# CUstom Built HEterogeneous Multi-Core ArCHitectures (CUBEMACH): Breaking the Conventions

**Nagarajan Venkateswaran**
**Director, Waran Research Foundation**
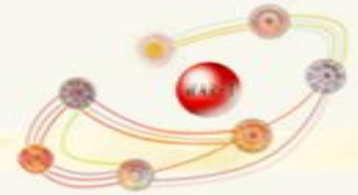
**Karthikeyan Palavedu Saravanan - Nachiappan Chidambaram Nachiappan**
**Research Trainees (2008 - 2010), Waran Research Foundation**

**Aravind Vasudevan - Balaji Subramaniam - Ravindhiran Mukundarajan**
**Former Research Trainees (2007 - 2009), Waran Research Foundation**

# Motivation : Heterogeneity Redefined

- Cost  Effective High Performance Custom Built Heterogeneous Multi-Core Node Design for wider class applications
  - Inter and Intra core heterogeneity
- Breaking the Conventions
  - Multiple User Multiple Application without Space-Time sharing in a Cluster : Cost sharing across users
  - Single User Multiple Application without Space-Timer Sharing (non-multiprogramming) : Cost sharing across applications
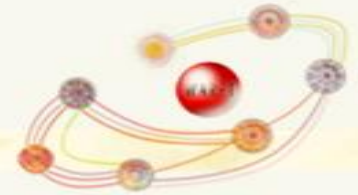
- Custom Built Heterogeneous Multi-Core Architectures (CUBEMACH)
- Design Space
  - Architectural Space
  - Optimization Space
    - Customer Vendor Interaction
  - Simulation Space
- CUBEMACH Design and Simulation Tool Framework
- Conclusion

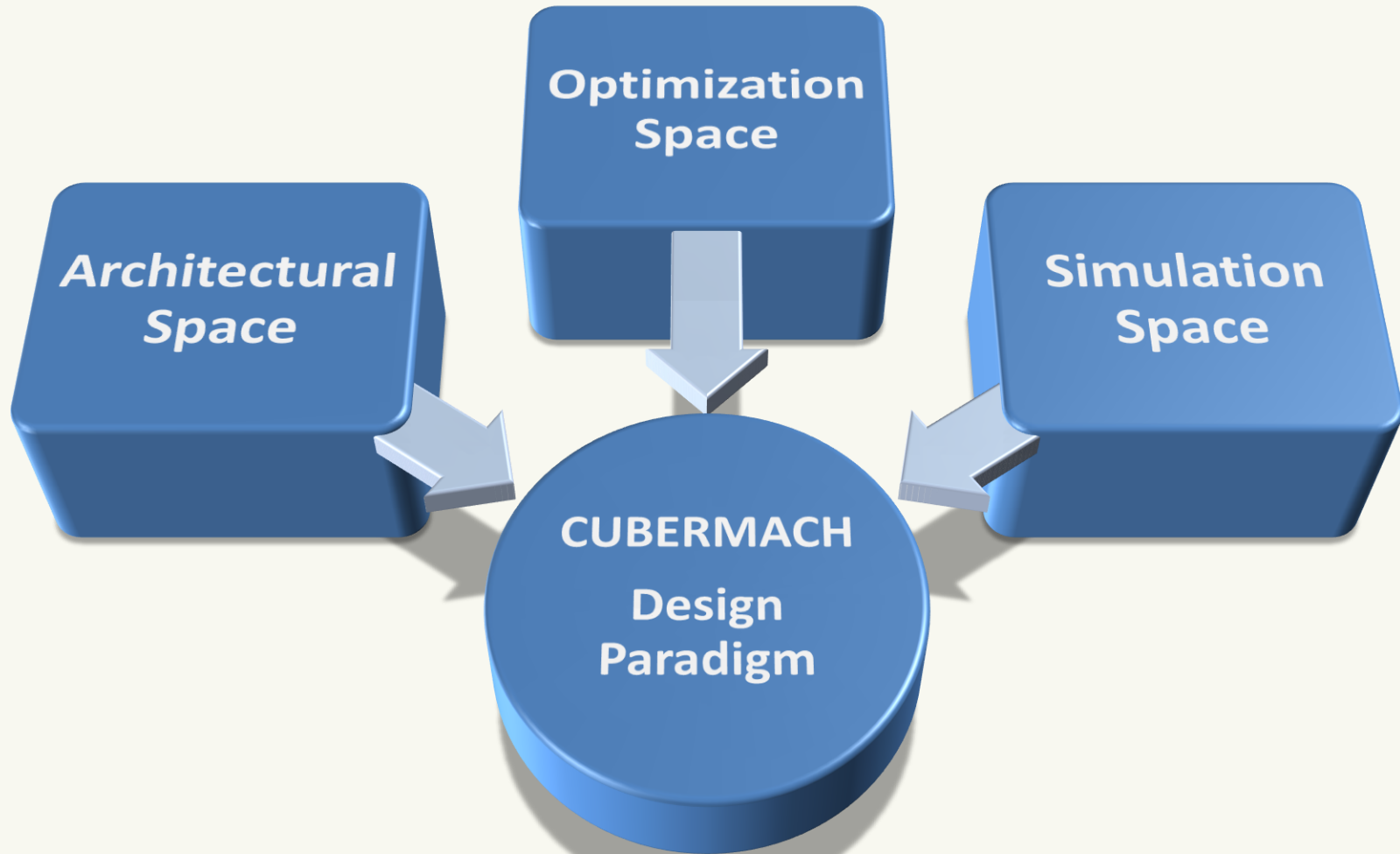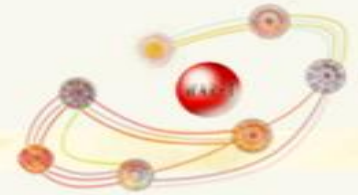# Custom Built Heterogeneous Multi-Core Architectures (CUBEMACH)

- ## CUBEMACH promises
  - Increased Resource Utilization
  - Multiple Application Flavored Architectures
  - Elimination of Space Time Sharing at the Quantum Level during Multiple Application Execution
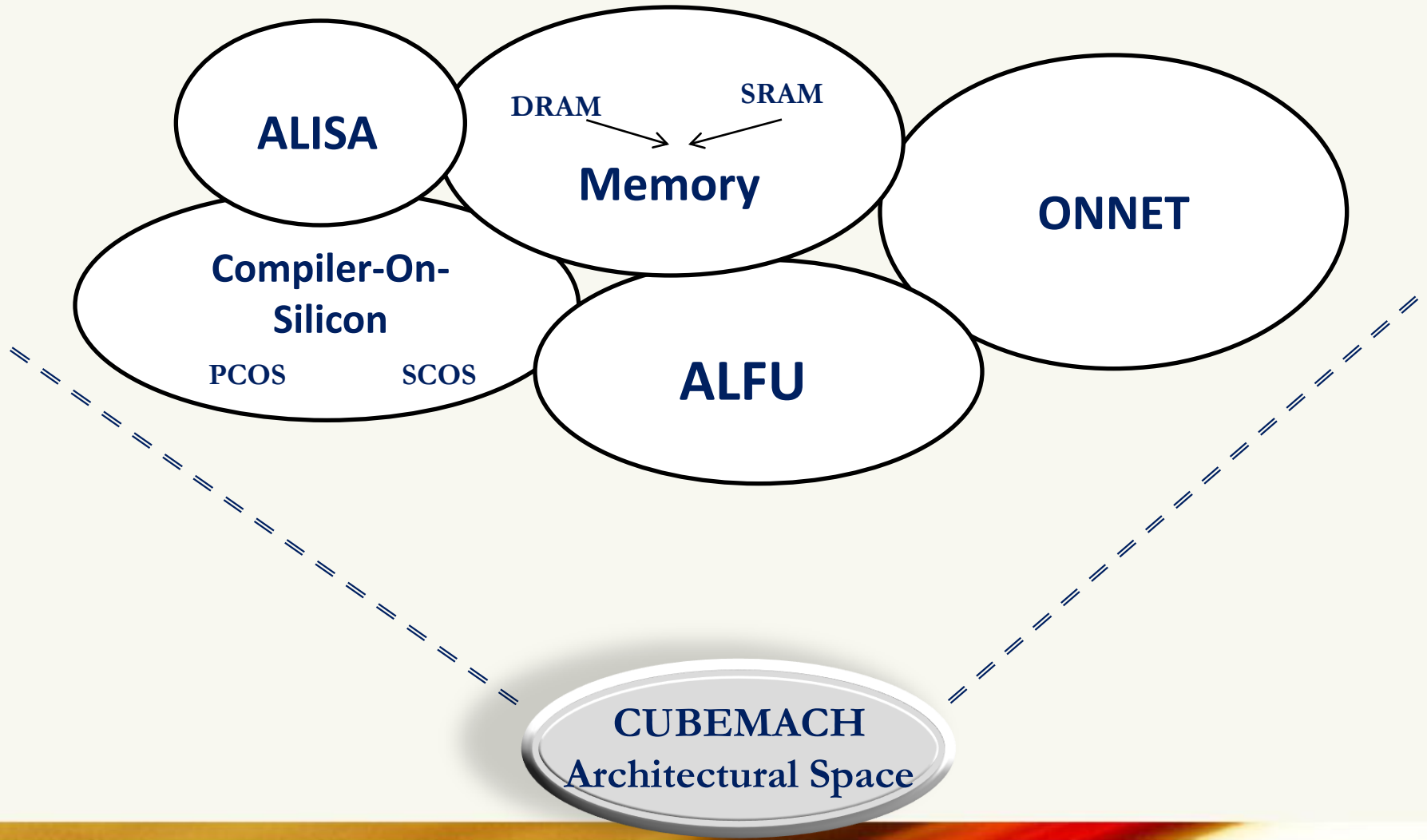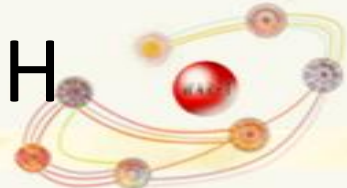  - Manufacturing and Operational Cost reduction

# Overview

- Custom Built Heterogeneous Multi-Core Architectures (CUBEMACH)

- Design Space

  - Architectural Space

  - Optimization Space

    - Customer Vendor Interaction

  - Simulation Space

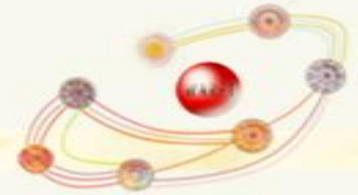- CUBEMACH Design and Simulation Tool Framework

- Conclusion

# CUBEMACH Design Paradigm

# Architectural Design Space - CUBEMACH

**ALISA**

DRAM ——→ ←—— SRAM

**Memory**

**ONNET**

**Compiler-On-Silicon**

PCOS          SCOS

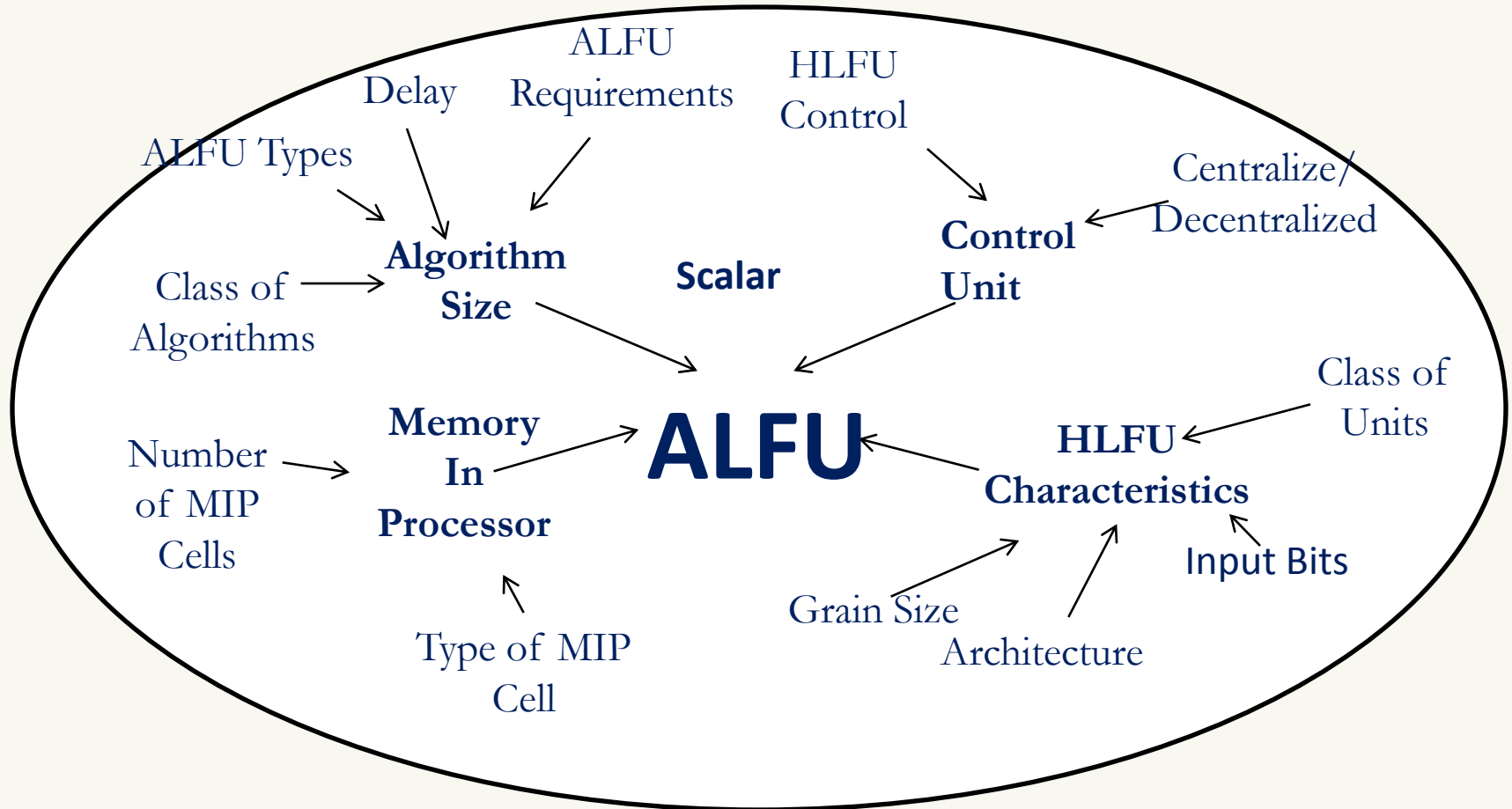**ALFU**
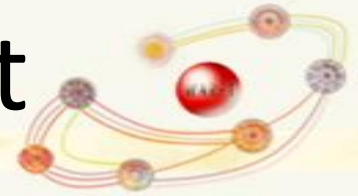
**CUBEMACH**
**Architectural Space**
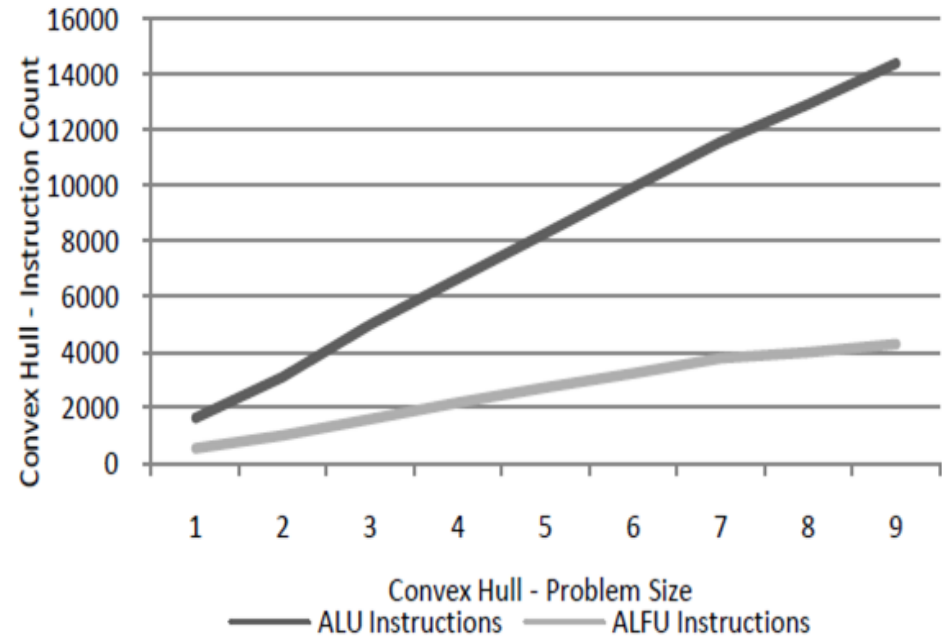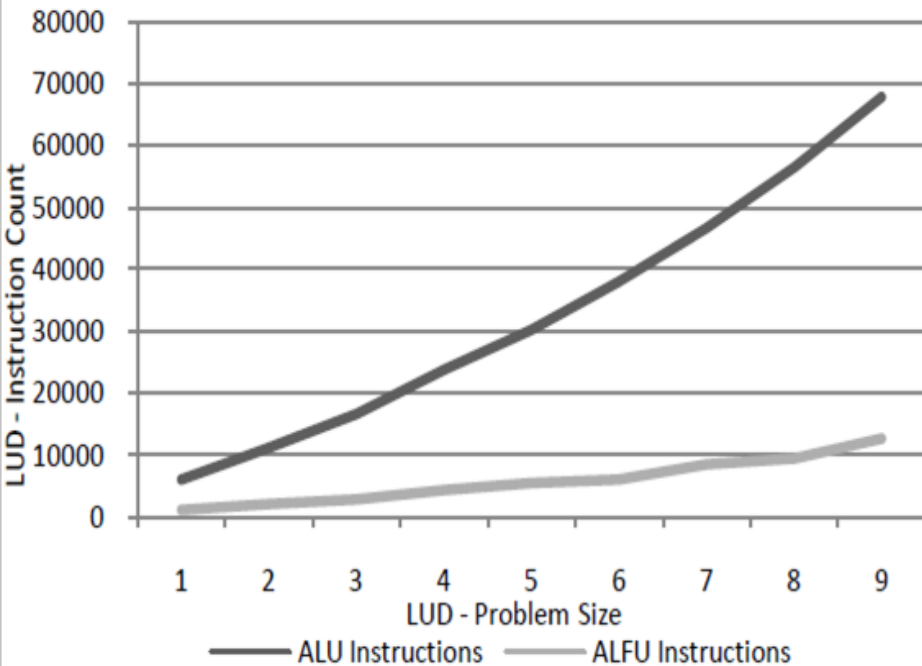
# Architectural Space

- Why ALU Why Not ALFU??
    - Hardwired units
    - Design : Homogeneously Structured
    - Reduced Instruction Generation & Fetches : Employ a Higher Level ISA
    - Reduced memory-functional unit interaction
    - Helps execute multiple applications without space & time sharing

# Algorithm Level Functional Unit

# ALU vs ALFU Instruction Generation Results

# Sample Algorithm Level Functional Units

- **Scalar Units**
  - Scalar Adder / Subtractor
  - Scalar Multiplier
  - Scalar Divider
  - Comparator
  - Sorter
  - Multiple Operand Adder
  - Min / Max Finder

- **Vector Units**
  - Inner Product

- **Matrix Centric Units**
  - Matmul
  - Matadd
  - Chain Matadd

- **Graph Theoretic Units**
  - Graph Traversal Unit – BFS, DFS
  - KL Graph Partitioning

## ALISA – Algorithm Level Instruction Set Architecture

- Algorithm Level Instructions

- Triggers ALFUS

- ALISA $\longrightarrow$ Multiple VLIWs
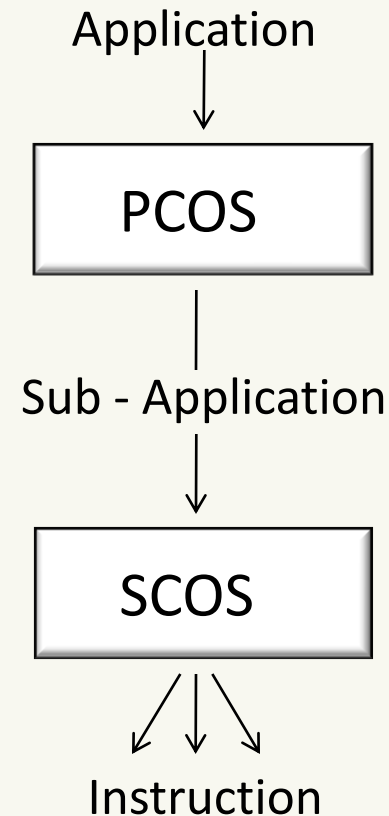
- ALISA for heterogeneous multi-cores

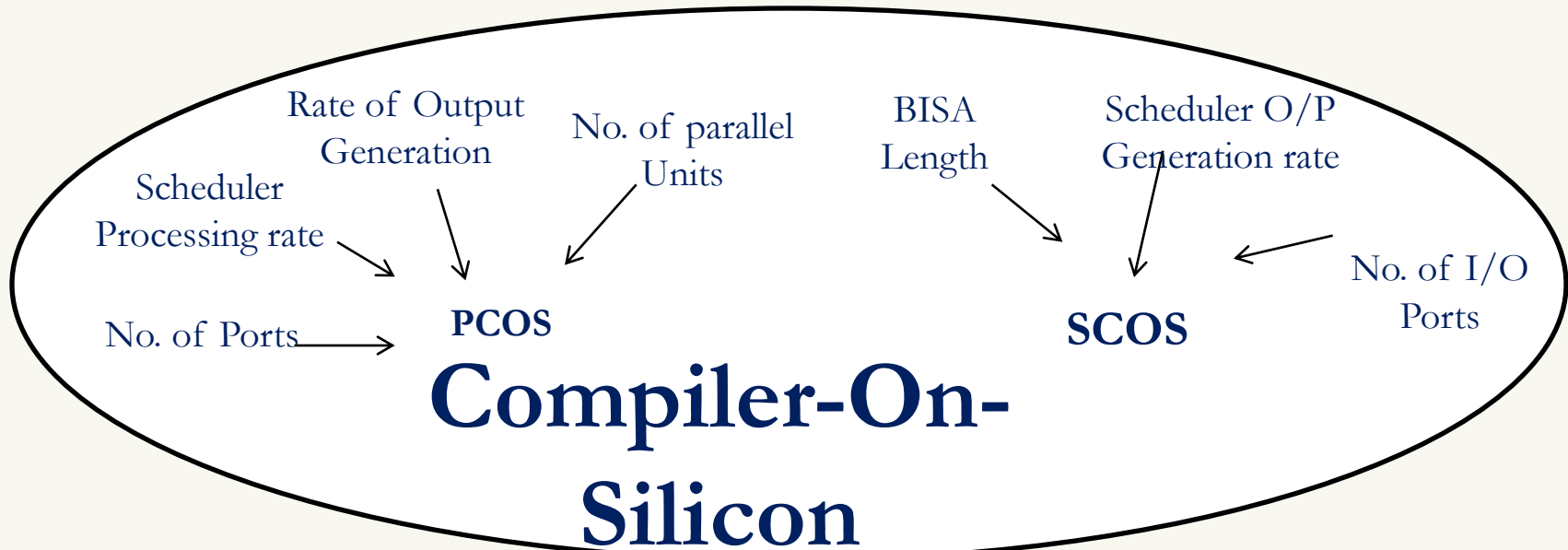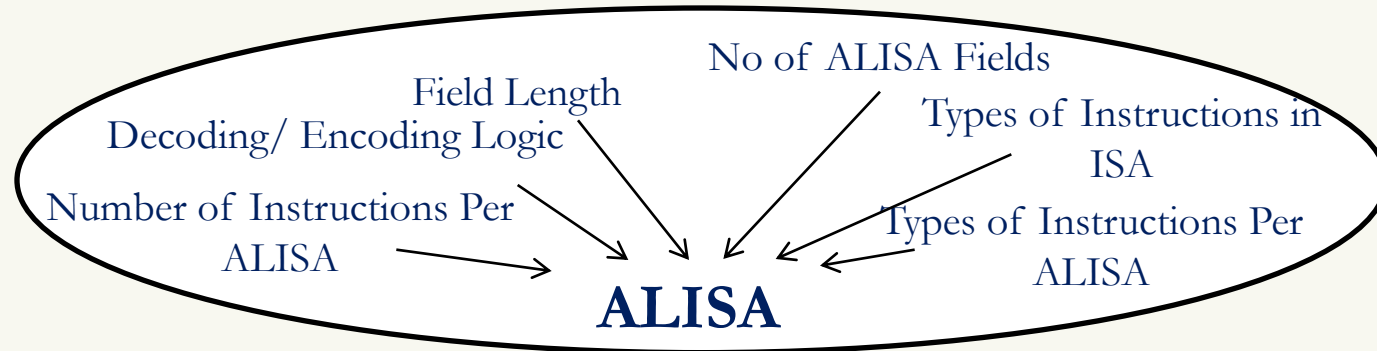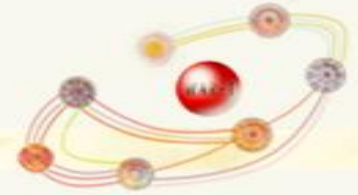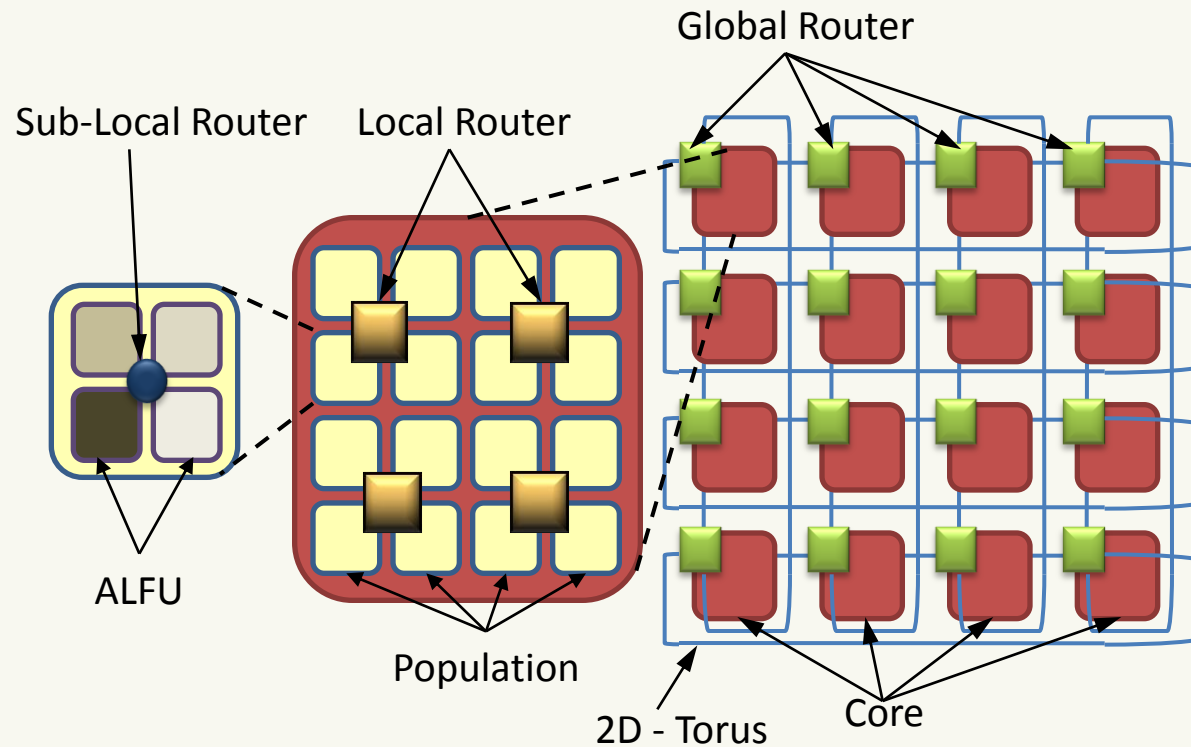| VLIW 1 | VLIW 2 | VLIW 3 | VLIW 4 |
|--------|--------|--------|--------|

**ALISA**

# Hierarchical Compilation Scheme

- PCOS Partitions A Problem Into Sub-Problems – Level 1

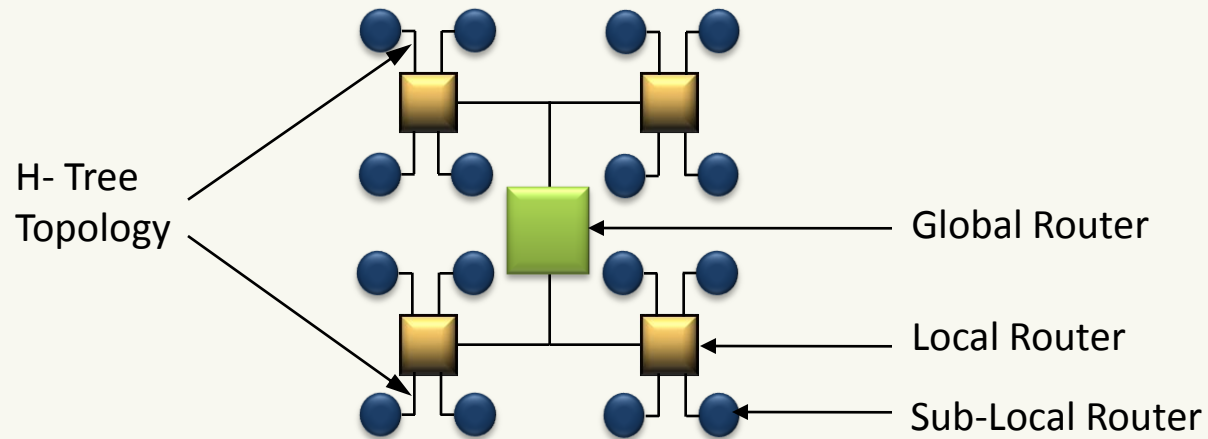- SCOS Partitions The Sub-Problems Into ALFU Level Instruction – Level 2

Application

↓

PCOS

↓

Sub - Application

↓

SCOS

↓ ↓ ↓

Instruction

# ALISA & Compiler On Silicon

No of ALISA Fields

Field Length

Decoding/ Encoding Logic

Types of Instructions in ISA

Number of Instructions Per ALISA

Types of Instructions Per ALISA

**ALISA**

Rate of Output Generation

No. of parallel Units

BISA Length

Scheduler O/P Generation rate

Scheduler Processing rate

No. of I/O Ports

No. of Ports

**PCOS**

**SCOS**

**Compiler-On-Silicon**

# ON-Node-Network Architecture

# ON-Node-Network Architecture

H- Tree
Topology

Global Router

Local Router

Sub-Local Router

# Comparison of Conventional NOCs with ONNET

|  | ONNET | Conventional NOCs |
|---|---|---|
| Type of Switch | MIN | Crossbar |
| Number of Routers | $N* \log_2 (N)$ | $N^2$ |
| Hierarchy | Yes | No |
| Switching Latency | $\log_2$(Number of Inputs) * Switch Delay | Number of Inputs * Switch Delay |

# On Node Network Architecture



Decoding Rate → Address Decoding
No. of Decoders → Address Decoding
I/O Port → Address Decoding
Length of Stack → Address Decoding
Data Rate → Routers
Route Location → Routers
Buffer/Stack Size → Routers
Logical Grouping → Organization
HLFU Count → Organization
No. of Buffers → Packetization
Packet Size → Packetization
I/P Data Size → Packetization
Word Length → Packetization
Output Data Size → Input/Output
Buffer Size
Destination ID
HLFU ID
Path Latency
Input Traffic → Packet Switching
Type of MIN → Packet Switching
Destination Router ID → Packet Switching
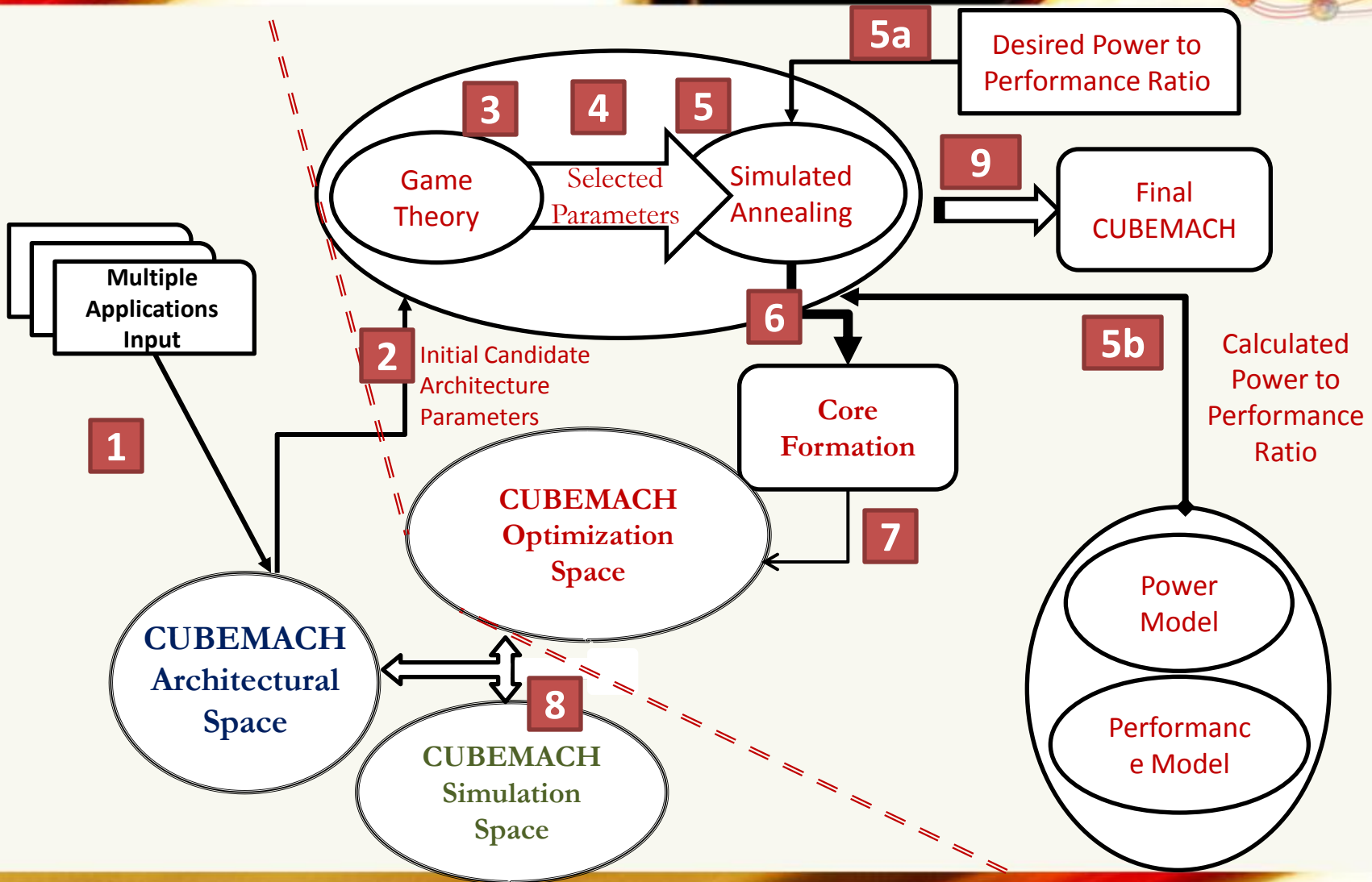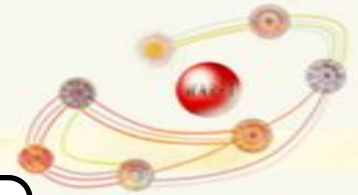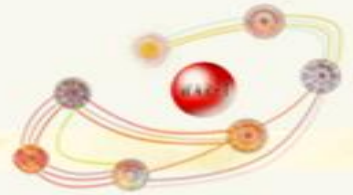
**ONNET**

# Overview

- Motivation
- Custom Built Heterogeneous Multi-Core Architectures (CUBEMACH)
- Design Space
  - Architectural Space
  - Optimization Space
    - Customer Vendor Interaction
  - Simulation Space
- CUBEMACH Design and Simulation Tool Framework
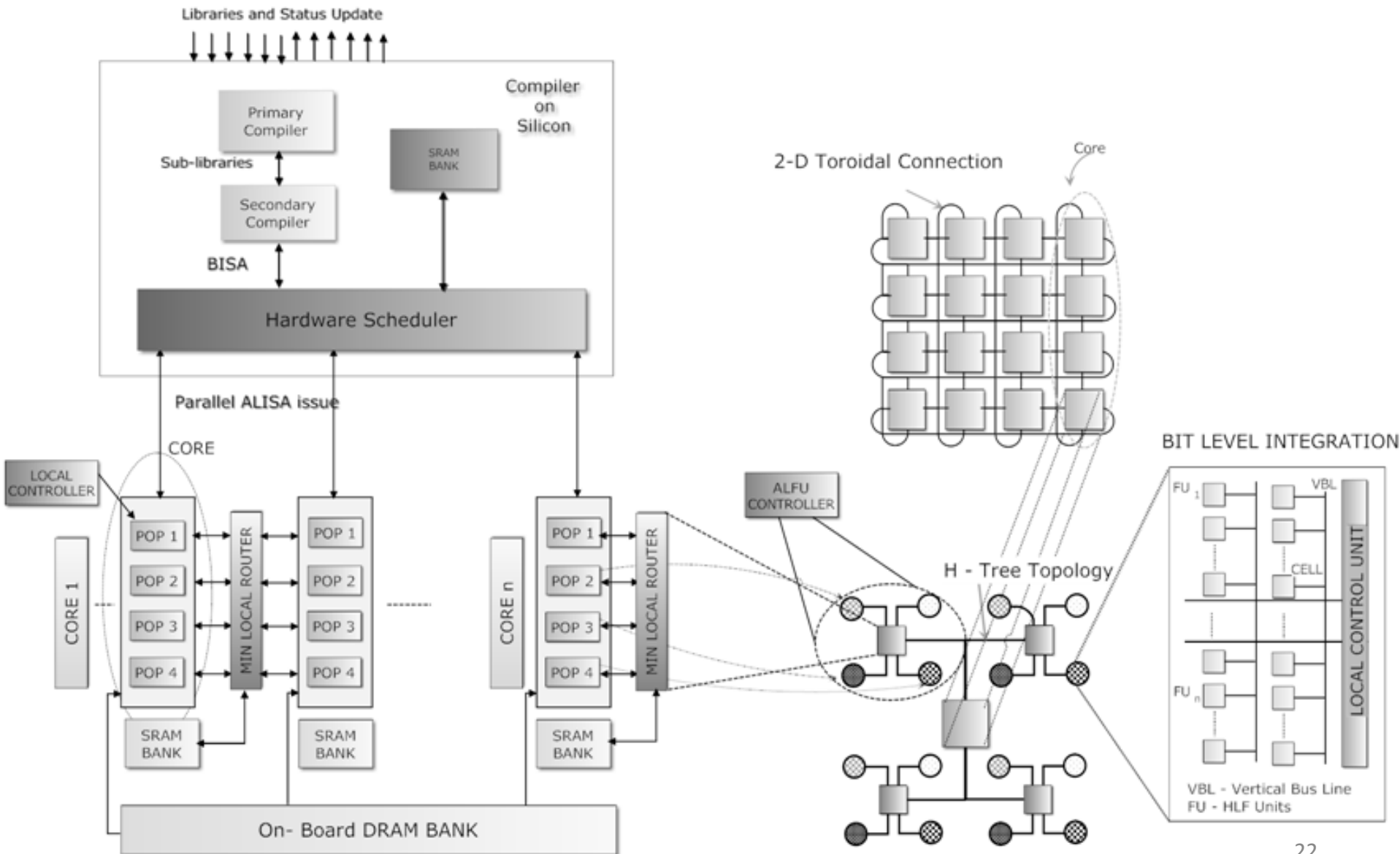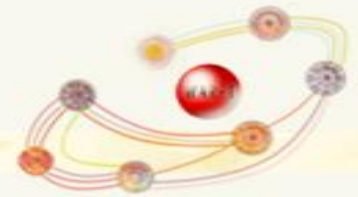- Conclusion

# Optimization Space

# Optimization Space

- Generates Optimized CUBEMACH for input specifications such as,
  - Power – Performance – Cost
  - Initial Architecture
- Power and Performance Model
- Uses GT and SA for optimization of Power and performance
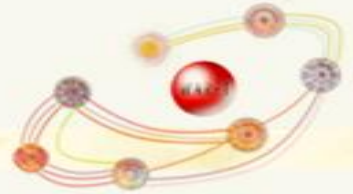- Uses KL For Core Grouping

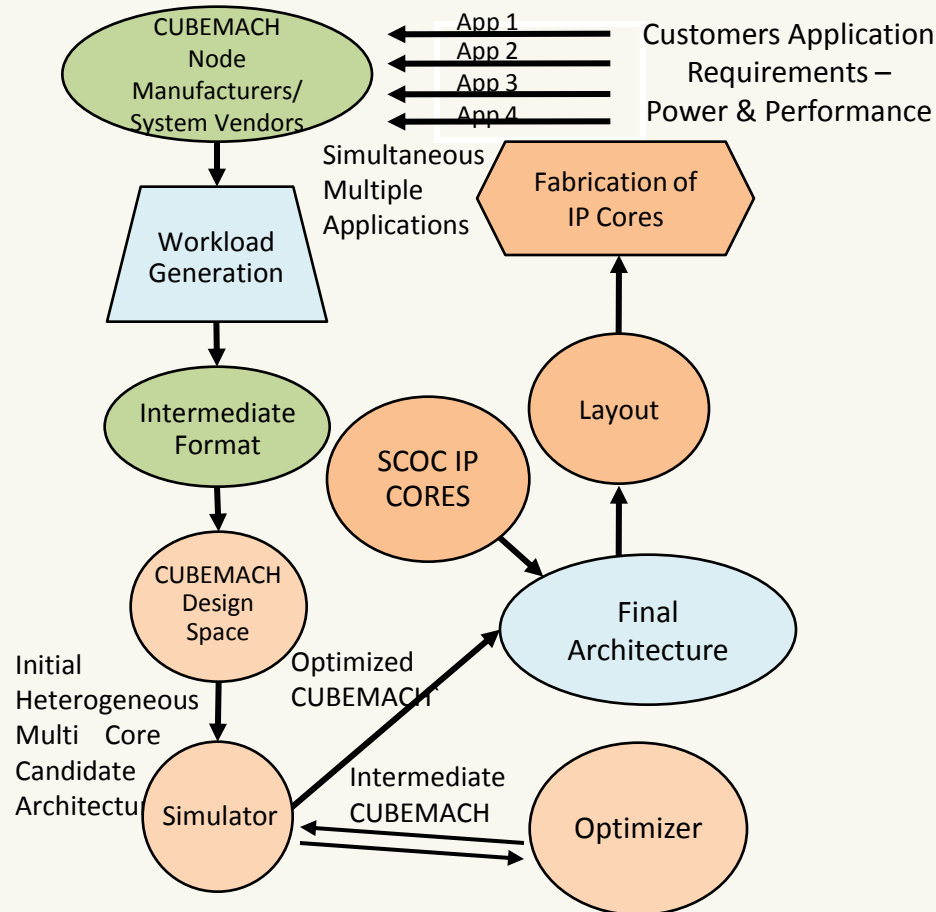# Sample CUBEMACH Architecture

# CUBEMACH Design Implementation : Supercomputer On Chip (SCOC) IP Cores
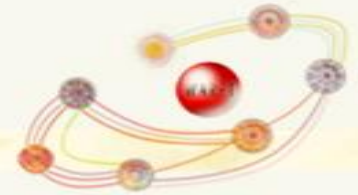
# SCOC IP Cores

- ALFUs  designed as SCOC  IP Cores

- Soft IP Core

- Coarse-grained Reusable Soft IP Cores
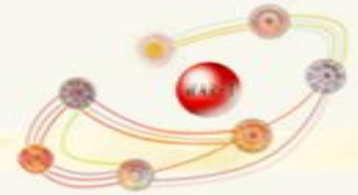
- Scalable IP Cores

# Customer Vendor Interaction
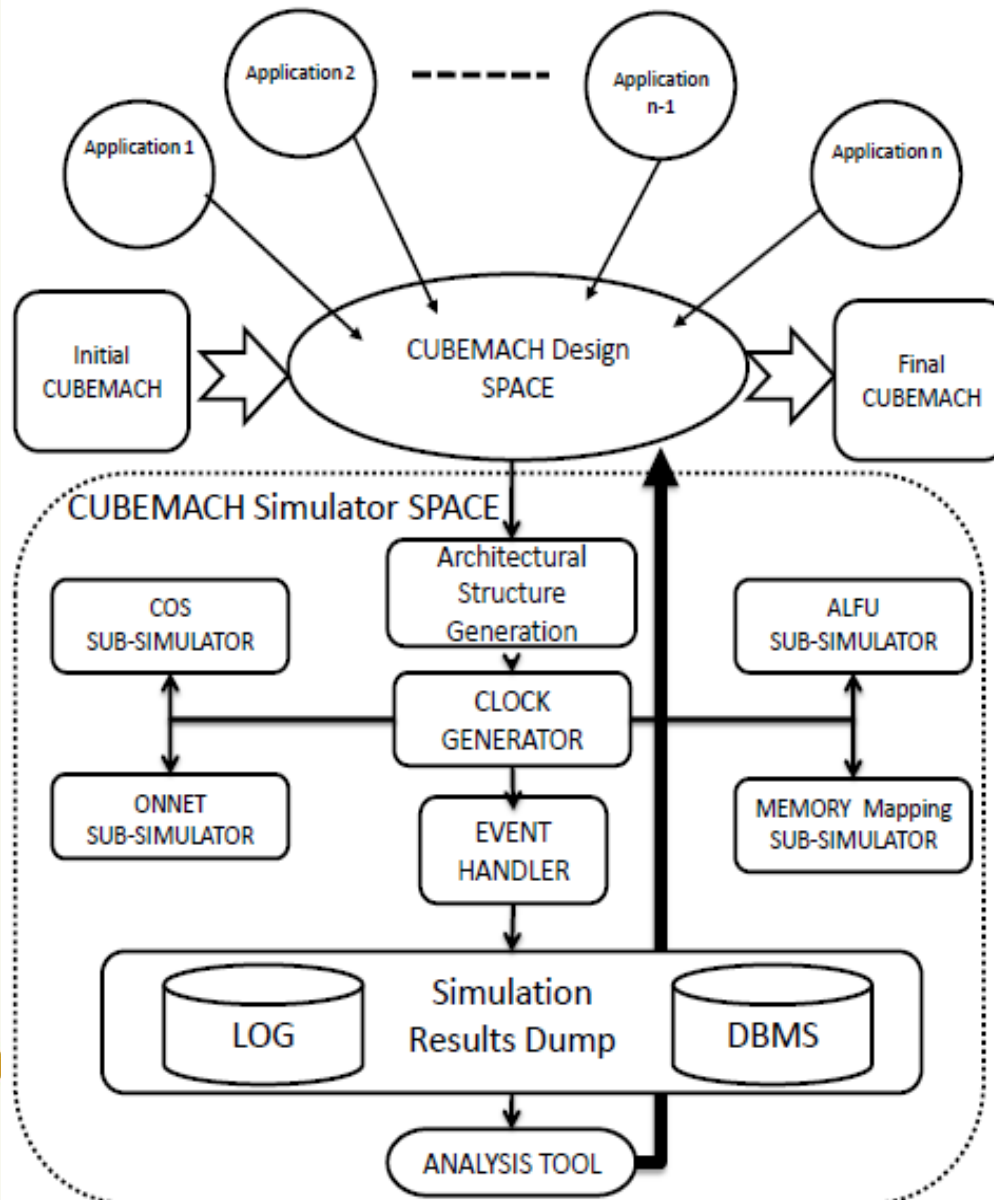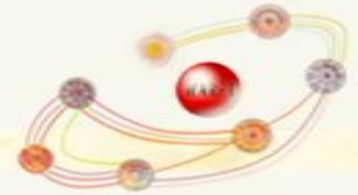
# Overview

- Motivation
- Custom Built Heterogeneous Multi-Core Architectures (CUBEMACH)
- Design Space
  - Architectural Space
  - Optimization Space
    - Customer Vendor Interaction
  - Simulation Space
- CUBEMACH Design and Simulation Tool Framework
- Conclusion

# CUBEMACH Simulator

- pThread based Simulator

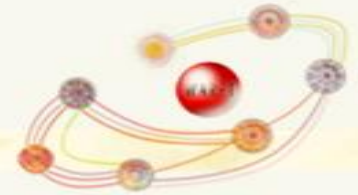- Evaluates candidate CUBEMACH Architecture

- Feed results to CUBEMACH Optimizer

- CUBEMACH Optimization Engine (COE) produces Optimized Architecture

- Simulation & Optimization : An iterative process

- Consists of

    ALFU Sub-Simulator       COS Sub-Simulator
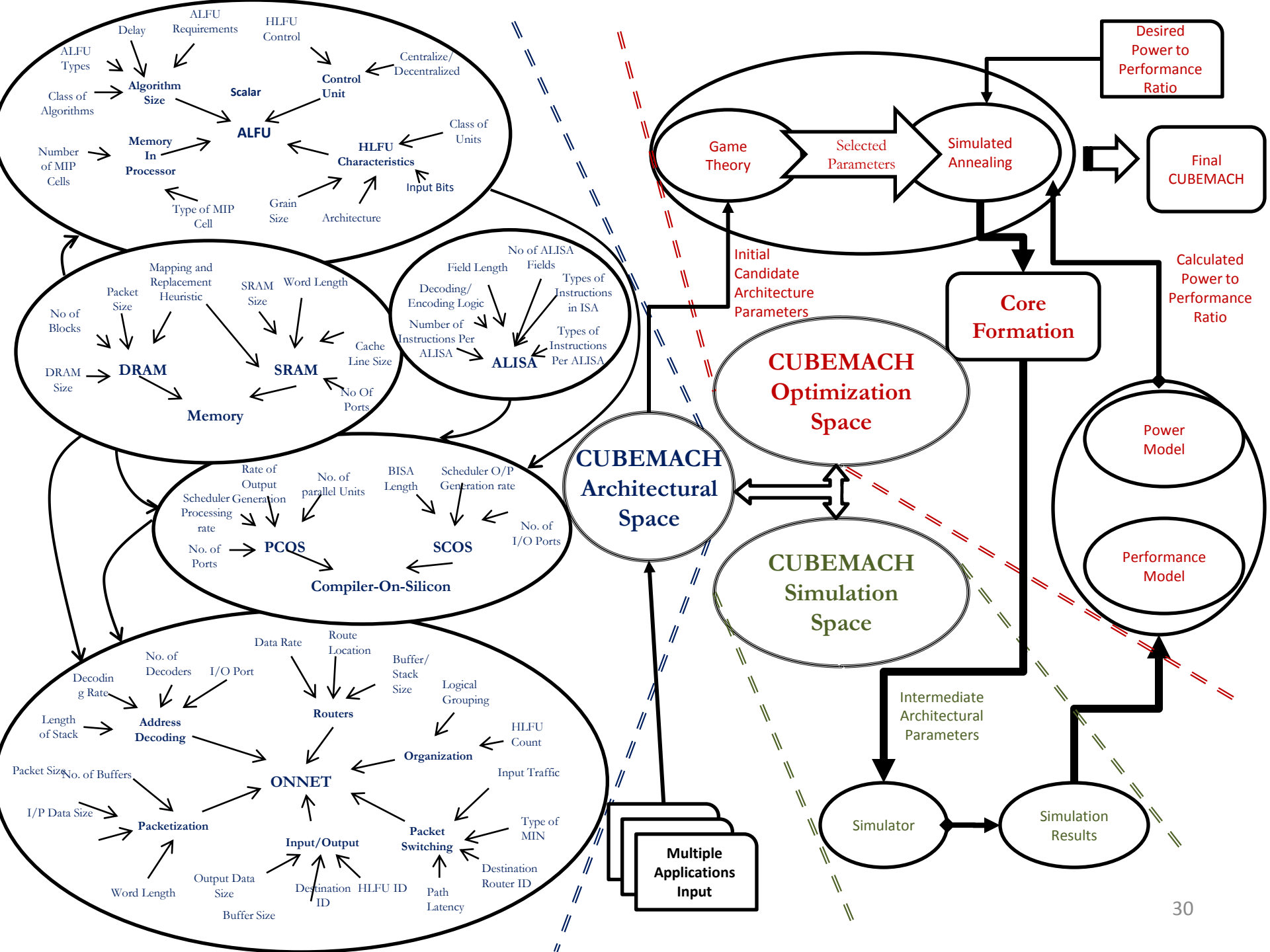
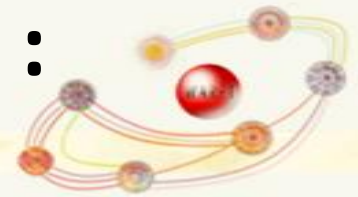    ONNET Sub-Simulator    Memory Sub-Simulator
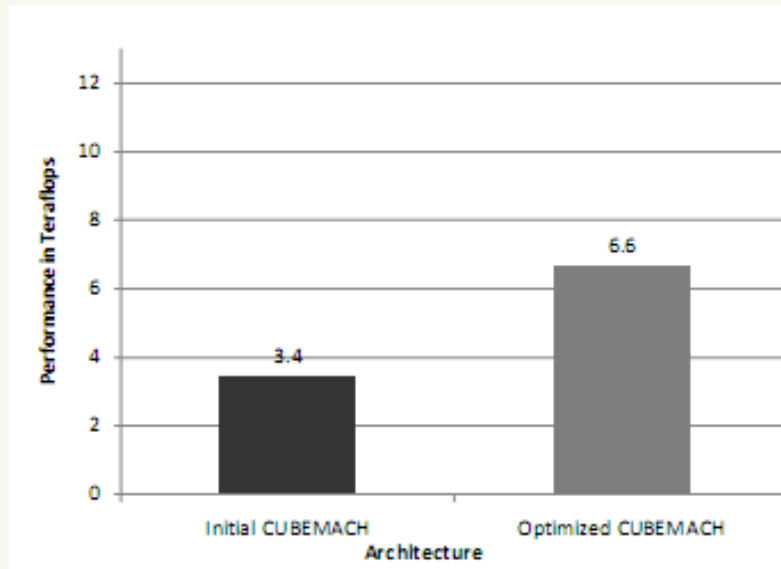
# CUBEMACH Simulator

# What we have seen . . .

## Integrated CUBEMACH Design Paradigm …

ALFU Types
Delay
ALFU Requirements
HLFU Control
Centralize/ Decentralized
Class of Algorithms
**Algorithm Size**
**Scalar**
**Control Unit**
Number of MIP Cells
**Memory In Processor**
**ALFU**
**HLFU Characteristics**
Class of Units
Type of MIP Cell
Grain Size
Architecture
**Input Bits**

Mapping and Replacement Heuristic
Packet Size
SRAM Size
Word Length
No of Blocks
DRAM Size
**DRAM**
**SRAM**
Cache Line Size
**Memory**
No Of Ports

No of ALISA Fields
Field Length
Types of Instructions in ISA
Decoding/ Encoding Logic
Number of Instructions Per ALISA
**ALISA**
Types of Instructions Per ALISA

Rate of Output
Scheduler Generation
No. of parallel Units
BISA Length
Scheduler O/P Generation rate
Processing rate
No. of Ports
**PCQS**
**SCOS**
No. of I/O Ports
**Compiler-On-Silicon**

Data Rate
Route Location
Buffer/ Stack Size
Decoding Rate
No. of Decoders
I/O Port
Logical Grouping
Length of Stack
**Address Decoding**
**Routers**
HLFU Count
Packet Size
No. of Buffers
**Organization**
Input Traffic
I/P Data Size
**ONNET**
**Packetization**
**Packet Switching**
Type of MIN
Word Length
Output Data Size
**Input/Output**
Destination Router ID
Buffer Size
Destination ID
HLFU ID
Path Latency

**CUBEMACH Architectural Space**

**CUBEMACH Optimization Space**

**CUBEMACH Simulation Space**

Game Theory
Selected Parameters
Simulated Annealing
Final CUBEMACH

Desired Power to Performance Ratio

Initial Candidate Architecture Parameters

**Core Formation**

Calculated Power to Performance Ratio

Power Model

Performance Model

Intermediate Architectural Parameters

Simulator
Simulation Results

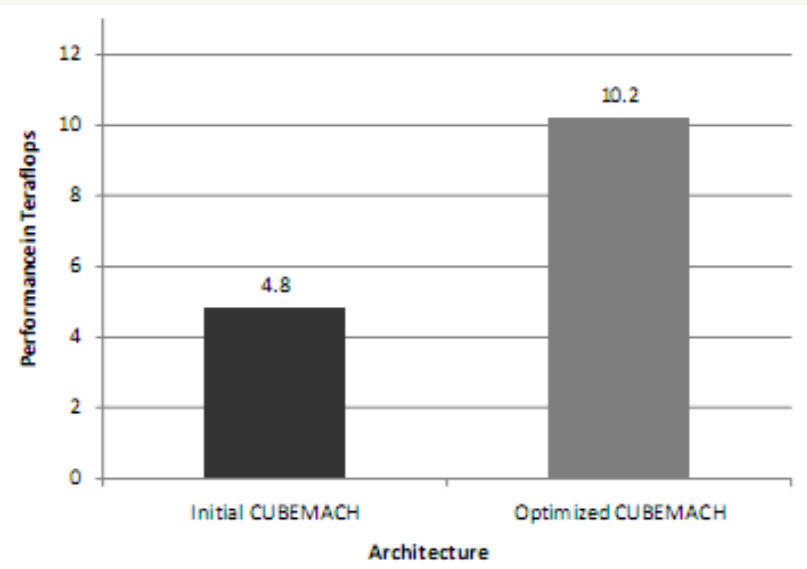**Multiple Applications Input**

30

# Sample CUBEMACH Architecture :
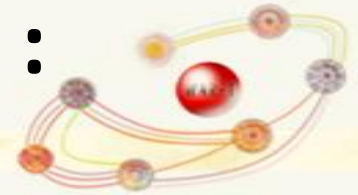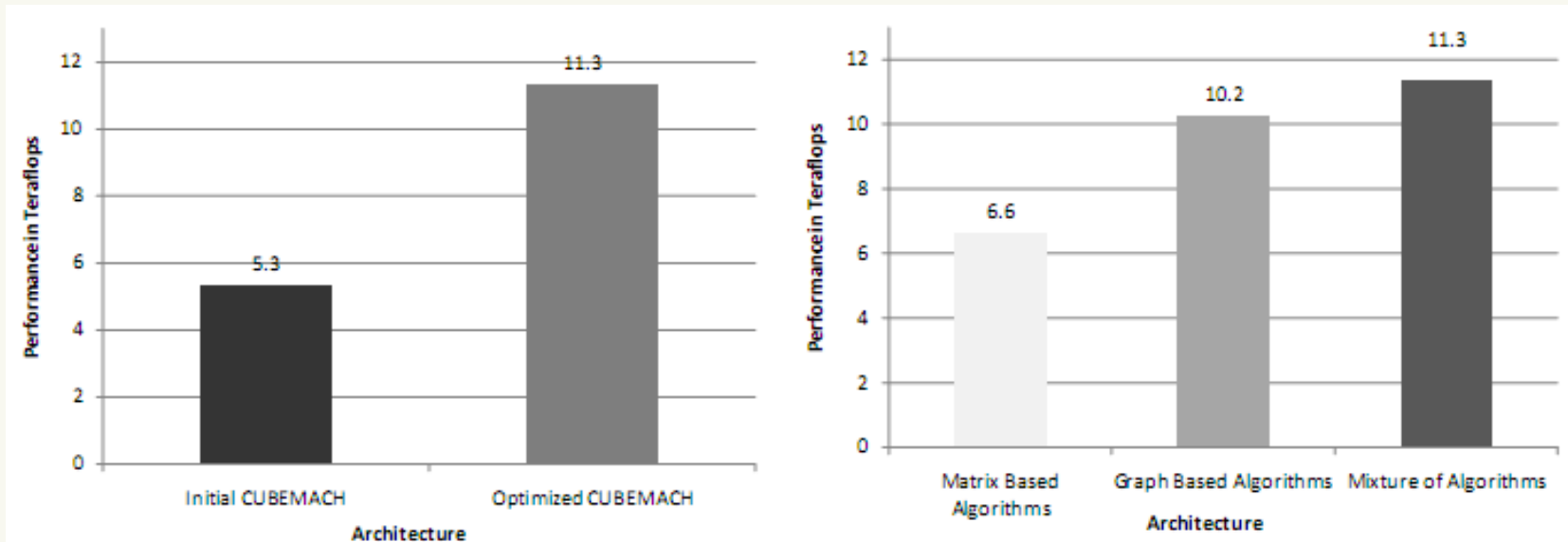# Simulation Results



Matrix Based Algorithms

Graph Based Algorithms

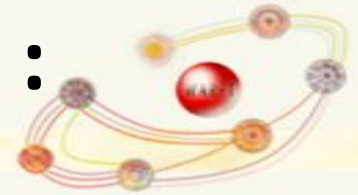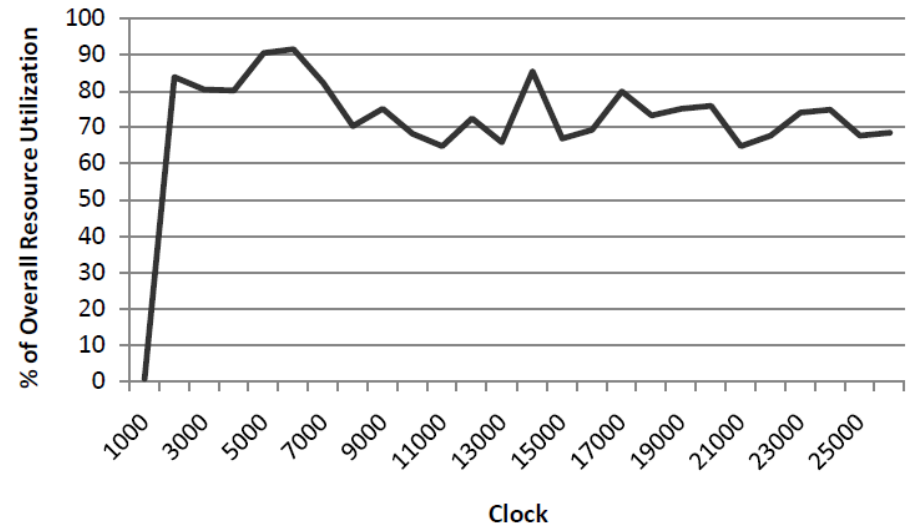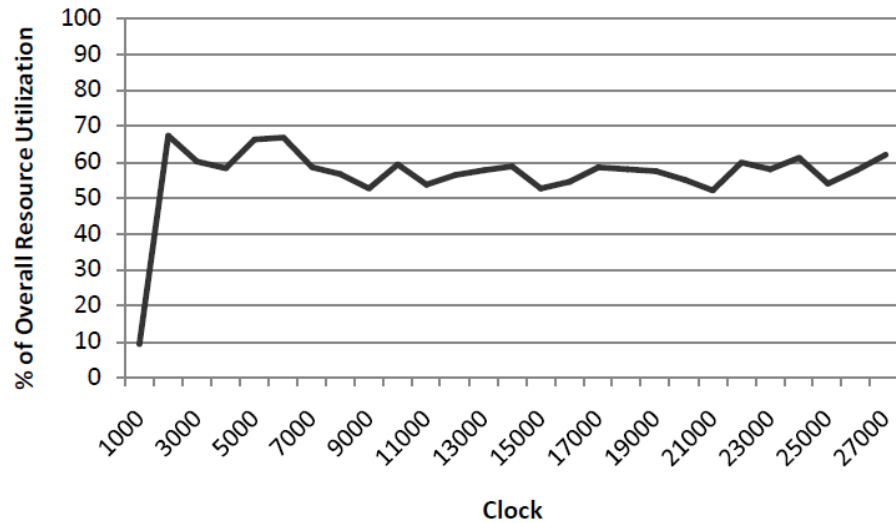# Sample CUBEMACH Architecture :

## Simulation Results



Mixture of Algorithms

Comparison of Performance delivered by Optimized Architectures for corresponding types of Algorithms
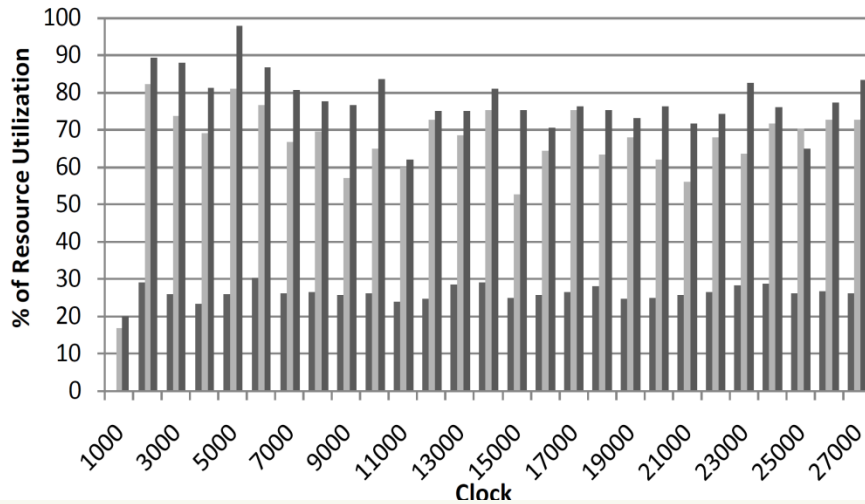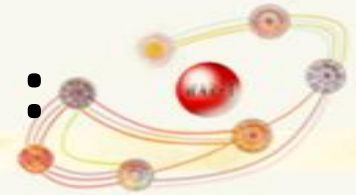
# Sample CUBEMACH Architecture :
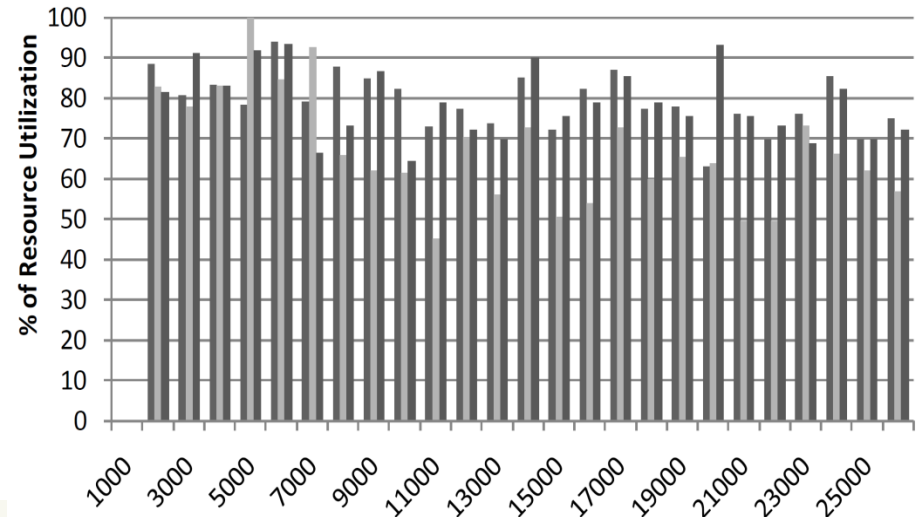# Simulation Results



Overall Resource Utilization of :
(i)   Initial CUBEMACH Architecture       : Mean = 59 %
(ii)  Optimized CUBEMACH Architecture : Mean = 74 %

# Sample CUBEMACH Architecture : Simulation Results
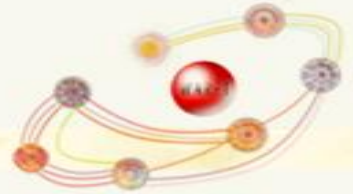


In Initial Candidate CUBEMACH Architecture,
- Matrix ALFUS – low usage
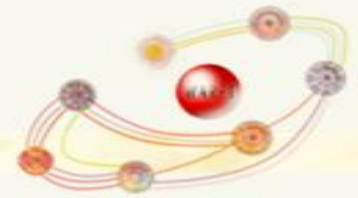- Scalar ALFUS – average usage
- Graph ALFUS – high usage

In Optimized Candidate CUBEMACH Architecture,
- Matrix ALFUS – high usage
- Scalar ALFUS – high usage
- Graph ALFUS – high usage

# Conclusion

- Custom Built Heterogeneous Multi-Core Architectures (CUBEMACH) promises,
  - Increased Resource Utilization
  - Multiple application flavored architectures
  - Elimination of Space Time Sharing at the Quantum Level during Multiple Application Execution (without multiprogramming)
  - Manufacturing and Running Cost reduction

# Thank You

# Questions??

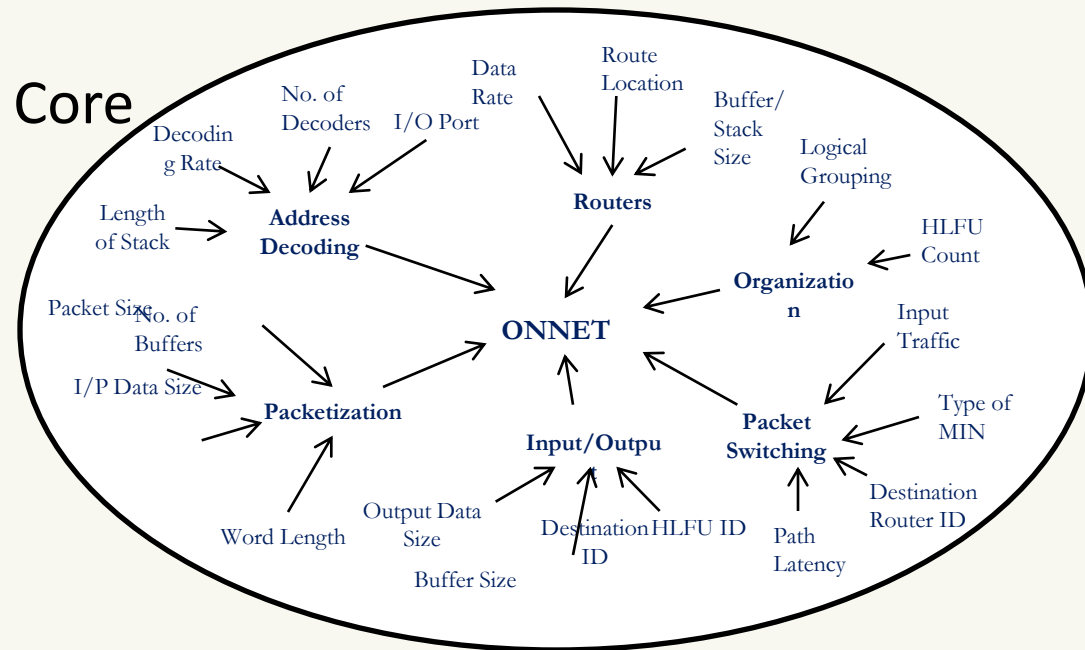# Customizable Compiler-On-Silicon

- What  Compiler-On-Silicon?

- Why do we  need Compiler-On-Silicon ?

- Why go for Customizable Compiler-On-Silicon ?
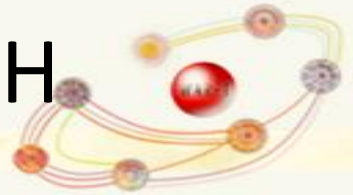
# ONNET

Architecture uses -

- Multistage Interconnect Network

- Hardware Packetization Unit

- ONNET Design Space
  - H-Tree Structure within a Core
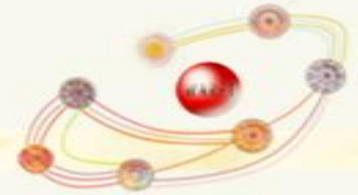  - 2D Torus Across Cores
  - MIN Type

# Architectural Design Space - CUBEMACH

- ALFU – Algorithm Level Functional Units

- BISA – Backbone Instruction Set Architecture

- COS – Compiler On Silicon

- ONNET – On Node Network

- Novel Cache Mapping Scheme

- SCOC IP Cores : Achieving cost effectiveness
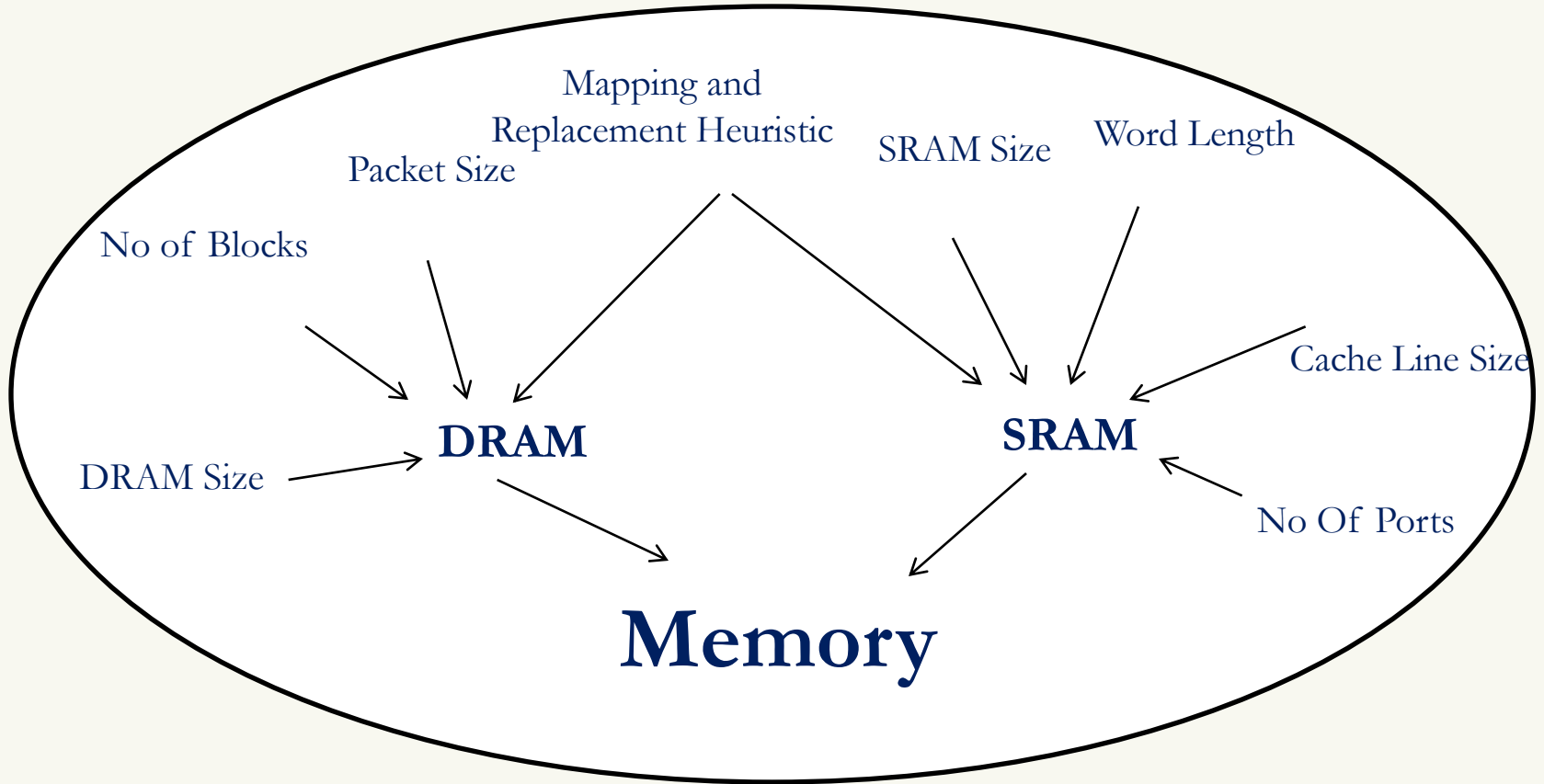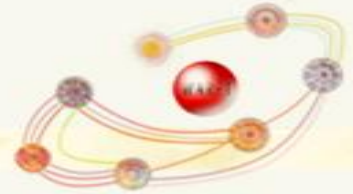
( Super Computer On Chip - IP Cores)

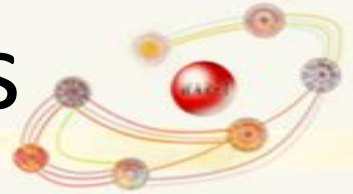# On Node Network Architecture

## Features  -

- Communication across heterogeneous  multi-cores

- Data requirements of diverse ALFUs

- High bandwidth

- Scalable

- Hierarchical Network-On-Chip

# Memory



No of Blocks

Packet Size

Mapping and Replacement Heuristic

SRAM Size

Word Length

**DRAM**

DRAM Size

**SRAM**

Cache Line Size

No Of Ports

**Memory**

# Advantages of SCOC IP Cores

- Fully Customizable

- Greatly reduces Design-Turnaround-Time

- Physically Design Friendly
  - Constraints of Area, Power and Performance

- Constrained & Rigid Design Methodology