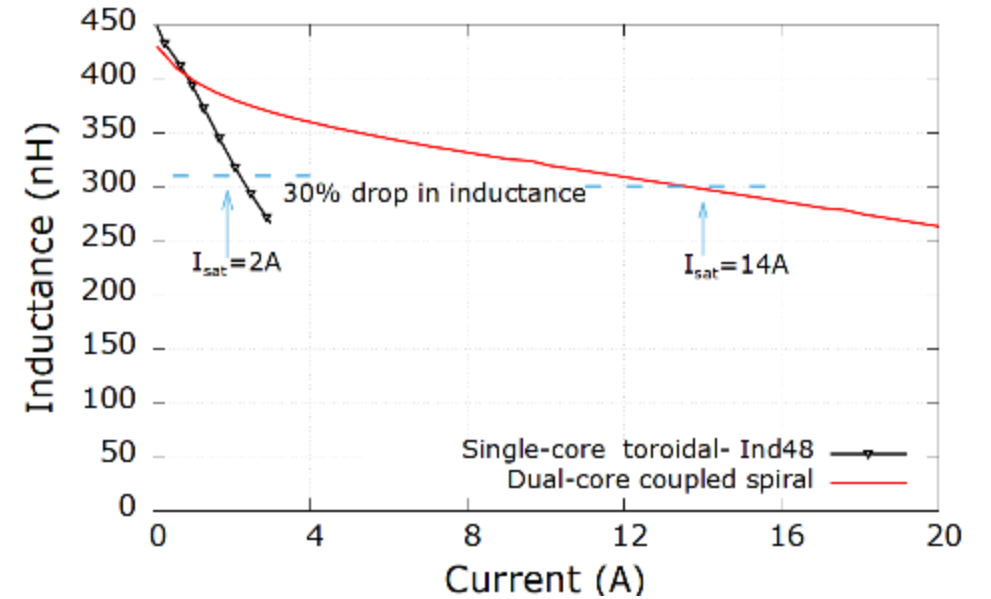
Surface
 $L = 25.1 \text{ nH @ } 10\text{MHz}$
 $Q\text{-Factor} = 35 \text{ @ } 10\text{MHz}$



Embedded
 $L = 16.1 \text{ nH @ } 10\text{MHz}$
 $Q\text{-Factor} = 21 \text{ @ } 10\text{MHz}$

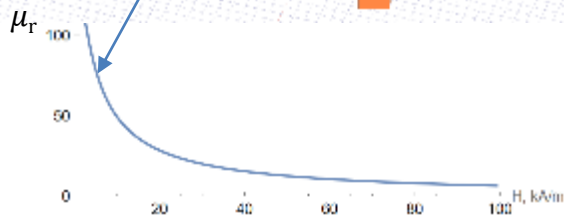
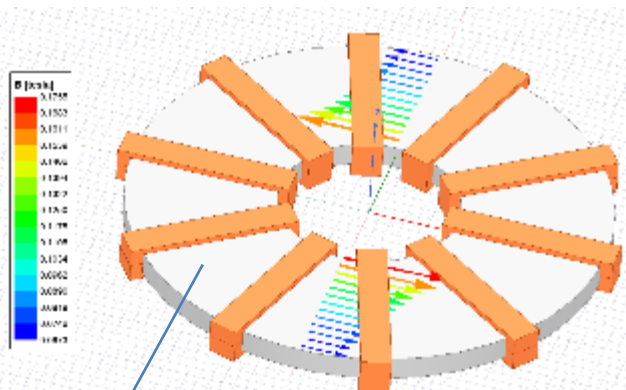
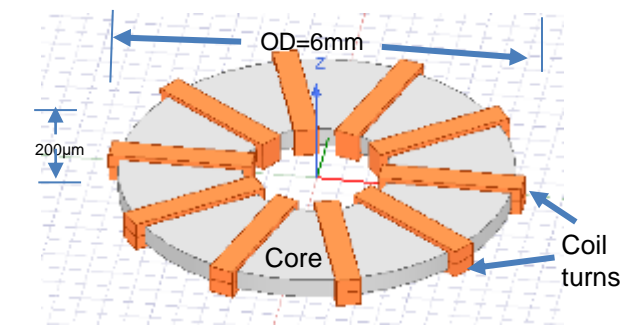


3. High Saturation Current

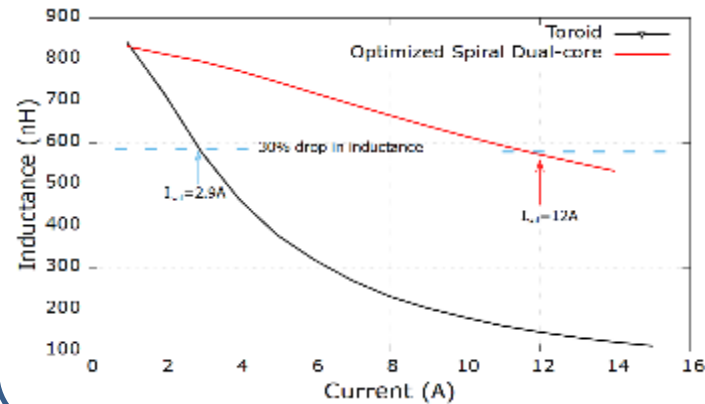
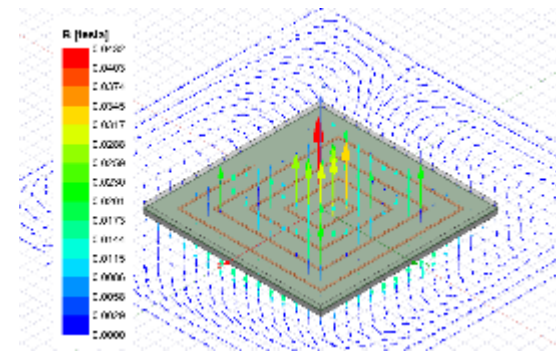
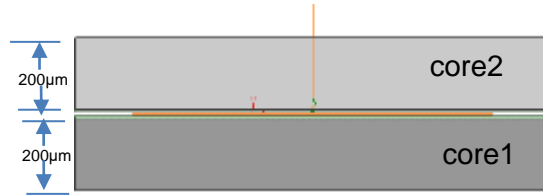


Toroid Vs Spiral Inductor Topology

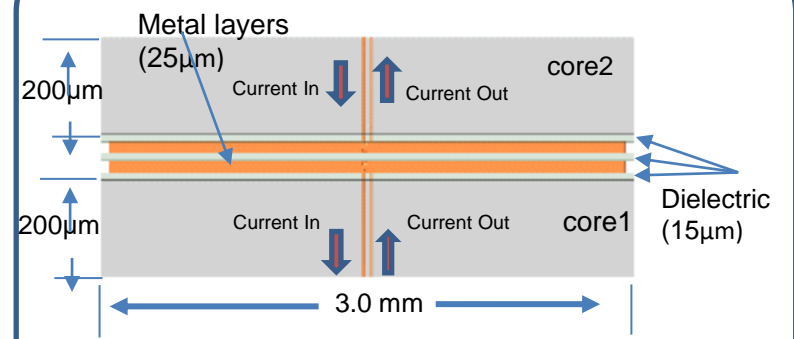
Traditional toroid inductor topology



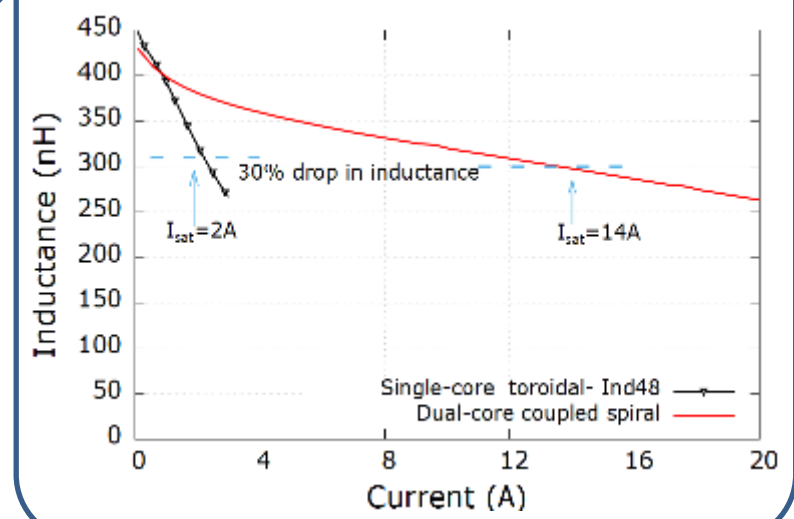
Dual-core spiral topology



Dual-core coupled spiral topology



$$L_{total} = L_1 + L_2 + 2M$$

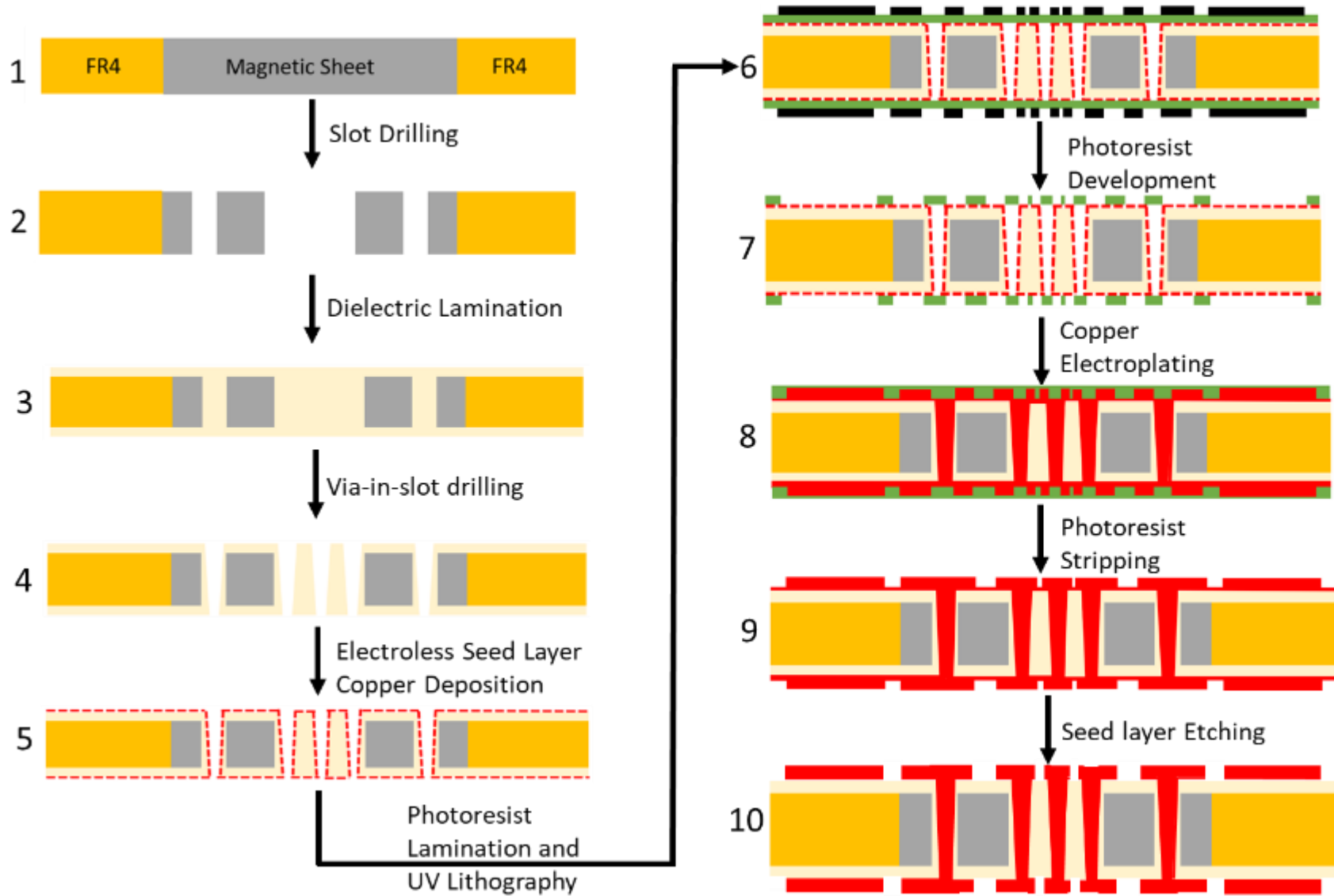


➤ Toroid inductor has a higher magnetic field present in its core; therefore, its saturation current is lower.

➤ Dual-core spiral topology improves saturation current performance of the toroidal inductors by 4x.

➤ Dual-core coupled spiral topology improves saturation current performance of the toroid inductors by 7x.

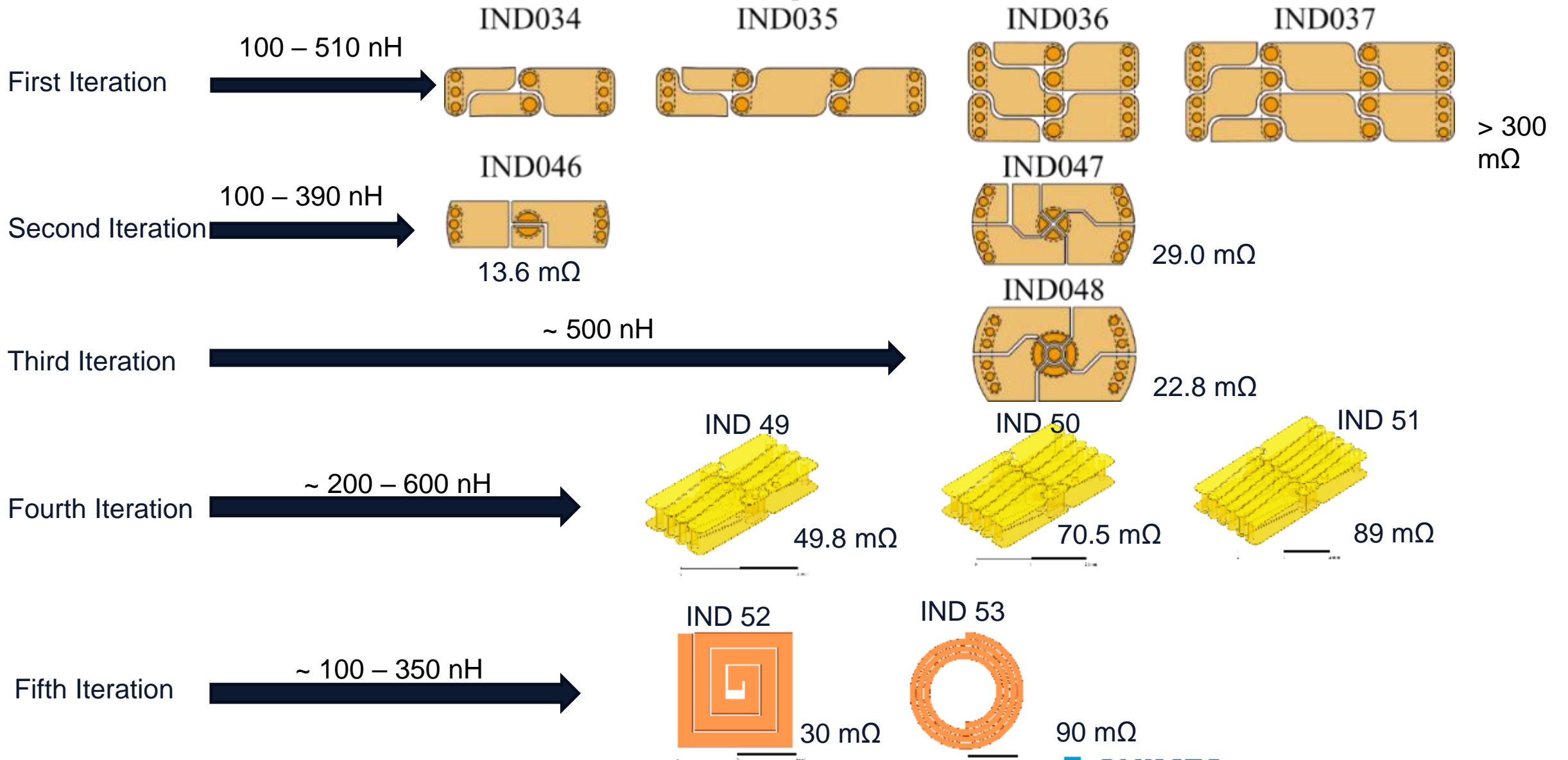
Fabrication Process Flow



Process flow for metal layer 1 of toroidal inductors using semi additive patterning process

Murali, P., et al., 2022 Elsevier *PEDC*

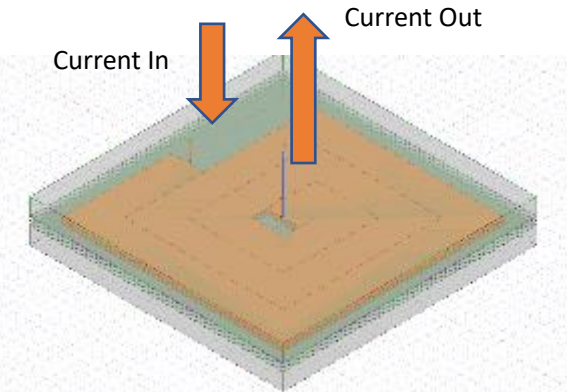
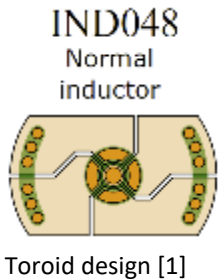
Toroidal and Spiral Inductors



Comparison between Toroid and Spiral Inductors

- Since core losses for sinusoidal excitations are proportional to square of magnetic field (B^2), as per Steinmetz's equation is, $P_L = K_m f^\alpha B^\beta$, the core losses in the dual-core topology is 49x lower.
- Therefore, large signal saturation current of 2 A to 3 A is expected for the dual-core topology.

Parameter	IVR requirements	Toroid HBS1 IND048 [1]	Coupled dual-core Spiral- IND 53 HPE1
Area (mm ²)		1.62	9
Thickness (mm)		0.6	0.5
Turns		5	3
Resistance(mΩ)	25	23	50
Inductance (nH)	400 nH	450	430
Inductance density (nH/mm ²)		96	48
I _{sat} (A) (from DC bias simulation)		2	14 (7x higher)
Core loss, P _L (W)			~49 x lower
I _{sat} (A) (from large signal measurement)	3 A	0.1	2 A to 3 A ^



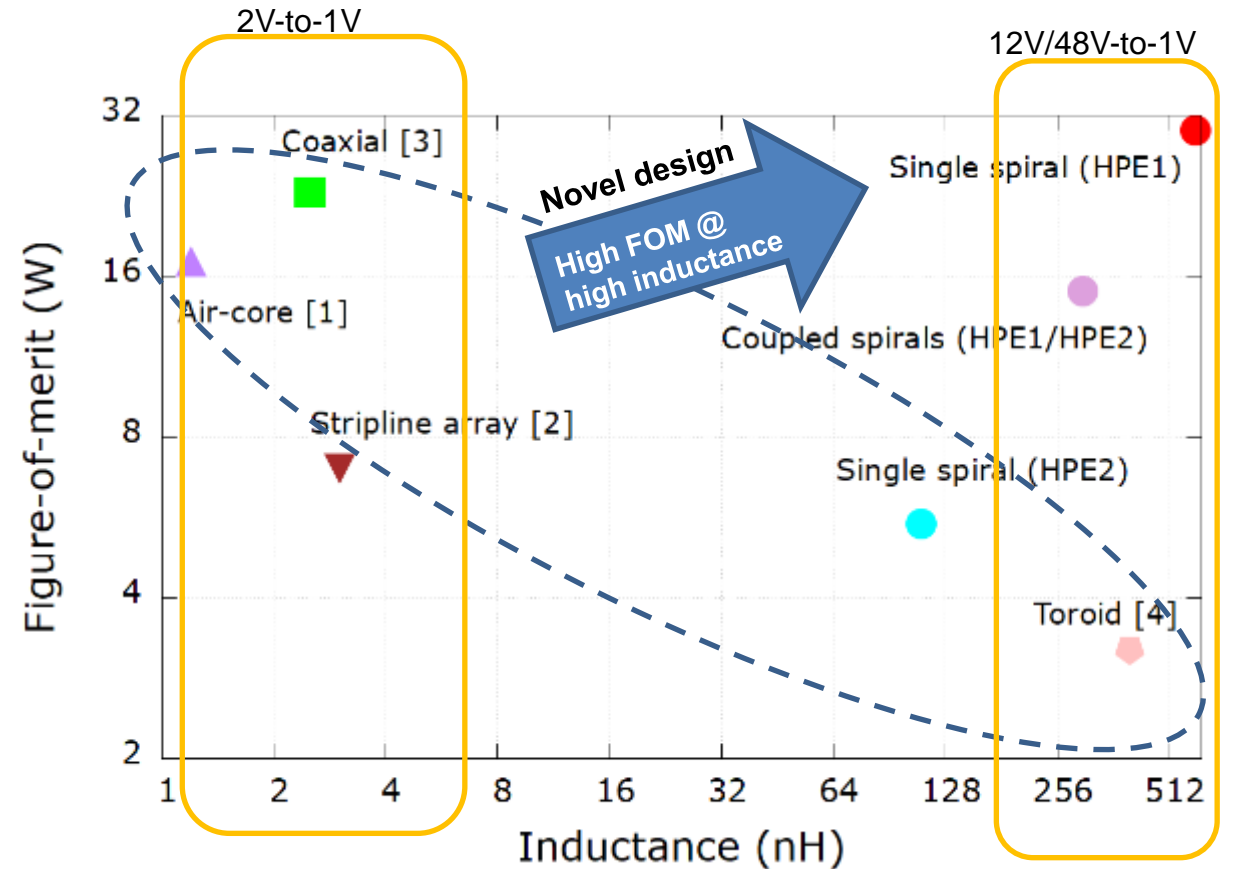
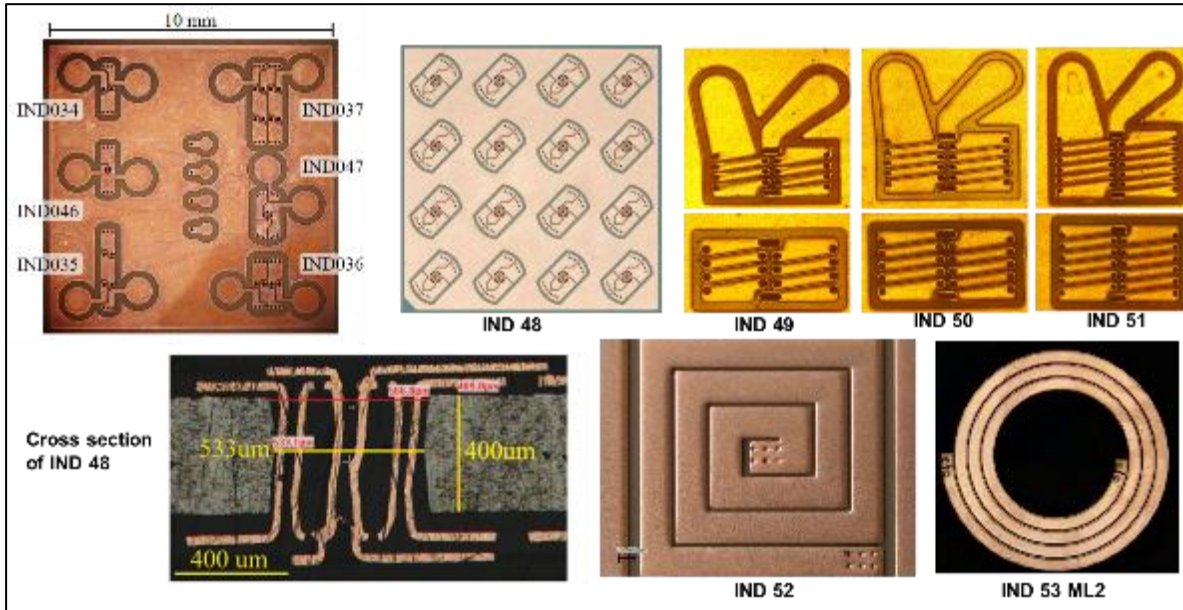
Dual magnetic core-coupled spirals- 3 turns each

^(expected due to 49x core loss reduction and 7x DC bias saturation improvement)

[1] Claudio Alvarez et al, Design and Demonstration of Embedded Inductors for High-Voltage Integrated Voltage Regulators, PhD Thesis, Georgia Tech, 2021

Figure of Merit for Inductors (FOM)

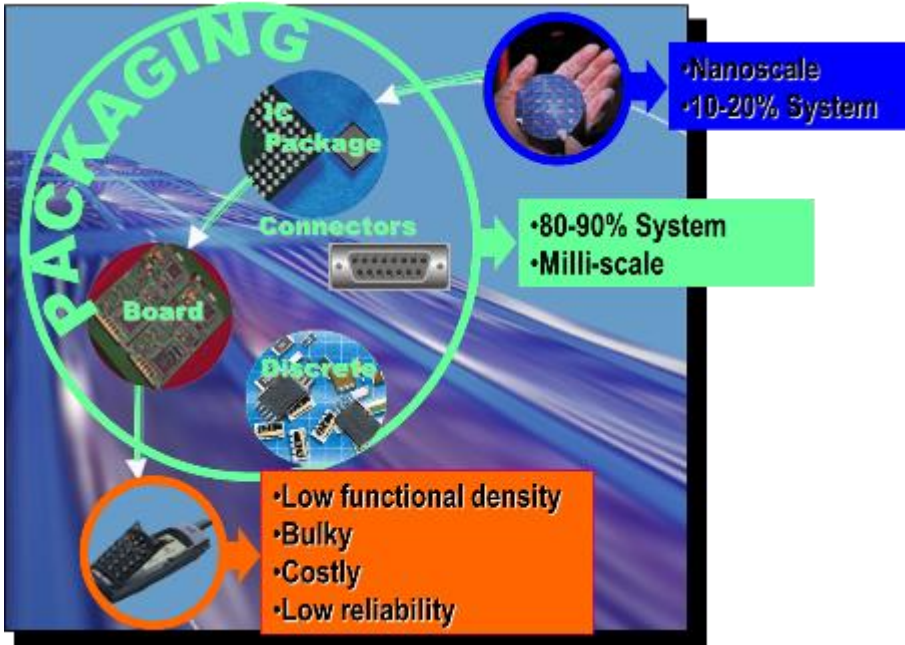
$$FOM = LI^2 f \left(\frac{\text{Energy Stored}}{\text{Cycle}} \right)$$



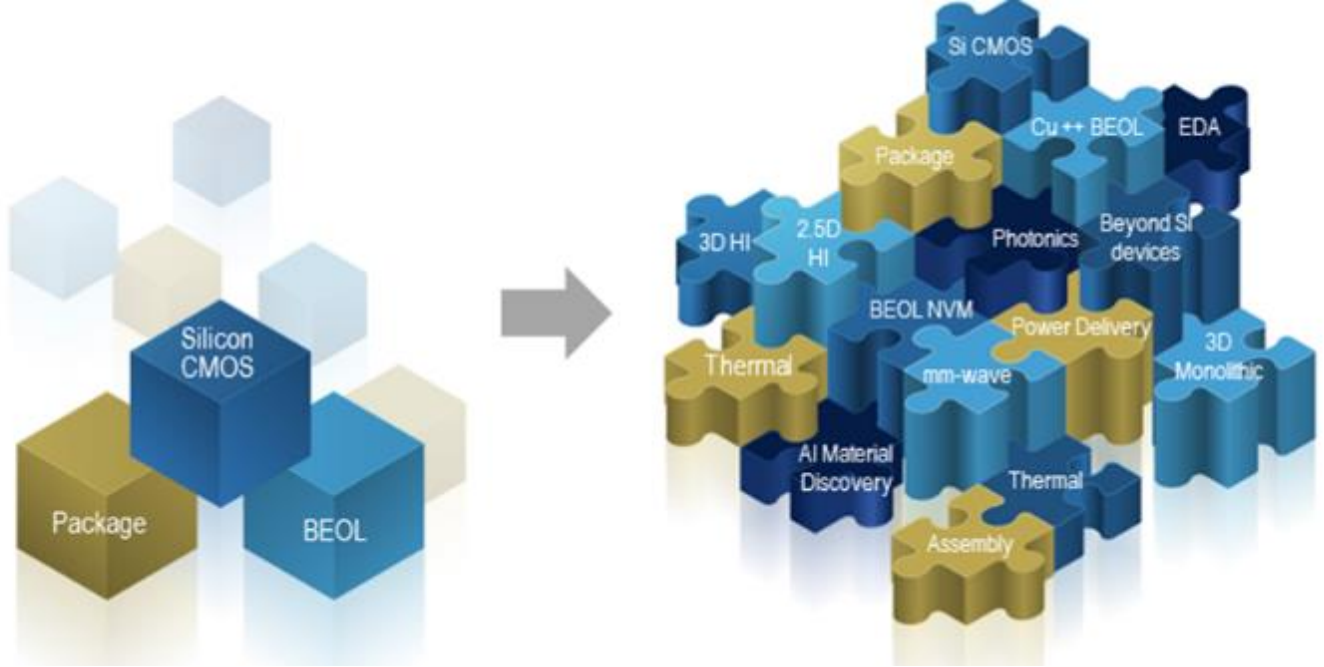
Comparison of package-embedded inductors

- [1] Burton et al, APEC 2014
- [2] M. Sankarasubramanian et al, ECTC2020
- [3] Krishna Bharath et al, ECTC 2021
- [4] Prahalad et al, ECTC 2022

Improving Efficiencies Further – A Case for Heterogeneous Integration



Packaging (Past)

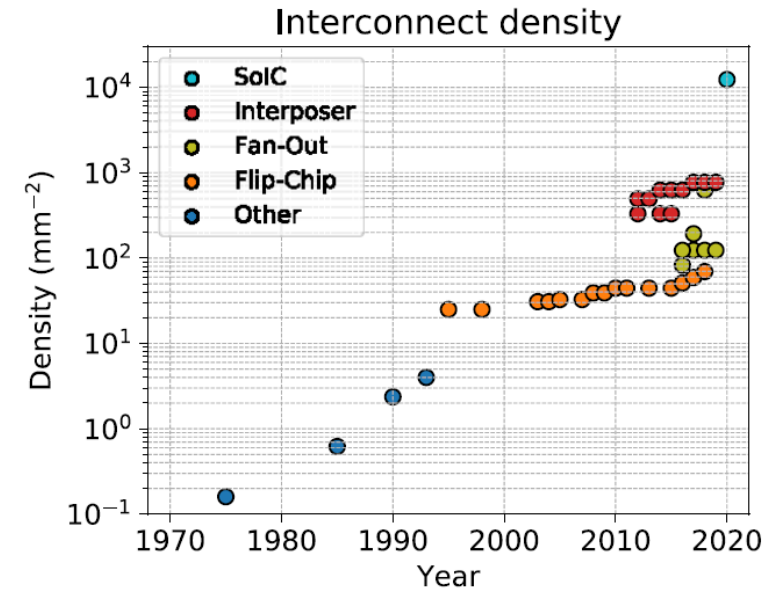
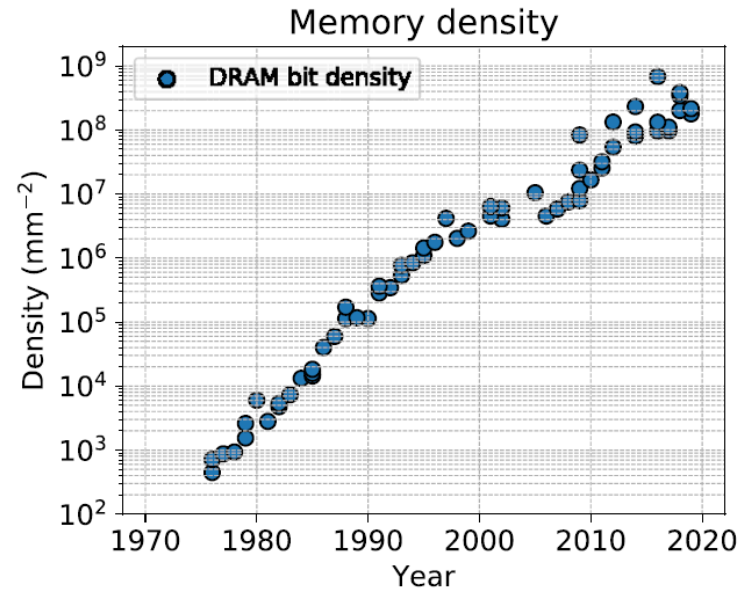
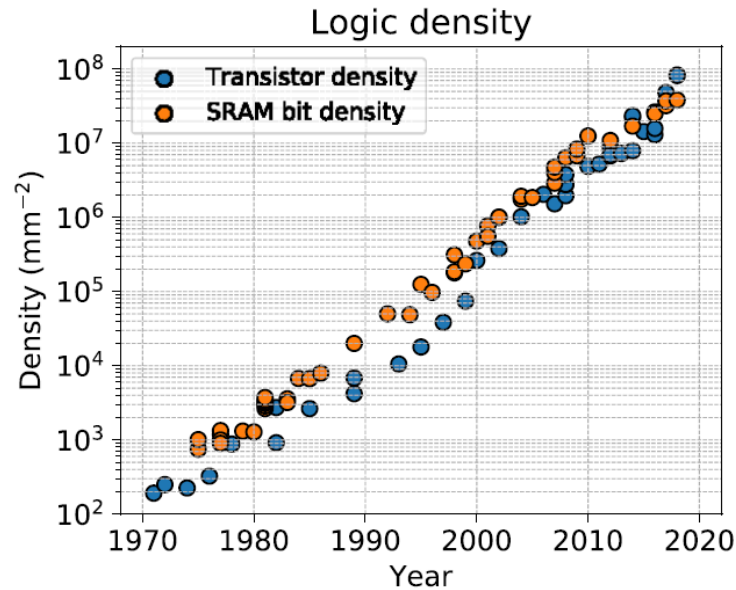


Advanced Packaging (Present)

Heterogeneous Integration (Future)

- ❑ **Past/Present:** FEOL Transistor, BEOL Wiring & Package individually developed and combined
- ❑ **Future:** New and transformative logic, memory, and interconnect technologies that overcome the inevitable slowdown of traditional dimensional scaling of CMOS by interconnecting a **diversity of transistors and integrated circuit components, blurring the line between what is on-chip and what is off-chip.**

Historical Trends & Future Needs



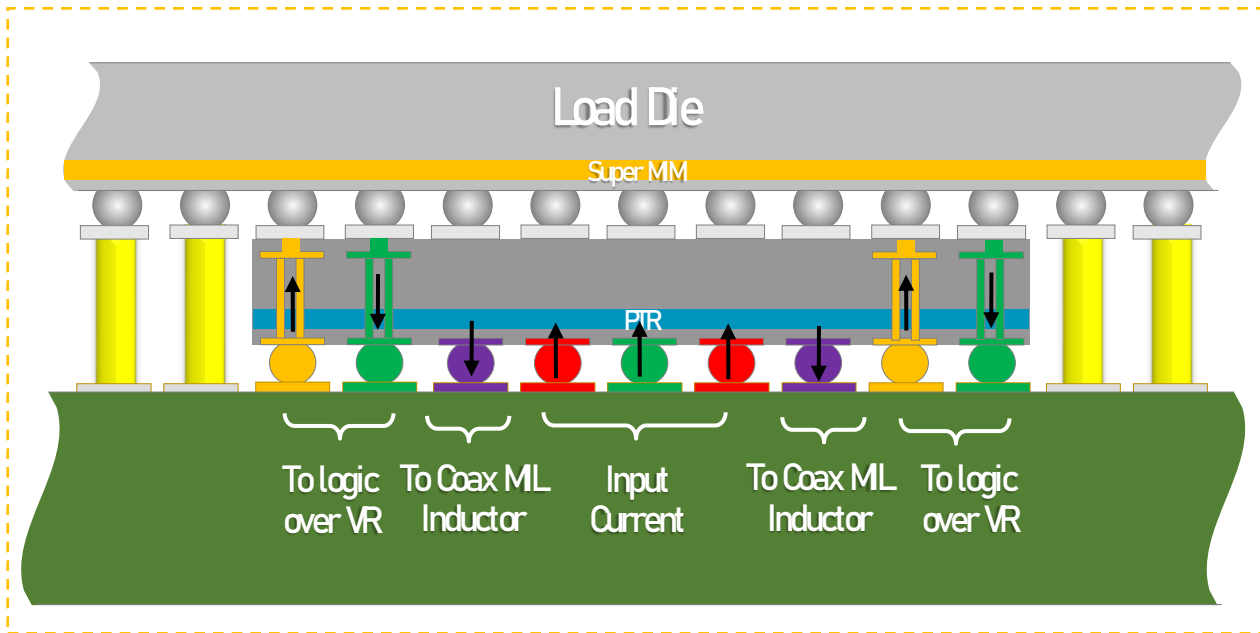
H.-S. Philip Wong, et al, "A Density Metric for Semiconductor Technology", Proceedings of the IEEE, April 2020

Current State of the Art

- Monolithic logic 10^8 transistors/mm²
- DRAM 10^9 transistors/mm²
- **IO density 10^4 IO/mm²**
- SRAM Access 20-50 TBps

10-100X increase in transistor densities
Interconnect densities 10^6 and higher (100X)
 Energy per bit (EPB) reduced to femto-joules/bit
 500TBps/mm² of bandwidth (10X increase)
 Wireless communication at 1Tbps

Emergence of 3D Technologies

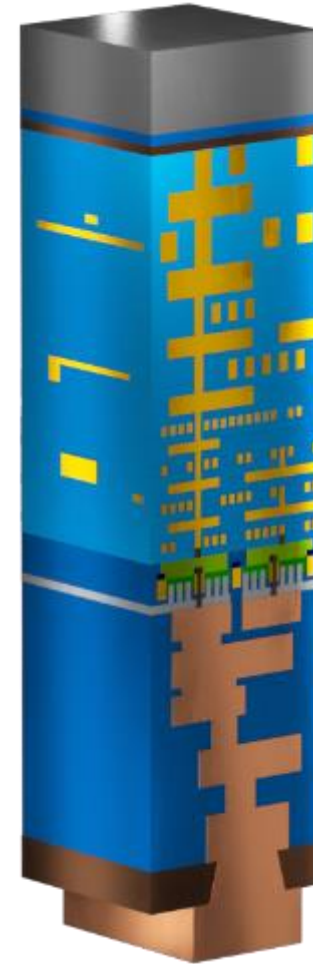


To logic over VR To Coax ML Inductor Input Current To Coax ML Inductor To logic over VR

- VCCIN (Input Voltage)
- VSS (Ground)
- VCCOUT (Output Voltage)
- VXBR (Switch node)

K. Radhakrishnan, EDAPS Keynote, 2021

Power Vias

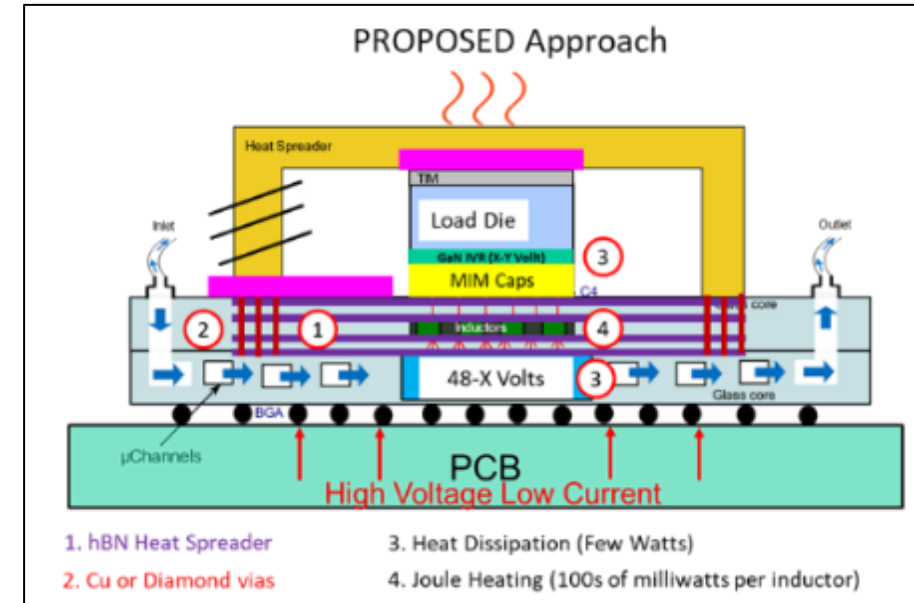
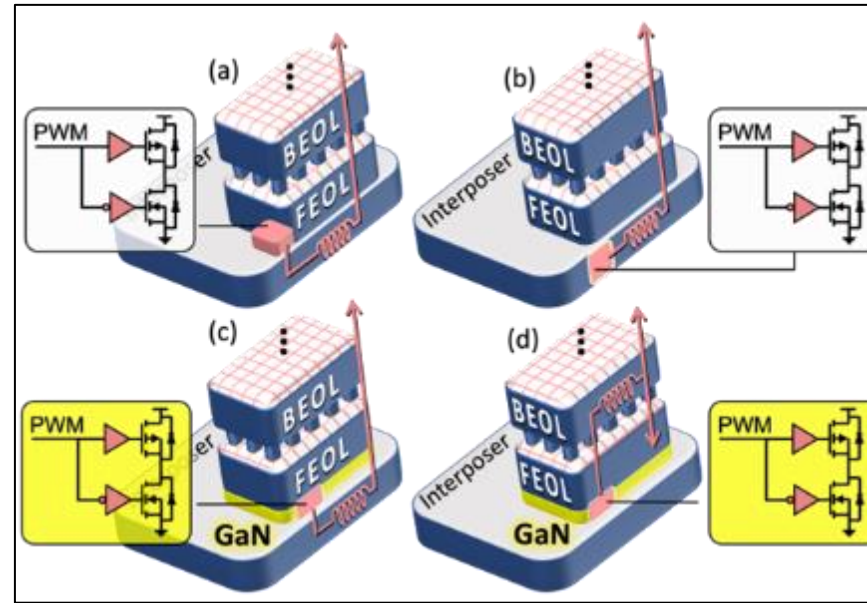
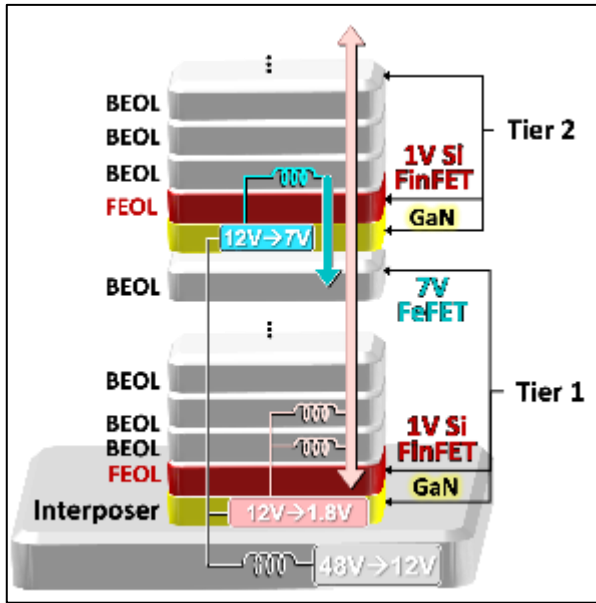


Front Side Interconnects: Signal Routing

Transistors
Nano TSV

Back Side Interconnects: Power Routing

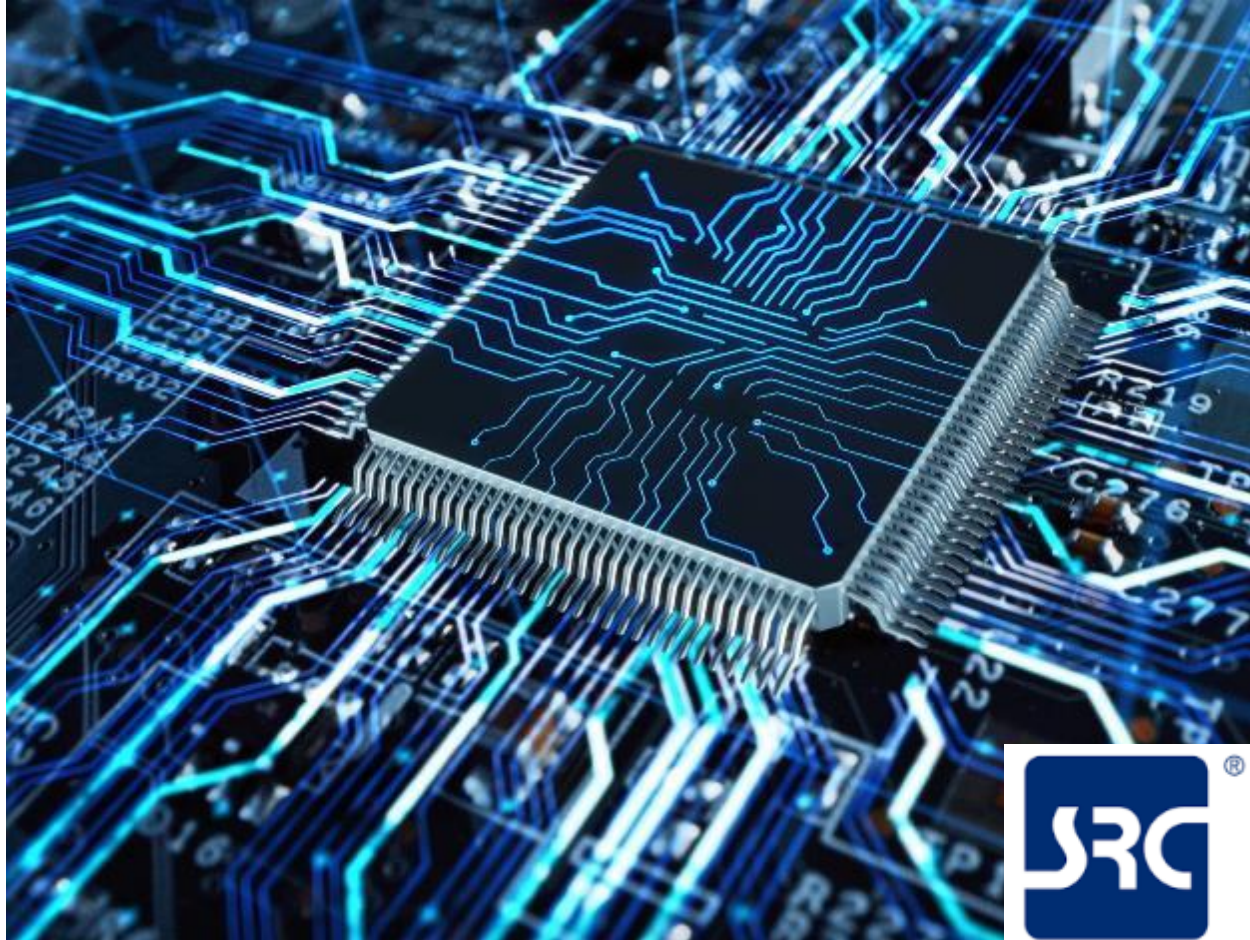
Further Integration



- ❑ Si VR on top of embedded inductor in interposer.
- ❑ GaN Voltage Regulator (VR) on backside of front-end-of-the-line (FEOL) with Power Vias.
- ❑ POL architecture with GaN VR on backside FEOL and inductors on back-of-the-line (BEOL).
- ❑ Multi-stage combination of the power architectures with intermediate voltage values.
- ❑ Integration of heat spreaders and micro-fluidics to remove heat from embedded ICs and inductors.

JUMP 2.0 Center @ Penn State (14 Univ. Partners)

Penn State leads semiconductor packaging, heterogeneous integration center

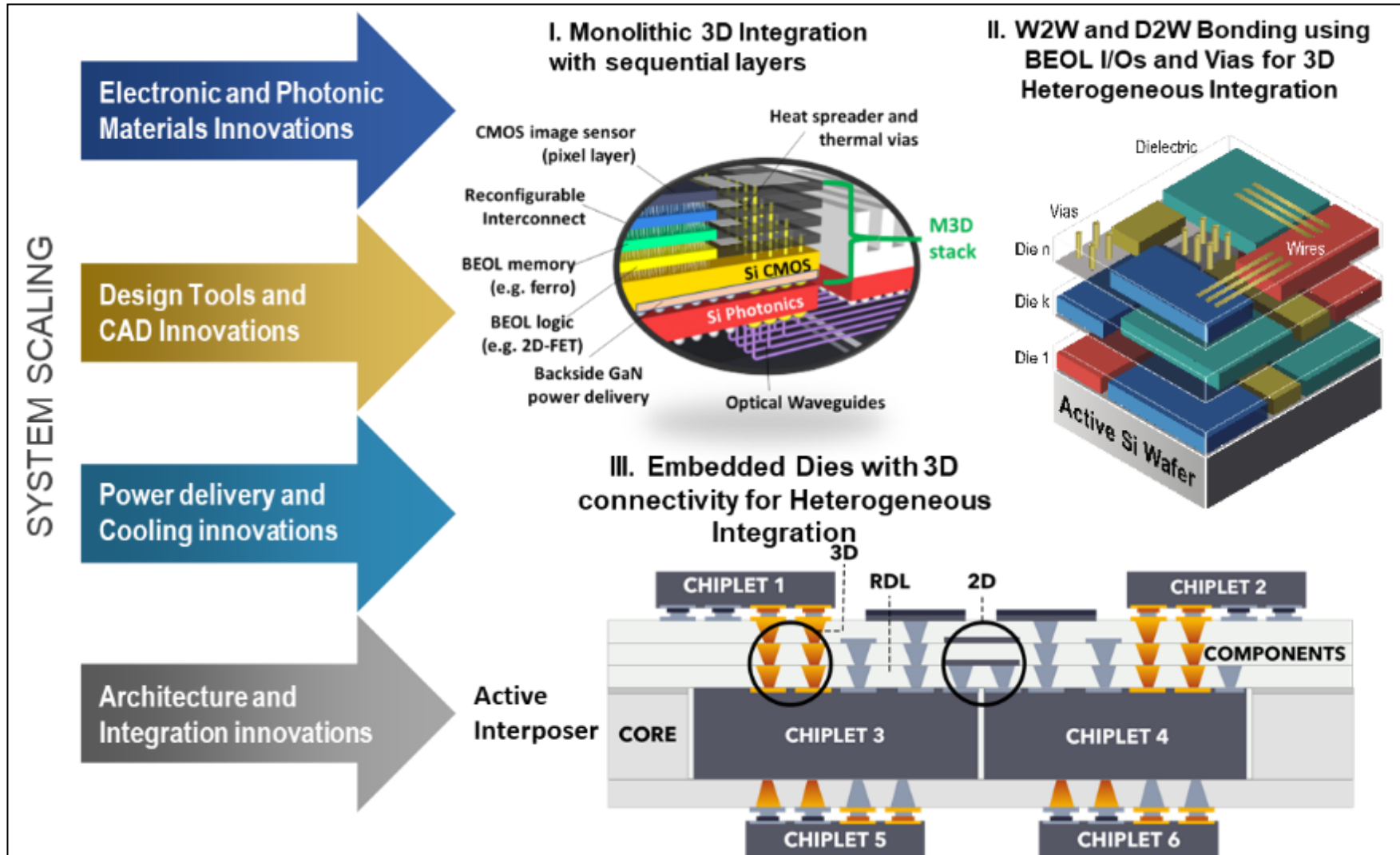


- Center for Heterogeneous Integration of Micro Electronic Systems (CHIMES)
- Supported by the Semiconductor Research Corporation (SRC)'s Joint University Microelectronics Program 2.0 (JUMP 2.0), a consortium of industrial partners in cooperation with the Defense Advanced Research Projects Agency (DARPA)



<https://www.psu.edu/news/engineering/story/penn-state-leads-semiconductor-packaging-heterogeneous-integration-center/>

CHIMES Focus



Summary

- ❑ System Scaling by using more GPUs for Deep Learning will reach an energy limit.
- ❑ Need to address Power Delivery and Thermal Management at Scale.
- ❑ Integrating Power Sources with new devices and storage elements is a possible solution. Major advances in the last 5 years.
- ❑ Heterogeneous Integration is key.
- ❑ Vertical power delivery structures is a necessity to reduce losses.
- ❑ New materials and technologies needed to enable 3D Integration. Not there yet!

